



# DAGSTUHL REPORTS

**Volume 12, Issue 6, June 2022**

Theories of Programming (Dagstuhl Seminar 22231) <i>Thomas D. LaToza, Amy Ko, David C. Shepherd, and Dag Sjøberg</i> .....	1
Efficient and Equitable Natural Language Processing in the Age of Deep Learning (Dagstuhl Seminar 22232) <i>Jesse Dodge, Iryna Gurevych, Roy Schwartz, and Emma Strubell</i> .....	14
Human-Game AI Interaction (Dagstuhl Seminar 22251) <i>Dan Ashlock, Setareh Maghsudi, Diego Perez Liebana, Pieter Spronck, and Manuel Eberhardinger</i> .....	28
Visualization Empowerment: How to Teach and Learn Data Visualization (Dagstuhl Seminar 22261) <i>Benjamin Bach, Sheelagh Carpendale, Uta Hinrichs, and Samuel Huron</i> .....	83
Human-Centered Artificial Intelligence (Dagstuhl Seminar 22262) <i>Wendy E. Mackay, John Shawe-Taylor, and Frank van Harmelen</i> .....	112

ISSN 2192-5283

*Published online and open access by*

Schloss Dagstuhl – Leibniz-Zentrum für Informatik GmbH, Dagstuhl Publishing, Saarbrücken/Wadern, Germany. Online available at <https://www.dagstuhl.de/dagpub/2192-5283>

*Publication date*

January, 2023

*Bibliographic information published by the Deutsche Nationalbibliothek*

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <https://dnb.d-nb.de>.

*License*

This work is licensed under a Creative Commons Attribution 4.0 International license (CC BY 4.0).



In brief, this license authorizes each and everybody to share (to copy, distribute and transmit) the work under the following conditions, without impairing or restricting the authors' moral rights:

- Attribution: The work must be attributed to its authors.

The copyright is retained by the corresponding authors.

*Aims and Scope*

The periodical *Dagstuhl Reports* documents the program and the results of Dagstuhl Seminars and Dagstuhl Perspectives Workshops.

In principal, for each Dagstuhl Seminar or Dagstuhl Perspectives Workshop a report is published that contains the following:

- an executive summary of the seminar program and the fundamental results,
- an overview of the talks given during the seminar (summarized as talk abstracts), and
- summaries from working groups (if applicable).

This basic framework can be extended by suitable contributions that are related to the program of the seminar, e. g. summaries from panel discussions or open problem sessions.

*Editorial Board*

- Elisabeth André
- Franz Baader
- Daniel Cremers
- Goetz Graefe
- Reiner Hähnle
- Barbara Hammer
- Lynda Hardman
- Oliver Kohlbacher
- Steve Kremer
- Rupak Majumdar
- Heiko Mantel
- Albrecht Schmidt
- Wolfgang Schröder-Preikschat
- Raimund Seidel (*Editor-in-Chief*)
- Heike Wehrheim
- Verena Wolf
- Martina Zitterbart

*Editorial Office*

Michael Wagner (*Managing Editor*)  
Michael Didas (*Managing Editor*)  
Jutka Gasiorowski (*Editorial Assistance*)  
Dagmar Glaser (*Editorial Assistance*)  
Thomas Schillo (*Technical Assistance*)

*Contact*

Schloss Dagstuhl – Leibniz-Zentrum für Informatik  
Dagstuhl Reports, Editorial Office  
Oktavie-Allee, 66687 Wadern, Germany  
[reports@dagstuhl.de](mailto:reports@dagstuhl.de)  
<https://www.dagstuhl.de/dagrep>

Digital Object Identifier: 10.4230/DagRep.12.6.i

# Theories of Programming

Thomas D. LaToza<sup>\*1</sup>, Amy Ko<sup>\*2</sup>, David C. Shepherd<sup>\*3</sup>,  
Dag Sjøberg<sup>\*4</sup>, and Benjamin Xie<sup>†5</sup>

- 1 George Mason University – Fairfax, US. [tlatoya@gmu.edu](mailto:tlatoya@gmu.edu)
- 2 University of Washington – Seattle, US. [ajko@uw.edu](mailto:ajko@uw.edu)
- 3 Virginia Commonwealth University – Richmond, US. [shepherdd@vcu.edu](mailto:shepherdd@vcu.edu)
- 4 University of Oslo, NO. [dagsj@ifi.uio.no](mailto:dagsj@ifi.uio.no)
- 5 University of Washington – Seattle, US. [bxie@uw.edu](mailto:bxie@uw.edu)

---

## Abstract

Much of computer science research focuses on techniques to make programming easier, better, less error prone, more powerful, and even more just. But rarely do we try to explain any of these challenges. Why is programming hard? Why is it slow? Why is it error prone? Why is it powerful? How does it do harm? These why and how questions are what motivated the Dagstuhl Seminar 22231 on Theories of Programming. This seminar brought together 28 CS researchers from domains most concerned with programming human and social activities: software engineering, programming languages, human-computer interaction, and computing education. Together, we sketched new theories of programming and considered the role of theories more broadly in programming.

**Seminar** June 6–10, 2022 – <http://www.dagstuhl.de/22231>

**2012 ACM Subject Classification** Social and professional topics → Computing education; Human-centered computing → Human computer interaction (HCI); Software and its engineering

**Keywords and phrases** computing education, human-computer interaction, programming languages, software engineering, theories of programming

**Digital Object Identifier** 10.4230/DagRep.12.6.1

## 1 Executive Summary


*Benjamin Xie (University of Washington – Seattle, US)*

*Amy Ko (University of Washington – Seattle, US)*

*Thomas D. LaToza (George Mason University – Fairfax, US)*

*David C. Shepherd (Virginia Commonwealth University – Richmond, US)*

*Dag Sjøberg (University of Oslo, NO)*

**License**  Creative Commons BY 4.0 International license  
© Benjamin Xie, Amy Ko, Thomas D. LaToza, David C. Shepherd, and Dag Sjøberg

Mature scientific disciplines are characterized by their theories, synthesizing what is known about phenomena into forms which generate falsifiable predictions about the world. In computer science, the role of synthesizing ideas has largely been through formalisms that describe how programs compute. However, just as important are scientific theories about how programmers write these programs. For example, software engineering research has increasingly begun gathering data, through observations, surveys, interviews, and analysis of artifacts, about the nature of programming work and the challenges developers face, and evaluating novel programming tools through controlled experiments with software developers.

---

\* Editor / Organizer

† Editorial Assistant / Collector



Except where otherwise noted, content of this report is licensed under a Creative Commons BY 4.0 International license

Theories of Programming, *Dagstuhl Reports*, Vol. 12, Issue 6, pp. 1–13

Editors: Thomas D. LaToza, Amy Ko, David C. Shepherd, and Dag Sjøberg



Dagstuhl Reports

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Computer science education and human-computer interaction research has done similar work, but for people with different levels of experience and ages learning to write programs. But data from such empirical studies is often left isolated, rather than combined into useful theories which explain all of the empirical results. This lack of theory makes it harder to predict in which contexts programming languages, tools, and pedagogy will actually help people successfully write and learn to create software.

Computer science needs scientific theories that synthesize what we believe to be true about programming and offer falsifiable predictions. Whether or not a theory is ultimately found to be consistent with evidence or discarded, theories offer a clear statement about our current understanding, helping us in prioritizing studies, generalizing study results from individual empirical results to more general understanding of phenomena, and offering the ability to design tools in ways that are consistent with current knowledge.

Dagstuhl Seminar 22231 on *Theories of Programming* explored the creation and synthesis of scientific theories which describe the relationship between developers and code within programming and social activities. The seminar brought together researchers from software engineering, human-computer interaction, programming languages, and computer science education to exchange ideas about potential theories of programming. We identified and proposed theories that arose from many sources: untested but strongly-held beliefs, anecdotal observations, assumptions deeply embedded in the design of languages and tools, reviews of empirical evidence on programming, and applications of existing theories from psychology and related areas. Our aim was to bridge this gulf: formulating deeply-held beliefs into theories which are empirically testable and synthesizing empirical findings in ways that make predictions about programming tools and languages.

To achieve this aim, the seminar had three specific goals. 1) Bring together researchers with diverse expertise to find shared understanding. 2) Create a body of theories which make testable predictions about the effects of programming tools, languages, and pedagogy on developer behavior in specific contexts. 3) Propose future activities which can advance the use of theories, including identifying studies to conduct to test theories and ways to use theories to communicate research findings to industry.

During this seminar, a few short talks first reviewed the nature, creation, and use of theories as well as existing evidence about developer behavior during programming activities. The main activity of the seminar was working in small groups to sketch new theories of programming.

## Seminar Overview

The seminar was divided into the following sessions across four days in June 2022:

- Tuesday: welcome, what is theory, describing theories, critiquing theories
- Wednesday: brainstorming unexplained programming phenomena, sketching theories, getting feedback on theories, and refining theories
- Thursday: presenting theory sketches, discussing ways of sharing theories, and skeptically examining whether developing theories of programming is really worth the time
- Friday: reflecting on takeaways and departure

The seminar was organized by Thomas Latoza, Amy J. Ko, Dag Sjøberg, David Shepherd, and Anita Sarma. Anita later had to drop out, leaving Thomas, Amy, Dag, and David as the four organizers who were able to attend.

## What is theory?

The goal of this opening session was to find common ground on what theory was. To achieve this, each organizer gave short presentations related to theories.

Thomas identified how researchers used theories to generate falsifiable predictions about the world. He described common characteristics of theories as abstract, explanatory, relevant, and operationalizable. An example of a theory of programming he provided was how violating constraints cause defects or reduces code quality.

Amy described an interpretivist framing of theories, where theories were cultural and experience-based. Some theories were folk theories (e.g. code is magic, Not Invented Here, and spaces vs tabs for white space). Some theories were personal, such as programming as common sense machines and tinkering towards correctness. Other theories came from research communities such as ICSE. For example, a theory of programming is that we can copy and adapt code from another location in a program to fix bugs.

Dag drew guidelines between what was and was not a theory. He identified multiple examples of what were not theories: scientific laws were not theories because they were missing the “why;” trivial statements were also not theories. The building blocks of theories included constructs, propositions, explanations, and scope. Theories can help us explain surprising empirical results, while empirical results can help us support or refute certain theories. Finally, Dag noted how premature theorizing is likely to be wrong, but can still be useful.

David emphasized the importance of keeping theories practical. He defined a relationship from theorems to corollaries to examples and applications. He provided an example of using different representations in music for different use cases and users.

In open discussions and breakout groups, attendees identified additional nuances to theories. We noted how it is useful for theories to enable ease of communication or shared understanding. But by defining a vocabulary, theories can also limit the scope of explanation. We can also use theories to understand what we observe or to justify interventions. Finally, there was discussion about creating theories inductively, deductively, and/or abductively.

Common themes that arose from discussion include how theories are seldom used to justify the design of programming languages and tools, and how programming is a social endeavor and drawing upon social science research (e.g. psychology) can support theory building.

## Expressing Theories using a Theory Template

The goal of this session was to try to express theories using a theory template developed by the organizers. While the goal of this template was to support the creation of new theories, attendees used it to describe existing theories for this session. Attendees broke into five groups to attempt to apply the theory template to the following existing theories of programming:

- Asking and answering questions [1]
- Program comprehension as fact finding [2]
- Leaky abstractions [3]
- Information hiding [4]
- Theory of programming instruction [5]

After considering feedback from attendees, organizers revised the theory template. The revised theory template’s section headers and helper text are as follows:

1. *Theory's name*: Choose a name that is memorable, short, and descriptive.
2. *Summary*: In a few sentences, summarize the phenomena, constructs, relationships, and a concrete example, hypothesis, and study.
3. *Contributors*: Who has contributed to this theory? Add your name here.
4. *Phenomena*: What programming phenomena is your theory trying to explain? And in what scope (people, expertise, contexts, tools, etc.)? This description should just describe what is being explained, should not offer an explanation; that below. (“Programming” includes any and all interactions between people and code, in any context (e.g., software engineering, learning, play, productivity, science, and all of the activities involved in creating programs, including requirements, architecture, implementation, verification, monitoring, and more).
5. *Prior Work*: What prior work offers an explanation of this phenomena, or might help generate an explanation of this phenomena? For the purposes of the seminar, this does not need to be complete, but a complete description of this theory would have an extensive literature review covering theories that inspired this theory, as well as conflicting theories.
6. *Concepts*: Describe the key concepts of the theory and some concrete examples of them, building upon the phenomena above. These might be variables, processes, people, aspects of people, structures, contexts or other phenomena that are essential to the theory's account of the phenomena. Note: concepts should be descriptions of ideas that give some structure and precision to describing the phenomena, not operationalizations or measurements – those belong in example hypotheses and/or studies.
7. *Relationships and Mechanisms*: Using the constructs described above, explain the causality of how the phenomena works. What causes what and how? Provide a few concrete examples to illustrate the idea.
8. *Example Hypotheses*: What testable claims do the constructs, relationships, and mechanisms imply?
9. *Example Studies*: What are existing or envisioned example study methods that might investigate the hypotheses above? How might the concepts be operationalized and measured? Describe details about populations, samples, tasks, contexts, tools, observations. Remember that studies can involve many forms of observation and data, both qualitative and quantitative and even design contributions. Studies do not have to be feasible to be proposed and can vary in scope, from single-study sized methods to long-term research agendas that might explore a theory over many years and many projects.
10. *Corollaries*: What follows from this theory, if true? Provide potential implications, concrete or otherwise.

## Unexplained phenomenon

After spending the first day discussing what theories were and applying a theory template, the goal of the second day was to identify unexplained phenomena related to programming and apply theories to explain them. After an informal voting process, attendees created groups to develop theories around the following phenomena:

- Debugging
- Types
- Neurodiversity in programming
- Data programming
- Code examples

- Developer tools
- Learning effects from code analysis

Groups spent all of Wednesday developing theories by filling out the theory template and then getting feedback from members of other groups. They then iterated and created presentations. See included abstracts of talks for descriptions of each presentation.

## Sharing Theories

On Thursday after the presentations, attendees had discussions about how to share theories of programming to broader audiences. Many ideas included written dissemination, such as publishing research, writing books, creating a wiki, adding to reviewer guidelines, creating a website, defining syllabi for reading groups, speaking on podcasts, and posting on social media sites. Other ideas featured opportunities for further interaction, such as workshops, special interest groups, demonstrations of theories for practitioners, and stickers/flair for engagement at conferences. Other ideas focused on incentive structures, such as creating a “best new theory” award at conferences.

Group-wide discussions about sharing theories identified some structural barriers and opportunities. A barrier to broader theory creation and/or use is that most computing researchers do not have much training in theories. Workshops, reading groups, or changes to undergraduate or graduate level coursework could help address this. Another structural barrier is that most conferences lack instructions about theory. Adding instructions in paper calls and reviewer instructions as well as “theory shepherds” could help address this systemic barrier.

## Do we really need theories?

The final session for Thursday was critically reflective about whether programming actually required theories. Given this session occurred after lunch on the final full day, this session got silly. After splitting into groups to discuss, groups shared eclectic presentations to reflect their discussions:

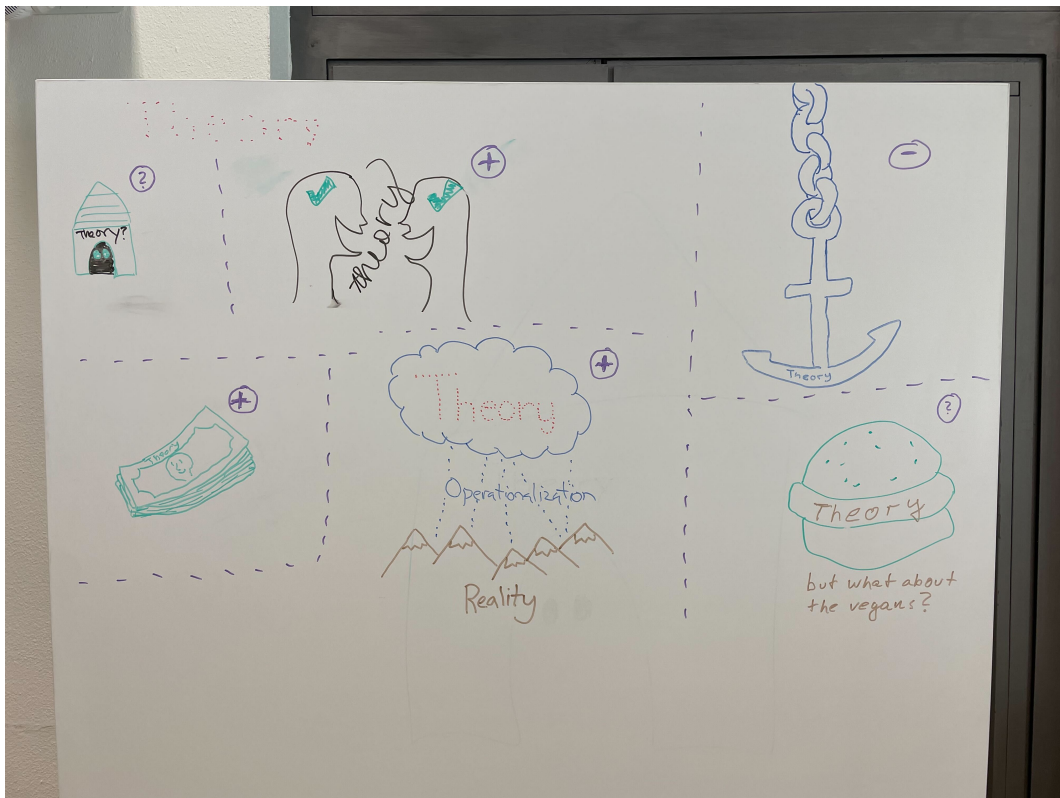
- a colorful whiteboard diagram about pros and cons of theories (Fig. 1)
- Another whiteboard diagram about whether to use theory (Fig. 2)
- A list of bullets about challenges of making changes in publishing
- An humorous improv skit about a conference Q&A session on a theory paper.

## Reflections on the week

The final session of this seminar asked attendees to reflect on the seminar as a whole. Attendees identified some high-level takeaways:

Attendees found theories useful for helping understand why things (e.g. languages or tools) do or do not work. They also found theories helpful for differentiating between how we think people work and how they actually work.

Attendees also felt that the engagement of computing researchers with theories of programming was often limited by the lack of interest and/or lack of expertise. Interdisciplinary



■ **Figure 1** Whiteboard sketches of the pros and cons of theories, as depicted by various diagrams.

research can help create the gestalt of expertise required to create theories of programming, but narrow conference and journal scopes often make this difficult. Specifically, many computing researchers lack expertise in empirical evaluations, making it difficult to develop rigorous evidence that is often foundational to theory building. Furthermore, much training in empirical evaluations focuses on lab settings, whereas most programming happens “in the wild.”

Multiple attendees also felt that theories were more implicitly prevalent in computing research than was explicitly discussed. Some conversation focused on “lower case ‘t’” theories, or theories that we not fully formalized, but provided use and explanation. Many attendees felt that theories implicitly existing in papers, but were unaware of explanations into this work.

A concluding consensus was that theories of programming have existed in the background. Through explicit engagement and discourse, this Dagstuhl Seminar could serve as a catalyst to augment existing theories and craft new ones.

## References

- 1 Sillito, Jonathan, Gail C. Murphy, and Kris De Volder. “Asking and Answering Questions during a Programming Change Task.” *IEEE Transactions on Software Engineering* 34, no. 4 (July 2008): 434–51. <https://doi.org/10.1109/TSE.2008.26>.
- 2 LaToza, Thomas D., David Garlan, James D. Herbsleb, and Brad A. Myers. “Program Comprehension as Fact Finding.” In *Proceedings of the the 6th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations*



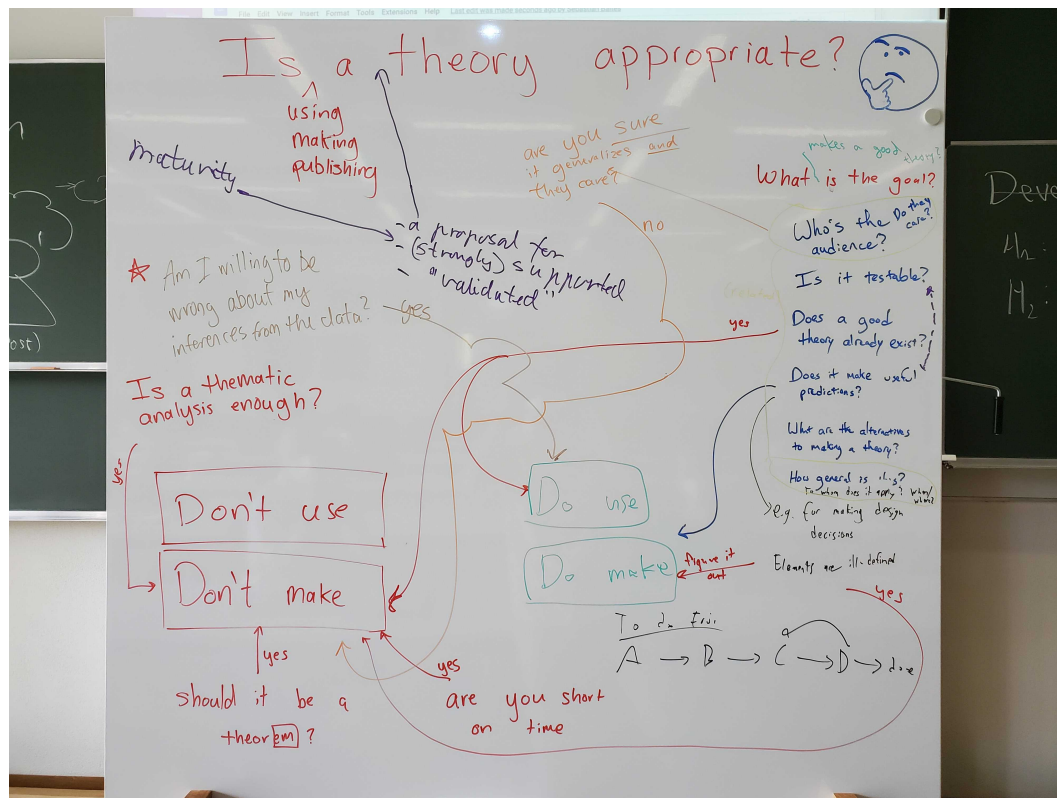


Figure 2 Whiteboard sketch of flowchart considering whether a theory is appropriate.

of Software Engineering, 361–70. ESEC-FSE '07. New York, NY, USA: Association for Computing Machinery, 2007. <https://doi.org/10.1145/1287624.1287675>.

- Spolsky, Joel. "The Law of Leaky Abstractions." Joel on Software, November 11, 2002. <https://www.joelonsoftware.com/2002/11/11/the-law-of-leaky-abstractions/>.
- Parnas, David L. "On the criteria to be used in decomposing systems into modules." *Pioneers and their contributions to software engineering*. Springer, Berlin, Heidelberg, 1972. 479-498.
- Xie, Benjamin, Dastyni Loksa, Greg L. Nelson, Matthew J. Davidson, Dongsheng Dong, Harrison Kwik, Alex Hui Tan, Leanne Hwa, Min Li, and Amy J. Ko. "A Theory of Instruction for Introductory Programming Skills." *Computer Science Education*, January 25, 2019, 1–49. <https://doi.org/10.1080/08993408.2019.1565235>.

## 2 Table of Contents

### Executive Summary

*Benjamin Xie, Amy Ko, Thomas D. LaToza, David C. Shepherd, and Dag Sjøberg* 1

### Overview of Talks

Formulating and Checking Hypotheses in Debugging

*Moritz Beller, Sebastian Balthes, Jonathan Bell, Brittany Johnson-Matthews, and Hila Peleg* . . . . . 9

Tool-Tainted Knowledge Guides Developer Decisions (TTKGDD)

*Thomas Fritz, Tudor Girba, Gail C. Murphy, Dag Sjøberg, and Kathryn T. Stolee* . 9

Theory of Code Examples

*Jun Kato, Gunnar Bergersen, Scott Fleming, Robert Hirschfeld, and Andreas Zeller* 10

Types as tools for structuring thought

*Sarah Lim, Michael Coblenz, Andrew Head, and Thomas D. LaToza* . . . . . 11

Learning Effects from Code Analysis

*Justin Lubin, Francisco Servant, Justin Smith, and Emma Söderberg* . . . . . 11

Neurofriction in Programming Tools

*Jeffrey Stylos, Amy Ko, and Lutz Prechelt* . . . . . 12

Narrative Data Programming: Narrative First, Program Second

*Benjamin Xie and David C. Shepherd* . . . . . 12

**Participants** . . . . . 13

## 3 Overview of Talks

### 3.1 Formulating and Checking Hypotheses in Debugging

*Moritz Beller (Facebook – Menlo Park, US), Sebastian Baltes (University of Adelaide, AU), Jonathan Bell (Northeastern University – Boston, US), Brittany Johnson-Matthews (George Mason University – Fairfax, US), and Hila Peleg (Technion – Haifa, IL)*

**License** © Creative Commons BY 4.0 International license  
© Moritz Beller, Sebastian Baltes, Jonathan Bell, Brittany Johnson-Matthews, and Hila Peleg

Particularly in large software systems, complex bugs may entirely stump developers who attempt to debug them solely through breakpoints and printf statements. One strategy for debugging failing test cases is “scientific debugging” – the process of formulating hypotheses that could explain the failure, and then investigating the code and its execution to see which hypotheses hold and refine them. Some programmers generate more accurate hypotheses than others; some programmers are more efficient at prioritizing and checking hypotheses than others; some may even intuitively jump from a failing test to the valid hypothesis. We have formulated a theory of scientific debugging to enable the externalization of the process that experts follow when debugging, providing a framework for future research in the growth and transfer of expertise. In our theory, developers use experiences from previous debugging to sort through patterns of symptoms and related hypotheses, which helps them navigate the hypothesis space. They iteratively select new working hypotheses that when checked either hold or are refuted, providing more information to diagnose the underlying root cause of the failure. Developers might use different strategies to navigate throughout the hypothesis space, likely without tool support.

### 3.2 Tool-Tainted Knowledge Guides Developer Decisions (TTKGDD)

*Thomas Fritz (Universität Zürich, CH), Tudor Girba (feenk – Wabern, CH), Gail C. Murphy (University of British Columbia – Vancouver, CA), Dag Sjøberg (University of Oslo, NO), and Kathryn T. Stolee (North Carolina State University – Raleigh, US)*

**License** © Creative Commons BY 4.0 International license  
© Thomas Fritz, Tudor Girba, Gail C. Murphy, Dag Sjøberg, and Kathryn T. Stolee

During the “Theories of Programming” seminar, we developed an initial version of a theory entitled “Tool-Tainted Knowledge Guides Developer Decisions (TTKGDD)”. This theory addresses observations that have been made about programming today, particularly that all too often, programmers make decisions about their program based on beliefs rather than evidence. Our theory is that basing decisions on evidence requires contextual tools that extract and present facts about the system in terms that the programmer can easily comprehend. Taking an evidence-based approach leads to higher quality decisions being made about the system.

The major concepts in the theory are developers, tools, decisions, code and hypotheses. There are many relationships between these concepts, such as how developers form hypotheses about their code and how developer’s knowledge guides the choice of a tool to answer a hypothesis. This theory suggests many hypotheses to test, including that “contextualized tools lead to better developer decisions”, “too much automation in tools reduces knowledge of the system” and “biased tools lead to biased knowledge and therefore biased decisions”.

Interestingly, with this last hypothesis, the choice of the term “biased” may indeed bias experiments conducted. Equally interesting hypotheses could be “opinionated tools lead to opinionated knowledge and therefore opinionated decisions” or “appropriate tools lead to appropriate knowledge and therefore appropriate decisions”.

We intend to propose a workshop at a conference that investigates the meaning of a “contextual tool”, the implications of applying them in concrete scenarios, and refinements of the theory.

### 3.3 Theory of Code Examples

*Jun Kato (AIST – Tsukuba, JP), Gunnar Bergersen (University of Oslo, NO), Scott Fleming (University of Memphis, US), Robert Hirschfeld (Hasso-Plattner-Institut, Universität Potsdam, DE), and Andreas Zeller (CISPA – Saarbrücken, DE)*

**License** © Creative Commons BY 4.0 International license  
© Jun Kato, Gunnar Bergersen, Scott Fleming, Robert Hirschfeld, and Andreas Zeller

We were concerned with concrete examples of both code and data in programming activities. Often such examples come in small sizes, as collections rather than single snippets, and are generally considered beneficial to help people understand and extend a codebase or system. Despite this common understanding, examples take many different forms, and there remain numerous challenges to ensuring their benefits.

The authors come from multiple disciplines including Human-Computer Interaction, Software Engineering, Computer Science Education, and Programming Languages, and use the same terminology “code examples” mainly in the following contexts.

First, there are code examples in tutorials and learning materials for computer science education. Several key issues that arise in the use of examples for teaching and learning. Authoring educational examples holds significant challenges in terms of producing correct (e.g., well-tested) code examples and of keeping examples up to date in the face of rapidly evolving platforms and APIs. Designing effective worked examples (i.e., which entail a problem statement, solution steps, and the final solution to the problem) similarly holds challenges with respect to helping learners gain transferable problem-solving knowledge and an understanding of the rationale that underlies the solution steps.

Second, there are code examples in exploratory programming. Exploratory programmers have an open-ended goal, learning about a new domain, working toward a specification and growing a system. To make sure their code works correctly and to find an appropriate next step for their exploration, they provide multiple examples and examine the results of their execution. The major challenge is the difficulty of writing good code examples that cover the test cases of current interest, run in a reasonable time frame, and provide informative feedback for the next steps. We foresee that addressing these issues will require further efforts on improving the liveness of the programming environment and adding guidance based on code understanding techniques, including static and dynamic code analysis.

While we saw differences in these contexts such as the “goodness” criteria for the code examples, we also found similarities like the need to support the authoring process of good code examples. Possible areas of study include how to keep examples relevant to the purpose, how to organize examples in the order that makes the most sense to the learners and programmers, and how to make examples more informative and explorable.

### 3.4 Types as tools for structuring thought

*Sarah Lim (University of California – Berkeley, US), Michael Coblenz (University of Maryland – College Park, US), Andrew Head (University of Pennsylvania – Philadelphia, US), and Thomas D. LaToza (George Mason University – Fairfax, US)*

License © Creative Commons BY 4.0 International license  
© Sarah Lim, Michael Coblenz, Andrew Head, and Thomas D. LaToza

Existing theories about types focus mainly on formal semantics, without considering how type systems actually influence the human practice of programming. Our theory aims to characterize how static type systems can shape the user experience of programming beyond simply surfacing type errors. In our theory, one key way that types can support programmers is by helping a programmer encode an ontology of the domain while planning their program, which will later support reasoning in terms of domain-specific constructs. Ultimately, this leads a type system to support a programmer during ideation, solution search, refactoring, and program comprehension. This encoding can be promoted or inhibited according to the features of the type system in which the work is done, and the expressivity of the type system affects programmers' success in encoding relevant constraints. Then, when programmers use the resulting encoding, their search for an appropriate solution can be guided or inhibited by the constraints in the encoding.

We theorize that well-designed types within a sufficiently expressive type system can (1) catch common mistakes and offload verification work to the computer; (2) help programmers identify good solutions to their problems; and (3) allow types to be expressed in a way that matches the problem domain in the way the programmer thinks about it. Importantly, rich type systems provide counterparts to the execution-based strategies common in dynamically-typed languages, such as defining example data, watching unit tests, or working heavily with a REPL during implementation. We propose future research to explore which design decisions around types support programmers in the tasks described above.

### 3.5 Learning Effects from Code Analysis

*Justin Lubin (University of California – Berkeley, US), Francisco Servant (King Juan Carlos University – Madrid, ES), Justin Smith (Lafayette College – Easton, US), and Emma Söderberg (Lund University, SE)*


License © Creative Commons BY 4.0 International license  
© Justin Lubin, Francisco Servant, Justin Smith, and Emma Söderberg

Developers typically use code analysis tools to improve code quality. We posit that these tools have an equally transformative impact on developers' mental models of a problem domain. For instance, reachability analysis tools may reveal incorrect models of possible states in a system, lifetime analysis in Rust may prompt developers to decide how long certain objects should live in their models, and dependency analysis may reveal circular reasoning in developers' mental models. We theorize that the extent to which the mental model of the problem domain and the embodiment in the code base and by extension the code analysis tools overlap affects the extent to which domain-specific learning effects can occur.

We spent several sessions mapping out an initial framework about the concepts, processes, relationships involved in code analysis tools affecting the mental model of developers as they write code. These discussions resulted in the formation of a team of researchers who will explore this topic further and in a plan to formalize and disseminate our findings in future publications.

### 3.6 Neurofriction in Programming Tools

*Jeffrey Stylos (Stylos Research – Northampton, US), Amy Ko (University of Washington – Seattle, US), and Lutz Prechelt (FU Berlin, DE)*

License  Creative Commons BY 4.0 International license  
© Jeffrey Stylos, Amy Ko, and Lutz Prechelt

Some people approach programming systematically, making plans and then implementing. Other people approach programming more opportunistically, working through examples and then building programs as they test and gather feedback. Others still may have different approaches to programming, shaped by how they learn, process information, and manage their attention.

Our ecosystem of tools, languages, and APIs are mostly created by people who are more systematic, disadvantaging those who do prefer to be more opportunistic, or have other problem solving approaches and preferences. This creates a mismatch between tools and programmers' needs that produces what we call "Neurofriction". Some of this friction is even framed as desirable by tool designers, describing some ways of programming as the "right" or "desirable" way, stigmatizing other ways of working. This is complicated by collaboration, where one team may need to agree on a particular set of languages, tools, and processes that further create Neurofriction.

By understanding Neurofriction better, and developing understanding amongst tool, language, and API designers about diverse needs, we may be able to create more universal tool designs.

### 3.7 Narrative Data Programming: Narrative First, Program Second

*Benjamin Xie (University of Washington – Seattle, US) and David C. Shepherd (Virginia Commonwealth University – Richmond, US)*

License  Creative Commons BY 4.0 International license  
© Benjamin Xie and David C. Shepherd

Computational notebooks are often criticized for their lack of adherence to traditional software engineering best practices. While this is certainly true, and often problematic, there may be good reasons for this departure from accepted norms. Because their purpose is to tell a story, we believe that computational notebooks should be seen as narratives first, and programs second. That is, while essential qualities like correctness are still important, the narrative that is woven from top to bottom of the notebook should take precedence over other non-essential qualities, such as efficiency and code reuse. Viewing notebooks in this way will allow us, as a community, to properly support users with essential best practices without over-burdening these often novice and end-user programmers with unnecessary complexity from practices, tools, and environments.

## Participants

- Sebastian Baltes  
University of Adelaide, AU
- Jonathan Bell  
Northeastern University –  
Boston, US
- Moritz Beller  
Facebook – Menlo Park, US
- Gunnar Bergersen  
University of Oslo, NO
- Michael Coblenz  
University of Maryland –  
College Park, US
- Scott Fleming  
University of Memphis, US
- Thomas Fritz  
Universität Zürich, CH
- Tudor Girba  
feenk – Wabern, CH
- Andrew Head  
University of Pennsylvania –  
Philadelphia, US
- Robert Hirschfeld  
Hasso-Plattner-Institut,  
Universität Potsdam, DE
- Brittany Johnson-Matthews  
George Mason University –  
Fairfax, US
- Jun Kato  
AIST – Tsukuba, JP
- Amy Ko  
University of Washington –  
Seattle, US
- Thomas D. LaToza  
George Mason University –  
Fairfax, US
- Sarah Lim  
University of California –  
Berkeley, US
- Justin Lubin  
University of California –  
Berkeley, US
- Gail C. Murphy  
University of British Columbia –  
Vancouver, CA
- Hila Peleg  
Technion – Haifa, IL
- Lutz Prechelt  
FU Berlin, DE
- Francisco Servant  
King Juan Carlos University –  
Madrid, ES
- David C. Shepherd  
Virginia Commonwealth  
University – Richmond, US
- Dag Sjøberg  
University of Oslo, NO
- Justin Smith  
Lafayette College – Easton, US
- Emma Söderberg  
Lund University, SE
- Kathryn T. Stolee  
North Carolina State University –  
Raleigh, US
- Jeffrey Stylos  
Stylos Research –  
Northampton, US
- Benjamin Xie  
University of Washington –  
Seattle, US
- Andreas Zeller  
CISPA – Saarbrücken, DE



# Efficient and Equitable Natural Language Processing in the Age of Deep Learning

Jesse Dodge<sup>\*1</sup>, Iryna Gurevych<sup>\*2</sup>, Roy Schwartz<sup>\*3</sup>, Emma Strubell<sup>\*4</sup>,  
and Betty van Aken<sup>†5</sup>

- 1 AI2 – Seattle, US. [jessed@allenai.org](mailto:jessed@allenai.org)
- 2 TU Darmstadt, DE. [gurevych@cs.tu-darmstadt.de](mailto:gurevych@cs.tu-darmstadt.de)
- 3 The Hebrew University of Jerusalem, IL. [roy.schwartz1@mail.huji.ac.il](mailto:roy.schwartz1@mail.huji.ac.il)
- 4 Carnegie Mellon University – Pittsburgh, US. [strubell@cmu.edu](mailto:strubell@cmu.edu)
- 5 (Berliner Hochschule für Technik, DE. [Betty.vanAken@bht-berlin.de](mailto:Betty.vanAken@bht-berlin.de))

---

## Abstract

This report documents the program and the outcomes of Dagstuhl Seminar 22232 “Efficient and Equitable Natural Language Processing in the Age of Deep Learning”. Since 2012, the field of artificial intelligence (AI) has reported remarkable progress on a broad range of capabilities including object recognition, game playing, speech recognition, and machine translation. Much of this progress has been achieved by increasingly large and computationally intensive deep learning models: training costs for state-of-the-art deep learning models have increased 300,000 times between 2012 and 2018 [1]. Perhaps the epitome of this trend is the subfield of natural language processing (NLP) that over the past three years has experienced even sharper growth in model size and corresponding computational requirements in the word embedding approaches (e.g. ELMo, BERT, openGPT-2, Megatron-LM, T5, and GPT-3, one of the largest models ever trained with 175B dense parameters) that are now the basic building blocks of nearly all NLP models. Recent studies indicate that this trend is both environmentally unfriendly and prohibitively expensive, raising barriers to participation in NLP research [2, 3]. The goal of this seminar was to mitigate these concerns and promote equity of access in NLP.

