

---

# 000 POST-AGI SCIENCE AND SOCIETY WORKSHOP

001

002

003

004

## 005 ABSTRACT

006

007

008

009

010

011

012

013

014

015

016

017

018 Artificial General Intelligence (AGI) has long seemed distant, but rapid advances  
019 in large-scale learning, autonomous reasoning, and open-ended discovery make  
020 its emergence increasingly plausible. The Post-AGI Science and Society Work-  
021 shop asks what comes next. If AGI becomes ubiquitous, reliable, and affordable,  
022 how will it reshape scientific inquiry, the economy of knowledge, and human  
023 society? Will humans remain central to discovery or become curators and inter-  
024 preters of machine-generated insights? The workshop brings together researchers  
025 from machine learning, philosophy of science, and policy to explore human-AI  
026 scientific coexistence. Topics include automated hypothesis generation, causal  
027 reasoning in AGI, collaborative discovery, epistemic alignment between humans  
028 and machines, and socio-economic shifts driven by pervasive intelligence. Through  
029 keynotes, talks, and a panel, we will examine how science and our understanding  
030 of knowledge might evolve in a post-AGI world.

031

[032 https://p-agi.netlify.app/](https://p-agi.netlify.app/)

033

034 **Workshop Summary** Never before has there been such a growing sentiment among researchers  
035 that artificial general intelligence (AGI), an AI system capable of performing most tasks at or beyond  
036 human level (Morris et al., 2023), may be within reach in the foreseeable future. Although uncertainty  
037 remains about precise timelines, an important shift has taken place: debates are no longer centered on  
038 whether AGI is possible, but increasingly on how to define it and what its implications might be for  
039 science, society, and culture.

040

041 This view is no longer limited to speculative debate; it is increasingly reflected in the expectations of  
042 the scientific community. A significant fraction of AI researchers estimate a 50% chance that unaided  
043 machines will surpass humans in every task by 2047 (Grace et al., 2024). The central disagreements  
044 are shifting away from whether AGI is possible to when it will emerge and what its implications will  
045 be for science and society.

046

047 Given this accelerating progress, one of the most critical and unresolved questions is how AGI will  
048 reshape the scientific process itself. When machines can autonomously formulate hypotheses, design  
049 experiments, and generate discoveries, what research remains for humans to do?

050

051 This workshop builds on that question. Rather than debating feasibility, we invite participants to  
052 explore how to make today's scientific research directions resilient and meaningful in the event that  
053 AGI becomes a pervasive tool. Our purpose is not to endorse specific timelines but to ask: **Assuming**  
054 **AGI becomes widely available, what is the future of science?** As scientists and innovators, it is our  
055 responsibility to address this question proactively, ensuring that as discovery becomes automated, it  
056 remains trustworthy, safe, and beneficial for humanity.

057

058 To that end, we bring together leading voices from AI, robotics, and safety to explore the following  
059 questions:

060

- 061 • What are the true limits of reasoning-capable systems, and could their capabilities scale  
062 beyond our control? We will hear from Noam Brown, co-creator of ChatGPT o1, one of the  
063 most advanced models to date.
- 064 • As these systems grow more capable, how will the scientific process evolve? Will machines  
065 generate discoveries for humans to validate, or will humans generate questions for machines  
066 to resolve? Jeff Clune and Francesco Locatello will examine this shift from complementary  
067 angles: the former through the lens of open-ended learning and automated discovery, and  
068 the latter through the role of causality in shaping how scientists will formulate and interpret  
069 causal questions in the age of AGI.
- 070 • As robotics and AI become more pervasive, so do challenges in adaptability, efficiency,  
071 and real-world integration. Daniela Rus will discuss how autonomous, reconfigurable, and

054 scalable robotic systems can advance energy-efficient, intelligent technologies that enhance  
055 daily life.  
056

- 057 • Finally, what role will humans play? Been Kim, a leader in interpretability and human-  
058 machine understanding, will explore how to build systems that remain cooperative and  
059 enable meaningful human-AI collaboration.

060 We invite submissions of original research that explore the future of science, technology, and society  
061 in an age where Artificial General Intelligence is pervasive. Our goal is to foster a discussion that  
062 looks beyond immediate technical hurdles to identify and address the fundamental questions that will  
063 remain relevant as general-purpose AI systems become ubiquitous.

064 Submissions are encouraged across two primary tracks: **Track1 - Technical Foundations for a**  
065 **Post-AGI World** and **Track2 - Socio-Economical and Future Visions**. All the submissions will be  
066 in the "Tiny Papers" format.  
067

068 **Track 1: Technical Foundations for a Post-AGI World.** This track focuses on the core technical  
069 challenges required to build, understand, and control highly capable AI systems. We are interested  
070 in work that addresses the safety, robustness, and scalability of models that approach and surpass  
071 human-level intelligence.

072 Suggested topics include, but are not limited to:  
073

- 074 • **Automated Scientific Discovery:**

- 075 – AI systems and algorithms designed to automate research, generate novel hypotheses,  
076 and accelerate scientific discovery. This includes areas like symbolic-neural theorem  
077 proving (Selsam et al., 2019; Bansal et al., 2021) and AI-guided program search (Shi  
078 et al., 2023; Rule et al., 2024).
- 079 – Explorations into the future of the scientific process when machines can independently  
080 design, execute, and interpret experiments (Hu et al., 2024; Mankowitz et al., 2023;  
081 Zhang et al., 2025).
- 082 – Research on “Deep Research” paradigms where AI drives discovery, including program  
083 synthesis and library learning (Bowers et al., 2023; Ellis et al., 2021).

