

Adventures with Agentic AI (in particle physics)

**11th LITP Spring Symposium: Theoretical Physics and AI
University of Michigan**

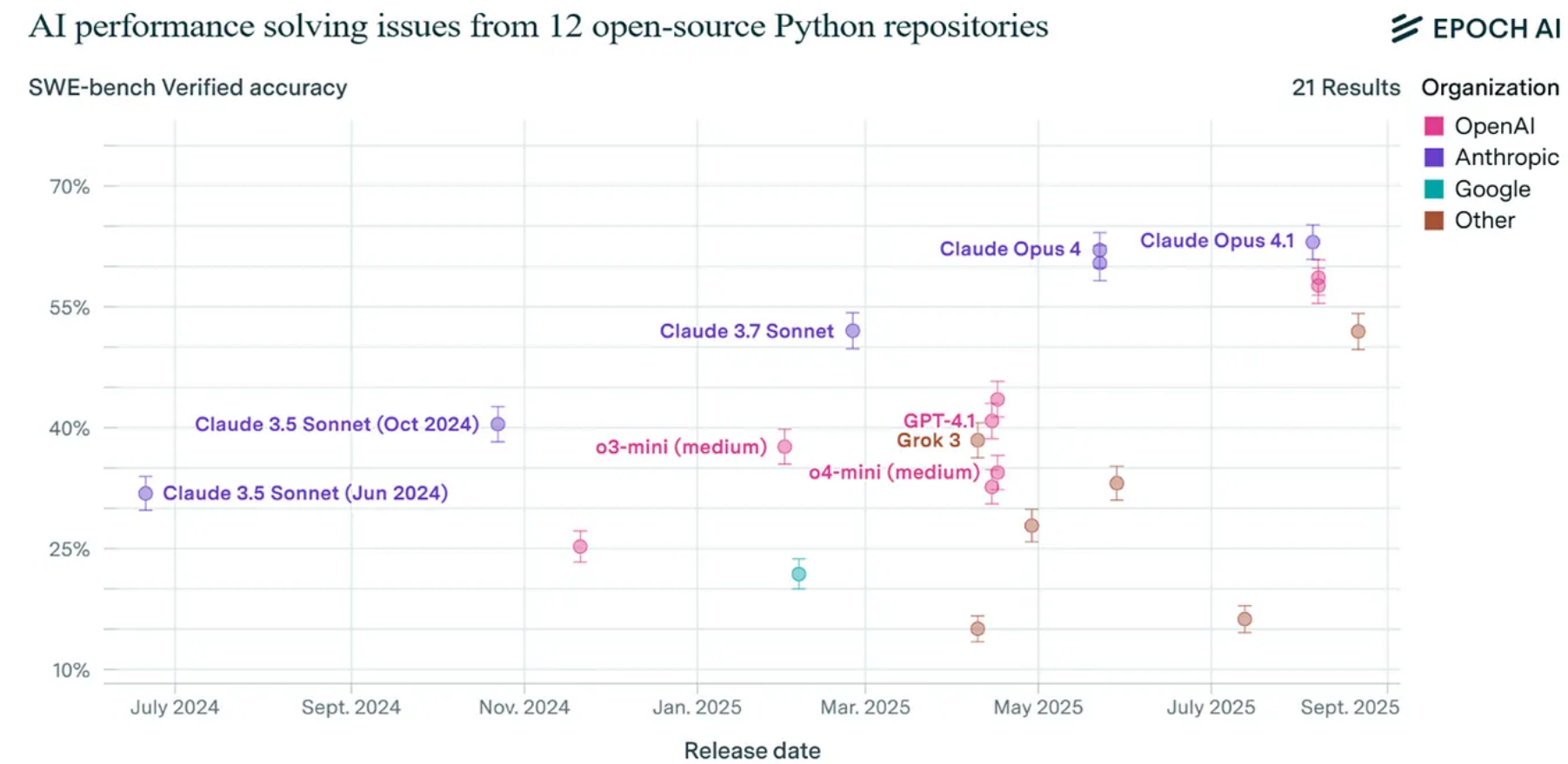
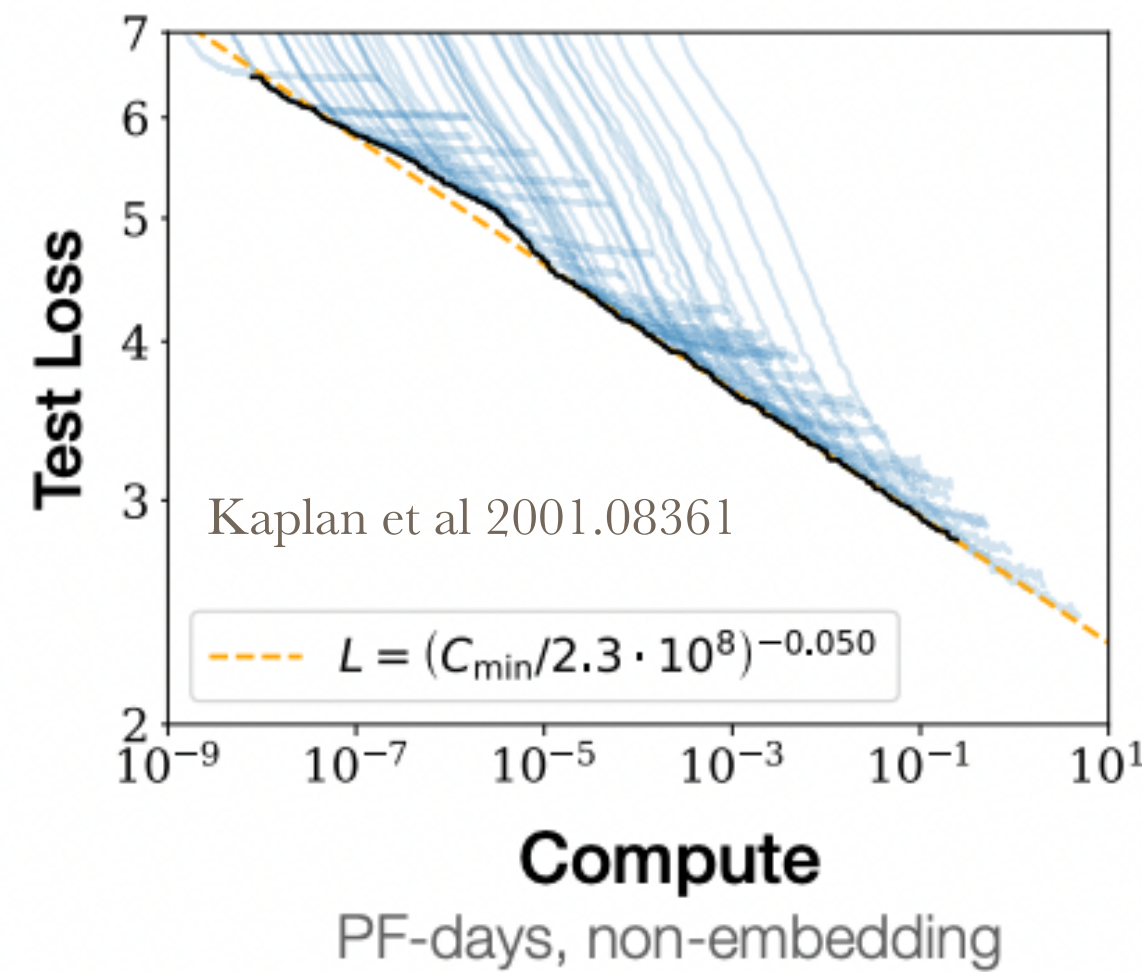
**David Shih
May 19, 2026**