## References

- 1 D. Amodei and D. Hernandez. 2018. AI and Compute. <https://openai.com/blog/ai-and-compute>
- 2 R. Schwartz, D. Dodge, N. A. Smith, and O. Etzioni. 2020. Green AI. Communications of the ACM (CACM)
- 3 E. Strubell, A. Ganesh, and A. McCallum. 2019. Energy and Policy Considerations for Deep Learning in NLP. In Proc. of ACL.

**Seminar** June 6–10, 2022 – <http://www.dagstuhl.de/22232>

**2012 ACM Subject Classification** Computing methodologies → Natural language processing; Computing methodologies → Neural networks; Social and professional topics → Sustainability

**Keywords and phrases** deep learning, efficiency, equity, natural language processing (nlp)

**Digital Object Identifier** 10.4230/DagRep.12.6.14

---

\* Editor / Organizer

† Editorial Assistant / Collector



Except where otherwise noted, content of this report is licensed under a Creative Commons BY 4.0 International license

Efficient and Equitable Natural Language Processing in the Age of Deep Learning, *Dagstuhl Reports*, Vol. 12, Issue 6, pp. 14–27

Editors: Jesse Dodge, Iryna Gurevych, Roy Schwartz, and Emma Strubell



Dagstuhl Reports

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany




## 1 Executive Summary

*Roy Schwartz (The Hebrew University of Jerusalem, IL)*

*Jesse Dodge (AI2 – Seattle, US)*

*Iryna Gurevych (TU Darmstadt, DE)*

*Emma Strubell (Carnegie Mellon University – Pittsburgh, US)*

License  Creative Commons BY 4.0 International license  
© Roy Schwartz, Jesse Dodge, Iryna Gurevych, and Emma Strubell

For this seminar, we brought together a diverse group of researchers and practitioners in NLP and adjacent fields to develop actionable policies, incentives and a joint strategy towards more efficient and equitable NLP. This Dagstuhl Seminar covered a range of related topics, which we summarize as follows.

### Efficient NLP models

A key method for mitigating the raised concerns is reducing costs by making models more efficient. We surveyed the different methods that exist for making NLP technology more efficient. We discussed their tradeoffs, prioritized them, and aimed to identify new opportunities to promote efficiency in NLP. During the seminar, we drafted a survey paper summarizing multiple methods for increasing the efficiency of NLP models. We aim to publish this work later this year.

### Systemic issues

We also addressed systemic issues in the field relating to the reporting of computational budgets in NLP research, and how we can use incentive structures such as the NLP Reproducibility Checklist [1] to motivate researchers throughout the field to improve reporting. We discussed the survey responses for the reproducibility checklist used at four major NLP conferences, and we plan to release a report of this data.

### Equity of access

A third topic of discussion was the equity of access to computational resources and state-of-the-art NLP technologies. Prior to the seminar, we conducted a survey of different stakeholders across the NLP community. During the seminar, we analyzed and discussed the results of this survey to better understand who is most affected and how, and developed informed strategies and policies to mitigate this inequity moving forward. We are currently working on a paper summarizing the results of this survey, which we hope to publish later this year.

### Measuring efficiency and equity

All of the above endeavors require establishing the right metrics and standards to measure our current status and progress towards efficiency and equity goals. We discussed multiple metrics and evaluation frameworks that capture the bigger picture of how different approaches compare in terms of energy efficiency not just in the research environment but in practice and over the entire ML model lifecycle (development, training and deployment), and that work under a wide range of computational budgets.

### References

- 1 Jesse Dodge, Suchin Gururangan, Dallas Card, Roy Schwartz, Noah A. Smith: Show Your Work: Improved Reporting of Experimental Results. EMNLP/IJCNLP (1) 2019: 2185-2194

## 2 Table of Contents

### Executive Summary

*Roy Schwartz, Jesse Dodge, Iryna Gurevych, and Emma Strubell* . . . . . 15

### Overview of Talks

Forays into Efficiency and Energy of NLP Models

*Niranjan Balasubramanian* . . . . . 18

Faster Neural Network Training, Algorithmically

*Jonathan Frankle* . . . . . 18

Evaluating Approximations is Hard; Efficient Machine Translation Shared Task

*Kenneth Heafield* . . . . . 18

ML Efficiency: Open Challenges and Opportunities.

*Sara Hooker* . . . . . 19

Neurosymbolic models in semantic parsing

*Alexander Koller* . . . . . 19

Investigating Rational Activation Functions to Train Transformer Models

*Ji-Ung Lee* . . . . . 20

Holistic model evaluation

*Alexandra Sasha Luccioni* . . . . . 20

Deep Patient Representation

*Alexander Löser* . . . . . 20

Is Sparsity a Path for Efficiency?

*André F. T. Martins* . . . . . 21

The Sweet Lesson

*Colin Raffel* . . . . . 21

On #Reviewer2 and paper-reviewer assignments

*Anna Rogers* . . . . . 21

BigScience Large LMs and small labs

*Thomas Wolf* . . . . . 22

### Working groups

Efficiency benchmarking

*Niranjan Balasubramanian, Leon Derczynski, Jesse Dodge, Jonathan Frankle, Iryna Gurevych, Kenneth Heafield, Sara Hooker, André F. T. Martins, Haritz Puerto, Colin Raffel, Roy Schwartz, Edwin Simpson, Noam Slonim, and Thomas Wolf* . . . . . 22

Implementing changes in NLP research

*Jesse Dodge, Niranjan Balasubramanian, Jessica Forde, Kenneth Heafield, Alexander Koller, Ji-Ung Lee, André F. T. Martins, Nils Reimers, Leonardo Ribeiro, Andreas Rücklé, and Betty van Aken* . . . . . 23

Breakout on “Making Change”

*Roy Schwartz, Leon Derczynski, Jonathan Frankle, Iryna Gurevych, Alexander Koller, Alexandra Sasha Luccioni, André F. T. Martins, Colin Raffel, Anna Rogers, Noam Slonim, Emma Strubell, and Thomas Wolf* . . . . . 23

**Panel discussions**

Panel on Equity in NLP research


*Colin Raffel, Iryna Gurevych, Alexandra Sasha Luccioni, Noah A. Smith, Emma Strubell, and Thomas Wolf* . . . . . 26

**Participants** . . . . . 27

### 3 Overview of Talks

#### 3.1 Forays into Efficiency and Energy of NLP Models


*Niranjana Balasubramanian (Stony Brook University, US)*

License  Creative Commons BY 4.0 International license  
© Niranjana Balasubramanian

This talk presents forays into efficient QA models, modeling sparsity for hardware acceleration, and issues in measuring energy consumption.

#### 3.2 Faster Neural Network Training, Algorithmically

*Jonathan Frankle (Harvard University – Allston, US)*

License  Creative Commons BY 4.0 International license  
© Jonathan Frankle

Training modern neural networks is time-consuming, expensive, and energy-intensive. As neural network architectures double in size every few months, it is difficult for researchers and businesses without immense budgets to keep up. In this talk, I describe one approach for managing this challenge: changing the training algorithm itself. While many companies and researchers are focused on building hardware and systems to allow existing algorithms to run faster in a mathematically equivalent fashion, there is nothing sacred about this math. To the contrary, training neural networks is inherently approximate, relying on noisy data, convex optimizers in nonconvex regimes, and ad hoc tricks and hacks that seem to work well in practice for reasons that elude us.

I discuss how we have put this approach into practice at MosaicML, including the dozens of algorithmic changes we have studied (which are freely available open source), the science behind how these changes interact with each other (the composition problem), and how we evaluate whether these changes have been effective. I will also detail several surprises we have encountered and lessons we have learned along the way. In the four months since we began this work in earnest, we have reduced the training times of standard computer vision models by 7x and standard language models by 2x on publicly available cloud instances, and we believe we are just scratching the surface.

#### 3.3 Evaluating Approximations is Hard; Efficient Machine Translation Shared Task

*Kenneth Heafield (University of Edinburgh, GB)*

License  Creative Commons BY 4.0 International license  
© Kenneth Heafield

Papers about a new approximation (i.e. faster for some loss in quality) often claim the quality loss is small, while better papers perform a Pareto comparison. Unfortunately, the baseline approximations used for the Pareto comparison are usually restricted to the same type of method, such as pruning. I argue the correct baseline is all approximations that already exist. Approximations are stackable, so the question is really whether the proposed

approximation belongs to a set of stacked approximations that advance the Pareto frontier. This is a high standard and difficult for the average paper to reach, so I present a partial solution. The efficient machine translation shared task establishes the state-of-the-art by soliciting competitive submissions and comparing them. Starting from a range of already efficient systems provides a much stronger baseline for evaluating a new approximation.

### 3.4 ML Efficiency: Open Challenges and Opportunities.

*Sara Hooker (Google – Mountain View, US)*

License  Creative Commons BY 4.0 International license  
© Sara Hooker

Our field is currently characterized by a “bigger is better” trend in the size of deep neural networks. This talk posits that this is an unsustainable recipe – akin to building a ladder to the moon. We discuss some important directions for revisiting the efficiency of our representation learning approaches.

### 3.5 Neurosymbolic models in semantic parsing

*Alexander Koller (Universität des Saarlandes, DE)*


License  Creative Commons BY 4.0 International license  
© Alexander Koller

There are many approaches to mapping natural-language sentences to symbolic meaning representations. The current dominant approach is with neural sequence-to-sequence models, which map the sentence to a string version of the meaning representation. Seq2seq models work well for many NLP tasks, including tagging and parsing, and deliver excellent accuracy on broad-coverage semantic parsing as well. However, it has recently been found that seq2seq models struggle with “compositional generalization”: They have a hard time generalizing from training examples to structurally similar unseen test sentences. I will show some new results that pinpoint this difficult more precisely, and discuss what this means for how to best evaluate semantic parsers.

I will then present our own research on compositional semantic parsing, which combines neural models with the Principle of Compositionality from theoretical semantics. Our semantic parser uses a neural supertagger to predict word meanings and a neural dependency parser to predict the compositional structure, and then evaluates this dependency structure in a graph algebra to obtain the meaning representation. We achieve state-of-the-art parsing accuracy across a number of graphbanks, at a speed of up to 10k tokens/second.

### 3.6 Investigating Rational Activation Functions to Train Transformer Models

*Ji-Ung Lee (TU Darmstadt, DE)*

**License**  Creative Commons BY 4.0 International license

© Ji-Ung Lee

**Joint work of** Haishuo Fang, Ji-Ung Lee, Nafise Sadat Moosavi, Iryna Gurevych

In this work, we explore rational activation functions for training transformer models. In contrast to activation functions such as GELU which remain fixed after initialization, rational activation functions are capable of approximating any arbitrary activation function during training. In preliminary experiments we find that using rational activation functions can lead to a faster convergence during pre-training as well as a higher performance on several downstream tasks.

### 3.7 Holistic model evaluation

*Alexandra Sasha Luccioni (Hugging Face – Paris, FR)*


**License**  Creative Commons BY 4.0 International license

© Alexandra Sasha Luccioni

In both research and industry, there are multiple factors to consider when comparing models. Our current ML benchmarks measure one aspect of this, e.g. NLI, NER, QA. How do we integrate different aspects of model performance when comparing models?

### 3.8 Deep Patient Representation

*Alexander Löser (Berliner Hochschule für Technik, DE)*

**License**  Creative Commons BY 4.0 International license

© Alexander Löser

**Joint work of** Betty van Aken, Jens-Michalis Papaioannou, Manuel Mayrdorfer, Klemens Budde, Felix A. Gers, Alexander Löser

**Main reference** Betty van Aken, Jens-Michalis Papaioannou, Manuel Mayrdorfer, Klemens Budde, Felix A. Gers, Alexander Löser: “Clinical Outcome Prediction from Admission Notes using Self-Supervised Knowledge Integration”, in Proc. of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 – 23, 2021, pp. 881–893, Association for Computational Linguistics, 2021.

**URL** <https://doi.org/10.18653/v1/2021.eacl-main.75>

Understanding clinical outcomes requires to integrate different modalities in a single latent representation. We present such operators that reuse clinical large language models in English language and integrate complementary medical latent representation from low resource languages, from ontologies, from time variant data and from set data. An example application is the differential diagnosis at <https://outcome-prediction.demo.dataxis.com>.

### 3.9 Is Sparsity a Path for Efficiency?

*André F. T. Martins (IST – Lisbon, PT)*

**License** © Creative Commons BY 4.0 International license  
© André F. T. Martins

Current NLP models are increasingly larger and data-hungry, which poses important environmental challenges. In this talk, I discuss several ways in which sparsity might lead to more efficient NLP models. The current life cycle of NLP models offers several opportunities to improve memory and runtime efficiency at different stages: during pretraining, during finetuning, and during inference. I first distinguish between model sparsity and activation sparsity. Then, I focus on adaptive sparse attention approaches for the latter, where the softmax transformation is replaced by sparse transformations – entmax – which maintain end-to-end differentiability and have a learnable parameter which controls their sparsity. I finish by asking several open questions and inviting discussion.

### 3.10 The Sweet Lesson

*Colin Raffel (University of North Carolina at Chapel Hill, US)*

**License** © Creative Commons BY 4.0 International license  
© Colin Raffel

Richard Sutton’s essay “The Bitter Lesson” argues that “general methods that leverage computation are ultimately the most effective”. In this talk, I will argue that the bitter lesson implies that, at a given point in time, it is often possible to outperform large-scale methods with methods that are more efficient and clever. Furthermore, actively working to develop more efficient methods has often uncovered new approaches that scale better. I call this perspective “the sweet lesson” and will present many examples of this principle. Finally, I will wrap up with some thoughts on how to internalize bitter and sweet lessons in NLP’s current era of scale.

### 3.11 On #Reviewer2 and paper-reviewer assignments

*Anna Rogers (University of Copenhagen, DK)*

**License** © Creative Commons BY 4.0 International license  
© Anna Rogers

**Joint work of** Terne Thorn Jakobsen, Anna Rogers

**Main reference** Terne Thorn Jakobsen, Anna Rogers: “What Factors Should Paper-Reviewer Assignments Rely On? Community Perspectives on Issues and Ideals in Conference Peer-Review”, in Proc. of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 4810–4823, Association for Computational Linguistics, 2022.

**URL** <https://doi.org/10.18653/v1/2022.naacl-main.354>

Some thoughts and new data on community preferences for how papers should be matched to reviewers.

### 3.12 BigScience Large LMs and small labs

*Thomas Wolf (Hugging Face – Paris, FR)*


License  Creative Commons BY 4.0 International license  
© Thomas Wolf

In this talk, I'll be presenting the BigScience (<https://bigscience.huggingface.co>) project. A collaborative experiment in building a multilingual large scale dataset as well as a multilingual large language model, inspired by other fields of research like the Large Hadron Collider.

## 4 Working groups

### 4.1 Efficiency benchmarking

*Niranjana Balasubramanian (Stony Brook University, US), Leon Derczynski (IT University of Copenhagen, DK), Jesse Dodge (AI2 – Seattle, US), Jonathan Frankle (Harvard University – Allston, US), Iryna Gurevych (TU Darmstadt, DE), Kenneth Heafield (University of Edinburgh, GB), Sara Hooker (Google – Mountain View, US), André F. T. Martins (IST – Lisbon, PT), Haritz Puerto (TU Darmstadt, DE), Colin Raffel (University of North Carolina at Chapel Hill, US), Roy Schwartz (The Hebrew University of Jerusalem, IL), Edwin Simpson (University of Bristol, GB), Noam Slonim (IBM – Haifa, IL), and Thomas Wolf (Hugging Face – Paris, FR)*

License  Creative Commons BY 4.0 International license  
© Niranjana Balasubramanian, Leon Derczynski, Jesse Dodge, Jonathan Frankle, Iryna Gurevych, Kenneth Heafield, Sara Hooker, André F. T. Martins, Haritz Puerto, Colin Raffel, Roy Schwartz, Edwin Simpson, Noam Slonim, and Thomas Wolf

This breakout session concerned questions about what we should measure and report for NLP experiments and how efficiency can be measured. A common problem with current practice in NLP is that efficiency is either not reported at all or that the used metrics are hard to compare. Different hardware environments often further require individual solutions. Shared tasks and benchmarks with fixed hardware were identified as one attempt to mitigate the problem of comparability. The MLPerf benchmarks [1] were mentioned as one positive example. However, participants raised the question whether such benchmarks lead to overfitting and distract from real life concerns. Also, different tasks require different constraints, e.g. the pre-training of a large language model entails different concerns than inferencing on this model. Use cases should therefore be viewed from different perspectives. The group agreed that pushing people to report and review efficiency measures can result in a culture shift and an acceleration of science in general.

#### References

- 1 Farrell, Steven, Murali Emani, Jacob Balma, Lukas Drescher, Aleksandr Drozd, Andreas Fink, Geoffrey Fox et al. “MLPerf™ HPC: A Holistic Benchmark Suite for Scientific Machine Learning on HPC Systems.” In 2021 IEEE/ACM Workshop on Machine Learning in High Performance Computing Environments (MLHPC), pp. 33-45. IEEE, 2021.



## 4.2 Implementing changes in NLP research

*Jesse Dodge (AI2 – Seattle, US), Niranjan Balasubramanian (Stony Brook University, US), Jessica Forde (Brown University – Providence, US), Kenneth Heafield (University of Edinburgh, GB), Alexander Koller (Universität des Saarlandes, DE), Ji-Ung Lee (TU Darmstadt, DE), André F. T. Martins (IST – Lisbon, PT), Nils Reimers (Hugging Face – Paris), Leonardo Ribeiro (TU Darmstadt, DE), Andreas Rücklé (Amazon – Berlin, DE), and Betty van Aken (Berliner Hochschule für Technik, DE)*

**License** © Creative Commons BY 4.0 International license

© Jesse Dodge, Niranjan Balasubramanian, Jessica Forde, Kenneth Heafield, Alexander Koller, Ji-Ung Lee, André F. T. Martins, Nils Reimers, Leonardo Ribeiro, Andreas Rücklé, and Betty van Aken

In this breakout session, the group discussed how to *implement* the ideas for more efficient and equitable NLP research within the community. A step toward that goal is the Checklist for Responsible NLP Research recently introduced in the ACL Rolling Review Submissions. The list contains questions, e.g. whether a submission discusses the risks of a work or mentions the computational budget of a solution. Jesse Dodge presented statistics from the first rounds of reviews using the checklist. The results showed that the acceptance rate was positively correlated with the number of items checked on the list, indicating that the initiative was successful in stimulating higher-quality research. Participants also agreed that templates for sections that mention limitations or reproducibility of experiments are a useful tool for early-stage researchers. The common view within the group was that current academic structures often encourage quantity over quality of publications, especially for junior researchers. Changing this culture and the incentives for doing research was identified as necessary for more carefully thought out and reproducible research.

## 4.3 Breakout on “Making Change”

*Roy Schwartz (The Hebrew University of Jerusalem, IL), Leon Derczynski (IT University of Copenhagen, DK), Jonathan Frankle (Harvard University – Allston, US), Iryna Gurevych (TU Darmstadt, DE), Alexander Koller (Universität des Saarlandes, DE), Alexandra Sasha Luccioni (Hugging Face – Paris, FR), André F. T. Martins (IST – Lisbon, PT), Colin Raffel (University of North Carolina at Chapel Hill, US), Anna Rogers (University of Copenhagen, DK), Noam Slonim (IBM – Haifa, IL), Emma Strubell (Carnegie Mellon University – Pittsburgh, US), and Thomas Wolf (Hugging Face – Paris, FR)*

**License** © Creative Commons BY 4.0 International license

© Roy Schwartz, Leon Derczynski, Jonathan Frankle, Iryna Gurevych, Alexander Koller, Alexandra Sasha Luccioni, André F. T. Martins, Colin Raffel, Anna Rogers, Noam Slonim, Emma Strubell, and Thomas Wolf

- Roy: Wrote a policy document to be adopted by the ACL exec
  - Iryna came to talk at Roy’s lab, mentioned she was part of ACL exec suggested going through ACL to try to influence things
  - Brought on Emma and Jesse and Andreas and had weekly meetings to think about what to do
  - Iryna knew next steps and how to promote – maybe not a general recipe since it relied on her expertise and position

- Wrote a document and took time to get it right, including feedback at the ACL business meetings, assembled an advisory panel for feedback too (including people you disagreed with)
- What were the three recommendations?
  - \* Add instructions to reviewers and authors and question for review form
  - \* Add efficiency track permanently to all future conferences
  - \* Encourage submission of code and data using the badge system
- ACL exec ultimately decides whether they will adopt it, and then it is technically a new set of recommendations for all of the rest of the conferences
- There is no set of centralized rules – PCs of individual conferences can ultimately choose to ignore the policy documents. May therefore need to also talk to individual PCs to get the changes implemented.
- There is also a conference handbook – also hoping to have the recommendations put in the conference handbook. The person responsible for it is part of the ACL exec – need to work with the person to say “this is the paragraph that needs to be included here, this is the paragraph that needs to be included there”.
- How do PCs get elected? They are invited by the exec.
- How does the ACL exec get elected? Candidates are nominated/selected and then the community votes.
- Anna: Since PCs rotate, there may be no continuity. Therefore things that stick are ones that are perceived as being a positive change.
- Jonathan: Given this, what are the most “durable” changes? The efficiency track?
- Alexander: Things are more durable thanks to ARR, which is controlled by the exec rather than the individual PCs
- Noam: Meta-point – make OKRs?
  - What are the objectives and key results required for that?
  - Need to figure out how we are going to measure whether we were successful.
  - André: Agree, seems necessary to separate the what from the how.
  - Leon: Super hard to truly do in a super exact way – “does the change in reviewing have an impact” – ultimately you have subjective judgements, you can’t really measure a lot of these changes.
  - Alexander: Certain things are easier to measure – e.g. are people releasing more code?
  - Emma: Survey every year to ask things like “are things getting better?”
  - Are other communities dealing with this? E.g. quantum computing.
- Jonathan: Split off?
  - Why are you fighting to change the communities you have rather than split off and create a new subcommunity?
  - Good example: FACCT.
  - Emma: Strongly disagree – sort of like checking out, would rather change the community rather than make a new community that cares about different issues. Worry that FACCT makes those issues not first-priority issues in the ML community.
  - Jonathan: FACCT is changing the machine learning community, through influence. For example, with mlsys. Arguably it’s even more impactful than having them be lost in the shuffle at NeurIPS.
  - Jonathan: Not creating a new community – the community exists, really about how to get the message out there.
  - Thom: Keep them connected so that there’s cultural exchange.

- Alexander: Meta-point – levels of change
  - There are multiple levels at which change can be caused – e.g. individual (blog posts, tweets, whatever) and top-down (e.g. policy doc) and community building (there is a community that agrees with us)
  - Thom: Build tools – easy to use, best thing to use – that causes change.
  - Anna: Importance of the human interface – what about putting “efficiency” badges on the HF hub?
- Leon: How much of the community are you reaching/representing?
- Thom: Peer review?
  - Less worried about efficiency, seems like a lot of people are interested in it.
  - More worried about carbon emissions especially in contrast to the desire for more GPUs.
  - Worried about the reviewing process because it can mean big groups leave the reviewing process.
  - Alexander: Super worrying development that people are circumventing the review process through arxiv and press releases.
  - What about connecting other communities? E.g. EleutherAI, can we do a way to connect the communities.
  - Anna: I want anonymous pre-prints.
  - Jonathan: Can we publicly peer review the non-peer reviewed papers?
  - André: Shouldn't assume peer review is the correct thing.
  - André: Don't make peer review too complicated – e.g. don't use super-structured review forms.
  - Emma: Structured review forms are a super effective tool for effecting change. TACL “Single box and you write what you want” vs. ARR prompting about the science specifically. Prompting is important, not complexity.
  - Colin: Any examples of successful public critiques?
  - Leon: Yes, about a stock market paper.
  - Roy: Yes, Yoav about a language generation paper.
  - André: Blog posts don't solve the problem because in the end it still depends on influential people.
  - Anna: Sort of like fake news and fact checking. E.g. post a negative results/can't reproduce paper.
  - Sasha: Use openreview.
- Roy: How do we create a community towards these initiatives/keep the community alive?
  - Jonathan: Start a non-archival workshop that accepts minimum standards to have a big tent.
  - Jonathan: Trying to find ways to sustain conversations can be hard, e.g. dead Slacks – needs to be an event to have people meet regularly.
  - Iryna: Reflecting on argument mining – also a Dagstuhl Seminar done several times. Longest-lasting effect were the people – students in the seminar and the students of the PIs in the seminar.
  - Iryna: Tutorials, summer schools.
  - Anna: Does this work for norms vs. research areas?
  - Iryna: For norms you also have the institutional factor and scientific debate.
  - André: Have there been any workshops or tutorials?
  - Roy: SustaiNLP.
  - André: Make it about research and not about policies.

- Sasha: These are longer-term things – what about short-term things? E.g. the fact that large LM papers keep winning best paper awards.
- Alexander: Are people in other communities interested in the same things?
- Alexander: Really like workshops and tutorials.
- Alexander: Transformative papers should be recognized as transformative.
- Alexander: Octopus paper was a nice example where a best paper award brought a lot of attention to the paper, beyond a large Twitter following. But it didn't change people's minds – just connected people who already thought the same thing.
- Iryna: Money!
  - Set up large-scale funding programs to support the work. In Germany, funding scheme where you can propose a special topic and they fund the faculty.

## 5 Panel discussions

### 5.1 Panel on Equity in NLP research

*Colin Raffel (University of North Carolina at Chapel Hill, US), Iryna Gurevych (TU Darmstadt, DE), Alexandra Sasha Luccioni (Hugging Face – Paris, FR), Noah A. Smith (University of Washington – Seattle, US), Emma Strubell (Carnegie Mellon University – Pittsburgh, US), and Thomas Wolf (Hugging Face – Paris, FR)*

**License** © Creative Commons BY 4.0 International license  
 © Colin Raffel, Iryna Gurevych, Alexandra Sasha Luccioni, Noah A. Smith, Emma Strubell, and Thomas Wolf

The panel, together with contributions from the audience, discussed a number of aspects relating to equity in NLP research. Two of the most prominent discussion items were as follows. (1) An unequal allocation of resources can lead to a misalignment between real-world problems and research work. Sharing of resources (HPC cluster usage, collaborations), and hiring of researchers with experience and passion for real-world problems may provide some mitigation. (2) The lack of diversity in the research community – e.g., geographically and institutionally – can lead to an over-exposure and hype for certain types of work and research agendas. This leaves little attention for progress being made in domains that deviate from the mainstream. Mitigation strategies can include changing the incentive structures for publication, actively endorsing research work on a personal level, promoting the inclusion of researchers in discussions with different backgrounds, and simplifying communication with people affected by real-world problems relating to NLP.

## Participants

- Yuki Arase  
Osaka University, JP
- Niranjan Balasubramanian  
Stony Brook University, US
- Leon Derczynski  
IT University of  
Copenhagen, DK
- Jesse Dodge  
AI2 – Seattle, US
- Jessica Forde  
Brown University –  
Providence, US
- Jonathan Frankle  
Harvard University –  
Allston, US
- Iryna Gurevych  
TU Darmstadt, DE
- Michael Hassid  
The Hebrew University of  
Jerusalem, IL
- Kenneth Heafield  
University of Edinburgh, GB
- Sara Hooker  
Google – Mountain View, US
- Alexander Koller  
Universität des Saarlandes, DE
- Ji-Ung Lee  
TU Darmstadt, DE
- Alexander Löser  
Berliner Hochschule für  
Technik, DE
- Alexandra Sasha Luccioni  
Hugging Face – Paris, FR
- André F. T. Martins  
IST – Lisbon, PT
- Haritz Puerto  
TU Darmstadt, DE
- Colin Raffel  
University of North Carolina  
at Chapel Hill, US
- Nils Reimers  
Hugging Face – Paris
- Leonardo Ribeiro  
TU Darmstadt, DE
- Anna Rogers  
University of Copenhagen, DK
- Andreas Rücklé  
Amazon – Berlin, DE
- Roy Schwartz  
The Hebrew University of  
Jerusalem, IL
- Edwin Simpson  
University of Bristol, GB
- Noam Slonim  
IBM – Haifa, IL
- Noah A. Smith  
University of Washington –  
Seattle, US
- Emma Strubell  
Carnegie Mellon University –  
Pittsburgh, US
- Betty van Aken  
Berliner Hochschule für  
Technik, DE
- Thomas Wolf  
Hugging Face – Paris, FR



# Human-Game AI Interaction

Dan Ashlock<sup>\*†1</sup>, Setareh Maghsudi<sup>†2</sup>, Diego Perez Liebana<sup>†3</sup>,  
Pieter Spronck<sup>†4</sup>, and Manuel Eberhardinger<sup>‡5</sup>

- 1 University of Guelph, CA
- 2 Universität Tübingen, DE. [setareh.maghsudi@uni-tuebingen.de](mailto:setareh.maghsudi@uni-tuebingen.de)
- 3 Queen Mary University of London, GB. [diego.perez@qmul.ac.uk](mailto:diego.perez@qmul.ac.uk)
- 4 Tilburg University, NL. [p.spronck@gmail.com](mailto:p.spronck@gmail.com)
- 5 Hochschule der Medien – Stuttgart, DE. [eberhardinger@hdm-stuttgart.de](mailto:eberhardinger@hdm-stuttgart.de)

---

## Abstract

People interact with semi-intelligent machines during their daily lives. They desire systems to respond intelligently to requests. While improvements to the interaction between humans and AI have been made over the years, these systems are a long way from responding like a human partner. Virtual (game) worlds are an ideal environment in which to experiment with the interaction between humans and AI, due to their similarity with real world environments and the presence of agents that represent “real people” that make decisions and interact among them.

In recent years, the number of ways in which players can interact with games have increased considerably: from the traditional mouse, keyboard, and controller, to responding to natural movements, facial expressions, voice, eye movements and brain signals, among others. This seminar brought together scientists, researchers, and industrial developers who specialize in intelligent interaction between humans and computer agents in virtual (game) environments. This report documents the program and its outcomes.

**Seminar** June 19–24, 2022 – <http://www.dagstuhl.de/22251>

**2012 ACM Subject Classification** Computing methodologies → Artificial intelligence; Human-centered computing → Human computer interaction (HCI); Applied computing → Computer games

**Keywords and phrases** Computational intelligence, artificial intelligence, games, modeling, interaction

**Digital Object Identifier** 10.4230/DagRep.12.6.28

## 1 Executive Summary

*Pieter Spronck (Tilburg University, NL)*

*Daniel Ashlock*

*Setareh Maghsudi (Universität Tübingen, DE)*

*Diego Perez Liebana (Queen Mary University of London, GB)*

**License** © Creative Commons BY 4.0 International license  
© Pieter Spronck, Daniel Ashlock, Setareh Maghsudi, and Diego Perez Liebana

Over the past decades, artificial intelligence has evolved from esoteric techniques used mainly in computer science research to an integral and ever-growing part of the daily lives of most humans. People regularly interact with semi-intelligent machines during their daily lives, whether it is via smartphone applications, embedded systems in cars and household electronics,

---

\* in memoriam; † April 5, 2022

† Editor / Organizer

‡ Editorial Assistant / Collector



Except where otherwise noted, content of this report is licensed under a Creative Commons BY 4.0 International license

Human-Game AI Interaction, *Dagstuhl Reports*, Vol. 12, Issue 6, pp. 28–82

Editors: Dan Ashlock, Setareh Maghsudi, Diego Perez Liebana, Pieter Spronck, and Manuel Eberhardinger



DAGSTUHL  
REPORTS Dagstuhl Reports

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

client support systems, or helpful technology installed on personal computers. People wish and expect systems to respond intelligently to their requests, and even to anticipate their actions. While improvements to the interaction between humans and intelligent systems in this respect have been made over the years, there is still a long way to go before these systems exhibit a level of understanding and intuition which can be expected from a human partner.

Human-computer interfaces (HCI) are a well-established scientific research domain. We noted that HCI research generally neglects the use of artificial intelligence as a integral part of an interface. Almost any person that uses computers can quickly recall multiple frustrating interactions with the current state of the art in artificial intelligence in interfaces. Since annoyance and apparent incompetence can derail the adoption of otherwise promising and potentially transformative technology, research into improving interfaces using AI is timely.

An AI assistant is broadly recognized as being a key factor in increasing human productivity, but it must be an AI assistant that the user either enjoys working with or that the user barely notices, not one that must be bludgeoned into useful behavior or constantly fought with. Perfection of assistants, companions, and even opponents that correctly anticipate and collaborate in the relatively controlled domain of games provides a smooth path to such developments in broader contexts.

We argue that virtual worlds, as found in computer games, are an ideal environment in which to experiment with the interaction between humans and artificial intelligence. There are at least three reasons for this. First, virtual worlds often approach the complexity of the “real world”, while still being under the control of the researcher and completely observable. Second, the agents in virtual worlds are supposed to represent “real people” and are approached as such by the humans who “play” with the virtual world. Third, the potential interactions that players have with the virtual worlds are highly diverse and wide-ranging, which presents a substantial challenge for artificial intelligence to respond to in a reasonable fashion.

In recent years, the number of ways in which human players can interact with games have increased considerably. While ten years ago interaction was almost exclusively through mouse and keyboard or controllers, nowadays games can potentially respond to natural movements and facial expressions captured by a camera, to spoken language, to eye movements, and to signals captured by a variety of sensors. Brain-computer interface (BCI) technology has become more mainstream, offering possibilities for games to respond to a users’ brain activity. Using VR technology, games can respond to movements of players in natural space. AI that can use all these interface elements to make game agents a natural and appreciated partner or opponent for humans can form the basis for advanced AI agents that interact with humans not only in games, but also in the real world.

The research area in which the proposed seminar is rooted is the interaction between humans and game AI, aiming for natural and appropriate responses of computational agents in virtual worlds to human behavior, making use of both traditional interaction technology as well as modern sensor and interaction technology.

The research area lends itself for a wide range of research topics. For the preparation of this seminar, we proposed the following set of sub-topics:

- **Personalized Human-Game AI Interaction:** Humans have different backgrounds, interests, and goals. As such, there is no “one-size-fits-all” interference and interaction form. Under this topic, we explore game adaptation as a type of automatic game design. The goal is to permit the AI to adapt the game environment to the player based on the observed features and received feedback. Instead of fully automatic game design, a

- sophisticated game design leaves scope for an AI to adapt to a broad variety of players. Such personalized adaptation could be extended to adaptation of the actual game interface – in games, usually complex interactions are possible, which novice players are not capable of employing. Therefore, automatically adapting the interface to the observed experience level of the player may be a valid approach to effective personalization.
- **Human-Game AI Interaction for People with Disabilities:** People with disabilities require special attention when designing interfaces, to mitigate adverse effects of disabilities, so that a suitable experience is ensured for everyone. Game AI can potentially help to diagnose disabilities, both physically and psychologically. There is also the potential for game AI to create awareness of issues faced by those with disabilities, by intelligently adapting the interface in such a way that the player experiences it as a person with disabilities would.
  - **Multimodal Interfacing and Interaction:** Multimodal systems offer a flexible and efficient interaction environment that consists of several input/output possibilities including text, speech, and vision. How to effectively use these possibilities in game design is still an open problem. A compelling application of artificial intelligence is to rapidly learn which modes a given player finds natural and enjoyable. The type of interface a user is comfortable with is likely to cross boundaries between different applications, meaning that an “interface fingerprint” may be derivable that can be carried with the user, permitting the re-use of information gained.
  - **Enhancing Human Creativity with Artificial Intelligence:** Computational Creativity is a field of AI where automatic AI systems design and create various forms of art, which may include images, drawings, poetry and music. In the broader sense, these systems create new content either completely by themselves, or with the human providing input at specific points. Research into this fusion of the creative skills of humans and AI systems would move the state of the art a step forward: from being inferior content creators the AI systems would become a tool for amplifying and augmenting the superior creative abilities of a human being, in a bi-directional collaboration process. AI systems should be able to learn from the human, anticipate what they intend to do, and understand the domain of discourse. They would provide advice on content creation and help when the user struggles with certain techniques or creative methods. By learning the skills of the human, AI systems would be able to propose alternatives that lie outside their expertise, allowing the humans to learn, refine and improve their capabilities. The users would experience a system that adapts to their skills, needs and pace, and becomes a personalized companion in their learning process.
  - **Trustful and Reliable Human-Game AI Interaction:** We often observe that humans feel uncomfortable with AI recommendations. Moreover, mistakes made by humans are deemed more tolerable than those made by an AI. While there is no objective rationale for this difference, it is hard to justify the use of AI for humans by arguing that AI offers a lower mistake probability compared to humans. It is therefore imperative to find new ways to convince humans to interact with the game AI and to take its advice seriously. Moreover, it is crucial to minimize any effect that might harm such trust, regardless of its origin.
  - **Information Flow in Human-Game AI Interaction:** A game AI must observe the human player and, in turn, provide players with information that they find helpful, valuable, or interesting. Even the most potentially helpful information is not actually helpful if the player cannot understand it or if it is not useful to their particular style of play. The flow of information is particularly important between the human player and an



AI companion. Reliable metrics that ascertain if the human uses information offered by the AI, that check if the AI fails to provide information that the human tries to find in other ways, and assessment of defects in the human's play that suggest which information is needed, are potential goals of research in this area.

