

INAUGURAL COHORT · FEBRUARY - MAY 2026

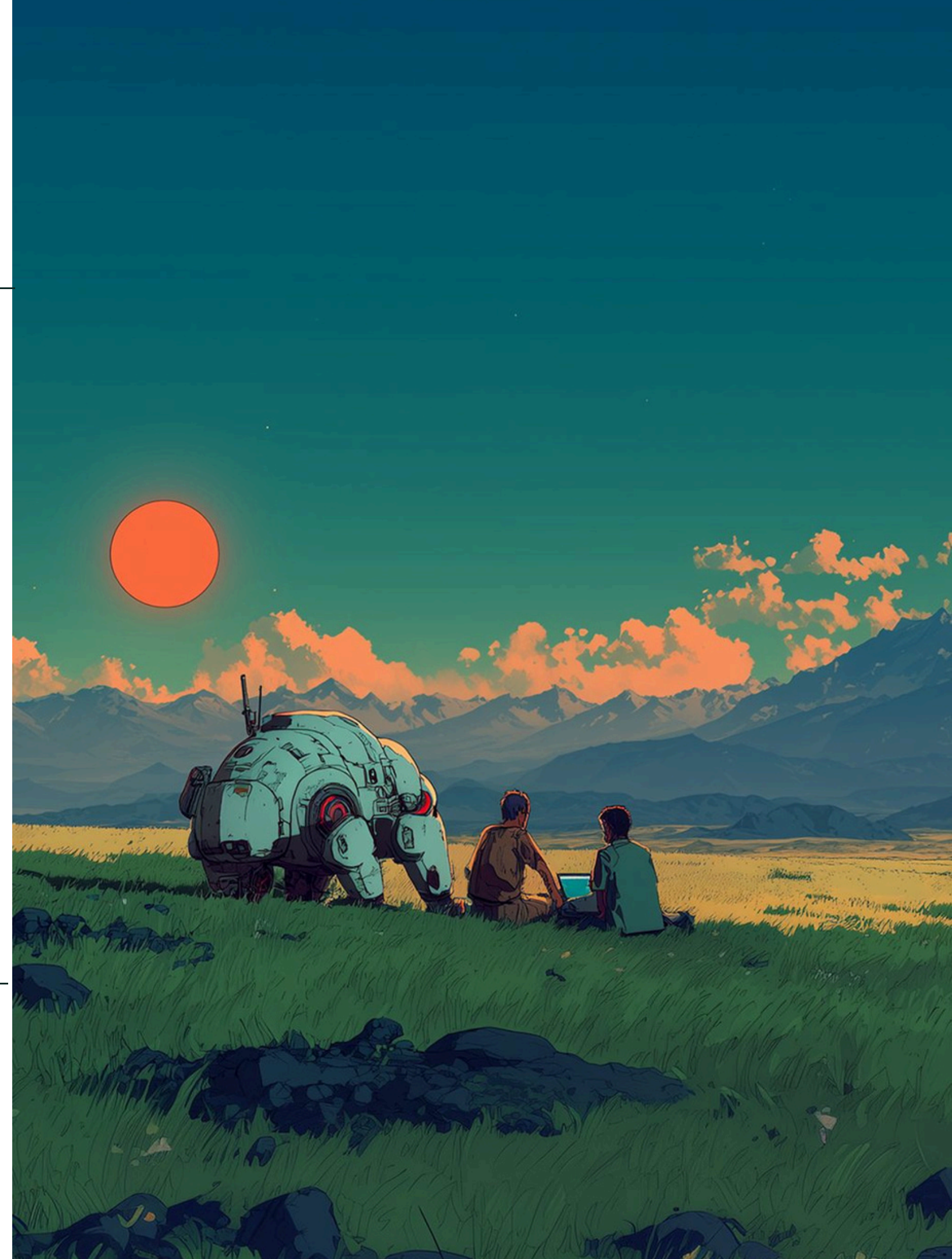
---

# Seeding a new cooperative AI ecosystem

---

COOPERATIVE AI RESEARCH FELLOWSHIP

An AI Safety South Africa Initiative



# Cohort & Mentors

The program brought 10 researchers from 8 countries to South Africa to work on Cooperative AI for 3 months in early 2026. Fellows were selected from 1100+ applicants across a five-stage application process and were paired with mentors from institutions such as Google DeepMind and MIT. The program started with a week-long retreat, which was foundational in developing project ideas and cohort bonding.