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

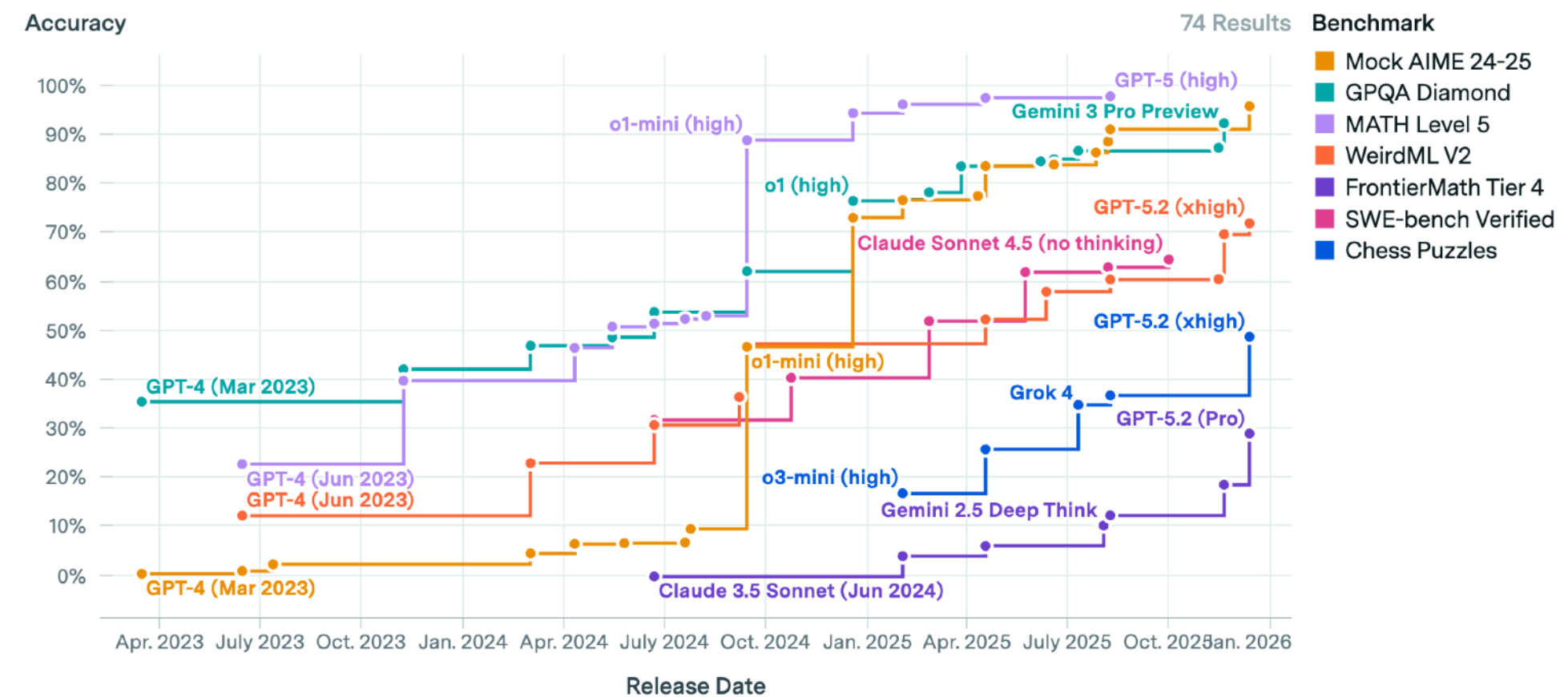
AI Revolution