- **Believable Human-Game AI Interaction:** In the last decade, contests have been held at several conferences where human judges voted on the “humanity” of both human game players and AI players in an effort to score the ability of the AI players to behave in a plausibly human manner. Attempts to make AIs interact in a way that is indistinguishable from human interaction are a natural way to structure research into human-game AI interaction. We note that the believability of game AI often suffers because it fails to recognize that it misunderstands the human player, or that the human player misunderstands the AI. How to recognize misunderstanding, followed by how to correct for misunderstanding, are important steps in making game AI more believable.
- **Ethics of Human-Game AI Interaction:** Several of the aforementioned research directions rely heavily on big data analysis. Acquiring such a massive amount of data is a challenging task. Perfect anonymization is hard to achieve, and often undesirable as multiple parties are involved in data collection and integration. To what extent is it ethical to collect personal interaction information? Are there ethical restrictions to the extent to which an AI is allowed to analyze a player's personality and demographics? These questions need answering even if a player gives permission to collect and use such data.
- **Novel Forms of Interaction and Interfaces in Game AI:** New technology gives rise to new possibilities in game interaction and interfacing. While developers often try to restrict themselves to small adaptations in tried-and-true forms of interaction, it makes sense to consider the interaction possibilities originating with novel technology, such as virtual reality and brain-computer interfacing. Beyond those, there may be ways for humans and AI to interact with each other that has not yet been imagined, or which can benefit from re-imagining. Player-AI interaction can be implemented in many forms, such as (1) cuing a player with environmental information from music to decor, (2) influencing a player by adjusting game elements such as local architecture, opponents, and rewards, and (3) making a player respond to the social tone of non-player characters. Such alternate forms of player-AI interaction warrant investigation.

This seminar was organized around workgroups, which worked in teams and topics proposed by the participants of the seminar in the areas outlined above. These workgroups were accompanied by plenary sessions for group formation, topic debate and discussions of the deliberation of each group. Workgroups were dynamic, so participants could move between them, and new groups were formed during the week. A Discord server was setup for coordination and announcements, and it was also used by the different groups for document and link sharing. This also has the benefit of providing a place for discussions after the seminar, easing the communication and further work among the members of each workgroup.

It is worthwhile mentioning the work carried out during the invitation process. Due to the COVID crisis, the changes in the political landscape, and the war in Ukraine, many declined the invitation, and many participants dropped out after originally having accepted the invitation. Thus, multiple rounds of invitations were run until two weeks before the seminar. We invited close to 100 people, the full list of invitations having a high diversity (a 50% male-female split, about half invitations for 'junior' people, and invitees hailing from all continents – including South America and Africa, which are usually highly underrepresented). In the end, just over 30 participants attended the seminar (out of the 45 possible). Size-wise this was a slight disappointment. We were fortunate, however, that those that did attend were highly enthusiastic and highly knowledgeable about the topics covered, which made the seminar a great success.

## 2 Table of Contents

### Executive Summary

*Pieter Spronck, Daniel Ashlock, Setareh Maghsudi, and Diego Perez Liebana . . . . .* 28

### Working groups

#### Language Models for Procedural Content Generation

*Maren Awiszus, Alexander Dockhorn, Amy K. Hoover, Antonios Liapis, Simon M. Lucas, Mirjam Palosaari Eladhari, Jacob Schrum, and Vanessa Volz . . . . .* 34

#### AI for Romantic comedies

*Michael Cook, Maren Awiszus, Duygu Cakmak, Alena Denisova, Alexander Dockhorn, Casper Hartevelde, Antonios Liapis, Mirjam Palosaari Eladhari, Diego Perez Liebana, Lisa Rombout, and Tommy Thompson . . . . .* 37

#### Pokegen

*Alexander Dockhorn, Manuel Eberhardinger, Daniele Loiacono, Diego Perez Liebana, and Remco Veltkamp . . . . .* 39

#### Program Synthesis for Explaining Strategies

*Manuel Eberhardinger, Cameron Browne, Jakob Foerster, Daniele Loiacono, Ana Matran-Fernandez, and Remco Veltkamp . . . . .* 43

#### Artificial Intelligence for Time-Travelling Games

*Ana Matran-Fernandez, Manuel Eberhardinger, Jakob Foerster, Simon M. Lucas, Paris Mavromoustakos Blom, and Pieter Spronck . . . . .* 45

#### Multiplayer Time Travel

*Jakob Foerster, Duygu Cakmak, Simon M. Lucas, Setareh Maghsudi, Ana Matran-Fernandez, Paris Mavromoustakos Blom, Diego Perez Liebana, Lisa Rombout, and Pieter Spronck . . . . .* 49

#### Artificial Intelligence for Audiences

*Antonios Liapis, Maren Awiszus, Alex J. Champandard, Michael Cook, Alena Denisova, Alexander Dockhorn, Tommy Thompson, and Jichen Zhu . . . . .* 50

#### Personalized Long-Term Game Adaptation Assistant AI

*Antonios Liapis, Guillaume Chanel, Alena Denisova, Casper Hartevelde, Mike Preuß, and Vanessa Volz . . . . .* 55

#### The Tabletop Board Games AI Tutor

*Diego Perez Liebana, Duygu Cakmak, Setareh Maghsudi, Pieter Spronck, and Tommy Thompson . . . . .* 60

#### Artificial Intelligence for Alternative Controllers

*Lisa Rombout, Alex J. Champandard, Ahmed Khalifa, Paris Mavromoustakos Blom, and Mark J. Nelson . . . . .* 65

#### Quality Diversity for Procedural Content Generation

*Jacob Schrum, Alex J. Champandard, Guillaume Chanel, Amy K. Hoover, Ahmed Khalifa, Mark J. Nelson, Mike Preuß, and Vanessa Volz . . . . .* 65

#### Benchmarking Coordination Games

*Pieter Spronck, Duygu Cakmak, Jakob Foerster, and Setareh Maghsudi . . . . .* 69

Explainable AI for Games  
*Jichen Zhu, Maren Awiszus, Michael Cook, Alexander Dockhorn, Manuel Eberhardinger, Daniele Loiacono, Simon M. Lucas, Ana Matran-Fernandez, Diego Perez Liebana, Tommy Thompson, and Remco Veltkamp . . . . .* 73

Human-AI Collaboration Through Play  
*Jichen Zhu, Guillaume Chanel, Michael Cook, Alena Denisova, Casper Hartevelde, and Mike Preuß . . . . .* 75

**Panel discussions**


Discussion and Evaluation  
*Pieter Spronck, Setareh Maghsudi, and Diego Perez Liebana . . . . .* 78

**Participants . . . . .** 82

### 3 Working groups

#### 3.1 Language Models for Procedural Content Generation

*Maren Awiszus (Leibniz Universität Hannover, DE), Alexander Dockhorn (Leibniz Universität Hannover, DE), Amy K. Hoover (New Jersey Institute of Technology, US), Antonios Liapis (University of Malta – Msida, MT), Simon M. Lucas (Queen Mary University of London, GB), Mirjam Palosaari Eladhari (Stockholm University, SE), Jacob Schrum (Southwestern University – Georgetown, US), and Vanessa Volz (modl.ai – Copenhagen, DK)*

**License**  Creative Commons BY 4.0 International license  
© Maren Awiszus, Alexander Dockhorn, Amy K. Hoover, Antonios Liapis, Simon M. Lucas, Mirjam Palosaari Eladhari, Jacob Schrum, and Vanessa Volz

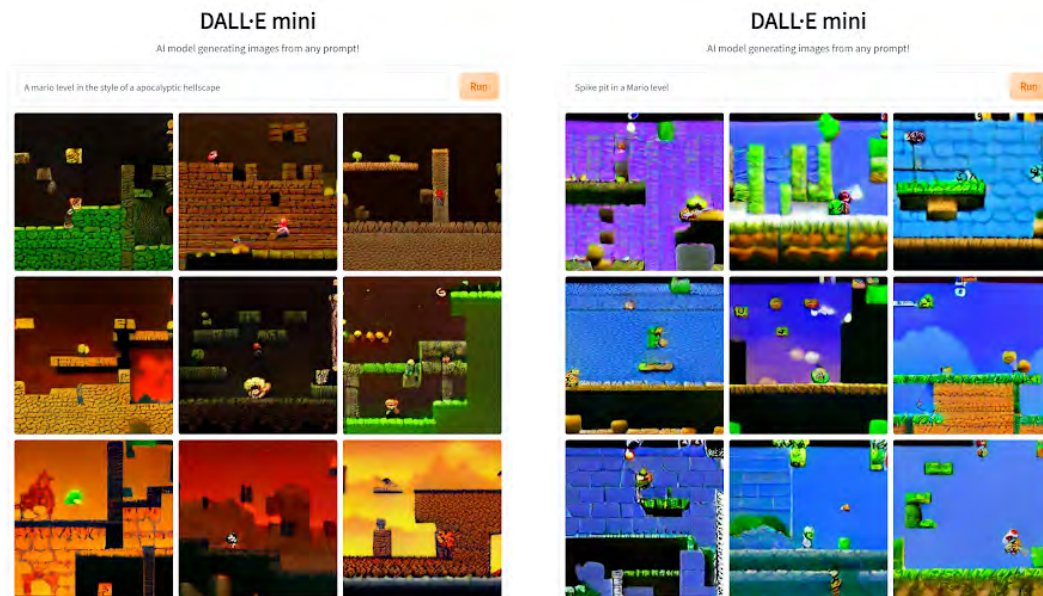
##### 3.1.1 Introduction and Motivation

Recent advances in ML-based image generation via systems like DALL-E [6] beg the question of whether similar tools can be used for generating game content. Specifically, it is desirable to generate game content based on simple text input. As our group was composed of researchers most familiar with level generation, we focused on video game level generation for 2D games first. However, other assets like textures also seem like great examples to use these generative methods on. Our general intuition was, that methods like DALL-E are able to generate impressive previously unseen images due to the strength of the diverse language processing learned with huge datasets of images and their descriptions. For games, in which only little data can be given for a domain like e.g. 2D Super Mario levels, such a large network can not easily be trained from scratch. Therefore, we wanted to investigate the capabilities of a pretrained DALL-E to generate content without any game specific training.

##### 3.1.2 Exploration of Applications

We ran tests using **DALL-E mini** [3] as an intermediate tool for generating content. Figure 1 shows some of those generated examples. Although DALL-E mini can create outputs that look like Mario levels given a prompt like “Mario level”, it has problems incorporating specific details suggested from prompts such as “spikes” or “pipes” in “Mario level with a spike pit” or “Mario level with pipes”. The problem seems to be that DALL-E’s concept of what spikes or pipes are is based on typical photo examples of these items rather than examples of these items in the context of a Mario game. However, prompts that only change the style of the level, like “apocalyptic”, can influence the output. We also did some small preliminary tests on content other than platformer levels, like generating images of new “Pokémon” from a textual description. These examples suffer from similar problems. DALL-E can generate a “Pokémon”, but adding additional descriptors is less likely to be successful. It will be interesting for future work to find out what kind of prompts can and cannot be mixed with DALL-E and why. Especially if this is a tool to be used by game designers, one needs to make sure that the method does not ignore additional descriptors that are a crucial part of the game’s design.

The model **CLIP** [5], which is part of the DALL-E pipeline, can also be used on its own to gauge how well a text description matches an image. We tried matching the images of some original Mario levels to certain prompts that describe a level in more detail, like “under ground Mario level”. For that, we used images of Mario levels provided by the Video Game Level Corpus (VGLC) [7]. The results indicate, that while CLIP does seem to distinguish



■ **Figure 1** Examples of images generated with DALL-E Mini [3]. The generated examples look like Mario levels and certain prompts, such as “apocalyptic” can change the style of the generated images. However, prompts like “spike pit” are ignored.

between “over ground” and “under ground” levels, objects such as pipes do not seem to be recognized as well, which likely explains why pipes cannot be easily generated either. Note, that these results are also without fine-tuning the model at all.

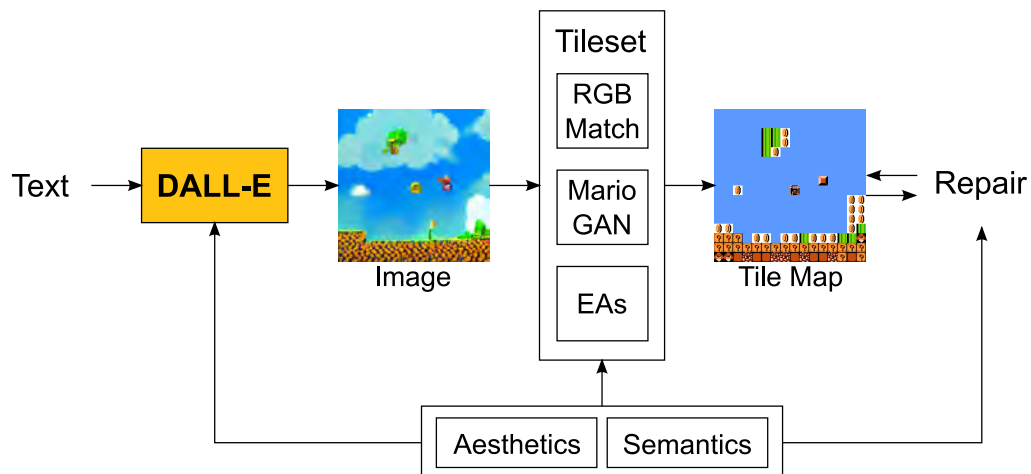
The results of these experiments indicate, that even without any fine-tuning of DALL-E and CLIP, the methods already show some understanding of video game levels. We identify creating a small data set of levels and their appropriate textual descriptions for fine-tuning as an important task to further research in this direction.

### 3.1.3 A Functional Pipeline

While creating images of levels with DALL-E can indicate whether or not the method can be used for level generation in general, this neglects the problem of creating a playable level from that image. Therefore, we established a prototype pipeline for getting playable levels from text, which is shown in Figure 2. Text can be sent to DALL-E to create a level image. Preexisting tools can derive a structured level segment from the image. For now, only the naive RGB tile matching method provided in [2] is applied to this task. Finally, a larger, more complete level can be made from that segment with TOAD-GAN [1]. As the example of a Mario level snippet in Figure 2 shows, the preliminary pipeline works and can create playable Mario level snippets from DALL-E mini. From here, the pipeline needs to be completed by implementing other options to create a tile map from an image, as well as assembling the implemented parts of the pipeline into one cohesive system.

### 3.1.4 Conclusion and Future Work

In this workgroup, we investigated the possibility of using current text to image methods like DALL-E for video game content generation. We show promising results for Super Mario level generation while identifying problems of the method ignoring certain prompts that



■ **Figure 2** The pipeline established at the end of the group session. As shown with the examples, we implemented the pipeline up to the point of being able to generate a tile map from a text prompt. The image is turned into the tile map with an RGB matching algorithm based on [2].

might be important for a game designer. Additionally, we tested if CLIP, a part of DALL-E, can match certain prompts with given Mario level images, and find a similar result: That it can only distinguish some prompts and might ignore others. This however, is only using the pretrained models as is, and we pose that fine-tuning will improve the results for both experiments. We also established a functional text to level pipeline, which can turn a text prompt into a playable Mario level snippet.

For future work, there are two distinct goals: creating a data set to allow for fine-tuning a pretrained DALL-E and completing the missing pieces of the pipeline. For the data set, detailed descriptions of level images need to be found or created and a way to convert them into a usable format needs to be found. Also, other kinds of data sets that deal with assets other than levels can be explored, like texture images. The missing pieces of the pipeline include include other tile set representations that generate a tile map from an image, such as Generative Adversarial Networks and Evolutionary Algorithms, and using the Tile-Pattern KL-Divergence [4] as a repair mechanism for the tile maps. Also, the currently still fragmented pieces need to be combined to form one cohesive system for ease of use.

## References

- 1 Maren Awiszus, Frederik Schubert, and Bodo Rosenhahn. *Toad-gan: Coherent style level generation from a single example*. In Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment, 2020
- 2 Eugene Chen, Christoph Sydora, Brad Burega, Anmol Mahajan, Abdullah Abdullah, Matthew Gallivan, and Matthew Guzdial. *Image-to-level: Generation and repair*. In Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment, volume 16, pages 189–195, 2020
- 3 Boris Dayma, Suraj Patil, Pedro Cuenca, Khalid Saifullah, Tanishq Abraham, Phúc Lê Khac, Luke Melas, and Ritobrata Ghosh. *Dall-e mini*, 7 2021
- 4 Simon M Lucas and Vanessa Volz. *Tile pattern KL-divergence for analysing and evolving game levels*. In Proceedings of the Genetic and Evolutionary Computation Conference, pages 170–178, 2019
- 5 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. *Learning trans-*

- ferable visual models from natural language supervision. In Proceedings of the International Conference on Machine Learning, pages 8748–8763. PMLR, 2021
- 6 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. *Zero-shot text-to-image generation*. ArXiv preprint, abs/2102.12092, 2021
  - 7 Adam James Summerville, Sam Snodgrass, Michael Mateas, and Santiago Onta n'on Villar. *The vglc: The video game level corpus*. Proceedings of the 7th Workshop on Procedural Content Generation, 2016

## 3.2 AI for Romantic comedies

Michael Cook (Queen Mary University of London, GB), Maren Awiszus (Leibniz Universität Hannover, DE), Duygu Cakmak (Creative Assembly – Horsham, GB), Alena Denisova (University of York, GB), Alexander Dockhorn (Leibniz Universität Hannover, DE), Casper Hartevelde (Northeastern University – Boston, US), Antonios Liapis (University of Malta – Msida, MT), Mirjam Palosaari Eladhari (Stockholm University, SE), Diego Perez Liebana (Queen Mary University of London, GB), Lisa Rombout (Tilburg University, NL), and Tommy Thompson (AI and Games – London, GB)

**License** © Creative Commons BY 4.0 International license  
 © Michael Cook, Maren Awiszus, Duygu Cakmak, Alena Denisova, Alexander Dockhorn, Casper Hartevelde, Antonios Liapis, Mirjam Palosaari Eladhari, Diego Perez Liebana, Lisa Rombout, and Tommy Thompson

### 3.2.1 Introduction and Work Process

Both romance and comedy are integral parts of human culture, yet despite the breadth of AI research into games and creativity, little work has been done to explore these themes in the context of games. In AI research, the best examples are games that deal with ‘social physics’ or human relationships, such as *Prom Week* [4] or *Façade* [3], where both romantic and comedic themes are hinted at. In the games industry, while romance is a key feature in many games (such as *The Sims*), it is often reduced to static linear narratives, while comedy is notoriously difficult to achieve in games and is often achieved unintentionally [2].

In this workgroup, we aimed to explore the possibility that these two things are connected. Due to a lack of AI research into topics such as romance and comedy, there are fewer systems and techniques available to support the exploration of these themes in game design. Our workgroup aimed to explore the potential for AI research in these areas, to think about the open questions and pitfalls ahead, and to collaboratively sketch out some ideas for work that we could act as inspiring examples for future AI research projects. The group began with a short presentation, including a series of tweets from @NightlingBug on Twitter, who made an observation that playing a game such as *Stardew Valley* from the perspective of a character competing for the player’s attention would be an interesting idea.

We began with an open discussion of the topic, encouraging perspectives from everyone present, covering both existing examples of technology and games, as well as concerns, questions, and ideas that arose as we thought about the topic. All of the topics that came out of this discussion were interesting and thought-provoking, but a few ideas stood out as something the groups were particularly excited to take forward during the day. The first was the idea of connecting existing AI narrative techniques, such as the Nemesis system in *Shadow of Mordor* [7] to large-group dynamics like the romantic NPCs in *Stardew Valley*. The second idea was to think about how information flow is often crucial in romantic stories, both within

the fiction and between the reader and the author. The third was to investigate unusual concepts such as discomfort, embarrassment, or “cringe” as a component of a narrative or social AI system. The workgroup split into three subgroups to explore these ideas separately, before reconvening at the end of the day.

### 3.2.2 Nemesis Island

The first group proposed an AI-driven spectator sport based on popular reality TV franchises such as *Love Island*. In their prototype, a network of AI agents compete both for the romantic attentions of other AI agents, and the real-world attention of people viewing the game on livestream services, such as Twitch. As a third role, a director can be introduced, whose task it is to steer the narrative by setting hidden internal goals for each agent. Agents respond to the internal social network of the game, the pursuit of their internal goals, as well as their meta-level understanding of the show they are in, intentionally creating drama or showing off to create interest in the audience, in the hope that they will survive rounds of voting and elimination.

### 3.2.3 JANE (Judicious Artificial Narrator Experience)

The second group proposed a game inspired by Jane Austen’s use of free indirect discourse [1], where the author disseminates information to the reader that could be biased by a particular viewpoint, or actual narrative fact [6]. In this approach, the reader always only has partial (and potentially misleading) information on the characters, and they about each other – which leads to both romantic and comedic situations. The setting for this game could be based on shows such as *Bridgerton* or *Gossip Girl*. The player takes the role of a pseudonymous gossip columnist, who must explore and learn about high society by attending events, engaging in gossip, and dealing favours. The columns written by the player impact the knowledge and social simulation of AI socialites, which in turn changes the situations the player finds themselves in. This creates a kind of participatory take on social simulations like *Bad News* [5], with the added complication of allowing the player to engage in high society themselves, potentially manipulating the social scene to help them achieve their personal goals.

### 3.2.4 #CringeFestival

The third group considered the role of embarrassment and negative emotions in romantic comedies. One issue that came up in our initial discussions was understanding the role of the player in such games. As the audience for a romantic comedy, we have a distance between us and the actions of the characters (“cringe” is defined as experiencing embarrassment on behalf of someone else). If the player is participating as a character then they might feel closer to the negative experience. This group explored the idea of games in which the player acts as an external force, either trying to set up artificially embarrassing moments for AI agents, or acting to save and rescue AI agents from embarrassing situations to gain catharsis.

### 3.2.5 Conclusion and Outcomes

Our group discussions have yielded a number of new directions to explore, both in terms of prototyping new systems, as well as exploring the affordances and applications of existing technology. We are hoping to pursue some of these ideas a little further and write the results up, and to continue to maintain the working group as an ongoing collaboration.



## References

- 1 Jane Austen. *The complete novels of Jane Austen*, volume 4. Chartwell Books, 2016.
- 2 Claire Dormann and Robert Biddle. Making players laugh: The value of humour in computer games. In *Proceedings of the 2007 conference on Future Play*, pages 249–250, 2007.
- 3 Michael Mateas and Andrew Stern. Procedural authorship: A case-study of the interactive drama Façade. In *Digital Arts and Culture*, 2005.
- 4 Josh McCoy, Mike Treanor, Ben Samuel, Aaron A. Reed, Michael Mateas, and Noah Wardrip-Fruin. Prom Week: Designing past the game/story dilemma. In *Proceedings of the Foundations of Digital Games Conference*, 2013.
- 5 Ben Samuel, James Ryan, Adam Summerville, Michael Mateas, and Noah Wardrip-Fruin. Bad news: An experiment in computationally assisted performance. In *Proceedings of the International Conference on Interactive Digital Storytelling*, 2016.
- 6 Carmen Smith and Laura Mooneyham White. Discerning voice through austen said: Free indirect discourse, coding, and interpretive (un) certainty. *Persuasions: The Jane Austen Journal On-Line*, 37(1), 2016.
- 7 Ryan Taljonick. Shadow of Mordor’s Nemesis system is amazing—here’s how it works. <https://www.gamesradar.com/shadow-mordor-nemesis-system-amazing-how-works/>, 2014. accessed 3 July 2022.

## 3.3 Pokegen

Alexander Dockhorn (Leibniz Universität Hannover, DE), Manuel Eberhardinger (Hochschule der Medien – Stuttgart, DE), Daniele Loiacono (Polytechnic University of Milan, IT), Diego Perez Liebana (Queen Mary University of London, GB), and Remco Veltkamp (Utrecht University, NL)

**License** © Creative Commons BY 4.0 International license  
 © Alexander Dockhorn, Manuel Eberhardinger, Daniele Loiacono, Diego Perez Liebana, and Remco Veltkamp

The generation of art assets plays a huge part in game development, costing both time and money. We explored how the process of generating game art can be supported using recent advances in generative art.

Machine learning models such as Dall-E 2 [1] and Imagen [2] have demonstrated powerful art generation capabilities. Starting from text prompts, they are able to combine concepts, attributes, and styles to generate artworks of generally high quality. Nevertheless, their usage is restricted and similar projects such as ruDall-E [4] and Mini-Dalle-E [3] do not produce results at the same level of detail, i.e. generating blurry images, struggling to include concepts that are not well represented in the training data, and sometimes creating stock image overlays (c.f. Figure 3). This often results in prompt engineering, a process in which the user adapts the text prompt to guide the black box model to produce the desired outcome [8]. Due to the black-box nature of deep learning models, this process can yield unstable results and is therefore hard to control, making it inefficient and unreliable for creating game assets.

Therefore, we have envisaged several pipelines that may support designers and artists during game development. Starting from possible inputs such as a designer’s textual descriptions of the required game asset, some image ideas, sketches, or even expected game mechanics, we have multiple ways to approach the problem of game asset generation. Simple text and image search models may guide the artistic exploration process and spawn new ideas.

Nevertheless, those can only return results that already exist. Given textual descriptions, we can apply text-to-image models for generating new assets. Alternatively, we may use style-transfer models to enforce characteristics described in the text to an existing image (e.g. CycleGAN [9]). The latter may also be used to adjust image characteristics such as drawing style or the choice of colors (e.g. Neural Style Transfer [10]). Especially interesting is the combination of such models, which may allow to tune each component of the processing chain separately.

In our working group, we have worked on implementing a toolchain to generate Pokemon-like creatures. A Pokemon often represents an animal or object and in terms of visual style, does only consist of a few colors as well as simple shapes and textures. Aiming to use existing models without retraining, we started our process by generating images of dragons using ruDall-E [4]. Generated images varied hugely in quality. Since due to our style constraints our final image does not need to include a lot of details we have chosen a rather blurry image of a dragon with a simple background. Having selected a generated image of a dragon we applied style-transfer as a combination of VQGAN [7] and CLIP [6]. Without retraining any of these components to our specific domain (due to time constraints), we were unable to achieve results of high visual quality (see Figure 4). Nevertheless, this process show-cases how mock-ups and ideas may be generated to guide the development process.

While having struggled to develop a multi-stage model for generating Pokemon-like creatures, it has helped us to better understand the main challenges for generating game assets in general. The following challenges have been identified by us and may guide further research in this domain:

- **Copyright:** Generating art from machine learning models poses the question of who owns the copyright of the final result. This may be a complicated question to answer since the result itself is likely to be a product of an enormous training corpus on which the model is based and the user's input. While there is no definitive answer to this question yet, the current suggestion seems to be an evaluation on a case-by-case basis [5].
- **Training data:** Depending on the stage of production, the amount of available training data may be minor in comparison to the variety of elements that need to be generated. Especially in the early stages of development, machine learning models may merely be used to generate interesting mock-ups or explore ideas. Later development stages may allow to train specialized models or refine existing models to produce desired results.
- **Costs:** Creating your own machine learning models or using the models provided by others can come with non-negligible costs. The required hardware, energy, and time for training and inference should be kept in mind while planning a pipeline. Reducing these costs is already a key aspect of machine learning research and further advancements may considerably reduce the related costs.
- **Usability and Explainability:** Each of the envisaged pipelines comes with its own unique challenges. Especially, the usability of black box models may become a problem in case the input space is not well understood. We have tackled this problem by splitting the asset generation into multiple sub-tasks which we were able to control independently with limited success. Better explaining a model's relation between in- and output as well as its parameter space may help in increasing the usability of such models.

While there are still many steps ahead of us, supporting the generation of game assets using machine learning models may have huge impact on the field. At the current stage, existing models may already be used to support the prototyping stage or generate mock-ups and ideas for the human-guided generation process. Having further advanced on the models' capabilities, it may be possible to learn from just a few examples and produce game assets

Dall-E 2: A dragon in the style of a pokemon, digital art



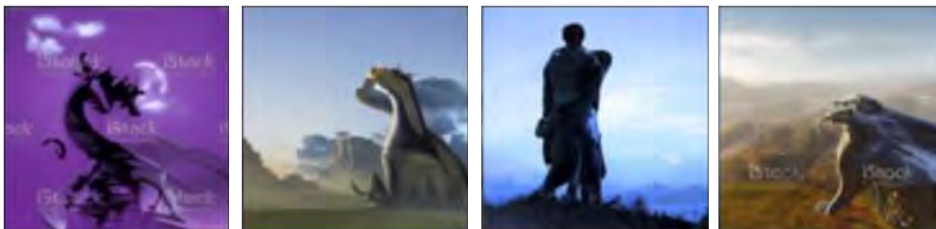
Dall-E 2: A dragon in the style of a pokemon, pixel art



Dall-E 2: dragon standing on a mountain

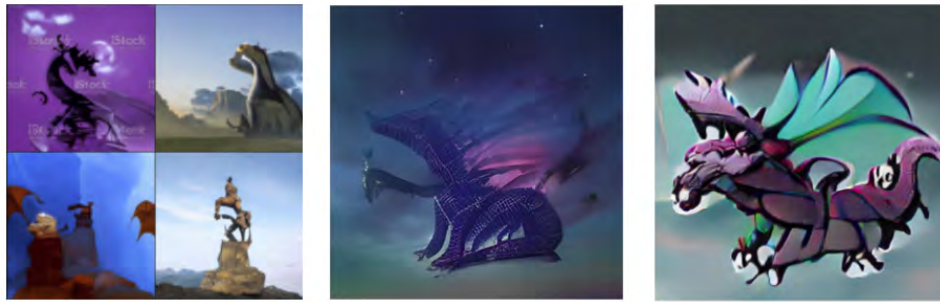


ruDall-E: dragon standing on a mountain



■ **Figure 3** Comparison of Dall-E 2 and the openly accessible ruDall-E for generating dragons and pokemon-like creatures. Dall-E 2 is able to produce images of higher quality. However, it requires an invitation from OpenAI to be used.

of matching styles. In the long run, combinations of machine learning models may even guide the development of whole game worlds and game mechanics, allowing us to generate complete game experiences given a user's queries.



■ **Figure 4** Demo pipeline for generating Pokemon-like creatures. First, generating images of dragons using ruDall-E[4] (left) discarded examples, (middle) chosen example, (right) given the text prompt “A dragon in the style of a pokemon” we used VQGAN [7] and CLIP [6] to produce the final result.

## References

- 1 Ramesh, A., Dhariwal, P., Nichol, A., Chu, C. & Chen, M. Hierarchical Text-Conditional Image Generation with CLIP Latents. (arXiv,2022), <https://arxiv.org/abs/2204.06125>
- 2 Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S., Ayan, B., Mahdavi, S., Lopes, R., Salimans, T., Ho, J., Fleet, D. & Norouzi, M. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. (arXiv,2022), <https://arxiv.org/abs/2205.11487>
- 3 Dayma, B., Patil, S., Cuenca, P., Saifullah, K., Abraham, T., Lê Khâc, P., Melas, L. & Ghosh, R. DALL E Mini. (2021,7), <https://github.com/borisdayma/dalle-mini>
- 4 Shonenkov, A. ruDall-E. (2021), <https://pypi.org/project/rudalle/>
- 5 Chiou, T. Copyright lessons on Machine Learning: what impact on algorithmic art?. (2019)
- 6 Radford, A., Kim, J., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G. & Sutskever, I. Learning Transferable Visual Models From Natural Language Supervision. (arXiv,2021), <https://arxiv.org/abs/2103.00020>
- 7 Esser, P., Rombach, R. & Ommer, B. Taming Transformers for High-Resolution Image Synthesis. (2020)
- 8 Liu, V. & Chilton, L. Design Guidelines for Prompt Engineering Text-to-Image Generative Models. *CHI Conference On Human Factors In Computing Systems*. (2022)
- 9 Zhu, J., Park, T., Isola, P. & Efros, A. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *Computer Vision (ICCV), 2017 IEEE International Conference On*. (2017)
- 10 Gatys, L., Ecker, A. & Bethge, M. Image Style Transfer Using Convolutional Neural Networks. *Proceedings Of The IEEE Conference On Computer Vision And Pattern Recognition (CVPR)*. (2016,6)

### 3.4 Program Synthesis for Explaining Strategies

*Manuel Eberhardinger (Hochschule der Medien – Stuttgart, DE), Cameron Browne (Maastricht University, NL), Jakob Foerster (University of Oxford, GB), Daniele Loiacono (Polytechnic University of Milan, IT), Ana Matran-Fernandez (University of Essex – Colchester, GB), and Remco Veltkamp (Utrecht University, NL)*

License  Creative Commons BY 4.0 International license  
© Manuel Eberhardinger, Cameron Browne, Jakob Foerster, Daniele Loiacono, Ana Matran-Fernandez, and Remco Veltkamp

#### 3.4.1 Introduction & Motivation

Artificial intelligence (AI) in games has attracted a lot of public attention by defeating world champions in board games such as Go or Chess [1, 2]. In eSports, OpenAI trained multiple agents simultaneously that defeated the world champion team in Dota 2, a real-time strategy multiplayer online game where two teams of five members compete against each other [3]. Additionally, DeepMind developed an AI model called AlphaStar that defeated the world champion in the real-time strategy game StarCraft II [4].

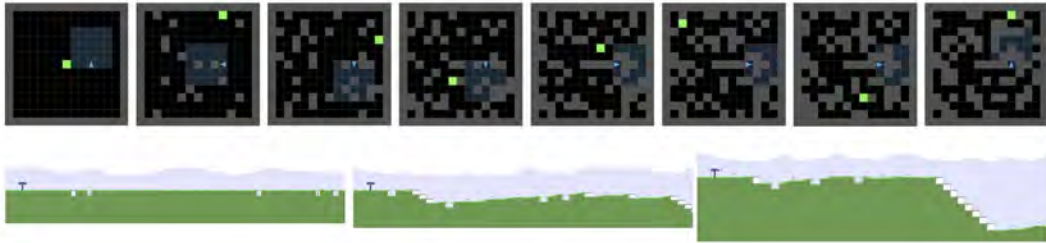
Nevertheless, most strategies of AI models are not explainable or interpretable because a trained neural network is a black box or the search space is too large. This makes it difficult for humans to follow the path of the agent as it traverses the search tree to choose the best action. In this work, we are investigating new ways to make agent behavior interpretable by using program synthesis to explain strategies of agents by distilling black-box policies into programmatic policies.

The rest of this abstract gives a brief introduction to program synthesis for generating programmatic policies and unsupervised environment design (UED), a new approach to providing agents with increasingly difficult environments to create a curriculum for the agent. We conclude this abstract by proposing a method how to combine UED with program synthesis to make game strategies interpretable.

#### 3.4.2 Program Synthesis & Programmatic Policies

In recent years, more and more work investigated program synthesis in reinforcement learning to create programmatic policies for making agent behavior in games or other environments interpretable [5, 6, 7]. Most methods use a form of imitation learning by trying to imitate the behavior of an oracle such as a neural network policy or a human demonstrator. The generation of programs is only possible if the domain-specific language (DSL) is specified and adapted to the task at hand. If the DSL is too general, like a normal programming language, the search for programs is not feasible and leads to no program being found at all [5]. Another way to increase the chances of finding a correct program is to reduce the search space by providing sketches of the program structure where only the missing gaps need to be filled [6]. The problem of finding programmatic policies without defining a task-specific DSL or given prior knowledge about the structure of the program is still an open problem.

However, recent work showed that it is possible to learn a library of functions from previously solved problems. These functions are then reusable in an updated DSL to solve more difficult problems [8, 9]. This leads to a form of curriculum learning by the agent, similar to self-play, as the agent is able to find programs for problems it could not solve before.