- 084 • **Scalable and Efficient Intelligence:**

- 085 – Hardware-AI co-design to drive scalable, energy-efficient intelligence (Mirhoseini  
086 et al., 2021; 2020; Fawzi et al., 2022).
- 087 – Research into resource-aware scaling and the fundamental limits of reasoning-capable  
088 systems (Costello et al., 2025; Shao et al., 2024).
- 089 – Methods for building and deploying powerful AI systems with reduced computational,  
090 financial, and environmental costs (Schwartz et al., 2020).

- 091 • **Safety, Robustness, and Alignment:**

- 092 – Novel techniques for aligning powerful AI systems with human values and intentions,  
093 including next-generation preference learning (Lee et al., 2024; Rafailov et al., 2023;  
094 Meng et al., 2024) and broader human-AI alignment strategies (Ellis, 2023; Li et al.,  
095 2025).
- 096 – Methods for ensuring safety and robustness, including superalignment and scalable  
097 oversight (Burns et al., 2024; Bowman et al., 2022; Belrose et al., 2023).
- 098 – Frameworks and benchmarks for rigorously evaluating AGI-level reasoning (Chollet  
099 et al., 2025a;b; Li et al., 2024).
- 100 – Advanced methods for transparency and human-AI collaboration, enabling meaningful  
101 human oversight (Tang et al., 2024; Yang & Deng, 2019).

104 **Track 2: Socio-Economical and Future Visions** This track invites contributions that analyze  
105 the broader societal, ethical, and economic implications of AGI. We encourage speculative works,  
106 position papers, and in-depth studies that grapple with the profound transitions society will face.  
107

Suggested Topics include, but are not limited to:

---

108 • **Economic and Societal Impact:**

109 – Studies analyzing AGI’s impact on key domains such as economics, law, education,  
110 and healthcare (Lee et al., 2025; Dunlop et al., 2024).

111 – Analyses of labor-market disruptions, the future of work, and the role of humans in a  
112 world with ubiquitous AGI (Shao et al., 2025; Tomlinson et al., 2025).

113 – Frameworks for evaluating the societal effects of AI, such as “GDPval” and other  
114 systemic evaluations (Patwardhan et al., 2025).

115 • **Governance, Regulation, and Risk:**

116 – Proposals for the regulation and governance of powerful, general-purpose AI systems  
117 (Raji et al., 2022; Nolte et al., 2025).

118 – Comprehensive risk assessments, including the study of existential risks and pathways  
119 from AGI to Artificial Superintelligence (ASI) (Anderljung et al., 2023; Green, 2020).

120 – Position papers on how society can best govern and benefit from systems that reshape  
121 every field (Bengio et al., 2023).

122 • **Foundational Questions and Position Papers:**

123 – Analyses and positions of long-term AI trajectories, focusing on potential further  
124 developments and their technical and societal implications (Morris et al., 2023; Gabriel,  
125 2020; Wen et al., 2025).

126 – Studies and position papers examining the theoretical and computational aspects of  
127 machine consciousness, models of computational creativity, and the evolving role of AI  
128 in creative practices and art production (Kamb & Ganguli, 2025; Blili-Hamelin et al.,  
129 2025).

130 – Ethical frameworks for the deployment and use of AGI for social good (Prabhakaran  
131 et al., 2022; Gabriel et al., 2024).

132

133 P-AGI invites the community to step beyond immediate technical challenges and to define the next set  
134 of fundamental questions that will matter in a world where general intelligence is no longer rare. We  
135 believe that by hosting this discussion now, ICLR will reaffirm its leadership not only in advancing  
136 the field of machine learning, but in preparing it for its most profound transitions.

137 **Workshop Format** P-AGI is structured to foster deep discussion and cross-disciplinary collabora-  
138 tion. The program will integrate three main elements: *invited talks*, a *panel discussion*, and a *poster*  
139 *session*.

140 The *invited talks* will bring together leaders from frontier AI, robotics, AI safety, philosophy, and  
141 social sciences. Each speaker will present a vision or provocation addressing the scientific, societal,  
142 or creative challenges in a world where AGI is widespread.

143 The *panel discussion* will gather experts from across disciplines to debate key questions raised  
144 during the workshop: Which research problems will endure? How will society adapt? What forms  
145 of creativity and governance will be most important? The panel will be guided by live audience  
146 questions to keep the discussion responsive and grounded.

147 The *poster session* will showcase work submitted to the two workshop tracks. We will organize the  
148 poster session to maximize interaction, including assigning experienced attendees to engage with  
149 selected posters. This will ensure that junior researchers receive substantive feedback and encourage  
150 rich discussions around forward-looking ideas. Extended coffee breaks and lunch will be designed to  
151 encourage informal conversations and collaborations across fields. Our goal is to create not just a  
152 series of presentations, but an environment where bold ideas about the future of research can emerge  
153 and take root.

154 The organizers’ prior experience running successful workshops at top venues will be instrumental in  
155 coordinating a smooth, high-impact event.

156 **Points of difference** Unlike workshops focused on advancing the current state of AI, P-AGI is  
157 designed to prepare the research community for the downstream effects of general intelligence. Our  
158 unique contribution is to create a dedicated space where researchers can shape a scientific agenda  
159 that is resilient, relevant, and responsible in the face of this transformative change. We will achieve  
160 this through three core objectives:

---

Tentative Schedule					
162	163	09:15	Opening Remarks	13:00	Lunch Break
164	165	09:30	Invited Talk 1	14:00	Invited Talk 4
166	167	10:00	Invited Talk 2	14:30	Invited Talk 5
168	169	10:30	Invited Talk 3	15:00	Invited Talk 6
170	171	11:00	Contributed Talks	15:30	Panel Discussion
172	173	11:30	Poster Session	16:30	Closing Remarks

---

- **Fostering Cross-Disciplinary Discussion:** The workshop is structured to bring together leading voices from diverse fields like AI, robotics, safety, and social sciences. Through invited talks, a cross-disciplinary panel, and interactive poster sessions, we will facilitate the deep discussions needed to tackle the multifaceted challenges of a post-AGI world.
- **Future-Proofing the Research Agenda:** Rather than pushing the immediate frontier, we seek to **future-proof the questions we ask**. By considering the implications of widely available AGI, our goal is to help the community identify and prioritize research directions that will hold fundamental value under a broad range of futures, ensuring today's work remains meaningful.
- **Defining the Next Set of Questions:** We aim to challenge participants to look beyond immediate technical hurdles and collectively **define the next set of fundamental questions** that will matter in a world where general intelligence is no longer rare. By hosting this forward-looking dialogue, P-AGI will help ensure the field is proactively and thoughtfully preparing for its most profound transitions.