We are witnessing a historic moment of technological advancement

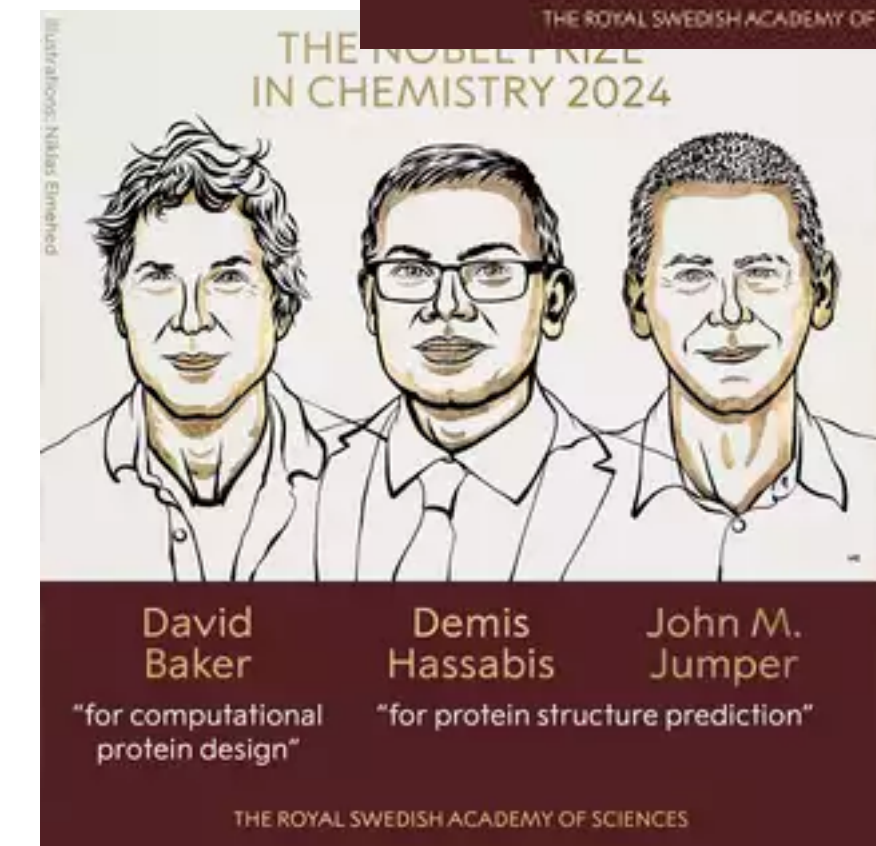
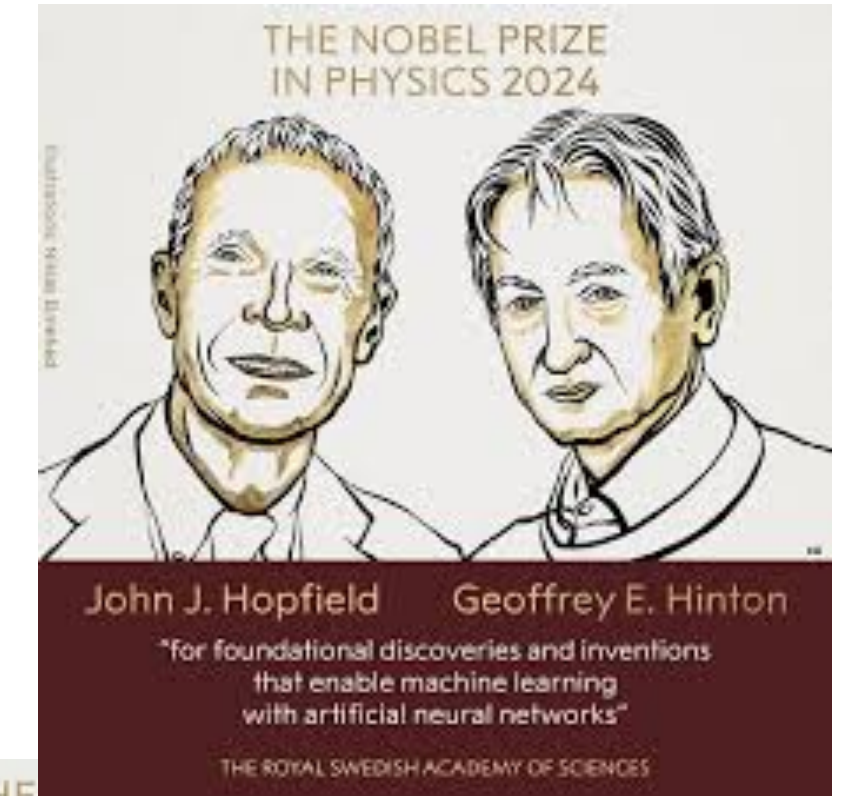