■ **Figure 5** Two examples of evolving environments with increasingly difficult levels using the ACCEL algorithm (from [11]).

### 3.4.3 Unsupervised Environment Design

Unsupervised environment design is a method for reinforcement learning in which the agent is given increasingly difficult environments that are still solvable for the agent, but also challenging enough so that the environment is not too easy to master. This discovers a curriculum for agents by always providing the agent with environments that are hardly solvable in the current training process. Dennis et al. [10] train two opposing agents with minimax regret, with one agent coupled to the environmental designer. Regret is the difference between the performance of the two agents, namely how good the agent could be and how good it actually is. This ensures that the levels generated are still solvable.

ACCEL [11] improves this method by using an evolutionary approach to adapt the difficulty of the environments and only trains a single agent for calculating the minimax regret. Figure 5 shows two examples of evolving environments that are increasingly difficult to solve for the agent. The upper environment is the MiniGrid environment [12], which is used to create mazes that the agent has to solve. The lower environment is the bipedal walker environment introduced in [13], where an agent must learn to run over obstacles.

### 3.4.4 Proposed Method

We propose to train a teacher agent that discovers a curriculum of increasingly hard problem sets to challenge a student agent in combination with a program synthesis system such as DreamCoder [8]. This should enable the program synthesis system to explain the behavior of the strategies found, while at the same time the student agent learns to solve increasingly difficult problems. As DreamCoder solves more and more levels by imitating the student agent, the DSL is updated with more representative functions for the environment. This will bootstrap the entire system and makes it possible to learn a custom DSL for the problem, which can be used by human experts to examine agent behavior.

In general, this method proposes a new idea for learning an end-to-end system that can explain game strategies or reinforcement learning policies by finding a tailored DSL for a given problem without using too much prior knowledge, since this knowledge should be found by the system itself. One challenge is the combination and interaction of all mentioned methods into a single system, that can generate programmatic policies and is trainable from scratch.

### References

- 1 Silver, D., Huang, A., Maddison, C., Guez, A., Sifre, L., Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V. & Others, Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*. **529**, 484-489 (2016)

- 2 Campbell, M., Hoane, A. & Hsu, F. Deep Blue. *Artif. Intell.* **134**, 57-83 (2002,1)
- 3 Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C. & Others, Dota 2 with large scale deep reinforcement learning. *ArXiv Preprint ArXiv:1912.06680*. (2019)
- 4 Vinyals, O., Babuschkin, I., Czarnecki, W., Mathieu, M., Dudzik, A., Chung, J., Choi, D., Powell, R., Ewalds, T. & Others, Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*. pp. 1-5 (2019)
- 5 Silver, T., Allen, K., Lew, A., Kaelbling, L. & Tenenbaum, J. Few-Shot Bayesian Imitation Learning with Logical Program Policies. *The Thirty-Fourth AAAI Conference On Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*. pp. 10251-10258 (2020)
- 6 Verma, A., Murali, V., Singh, R., Kohli, P. & Chaudhuri, S. Programmatically Interpretable Reinforcement Learning. *Proceedings Of The 35th International Conference On Machine Learning*. **80** pp. 5045-5054 (2018,7,10)
- 7 Inala, J., Bastani, O., Tavares, Z. & Solar-Lezama, A. Synthesizing Programmatic Policies that Inductively Generalize. *International Conference On Learning Representations*. (2020)
- 8 Ellis, K., Wong, C., Nye, M., Sablé-Meyer, M., Morales, L., Hewitt, L., Cary, L., Solar-Lezama, A. & Tenenbaum, J. DreamCoder: bootstrapping inductive program synthesis with wake-sleep library learning. *PLDI '21: 42nd ACM SIGPLAN International Conference On Programming Language Design And Implementation*. pp. 835-850 (2021)
- 9 Hewitt, L., Anh Le, T. & Tenenbaum, J. Learning to learn generative programs with Memoised Wake-Sleep. *Proceedings Of The 36th Conference On Uncertainty In Artificial Intelligence (UAI)*. **124** pp. 1278-1287 (2020)
- 10 Dennis, M., Jaques, N., Vinitzky, E., Bayen, A., Russell, S., Critch, A. & Levine, S. Emergent Complexity and Zero-shot Transfer via Unsupervised Environment Design. *Advances In Neural Information Processing Systems*. **33** pp. 13049-13061 (2020)
- 11 Parker-Holder, J., Jiang, M., Dennis, M., Samvelyan, M., Foerster, J., Grefenstette, E. & Rocktäschel, T. Evolving Curricula with Regret-Based Environment Design. *International Conference On Machine Learning*. (2022)
- 12 Chevalier-Boisvert, M., Willems, L. & Pal, S. Minimalistic Gridworld Environment for OpenAI Gym. *GitHub Repository*. (2018), <https://github.com/maximecb/gym-minigrid>
- 13 Wang, R., Lehman, J., Clune, J. & Stanley, K. POET: open-ended coevolution of environments and their optimized solutions. *Proceedings Of The Genetic And Evolutionary Computation Conference*. (2019)

### 3.5 Artificial Intelligence for Time-Travelling Games

Ana Matran-Fernandez (University of Essex – Colchester, GB), Manuel Eberhardinger (Hochschule der Medien – Stuttgart, DE), Jakob Foerster (University of Oxford, GB), Simon M. Lucas (Queen Mary University of London, GB), Paris Mavromoustakos Blom (Tilburg University, NL), and Pieter Spronck (Tilburg University, NL)

**License** © Creative Commons BY 4.0 International license  
 © Ana Matran-Fernandez, Manuel Eberhardinger, Jakob Foerster, Simon M. Lucas, Paris Mavromoustakos Blom, and Pieter Spronck

Time-travel has long been explored as a mechanism in science-fiction, particularly books, movies, and, to a lesser extent, video games, all with different degrees of success by reviewers and consumers. In this report we explore different types of time-travel mechanisms that could be offered in video games and discuss some design aspects that need to be considered in such

video games, as well as the components that might need artificial intelligence implementations and the considerations for a successful time-travel experience. A companion report in this collection considers practical implementation and multi-player aspects of time travel.

### 3.5.1 Introduction

There are many types of time-travelling operations that can enhance a player's experience of a video game. Many aspects of time travel have been explored in film and literature, with time-travel paradoxes being fundamental to the plot in films such as *Terminator*, *Back to the Future* and *Source Code*. As far as we are aware only a handful of games have significant time-travel elements, beyond the *Save Game* facility which we discuss below. A key difference between games compared with other narrative media (and indeed normal lived experience), is the assumption of a linear timeline and single common reality in the latter, whereas it is standard to analyse games in terms of game trees (or more generally graphs). Hence what appears as a paradox in a film may simply be an alternative branch of a game-tree in a game.

Notable examples of time-travel games include the following – here we just comment on the time-travel aspects of them, they are all covered extensive Web coverage including Wikipedia entries which we recommend for further reading.

**5D Chess with Multiverse Time Travel:** a mind-bending game where play proceeds on multiple time-lines and moves are made in four dimensions: the normal x-y of a chess-board, plus time and time-line.

**Life is Strange:** The player can rewind recent actions to play them out differently.

**Millennia: Altered Destinies:** The interplanetary civilisation game plays out through 10,000 years with the player time travelling to nurture four races throughout this period with multiple time-related mechanics to help foster their development.

**Deathloop:** Like *Ground Hog Day* each day resets for most characters, and depending which character we play as the aim is to either break this loop, or to maintain it.

Although limited in number, time travel games have a distinct appeal with very positive reviews.

When talking about travelling to the past, the simplest, and least interesting, form of time-travelling would be equivalent to loading a previously saved game, which can be seen as a trip to the past in which a player retains only their memories, but their inventories, etc., remain as they were at the point where the game was saved. For example, in a poker game, this would be equivalent to loading the game you just lost with the knowledge you have from having played it, and using it to change your own strategy and win this second time.

More interesting approaches are those in which a player can travel back in time with items from the future that can be used in the past, or use the trip to the past to influence the further progress of the game, while still being only one copy of the avatar. Following the poker example above, this could mean that one can travel to the past and manipulate the deck of cards before playing the game, so that the player wins.

In a further level of complexity, we explore also the cases in which the trip to the past (in any of the cases mentioned above) involve the cloning of the avatar, so that there are at least two copies of the player's avatar in the game. Whether or not both can be manipulated, and if one is terminated after reaching the point where the trip to the past was triggered, are further design decisions that influence the game play. In the poker scenario, an example of this level would be one in which the player travels back in time and joins the game with their clone, being able to manipulate it in some way to enable their clone to win.



### 3.5.2 Time-travel as a state transition

One of the reasons why time-travelling content is tricky to execute well is the time-travel paradox that could arise when a character travels back in time and leads to logical inconsistencies in the future [1].

To illustrate the time-travel paradox options, let us assume that we are currently in Universe A, which contains a diamond  $D_A$ , and that we can travel back in time and carry it with us to a previous time. If the time-travel operation is implemented so that we are in the same universe A, the options are that we either have two copies of this diamond (which may be a paradox), or that the original diamond  $D_A$  that was in the universe at that time is in some way destroyed, changed, or teleported so that only one  $D_A$  exists. If instead the time-travel operation transfers our character to a different universe B, which contains its own diamond  $D_B$ , then there is no paradox, as the two diamonds can exist in this universe.

Time-travel in a game is an operation that can be implemented as a state transition which needs an additional state. Namely, this is the state we want to travel to. Let a typical game operator be represented as  $S_{n+1} = f(S_n, a_n)$ , where  $S_n$  is the state of the game at time  $n$ , and  $a_n$  is the action taken at that time. Then, we can define a low-level time-travel operator as  $S_k = f(S_n, k, a_n)$ , where  $k$  is a time in the past (we will discuss time-travelling to the future below). Here we envisage that  $k$  is an absolute time, but an alternative is to have  $k$  as a delta to the current time. The choice of types of time-travel that are allowed to a player are therefore embedded in the state transitions, and these specify the rules for the universe of the game. The time travel action  $a_n$  could indicate the inventory the avatar takes back with them to time  $k$ , for example. Therefore, even when visiting a state in the past (i.e.  $k < n$ ), the state may be different as it includes the inventory (and perhaps avatar state) from the future.

Although we are considering time-travel in the most obvious sense of the player's avatar travelling through time, there are other possible time-warping mechanisms such as sending objects or messages through time. All of these have interesting game-play possibilities.

One should also consider which aspects of the game should be fixed and which should be stochastic. For example, a gamer could go back in time once they know the winning combination of the lottery and buy the correct ticket. However, this might also be an unintended consequence of bad design if the seed for certain simulations remains fixed. This could potentially be offset if there was a cost to using time-travelling as an action. For example, there could be a limit to how many times an avatar can time-travel (thus being a limited resource), or perhaps the avatar could be slower or older after each trip in time.

### 3.5.3 Artificial Intelligence in Time-travel

First we consider the role AI can play in the more mechanical aspects of time travel i.e. the enabling of the time-travelling state transition operator, when realistic agent actions are required to reach the required state.

Whereas travelling in time to the past does not necessarily require AI agents, the need for these is clear when the time-travel is in the forward direction, so we will look at these two cases (briefly) separately, and then consider other applications of AI.

#### 3.5.3.1 AI when travelling to a previous state

One of the clear cases where AI agents are needed when time-travelling to the past is when a clone of the player's avatar is created in the trip, and there are now two versions of the avatar during gameplay (the original and the one that has travelled to the past, which is

typically the one that will be controlled by the player). In this case, the actions of the original avatar need to be recreated (no AI needed), but if the clone is acting in a way that interferes with actions already taken by the avatar in the first gameplay, logical inferences about what the player would have done are needed. This highlights the need for AI agents that can model the player's behaviour and act in similar ways. This also raises the question of how intelligent the AI is required to be – for example, should the AI controlled avatar show surprise on encountering a twin they never knew they had?

Furthermore, if the original avatar persists beyond the point where the trip to the past started, their new actions need to be inferred.

### 3.5.3.2 AI when travelling to a future time

The main motivation for implementing AI in a game when time-travelling to the future is the need for extrapolating a new future from the current state in a way that seems plausible. Note that even if intermediate steps are not observed, the generated future should still appear plausible to the person playing the game.

The complexity of this problem depends on the nature of the game – for games with rich narratives it could be both complex and interesting. For example, if the actions of the player have long term effects, then ensuring those actions are in line with the human player's personality is important in order to present compelling visions of the future. To stay with the lottery example, winning the lottery should be unlikely for a character that shuns any form of gambling.

When travelling to a future time, unless we specify all the actions of all agents, or fix all random seeds, then it's reasonable to have a distribution over possible future states, which can be filtered to meet certain criteria before being presented to the player.

### 3.5.3.3 Other applications of AI in Time Travel Games

Beyond enabling the mechanics of achieving plausible states, there is great potential for AI in play-testing time travel games. While AI agents can in principle be used to play-test any game, time-travel games can be especially confusing, making it hard for human designers and players to spot game-breaking loopholes.

AI for play-testing serves two main roles: one is to find bugs, included crashes, and the other is to check the quality of the game-play. The latter is harder to do well, and often relies on having agents of sufficient intelligence to explore the richness of the gameplay and strategic depth. The extra challenges posed by time-travel for AI are as yet unclear.

Adding time travel actions to a game would most likely increase its complexity, as we are increasing the action space. However, this is not necessarily so, as time travel could also break a game, to the point of rendering it trivial from a competitive viewpoint.

## 3.5.4 Conclusions

Time-travelling in games has the potential to be a fun mechanic that could be added to many games, but there are many considerations that need to be taken into account when designing the time-travelling component of the game, particularly where stochasticity is involved. Multi-player games also require special attention, which we cover in a companion report.

## References

- 1 Tobar, G. & Costa, F. Reversible dynamics with closed time-like curves and freedom of choice. *Classical And Quantum Gravity*. **37**, 205011 (2020).

### 3.6 Multiplayer Time Travel

*Jakob Foerster (University of Oxford, GB), Duygu Cakmak (Creative Assembly – Horsham, GB), Simon M. Lucas (Queen Mary University of London, GB), Setareh Maghsudi (Universität Tübingen, DE), Ana Matran-Fernandez (University of Essex – Colchester, GB), Paris Mavromoustakos Blom (Tilburg University, NL), Diego Perez Liebana (Queen Mary University of London, GB), Lisa Rombout (Tilburg University, NL), and Pieter Spronck (Tilburg University, NL)*

**License** © Creative Commons BY 4.0 International license

© Jakob Foerster, Duygu Cakmak, Simon M. Lucas, Setareh Maghsudi, Ana Matran-Fernandez, Paris Mavromoustakos Blom, Diego Perez Liebana, Lisa Rombout, and Pieter Spronck

Being able to travel in time is one of the ancient dreams of humanity and a topic that is explored broadly in popular culture and science-fiction. However, due to the first law of thermodynamics (“the entropy always increases”) it is unlikely that time travel in the real world will ever be possible. In contrast, simulated worlds like computer games do not need to obey the laws of physics and thus, in principle, can offer the ability to time travel. Indeed, there are a number of examples of games in which players can use time travel as part of the gameplay. Crucially though, currently time travel is both limited to specific games and to the single player case. In this report we put forward a proposal for an API that allows extending time travel to arbitrary games and to the multi-player case.

#### 3.6.1 Introduction

It is striking that time travel (TT) has captured the imagination of writers and scientists for centuries and is yet only rarely present in computer games, where it is *actually* possible. Of course, there are exceptions to the rule, but by and large existing computer games do not take advantage of TT. In this proposal we investigate what it would take to “time-travel-fy” arbitrary games. We will also investigate how we can extend this idea to the multi-player case. Imagine a world in which a game designer can easily import the “time-travel package” and use standard functions to deal with the state-keeping, game play logic etc. associated with time travel. In particular, we focus on two aspects of this logic: First of all, we investigate the role of randomness and, secondly, we address multi-player time travel.

#### 3.6.2 Of Lottery Tickets and Dice

One of the crucial issues is that TT allows a player to potentially “hack” the game logic as long as there is any randomness in the game. There are two different potential problems with opposite impact: The first case is the “lottery ticket”, whereby a player could travel back in time and guess the correct lottery ticket which they had observed in the future. In this instance a simple fix is to *re-randomise* the lottery draw during each instance of time travel. However, re-randomisation has a separate issue: It allows the player to “keep trying” until they obtain the outcome that they want and continue the gameplay from there. For example, when a unit attacks a stronger unit it might still have a finite probability of success and the player could travel back in time until they “get lucky”. To mitigate this, some circumstances require *freezing the randomness* across time, rather than re-running random events on every path forward through time.

Finally, addressing both issues requires a higher-level “semantic understandig” of outcomes. A player should not be able to *ceteris paribus* obtain a better outcome by traveling back in time and *hacking the randomness*. In other words, if the player didn’t win the lottery on the initial travel through time, they should not be able to do so on the second attempt. How to implement this using game AI is an open problem that we hope to address in future work.

### 3.6.3 Multi-Player Time Travel

Time-travel in the single player case closely resembles saving the game state and reloading past checkpoints later on. In contrast, in the multiplayer case things get a lot more interesting. When a number of players co-exist in the same environment it is unclear how time travel of one player should change the current time step of other players. A naive approach is to simply use the single-player option, whereby all players are “dragged through time” by the time travel decisions of any player. However, this will likely make for a confusing playing experience and also break game dynamics since any player might be incentivised to travel in time when things are not working well for them.

Instead, we suggest a new approach for multi-player TT which relies on *branching timelines*: Any player can travel back in time *independently* while all other players have the option to continue playing on their current timeline or TT independently. This leaves a crucial question: What are the characters of other players doing on branches that the player is not currently playing? We suggest to use “zombie-actors”, i.e. machine learning models that predict the actions of a player in the *alternate* reality given their realised actions in the played reality. This problem is similar to the issues caused by time-delay in multi-player games, which are solved e.g. with *Rollback* which is now being improved using machine learning [1].

To reduce computational overhead and avoid pure “zombie-games”, branches that have been abandoned by all players get frozen in time until a player rejoins said branch.

#### References

- 1 Anton Ehlert. *Improving input prediction in online fighting games*. 2021

## 3.7 Artificial Intelligence for Audiences

*Antonios Liapis (University of Malta – Msida, MT), Maren Awiszus (Leibniz Universität Hannover, DE), Alex J. Champandard (creative.ai – Wien, AT), Michael Cook (Queen Mary University of London, GB), Alena Denisova (University of York, GB), Alexander Dockhorn (Leibniz Universität Hannover, DE), Tommy Thompson (AI and Games – London, GB), and Jichen Zhu (IT University of Copenhagen, DK)*

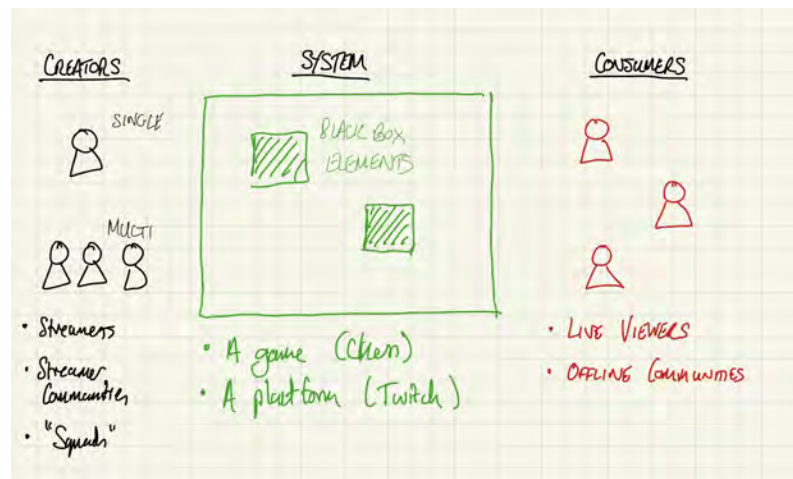
**License** © Creative Commons BY 4.0 International license

© Antonios Liapis, Maren Awiszus, Alex J. Champandard, Michael Cook, Alena Denisova, Alexander Dockhorn, Tommy Thompson, and Jichen Zhu

Artificial Intelligence (AI) has been leveraged for assisting individual players [20, 12] and individual designers or creators [9], but the rise of for-profit content creation platforms [3], and games as a spectacle [1] opens a new and exciting opportunity for AI support. In this working group, we explore applications, algorithms, and interfaces for *AI for audiences*.

The simplest inception of an AI application in this vein would be as mediator between a content creator (e.g. a YouTuber or a Twitch streamer) and the consumers that may be enjoying this content in real-time (e.g. during a stream) or asynchronously (e.g. watching a YouTube video). Focusing on the communication between audience and content, the working group identified the following non-exhaustive list for possible AI roles:

- **AI as mediator.** For instance, the AI may inform a viewer when the content changes (e.g. a new game area is entered or the creator changes the discussion topic), or inform a live-streamer when audience engagement shifts (in tone, volume, or discussion topic).



■ **Figure 6** Envisioned AI as mediator between an audience and one or many content creators.

- **AI as entertainer.** For instance, the AI can add a (textual) commentary to a playthrough in real-time. In this role, the AI may act as an *unreliable narrator*, in which case the state of the game need not be described reliably in order to increase engagement through uncertainty and curiosity. Similar patterns are observed in e.g. e-sport competitive matches, where (human) casters give more “optimistic” predictions for a comeback of the currently losing team.
- **AI for hype.** For instance, the AI can algorithmically generate audio, visual, or text assets to promote content scheduled in the future by connecting it with past content from the same creator or a broader context. Similarly, the AI can promote existing content to the audience based on more in-depth patterns (e.g. gameplay progression) and player/viewer models than current recommender systems.
- **AI as tutor.** For instance, when requested by a viewer an AI could explain game mechanics and their interactions as relevant to the current context. The issue of personalisation is pertinent here, as modeling the viewer’s expertise (based on the number of similar content they have viewed or games they have played, as well as questions they have asked the AI) could impact the level of explanation and possible examples or anchor points to scaffold the explanation.
- **AI as filter of needless data.** For instance, an on-demand AI can jump to the highlights in the video, or an always-on AI can remove uninteresting or toxic chat between audience members.

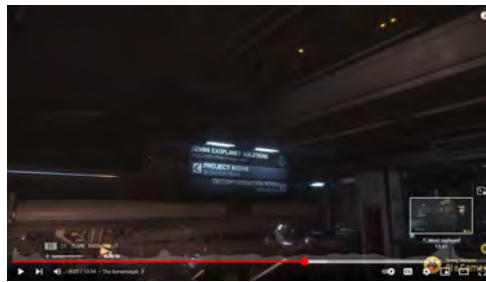
The issue of synchronous versus asynchronous engagement can heavily impact the affordances and constraints for both the AI algorithms and the user interfaces. Beyond the obvious fast-response and low-latency requirements, the issue is pertinent because synchronous viewing may foster shorter but more direct interactions between content creator and audience and between members of the audience (e.g. chat). Synchronous viewing opens additional opportunities for AI assistance, such as a personalized recap of the stream so far in case a viewer joins late, or a recap of events while the user was away in case they leave and rejoin. On the other hand, asynchronous viewing allows for more thoughtful discussions to emerge in comments; at the same time interaction with the content is more granular and controlled as viewers can choose which parts of the video to view, rewind, etc.



(a) Lichess turn-by-turn replays with predicted wins and suggested moves.



(b) DotA 2 real-time match progression with gold, experience, deaths, and predicted win chances.



(c) YouTube viewership analytics, including the a “most replayed” label for popular video sections.

■ **Figure 7** Current examples of visualizations, analytics, and predictions intended for audiences.

Note, that the data format of the content that is made available to the AI should ideally not be simply the end-product (e.g. a video) but additional meta-data regarding game actions, context, and potentially even game-specific AI game players. An example of such rich data is provided in *lichess*<sup>1</sup> where viewers (or players after the game is completed) can watch replays of chess matches along with AI-based predictions of win versus loss after every move, as well as suggested moves instead of the one played. Beyond chess, having access to such granular game data could allow for highlight detection (e.g. at points where the predictions shift dramatically between players), summarization (e.g. grouping similar moves together and focusing on highlights), or tutoring (e.g. showing the causal links between early choices and later outcomes). To maximize the potential of such an approach, however, the game developers would need to provide not only game state and action events but also ideally some game-specific AI that could provide nuanced context-specific metrics such as predicted win probability or chosen next moves. Such meta-data and AI-predicted game metrics are already made available for certain games that embrace the game as spectacle philosophy, especially e-sports such as *Dota 2* (Valve, 2013).

However, AI for audiences need not rely on the assumption of a *one-to-many* interaction, or the implicit assumption that the audience consists of passive consumers with no agency over the content or how they interact with it. AI for audiences can be used to promote and support *augmented communities*, where some or all of the audience members can take more proactive roles (indicatively, live commentators with AI visualization assistance or cinematographers by creating custom camera positions in live or replay game data). Audience interactions

<sup>1</sup> <https://lichess.org/>

with the AI itself can also lead to improved computational models, including player models [26, 19] that can provide personalized tutoring (based on detected expertise level) but also for matchmaking between audience members (especially those with proactive roles). Similarly, the AI can operate on a *many-to-many* assumption and find similar content with similar game-states from other streamers to propose to viewers, but also for matchmaking between content creators. The simplest form of AI for content creators could suggest scheduling clashes with popular content creators in the same genre (or followed by the same audience) or niche topics that have not been explored by other content creators. A more proactive AI could also act as a matchmaker between content creators, suggesting ideas on how and on what topic this collaboration could be built on. Algorithms and interfaces for this type of AI assistance can have broader ramifications, as similar many-to-many relationships can be found in crowdfunding platforms (e.g. Kickstarter), virtual crowd working platforms (e.g. Fiverr or creative.ai), and service providers more broadly (e.g. Uber, Wolt).

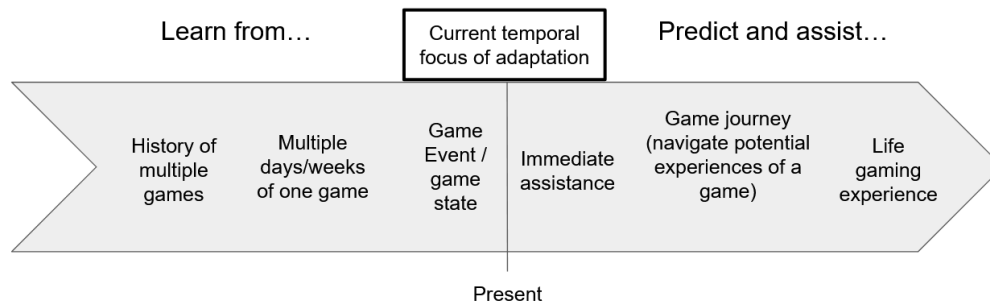
Several existing algorithmic advancements can be leveraged towards the goals laid out above, including recommender systems [19, 12], text summarisation [13, 21], personalisation [10] and personas [6], highlight detection [12], video indexing and matching [22], viewership analytics [7], coordination and scheduling [2], monetisation and churn prediction [8], expressive range analysis [16] and quality-diversity search [4], AI directors [11, 17], and more. However, novel AI research will be warranted in this vein tailored to the format (video, speech, and game meta-data) and user requirements of such applications. Example directions for AI research include question-answering systems (including natural language processing), text summarisation of real-time expanding datasets (of comments or gameplay), context-aware detection of video segments (e.g. based on text mentions in the comments), or causal models [14] based on audio, visual, video, gameplay, and comment/chat data.

## References

- 1 Dave Boling. How The International became a global “Super Bowl for nerds”. [https://www.espn.com/esports/story/\\_/id/20343989/super-bowl-nerds-dota-2-fans-globe-lured-spectacle-camaraderie-international-7](https://www.espn.com/esports/story/_/id/20343989/super-bowl-nerds-dota-2-fans-globe-lured-spectacle-camaraderie-international-7), 2017. Accessed 27 July, 2022.
- 2 Elisabeth Crawford and Manuela M. Veloso. Learning to select negotiation strategies in multi-agent meeting scheduling. In *Proceedings of the Portuguese Conference on Artificial Intelligence*, 2005.
- 3 Cecilia D’Anastasio. Amazon’s Twitch seeks to revamp creator pay with focus on profit. <https://www.bloomberg.com/news/articles/2022-04-27/amazon-s-twitch-seeks-to-revamp-creator-pay-with-focus-on-profit>, 2022. Accessed 27 July, 2022.
- 4 Daniele Gravina, Ahmed Khalifa, Antonios Liapis, Julian Togelius, and Georgios N. Yannakakis. Procedural content generation through quality-diversity. In *Proceedings of the IEEE Conference on Games*, 2019.
- 5 Fabian Hadiji, Rafet Sifa, Anders Drachen, Christian Thureau, Kristian Kersting, and Christian Bauckhage. Predicting player churn in the wild. In *Proceedings of the IEEE Conference on Computational Intelligence in Games*, 2014.
- 6 Christoffer Holmgård, Michael Cerny Green, Antonios Liapis, and Julian Togelius. Automated playtesting with procedural personas through MCTS with evolved heuristics. *IEEE Transactions on Games*, 11(4):352–362, 2019.
- 7 Andrew Hutchinson. Youtube rolls out activity graph to all videos, ups the maximum price of channel memberships. <https://www.socialmediatoday.com/news/youtube-rolls-out-activity-graph-to-all-videos-ups-the-maximum-price-of-ch/624036/>, 2022. Accessed 27 July, 2022.
- 8 Erik Johnson. A deep dive into Steam’s Discovery Queue 2. <https://www.gamedeveloper.com/business/a-deep-dive-into-steam-s-discovery-queue>, 2019. Accessed 6 July, 2022.

- 9 Antonios Liapis, Gillian Smith, and Noor Shaker. Mixed-initiative content creation. In Noor Shaker, Julian Togelius, and Mark J. Nelson, editors, *Procedural Content Generation in Games: A Textbook and an Overview of Current Research*, pages 195–214. Springer, 2016.
- 10 Santiago Ontanon and Jichen Zhu. The personalization paradox: The conflict between accurate user models and personalized adaptive systems. In *Companion Proceedings of the International Conference on Intelligent User Interfaces*, page 64–66, 2021.
- 11 Mark O. Riedl, H. Chad Lane, Randall Hill, and William Swartout. Automated story direction and intelligent tutoring: Towards a unifying architecture. In *Proceedings of the AIED Workshop on Narrative Learning Environments*, 2005.
- 12 Charlie Ringer and Mihalis A. Nicolaou. Deep unsupervised multi-view detection of video game stream highlights. In *Proceedings of the International Conference on the Foundations of Digital Games*, 2018.
- 13 Alexander M. Rush, Sumit Chopra, and Jason Weston. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2015.
- 14 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- 15 Adam M. Smith, Chris Lewis, Kenneth Hullet, Gillian Smith, and Anne Sullivan. An inclusive taxonomy of player modeling. Technical Report UCSC-SOE-11-13, 2011, University California Santa Cruz, 2011.
- 16 Gillian Smith and Jim Whitehead. Analyzing the expressive range of a level generator. In *Proceedings of the FDG workshop on Procedural Content Generation in Games*, 2010.
- 17 Tommy Thompson. In the directors chair: The AI of Left 4 Dead. <https://medium.com/@t2thompson/in-the-directors-chair-the-ai-of-left-4-dead-78f0d4fbf86a>, 2014. Accessed 27 July, 2022.
- 18 Tommy Thompson. How Forza’s Drivatar actually works. <https://www.gamedeveloper.com/design/how-forza-s-drivatar-actually-works>, 2021. Accessed 6 July, 2022.
- 19 Hao Wang, Naiyan Wang, and Dit-Yan Yeung. Collaborative deep learning for recommender systems. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.
- 20 Georgios N. Yannakakis, Pieter Spronck, Daniele Loiacono, and Elisabeth André. Player modeling. In Simon M. Lucas, Michael Mateas, Mike Preuss, Pieter Spronck, and Julian Togelius, editors, *Artificial and Computational Intelligence in Games*, volume 6 of *Dagstuhl Follow-Ups*, pages 45–59. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2013.
- 21 Kevin Yauris and Masayu Leylia Khodra. Aspect-based summarization for game review using double propagation. In *Proceedings of the International Conference on Advanced Informatics, Concepts, Theory, and Applications*, 2017.
- 22 Xiaoxuan Zhang, Zeping Zhan, Misha Holtz, and Adam M. Smith. Crawling, indexing, and retrieving moments in videogames. In *Proceedings of the Foundations of Digital Games Conference*, 2018.





■ **Figure 8** Horizon of past user activities and horizon of predicted future trends, used for choosing actions to take on the part of the assistive AI.

### 3.8 Personalized Long-Term Game Adaptation Assistant AI

*Antonios Liapis (University of Malta – Msida, MT), Guillaume Chanel (University of Geneva, CH), Alena Denisova (University of York, GB), Casper Hartevelde (Northeastern University – Boston, US), Mike Preuß (Leiden University, NL), and Vanessa Volz (modl.ai – Copenhagen, DK)*

**License** © Creative Commons BY 4.0 International license  
© Antonios Liapis, Guillaume Chanel, Alena Denisova, Casper Hartevelde, Mike Preuß, and Vanessa Volz

Artificial Intelligence (AI) in games has already been extensively used for the purposes of modeling players [26, 19]. This working group viewed issues of player modeling through the lens of **personalized assistance**, focusing on the horizon(s) that such models could use to **learn from the past** and to **predict the future**. While the scope and purpose of player models can vary [19], this working group focused on models of an individual person/player which can be generative in scope, i.e. generating “data where a human player could otherwise be consulted” [19].

#### What constitutes long-term personalization?

As a central issue of the topic tackled by the working group was the “long-term” aspect, it is important to define the scope of such temporal and contextual information. As depicted in Figure 8, the main questions on this aspect revolve around (a) the horizon of past user actions (and their context) that the model will learn from, and (b) the horizon of future expectations that the AI can predict and assist towards.

The working group identified that a personalized computational model could learn patterns from a very long-term history of player behavior at low granularity (with metrics such as game purchase behavior and/or playtime), a mid-term history within one game (spanning e.g. multiple days or weeks), or short-term history spanning a few actions or game-states within the current game context. In terms of what predictions the model could make, the working group similarly identified short-term predictions regarding actions within the current game session (e.g. whether the player would fail in an upcoming challenge), mid-term predictions regarding behaviors within the same game (e.g. which parts of future game content the player would enjoy and how), or longer-term predictions (e.g. when the player would quit this game, or which games they would pick up after it).

We note here that an assistant AI does not necessarily require prediction of future states in order to provide assistance, as it can operate without output [27] by detecting patterns between this player and a broader player corpus through unsupervised learning (as would be the case for recommender systems, for instance). However, the context of the assistance similarly fits the same time-scales as player predictions, from short-term assistance regarding e.g. a current problem the player is facing in this phase/location of the game versus long-term assistance in terms of e.g. similar games they can play once they finish this game.

At this level of granularity, there is an abundance of examples to draw from in commercial and research applications modeling past and future horizons. Indicatively, player profiles on Steam take into account game purchases in order to provide similar games to recommend in the player’s discovery queue [12] (long-term past for long-term future). On the other end of the spectrum, AI replicants [16] that follow the low-level decision-making of a single player in a singular context (e.g. level layout and opponent) can be used to simulate the next in-game actions (short-term past for short-term future). On a more realistic mid-term assistance, the Drivatar models [20] learn how to drive as a player would and can generate complete playthroughs through sequences of short-term decisions in the style of the player, even in unseen tracks (mid-term past for short-term future). While not explicitly aimed to assist players, similar work on churn prediction [8] focuses on learning patterns from a player’s long- or mid-term gameplaying and purchase history in order to predict when players may quit playing the game (long-/mid-term past for mid-term future); these predictions are often used to provide players with interesting content or power-ups in the short-term in order to delay players from dropping out.

### How can the AI assist a player?

While to a large extent the algorithms for player modeling are already mature, a more pertinent issue relates to the type of assistance that such models can offer to players. Through extensive brainstorming, the working group identified the following non-exhaustive list of assistant AI actions:

- **Game Selection:** The assistant AI can suggest new games for the player to explore. In this level of granularity, the AI does not adapt any game content and relies on human-authored games that are better suited for this player.
- **Modification of an initial game state:** The assistant AI provides the player with new content within the same game, modifying the initial state via e.g. a new game level to explore [25], making a new mechanic available [1, 3], or new opponent abilities [13, 18]. The distinction here is that the assistant AI adapts the *possibility space* offered to the player without hand-holding the player on how they should take advantage of these possibilities. This assistance is also highly relevant in terms of generating *end-game content* through e.g. recombining existing hand-authored content in novel ways that match player preferences, performance, or expectations.
- **Mechanics adaptation:** The assistant AI interferes in a more granular manner on the moment-to-moment playthrough by adjusting the game mechanics themselves. This could for example take the form of aim assistance by increasing the leniency on what constitutes a hit in a shooter game (e.g. [23]). This same adaptation, due to its subtle nature, could also be used for increasing accessibility in games that would normally require fast reflexes [21].
- **Adapting the player’s behavior towards a normative gameplay goal:** Rather than changing the game according to the player’s preferences, this assistant AI shoehorns the player into playing the game as-is according to the designers’ (rather than the players’) intentions. This assistive AI can take two complementary roles, guiding players during

their playthrough towards intended outcomes by *scaffolding and mediating their learning* and by *nudging them towards desirable behaviors*. Scaffolding can be done through generated tutorials on overlooked game mechanics [7] or on-demand hints regarding actions – or in-game knowledge – needed to overcome a current challenge (e.g. a puzzle). Nudging and priming on the other hand can be achieved by making certain decisions seem more appealing, and has been used extensively in both advertising [24] and teaching [4]. In order to guide a player towards a specific level traversal path, for example, a specific path may be adapted to have less clutter in order to guide players through it [22], sound emitters could be used to guide players towards their source [14], or user interface elements (e.g. quest markers) could be adapted to be more or less prominent. This type of assistance is perhaps the least explored academically in the context of games while also the most promising and realistic from the perspective of the game industry. This is relevant to the game industry as game developers can thus streamline a singular play experience, while taking advantage of existing – carefully crafted – assets without requiring unpredictable generation or adaptation.