187 To our knowledge, no other workshop takes this long-term perspective, placing future relevance at  
188 the core of its mission.

189  
190 **Diversity and Inclusivity** We believe that diversity and inclusivity are essential to fostering  
191 a forward-thinking research community, especially when confronting the broad implications of  
192 AGI. P-AGI is committed to welcoming participants from all backgrounds and identities, and this  
193 commitment is reflected across every aspect of the event. We have prioritized diversity in seniority,  
194 gender, nationality, and institutional affiliation. Our organizing team and invited speakers span PhD  
195 students, early-career researchers, and senior academics, with an intentional 50/50 gender balance  
196 in both groups. The committee also represents a wide geographic spread, including Europe, North  
197 America, and Asia. To further broaden participation, we will actively engage affinity groups such as  
198 Black in AI, Women in Machine Learning (WiML), Queer in AI, and Latinx in AI, distributing calls  
199 for program committee roles and general participation among these corresponding networks. Finally,  
200 pending sponsorship, we intend to offer travel support for members of underrepresented groups to  
201 help reduce financial barriers to attendance. Through these efforts, we aim to cultivate an inclusive,  
202 interdisciplinary environment where all voices are heard, supported, and empowered to shape the  
203 future of AGI research.

203  
204 **Attendance** We expect a vibrant and diverse audience, with over 150 in-person attendees and an  
205 additional 200 joining virtually. This strong interest is driven by the workshop's themes, which  
206 are widely felt across the AI community yet often confined to social media threads and informal  
207 conversations. P-AGI offers a space to bring these discussions into the open, with the rigor, depth, and  
208 interdisciplinary engagement they deserve. Participants will include students, academic researchers,  
209 and industry professionals from a wide range of backgrounds and cultures. To encourage interaction  
210 and community building, all attendees will have access to an open-source communication workspace.  
211 This shared space allows newcomers to introduce themselves, share their work, and stay informed  
212 about upcoming opportunities, events, and research highlights. A dedicated social channel will  
213 support informal meetups and casual conversations, helping participants connect beyond the formal  
214 sessions. Our goal is to create a dynamic, inclusive environment where everyone can contribute,  
215 learn, and remain engaged long after the workshop concludes.

215 To maximize the workshop's reach and impact for a global audience, all talks and presentations  
will be recorded (pending speaker permission) and made publicly available on the official workshop

216 website. This digital archive will also feature a complete collection of the accepted papers and posters,  
217 ensuring that the content remains accessible to the entire research community long after the event.  
218

219 **INVITED SPEAKERS AND PANELISTS**  
220

221 **Noam Brown** *Research Scientist, OpenAI* [Confirmed]  
222 **Noam Brown** leads the development of advanced reasoning models at OpenAI, including o1,  
223 designed for general problem-solving in math, science, and code. He was named one of MIT's 35  
224 Innovators Under 35 and received the Marvin Minsky Medal for his contributions to AI. Having  
225 led work on one of the most capable reasoning AIs to date, Noam offers a rare perspective on the  
226 capabilities, limits, and implications of general-purpose systems.  
227

228 **Jeff Clune** *Full Professor, University of British Columbia* [Confirmed]  
229 **Jeff Clune** is a Full Professor at the University of British Columbia, a Senior Research Advisor  
230 at DeepMind, and a Canada CIFAR AI Chair at the Vector Institute. His research focuses on  
231 open-ended learning, AI-generating algorithms, and deep reinforcement learning – core ideas for  
232 scaling intelligence. As a pioneer in evolving increasingly general and autonomous AI systems, Jeff  
233 brings a visionary yet technically grounded perspective, making him an ideal contributor to the  
234 workshop themes.  
235

236 **Daniela Rus** *Andrew & Erna Viterbi Professor, MIT | Director, MIT CSAIL* [Confirmed]  
237 **Daniela Rus** leads MIT's Computer Science and Artificial Intelligence Laboratory (CSAIL),  
238 advancing robotics and AI for real-world autonomy, from soft robots to modular and swarm systems.  
239 A MacArthur Fellow and elected member of the National Academy of Engineering and National  
240 Academy of Sciences, she has received several awards, including the IEEE Edison Medal and the  
241 John Scott Medal. Her work focuses on scalable, energy-efficient AI–robotics integration that  
242 enhances daily life.  
243

244 **Francesco Locatello** *Assistant Professor, ISTA* [Confirmed]  
245 **Francesco Locatello** is an assistant professor at ISTA and an AI Resident at the Chan Zuckerberg  
246 Initiative. His research focuses on causal representation learning, developing models that understand  
247 cause and effect, adapt to interventions, and generalize beyond fixed conditions. At P-AGI, he  
248 will discuss how AGI may transform the way scientists approach causal questions, exploring how  
249 advances in causal reasoning could redefine discovery and explanation in a post-AGI scientific  
250 landscape.  
251