CC-BY

Frontier performance across benchmarks



CC-BY

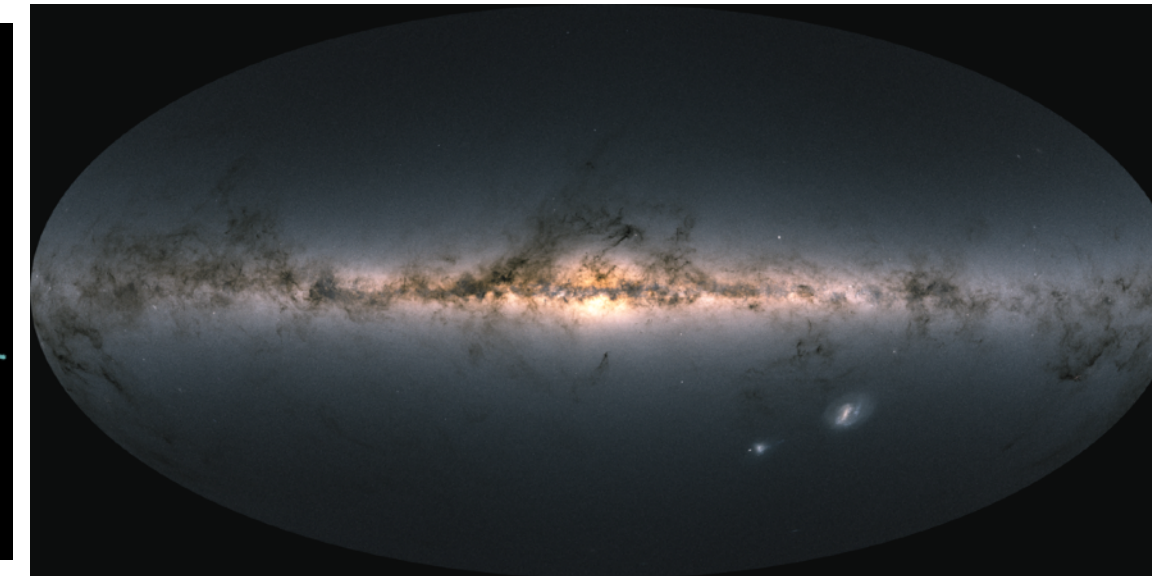