- **Providing Reflection and Explainability:** The assistant AI provides feedback to the player regarding their performance or playstyle, allowing players themselves to reflect on how to improve the former or diversify the latter. While post-game summaries abound in games, the AI aspect can be leveraged to provide personalized feedback (e.g. focusing on presenting metrics that are important to this player, based on their personal model) or for highlighting aspects or portions of the playthrough where alternative decisions or actions could have led to better results – as a form of post-game scaffolding. Moreover, post-game visualizations can also serve to explain certain AI decisions or prompts during the game covered in other AI actions above. For example, in a racing game the player's trajectory on the track is shown in a post-game summary, juxtaposed with an AI driver's trajectory and highlighting the points where the AI detected (and verbalized) a hint towards course correction. Therefore this AI action can be a standalone component or an accompanying explanation for other AI actions.

The issue of assistance was central in this particular envisioned application of AI, specifically regarding how (in)visible the assistant would be (e.g. performing difficulty adjustment or aim assistance behind the scenes versus coaching players to better handle the game's challenges). Relatedly, whether the assistance would be on-demand by the player or always in effect would impact the type of assistance the AI can provide as well as issues of players' perception of the AI and explainability requirements. Based on the type of assistance and how it is presented to the player, such AI could take the role of salesperson, tutor, gamemaster, commentator, tour guide, and even as general on-demand virtual assistant similar to Siri or Google Assistant but within the game.

### What can the AI assist towards?

As a final dimension regarding the goals of the assistant, the player model could be trained to focus on a variety of metrics or key performance indicators (KPIs) of the player. The following non-exhaustive list covers some KPIs of interest:

- **Emotional state:** a player's emotional state in the game. AI assistance that keeps track of and aims to improve such a KPI could tailor content towards the intended emotional state (e.g. fear [6] or stress [15] in horror games) or in order to course-correct in case experienced emotions are overwhelming (e.g. in games for rehabilitation [9]).

- **Performance:** the difficulty or challenge a player faces in the game and how they overcome it (e.g. number of retries or game score). AI assistance can be used to tailor the experience to the player’s skills. This is the most traditional application through dynamic difficulty adjustment [10], but can be enhanced beyond invisible rubber-banding through e.g. on-demand coaching and personalized assistance.
- **Coverage:** how much of the game a player has explored (or tends to explore). Coverage can refer to spatial coverage (e.g. heatmap of the level), action coverage (e.g. whether the player makes use of all mechanics and dynamics [11] available to them), narrative coverage (e.g. which non-player-character relationships the player has focused on and how), or temporal coverage (e.g. build orders in strategy games).
- **Learning:** how much the player has mastered the game’s mechanics (or concepts) and improved their repertoire. This KPI could be especially meaningful for assistance on mechanics that the player has overlooked (e.g. via hints and tutorials) or has trouble with (e.g. via aim assist that progressively becomes less pronounced). The aspects of the game that the player has mastered can also inform recommendations for future games that specifically offer additional challenges (or content) in these specific aspects.
- **Intentions:** why and how a player likes to engage with the game. This is particularly meaningful for detecting cases where certain play behaviors are not due to poor performance but due to conscious decisions to play subversively or towards the player’s own goals. A broader example of this would be speedrun challenges [17], where assistance should be tailored not to guide players towards maximizing spatial coverage but towards shortcuts or skill-based shorter traversal paths.
- **Social experience:** how a player interacts with other players in the game. This KPI is mostly relevant to multi-player games with a social component, such as massively multiplayer online games rather than competitive brawlers or racing games. The AI model could capture cases of toxic behavior or interactions within and outside the game, and its “assistance” could include warnings in cases of toxic behavior [2] either to the perpetrator or to new players interacting with them.
- **Monetization:** how a player spends their real-world money in the game, and towards which content. This KPI is less relevant for the research purposes of this working group, but could be relevant for industrial use cases. Moreover, extensive AI research on churn prediction and analytics [5, 8] has been motivated by monetization and thus can not be overlooked. Ethical assistance in terms of this KPI could focus on the explainability and reflection aspect, highlighting to the player which purchase behaviors they tend to make and coaching them towards diversifying or reducing their in-game purchases.

### Envisioned use-cases of assistant AI

After these high-level implications of personalized AI assistance were laid out, the working group focused on two practical examples, both including nudging player behaviors to better experience a game as-is (without new content being generated for it). The first example was focused on narrative-based games with explicit role-playing decision points. Speculative work under this more narrow use-case explored different ways of nudging the player’s decisions in terms of their *invasiveness* (i.e. how much the AI takes over the decision-making) and the *subconscious nature* of the nudging (i.e. whether the player might understand that they are being manipulated). The second example was focused on assisting players to brake within a racing game, which relies on kinesthetic player behavior and the AI aims to reduce the challenge of an existing game. Speculative work to address this issue identified audiovisual feedback in real-time as the most intuitive way of AI assistance, exploring how audio or visual

feedback could be more or less invasive (e.g. a popup foreshadowing the type of upcoming turn, versus a ghost trail of the ideal trajectory for taking the turn). The two practical examples allowed for some more in-depth discussion on the specific challenges that would need to be addressed when developing personalized AI assistants.

## References

- 1 Eric Butler, Adam M. Smith, Yun-En Liu, and Zoran Popovic. A mixed-initiative tool for designing level progressions in games. In *Proceedings of the ACM Symposium on User Interface Software and Technology*, page 377–386, 2013.
- 2 Alessandro Canossa, Dmitry Salimov, Ahmad Azadvar, Casper Harteveld, and Georgios Yannakakis. For honor, for toxicity: Detecting toxic behavior through gameplay. *Proceedings of the ACM CHIPLAY Conference*, 2021.
- 3 Michael Cook, Simon Colton, Azalea Raad, and Jeremy Gow. Mechanic Miner: Reflection-driven game mechanic discovery and level design. In *Applications of Evolutionary Computation*, volume 7835, LNCS. Springer, 2012.
- 4 Mette Trier Damgaard and Helena Skyt Nielsen. Nudging in education. *Economics of Education Review*, 64:313–342, 2018.
- 5 Magy Seif El-Nasr, Anders Drachen, and Alessandro Canossa. *Game Analytics: Maximizing the Value of Player Data*. Springer, 2013.
- 6 Magy Seif El-Nasr, Simon Niedenthal, Igor Kenz, Priya Almeida, and Joseph Zupko. Dynamic lighting for tension in games. *Game Studies*, 7(1), 2007.
- 7 Michael Cerny Green, Ahmed Khalifa, Gabriella A. B. Barros, and Julian Togelius. “Press Space to Fire”: Automatic Video Game Tutorial Generation. In *Proceedings of the AIIDE workshop on Experimental AI in Games*, 2018.
- 8 Fabian Hadiji, Rafet Sifa, Anders Drachen, Christian Thureau, Kristian Kersting, and Christian Bauckhage. Predicting player churn in the wild. In *Proceedings of the IEEE Conference on Computational Intelligence in Games*, 2014.
- 9 Christoffer Holmgård, Georgios N. Yannakakis, Karen-Inge Karstoft, and Henrik Steen Andersen. Stress detection for PTSD via the StartleMart Game. In *Proceedings of the Conference on Affective Computing and Intelligent Interaction*, page 523–528, 2013.
- 10 Robin Hunicke. The case for dynamic difficulty adjustment in games. In *Proceedings of the International Conference on Advances in Computer Entertainment Technology*, 2005.
- 11 Robin Hunicke, Marc Leblanc, and Robert Zubek. MDA: A formal approach to game design and game research. In *Proceedings of AAAI Workshop on the Challenges in Games AI*, 2004.
- 12 Erik Johnson. A deep dive into Steam’s Discovery Queue 2. <https://www.gamedeveloper.com/business/a-deep-dive-into-steam-s-discovery-queue>, 2019. Accessed 6 July, 2022.
- 13 Ahmed Khalifa, Scott Lee, Andy Nealen, and Julian Togelius. Talakat: Bullet hell generation through constrained map-elites. In *Proceedings of the Genetic and Evolutionary Computation Conference*, page 1047–1054, 2018.
- 14 Daryl Marples, Duke Gledhill, and Pelham Carter. The effect of lighting, landmarks and auditory cues on human performance in navigating a virtual maze. In *Proceedings of the Symposium on Interactive 3D Graphics and Games*, 2020.
- 15 Paraschos Moschovitis and Alena Denisova. Keep calm and aim for the head: Biofeedback-controlled dynamic difficulty adjustment in a horror game. *IEEE Transactions on Games*, 2022.
- 16 Johannes Pfau, Antonios Liapis, Georg Volkmar, Georgios N. Yannakakis, and Rainer Malaka. Dungeons & Replicants: Automated game balancing via deep player behavior modeling. In *Proceedings of the IEEE Conference on Games*, 2020.

- 17 Tom Phillips. World record Portal speedrun completed in 8 minutes. <https://www.eurogamer.net/world-record-portal-speedrun-completed-in-8-minutes>, 2012. Accessed 6 July, 2022.
- 18 Kristin Siu, Eric Butler, and Alexander Zook. A programming model for boss encounters in 2d action games. In *Proceedings of the AIIDE workshop on Experimental AI in Games*, 2016.
- 19 Adam M. Smith, Chris Lewis, Kenneth Hullet, Gillian Smith, and Anne Sullivan. An inclusive taxonomy of player modeling. Technical Report UCSC-SOE-11-13, 2011, University California Santa Cruz, 2011.
- 20 Tommy Thompson. How Forza’s Drivatar actually works. <https://www.gamedeveloper.com/design/how-forza-s-drivatar-actually-works>, 2021. Accessed 6 July, 2022.
- 21 Tommy Thompson and Matthew Syrett. “play your own way”: Adapting a procedural framework for accessibility. In *Proceedings of the FDG Workshop on Procedural Content Generation*, 2018.
- 22 Christopher W. Totten. *An Architectural Approach to level Design*, chapter Teaching in Levels through Visual Communication. CRC Press, 2014.
- 23 Rodrigo Vicencio-Moreira, Regan L Mandryk, and Carl Gutwin. Balancing multiplayer first-person shooter games using aiming assistance. In *2014 IEEE Games Media Entertainment*, pages 1–8. IEEE, 2014.
- 24 Markus Weinmann, Christoph Schneider, and Jan vom Brocke. Digital nudging. *Business & Information Systems Engineering*, 58:433–436, 2016.
- 25 Georgios Yannakakis and Julian Togelius. Experience-driven procedural content generation. *IEEE Transactions on Affective Computing*, 2(3):147–161, 2011.
- 26 Georgios N. Yannakakis, Pieter Spronck, Daniele Loiacono, and Elisabeth André. Player modeling. In Simon M. Lucas, Michael Mateas, Mike Preuss, Pieter Spronck, and Julian Togelius, editors, *Artificial and Computational Intelligence in Games*, volume 6 of *Dagstuhl Follow-Ups*, pages 45–59. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2013.
- 27 Georgios N. Yannakakis and Julian Togelius. *Artificial Intelligence and Games*. Springer, 2018. <http://gameaibook.org>.

### 3.9 The Tabletop Board Games AI Tutor

*Diego Perez Liebana (Queen Mary University of London, GB), Duygu Cakmak (Creative Assembly – Horsham, GB), Setareh Maghsudi (Universität Tübingen, DE), Pieter Spronck (Tilburg University, NL), and Tommy Thompson (AI and Games – London, GB)*

License © Creative Commons BY 4.0 International license

© Diego Perez Liebana, Duygu Cakmak, Setareh Maghsudi, Pieter Spronck, and Tommy Thompson

#### 3.9.1 Introduction and Motivation

The aim of game tutorials is to teach a player how to play a game. The design and implementation of tutorials is no easy task, as they will be used by different types of players, with different experiences and prior knowledge of games. It’s also hard to make them adaptive and interactive, especially in a way that identifies the actual capabilities and needs of the player as the learning progresses. In recent years, video-game tutorials have been object of study by game AI researchers who approach this problem from multiple angles. Recently, L. Poretski et al. [7] analyzed tutorials from 40 contemporary video games to identify common patterns and strategies used to generate these elements. Previously, E. Andersen et al. [1] studied the effect of tutorials on player engagement and game retention.

Later on, B. Aytemiz [2] investigated the design of tutorials that considered the skill of the player. Here, the authors implemented a Unity plug-in that provided information to the player only when this is needed. The automatic generation of tutorials has also been recently explored by M. Green and colleagues [5, 4] in AtDELFI, a system that procedurally generates instructions to teach players how to play 2D arcade games within the GVGAI framework [6].

While some research has been put into video-game tutorials, to our knowledge, no work has been carried out for table-top board game tutorials. Table top games (TTG) are normally played on a physical surface and involve the use and manipulation of certain elements, such as tokens, cards, dice or counters. TTG can feature relatively complex rules which may be hard to explain and understand, once the relationships and connections between the different game components and how are they used become intricate. Novice TTG players may experience a steep learning curve when facing some complex board games for the first time. Even some expert players may spend a big proportion of their playing session learning the game rules – even if some of the game mechanics are not new to them. Not surprisingly, many TTG players prefer to learn game rules via video tutorials or live explanations given by a colleague, rather than reading long rule-books. An interesting difference with video games is that TTG can't prevent a player *programmatically* to take invalid actions, sum victory points incorrectly or move your tokens not respecting the rules. Video games can enforce this by design and implementation (an avatar is a subjects to the game's physics model), so the tutorials teach the player to play the game *well*. TTG tutorials, however, must ensure that a player understands and applies rules in a context different to that of a digital tutorial. Current TTG frameworks, such as TAG [3], blur the line between TTG and their *digital twins*, but the game's user interface needs to be implemented with this aspect in mind.

Therefore, we propose research on tutorial generation for TTG as an area worth exploring: it would not only help players learn how to play games more effectively, but it will also provide interesting insights in decision making, game interfaces and game analysis. During this workgroup, we discussed the different challenges and opportunities offered by the generation of AI tutors. This chapter summarizes the discussions and considerations that arose from this group during the seminar.

### 3.9.2 Learning to Play

Once we ask ourselves the question “*how can an AI teach a human to play a Tabletop Board Game?*”, it's important to understand the different dimensions of the problem. One of them is the actual player: what does it mean for a player to know how to play a game? Is it enough that the player understands the rules so they can play the game? Should they be able to explain it to others, or is it sufficient that they make no mistakes when playing? Do players need to know just the valid actions and dynamics of the game, or do they need to have certain skill level to be able to devise good strategies and avoid flawed moves? Do we take into consideration that the player may know other similar games, or their general experience level with TTG?

Additionally, from the point of view of the AI, we may also ask what does the AI need to teach, how is this information conveyed, and when does the AI intervene (if at all) to explain a rule or to correct the player. What sort of capabilities does the AI tutor need to have, and does it build a model of the human to drive the teaching process? Can it analyze the player's gameplay traces to identify crucial decisions and does it take into account any external knowledge – for instance, the game rule-book, or other games that are known to the player? The following categorization analyzes the identified dimensions of this problem, which often overlap each other:

- **Introducing the Rules:** the main objective of a tutorial system is to explain the rules of the game. An important point regarding teaching these rules is *when* are they introduced. One option for the rules to be provided is to enumerate them all at front. This is the simplest scenario, which may suggest shouldn't probably deserve much consideration. However, some TTG *are* simple; if the number and/or the complexity of the rules is small, this approach may actually be the most appropriate. Another option is to provide, at the start, only information about the necessary rules, to then progressively explain new rules when they are needed. A final option would entail not providing any rules beforehand, letting the player *blindly* play the game. An explanation could then appear when a rule is broken, and subtle notifications could reinforce the player when a new rule has been successfully applied. This would have the added benefit of correcting the player mistakes and reinforcing their successes.
- **Contextualizing the Rules:** rules can be provided with different types of context. The simplest option would be, naturally, providing *no* context. For instance, in *Terraforming Mars* (TM; FryxGames, 2016), some cards can only be played if certain pre-requisites have been fulfilled. A rule that explains how this works with no context would just indicate the part in the card to look at when these pre-requisites exist. However, rules can also be explained with a *thematic* context: the card “Anti-gravity Technology”, in TM, can only be played if the player has previously played 7 cards with a Science tag. Context can be given in the form of the theme of the game (many scientific projects are needed before Anti-gravity Technology can be discovered) and even hinting a *strategic* context (the game can be played with a Science strategy in mind to allow playing powerful cards later on). Finally, a useful context to certain rules is providing *distribution* information. In this same example, it is useful to indicate that only one card in the deck has such a restrictive requirement, but a given percentage of cards have requirements of having played different Science tags earlier. While the former examples concentrate in the rules per se, the latter are more useful to convey stronger play tactics.  
As hinted earlier, different approaches may be used for different players, or even at different times in the learning process. In complex games, the role of the AI tutor may help determine which rules should be explained when, or which sort of feedback or context should be given to the player. This is something that could be *learned* from data, for instance from previous tutoring sessions or from similar games.
- **Rules as Rule-sets:** some TTG, such as *Gloomhaven* (Cephalofair Games, 2017) or *Mage Knight* (WizKids, 2011), have a large collection of rules. Teaching or learning all these rules at once is challenging, thus in some cases learning is conveyed through different rule-sets. Initially, a version of the game (often simplified) is played using a simple subset of rules, easier to learn. More rules are added to the rule-set for consecutive game sessions, increasing the complexity of the game but providing a smoother learning curve than using the complete rule-set at first<sup>2</sup>. Note that this is different to progressively introducing rules (mentioned in a previous point), as this allows a complete game being played. An interesting line of research could be to investigate if an AI tutor can automatically find and compose the rule-sets and scenarios required for this incremental teaching approach.
- **Level of Play:** an AI tutor may be configured to teach the game at different proficiency levels. The simplest one is to only explain the rules: what's possible and what's not,

---

<sup>2</sup> *Gloomhaven: Jaws of the Lion* (Cephalofair Games, 2020) used this system to progressively explain the rule-set through 5 consecutive game scenarios.



how is the game played and how the winner is decided. An intermediate approach would require the AI to provide advice and reasoning behind certain decisions at specific points during the game; for instance, that is important to play a certain card because it gives the player further maneuvering in subsequent turns. Finally, an approach that would require a higher skill is to teach the player good strategies, which inform the play-through from start to the end. In this case, the ability of an AI tutor to teach at one or another level can be closely related to its own ability to play the game as an AI player, thus the research into this area could benefit from previous works on AI for game playing.

- **Scenarios:** Teaching can be achieved through complete playthroughs, from game setup to game end, or via different scenarios. In game AI terms, the AI tutor could choose interesting mid-point game states for the player to play in. This is similar to puzzles in Chess, where the player needs to identify the correct move to get to a chess mate or escape from one. An AI capable of automatically generating interesting scenarios would be able to produce valuable *teachable moments*, where expected actions or common mistakes can be identified and corrected.
- **AI Interactions:** When the AI tutor monitors the actions and progress of the learner who is playing the game, a consideration needs to be made about how does the tutor interact with the player when they infringe the rules or take a particularly good or bad action. For the former, it is sensible to think that the AI tutor should intervene immediately to stop a rule for being broken, probably providing extra information about the rule itself and how it was violated. For the latter, however, different approaches could be considered. We could, again, interrupt the game when a mistake is made, or to provide reinforcement for a good move being executed. This could, however, turn the playing session tedious if multiple mistakes are made, so an alternative would be to provide a retrospective summary of good and bad decisions. The AI tutor would recap the “highlights” where the player made a particularly interesting choice (either good or bad) and provide extra information on the quality of the action. Nevertheless, it is important to not lose sight of the player’s psychology at this point – a recapitulation of every time a novice player makes a mistake may not be received kindly by some learners.
- **Learning a Human-model:** Following up from the previous point, it is interesting to discuss if the AI tutor should build a model of the human learner. An important aspect to consider could be the prior knowledge of the player about similar and dissimilar games. This could be useful to draw similarities to dynamics found in other games<sup>3</sup>. Additionally, the AI tutor may also require to model the learning process of the player, for instance to identify which rules have already been learned (therefore no more emphasis should be put on them) and which ones have not. An interesting consequence of this is to estimate the learning progress of the player through the session, to change strategy in case the learning is too slow (provide simpler or fewer rules at a time) or too fast (increase the cognitive *load* required so the full game can be learned earlier).
- **AI Priors:** The previous point identifies prior knowledge from the point of view of the player, but the AI system may also have a source of prior knowledge to use. This can come with regards to the game being explained, either as an expert system that provides a series of rules, information taken from the manual, or outsourced from online resources such as forums or Fandom Wikis. If the AI tutor has been used previously to teach

---

<sup>3</sup> For instance, saying that “the phases in *Terraforming Mars: Expedition Ares* (FryxGames, 2021) are chosen using a similar mechanic to the ones in *Roll for the Galaxy* (Rio Grande Games, 2014)” can be a useful teaching tactic.

the same game, information about previous teaching sessions can be used as data to signal common mistakes or interesting teaching moments identified in earlier attempts. Finally, one could also attempt to extract useful data from other similar games that share common characteristics: for example, the use of *workers* in EuroGames, which is a similar feature in games like *Everdell* (Starling Games, 2018), *Village* (Eggertspiele, 2011), *Istanbul* (Pegasus Spiele, 2014), and many others.

- **Game-play Traces:** An interesting possibility that can aid teaching is extracting information from old games, for instance through the analysis of game-play traces. These can be taken from proficient players, AI agents or the actual human learner. This analysis can be done both in the short-term (actions taken in particular game states with an immediate effect) or in the long-term (actions taken in a game state that have a noticeable effect several turns later). Especially interesting decisions would provide valuable teaching moments for the learner, such as identifying particularly strong moves or missed opportunities where a better action could have been taken. By means of self-play, these missed opportunities can be more illustrative by simulating alternative future states that could have been reached if the better action had been chosen.

In summary, we observe a great possibility for Game AI research in this domain. Stemming from the previous enumeration, we can identify several areas of research that can benefit from work in this area, such as question-answering, generation of interesting situations for teaching and training (game states), identification of teachable moments (game states and actions), detection of weak and strong moves, and the design of interfaces for the communication between the two parts. All this should be investigated in conjunction with building models of the human learner, while adapting the teaching mechanisms to the prior knowledge, progress rate and habits of the player. Research in this area would also shed some light in the capabilities of the current state-of-the-art AI methods, especially in what refers to the Human-AI interaction, and whether it is necessary to design different AI tutor systems for distinct types of games and players.

## References

- 1 Erik Andersen, Eleanor O’rourke, Yun-En Liu, Rich Snider, Jeff Lowdermilk, David Truong, Seth Cooper, and Zoran Popovic. *The impact of tutorials on games of varying complexity*. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pages 59–68, 2012.
- 2 Batu Aytemiz, Isaac Karth, Jesse Harder, Adam M Smith, and Jim Whitehead. *Talin: a framework for dynamic tutorials based on the skill atoms theory*. In Fourteenth Artificial Intelligence and Interactive Digital Entertainment Conference, 2018.
- 3 Raluca D. Gaina, Martin Balla, Alexander Dockhorn, Raul Montoliu, and Diego Perez-Liebana. *TAG: A Tabletop Games Framework*. In Experimental AI in Games (EXAG), AIIDE 2020 Workshop, 2020.
- 4 Michael Cerny Green, Ahmed Khalifa, Gabriella AB Barros, Tiago Machado, Andy Nealen, and Julian Togelius. *Atdelfi: automatically designing legible, full instructions for games*. In Proceedings of the 13th International Conference on the Foundations of Digital Games, pages 1–10, 2018.
- 5 Michael Cerny Green, Ahmed Khalifa, Gabriella AB Barros, and Julian Togelius. “Press space to fire”: *Automatic video game tutorial generation*. In Thirteenth Artificial Intelligence and Interactive Digital Entertainment Conference, 2017.

- 6 Diego Perez-Liebana, Jialin Liu, Ahmed Khalifa, Raluca D Gaina, Julian Togelius, and Simon M Lucas. *General video game ai: A multitrack framework for evaluating agents, games, and content generation algorithms*. IEEE Transactions on Games, 11(3):195–214, 2019.
- 7 Lev Poretzki and Anthony Tang. *Press a to jump: Design strategies for video game learnability*. In CHI Conference on Human Factors in Computing Systems, pages 1–26, 2022.

### 3.10 Artificial Intelligence for Alternative Controllers

*Lisa Rombout (Tilburg University, NL), Alex J. Champandard (creative.ai – Wien, AT), Ahmed Khalifa (University of Malta – Msida, MT), Paris Mavromoustakos Blom (Tilburg University, NL), and Mark J. Nelson (American University – Washington, US)*

**License** © Creative Commons BY 4.0 International license

© Lisa Rombout, Alex J. Champandard, Ahmed Khalifa, Paris Mavromoustakos Blom, and Mark J. Nelson

Most popular game controllers are relatively straightforward – a collection of buttons, perhaps some touchpads and scroll wheels, a joystick. These inputs give clear, unambiguous signals to the system, and are generally designed to be as unobtrusive as possible. The user should be able to interact so naturally with the controller that they are almost able to forget it exists.

However, controllers can also be an integral part of the gameplay experience, facilitating play that is a combination of physical and digital. Popular examples include the guitar from Guitar Hero and the dance pad from Dance Dance Revolution. Currently widely available hardware platforms, such as Arduino, have made it very easy and accessible to extend a digital game with real-world sensors and actuators of various kinds. As basically anything conductive can be made into a sensor (including bowls of custard, bananas, and people), it follows that anything conductive can be made into a controller. The challenge, then, is to make sense of the data coming in to the system, and make sure the interaction can run smoothly.

We explored the idea of using an Arduino device to a) monitor a player’s heart rate in real time and b) translate the heart rate measurement into the rhythm (beats per minute) of a custom-made music tune. This concept could be employed in an adaptive stress management game, where the player is rewarded for maintaining their heart rate at a rest-state level.

### 3.11 Quality Diversity for Procedural Content Generation

*Jacob Schrum (Southwestern University – Georgetown, US), Alex J. Champandard (creative.ai – Wien, AT), Guillaume Chanel (University of Geneva, CH), Amy K. Hoover (New Jersey Institute of Technology, US), Ahmed Khalifa (University of Malta – Msida, MT), Mark J. Nelson (American University – Washington, US), Mike Preuß (Leiden University, NL), and Vanessa Volz (modl.ai – Copenhagen, DK)*

**License** © Creative Commons BY 4.0 International license

© Jacob Schrum, Alex J. Champandard, Guillaume Chanel, Amy K. Hoover, Ahmed Khalifa, Mark J. Nelson, Mike Preuß, and Vanessa Volz

Quality diversity algorithms [8, 1, 5, 6] are well-suited for generating game content [4], because most often there is no single optimal piece of content: level, texture, character design, etc. Rather, games need a variety of content, but all of that content needs to be of a reasonably

high quality, hence the usefulness of QD algorithms. One of the more popular QD algorithms is the Multi-dimensional Archive of Phenotypic Elites (MAP Elites [6]), which collects an organized archive of diverse but high-quality solutions. Many variants of MAP-Elites exist [3, 7], but all rely on some sort of archive with a user-specified structure.

Unfortunately, there are unresolved questions that make variants of MAP Elites difficult to use, even for skilled practitioners. For any given domain, expert knowledge is required to define potentially beneficial behavior characteristic dimensions, which determine the number of dimensions in the archive. It is unclear how many dimensions should be used, and how many intervals should exist along each dimension. Several experiments were proposed to help guide practitioners in dealing with these issues, particularly for games.

### 3.11.1 Proposed Experiments

Here are a list of specific and detailed experiments considered as part of the seminar.

#### 3.11.1.1 Restrict Archive Based on Real Game Data

It is assumed that good games already feature a good variety of content. Although interesting content could be discovered by searching outside the bounds of the existing content, there is also a risk of wasting considerable computational resources discovering a wide variety of uninteresting content. Not only is it computationally expensive to search a larger space, but the cost in human effort to analyze an unnecessarily large archive can be prohibitive. Therefore, we propose a way of restricting an archive to focus on areas likely to be relevant to designers, by using existing game content as a basis. This experiment can readily extend previous work applying MAP Elites to evolve levels for Super Mario Bros. and The Legend of Zelda [9], but can be applied to any game or content.

First, behavior characteristics and archive structures taken from previous literature can be used to store original levels from these games. If at least a majority of the levels occupy distinct bins in the archive, it means that the archive is able to capture the diversity of content in the original game. Intervals along different dimensions of the archive could be adjusted to accommodate more of the original game content if levels that share a bin are deemed to be sufficiently different.

Second, the archive storing the original game content is contracted, so that bins outside the boundaries of what is present in the original game are excluded. Specifically, for each archive dimension, the minimum and maximum values are calculated, and the empty bins outside this range are deemed unreachable.

Third, this restricted archive is used to evolve new content with MAP-Elites or one of its variants. Whenever a level is generated that would fall into one of the unreachable bins, it can simply be discarded, or penalized in some way. This restricts the range of what evolution can produce, but it also means that there are no levels in the unreachable bins that can be chosen as a parent for a new level. Therefore, the search will be more focused on what are presumed to be the most relevant areas of the archive.

Finally, results with such restricted archives would need to be compared with results from unrestricted archives. Although an unrestricted archive will likely contain many more occupied bins, we can ask several pertinent questions: 1) which approach does a better job of filling the restricted/contracted portion of the archive with quality solutions? 2) which approach fills this restricted/contracted area faster? 3) How diverse and interesting are levels from outside of the restricted/contracted area? This last question is perhaps the most difficult to answer due to a lack of widely accepted definitions for *diverse* and *interesting*,

though a human subject study could help in this assessment. The first two questions can be directly answered using objective measures such as the difference in QD score and archive occupancy over time.

The results of such experiments would be informative, regardless of how they turn out. If restricting the archive allows a pertinent area to be filled better or more quickly, then it means that researchers and practitioners can stop wasting effort on large archives. However, if searching a larger archive actually makes it easier to fill a restricted portion of the archive, it would imply that the extra bins are providing useful stepping stones for searching the evolutionary space. It is possible that these outcomes will be dependent on the domain and archive structure being used, but this outcome would also be informative.

#### 3.11.1.2 Evolving Levels With a Wide Range of Archive Dimensions

There are many properties of content that can be objectively measured, and in principle, any such property can be used as a dimension in a behavior characterization. For a content designer, it is unclear what properties to incorporate into a behavior characterization, and how many, but a designer will presumably have access to many candidate properties worth considering.

This experiment proposes a way of assessing how useful different numbers of dimensions in a behavior characterization are. We can start with a game like Super Mario Bros. for which past work [11] has provided a wide range of level properties that can be measured, and which seem to impact the user experience of a level in meaningful ways.

Given some large set of level properties, a way of using each property in a behavior characterization and within an archive is needed. Given this starting point, all combinations of the different properties within a behavior characterization can be considered. Separate experiments can be conducted with each behavior characterization.

The challenge is in comparing these results in order to produce recommendations for designers on which types of properties and how many to include in a behavior characterization. For any given archive of evolved MAP Elites solutions, one can take those archive's members and transfer them to another archive. If the dimensions used in the target archive are a strict subset of those from the source archive, one would expect that some number of solutions are lost as bins along now missing dimensions are collapsed into one bin within the target archive. The degree of this loss can be measured, and the discarded solutions can be analyzed to see if anything of value was actually lost from using a simpler archive.

Also, results evolved with smaller archives can be compared against those transferred to the same archive structure from a larger archive to see if focusing on fewer behavior dimensions leads to higher quality, at least within that particular range.

A more daunting task is comparing archives whose defining behavior characteristics are made up of disjoint sets of properties. One can still transfer solutions between the different archive structures to track how many solutions are lost in the transfer, but it is less clear what to make of such results. If one archive can contain all of the solutions of another, then it implies that some aspect of the source archive might be superfluous, but it is easily possible for significant loss to occur when transferring in both directions, so further analysis will be needed to see what this information can actually teach prospective users of MAP Elites for games.

#### 3.11.1.3 Diversity of Content and of Playing Styles

The behavior characteristics used to evolve diverse game content using QD methods often include measures of playing style, estimated by simulated AI players. For example, we can evolve levels that require a lot of jumping to clear, and little to no jumping to clear (and

everything in between). The result of such evolution is a set of diverse content that supports diverse playing styles – diverse according to whatever measures of play style were chosen for the behavior characterizations.

This approach assumes a 1-to-1 mapping between specific pieces of content and the playing styles they support (at least in the sense of reducing any variation to a single numerical estimate per behavior characterization). Therefore, we evolve a range of content to support a range of different playing styles. However, a range of playing styles can often be supported by a single piece of content, such as a level that is replayable in several very different ways. This leads to a second place we can apply QD methods: to investigate diversity of playing styles supported by fixed pieces of content, and what implications that has for PCG.

To run concrete experiments on diverse playing styles, it is most straightforward to choose a specific parameterized representation for the simulated player. For example, one of the demos of the `pyribs` implementation of CMA-ME [2] evolves agents for Lunar Lander that complete the level with a range of  $x$  positions and  $y$  velocities.<sup>4</sup> That demo uses a linear policy representation. Linear policies are a good place to start, although evolving weights of a fixed neural network is another option.

The result of running QD on Lunar Lander shows that the game supports a wide range of playing styles (types of landing), which in this case is probably key to the game’s popularity. Should this be an explicit target for PCG? In addition to evolving diverse content, should individual pieces of content be generated with a goal that each one also supports diverse playing styles? That is less clear; it may in fact be better in some cases to retain a clearer mapping between individual pieces of content and desired playthroughs. For example, in an educational game it may explicitly *not* be desirable, if the goal of a level is to force the use of a specific concept being introduced [10].

The open question here is: What is the relationship between diverse game content and diversity of play styles supported by a single piece of game content?

#### 3.11.1.4 Conclusion

Although MAP Elites and its variants are powerful tools for creating diverse collections of quality content for games, more research is needed to understand how to apply these tools and how to evaluate the artifacts they produce. In particular, the relation of produced to existing game content is important here. This working group produced several intriguing ideas for future research along these lines.

#### References

- 1 Konstantinos Chatzilygeroudis, Antoine Cully, Vassilis Vassiliades, and Jean-Baptiste Mouret. Quality-diversity optimization: A novel branch of stochastic optimization. In *Black Box Optimization, Machine Learning, and No-Free Lunch Theorems*, pages 109–135. Springer, 2021.
- 2 Matthew C. Fontaine, Scott Lee, Lisa B. Soros, Fernando de Mesentier Silva, Julian Togelius, and Amy K. Hoover. Mapping hearthstone deck spaces through map-elites with sliding boundaries. In Anne Auger and Thomas Stützle, editors, *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO 2019, Prague, Czech Republic, July 13-17, 2019*, pages 161–169. ACM, 2019.
- 3 Matthew C. Fontaine, Julian Togelius, Stefanos Nikolaidis, and Amy K. Hoover. Covariance Matrix Adaptation for the Rapid Illumination of Behavior Space. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, New York, NY, USA, 2020. ACM.


---

<sup>4</sup> [https://docs.pyribs.org/en/stable/tutorials/lunar\\_lander.html](https://docs.pyribs.org/en/stable/tutorials/lunar_lander.html)

- 4 Daniele Gravina, Ahmed Khalifa, Antonios Liapis, Julian Togelius, and Georgios N Yannakakis. Procedural content generation through quality diversity. In *2019 IEEE Conference on Games (CoG)*, pages 1–8. IEEE, 2019.
- 5 Joel Lehman and Kenneth O. Stanley. Evolving a diversity of virtual creatures through novelty search and local competition. In Natalio Krasnogor and Pier Luca Lanzi, editors, *13th Annual Genetic and Evolutionary Computation Conference, GECCO 2011, Proceedings, Dublin, Ireland, July 12-16, 2011*, pages 211–218. ACM, 2011.
- 6 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *CoRR*, abs/1504.04909, 2015.
- 7 Olle Nilsson and Antoine Cully. Policy gradient assisted map-elites. In Francisco Chicano and Krzysztof Krawiec, editors, *GECCO '21: Genetic and Evolutionary Computation Conference, Lille, France, July 10-14, 2021*, pages 866–875. ACM, 2021.
- 8 Justin K Pugh, Lisa B Soros, and Kenneth O Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, page 40, 2016.
- 9 Jacob Schrum, Vanessa Volz, and Sebastian Risi. CPPN2GAN: Combining Compositional Pattern Producing Networks and GANs for Large-scale Pattern Generation. In *Genetic and Evolutionary Computation Conference*, New York, NY, USA, 2020. ACM.
- 10 Adam M. Smith, Eric Butler, and Zoran Popovic. Quantifying over play: Constraining undesirable solutions in puzzle design. In *Proceedings of the Foundations of Digital Games Conference*, pages 221–228, 2013.
- 11 Vanessa Volz. *Uncertainty Handling in Surrogate Assisted Optimisation of Games*. PhD thesis, TU Dortmund University, Germany, 2019.