## FELLOW

## MENTOR(S)

<b>Bhavyesh Sajja</b>	Max Kleinman-Weiner (University of Washington) & Tan Xhi Xuan (National University of Singapore)
<b>Joseph Low &amp; Oscar Duys</b>	Michiel Bakker (MIT, Google DeepMind) & Lewis Hammond (Cooperative AI Foundation, University of Oxford)
<b>Akash Kundu</b>	Vincent Conitzer & Emanuel Tewelde (Foundations of Cooperative AI Lab, Carnegie Mellon University)
<b>Omer Kamal Ali Ebead</b>	Joel Z. Leibo (Google DeepMind)
<b>Qi Guo</b>	Lewis Hammond (Cooperative AI Foundation, University of Oxford)
<b>Pramod Kaushik</b>	Sahar Abdelnabi (ELLIS Institute Tübingen & MPI-IS, COMPASS Lab)
<b>Yves Bicker</b>	Zhijing Jin & David Guzman Piedrahita (University of Toronto, Jinesis Lab)
<b>Van Quynh Thi Truong</b>	Zhijing Jin & David Guzman Piedrahita (University of Toronto, Jinesis Lab)
<b>Mariana Meireles</b>	Zhijing Jin & David Guzman Piedrahita (University of Toronto, Jinesis Lab)

# Key Outcomes
















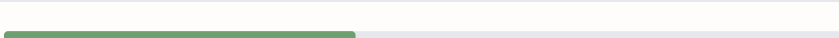
<p>FELLOWS</p> <p><b>10</b></p> <p>Paired with mentors at MIT, Oxford, DeepMind and more.</p>	<p>FELLOW NPS</p> <p><b>9.1/10</b></p> <p>Participant recommendation score. No fellow rated the program below 7.</p>	<p>STAKEHOLDER NPS</p> <p><b>9.1/10</b></p> <p>Mentor and partner recommendation score. No respondent rated below 8.</p>	<p>MENTORSHIP RATING</p> <p><b>8.5/10</b></p> <p>The highest-rated program component rated by fellows.</p>
<p>NEURIPS SUBMISSIONS</p> <p><b>4+</b></p> <p>Additional submissions at ICML.</p>	<p>CONTINUED FUNDING</p> <p><b>6/10</b></p> <p>Fellows with secured continuation funding post-fellowship.</p>	<p>PRODUCT OUTCOMES</p> <p><b>1</b></p> <p>Habermolt launched, paper to arXiv, acquired by Change.org mid-fellowship.</p>	<p>CAREER IMPACT</p> <p><b>10/10</b></p> <p>All fellows reported the program confirmed, redirected or accelerated their trajectory.</p>

# Research Outputs

FELLOW	PROJECT	HEADLINE OUTCOME
<b>Bhavyesh Sajja</b>	Evaluating Rationality in Natural-Language Contracting	NEURIPS SUBMITTED
<b>Joseph Low &amp; Oscar Duys</b>	Delegating Deliberation to AI Agents (Habermolt)	LAUNCHED · ACQUIRED · ARXIV
<b>Akash Kundu</b>	Similarity as a Signal	NEURIPS SUBMITTED · EXTENSION FUNDED
<b>Omer Kamal Ali Ebead</b>	Contextual Integrity in Multi-Agent LLM Systems	EXTENSION FUNDED
<b>Qi Guo</b>	Inter-Agent Influence	NEURIPS SUBMITTED · HIRED BY CAIF
<b>Van Quynh Thi Truong</b>	Social Strategies for Cooperation in Multi-Agent Societies	COLM / NEURIPS TARGET
<b>Yves Bicker</b>	MoralGym: Training Cooperative LLM Agents	RAPID GRANT · ICLR TARGET
<b>Pramod Kaushik</b>	Strategic Indirect Speech in LLMs	EXPERIMENTS COMPLETE
<b>Mariana Meireles</b>	Toward Collective Intelligence	EXPERIMENTS COMPLETE

# Fellow Ratings

Program components, scored out of 10.

Mentorship		8.5
AISSA infrastructure		8.5
Co-working (Workshop 17)		8.4
Stellenbosch workshop		8.3
Operational support		8.3
Opening retreat		8.1
\$500 compute budget		7.9
Opportunities for peer connection		7.9
ShockLabs (UCT AI lab)		7.9
Seminars		7.6
Team meals		7.2
Housing (Curiosity)		6.9
Research management		6.2
Peer review session		6.1
Pre-fellowship course		6.0
GPU compute (n=5)		4.2

FELLOW NPS

9.1/10

## Peer bonds were the most-cited value

*They're all super sharp and generous and I've learnt an enormous amount from them on everything from technical questions to how to think about my career. Three months is a short time but I genuinely believe these are friendships for life.*

FELLOW

*Having people to lean on, did have a huge impact on me, especially during the middle when there are no quick results to look at. I think they helped me push further. Moreover, I was amazed by how diligent everyone was at work. They somehow passed on their discipline which made me push slightly harder than I would have, because everyone around me, was.*

FELLOW

# Mentor & Stakeholder Feedback

10 respondents from MIT, NUS, DeepMind, CAIF, PIBBSS, CMU, Oxford, UCT, UW and Lambda. Aggregate NPS 9.1.