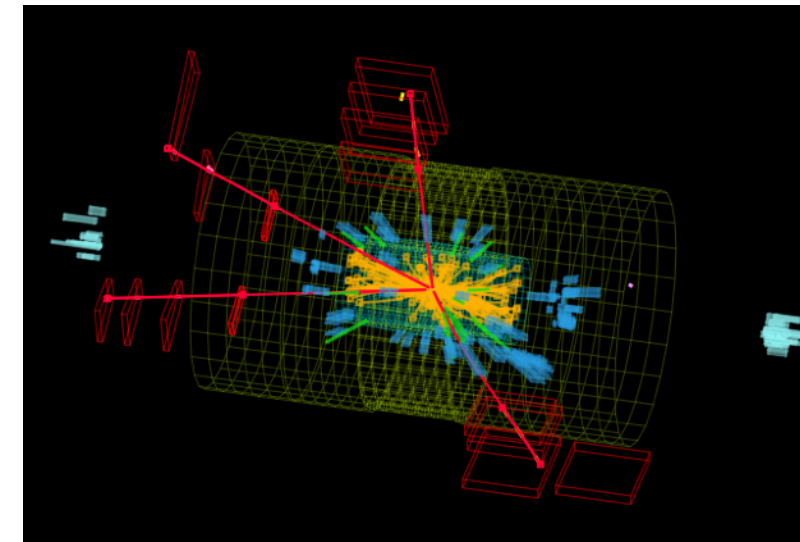
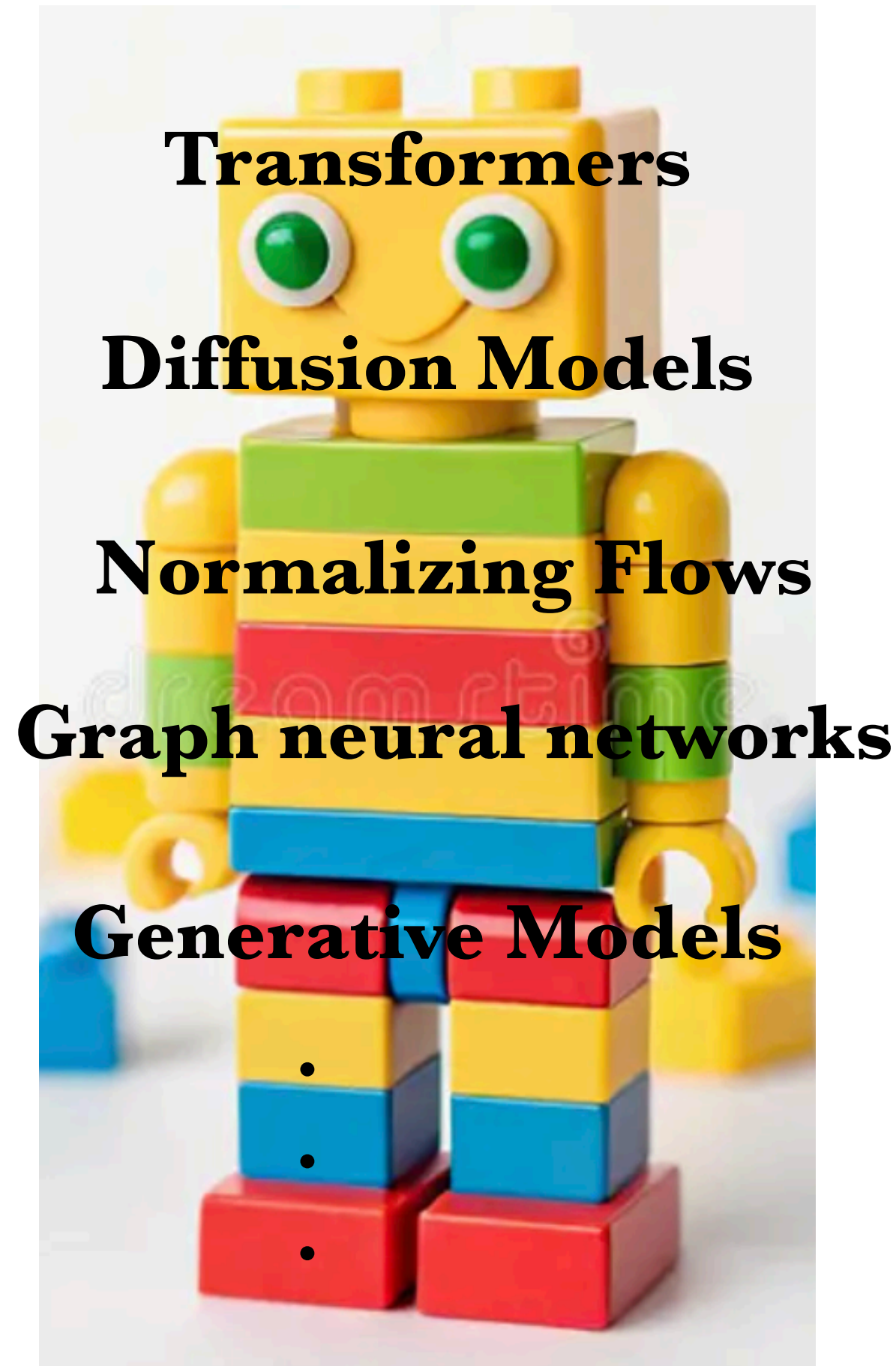