252 **Been Kim** *Senior Staff Research Scientist, Google DeepMind* [Confirmed]  
253 **Been Kim** is a leading voice in interpretable AI and creator of TCAV, a concept-based explanation  
254 method that won the UNESCO Netexplo Award. Her work focuses on enabling human–machine  
255 collaboration, detecting model failures, and using AI knowledge for human benefit. As humans  
256 and AI systems grow increasingly entangled, her perspective is not only advisable, but essential to  
257 ensuring reliable, understandable, and aligned intelligence.  
258

260 **ORGANIZERS**  
261

263 **Organizational experience** The P-AGI organizing committee brings deep and proven experience  
264 in curating high-impact scientific events. Three of its members were among the founding organizers  
265 of UniReps, one of the most attended workshops at NeurIPS in both 2023 and 2024. Collectively, the  
266 committee has organized a total of more than 20 workshops across major conferences in machine  
267 learning, natural language processing, robotics, and computer vision. This track record reflects not  
268 only logistical capability but also a strong ability to attract diverse, interdisciplinary communities and  
269 sustain meaningful discussion around frontier research questions, precisely the kind of engagement  
P-AGI aims to foster.

---

270 **Emanuele Rodolà**

*Sapienza University of Rome*

271 **Emanuele** is Full Professor of Computer Science at Sapienza University of Rome, where he leads  
272 the GLADIA group of learning and applied AI, funded by an ERC Grant, a FIS Grant, and a  
273 Google Research Award. Previously, he was Assistant and then Associate Professor at Sapienza  
274 (2017-2020), a postdoc at USI Lugano (2016-2017), an Alexander von Humboldt Fellow at TU  
275 Munich (2013-2016), and a JSPS Research Fellow at The University of Tokyo (2013). He is a fellow  
276 of ELLIS and the Young Academy of Europe, has received a number of research prizes, has been  
277 serving in the program and organizing committees of the top rated conferences in computer vision  
278 and machine learning, founded and chaired several successful workshops including UniReps at  
279 NeurIPS 2023 and 2024. His research interests lie at the intersection of representation learning,  
280 model merging, graph / geometric deep learning, language and learning for audio, and has published  
281 more than 170 papers in these areas. Previously, he has organized and lectured at 15 tutorials, and  
282 has co-organized and chaired more than 10 workshops co-located with ECCV, ICCV, and NeurIPS.

283 **Pratyusha Sharma**

*Microsoft & NYU*

284 **Pratyusha Sharma** is a Senior Research Scientist at Microsoft Research and an incoming Assistant  
285 Professor at the Courant Institute and Center for Data Science at New York University. She did her  
286 PhD from MIT in the Computer Science and Artificial Intelligence Lab. She studies the interplay  
287 between language, sequential decision making and intelligence in natural and AI systems. Her  
288 research has also been featured in articles in the New York Times, National Geographic Magazine,  
289 BBC, etc. She was recently a speaker at TED AI and was selected as a Rising Star in EECS, Data  
290 Science, and GenAI. She has previously organized the Language and Reinforcement Learning  
291 Workshop (LaReL) NeurIPS 2022, the workshop on Language and Robotics at Conference on Robot  
292 Learning in 2022 and Social Intelligence in Humans and Robots Workshop, Robotics Science and  
Systems 2023.

293 **Andrea Santilli**

*Independent Researcher*

294 **Andrea Santilli** is a Research Scientist at Nous Research, where he focuses on post-training large  
295 language models (LLMs) to improve their robustness, reliability, and alignment. He holds a PhD in  
296 Computer Science (2025) from Sapienza University of Rome, with a thesis on building “Effective,  
297 Efficient, and Reliable Large Language Models.” Previously, he was a Research Scientist at Apple  
298 MLR (2024) and collaborated with leading institutions like Hugging Face (2021). Despite being early  
299 in his career, Andrea’s work has already had a significant impact: his research has been published at  
300 top-tier conferences (ICLR, ICML, ACL, EMNLP, CVPR, SIGIR, AAAI) and has accumulated over  
301 7,000 citations. In addition to his research contributions, he actively serves on program committees  
302 for major conferences and has received multiple awards, including the prestigious Imminent Research  
303 Grant.

304 **Valentina Pyatkin**

*Allen Institute for AI & University of Washington*

305 **Valentina Pyatkin** is a Postdoctoral Researcher at the Allen Institute for AI and the University of  
306 Washington, working with Prof. Yejin Choi and Prof. Hannaneh Hajishirzi. She completed her  
307 PhD in Computer Science at the NLP lab of Bar Ilan University. Her work has been awarded an  
308 ACL Outstanding Paper Award and the ACL Best Theme Paper Award. Her research focuses on  
309 language model alignment, with a focus on steerability, values, pragmatics and context. Valentina  
310 was a co-organizer of the 2nd and 3rd UnImplicit workshop at NAACL 2022 and EACL 2024, and a  
311 co-organizer of the 2nd and 3rd SoLaR workshops collocated with NeurIPS 2024 and CoLM 2025.

312 **Donato Crisostomi**

*Sapienza University of Rome*

313 **Donato Crisostomi** is an ELLIS PhD student at Sapienza University of Rome and University  
314 of Cambridge, focusing on model merging and representational alignment. He currently leads  
315 the “Model Reuse” work package for the 1.5M€ project “*NEXUS: Interoperable Machine*  
316 *Learning with Universal Representations*”. He previously held roles as a visiting researcher at  
317 the University of Cambridge, a Research Scientist at Amazon Alexa, and an Applied Scientist  
318 at Amazon Search. His research has been featured in top-tier AI conferences and journals,  
319 including CVPR, NeurIPS, ACM, ACL, and LoG. In addition to his scientific contributions, he has  
320 played an active role in the research community as the organizer of several workshops, including  
321 the UniReps workshop at NeurIPS (one of the most attended in NeurIPS 2023 and 2024) and  
322 as a program committee member for leading conferences such as CVPR, ICML, NeurIPS, ICLR, etc.