### 3.12 Benchmarking Coordination Games

Pieter Spronck (Tilburg University, NL), Duygu Cakmak (Creative Assembly – Horsham, GB), Jakob Foerster (University of Oxford, GB), and Setareh Maghsudi (Universität Tübingen, DE)

License  Creative Commons BY 4.0 International license  
© Pieter Spronck, Duygu Cakmak, Jakob Foerster, and Setareh Maghsudi

Games are often used in AI research as benchmarks for new technologies and developments. For instance, a breakthrough application of deep convolutional neural networks and Monte Carlo tree search was in the game of Go, where the program AlphaGo was able to defeat the human world champion by using the aforementioned techniques [3]. Competitions between different AI approaches to play particular games are common, not only for classic deterministic two-player board games such as Chess, but also for modern board games and video games.

Commonly, the games used are competitive games. One reason why this is the case, is that for competitive games it is relatively easy to determine which AI is “the best,” namely the one that wins the most or scores the highest. However, competition only covers part of what AI needs to do. In the modern world, where AI get increasingly integrated in people’s lives, the ability of AI to cooperate effectively (with humans or with other AIs) is of great interest to researchers.

If the environment in which cooperation is required is fully observable, the solution to problems of cooperation is, in general, to let the smartest AI propose the most effective approach, and let every participant follow that approach. However, in practice cooperation

problems are often only partially observable, either because the participants have only partial information of the environment or because participants have individual goals.

To research AI which is able to cooperate, we therefore need partially observable, fully cooperative benchmarks. Games can be such benchmarks.

### 3.12.1 Examples of coordination games

A typical example of a coordination game is Hanabi, a game in which players cooperatively need to achieve a goal by playing cards, whereby they have knowledge of the cards of the other players but not of their own cards, while their ability to communicate is limited to very specific statements that they can make [1]. If AIs learn to play Hanabi with only other AIs in the mix, they will discover strong strategies to play the game, reaching a very high cooperative score. However, such strategies are not used by humans, so different approaches must be used to create AIs which are able to cooperate with human players.

Hanabi is a partially observable, turn-based game for 3 to 5 players, who all have a different observation space. In the previous iteration of this series of Dagstuhl Seminars, the “Guess 100” game was developed, which is a fully observable, simultaneous-move game for 2 players. In the Guess 100 game, the two players simultaneously choose a number between 1 and 100; the two numbers are then revealed to the two players, after which they choose a new number. There is no other communication between the players. They continue doing this until they chose the same number, with the ultimate goal to keep the number of turns needed as low as possible.

Another example of a well-known, popular coordination game is Codenames, where a turn for one team with a team leader can be considered a coordination task which is partially observable for the team and fully observable for the team leader, whereby the team leader attempts to give a one-word hint which directs the team to select from a grid of cards those cards which belong to the team (without the team being able to see which cards are theirs).

During the pandemic, a digital version of the board game Wavelength became highly popular. In this game one player, the Psychic, gets two words which are extremes on a range, such as “hot” and “cold.” The Psychic also gets a “bullseye”, a spot on the line between the extremes. The Psychic needs to give a clue (with certain limitations to what the clue can be) which indicates where the bullseye is. The players then indicate a spot on the line between the two extremes, and the closer they are to the bullseye, the more points they score. Like Codenames, this is a partially observable, turn-based game.

### 3.12.2 Axes of coordination

The goal of the workgroup was to determine different coordination games, preferably games that are simple to implement, play, and understand, which cover different coordination problems. We therefore started by determining the different “axes of coordination.” We arrived at the following list:

- Fully observable vs. partially observable
- Simultaneous moves vs. sequential moves
- Labeled vs. unlabeled states and actions
- Type of labels (e.g., ordered, unordered, semantic)
- Turn-based vs. real-time
- Iterated vs. single-shot
- Shared vs. different observation spaces
- 2-player vs. 3+ players



### 3.12.3 Variations of the Guess 100 game

As the Guess 100 game is close to the simplest implementation of a fully-observable, simultaneous-move, 2-player game with ordered labels for actions, it makes sense to use it as a template for game variations which touch different axes of coordination.

To turn Guess 100 into a partially observable game (leaving the other axes of coordination unchanged), the two players do not aim to find the same number, but to land on a particular target number that is different for each of them, whereby the target number for each player is known to the other player, while unknown to themselves. The players are not allowed to select the target number of the other player. The game ends when they simultaneously land on their assigned target number.

To turn Guess 100 into a game without ordering, the players would not select numbers but icons, and the icons are ordered differently for the two players, so that even the location on the grid cannot be considered an ordering.

To turn Guess 100 into a game without labels or ordering at all, the playing field can consist of a number of moving balls which are unmarked. The players' goal is still to select the same ball, and they get to see which ball was selected by the other player when both have selected a ball.

### 3.12.4 The Color-coder game

Since an integral part of the Guess 100 game is that the players move simultaneously, we continued by designing a game that is partially observable, has sequential moves, and no ordered labels. The aforementioned Codenames is such a game, but too complex to be a good basis for fundamental research. We wanted this game to be as simple as possible. We came up with the Color-coder game.

In the Color-coder game there are two players, one of which has the role of “hinter” and the other one the role of “guesser.” The hinter observes a “code” which consists of two colored tokens, randomly selected from four colors, e.g., RED-BLUE. The hinter and guesser now take turns, starting with the hinter. The turn of the hinter consists of providing a sequence of black and white tokens, e.g., BLACK-BLACK-WHITE-WHITE-BLACK. The turn of the guesser consists of producing a guess, which consists of two colored tokens, e.g., YELLOW-GREEN. No other communication between the hinter and guesser is allowed. The game ends when the guesser reproduces the code that the hinter observed. The goal is to minimize the number of turns needed for that.

The Color-coder game sounds like Mastermind, but it is substantially different. In Mastermind the meaning of the black and white tokens is predetermined, and the hinter's role is not to cooperate with the guesser, but to simply provide the predetermined hint. In the Color-coder game, the hinter actively wants to lead the guesser to the correct code, and can try to supply the guesser with more information. The hinter can decide upon a particular interpretation of the tokens and stick to it, but might also change the meaning of the tokens while playing.

We tested the Color-coder game several times, and found that different strategies were employed. Often the hinter chose a strategy before the game started, and did not diverge from it, hoping that the guesser would pick it up. Sometimes the hinter changed their way of hinting during the game. In principle there is no reason for the game to last longer than four turns, if the hinter simply always uses black to indicate correct and white to indicate incorrect, and then always places two tokens. However, with smart hinting fewer turns are needed.

One hinker came up with the idea to hint at the correct code by using black tokens for the left color and white tokens for the right color, and then placed as many of the respective tokens as there are letters in the color name. Considering that in our default game we used red, blue, green, and yellow, the colors were encoded in a unique way. If the guesser would pick up on the approach, one guess would suffice.

The Color-coder game felt as slightly too simple to give rise to interesting communication strategies. It can be made more interesting without increasing complexity too much by increasing the number of colors.

### 3.12.5 The Convergence and Divergence games

We also wanted to define a game that, like the Guess 100 game, is fully observable and has simultaneous moves, but has no numerically ordered labels, because with a numerical ordering players will use calculations to try to reach the same number.

In the Convergence game the players are presented with a list of ten words. The words are randomly selected from a dictionary, e.g., zoom, understood, women, income, joke, scrawny, waiting, bucket, picayune, camera. The list may be presented to the players in a different order, so that the place on the list is not part of the ordering; there is, of course, a semantic ordering. The players simultaneously select a word. If the word is the same, the game ends. If not, they select a different word, whereby the words that were selected before cannot be selected again. The goal is to select the same word as fast as possible.

The players can use the semantic ordering of the words to decide which next word to select. E.g., if from the previous list one player selected “understood” and the other one “women”, then they might decide to select “waiting” as their next guess as that is the only word which is between the two previously selected ones. However, with numbers such an approach is far more “natural” than it is with words, and players are more likely to use the “meaning” of words to direct their guesses.

However, we found a slight change to the Convergence game made it much more interesting: in the Divergence game the players undertake the same actions, but they try to avoid selecting the word that the other player selects. As soon as they land on the same word, the game ends, and they try to postpone this as long as possible. With ten words, the game can last no more than five turns, but the game can easily be extended by making the list of words longer, which would also allow the players to try to communicate more as they can select more words before the list becomes so short that the risk of selecting the same word is high.

### 3.12.6 Next steps

We implemented a digital version of the Guess 100 game during and after the previous Dagstuhl Seminar. We want to use it to collect data on game plays. It can also be used to implement variations of the Guess 100 game, to work on different axes of coordination. The datasets developed this way should be open-sourced. We then want to develop AIs which play these games, both with other AIs and with human players. The most interesting games can form the basis for a challenge paper.

#### References

- 1 N. Bard, J.N. Foerster, S. Chandar, et al. *The Hanabi challenge: A new frontier for AI research*. *Artificial Intelligence* 280, 2020
- 2 K. Hofmann, D. Cakmak, P. Cowling, et al. Human-AI Coordination. *Artificial and Computational Intelligence in Games: Revolutions in Computational Game AI*, Dagstuhl Seminar 19511, 2020
- 3 D. Silver, A. Huang, C. Maddison, et al. *Mastering the game of Go with deep neural networks and tree search*. *Nature* 529, 484–489, 2016. <https://doi.org/10.1038>

### 3.13 Explainable AI for Games

*Jichen Zhu (IT University of Copenhagen, DK), Maren Awiszus (Leibniz Universität Hannover, DE), Michael Cook (Queen Mary University of London, GB), Alexander Dockhorn (Leibniz Universität Hannover, DE), Manuel Eberhardinger (Hochschule der Medien – Stuttgart, DE), Daniele Loiacono (Polytechnic University of Milan, IT), Simon M. Lucas (Queen Mary University of London, GB), Ana Matran-Fernandez (University of Essex – Colchester, GB), Diego Perez Liebana (Queen Mary University of London, GB), Tommy Thompson (AI and Games – London, GB), and Remco Veltkamp (Utrecht University, NL)*

**License** © Creative Commons BY 4.0 International license

© Jichen Zhu, Maren Awiszus, Michael Cook, Alexander Dockhorn, Manuel Eberhardinger, Daniele Loiacono, Simon M. Lucas, Ana Matran-Fernandez, Diego Perez Liebana, Tommy Thompson, and Remco Veltkamp

In recent years more and more research has been invested into eXplainable artificial intelligence (XAI) to make machine learning (ML) and AI models more trustworthy and understandable for users. In an earlier vision paper, a new research area for designers and game designers was proposed called XAI for Designers (XAID) [1], which focused on *mixed-initiative co-creation* [2] approaches to help designers better leverage AI methods through co-creation in their respective design tasks. Since then, much development has been made in XAI. In this working group, we investigate whether and how these new methods for XAI can also be used for games.

#### 3.13.1 What is XAI for Games?

There are a large variety of possible use cases for XAI in games or game development, and this largely depends on what one wants to achieve. Some salient use cases include:

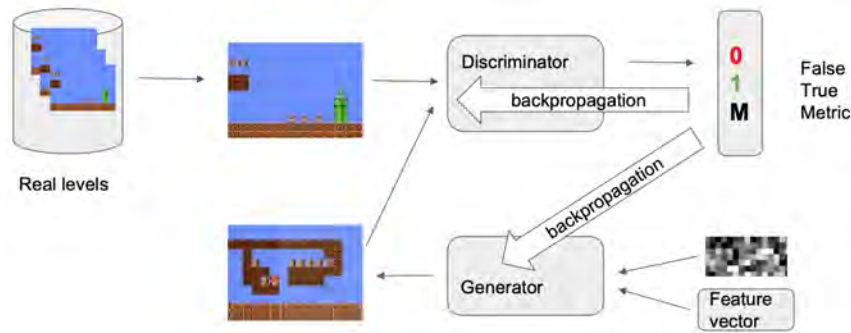
- Increasing the transparency for game AI decisions so that these decisions can be understood and trusted by humans.
- Explanations of key game AI decisions can be used as a feedback mechanism for how well a player is performing. For instance, a PCG-based educational game can explain to a player that a new level is generated based on her previous gameplay so that she can continue to practice a certain skill that she has not mastered. This type of explanation can be used as a feedback mechanism to foster player reflection and learning[6].
- Tools for framing in computational creativity and improving the design experience with mixed-initiative co-creativity systems.
- Highlighting to players why a given strategy is relevant, optimal, or exciting.

To further narrow the focus on the different use cases, in this report, we will focus on procedural content generation with ML (PCGML).

#### 3.13.2 Case Study: Mario Level Generation

First, we looked at different possibilities to generate Super Mario levels. TOAD-GAN [3] can be trained using only one example. This method also makes it possible for users to control the output of the generation process by changing the noise vector that represents the input of the generator network. Since noise vectors cannot be interpreted by designers, designers still do not have the ability to design content according to their needs. To accomplish this, one must make the noise vector explainable to designers and map the different areas of the noise vector to the content that would result from a change in the noise vector.

Another method for generating Super Mario levels uses an evolutionary algorithm with tilesets [4]. The tilesets enforce consistency of the output, and the Kullback-Leiber Divergence



■ **Figure 9** An overview of DOOMGAN architecture.

enables for control of variation and novelty. This method is explainable by design as the history of the gene values and the time steps of the mutation operators could be used to identify *when* something occurred, and *why* it was picked to be modified.

### 3.13.3 Case Study: DOOMGAN – Improving the PCGML Interpretability by Incorporating Metrics

PCGML[5] has been successfully applied to several kinds of game content. However, it generally has low interpretability to human designers because how the input (e.g., parameter/feature vectors) leads to generated content (or corresponding gameplay metrics) is often opaque. Recent Deep Learning-based generative models exacerbate this problem due to their complexity and blackbox nature. As a relevant case of study, we focused on GAN-based PCGML approaches and proposed to incorporate gameplay metrics (e.g., completion time, win rate) in part of the GAN architecture at the level of the discriminator. Figure 9 provides an overview of the proposed GAN architecture, dubbed DOOMGAN, where the discriminator is extended by adding one or more gameplay metrics as additional outputs. Our research hypotheses are that this method will 1) improve the interpretability of the system by providing meaningful intermediate output to designers and 2) improve the performance of the generative model (e.g., better data quality and data efficiency). Moreover, with the proposed method, existing XAI techniques, such as Saliency Map, LIME, and DeepSHAP, can be used to further open the blackbox of PCGML. An ideal testbed to investigate our ideas would be to extend one of the Mario level generators based on GAN previously introduced in the literature (e.g., TOAD-GAN[3]). To the best of our knowledge, this is among the first approach that connects XAI to PCGML methods.

### 3.13.4 Open Problems

Explainable AI for games is still a nascent research area. Below we summarize some of the key open problems in this area:

- How to turn explainability into explanation and actionable explanations to players and/or designers?
- How does content representation affect explainability? (e.g., representing a Mario level as tiles vs. objects)
- Whom do we design the explainable system for? What do the human players, designers, or other stakeholders need? Current XAI methods only explain predictive models but not generative models.
- How to capture functionality/playability of a level in XAI, which is absent in image generation?

### 3.13.5 Conclusion

In summary, this working group found eXplainable AI to be a rich research topic to explore in the context of computer games. Making the underlying AI process more transparent can benefit a wide range of stakeholders, including players, game designers, game analytics/user researchers, and game producers. Since computer games are end user-facing, we believe exploring eXplainable AI in the context of games will expedite the transition from technical explainability to usable human-centered explanations.

#### References

- 1 Zhu, J., Liapis, A., Risi, S., Bidarra, R. & Youngblood, G. Explainable AI for Designers: A Human-Centered Perspective on Mixed-Initiative Co-Creation. *2018 IEEE Conference On Computational Intelligence And Games (CIG)*. pp. 1-8 (2018)
- 2 Yannakakis, G., Liapis, A. & Alexopoulos, C. Mixed-initiative co-creativity. *FDG*. (2014)
- 3 Awiszus, M., Schubert, F. & Rosenhahn, B. TOAD-GAN: Coherent Style Level Generation from a Single Example. *AAAI Conference On Artificial Intelligence And Interactive Digital Entertainment Best Student Paper Award*. (2020,10)
- 4 Lucas, S. & Volz, V. Tile Pattern KL-Divergence for Analysing and Evolving Game Levels. *Proceedings Of The Genetic And Evolutionary Computation Conference*. pp. 170-178 (2019), <https://doi.org/10.1145/3321707.3321781>
- 5 Summerville, A., Snodgrass, S., Guzdial, M., Holmgård, C., Hoover, A., Isaksen, A., Nealen, A. & Togelius, J. Procedural Content Generation via Machine Learning (PCGML). *IEEE Transactions On Games*. **10**, 257-270 (2018)
- 6 Zhu, J. & El-Nasr, M. Open player modeling: Empowering players through data transparency. *ArXiv Preprint ArXiv:2110.05810*. (2021)

## 3.14 Human-AI Collaboration Through Play

*Jichen Zhu (IT University of Copenhagen, DK), Guillaume Chanel (University of Geneva, CH), Michael Cook (Queen Mary University of London, GB), Alena Denisova (University of York, GB), Casper Harteveld (Northeastern University – Boston, US), and Mike Preuß (Leiden University, NL)*

License  Creative Commons BY 4.0 International license

© Jichen Zhu, Guillaume Chanel, Michael Cook, Alena Denisova, Casper Harteveld, and Mike Preuß

### 3.14.1 Motivation

Human-AI collaboration is a rapidly growing research area. As AI becomes an integral part of the workplace as well as home, developing technology that can efficiently collaborate with humans is essential.

Existing psychology research found that successful collaborations between humans need the foundation of 1) a relational interaction (conflict, small talk, emotional exchanges, relationship construction) and 2) efficient cognitive interaction (e.g., building on others' ideas – transactivity, synthesis, building a common ground) [1]. However, in current *Human-AI interaction* (HAI) research, this social-cognitive element and the social experience between human users and the AI is under-explored. This is problematic because since most users, especially novel users of AI, tend to approach AI based on their knowledge of similar human interactions [9, 8].

We argue that computer games, and playful interactions around games, provide a platform to explore new forms of collaborations and social interactions that consider relational aspects and are more cognitively nuanced. Such AI agents can potentially lead to better outcomes such as deeper social engagement when a player collaborates with it towards a common goal (e.g., solving a puzzle or mixed-initiative co-creative game design [10]).

Many games are social by nature as they propose several mechanisms for collaborative and competitive play. In addition, it appears that watching a game together is already sufficient to provide a significant social experience as demonstrated by streamed games [11], probably being not that different from groups of people watching sports games. This implies that there is huge potential to explore human-AI collaboration through play.

### 3.14.2 Background on Collaborative Game AI

A particularly interesting collaborative AI player is OpenAI Five [7]. It is able to play the MOBA game Dota 2 at the human professional level. Teams in this game consist of 5 players, and close collaboration is mandatory for winning the game. The AI has learned a specific type of collaborative behavior that may be called generous self-sacrifice. It is able to play extremely well when playing as a pure AI team but did not manage to collaborate well with human players, most likely because humans play more selfishly.

In addition, [12] used intention recognition to infer the task the human player was performing at the moment so that the AI could provide the appropriate assistance to the player. In board games, [13] explored agents for games built on collaborative game mechanics between human players. The authors used Rolling Horizon Evolutionary Algorithm to develop an artificial agent to balance gameplay in *Pandemic*. [14] explored how procedural content generation, such as character generation, story sifting, and social simulation, can be used to facilitate collaborative storytelling among human players.

Finally, HCI researchers have used computer games as a platform to study how human perception and gameplay outcomes may be affected when interacting with an AI [15].

### 3.14.3 Design space of human-AI interactions

Social interaction can take several forms [2, 3], including:

- Competition: the various agents are competing on a limited amount of resources to accomplish goals that are generally orthogonal.
- Collaboration: collaboration is defined as the action of working together on a single shared goal which generally leads to joint and strongly coordinated behaviors.
- Cooperation: during cooperation, some goals might be distinct, and the agents generally dispatch sub-tasks among the group to only assemble the results at the end of the task.
- Mediation: during mediation, a single agent has the goal of reducing the amount of conflict between at least two other agents.

Researchers have found that different AI techniques support different collaborative interaction mechanisms. For example, [16] noticed that generative adversarial networks (GANs) require a modified set of interaction patterns compared to other generative models.

In the context of video games, most artificial agents are designed to be competitive, while some are able to collaborate with other artificial agents. However, AIs that are able to collaborate or cooperate with human players are more scarce. Cooperation would necessitate understanding which goals are shared, dividing the objectives into sub-tasks, and reaching a common agreement on the distribution of those tasks. Collaboration is more complex as there is a need to synchronize actions between the AI and the players at any time. This is achievable

only by having a common understanding of the situation and reaching an agreement on which steps to take next. Finally, a very under-studied type of social interaction in games is mediation. Artificial agents could help humans to collaborate better by assisting them in organizing and avoiding conflicts. In all these interaction types, there is a need to measure the interaction to determine the best approach and actions to take next.

#### 3.14.4 Measuring social interactions

Social interactions can be measured in real-time or a posteriori – after the actual interaction has taken place. The former might be useful to provide feedback to AI agents so they can adapt their behavior to the social situation, for instance, by being included as a reward in reinforcement learning. The latter would allow for evaluating the efficiency and outcomes of the proposed approach.

Questionnaires can be used to evaluate the interaction a posteriori but are more difficult to be administered during gameplay. The importance of players' social experience was acknowledged by the Game Experience Questionnaire [4], which includes a social presence module mostly inspired by the social presence theory. The competitive and cooperative presence in gaming questionnaire [5] allows researchers to capture the sense of social presence in different types of interaction.

Measuring a social interaction can also be achieved by measuring in-game (position, firing rate) and reactions outside of the game (for example, players' facial expressions, eye movements, and physiological signals). The advantage of this method is that it can provide insights both during and after the game. Several publications, including [6], have studied the possibility of using joint reactions to identify the type of interaction and collaborative processes. However, the use of in-game features to characterize game social interaction remains an under-studied research area. In addition, there is a need to investigate if these methods, including the usage of questionnaires, can be transferred from the context of human-human interactions to human-AI interactions.

#### 3.14.5 Conclusions and Future Work

In conclusion, computer games and playful interactions are a particularly rich domain to explore and advance human-AI collaboration research. This includes both the technical research of how to build better collaborative AIs and the HCI research of how to design new forms of collaborative interaction between humans and agents.

#### References


- 1 Avry, S., Chanel, G., Bétrancourt, M., & Molinari, G. (2020). Achievement appraisals, emotions and socio-cognitive processes: How they interplay in collaborative problem-solving?. *Computers in Human Behavior*, 107, 106267.
- 2 Arnold, N., Ducate, L. & Kost, C. (2012). Collaboration or cooperation? Analyzing group dynamics and revision processes in wikis. *Calico Journal*, 29(3), 431-448.
- 3 Bogacz, F., Pun, T. & Klimecki, O.M. Improved conflict resolution in romantic couples in mediation compared to negotiation. *Humanities and Social Sciences Communications* 7, 131 (2020)
- 4 Jsselsteijn, W. A., de Kort, Y. A. W. & Poels, K. (2013). The Game Experience Questionnaire. *Technische Universiteit Eindhoven*.
- 5 Hudson, M., & Cairns, P. (2014). Measuring social presence in team-based digital games. *Interacting with Presence: HCI and the Sense of Presence in Computer-mediated Environments*, 83.

- 6 Chanel, G., & Mühl, C. (2015). Connecting brains and bodies: applying physiological computing to support social interaction. *Interacting with Computers*, 27(5), 534-550.
- 7 Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., Józefowicz, R., Gray, S., Olsson, C., Pachocki, J., Petrov, M., Pondé de Oliveira Pinto, M., Raiman, J., Salimans, T., Schlatter, J., Schneider, J., Sidor, S., Sutskever, I., Tang, J., Wolski, F. & Zhang, S. (2018) Dota 2 with Large Scale Deep Reinforcement Learning. <http://arxiv.org/abs/1912.06680>.
- 8 Zhang, R., McNeese, N., Freeman, G. & Musick, G. “An Ideal Human” Expectations of AI Teammates in Human-AI Teaming. *Proceedings Of The ACM On Human-Computer Interaction*. 4, 1-25 (2021)
- 9 Myers, C., Furqan, A., Nebolsky, J., Caro, K. & Zhu, J. Patterns for how users overcome obstacles in voice user interfaces. *Proceedings Of The 2018 CHI Conference On Human Factors In Computing Systems*. pp. 1-7 (2018)
- 10 Zhu, J., Liapis, A., Risi, S., Bidarra, R. & Youngblood, G. Explainable AI for designers: A human-centered perspective on mixed-initiative co-creation. *2018 IEEE Conference On Computational Intelligence And Games (CIG)*. pp. 1-8 (2018)
- 11 Sjöblom, M., & Hamari, J. (2017). Why do people watch others play video games? An empirical study on the motivations of Twitch users. *Computers in human behavior*, 75, 985-996.
- 12 Nguyen, T., Hsu, D., Lee, W., Leong, T., Kaelbling, L., Lozano-Perez, T. & Grant, A. Capir: Collaborative action planning with intention recognition. *Seventh Artificial Intelligence And Interactive Digital Entertainment Conference*. (2011)
- 13 Sfikas, K. & Liapis, A. Collaborative agent gameplay in the pandemic board game. *International Conference On The Foundations Of Digital Games*. pp. 1-11 (2020)
- 14 Kreminski, M., Acharya, D., Junius, N., Oliver, E., Compton, K., Dickinson, M., Focht, C., Mason, S., Mazeika, S. & Wardrip-Fruin, N. Cozy Mystery Construction Kit: prototyping toward an AI-assisted collaborative storytelling mystery game. *Proceedings Of The 14th International Conference On The Foundations Of Digital Games*. pp. 1-9 (2019)
- 15 Ashktorab, Z., Liao, Q., Dugan, C., Johnson, J., Pan, Q., Zhang, W., Kumaravel, S. & Campbell, M. Human-ai collaboration in a cooperative game setting: Measuring social perception and outcomes. *Proceedings Of The ACM On Human-Computer Interaction*. 4, 1-20 (2020)
- 16 Grabe, I., González-Duque, M., Risi, S. & Zhu, J. Towards a Framework for Human-AI Interaction Patterns in Co-Creative GAN Applications. (2022)

## 4 Panel discussions

### 4.1 Discussion and Evaluation

Pieter Spronck (Tilburg University, NL), Setareh Maghsudi (Universität Tübingen, DE), and Diego Perez Liebana (Queen Mary University of London, GB)

License  Creative Commons BY 4.0 International license  
© Pieter Spronck, Setareh Maghsudi, and Diego Perez Liebana

On the last day of the seminar only, the participants gathered in the common room to have an evaluation and discussion of the seminar, and to look forward to a possible follow-up seminar. Multiple topics came up, which are discussed below.



#### 4.1.1 Seminar setup

The setup of the seminar was as follows: every day, immediately after breakfast, participants gathered in the common lecture room, where workgroups were formed around themes proposed by the participants. These workgroups consisted of at least three and at most eight people (usually four or five). They worked on their respective themes for the day. This work could consist of a discussion, the building of a prototype, or the running of an experiment. Around five o'clock everyone gathered in the common lecture room for a plenary session, where each workgroup presented their achievements. Usually a workgroup ended after one day, but a few ran a second day, sometimes with different participants, sometimes with a variation on the theme.

After dinner, on Tuesday, Wednesday, and Thursday, an extra activity was planned: on Tuesday and Thursday this concerned a lecture/discussion led by one of the participants on a topic which was of general interest, while on Wednesday this concerned the playing of a game. This game was a roleplaying game to commemorate Daniel Ashlock, one of the organizers and a rather prolific member of the community, who passed away two months before the seminar took place. The game was based on a collection of ideas that Dan developed for roleplaying games (the game has been made available for free from <https://www.drivethrurpg.com/product/400792/Ashlocks-Maze>). Thursday evening also a pub quiz was held, and VR games were made available on several other evenings.

Dagstuhl recommends having a longer walk on one of the days. As was suggested during the previous seminar, we replaced that with a shorter, 45-minute walk, every day after lunch. Still, as many participants also desired a longer walk, time was set aside on Wednesday afternoon to have that.

Before, during, and after the seminar, we used Discord to communicate between participants.

#### 4.1.2 Workgroups

For previous seminars we had participants approach the blackboard to write down ideas for workgroups, after which participants wrote their names next to some of those ideas, to decide which workgroups would start. A problem that was recognized with this approach, is that newcomers to the seminar or the research field might feel intimidated by this process and thus reluctant to bring their ideas forward. In an attempt to resolve this problem, we let participants write their ideas on sheets of paper which we collected, and then anonymously wrote down on the blackboard.

The big disadvantage of this approach was that there were many more ideas formulated than with the old approach (which is good), but with a lot of overlap between them. This made it harder to come up with a good set of workgroups to run, especially since we did not want to have ideas “get lost”. In the end everything worked out, but that was also because the number of participants was smaller than for previous seminars. We had already more ideas than fit on six blackboards with this relatively small group.

It was noted that rather than using the blackboard to form workgroups, we could have used Discord for that. It may even be possible to let then people sign up for multiple workgroups and have a semi-automated process divide the workgroups over the seminar days so that every participant can attend an optimal number of workgroups that hold their interest.

One point of stress that was recognized is that participants found it hard to choose between workgroups as there were multiple that they wanted to be part of. When too many people wanted to be part of one workgroup, it was split up along lines of interest, where each “subgroup” worked on their own perspective on the theme. This led to a suggestion that

we could actually spend part of the seminar on having most or all participants work on the same theme, but in randomly constituted workgroups, which at the end of the day would all present their own view on the theme. A possible enhancement to this idea is that the groups would reform halfway through the day to stimulate cross-pollination. This idea needs some consideration, as a disadvantage of it would be less freedom in choosing workgroups.

### 4.1.3 Evening program

One general remark was that the evening program was for many participants “a bit much.” A possible reason why many felt this way, may be found in the effects of the recent COVID crisis, where people became less used to interacting with bigger groups for a long period of time. Especially the fact that the talks were immediately after dinner (which was deliberate, to still have social time afterwards) was deemed “intense.” The organized games, however, were evaluated positively.

It was suggested that perhaps the Sunday evening could also be used for some activity, though this should be an “unofficial” activity as the seminar officially starts on Monday.

### 4.1.4 Relaxation time

The short walks after lunch were appreciated by many participants, although for most of the week it was rather hot and therefore not ideal for walks. Some participants asked for extra time to exercise. Ideally, such time should be at the end of the day, before dinner (because exercising right after lunch or dinner is not a good idea). This would shorten the time available for workgroups; therefore, a possible approach could be to leave out the short walk after lunch, and instead have the plenary session between four and five o'clock, and have the time after that available for a walk or exercising.

Optionally, the short walk after lunch could still take place, but it should start around 12.45, as we found that by that time most participants had already finished lunch. That way, not much time is lost from the afternoon sessions.

### 4.1.5 Recording

Considerable discussion time was spent on “recording the seminar.” This discussion started with the suggestion that the plenary sessions at the end of the day could be recorded and made available to people who do not attend the seminar. The suggestion was extended with the idea that a short video which shows off the seminar as a whole (including morning sessions, workgroups, plenary sessions, and social interaction) could be used as a way to show to invitees who have not been at Dagstuhl before what the seminar is like, and make them enthusiastic about accepting an invitation.

There are, however, issues with this. A major issue is privacy: not everyone might agree to being recorded, and even if they agree, they might feel uncomfortable with it. Moreover, due to the friendly atmosphere, people tend to be open about their ideas and how they talk about them, but when a camera is present they may feel guarded and cautious.

Recording the plenary session may be a bridge too far, but there would be a lot of value in a 15-20 minute documentary on a next seminar as an advertising tool. This could, for instance, consist of some soundless recordings (with music and commentary) of moments during one day of the seminar, interspersed with brief interviews with participants and explanations about the seminar's setup. Participants who do not want to be recorded can wear visual labels which indicate that they “opt-out.”

#### 4.1.6 Invitation process

The invitation process was rather involved this time around. For previous seminars, we needed at most two rounds of invitations before the seminar filled up. This time, due to the COVID crisis, we had to be cautious in sending invitations, so that early in the process we were very limited in the number of invited participants. Later on this was expanded, and we could invite more participants. At some point the seminar was close to being completely full, but then we ran into a second problem: due to changes in policies at many universities and institutes (particularly in Asia and the US), rising fears of people traveling, and the war in Ukraine, many participants started to drop out again. We frantically invited new people late in the process, up to two weeks before the seminar took place, but that was so late that almost no one could make it.

We invited mostly people from Europe, but we also sent a good number of invitations to Asia and North America, and a few to South-America, Oceania, and Africa. We ended up with just over 30 participants, out of 45 that would have been possible. Over the course of the invitation process, we invited close to 100 people, of which more than half were women, and about half were ‘junior’ people. While in the end the majority of the participants were people who had visited an earlier seminar, we managed to bring in multiple new faces. Often these were people who had no idea what Dagstuhl Seminars were about, but got a recommendation from someone who was aware of the event.

#### 4.1.7 Organization

As always, from the side of Dagstuhl the organization and support were excellent. We noted a few possible improvements, which were communicated to Dagstuhl. As for the scientific organization, we got one recommendation for a follow-up seminar, which was to place a ballot box in the common room where people can deposit ideas or potential complaints – not that there appeared to be anything to complain about, but the existence of such a box would take away hesitation in reporting complaints as it offers the possibility of anonymity.

#### 4.1.8 Topics for a follow-up seminar

Three topic ideas were brought forth for a potential follow-up seminar: (1) multi-agent social games; (2) benchmarks for game AI; and (3) creativity for games. Clearly, the current group of participants would enthusiastically support a follow-up seminar.

## Participants

- Maren Awiszus  
Leibniz Universität  
Hannover, DE
- Cameron Browne  
Maastricht University, NL
- Duygu Cakmak  
Creative Assembly –  
Horsham, GB
- Alex J. Champandard  
creative.ai – Wien, AT
- Guillaume Chanel  
University of Geneva, CH
- Michael Cook  
Queen Mary University of  
London, GB
- Alena Denisova  
University of York, GB
- Alexander Dockhorn  
Leibniz Universität  
Hannover, DE
- Manuel Eberhardinger  
Hochschule der Medien –  
Stuttgart, DE
- Jakob Foerster  
University of Oxford, GB
- Casper Hartevelde  
Northeastern University –  
Boston, US
- Amy K. Hoover  
New Jersey Institute of  
Technology (NJIT) – Newark, US
- Ahmed Khalifa  
University of Malta – Msida, MT
- Antonios Liapis  
University of Malta – Msida, MT
- Daniele Loiacono  
Polytechnic University of  
Milan, IT
- Simon M. Lucas  
Queen Mary University of  
London, GB
- Setareh Maghsudi  
Universität Tübingen, DE
- Ana Matran-Fernandez  
University of Essex –  
Colchester, GB
- Paris Mavromoustakos Blom  
Tilburg University, NL
- Mark J. Nelson  
American University –  
Washington, US
- Mirjam Palosaari Eladhari  
Södertörn University –  
Huddinge, SE
- Diego Perez Liebana  
Queen Mary University of  
London, GB
- Mike Preuß  
Leiden University, NL
- Lisa Rombout  
Tilburg University, NL
- Jacob Schrum  
Southwestern University –  
Georgetown, US
- Pieter Spronck  
Tilburg University, NL
- Tommy Thompson  
AI and Games – London, GB
- Remco Veltkamp  
Utrecht University, NL
- Vanessa Volz  
modl.ai – Copenhagen, DK
- Jichen Zhu  
IT University of  
Copenhagen, DK



# Visualization Empowerment: How to Teach and Learn Data Visualization

Benjamin Bach<sup>\*1</sup>, Sheelagh Carpendale<sup>\*2</sup>, Uta Hinrichs<sup>\*3</sup>, and  
Samuel Huron<sup>\*4</sup>

1 University of Edinburgh, GB. [bbach@ed.ac.uk](mailto:bbach@ed.ac.uk)

2 Simon Fraser University, Vancouver, CA. [sheelagh@sfu.ca](mailto:sheelagh@sfu.ca)

3 University of Edinburgh, GB. [uhinrich@ed.ac.uk](mailto:uhinrich@ed.ac.uk)

4 Télécom Paris, Institut Polytechnique de Paris, CNRS i3 (UMR 9217),  
Palaiseau, FR. [shuron@enst.fr](mailto:shuron@enst.fr)

---

## Abstract

Data visualization is becoming an important asset for a data-literate, informed, and critical society. Despite the variety of existing resources to teach theories and practical skills in this domain, little is known about 1) how learning processes in the context of visualization unfold and 2) best practices for engaging and teaching data visualization to diverse audiences and in different contexts. This Dagstuhl Seminar invited practitioners, researchers, and teachers from the areas of visualization, design, education and cognitive psychology to explore these questions from multiple perspectives. Through a range of practical activities, talks, and discussions, we have begun characterizing and classifying teaching methodologies. We have redacted a pedagogical manifesto, and started formalizing the concept of improvisation with visualization in the context of teaching and learning. We have also interrogated creativity as an important aspect of visualization teaching and learning and explored links between data physicalization and visualization teaching activities. Across these different themes, we have begun to map out the challenges of visualization teaching and learning and the opportunities for research and practice in this area.