Google A.I. System Wins Gold Medal in International Math Olympiad

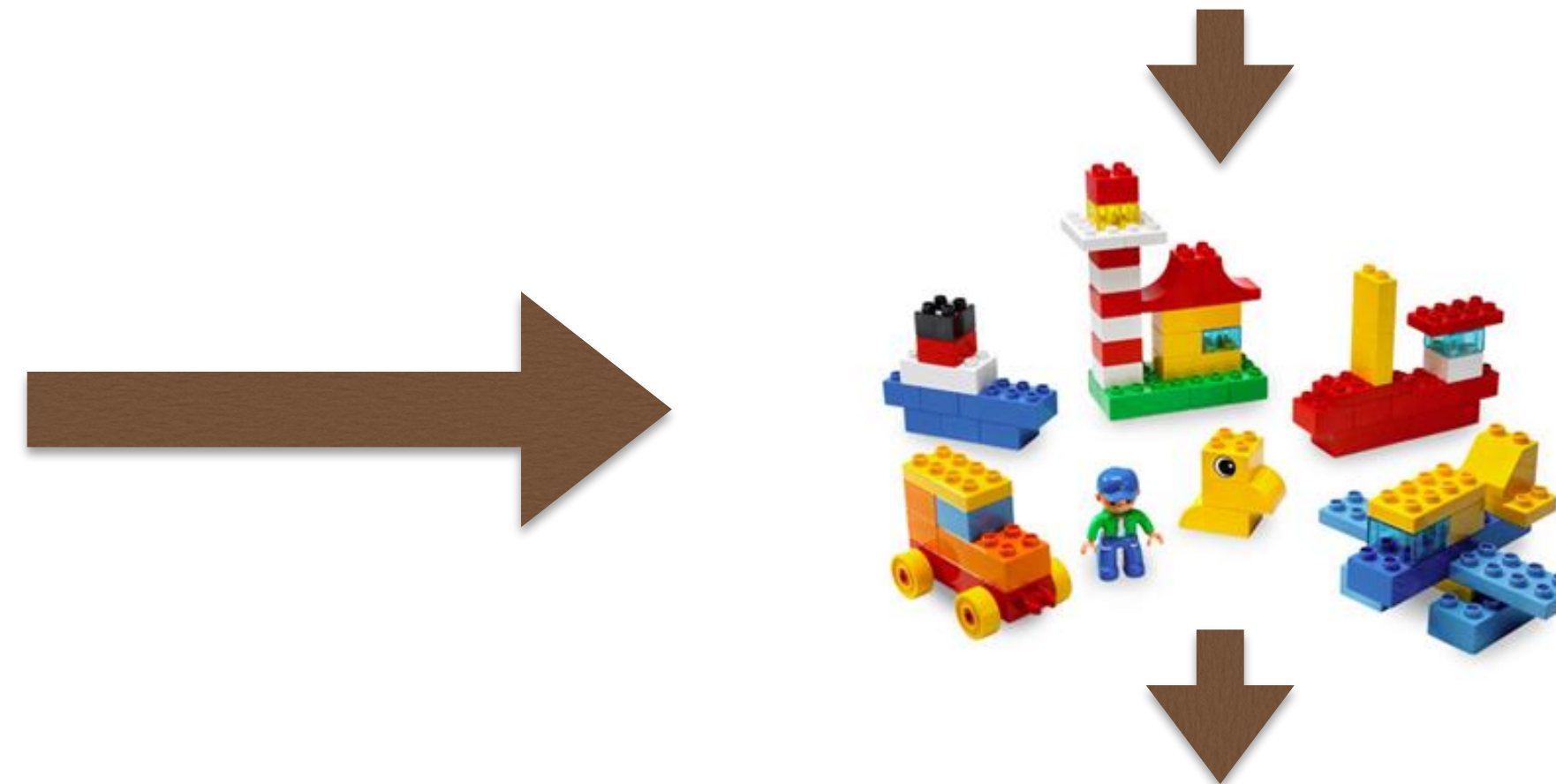
OpenAI said it, too, had built a system that achieved similar results.

ML4HEP

Exciting progress over the last ~10 years



$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + \chi_i y_{ij} \chi_j \phi + h.c. + |D_\mu \phi|^2 - V(\phi)$$



SciPost Physics Submission

Normalizing Flows for High-Dimensional Detector Simulations

Florian Ernst^{1,2}, Luigi Favaro^{1,3}, Claudius F. Romker^{1,2}

¹ Institut für Theoretische Physik,
² Experimental Physics Department,
³ CP3, Université catholique de Louvain (UCL), 1348 Louvain-la-Neuve, Belgium
⁴ Institut für Hochenergiephysik (HEPHY), ÖAW, Wien
⁵ NHETC, Department of Physics & Astronomy, Rutgers University, Piscataway, NJ 08854, USA

January 2024

Abstract

Whenever invertible generative networks are used for detector simulations, they can provide excellent performance. In this work, we investigate the performance of normalizing flows in shower simulations with increasing phase space dimensionality. By using spline transformations applied to the Cauchy distribution, we also employ a VAE to compare the performance of normalizing flows to generative models.