---

324 **Zorah Lähner**

325 *University of Bonn and Lamarr Institute*

326 **Zorah Lähner** is assistant professor and head of the Geometry in Machine Learning group at the  
327 University of Bonn and the Lamarr Institute for Machine Learning and Artificial Intelligence in  
328 Germany. She previously worked at the University of Siegen, Technical University of Munich, Meta  
329 Reality Labs and Toshiba Research Europe. Her research focuses on geometric deep learning and 3D  
330 vision, and has been published in major machine learning and computer vision conferences, including  
331 NeurIPS, ICLR, CVPR, and ICCV. Beyond that she has been involved in the organization of the  
332 UniReps workshop at NeurIPS, a section of Women in Vision, been program chair of the German  
333 Conference on Pattern Recognition, and received several outstanding reviewer awards.

334 **PROGRAM COMMITTEE**

335 The following is a preliminary list of 50 *confirmed* Program Committee members. Based on an  
336 anticipated 100 submissions and a target load of no more than four short papers per reviewer, we will  
337 need roughly 75–100 reviewers. We are therefore actively recruiting additional members, especially  
338 from diverse and historically under-represented communities in AI, and expect to reach the required  
339 committee size within the next few weeks.

1. Tommaso Mencattini (EPFL) [\[Confirmed\]](#)
2. Andrea Caciolai (Meta) [\[Confirmed\]](#)
3. Simone Antonelli (CISPA) [\[Confirmed\]](#)
4. Stefano Esposito (U Tübingen) [\[Confirmed\]](#)
5. Irene Tallini (Area Science Park, Trieste) [\[Confirmed\]](#)
6. Pietro Barbiero (USI) [\[Confirmed\]](#)
7. Lorenzo Giusti (CERN) [\[Confirmed\]](#)
8. Eleonora Gualdoni (Apple) [\[Confirmed\]](#)
9. Irene Cannistraci (ETH) [\[Confirmed\]](#)
10. Marco Fumero (ISTA) [\[Confirmed\]](#)
11. Daniele Baieri (U Milano-Bicocca) [\[Confirmed\]](#)
12. Luca Moschella (Apple) [\[Confirmed\]](#)
13. Luca Cosmo (Ca' Foscari U Venice) [\[Confirmed\]](#)
14. Robert Adrian Minut (Sapienza) [\[Confirmed\]](#)
15. Adam Goliński (Apple) [\[Confirmed\]](#)
16. Daniele Solombrino (Sapienza) [\[Confirmed\]](#)
17. Miao Xiong (NUS) [\[Confirmed\]](#)
18. Maks Ovsjanikov (Ecole polytechnique) [\[Pending\]](#)
19. Federico Danieli (Apple) [\[Confirmed\]](#)
20. Gabriele Sarti (U Groningen) [\[Confirmed\]](#)
21. Arianna Rampini (Autodesk) [\[Confirmed\]](#)
22. Davide Marincione (Sapienza) [\[Confirmed\]](#)
23. Michele Miranda (Translated) [\[Confirmed\]](#)
24. Emanuele Rossi (VantAI) [\[Confirmed\]](#)
25. Michele Mancusi (Moises) [\[Confirmed\]](#)
26. Emily Cheng (UPF) [\[Confirmed\]](#)
27. Roberto Dessì (Samaya AI) [\[Confirmed\]](#)
28. Tal Remez (Meta) [\[Pending\]](#)
29. Giorgio Mariani (U Milano-Bicocca) [\[Confirmed\]](#)
30. Antonio Norelli (MIT) [\[Confirmed\]](#)
31. Giorgio Strano (Sapienza) [\[Confirmed\]](#)
32. Florian Bernard (U Bonn) [\[Pending\]](#)
33. Alessio Devoto (NVIDIA) [\[Confirmed\]](#)
34. Lucas Weber (Fraunhofer IIS) [\[Confirmed\]](#)
35. Or Litany (Technion) [\[Confirmed\]](#)
36. Riccardo Marin (TUM) [\[Confirmed\]](#)
37. Giovanni Trappolini (Sapienza) [\[Confirmed\]](#)
38. Emtiyaz Khan (Riken) [\[Pending\]](#)
39. Emilian Postolache (IRIS Audio) [\[Confirmed\]](#)
40. Simone Melzi (U Milano-Bicocca) [\[Confirmed\]](#)
41. Asako Kanezaki (Tokyo Tech) [\[Pending\]](#)
42. Zorah Lähner (U Bonn) [\[Confirmed\]](#)
43. Federico Bomba (Sineglossa) [\[Confirmed\]](#)
44. Silvia Zuffi (CNR) [\[Confirmed\]](#)
45. Alex Bronstein (ISTA) [\[Confirmed\]](#)
46. Fabrizio Frasca (Technion) [\[Confirmed\]](#)
47. Thomas Möllenhoff (Riken) [\[Confirmed\]](#)
48. Michael Möller (U Siegen) [\[Confirmed\]](#)
49. Laura Leal-Taixé (NVIDIA) [\[Pending\]](#)
50. Antonio Gargiulo (Sapienza) [\[Confirmed\]](#)
51. Haggai Maron (Technion) [\[Confirmed\]](#)
52. Silvio Severino (Amazon) [\[Confirmed\]](#)
53. Jacob Morrison (Ai2) [\[Confirmed\]](#)
54. Ethan Shen (U Washington) [\[Confirmed\]](#)
55. Osamu Hirose (U Kanazawa) [\[Confirmed\]](#)
56. Clementine Domine (UCL) [\[Confirmed\]](#)
57. Pritish Chakraborty (IIT Bombay) [\[Pending\]](#)
58. Viktor Stenby Johansson (DTU) [\[Pending\]](#)
59. Akshit Acharya (King's College London) [\[Pending\]](#)
60. Hugo Daniel Monzón Maldonado (Riken) [\[Pending\]](#)