**Seminar** June 26–July 1, 2022, Dagstuhl Seminar 22261 – <http://www.dagstuhl.de/22261>

**2012 ACM Subject Classification** Human-centered computing → Visualization theory, concepts and paradigms; Human-centered computing → Visualization design and evaluation methods

**Keywords and phrases** Information Visualization, Visualization Literacy, Data Literacy, Education

**Digital Object Identifier** 10.4230/DagRep.12.6.83

## 1 Executive Summary

*Benjamin Bach (University of Edinburgh, GB)*

*Samuel Huron (Télécom Paris, Institut Polytechnique de Paris, CNRS, FR)*

*Uta Hinrichs (University of Edinburgh, GB)*

*Sheelagh Carpendale (Simon Fraser University – Burnaby, CA)*

**License**  Creative Commons BY 4.0 International license

© Benjamin Bach, Samuel Huron, Uta Hinrichs, and Sheelagh Carpendale

This seminar set out to discuss timely issues and approaches to teaching and education in data visualization. The topic is of growing importance in a world where more and more content is being shared through online news and social media. Our mission as researchers, practitioners and educators in data visualization is to assure quality education for everyone

---

\* Editor / Organizer



Except where otherwise noted, content of this report is licensed under a Creative Commons BY 4.0 International license

Visualization Empowerment, *Dagstuhl Reports*, Vol. 12, Issue 6, pp. 83–111

Editors: Benjamin Bach, Sheelagh Carpendale, Uta Hinrichs, and Samuel Huron



DAGSTUHL  
REPORTS

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

engaging with visualization; this ranges from visualization designers, data scientists, school teachers, journalists, working professionals, students, as well as general public audiences. Teaching visualization is tricky for a range of reasons:

- Data visualization is a skill that is only slowly starting to make its way into school curricula (at least in some countries);
- While the **range of visualization tools** available makes it easy for almost anyone to create visualizations regardless of their technical background, it can be overwhelming to know where to start and to navigate this ever-growing and changing tool landscape;
- Data visualization is a highly interdisciplinary field, influenced and moved forward by psychology, cognitive science, design, computer science, data science, art, and many more disciplines. As a result, **learning objectives and teaching practices greatly vary**;
- There is currently **no defined agreement on the learning goals and criteria** for visualization literacy. For example, what defines a beginner, intermediate, professional in data visualization? What aspects of data visualization should be taught at different levels? And: how can we assess visualization skills?;
- From a learner perspective, the motivation to pick up visualization as a skill is broad: some people “just” want to use a specific tool to get things done quickly, others pursue a design approach (no coding language required), others want to build systems for visualization (Computer Sciences), others go on and become educators or researchers;
- Visualization is important in many domains and knowledge and specific solutions might be specific to these domains, rather than valid universally (e.g., color choices, symbolic conventions, level of interactivity);
- There are a lot of **tacit knowledge and skills** involved in visualization which can be difficult to pin down and transform into learning activities.

In order to discuss these challenges and how to navigate them, we invited participants from academia and industry, including senior and junior thinkers.

### Participants & Seminar Format

Given the highly interdisciplinary field of data visualization and visualization literacy, participants covered a range of expertises including the fields of design, computer science, human-computer interaction, education, graphics, and cognitive psychology. The 5-day seminar was run in a hybrid format with 28 participants joining us at Schloss Dagstuhl in-person, 7 participants joining us online synchronously from Europe, and 6 participants joining us asynchronously from North America. Two organizers were on-site at Schloss Dagstuhl, while two joined the seminar remotely (one synchronously and one asynchronously). One of the online organizers led the asynchronous North America group from Canada. All seminar participants (synchronous and asynchronous) met for a daily debriefing session at 5pm local (Dagstuhl) time to share their progress and discussions. The synchronous remote participants (Europe) joined different local discussion groups through online calls, which did work out surprisingly well – *special thanks to the Dagstuhl technical team for the amazing help with the hybrid setup.*

### Seminar Structure & Activities

The seminar followed an open-ended approach with respect to the possible outcomes, to allow discussion topics to emerge and develop, based on participants’ expertise and interests. Discussions were sparked by brief talks and visualization activities led by selected participants.

The **seminar talks** included presentations on visualization teaching and learning with children, a syntactic analytical framework for visualization, engaging new students with visualization, using forums to engage students with visualization content, how to approach and streamline large-scale assessment of university students' visualization projects, as well as an overview over a book project from a past Dagstuhl Seminar (find the complete list and abstracts of talks in Section 5).

From a practical end, the **visualization activities** invited seminar participants to actively engage in and experience a number of visualization teaching methods and techniques (see section 6 for more details). One activity invited participants to sketch their relation to the seminar topic in order to introduce participants to each other and to start immersing them into the seminar topic. Another activity asked participants to analyze a given visualization systematically. In one activity, we classified existing visualization activities that were submitted by participants prior to the seminar. Another activity took a speculative approach to visualization, inspiring critical visualization scenarios and designs through a card game.

There was ample time to discuss topics of interests through **breakout groups** which focused on topics related to

- Teaching methods and taxonomies for educational activities;
- Teaching creativity and criticality for visualization;
- Data physicalization and how corresponding methods can be used for education and engagement;
- Practical approaches to teaching visualization and the politics involved in teaching visualization;
- Approaches to visualization teaching and creation inspired by improvisation in the arts, and eventually;
- Grand challenges in visualization education.

From an organizer perspective, the seminar was a great success. All participants – both on-site and online – were extremely engaged, and we obtained very positive feedback. Participants appreciated the creative and open-ended nature of this seminar that invited for sharing and reflection of practices from different disciplines and perspectives. The seminar produced a long list of outcomes ranging from paper outlines and book projects, to collecting teaching manifestos and taxonomies, to grant projects and platforms for sharing teaching tools and resources. The plan emerged to establish a reoccurring international symposium around visualization education as part of the IEEE VIS conference, the largest annual conference on visualization with over 1000 participants. The individual working groups will move their individual goals forwards after the seminar. As organizers, we will coordinate between groups and support each of the projects as best as we can, e.g., through regular check-ins with the workgroup leaders as well as townhouse meetings with all Dagstuhl participants, e.g., once a semester. We all believe strongly that this Dagstuhl Seminar – the first formal event on visualization education besides smaller conference workshops – has created a strong momentum for visualization empowerment and education, and we are looking forward to sharing our outcomes on a dedicated website soon.

## 2 Table of Contents

### Executive Summary

*Benjamin Bach, Samuel Huron, Uta Hinrichs, and Sheelagh Carpendale* . . . . . 83

### The week at a glance

Monday . . . . . 88

Tuesday . . . . . 90

Wednesday . . . . . 91

Thursday . . . . . 91

Friday . . . . . 92

### Working Groups

Working group on creativity

*Jonathan C. Roberts, Fateme Rajabiyazdi, Rebecca Noonan, Christina Stoiber, Andy Kirk, Fanny Chevalier, Nathalie Riche, Magdalena Boucher, Alexandra Diehl, Benjamin Bach, and Samuel Huron* . . . . . 92

Working Group on Improvisation with Visualization

*Émeline Brulé, Sheelagh Carpendale, Dietmar Offenhuber, and Charles Perin* . . . . . 95

Working Group on Democratization & Manifesto

*Georgia Panagiotidou, Jagoda Walny, Soren Knudsen, Uta Hinrichs, Wesley Willett, Jason Dyke, Tatiana Losev, Doris Kosminsky, and Samuel Huron* . . . . . 95

Working Group on Physicalization

*Uta Hinrichs, Wolfgang Aigner, Peter Chen, Georgia Panagiotidou, Sarah Hayes, Trevor Hogan, Tatjana Losev, Andrew Manches, Luiz Morais, Till Nagel, and Rebecca Noonan* . . . . . 96

Working Group on Teaching methods

*Isabel Meirelles, Jan Aerts, Wolfgang Aigner, Mashaël Alkadi, Magdalena Boucher, Alexandra Diehl, Christoph Huber, Mandy Keck, Christoph Kinkeldey, Søren Knudsen, Robert S Laramée, Areti Manataki, Till Nagel, and Laura Pelchmann* . . . . . 97

Working Group on Challenges

*Benjamin Bach, Jan Aerst, Andy Kirk, Madny Keck, Till Nagel, Areti Manataki, Soren Knudsen, Georgia Panagiotidou, Wesley Willet, Bob Laramée, Uta Hinrichs, Isabel Meirelles, Benjamin Bach, Doris Kosminsky, Tatiana Losev, Jagoda Walny, Luiz Morais, Fateme Rajabiyazdi, Alexandra Diehl, Wolfgang Aigner, Samuel Huron, and Peter Cheng* . . . . . 98

### Overview of Talks

The potential of a more embodied approach to supporting children's data understanding

*Andrew Manches* . . . . . 98

How I can help you? How can you help me?

*Andy Kirk* . . . . . 99

Cognitive Science of Representational Systems

*Peter Cheng* . . . . . 100



Breaking the Monolith	
<i>Isabel Meirelles</i> . . . . .	100
Jason Dykes tips about assessment	
<i>Jason Dykes</i> . . . . .	101
From visioning to solution, via sketching	
<i>Jonathan C. Roberts</i> . . . . .	101
VisGuides	
<i>Alexandra Diehl</i> . . . . .	103
Teaching visualization Free Form	
<i>Fateme Rajabiyazdi</i> . . . . .	103
Making with Data – Using an open structured template to document the creation of physical data objects	
<i>Till Nagel</i> . . . . .	104
Visual Robot Glyphs	
<i>Jason Dykes</i> . . . . .	104
Reflective on VISualization Learning Outcomings (VISLOs)	
<i>Jason Dykes</i> . . . . .	105
<b>Overview of Activities</b>	
Sketching Introductions: An ice-breaker	
<i>Tatiana Losev</i> . . . . .	106
Visualization Futures Cards	
<i>Wesley Willett</i> . . . . .	106
Classifying teaching activity in a design space	
<i>Samuel Huron</i> . . . . .	107
Activity : Cognitive Science of Representational Systems	
<i>Peter Cheng</i> . . . . .	109
<b>Summary</b> . . . . .	109
<b>Participants</b> . . . . .	110
<b>Remote Participants</b> . . . . .	111

PST	EST	CEST	Monday	Tuesday	Wed	Thu	Fri
00:00	03:00	7:30 - 9:00	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast
		09:00	Welcome!	Morning Meeting, Logistics	Morning Meeting, Logistics	Morning Meeting, Logistics Pre-approved VIS workshop	Wrap up
00:15	03:15	09:15	Icebreaker activity	3 x 5min talks + discussion - Andrew Manches - Andy Kirk - Peter Cheng	3 x 5min talks + discussion - Isabel Mirelles - Jonathan Roberts - Jason Dykes	3 x 5min talks + discussion - Fateme Rajabi - Doris Kosminsky - Alexandra Diehl - Wes & Till	
		09:45		Breakout Groups	Breakout Groups	Breakout Groups	
01:30	04:30	10:30 - 10:45	Coffee	Coffee	Coffee	Coffee	Coffee
			Brainstorming	Breakout Groups Plenary group discussion	Categorizing the Vis Activities cards (Samuel Huron) Framing learning objectives	Breakout Groups	people leave
03:00	06:00	12:15 - 14:00	Lunch *starts promptly	Lunch *starts promptly	Lunch *starts promptly	Lunch *starts promptly	
05:00	08:00	14:00	Breakout groups	VIS Activities: Wesley Willett (30 minutes) Breakout groups	Social activity (for on-site participants)	1:30pm - Breakout groups / VIS Activities Peter Cheng (60min) Coffee & Cake	
		15:30 - 16:00	Coffee & Cake	Coffee & Cake			
07:00	10:00	16:00	Breakout groups	Breakout groups			Breakout groups
08:00	11:00	17:00-00 (ALL)	Informal report back (ALL) Ideas & questions for the larger group	Informal report back (ALL) Ideas & questions for the larger group - Northern Americans introduce themselves to the group			Informal report back (ALL) Ideas & questions for the larger group
09:00	12:00	18:00	Dinner *starts promptly	Dinner *starts promptly		Dinner *starts promptly	

■ Figure 1 Seminar schedule.

### 3 The week at a glance

#### 3.1 Monday

After introducing the seminar and presenting various organizational matters, the first day started with a **sketching activity** organized by Tatiana Losev from Simon Fraser University in Vancouver, Canada (Section 6.1). This activity acted as an icebreaker for participants start to get to know each other, and to initiate reflections on the seminar topic in a playful way. Participants were invited to sketch their relationship to the seminar topic, and to then use this sketch to introduce themselves in the form of a 1-minute presentation. Some participants sketched their relation to teaching and research in a literal, metaphorical or abstract way, others focused on the variety of themes and questions to be discussed, and again others created visual representations of their past curricula (see Figure 2). The activity brought to the fore an exciting diversity of viewpoints as well as the coverage of interest toward the topics of the seminar.

The sketching activity was followed by a **brainstorming session** that invited participants to identify high- and low-level topics around visualization empowerment and teaching, that they wanted to discuss during the seminar. This session was intended to initiate and fuel discussions that would take place in the form of smaller working groups throughout the week. Participants noted topics on sticky notes that we then collaboratively reviewed and grouped (see Figure 3). We identified a great diversity of themes including design creativity, physicalization, ethics, democratization, scalability of teaching, humanism, tools, hybrid and online teaching, teaching methods, community building, success stories and inspiration, learning goals, planning, contexts, and barriers, audiences (from practitioners to children), evaluation and assessment, measuring learning progress, interdisciplinarity, critical thinking, inclusivity, resources, and cataloging educational material.

We ranked topics in a voting activity, based on participants' interests to discuss them. This led to the formation of initially four working groups: teaching methods, democratization, creativity, and physicalization.



■ **Figure 2** Sketches produce by the participant for their introduction.

- The **Teaching Methods Group** (see Section 4.5) first focused on the diversity of challenges in visualization education. They explored possibilities of formalizing a multidimensional problem space to capture these challenges.
- The **Democratization Group** (see Section 4.3) discussed practices, beliefs, intentions and biases that influence visualization teaching and creating with the aim to make visualization more accessible. Based on these discussions they decided to question the value behind our teaching activities.
- The **Creativity Group** (see Section 4.1) focused on how to teach creativity and criticality in visualization; how one can be creative in teaching visualization, how one can teach creativity through visualization, and eventually focused on an activity book for novice visualization designers.
- The **Physicalization Group** (see Section 4.4) discussed how data physicalizations could be a mediator for teaching and learning activities, but also how it can be used to breach disciplines, and also how inclusive the physicalization can be for teaching and learning.

These working groups continued discussions throughout the week in different participant constellations. A number of participants shifted between groups to absorb different discussions, which proved to be useful for cross-dissemination across working groups.

At the end of the day, i.e., after some initial discussions and topic finding within the individual groups, each group briefed the entire seminar on their discussion and focus. The North America group joined to get updated on the European groups.



■ **Figure 3** Identification of thematic thought creating an affinity diagram of post it.



■ **Figure 4** A photo of the visualization future card game.

### 3.2 Tuesday

This day was mainly reserved for discussions within the individual working groups. “Over night”, the one organizer based in North America lead the North America discussion group which decided on the topic of Improvisation in visualization and what can be learned from improvisation in art for how to approach visualization (design). We started the day with a short briefing into the day schedule and asked people if they wanted to switch or split groups. Then, we had a series of short 5min talks from seminar participants Andrew Manches (education) providing a learning science perspective in his talk *The potential of a more embodied approach to supporting children’s data understanding* (Section 5.1); Andy Kirk (freelance visualization designer and educator), provided a non academic visualization trainer perspective “How I can help you? How can you help me?”(Section 5.2); eventually Peter Cheng brought a cognitive psychology perspective through his talk “Cognitive Science of Representational Systems” (Section 5.3). A joined question and answer session followed these talks. Then, participants broke out into their groups. At the beginning of the afternoon Wesley Willet ran a visualization activity “Visualization future card game”(Section 6.2).

Again, at the end of the day, all working groups met, including the North America group on *ImproVISation* to brief the other seminar participants.

### 3.3 Wednesday

We started the day, again, with four short 5min talks: Isabel Meirelles (visualization design) “Breaking the Monolith” (Section 5.4); Jason Dykes (cartography and visualization) about assessment of visualization teaching (Section 5.5); Jonathan Roberts (visualization) “From visioning to solution via sketching” (Section 5.6). Following these presentations a vibrant and spontaneous conversation happened among participants about the different assessment strategies in various teaching constraints. The discussion was so spontaneous and interesting it took most of the morning.

Before the afternoon socialization, Samuel Huron ran a visualization activity with all participants, aiming to classify visualization activities. Prior to the seminar, Samuel had invited participants to submit short descriptions of activities they do in their classes with their students. This collection was printed on small cards, one activity per card, and distributed among each working group. Each working group came up with a different classification scheme which is currently informing an ongoing discussion. There was no evening briefing with the North America group due to the socialization activity.

### 3.4 Thursday

During the morning talks, Fateme Rajabiyazdi shared some of the lessons learned in her class during the talk “Teaching visualization Free Form ”(Section 5.8). Doris Kosminsky discuss how can we empower mother with health data in Brazil in her talk “Reflections on learning and empowerment of those represented in health visualization”. Then Alexandra Diehl presented how Visguides <sup>1</sup> could be use for education. Visguides is an open community project to provide guidance and support on visualization design in a web forum (Section 5.7). Last Till Nagel presented the goal, reflect on the process and the design of the book *Making with data* (Section 5.9).

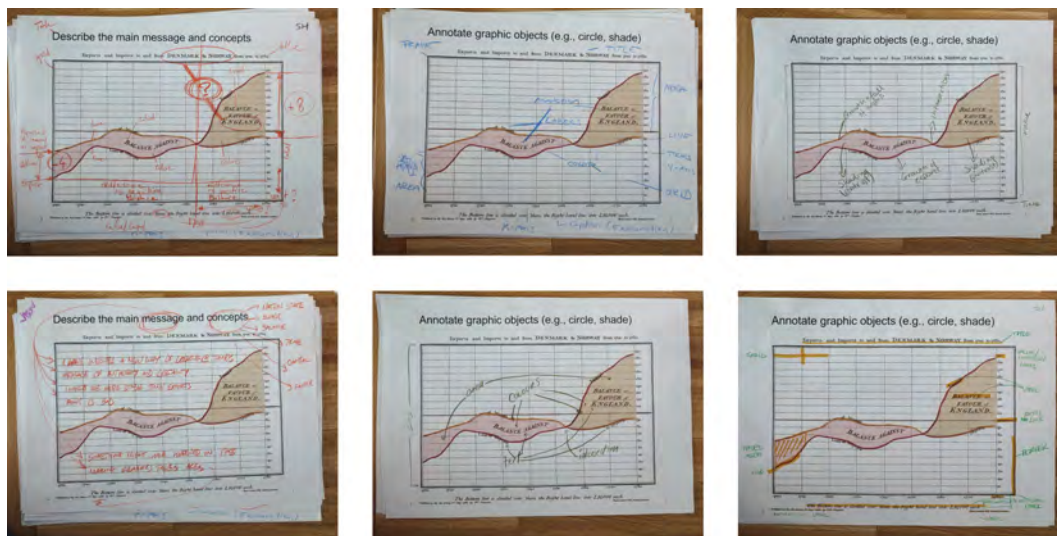
Mandy Keck proposed to create a pre-approved symposium at IEEE VIS conference, the major international forum for visualization with over 1000 attendees. The symposium would become a major outlet and forum for research around visualization education and learning. It would be a place for research, reflection, creation, and discussion of learning / teaching but also discussing higher-level issues in regards to human-centered approaches to visualization education and design, and build a permanent forum and community around these topics.

Later in the afternoon, Peter Cheng ran an activity to model the cognitive process of reading a visualization through annotating a visualization and then modeling the different cognitive steps of our reading procedure (Section 6.4).

The remainder of this day (morning, afternoon) was reserved for discussions within the working groups. The organizers encouraged goal-oriented thinking and to list and plan the different outcomes of each working group to be reported in the pre-dinner briefing session. Since some participant had to leave early on Friday morning, we started a general discussion about the individual outcomes of this seminar and how to organize working groups beyond the seminar.

---

<sup>1</sup> <https://visguides.org/>



■ **Figure 5** A visualization annotated by a participant from Cheng’s activity.

### 3.5 Friday

We started the day by synthesizing and presenting on one page all the potential outcomes that were planned by the individual working groups. Then, Jason Dykes gave a short talk to introduce data visualization to an audience (Section 5.10). Then the group started a discussion to reflect on the seminar experience and outcomes. This discussion push up to open two other topics thread, one on grants, and one on writing a paper about the main challenges in information visualization teaching and learning. After the coffee break the remaining participants decided to outline collectively a paper on grand challenges in visualization education, effectively forming a sixth working group at the seminar.

## 4 Working Groups

### 4.1 Working group on creativity

*Fateme Rajabiyazdi (Carleton University – Ottawa, CA), Rebecca Noonan (Munster Technological University – Cork, IE), Jonathan C. Roberts (Bangor University, GB), Christina Stoiber (FH – St. Pölten, AT), Andy Kirk (Visualising Data – Leeds, GB), Fanny Chevalier (University of Toronto, CA), Nathalie Riche (Microsoft Research – Redmond, US), Magdalena Boucher (FH – St. Pölten, AT), Alexandra Diehl (University of Zurich, CH), Benjamin Bach (University of Edinburgh, GB), Samuel Huron (Institut Polytechnique de Paris, FR)*

License © Creative Commons BY 4.0 International license

© Jonathan C. Roberts, Fateme Rajabiyazdi, Rebecca Noonan, Christina Stoiber, Andy Kirk, Fanny Chevalier, Nathalie Riche, Magdalena Boucher, Alexandra Diehl, Benjamin Bach, and Samuel Huron

This working group focused on creativity in visualization empowerment. Creativity is the use of people’s imagination to engender original ideas, to make and create something, be inventive, design new ideas or make different designs.

#### 4.1.1 Discussed Problems

After introductions, we discussed our backgrounds. The creativity group brought together people with different experiences and situations. Our backgrounds are diverse: ranging from PhD students, researchers, early career academics, company directors, to senior researchers and academics. We teach to undergraduates, postgraduates and the public. Most of us sit in computing, mathematics, and engineering schools; while some are in industry. But more and more we are teaching and discussing with people from broader backgrounds, and noted specifically those in arts, design, humanities, social science and psychology.

We focused on three questions, around teaching creativity, being creative in the pedagogic process, and creativity in data-visualisation.

- **What is creativity and how can we teach it?** This question focuses on the *learner*. Teaching creativity is not necessarily easy. We discussed many different ideas, from creativity and innovation, to inspiration. How (as a teacher) can we encourage, and help other people to be creative? What do we need to teach? What strategies do people need to learn? How can we get learners to be more creative and innovative? What processes can we use to help people to become creative? We discussed many skills that people can learn, from creating, design, elegance, imagination, aesthetics, elegance, harmony, flow, balance, beauty to storytelling.
- **How can teachers be creative when teaching data visualisation?** This question focuses on the *process*. Being innovative and creative in teaching can help to engender excitement, it can encourage people to be creative (pushing them out of their comfort zone) and can help to improve the relationship between learner and teacher. What new ideas can we use in teaching? What tools, technologies and resources can we use? For instance, it is possible to teach creative thinking through sketching, use of LEGO, modeling clay, and so on.
- **What is creativity in data visualization?** This question focuses on the broad challenges in data visualisation as a *domain*. Creativity can be applied to every part of data visualisation process, not just in pedagogic terms. For instance, it is possible to be creative in understanding and using data, in how we approach the visualisation design process, or how we interact with clients. Creativity can be achieved through any part of: research, specification, design, client-interaction, implementation, evaluation, maintenance, and so on.

#### 4.1.2 Possible Approaches

We approached these challenges through discussion using shared online documents. We broadly worked through each of the questions in order. We used zoom, shared Google documents, and Miro board as a virtual white board. We approached the challenge in portions of a few hours. First, we discussed different ideas, took notes in the shared documents, placed sticky-notes on the Miro board (See Figure 6), and added links to external resources in the shared document. We shared our experiences, gave examples of how we used creative activities in our teaching, and bounced off ideas from each other. Second we summarized our ideas, created a short report and reported back to the other Dagstuhl participants.

There were several important discussions and outcomes, and ideas that we will work on after the Dagstuhl Seminar. The group discussed and proposed that there is a huge need for resources. Resources, ideas, inspirational creative activities, and so on, that can help teachers, learners, educators, researchers and developers be creative in visualization.



■ **Figure 6** A screenshot of a part of the creativity group Miro board.

We discussed different ways to collage resources, perhaps to write a book, create an online resource of difference creative recipes, organize a workshop, interview experts, and so on. We also realized that our collective knowledge and experience was important, and that we had many creative ideas that we felt would be useful for others to view.

### 4.1.3 Conclusions

Creative visualization, teaching creativity in visualisation, is an exciting area. People can be creative by creating different assets, videos/illustrations, can be “creative” in how they approach thing (e.g., exploring data, defining audiences), and people can be creative in how they approach and imagine new ideas in data visualization. The group discussion ended with a two stage plan. In the short instance the group wrote a long paper that summarizes activities (demonstrating the collective experience and shared examples that were discussed at Dagstuhl) [1]. In the long term, the group proposed to consider summarizing a broader set of creative data visualization activities, as a larger resource, such as a book and website.

### References

- 1 Jonathan C. Roberts, et al. (2022, October). Reflections and Considerations on Running Creative Visualization Learning Activities. 4th IEEE Workshop on Visualization Guidelines in Research, Design, and Education 2022.



## 4.2 Working Group on Improvisation with Visualization

*Émeline Brulé (University of Sussex, GB), Sheelagh Carpendale (University of Vancouver, CA), Dietmar Offenhuber (Northeastern University – Boston, US), Charles Perin (University of Victoria, CA)*

**License** © Creative Commons BY 4.0 International license

© Émeline Brulé, Sheelagh Carpendale, Dietmar Offenhuber, and Charles Perin

Resources for data-visualization education often emphasize toolkits, frameworks, guidelines, outlining what a good visualization is and the building blocks to develop one. With this group, we asked: what if we were instead emphasizing improvisation and practice? We build on case-studies drawn from our own experiences, art education and performance to outline what *impro-vis* practices look like. We argue for centering such practices has the potential to widen what diverse audiences consider as data; the aesthetic and representation repertoire of students and data-viz practitioners; and strengthen research on the situated and improvisational aspects of visualization.

## 4.3 Working Group on Democratization & Manifesto

*Georgia Panagiotidou (University College London, GB), Jagoda Walny (Canada Energy Regulator – Calgary, CA), Soren Knudsen (IT University of Copenhagen, DK), Uta Hinrichs (University of Edinburgh, GB), Wesley Willett (University of Calgary, CA), Jason Dykes (City University London, GB), Tatiana Losev (Simon Fraser University – Burnaby, CA), Doris Kosminsky (University of Rio de Janeiro, BR), Samuel Huron (Institut Polytechnique de Paris, FR)*

**License** © Creative Commons BY 4.0 International license

© Georgia Panagiotidou, Jagoda Walny, Soren Knudsen, Uta Hinrichs, Wesley Willett, Jason Dyke, Tatiana Losev, Doris Kosminsky, and Samuel Huron

We are a group of visualization researchers from different countries, disciplines and generations. We came together to discuss “democratisation” in the context of teaching and learning in visualization. We found that “democratisation” existed in our common desire to empower the people with whom we interact to understand and use data in their lives. We thus set out to develop a shared vision: a manifesto of sorts that would guide us towards strategies to broaden data visualization skills, make them more common and accessible, and enable this empowerment. Instead of creating one common manifesto however, we ended up taking a different, more personal approach of what we understood by empowerment. Our perspectives highlighted our situated understandings, ranging from constructivist teaching and physicalization, to co-design and policy intervention. The variety in our approaches reflected the variety in our backgrounds, and the different situations through which we personally felt we could approach the task of strengthening visualization empowerment in others.

Inspired by this process, we created an exercise that helps visualization educators to elicit their personal reflections and make commitments for their teaching and learning. This exercise, which we named a “me-ifesto”, was eventually supported and co-authored by over 25 researchers present at the Dagstuhl. A “me-ifesto” paper, which described the exercise

and our process, was then presented at the alt. VIS workshop collocated with the IEEE VIS 2022 in Oklahoma [1]. This working group moreover, has since transformed into a recurring meeting in which we, as visualization teachers (and learners), continue to reflect on the values we embed in our teaching both consciously and not.

#### References

- 1 Walny et al, “Me-ifestos for Visualization Empowerment in Teaching (and Learning?)”, alt.VIS workshop, IEEE VIS, 2022, Oklahoma City.

#### 4.4 Working Group on Physicalization

*Wolfgang Aigner (St. Pölten University of Applied Sciences), Peter Chen (University of Sussex, GB), Georgia Panagiotidou (UCL, GB), Sarah Hayes (Cork Institute of Technology, IR), Uta Hinrichs (University of Edinburgh, GB), Trevor Hogan (Cork Institute of Technology, IR), Tatjana Losev (Simon Fraser University, CA), Andrew Manches (University of Edinburgh, GB), Luiz Morais (Inria, FR), Till Nagel (Mannheim University of Applied Sciences), Rebecca Noonan (Cork Institute of Technology, IR)*

License © Creative Commons BY 4.0 International license

© Uta Hinrichs, Wolfgang Aigner, Peter Chen, Georgia Panagiotidou, Sarah Hayes, Trevor Hogan, Tatjana Losev, Andrew Manches, Luiz Morais, Till Nagel, and Rebecca Noonan

This working group focused on distinguishing data physicalization as an activity to support teaching and learning about data visualization – both computer-supported physical representations of data, and hand-made constructions of data. They explored the role of physicalization in learning settings and developed questions, identified gaps and ethical considerations for further research: *How can we, as educators in data VIS, evaluate physicalization activities for classroom settings and public community settings? How might data physicalization, as both an activity and an output of a tangible artifact, facilitate teaching and learning? What are the benefits of using physicalizations as a mediator to bridge disciplines and connect different people and perspectives?* They identified a need to determine learning outcomes for teaching physicalization with different audience groups ranging from children, post-secondary students, the public, and diverse communities of practice. This is important because the benefits and limitations of using physicalization for learning data visualization with sustainability, inclusivity and accessibility are underexplored in computer sciences. By mapping the space of data physicalization in the learning context, the physicalization research group aims to explore the potential and limitations of physicalization as an interactive activity, a tool, data output, and a process for learning and teaching.

## 4.5 Working Group on Teaching methods

*Jan Aerts (Amador Bioscience – Hasselt, Hasselt University & KU Leuven, BE), Wolfgang Aigner (FH St. Pölten, AT), Mashael Alkadi (University of Edinburgh, GB), Magdalena Boucher (FH St. Pölten, AT), Alexandra Diehl (Universität Zürich, CH), Christoph Huber (Hochschule Mannheim, DE), Mandy Keck (Univ. of Applied Sciences – Hagenberg, AT), Christoph Kinkeldey (HAW – Hamburg, DE), Søren Knudsen (IT University of Copenhagen, DK), Robert S Laramée (University of Nottingham, GB), Areti Manataki (University of St Andrews, GB), Isabel Meirelles (OCAD University, CA), Till Nagel (Hochschule Mannheim, DE), Laura Pelchmann (Universität Köln, DE)*

**License** © Creative Commons BY 4.0 International license

© Isabel Meirelles, Jan Aerts, Wolfgang Aigner, Mashael Alkadi, Magdalena Boucher, Alexandra Diehl, Christoph Huber, Mandy Keck, Christoph Kinkeldey, Søren Knudsen, Robert S Laramée, Areti Manataki, Till Nagel, and Laura Pelchmann

The Teaching Methods group worked on:

- **Identifying challenges** that the group participants faced in their own teaching experiences. This resulted in an initial list that was later expanded given that a separate group formed on the last day to focus exclusively on challenges, now called Grand Challenges (Section 4.6). Our group's initial list of challenges can include challenges about Learning & Teaching Resources, Self-guided Learning, Learning participants, Implementation and development, Vis prototyping, measuring, marking and evaluation, teaching methods, story-telling, critical-thinking skills, technology, online (remote) teaching, combining teaching and research, and sharing teaching materials and resources with the wider visualization teaching community. When conducting the activity on Activities, we marked our initial list of challenges based on Moon's Handbook of Reflective and Experiential Learning [1].
- **Creating a multidimensional problem space** towards identifying existent resources and gaps in teaching and learning methods. The discussion consisted in identifying key components in teaching and learning to define topics and dimensions. The work is currently in progress. During Dagstuhl, we created a framework/taxonomy for future use in identifying literature, activities, gaps, etc.
- **Systematizing the Role of development in Data Visualization Teaching.** A subgroup in this working group worked on identifying key components required in implementation and evaluation as related to development in data vis.

### References

- 1 Moon, Jennifer: A handbook of reflective and experiential learning: Theory and practice, Routledge, 2006

## 4.6 Working Group on Challenges

*Benjamin Bach (University of Edinburgh), Jan Aerst, Andy Kirk, Madny Keck, Till Nagel, Areti Manataki, Soren Knudsen, Georgia Panagiotidou, Wesley Willet, Bob Laramée, Uta Hinrichs, Isabel Meirelles, Benjamin Bach, Doris Kosminsky, Tatiana Losev, Jagoda Walny, Luiz Morais, Fateme Rajabiyazdi, Alexandra Diehl, Wolfgang Aigner, Samuel Huron, Peter Cheng*

**License** © Creative Commons BY 4.0 International license  
 © Benjamin Bach, Jan Aerst, Andy Kirk, Madny Keck, Till Nagel, Areti Manataki, Soren Knudsen, Georgia Panagiotidou, Wesley Willet, Bob Laramée, Uta Hinrichs, Isabel Meirelles, Benjamin Bach, Doris Kosminsky, Tatiana Losev, Jagoda Walny, Luiz Morais, Fateme Rajabiyazdi, Alexandra Diehl, Wolfgang Aigner, Samuel Huron, and Peter Cheng

This group emerged spontaneously on Friday morning at the seminar closing session. It started after a question to keep collecting challenges. The group, comprising almost all of the seminar participants, collected around 30 challenges. Some challenges were based on the challenges already collected by the Teaching Methods working group. Others were entirely new. In the months after the seminar, we are still re-organizing these challenges and trying to come up with a suitable structure to describe these challenges. Many challenges are interwoven and otherwise related, e.g., teaching different audiences and hybrid teaching, or learning goals.

## 5 Overview of Talks

### 5.1 The potential of a more embodied approach to supporting children’s data understanding

*Andrew Manches (University of Edinburgh – Edinburgh, GB, a.manches@ed.ac.uk)*

**License** © Creative Commons BY 4.0 International license  
 © Andrew Manches

Early Education has over two centuries’ experience of designing materials to help children learn abstract concepts – such as colored rods to help children learn numerical relationships. Yet the representational transparency of these materials depends upon existing knowledge of the learner. Representations often integrate a range of conceptual and cultural metaphors that are often known but taken for granted by adulthood.

Simply looking at colored rods will not enable a child to just “get” numbers. Pedagogy, clearly, is key. This captures many things (e.g., narrative, construction activities); my work focuses on interaction – how adults scaffold children’s interaction and learning with materials. In particular, I attend to the multimodality and bidirectionality of scaffolding. It is not just words but other modes such as facial expression, body posture, gaze and gesture that educators employ. Gestures are notably powerful in their capacity to provide schematic, dynamic, visuo-spatial representations coproduced with speech to bridge communication with our environment. And importantly, scaffolding is two-way: children also employ a range of modes to structure and manipulate the support they need. The importance of multimodality is further accentuated when considering emerging theories of what it means to know (and hence what children have “learnt”). Increasing evidence points to the embodied nature of cognition and how learning involves the internalization of body-based experiences, emphasizing the interwoven nature of emotional, social, physical, and cognitive dimensions.

When communicating our understanding, we can activate these prior experiences – evident in the emotions and gestures we commonly produce in explanations (makes for a good observational activity at seminars like Dagstuhl).

Increasing attention to the significance of multimodality in how we think and interact has important implications for design as well as pedagogy. This may be greater recognition of existing activities and media (e.g., physicalization) or the potential for more body-based interaction with digital representations (e.g. tangibles, haptics, gesture recognition). More recently, my work has asked how we can tap into emerging theories of cognition and digital tools to help young children (3 years+) understand concepts of data. This is not just a conceptual challenge – children’s worlds are increasingly datafied- from how we measure their “progress” to the (smart) toys we give them. Here there is much potential. Young children understand and often articulate themselves and their interaction with the world (e.g., how old, noisy, tall, active, sleepy, or happy they are) – hence offering a design and learning opportunity through appropriately representing this personal information. Educators already do – colored rods to compare ages, stickers to quantify good behavior or classroom “noisemeters” to maintain sanity. Experts in the field of visualisation/physicalisation have much potential to create a new generation of embodied designs and activities that build upon this foundation.