STAKEHOLDER NPS

9.1/10

MENTOR RATINGS

- **Weekly meetings with fellows** — 9.3 / 10
- **Mentor matching experience** — 9.2 / 10
- **Program prepared mentee for career in coop AI** — 9 / 10
- **Operational support & comms** — 8.7 / 10
- **Pre-fellowship project development** — 8 / 10
- **Value from Research Manager (mentor-rated)** — 7.3 / 10

*CAIRF is launching the next generation of scholars in Cooperative AI from around the world. Highly recommend this well run program.*

MAX KLEIMAN-WEINER · UNIVERSITY OF WASHINGTON

*I loved working with Oscar and Joseph, the two CAIRF fellows I mentored. Having everyone together in person in Cape Town (except for me, sadly) made a huge difference to how quickly things moved and what the fellows were able to achieve. I'm also very grateful for all the support from Claude Formanek and the rest of the team!*

MICHIEL BAKKER · MIT

*Well-structured, made tons of research progress, funding is providing for experiments, and it seems like my student benefited a lot from the community and support too :)*

TAN ZHI XUAN · NUS

# Outcomes vs. Aims

The aims we set when we launched, and where we landed by the end of the fellowship in May 2026.

6 MONTHS	12 high-quality research projects in cooperative AI	10 projects produced, two fellows removed mid-program. Quality strong (NeurIPS submissions, acquisition, ongoing research); target nearly met.	PARTIAL
6 MONTHS	3–5 fellows identified for permanent positions	6 fellows with continuing funded research — Qi at CAIF, Joseph & Oscar at Change.org, Omer & Akash on CAIRF extension, Yves on rapid grant.	SUCCESS
6 MONTHS	2–3 South African academics shifting research focus toward AI safety	Had many professors express interest in getting involved in our work, either by providing students for projects, attending events, or working on research together.	STARTED
1–2 YEARS	2–3 permanent cooperative AI research positions at UCT AI Safety Hub	Submitted a joint proposal with UCT for a multi-agent safety lab with Claude Formanek with the PI, and Prof. Shock as CO-PI	STARTED
1–2 YEARS	Established talent pipeline from Africa to global AI safety organisations	Qi Guo's hire by CAIF is the first concrete instance. Mentor network now reaches Jin lab, FOCAL, COMPASS, CoSI lab, Computational Minds & Machines, MIT and DeepMind.	STARTED

# Going Forward

Six recommendations the team is committing to for the next iteration, synthesised from fellow, mentor and team feedback.

## Co-working space

### Customised office

Despite receiving a high rating out of 10, many written responses that highlighted improvements for the office setup. Future planning will prioritise securing a private, dedicated facility that separates quiet work zones from social or collaborative areas and caters to a wider variety of work styles and social needs.

## Community Manager

### Hire community manager

A community manager will serve as the dedicated bridge between the fellows and the AISSA team, so our presence is more visible and approachable. This role would also support preventing and resolving community health issues.

## Research management

### Clarify RM role & reporting

Provide a clearer brief to RMs and fellows about what is expected of them. Also, we would clarify reporting and feedback systems from the start of the fellowship. Furthermore, we think we could improve outcomes by scoping this role better in terms of hourly demand.

## Compute

### Expand the compute pool.

Combine GPU + API budgets into a single flexible per-fellow allocation. Raise default API budget for compute-heavy projects. Budget \$1500-3500/fellow.

## Timeline

### Secure funding & housing earlier

Funding transferred 6 months prior, housing at least 4 months before fellowship start. This would prevent significant uncertainty that blocked further decisions early on in the fellowship.

## Due Diligence

### Screening & code of conduct

Contact candidate's references and screen them through CEA. Require signing of code of conduct prior to start of fellowship and have a talk about norms and appropriate conduct on day 1.

END

# Thank you

CONTACT

[info@cai-research-fellowship.com](mailto:info@cai-research-fellowship.com)

MORE ON THE PROGRAM

[www.cai-research-fellowship.com](http://www.cai-research-fellowship.com)