---

## 378 REFERENCES

## 379

380 Markus Anderljung, Joslyn Barnhart, Anton Korinek, Jade Leung, Cullen O’Keefe, Jess Whittlestone,  
381 Shahar Avin, Miles Brundage, Justin Bullock, Duncan Cass-Beggs, Ben Chang, Tantum Collins,  
382 Tim Fist, Gillian Hadfield, Alan Hayes, Lewis Ho, Sara Hooker, Eric Horvitz, Noam Kolt, Jonas  
383 Schuett, Yonadav Shavit, Divya Siddarth, Robert Trager, and Kevin Wolf. Frontier ai regulation:  
384 Managing emerging risks to public safety, 2023. URL <https://arxiv.org/abs/2307.03718>. 3

385

386 Kshitij Bansal, Christian Szegedy, Markus Norman Rabe, Sarah M. Loos, and Viktor Toman. Learning  
387 to reason in large theories without imitation, 2021. URL <https://openreview.net/forum?id=qbRv1k2AcH>. 2

388

389 Nora Belrose, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and Stella  
390 Biderman. Leace: Perfect linear concept erasure in closed form. *Advances in Neural Information  
391 Processing Systems*, 36:66044–66063, 2023. 2

392

393 Joshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yu-  
394 val Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian K. Hadfield, Jeff Clune,  
395 Tegan Maharaj, Frank Hutter, Atilim Gunes Baydin, Sheila A. McIlraith, Qiqi Gao, Ashwin  
396 Acharya, David Krueger, Anca Dragan, Philip Torr, Stuart Russell, Daniel Kahneman, Jan Markus  
397 Brauner, and Sören Mindermann. Managing extreme ai risks amid rapid progress. *Science*, 384:842  
398 – 845, 2023. URL <https://api.semanticscholar.org/CorpusID:269929051>. 3

399

400 Borhane Blili-Hamelin, Christopher Graziul, Leif Hancox-Li, Hananel Hazan, El-Mahdi El-Mhamdi,  
401 Avijit Ghosh, Katherine A Heller, Jacob Metcalf, Fabricio Murai, Eryk Salvaggio, Andrew J Smart,  
402 Todd Snider, Mariame Tighanimine, Talia Ringer, Margaret Mitchell, and Shiri Dori-Hacohen.  
403 Position: Stop treating ‘AGI’ as the north-star goal of AI research. In *Forty-second International  
404 Conference on Machine Learning Position Paper Track*, 2025. URL <https://openreview.net/forum?id=1RlrtH6ydW>. 3

405

406 Matthew Bowers, Theo X. Olausson, Lionel Wong, Gabriel Grand, Joshua B. Tenenbaum, Kevin  
407 Ellis, and Armando Solar-Lezama. Top-down synthesis for library learning. *Proc. ACM Program.  
408 Lang.*, 7(POPL), January 2023. doi: 10.1145/3571234. URL <https://doi.org/10.1145/3571234>. 2

409

410 Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė  
411 Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron  
412 McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-  
413 Johnson, Jackson Kernion, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal  
414 Ndousse, Liane Lovitt, Nelson Elhage, Nicholas Schiefer, Nicholas Joseph, Noemí Mercado, Nova  
415 DasSarma, Robin Larson, Sam McCandlish, Sandipan Kundu, Scott Johnston, Shauna Kravec,  
416 Sheer El Showk, Stanislav Fort, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan  
417 Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, and Jared Kaplan. Measuring progress on  
418 scalable oversight for large language models, 2022. URL <https://arxiv.org/abs/2211.03540>. 2

419

420 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner,  
421 Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization:  
422 eliciting strong capabilities with weak supervision. In *Proceedings of the 41st International  
423 Conference on Machine Learning*, pp. 4971–5012, 2024. 2

424

425 Francois Chollet, Mike Knoop, Gregory Kamradt, and Bryan Landers. Arc prize 2024: Technical  
426 report, 2025a. URL <https://arxiv.org/abs/2412.04604>. 2

427

428 Francois Chollet, Mike Knoop, Gregory Kamradt, Bryan Landers, and Henry Pinkard. Arc-agi-2:  
429 A new challenge for frontier ai reasoning systems, 2025b. URL <https://arxiv.org/abs/2505.11831>. 2

430

431 Caia Costello, Simon Guo, Anna Goldie, and Azalia Mirhoseini. Think, prune, train, improve: Scaling  
reasoning without scaling models, 2025. URL <https://arxiv.org/abs/2504.18116>. 2

---

432 Connor Dunlop, Weiwei Pan, Julia Smakman, Lisa Soder, Siddharth Swaroop, and Noam Kolt.  
433 Position: AI agents & liability – mapping insights from ML and HCI research to policy.  
434 In *Workshop on Socially Responsible Language Modelling Research*, 2024. URL <https://openreview.net/forum?id=pa80BLEavx>. 3

435

436 Kevin Ellis. Human-like few-shot learning via bayesian reasoning over natural language. In  
437 *Proceedings of the 37th International Conference on Neural Information Processing Systems*,  
438 NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc. 2

439

440 Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt,  
441 Luc Cary, Armando Solar-Lezama, and Joshua B. Tenenbaum. Dreamcoder: bootstrapping  
442 inductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd ACM*  
443 *SIGPLAN International Conference on Programming Language Design and Implementation*,  
444 PLDI 2021, pp. 835–850, New York, NY, USA, 2021. Association for Computing Machinery.  
445 ISBN 9781450383912. doi: 10.1145/3453483.3454080. URL <https://doi.org/10.1145/3453483.3454080>. 2