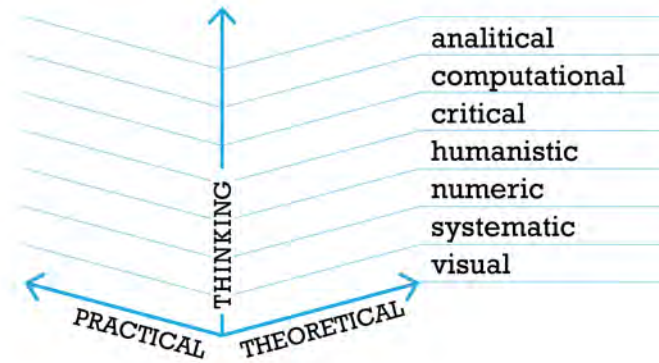
## 5.2 How I can help you? How can you help me?

*Andy Kirk (Visualising Data Ltd – GB)*

License © Creative Commons BY 4.0 International license  
© Andy Kirk

In this lightning talk I introduced myself to members of the seminar, especially as I’m a non-academic participant so my work may be less than familiar to most. I presented an outline of how I may be able to help my fellow attendees. I opened with an overview of my “boundary spanning” activities, as a freelancer: I publish via my website and podcast series, I teach (academically) and research, I consult and design, I author and present. Above all, perhaps, I train professionals, outside of academia and across a diverse array of client organization types and industries. I described some of the key objectives and approaches I take to the challenge of teaching and learning, and how delighted I was to observe alignment with the seminar’s theme of “Visualization Empowerment”. Given my activities and career experiences, could I be of service to offer any guidance to others?

I then switched over to introduce some matters of interest that I am particularly keen to learn about over the seminar, and maybe get some help from others. Firstly, listing some of the current challenges I experience in the forum of public training: including carving out distinctions in teaching levels (basic » advanced) for the same and different cohorts, how to teach the concept of elegance and instill journalistic curiosities. Secondly, and finally, issues that specifically affect training in private/client settings, including how to demonstrate (maybe prove?) success of visualisation in terms like ROI, how to encourage organizational readiness for cultural change, ambitions vs. reproducible pragmatism, and the challenges of educating across such a multi-disciplinary skillset.



■ **Figure 7** Framework of possible pedagogical goals.

### 5.3 Cognitive Science of Representational Systems

*Peter Cheng (Sussex University, GB)*

License Creative Commons BY 4.0 International license  
© Peter Cheng

My approach to the design of information visualisations, and representational systems more generally, combines cognitive science with the analysis of the conceptual structure of the to-be visualized topic. To obtain theoretical and empirical leverage to formulate and test principles of representation design, I create novel diagrammatic systems for conceptually challenging topics and interactive graphical user-interfaces for information intensive decision-making systems. Generalizing over the creation and evaluation of many such systems, I make four claims: (1) STEM topics should be easy to learn; (2) compared to extant conventional visualizations, factor of 2 improvements in problem solving and learning are feasible when representations effectively re-codify knowledge; (3) knowledge re-codification should attempt to capture the conceptual structure of a topic in the graphical structure of the representation; (4) this yields representations that possess semantic transparency and syntactic plasticity.

### 5.4 Breaking the Monolith

*Isabel Meirelles (OCAD University, CA)*

License Creative Commons BY 4.0 International license  
© Isabel Meirelles

The talk invited a conversation about the challenges we encounter preparing students to contribute to data visualization practices. There is an unbalance of how skills are taught across disciplines (sciences, arts, humanities, etc). This unbalance affects education at all levels and disciplines. I would like to suggest that we tailor data visualization pedagogical strategies in a situated manner. For that I proposed a scaffold built around thinking processes positioned along the theoretical-practical axis (7). The thinking processes are derived from the literacies needed for data visualization and dependent on the setting and pedagogical goals of the course and needs of our learners (listed alphabetically): analytical, computational,

critical, humanistic, numerical, systems, and visual (7). Educators can use the scaffold to identify areas of focus and whether they will approach through practical and/or theoretical means.

## 5.5 Jason Dykes tips about assessment

*Jason Dykes (City University London, GB)*

License © Creative Commons BY 4.0 International license  
© Jason Dykes

I made passionate plea to design assessments so that the teacher gets the information they need to make the best judgment that they can in the time they have. Few of my colleagues do this. Aim to use as little time as you (teacher) can to assess, in ways that are as efficient as possible. Aim to make it fun. Yes, really FUN and INFORMATIVE – how well are you teaching? What can your students do? Assessment is a creative design exercise and an informative teaching diagnostic – it helps to see it that way. What would you really like to see? You control this, so if you ask them to give you 400 hours of text to read, well, that’s your fault!

## 5.6 From visioning to solution, via sketching

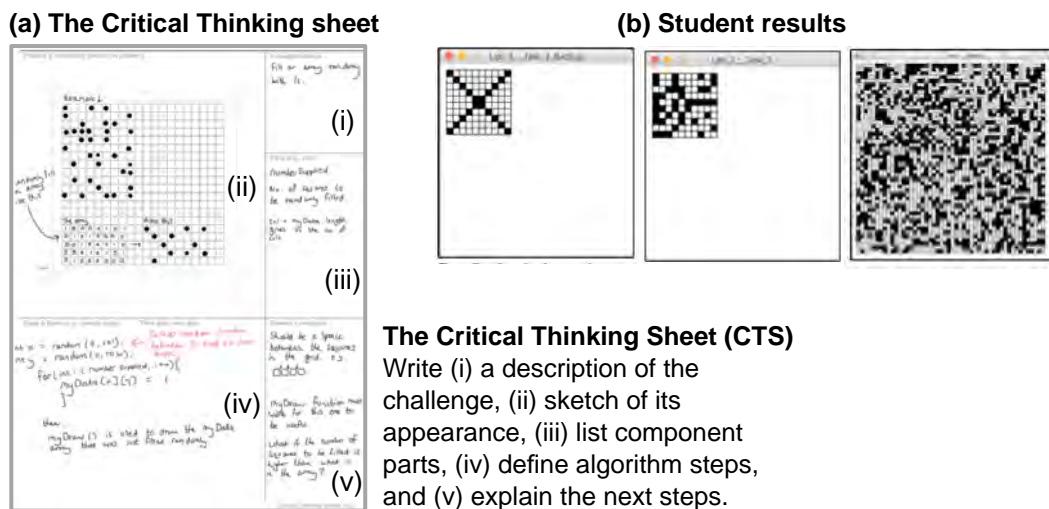
*Jonathan C. Roberts (Bangor University, GB, j.c.roberts@bangor.ac.uk)*

License © Creative Commons BY 4.0 International license  
© Jonathan C. Roberts

What will your visualization look like? What will it do? How will it work? These are important questions for designers and especially learners to ask. Far too often learners, and experienced researchers, create visualisations without thinking what they are doing. Students are often, far too keen to just get coding. When they jump into their code they create solutions that may not be fit for purpose. And when they realize what they have created, it is too close to the deadline to change their mind and adapt it. This early enthusiasm is admirable, and should be encouraged and tapped. Indeed, with some forethought – by becoming more reflective at an early stage, and performing critical thinking – they will create something better. Through thinking and sketching they will be able to contemplate how their solution will work, think who will use it, even imagine a specific person using their tool.

In this lightning talk I presented the need for “visioning”. I proposed that sketching solutions can help people think through their ideas, externalize their thoughts, and project their minds to imagine people using their solution for real, and for its intended purpose.

So to achieve this act of visioning, I proposed that people need to understand their “goal”. The goal then helps to frame the challenge and define the focus. Their goal could be a task they want to fulfill or challenge to solve. In addition, and especially in an education setting, I proposed that these tasks (assessments) should be **authentic** in their design [1]. In other words, that the tasks should be challenges that they could find in their real life (perhaps when the students have a job after they graduate). Furthermore, the task should be individual to each student.



■ **Figure 8** Results of a student using the *Critical Thinking Sheet (CTS)* method [4], to design and create a random pattern generator (using Processing.org). Starting with the CTS they sketch and plan the work, and then iterate better implementations of their solution.

But learners need structure. They need methods to follow. To **scaffold** these vision sketches, I proposed a few techniques. The Five Design-Sheets [2, 3] method uses five sheets of paper, with five stages to help drive critical thinking. Alternatively, for specific tasks a single sheet of sketching and planning could be used. One method is the Critical Thinking sheet [4], which gets students to think about the goal, sketch what the solution would look like, list component parts, articulate algorithmic steps, and list tasks that they need to achieve in order to implement it. Figure 8 shows how a student, thinking about a random pattern generator, starts with a sketch that presents the vision of their solution, before implementing and iterating towards their solution in code.

## References

- 1 Roberts J. C., Ritsos P. D., Jackson J. R., Headleand C.: The explanatory visualization framework: An active learning framework for teaching creative computing using explanatory visualizations. *IEEE Transactions on Visualization and Computer Graphics* 24, 1 (2018), 791–801. doi:10.1109/TVCG.2017.2745878.
- 2 Roberts J. C., Headleand C., Ritsos P. D.: Sketching designs using the five design-sheet methodology. *IEEE Transactions on Visualization and Computer Graphics* 22, 1 (Jan 2016), 419–428. doi:10.1109/TVCG.2015.2467271.
- 3 Roberts J. C., Headleand C. J., Ritsos P. D.: *Five Design-Sheets – Creative design and sketching in Computing and Visualization*. Springer International Publishing, 2017. doi:10.1007/978-3-319-55627-7.
- 4 Roberts J. C., Ritsos P. D.: Critical Thinking Sheet (CTS) for Design Thinking in Programming Courses. In *Eurographics 2020 – Education Papers (2020)*, Romero M., Sousa Santos B., (Eds.), The Eurographics Association, pp. 17–23. doi:10.2312/eged.20201029.



## 5.7 VisGuides

*Alexandra Diehl (University of Zürich, CH)*

License  Creative Commons BY 4.0 International license  
© Alexandra Diehl


VisGuides, engaging the VIS community on democratic discussions

Building a community platform that is open, democratic, and inviting is a big challenge. I presented VisGuides, a democratic discussion forum co-created by several colleagues over Europe. The main goal of VisGuides is to create an open space for evidence-based discussions where we can explore and contest well-known practices and guidelines.

We have been using VisGuides as a collaborative educational tool to collect resources, experiences, and educational material. We want to invite the VIS community to join us in these efforts, share with them resources and collected material, and find new ways of rewarding them for their contributions.

## 5.8 Teaching visualization Free Form

*Fateme Rajabiyazdi (Carleton University – Ottawa, CA, fateme.rajabiyazdi@carleton.ca)*

License  Creative Commons BY 4.0 International license  
© Fateme Rajabiyazdi

In this presentation, I share my experiences and questions as a first-time data visualization instructor.


I offered a data visualization course to 14 graduate students. Students from different disciplines or backgrounds could join the class, pick their own dataset and choice of audience, and their developing platform. I used “Design Study “Lite” Methodology” [1] to outline the course, and here are some lessons learned. Having students from different disciplines helped students learn from other areas of study. It was rather difficult to tailor the content to this heterogeneous group. Having different datasets for visualizing was valuable as students could teach others from other disciplines about their world. The conversations between students helped share knowledge beyond the course outcomes. However, that required me to learn and assess different datasets which are not scalable! By having the option to choose their audience, students said they could target audiences beyond class. Visualization empowered students to better articulate their idea and communicate important insights about their data to their supervisors and peers. Flexibility in selecting tools for creating the visualization ensured that students could learn and apply visualization techniques regardless of programming knowledge. The assessment focused on evaluating the understanding of visualization techniques. However, it was difficult to deal with students switching between platforms halfway through the semester? From my perspective, students did not have a full understanding of the difficulty of learning to program or use a new (visualization) tool.

### References

- 1 Uzma Haque Syeda, Prasanth Murali, Lisa Roe, Becca Berkey, and Michelle A. Borkin Design Study “Lite” Methodology: Expediting Design Studies and Enabling the Synergy of Visualization Pedagogy and Social Good. *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, 1–13. doi.org/10.1145/3313831.3376829.

## 5.9 Making with Data – Using an open structured template to document the creation of physical data objects

*Till Nagel (Hochschule Mannheim, DE, t.nagel@hs-mannheim.de)*

License  Creative Commons BY 4.0 International license  
© Till Nagel

*Making with Data* brings together a series of practical examples that highlight the diverse range of different ways in which people create physical data objects, showcasing the myriad considerations and decisions that are required to translate data into physical form. With this, our book introduces physicalization to a broad audience of learners, educators, makers, and researchers. Rather than illustrating one correct approach, the collection showcases the many ways in which people today are making with data – in the hope that these processes might inspire readers to make something new.

We started our collection process by interviewing participants at the Dagstuhl Seminar #18441 which informed the general direction of documenting practices and processes. Next, we created a template and over the years iteratively refined it by asking for a broader set of descriptive metadata, specifying the expected text length for each section, and most importantly giving more explicit prompts. We also asked the authors to document their projects in a very rich and visual way by providing high-resolution images from all steps of the creation process.

After we clustered the submissions into five thematic sections, we invited academic experts to write introductions for each in which they provide a personal and unique take on the value of creating physicalizations and help anchor the act of making with data in a different set of artistic, technical, and social practices.

Our approach for documenting the ideation and construction process of physical data objects can be adapted to related fields. We can imagine this open-structured template working similarly well for other forms of visualization creation. Furthermore, it resulted in a detailed dataset collecting diverse approaches to creation providing a range of opportunities for research, analysis, and sharing.

*Making with Data* is edited by Samuel Huron, Till Nagel, Lora Oehlberg, and Wesley Willet. The book will be published in fall 2022 as part of Routledge’s AK Peters Visualization Series edited by Tamara Munzner and Alberto Cairo.

## 5.10 Visual Robot Glyphs

*Jason Dykes (University City of London, London, GB)*

License  Creative Commons BY 4.0 International license  
© Jason Dykes

I showed some data robots (9). These multi-channel glyphs show characteristics of the names of a class of students, or the participants in a Dagstuhl Seminar. They are good for getting people to think about visual channels – how we can encode, the people behind the data – who are we representing, and designs that do not work – encoding is not enough. They also help introduce some issues associated with the ethics of visualization.



■ **Figure 9** Robot glyphs: visually encoding attributes of underlying data.



■ **Figure 10** Sketches to VISualization Learning Outcomings (VISLOs).

## 5.11 Reflective on VISualization Learning Outcomings (VISLOs)

Jason Dykes (University City of London, GB)

License © Creative Commons BY 4.0 International license  
© Jason Dykes

This talk was on intended module learning outcomes and summarized some experiences of using approaches introduced by Jenny Moon in the 2000s. See Moon (2004). The key message, leading to VISualization Learning Outcomings (VISLOs, 10) is:



- write base level learning outcomes that everyone must achieve;
- then write aspirational outcomes that you hope your best students will achieve;
- then work out how you can assess these
- then put all of your effort into helping students achieve the outcomes and do well in the assessment through your plan for teaching. #curriculumLast.

I suggest writing three part outcomes that involve: an **action** – **on a thing** – **at a level**. This works well for me, and helps create grading criteria as I have two points of reference. I asked people to log examples of VISLOs here – and was absolutely, totally, shockingly and painfully unsuccessful. But the opportunity remains: <http://bit.ly/dagVISLO>

## 6 Overview of Activities

### 6.1 Sketching Introductions: An ice-breaker

*Tatiana Losev (Simon Fraser University, Vancouver, CA, Tatiana.Losev@sfu.ca)*

License  Creative Commons BY 4.0 International license  
 Tatiana Losev

On the first morning of the seminar, I facilitated Sketching Introductions, a 90-minute icebreaker activity that invites people to make personal visualizations through simple drawing. I invited everyone in the group to draw a quick sketch, then introduce themselves using their sketch – people drew with pens or colored pencils on paper to sketch their responses to the question, “How do I see myself in relation to the topic Visualization Empowerment: How to Teach and Learn Data Visualization?”

The group sketched for 7 minutes accompanied by background music, then everyone introduced their sketches in 1-minute introductions. The remote attendees presented their sketches on a shared digital whiteboard via videoconferencing. The in-person attendees projected their sketches to both the in-person and remote attendees. Though some people could not finish their introduction in 1 minute and required more time, everyone completed the activity, and the sketches were as distinct as the personal experiences that they depicted. Each introduction proposed personal approaches to teaching and learning in the VIS context. This icebreaker was a creative visualization activity that enabled group members to share and be introduced to the many perspectives to teaching and learning.

### 6.2 Visualization Futures Cards

*Wesley Willett (University of Calgary, CA, wesley.willett@ucalgary.ca)*

License  Creative Commons BY 4.0 International license  
 Wesley Willett

This Vis Futures activity demonstrated a sketching exercise that uses design fiction, collaboration, and creative ideation to encourage players to envision opportunities, use cases, and designs for future visualizations.

The activity uses a set of themed playing card prompts, which emerged from a 2020 workshop on Vis Futures <sup>2</sup> at IEEE VIS. Using the Situation Lab’s *The Thing From the Future* <sup>3</sup> (a similar sketching game focused on more general future ideation) as a template, the attendees a set of roughly 40 attendees spent several hours proposing, designing, and playtesting a diverse set of different visualization-specific design prompts and cards. Both the original workshop and the resulting game were designed to encourage the use of design futuring to envision the next generation of vis tools and applications.

Since the conclusion of the Vis Futures workshop in 2020, a team of collaborators at the University of Calgary, University of Victoria, and Simon Fraser University have collaborated to develop and refine decks of playable cards to support the activity. These include an online version of the game developed at the University of Victoria <sup>4</sup> as well as a physical card deck currently under production at Calgary.

---

<sup>2</sup> <https://visfutures.github.io/>

<sup>3</sup> <http://situationlab.org/project/the-thing-from-the-future/>

<sup>4</sup> <https://observablehq.com/d/51a981cf418ab2ac>



■ **Figure 11** Left: A set of visualization futures cards. Middle: Attendees sketching possible future visualizations. Right: Attendees share and discuss future visualization designs and their implications.

During the Dagstuhl Seminar, attendees participated in a play test using an in-progress version of the physical cards. In this version, players form teams of 3-5 players. A dealer then composes a sketching prompt by dealing one card from each of four decks – *Audience*, *Data Type*, *Data Characteristics*, *Utopia/Dystopia*. For example the prompts dealt at the beginning of the Dagstuhl activity included “Student” (Audience), “3D” (Data Type), “Cherry Picked” (Data Characteristics), and “Dystopian Football” (Dystopia). Players then have 5 minutes to independently imagine and sketch possible visualization designs based on the prompts. Afterwards, players share and discuss their designs.

At the end of the short session attendees shared a number of their creative designs and reflections, and offered a variety of constructive suggestions for adapting the gameplay to different audiences and settings. Attendees also voiced considerable enthusiasm for the arrival of the complete, playable card game in Fall 2022.

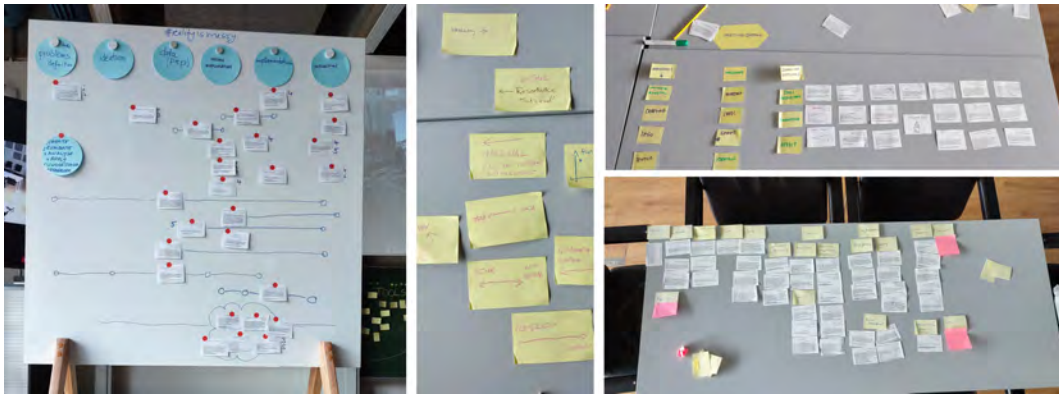
### 6.3 Classifying teaching activity in a design space

*Samuel Huron (Institut Polytechnique de Paris, i3 – Palaiseau, France, samuel.huron@enst.fr)*

License © Creative Commons BY 4.0 International license  
© Samuel Huron

As information visualization educators and teachers, we are all using a variety of activities to introduce, teach, familiarize students and other audiences with concepts relative to information visualization. Since now we have a poor overview of these types of information visualization activities teacher and workshop organizer are using. We wanted to reflect on these practices with the Dagstuhl Seminar participant in order, to have a better idea of the design space, of what are the meaning full categories to describe this space, what have been deeply explored, what has not been explored, and last what could be generated.

We collected 42 activities from more than 13 different authors gathered from three different sources: 1) the seminar participant and 2) IEEE VIS – VIS activities workshop 2020 [1], 2021 [1] and 3) few other ones we know. On this basis we generate a spreadsheet listing all



■ **Figure 12** A white board containing a design space of data visualization teaching activities.

these activities and created a card deck in which all activities are represented by one card. This card includes (title, author, description, keywords and URL to a document describing it). Each working group received a deck of 42 cards and the link to the spreadsheet.

The activity happened in two main steps 1) design space creation, 2) presentation & discussion. The prompts for the first step were the following “Categorize these activities in a design space from the focus of your group.” It will be in the room for local participants and on a Miro Board <sup>5</sup> for remote participants. The prompt for the second step was “Present the design space and the dimensions, maybe the case that was problematic, and the rationale behind each dimensions”. They were asked that it would be totally fine to remove some cards if they do not fit the focus of their group, or even complete the card with activities that were not described in the data set by using a post-it notes or other papers. Last the participant was asked to use the cards as tokens and look for more details in the spreadsheet (images, descriptions, paper, links).

The activity last one hour and 30 minutes, the participant spend 40 minutes to create one design design space by the five groups, and we spent the rest of the time for presenting theses spaces and discussing them. Each group have been able to create meaningful dimensions.

## References

- 1 Huron, S., Bach, B., Panagiotidou, G., Keck, M., & Roberts, J., Carpendale (2021, October). 2nd IEEE VIS Workshop on Data Vis Activities to Facilitate Learning, Reflecting, Discussing and Designing In IEEE VIS 2021.
- 2 Huron, S., Bach, B., Hinrichs, U., Keck, M., & Roberts, J. (2020, October). IEEE VIS Workshop on Data Vis Activities to Facilitate Learning, Reflecting, Discussing, and Designing. In IEEE VIS 2020.

<sup>5</sup> [https://miro.com/app/board/uXjV0p4UJmI=?share\\_link\\_id=161478137131](https://miro.com/app/board/uXjV0p4UJmI=?share_link_id=161478137131)



■ **Figure 13** After-hours music session.

## 6.4 Activity : Cognitive Science of Representational Systems

*Peter Cheng (University of Sussex, GB, p.c.h.cheng@sussex.ac.uk)*

License © Creative Commons BY 4.0 International license  
© Peter Cheng

We contend that understanding users' interpretations of visualizations is essential for anyone who wishes to teach about visualizations or to design visualizations for learning. Interpretations are the memory structures that users build as they read and interact with visualizations. This tutorial activity introduces an approach to modeling interpretations based on Representation Interpretive Structure Theory (RIST), which claims that interpretations of a representations depend upon four types of cognitive schemes and a small number of relations among them. Participants in the activity will learn the graphical notation (RISN) for building network models of such interpretations using a web-based graphical editor (RISE). We imagine that interpretation models may be built for many purposes. Instructors may construct a model of the conceptual structure of a representation in order to devise better explanations for learners of how a representation works. A visualization designer may build models to explore the consequences of expressiveness and cognitive demands of alternative visualization formats for a particular dataset. For researchers specializing in visualization, the approach potentially provides a coherent and rigorous approach for comparisons of representations across graphical formats and knowledge domains.

## 7 Summary

The week was an extraordinary energetic moment of encounters and intense discussion. Data visualization teaching and learning is an emerging domain that will need proper addressing in the years and decades to come. The Dagstuhl Seminar gave us the opportunity to place this topic onto the map and to create an early community with a strong agenda that will be remembered by the participants and organizers. The participants of the seminar generated a myriad of possible outcomes including books, scientific papers, workshop papers, online platforms, grant collaborations, a dedicated symposium proposal at our main conference IEEE VIS, and a potential follow up seminar in a few years time.

We thank Dagstuhl and its staff for providing the stage and the services in which incredible moment of fruitful collaboration can happen.

## Participants

- Jan Aerts  
Amador Bioscience – Hasselt,  
Hasselt University &  
KU Leuven, BE
- Fearn Bishop  
BBC – Salford, GB
- Peter C.-H. Cheng  
University of Sussex –  
Brighton, GB
- Alexandra Diehl  
University of Zurich, CH
- Jason Dykes  
City – University of London, GB
- Sarah Hayes  
Munster Technological University  
– Cork, IE
- Uta Hinrichs  
University of Edinburgh, GB
- Trevor Hogan  
Munster Technological University  
– Cork, IE
- Christoph Huber  
Mannheim University of Applied  
Sciences, DE
- Samuel Huron  
Institut Polytechnique de  
Paris, FR
- Mandy Keck  
Univ. of Applied Sciences –  
Hagenberg, AT
- Christoph Kinkeldey  
HAW – Hamburg, DE
- Søren Knudsen  
IT University of  
Copenhagen, DK
- Doris Kosminsky  
University of Rio de Janeiro, BR
- Tatiana Losev  
Simon Fraser University –  
Burnaby, CA
- Areti Manataki  
University of St Andrews, GB
- Isabel Meirelles  
The Ontario College of Art and  
Design University, CA
- Luiz Morais  
INRIA – Bordeaux, FR
- Till Nagel  
Mannheim University of Applied  
Sciences, DE
- Rebecca Noonan  
Munster Technological University  
– Cork, IE
- Georgia Panagiotidou  
University College London, GB
- Laura Pelchmann  
University of Cologne, DE
- Fateme Rajabiyazdi  
Carleton University –  
Ottawa, CA
- Jonathan C. Roberts  
Bangor University, GB
- Christina Stoiber  
FH – St. Pölten, AT
- Yagoda Walny  
Canada Energy Regulator –  
Calgary, CA
- Wesley J Willett  
University of Calgary, CA





## Remote Participants

- Wolfgang Aigner  
FH – St. Pölten, AT
- Benjamin Bach  
University of Edinburgh, GB
- Magdalena Boucher  
FH – St. Pölten, AT
- Robert S. Laramée  
University of Nottingham, GB
- Andrew Manches  
University of Edinburgh, GB
- Alison Powell  
London School of Economics, GB
- Mashaël Alkadi  
University of Edinburgh, GB
- Mine Çetinkaya-Rundel  
DGBe University – Durham, US
- Andy Kirk  
Visualising Data – Leeds, GB
- Fanny Chevalier  
University of Toronto, CA
- Nathalie Henry Riche  
Microsoft Research –  
Redmond, US
- Dietmar Offenhuber  
Northeastern University –  
Boston, US
- Sheelagh Carpendale  
Simon Fraser University –  
Burnaby, CA
- Marti Hearst  
University of California –  
Berkeley, US
- Charles Perin  
University of Victoria, CA
- Emeline Brulé  
University of Sussex –  
Brighton, GB



# Human-Centered Artificial Intelligence

Wendy E. Mackay<sup>\*1</sup>, John Shawe-Taylor<sup>\*2</sup>, and Frank van Harmelen<sup>\*3</sup>

1 INRIA Saclay - Orsay, FR. [wendy.mackay@inria.fr](mailto:wendy.mackay@inria.fr)

2 University College London, GB. [jst@cs.ucl.ac.uk](mailto:jst@cs.ucl.ac.uk)

3 VU University Amsterdam, NL. [frank.van.harmelen@vu.nl](mailto:frank.van.harmelen@vu.nl)

---

## Abstract

This report documents the program and the outcomes of Dagstuhl Perspectives Workshop 22262 “Human-Centered Artificial Intelligence”.

The goal of this Dagstuhl Perspectives Workshops is to provide the scientific and technological foundations for designing and deploying hybrid human-centered AI systems that work in partnership with human beings and that enhance human capabilities rather than replace human intelligence. Fundamentally new solutions are needed for core research problems in AI and human-computer interaction (HCI), especially to help people understand actions recommended or performed by AI systems and to facilitate meaningful interaction between humans and AI systems. Specific challenges include: learning complex world models; building effective and explainable machine learning systems; developing human-controllable intelligent systems; adapting AI systems to dynamic, open-ended real-world environments (in particular robots and autonomous systems); achieving in-depth understanding of humans and complex social contexts; and enabling self-reflection within AI systems.

**Seminar** June 26–July 1, 2022 – <http://www.dagstuhl.de/22262>

**2012 ACM Subject Classification** Human-centered computing → Human computer interaction (HCI); Human-centered computing → Interaction design; Mathematics of computing → Probability and statistics; Theory of computation → Probabilistic computation; Theory of computation → Automated reasoning; Theory of computation → Constraint and logic programming; Theory of computation → Machine learning theory; Theory of computation → Algorithmic game theory and mechanism design; Computing methodologies → Artificial intelligence; Computing methodologies → Machine learning

**Keywords and phrases** Human-centered Artificial Intelligence, Human-Computer Interaction, Hybrid Intelligence

**Digital Object Identifier** 10.4230/DagRep.12.6.112

## 1 Executive Summary

*Wendy E. Mackay*

*John Shawe-Taylor*

*Frank van Harmelen*

This workshop brought together 22 participants with a diverse background from AI, Robotics, HCI, Ubiquitous Computing, Business and Sociology, from across Europe and North America laid the groundwork for a manifesto on Hybrid Human-centered AI systems.

Informed by currently ongoing large initiatives such as the EU-funded Humane AI Net, the Dutch Hybrid Intelligence Center, the Danish Centre for Hybrid Intelligence, and OECD AI policy framework, four pillars of the manifesto emerged: (a) Collaboration and Cooperation,

---

\* Editor / Organizer



Except where otherwise noted, content of this report is licensed under a Creative Commons BY 4.0 International license

Human-Centered Artificial Intelligence, *Dagstuhl Reports*, Vol. 12, Issue 6, pp. 112–117

Editors: Wendy E. Mackay, John Shawe-Taylor, and Frank van Harmelen



DAGSTUHL  
REPORTS

Dagstuhl Reports  
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

(b) Control & Adaptivity, (c) Transparency & Explainability, and (d) Societal dimensions. For each of these pillars, the workshop resulted in (i) key terminology, (ii) key research questions, (iii) metrics and methodologies, and (iv) benchmarks and challenges.

The above resulted in a solid framework for the Hybrid Human-Centered AI manifesto to be written by the participants in the months following the Dagstuhl workshop.

**2 Table of Contents**

**Executive Summary**  
*Wendy E. Mackay, John Shawe-Taylor, and Frank van Harmelen* . . . . . 112

**Overview of Talks** . . . . . 115

**Plenary discussions** . . . . . 115

**Working groups** . . . . . 115

**Participants** . . . . . 117

### 3 Overview of Talks

Wendy Mackay presented the **Human-Computer Interaction approach** to taking a human perspective when interacting with an intelligent system. She explained that the focus is not on the human alone, or the system alone, but rather the interaction between them, where interaction is treated as a phenomenon in its own right that can be designed, studied and controlled. She explained the cycle in the natural and social sciences between theoretical and empirical research, and how to incorporate human-designed interactive systems into the research process. She briefly described generative theory for creating novel interactive systems, as well as the possible relationships between users and interactive systems, and how they blur when designing human-centered AI systems. Finally, she described a number of human-centered participatory design methods, and the importance of prototyping how humans will interact with a proposed intelligent system, including explicitly identifying, illustrating and reflecting on possible breakdowns.

**Marko Grobelnik presented the OECD AI policy framework.** He discussed issues of human vs. AI autonomy, and questioned whether the OECD framework is too reductionist, placing the human within a machine framework, rather than the other way around. He also discussed the AI system lifecycle and AI system classification.

**Frank van Harmelen presented the research agenda for the Dutch Hybrid Intelligence Center,** with 80 Ph.D. students over 10 years (slides here). He explained the project's CARE goals: Collaborative (synergy with humans), Adaptive (adapt to humans and the environment), Responsible (ethical performance), and Explainable (share and explain awareness, goals and strategies).

**Paul Lukowicz presented the HumanE-AI European research network.**

**Paul Lukowicz presented two overall frameworks.** The first framework includes: Communication; Oversight/Control; and Frame of reference. The second organizes the Outcome space in terms of physical effect, impact on humans, and impact on society.

### 4 Plenary discussions

After the opening sessions on Monday morning, each participant gave a three-minute presentation of their background, goals for the workshop, and key issues they would like to address in the workshop. The group is highly diverse, with representatives from AI, Robotics, HCI, Ubiquitous Computing, Business, Sociology, from across Europe and North America.

In the Monday afternoon session, the group identified six major issues to address, which were compressed into four topics for breakout groups: Collaboration and Cooperation, Control & Adaptivity, Transparency & Explainability, and Societal dimensions, to be discussed in greater detail in four breakout sessions during the week. Each breakout session was encouraged to discuss the following topics for the manifesto: Key Terminology; Research Questions; Methods, Competences and Frameworks; Benchmarks, Moonshots and Challenges; and Policy.

### 5 Working groups

In working groups on Monday afternoon, the participants addressed the following assignment:

1. Choose a target user and a context where a user or (users) interact with an intelligent system, based on one of the four morning session topics;
2. Develop a scenario and describe prototypes (on paper, or otherwise) to illustrate the interaction;
3. Shoot a video prototype of a realistic scenario that illustrates how the user(s) interact with the system, including at least one situation where the system breaks down.

Each of the five groups then presented and discussed their video and the challenge presented by the breakdown.

In working groups on Tuesday afternoon, participants discussed the issues related to Collaboration and Communication in Human-centered AI. At the end of the afternoon, participants met in a plenary session to discuss the results of each breakout session.

In working groups on Wednesday morning, participants formed four breakout groups with specific examples related to the topic of Transparency and Explainability. At the end of the morning, participants met in a plenary session to discuss the results of each breakout session.

In working groups on Thursday morning, participants formed four breakout groups to discuss specific examples with respect to the dimensions outlined in Paul Lukowicz' presentation (mentioned above). At the end of the morning, participants met in a plenary session to discuss the results of each breakout session.

In working groups on Thursday afternoon, participants formed three breakout groups, one focused on policy issues, the remaining three groups worked on writing the following:

1. Definition of HCAI,
2. List of provocative statements for the manifesto,
3. Identification of leading examples
4. Suggested Ph.D. topics

In a plenary session on Friday morning, the group discussed creating four concrete examples that illustrate the design space of human-centric AI and settled on four each with an assigned author: Interactive collaborative music, Social Media, Decision support for Health, and Crisis Management.

After organising the collected breakout notes into the key components of the manifesto (terminology, research questions, methods and metrics, benchmarks and moonshots, and policy), the participants then broke up into working groups and drafted the key sections of the manifesto

In a final plenary session the group then discussed the results and future directions, including transforming the manifesto into a 10-20 page journal article, with scientists as the target audience, contributing to the HHAI conference in Munich next year, and establishing a Human-centered AI Master's program.

## Participants

- Michel Beaudouin-Lafon  
University Paris-Saclay – Orsay, FR
- Mehul Bhatt  
University of Örebro, SE
- Stefan Buijsman  
TU Delft, NL
- Mohamed Chetouani  
Sorbonne University – Paris, FR
- Ulises Cortés  
UPC Barcelona Tech, ES
- Adam Dahlgren Lindström  
University of Umeå, SE
- Emmanuelle Dietz  
Airbus – Hamburg, DE
- Marko Grobelnik  
Jozef Stefan Institute – Ljubljana, SI
- Alípio Jorge  
University of Porto & INESC TEC – Porto, PT
- Antonis C. Kakas  
University of Cyprus – Nicosia, CY
- Samuel Kaski  
Aalto University, FI
- Janin Koch  
INRIA Saclay – Orsay, FR
- Helena Lindgren  
University of Umeå, SE
- Paul Lukowicz  
DFKI – Kaiserslautern, DE
- Wendy E. Mackay  
INRIA Saclay – Orsay, FR
- Ozlem Ozmen Garibay  
University of Central Florida – Orlando, US
- Janet Rafner  
Aarhus University, DK
- Laura Sartori  
University of Bologna, IT
- John Shawe-Taylor  
University College London, GB
- Jacob Sherson  
Aarhus University, DK
- Marija Slavkovic  
University of Bergen, NO
- Philipp Slusallek  
DFKI – Saarbrücken, DE
- Frank van Harmelen  
VU University Amsterdam, NL
- Katharina A. Zweig  
TU Kaiserslautern, DE

