

# Felix Jedidja Binder

me@felixbinder.net

ac.felixbinder.net

+1 (858) 291-2056

[Full CV](#)

## Education

2019-2025	<b>University of California San Diego</b>	PhD in Cognitive Science
2024-2025	<b>Stanford University</b>	Visiting Researcher
2013-2019	<b>Freie Universität Berlin</b>	Bachelor of Arts in Philosophy & Computer Science

## Experience

2025 Menlo Park	<b>Research Scientist</b>   Meta   <i>Meta Superintelligence Labs</i> <ul style="list-style-type: none"><li>AI Safety &amp; Alignment for TBD Labs</li></ul>
2025 San Francisco	<b>Research Scientist</b>   Scale AI   <i>Safety, Evaluations &amp; Alignment Lab</i> <ul style="list-style-type: none"><li>Leading development of benchmark measuring active value learning in LLMs for public evaluation.</li><li>Contributing to economic impact evaluation of Computer Use Agents.</li></ul>
2019-2024 San Diego	<b>Graduate Student Researcher</b>   University of California San Diego   <i>Cognitive Science Department</i> <ul style="list-style-type: none"><li>Created and maintained a full stack setup for running web experiments evaluating human and AI behavior on a range of cognitive tasks (<a href="#">Cognitive AI Benchmarking</a>).</li><li>Led a study comparing humans and planning algorithms on a simulated physical construction task.</li><li>Created a dataset for a large benchmarking study of physical understanding in humans &amp; AI (<a href="#">Physion</a>) with NeuroAILab (Stanford) and Computational Cognitive Science lab (MIT).</li><li>Evaluated a broad suite of state-of-the-art vision &amp; particle-based AI models on the Physion dataset. Found that AI models do not yet meet human performance in physical understanding.</li><li>Created public outreach videos on neural networks and AI ethics for high school students with <a href="#">pathways2AI</a>.</li><li>Taught undergrad &amp; graduate courses, including <i>Reinforcement Learning</i> and <i>Data Science</i>.</li><li>Organized the Cognitive AI Benchmarking workshop at the 45th Annual Meeting of the Cognitive Science Society.</li></ul>
2024 Berkeley	<b>AI Safety Research Fellow</b> with Owain Evans   Constellation Astra Fellowship <ul style="list-style-type: none"><li>Developed <a href="#">novel experimental framework</a> to train and evaluate introspection in large language models (LLMs).</li><li><a href="#">Demonstrated</a> that frontier LLMs (GPT-4, GPT-4o, Llama 3 70B) can acquire knowledge about themselves through introspection, not just from training data.</li></ul>
2023 Cambridge, MA	<b>AI Research Scientist Intern</b>   Cambria Labs <ul style="list-style-type: none"><li>Oversaw creation of multimodal video dataset for physical understanding and prediction.</li><li>Built a data pipeline for data management &amp; model training; implemented and trained a suite of vision transformer based models on the dataset.</li></ul>
2023	<b>Artificial General Intelligence Safety Fundamentals Course</b>   BlueDot Impact <ul style="list-style-type: none"><li>Developed an <a href="#">evaluation protocol</a> that isolates causal effects of context for analyzing steganographic tendencies (covert information encoding) in large language models.</li><li>Conducted an investigation into potential steganographic behavior in current LLMs, utilizing the aforementioned evaluation protocol.</li></ul>

## Skills

**Programming & AI** Python & PyTorch, AI Safety & Alignment, RL, LLMs, Interpretability, Planning & Reasoning, [Evals](#)  
**Statistics** Experiment Design, Model Fitting & Analysis, Hypothesis Testing, Bayesian Statistics  
**Communication** Scientific Writing, Public Science Communication, Data Visualization, Cross-Field Communication

## Selected Publications

\* indicates equal contribution.

2025 **Binder, F.**, Mattar, M., Kirsh, D., & Fan, J. Humans Select Subgoals That Balance Immediate and Future Cognitive Costs During Physical Assembly. *Cognitive Science*.

2025 **Binder, F.** Thinking Through Action: Prediction, Planning, and Metacognition in Problem-Solving. *Doctoral dissertation, University of California, San Diego*. [Dissertation](#)

2024 **Binder, F.\***, Chua, J.\*., Korbak, T., Sleight, H., Hughes, J., Long, R., Perez, E., Turpin, M., & Evans, O. Looking Inward: Language Models Can Learn About Themselves by Introspection. *ICLR 2025*. | [Code & Paper](#)

2021 Bear, D.\*., Wang, E.\*., Mrowca, D.\*., **Binder, F.\***, Tung, H., Pramod, R. T., Holdaway, C., Tao, S., Smith, K., Sun, F., Fei-Fei, L., Kanwisher, N., Tenenbaum, J., Yamins, D.\*\* & Fan, J.\*\* Physion: Evaluating Physical Prediction from Vision in Humans and Machines. *NeurIPS 2021 (Datasets & Benchmarks track)* | [Code & Paper](#), [NeurIPS Presentation](#)