446

447 Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mo-  
448 mhammadim Barekatain, Alexander Novikov, Francisco J. R. Ruiz, Julian Schrittwieser, Grze-  
449 gorz Swirszcz, David Silver, Demis Hassabis, and Pushmeet Kohli. Discovering faster ma-  
450 trix multiplication algorithms with reinforcement learning. *Nature*, 610:47 – 53, 2022. URL  
451 <https://api.semanticscholar.org/CorpusID:252717185>. 2

452

453 Iason Gabriel. Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437,  
454 2020. 3

455

456 Iason Gabriel, Arianna Manzini, Geoff Keeling, Lisa Anne Hendricks, Verena Rieser, Hasan Iqbal,  
457 Nenad Tomašev, Ira Ktena, Zachary Kenton, Mikel Rodriguez, Seliem El-Sayed, Sasha Brown,  
458 Canfer Akbulut, Andrew Trask, Edward Hughes, A. Stevie Bergman, Renee Shelby, Nahema  
459 Marchal, Conor Griffin, Juan Mateos-Garcia, Laura Weidinger, Winnie Street, Benjamin Lange,  
460 Alex Ingerman, Alison Lentz, Reed Enger, Andrew Barakat, Victoria Krakovna, John Oliver  
461 Siy, Zeb Kurth-Nelson, Amanda McCroskery, Vijay Bolina, Harry Law, Murray Shanahan, Lize  
462 Alberts, Borja Balle, Sarah de Haas, Yetunde Ibitoye, Allan Dafoe, Beth Goldberg, Sébastien Krier,  
463 Alexander Reese, Sims Witherspoon, Will Hawkins, Maribeth Rauh, Don Wallace, Matija Franklin,  
464 Josh A. Goldstein, Joel Lehman, Michael Klenk, Shannon Vallor, Courtney Biles, Meredith Ringel  
465 Morris, Helen King, Blaise Agüera y Arcas, William Isaac, and James Manyika. The ethics of  
466 advanced ai assistants, 2024. URL <https://arxiv.org/abs/2404.16244>. 3

467

468 Katja Grace, Harlan Stewart, Julia Fabienne Sandkühler, Stephen Thomas, Ben Weinstein-Raun, and  
469 Jan Brauner. Thousands of ai authors on the future of ai. *arXiv preprint arXiv:2401.02843*, 2024. 1

470

471 Ben Green. The false promise of risk assessments: epistemic reform and the limits of fairness.  
472 In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT\*  
473 '20, pp. 594–606, New York, NY, USA, 2020. Association for Computing Machinery. ISBN  
474 9781450369367. doi: 10.1145/3351095.3372869. URL <https://doi.org/10.1145/3351095.3372869>. 3

475

476 Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *ArXiv*, abs/2408.08435,  
477 2024. URL <https://api.semanticscholar.org/CorpusID:271892234>. 2

478

479 Mason Kamb and Surya Ganguli. An analytic theory of creativity in convolutional diffusion mod-  
480 els. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=ilpL2qACla>. 3

481

482 Gyeonggeon Lee, Lehong Shi, Ehsan Latif, Yizhu Gao, Arne Bewersdorff, Matthew Nyaaba, Shuchen  
483 Guo, Zhengliang Liu, Gengchen Mai, Tianming Liu, et al. Multimodality of ai for education:  
484 Towards artificial general intelligence. *IEEE Transactions on Learning Technologies*, 2025. 3

485

486 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu,  
487 Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling  
488 reinforcement learning from human feedback with ai feedback. In *International Conference on*  
489 *Machine Learning*, pp. 26874–26901. PMLR, 2024. 2

---

486 Belinda Z. Li, Alex Tamkin, Noah Goodman, and Jacob Andreas. Eliciting human preferences with  
487 language models. In *The Thirteenth International Conference on Learning Representations*, 2025.  
488 URL <https://openreview.net/forum?id=LvDwwAgMEW>. 2  
489

490 Wen-Ding Li, Keya Hu, Carter Larsen, Yuqing Wu, Simon Alford, Caleb Woo, Spencer M. Dunn,  
491 Hao Tang, Michelangelo Naim, Dat Nguyen, Wei-Long Zheng, Zenna Tavares, Yewen Pu, and  
492 Kevin Ellis. Combining induction and transduction for abstract reasoning, 2024. URL <https://arxiv.org/abs/2411.02272>. 2  
493

494 Daniel J. Mankowitz, Andrea Michi, Anton Zhernov, Marco Gelmi, Marco Selvi, Cosmin Paduraru,  
495 Edouard Leurent, Shariq Iqbal, Jean-Baptiste Lespiau, Alex Ahern, Thomas Köppe, Kevin Millikin,  
496 Stephen Gaffney, Sophie Elster, Jackson Broshear, Chris Gamble, Kieran Milan, Robert Tung,  
497 Minjae Hwang, Taylan Cemgil, Mohammadamin Barekatain, Yujia Li, Amol Mandhane, Thomas  
498 Hubert, Julian Schrittwieser, Demis Hassabis, Pushmeet Kohli, Martin Riedmiller, Oriol Vinyals,  
499 and David Silver. Faster sorting algorithms discovered using deep reinforcement learning. *Nature*,  
500 618(7964):257–263, Jun 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06004-9. URL  
501 <https://doi.org/10.1038/s41586-023-06004-9>. 2

502 Yu Meng, Mengzhou Xia, and Danqi Chen. SimPO: Simple preference optimization with a reference-  
503 free reward. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,  
504 2024. URL <https://openreview.net/forum?id=3Tzcot1Lkb>. 2  
505

506 Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Jiang, Ebrahim Songhori, Shen Wang, Young-  
507 Joon Lee, Eric Johnson, Omkar Pathak, Sungmin Bae, et al. Chip placement with deep reinforce-  
508 ment learning. *arXiv preprint arXiv:2004.10746*, 2020. 2

509 Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Wenjie Jiang, Ebrahim Songhori, Shen Wang,  
510 Young-Joon Lee, Eric Johnson, Omkar Pathak, Azade Nova, et al. A graph placement methodology  
511 for fast chip design. *Nature*, 594(7862):207–212, 2021. 2

512 Meredith Ringel Morris, Jascha Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksan-  
513 dra Faust, Clement Farabet, and Shane Legg. Levels of agi for operationalizing progress on the  
514 path to agi. *arXiv preprint arXiv:2311.02462*, 2023. 1, 3

515 Henrik Nolte, Miriam Rateike, and Michèle Finck. Robustness and cybersecurity in the EU ar-  
516 tificial intelligence act. In *NeurIPS 2024 Workshop on Regulatable ML*, 2025. URL <https://openreview.net/forum?id=1m1XKCGerV>. 3  
518

519 Tejal Patwardhan, Rachel Dias, Elizabeth Proehl, Grace Kim, Michele Wang, Olivia Watkins,  
520 Simón Posada Fishman, Marwan Aljubeh, Phoebe Thacker, Laurance Fauconnet, Natalie S. Kim,  
521 Patrick Chao, Samuel Miserendino, Gildas Chabot, David Li, Michael Sharman, Alexandra Barr,  
522 Amelia Glaese, and Jerry Tworek. Gdpval: Evaluating ai model performance on real-world  
523 economically valuable tasks, 2025. URL <https://arxiv.org/abs/2510.04374>. 3  
524

525 Vinodkumar Prabhakaran, Margaret Mitchell, Timnit Gebru, and Iason Gabriel. A human rights-based  
526 approach to responsible ai, 2022. URL <https://arxiv.org/abs/2210.02667>. 3

527 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
528 Finn. Direct preference optimization: Your language model is secretly a reward model. In  
529 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HPuSIXJaa9>. 2  
530

531 Inioluwa Deborah Raji, Peggy Xu, Colleen Honigsberg, and Daniel Ho. Outsider oversight: Designing  
532 a third party audit ecosystem for ai governance. In *Proceedings of the 2022 AAAI/ACM Conference*  
533 *on AI, Ethics, and Society*, AIES ’22, pp. 557–571, New York, NY, USA, 2022. Association  
534 for Computing Machinery. ISBN 9781450392471. doi: 10.1145/3514094.3534181. URL  
535 <https://doi.org/10.1145/3514094.3534181>. 3  
536

537 Joshua S. Rule, Steven T. Piantadosi, Andrew Cropper, Kevin Ellis, Maxwell Nye, and Joshua B.  
538 Tenenbaum. Symbolic metaprogram search improves learning efficiency and explains rule learning  
539 in humans. *Nature Communications*, 15(1):6847, 2024. doi: 10.1038/s41467-024-50966-x. URL  
540 <https://doi.org/10.1038/s41467-024-50966-x>. 2

---

540 Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green ai. *Commun. ACM*, 63(12):  
541 54–63, November 2020. ISSN 0001-0782. doi: 10.1145/3381831. URL <https://doi.org/10.1145/3381831>. 2

542

543 Daniel Selsam, Matthew Lamm, Benedikt Bünz, Percy Liang, Leonardo de Moura, and David L.  
544 Dill. Learning a SAT solver from single-bit supervision. In *International Conference on Learning  
545 Representations*, 2019. URL [https://openreview.net/forum?id=HJMC\\_iA5tm](https://openreview.net/forum?id=HJMC_iA5tm). 2

546

547 Rulin Shao, Jacqueline He, Akari Asai, Weijia Shi, Tim Dettmers, Sewon Min, Luke Zettlemoyer,  
548 and Pang Wei Koh. Scaling retrieval-based language models with a trillion-token datastore. In  
549 *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL  
550 <https://openreview.net/forum?id=iAkhPz7Qt3>. 2

551

552 Yijia Shao, Humishka Zope, Yucheng Jiang, Jiaxin Pei, David Nguyen, Erik Brynjolfsson, and Diyi  
553 Yang. Future of work with ai agents: Auditing automation and augmentation potential across the  
554 u.s. workforce, 2025. URL <https://arxiv.org/abs/2506.06576>. 3

555

556 Kensen Shi, Hanjun Dai, Wen-Ding Li, Kevin Ellis, and Charles Sutton. Lambdabeam: neural  
557 program search with higher-order functions and lambdas. In *Proceedings of the 37th International  
558 Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023.  
Curran Associates Inc. 2

559

560 Hao Tang, Darren Yan Key, and Kevin Ellis. Worldcoder, a model-based LLM agent: Building  
561 world models by writing code and interacting with the environment. In *The Thirty-eighth Annual  
562 Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=QGJSXMhVaL>. 2

563

564 Kiran Tomlinson, Sonia Jaffe, Will Wang, Scott Counts, and Siddharth Suri. Working with ai:  
565 Measuring the applicability of generative ai to occupations, 2025. URL <https://arxiv.org/abs/2507.07935>. 3

566

567 Ying Wen, Ziyu Wan, and Shao Zhang. Language games as the pathway to artificial superhuman  
568 intelligence, 2025. URL <https://arxiv.org/abs/2501.18924>. 3

569

570 Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants. In  
571 *International Conference on Machine Learning (ICML)*, 2019. 2

572

573 Jenny Zhang, Shengran Hu, Cong Lu, Robert Lange, and Jeff Clune. Darwin godel machine: Open-  
574 ended evolution of self-improving agents, 2025. URL <https://arxiv.org/abs/2505.22954>. 2

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593