I. INTRODUCTION

The search for new physics at the Large Hadron Collider (LHC) requires powerful and model-agnostic anomaly detection methods [1, 2]. Weakly supervised approaches (BDTs) as a tool for their strength in tabular data analysis. Our results show that they not only perform well in tabular data analysis, but they are also robust to a large number of features. By using advanced gradient boosted decision trees in combination with an extended set of features, we significantly improve the performance of anomaly detection at the LHC. This advance is a crucial step towards the search for new physics.

PREPARED FOR SUBMISSION TO JCAP

Mapping Dark Matter Way using Normalized Gaia DR3

Sung Hak Lim^{a,b}, Eric Putney^b, Matthieu Shih^c

^aParticle Theory and Cosmology Group, Center for Theoretical Physics, Institute for Basic Science (IBS), 55 Expo-ro, Yuseong-gu, Daejeon 34126, Republic of Korea
^bNHETC, Department of Physics and Astronomy, Rutgers, the State University of New Jersey, Piscataway, NJ 08854, USA
^cInstitut für Experimentelle Physik, Universität Hamburg, 22761 Hamburg, Germany

CTPU-PTC-25-01

Residual ANODE

Ranit Das^{1,*}, Gregor Kasieczka^{2,†} and David Shih^{1,‡}

¹NHETC, Dept. of Physics and Astronomy, Rutgers University, Piscataway, NJ 08854, USA
²Institut für Experimentelle Physik, Universität Hamburg, 22761 Hamburg, Germany

We present R-ANODE, a new method for data-driven, model-agnostic resonant anomaly detection that raises the bar for both performance and interpretability. The key to R-ANODE is to enhance the inductive bias of the anomaly detection task by fitting a normalizing flow directly to the small and unknown signal component, while holding fixed a background model (also a normalizing flow) learned from sidebands. In doing so, R-ANODE is able to outperform all classifier-based, weakly-supervised approaches, as well as the previous ANODE method which fit a density estimator to all of the data in the signal region instead of just the signal. We show that the method works equally well whether the unknown signal fraction is learned or fixed, and is even robust to signal fraction misspecification. Finally, with the learned signal model we can sample and gain qualitative insights into the underlying anomaly, which greatly enhances the interpretability of resonant anomaly detection and offers the possibility of simultaneously discovering and characterizing the new physics that could be hiding in the data.

I. INTRODUCTION

Despite countless searches at the LHC [1–7], so far none have turned up any definitive evidence for new physics beyond the Standard Model yet. Since the vast majority of these searches have been model-specific, there has been increasing interest [8–10] in developing new, model-agnostic search strategies powered by modern machine learning in recent years. The hope is that these density estimators (normalizing flows in practice) on the signal region and sideband events, interpolates the latter into the SR, and constructs R_{optimal} by taking their ratio directly. Since this approach is solely based on unsupervised density estimation and does not involve any classifiers, it is technically an unsupervised approach to resonant anomaly detection.

It has been recognized [13, 17] that since density estimation is much more difficult than classification. An-

Agentic AI — a watershed moment

Claude Code, Codex — early 2025

Giving *agency* to LLMs — from question answering to reading/writing/executing code



Home News Sport Business Technology Health Culture Arts Travel Earth Audio Video Live

'Vibe coding' named word of the year by Collins Dictionary

☰ **The New York Times** 👤

Artificial Intelligence > Apple Settlement Job Losses Meta Sued Vetting A.I. Models A.I. Spending Rec

SHOP TALK

With 'Vibecoding,' A.I. Can Help Anyone Build an App

Bringing on artificial intelligence as a collaborator can make coding feel more accessible to those with little training in it, but there are trade-offs.

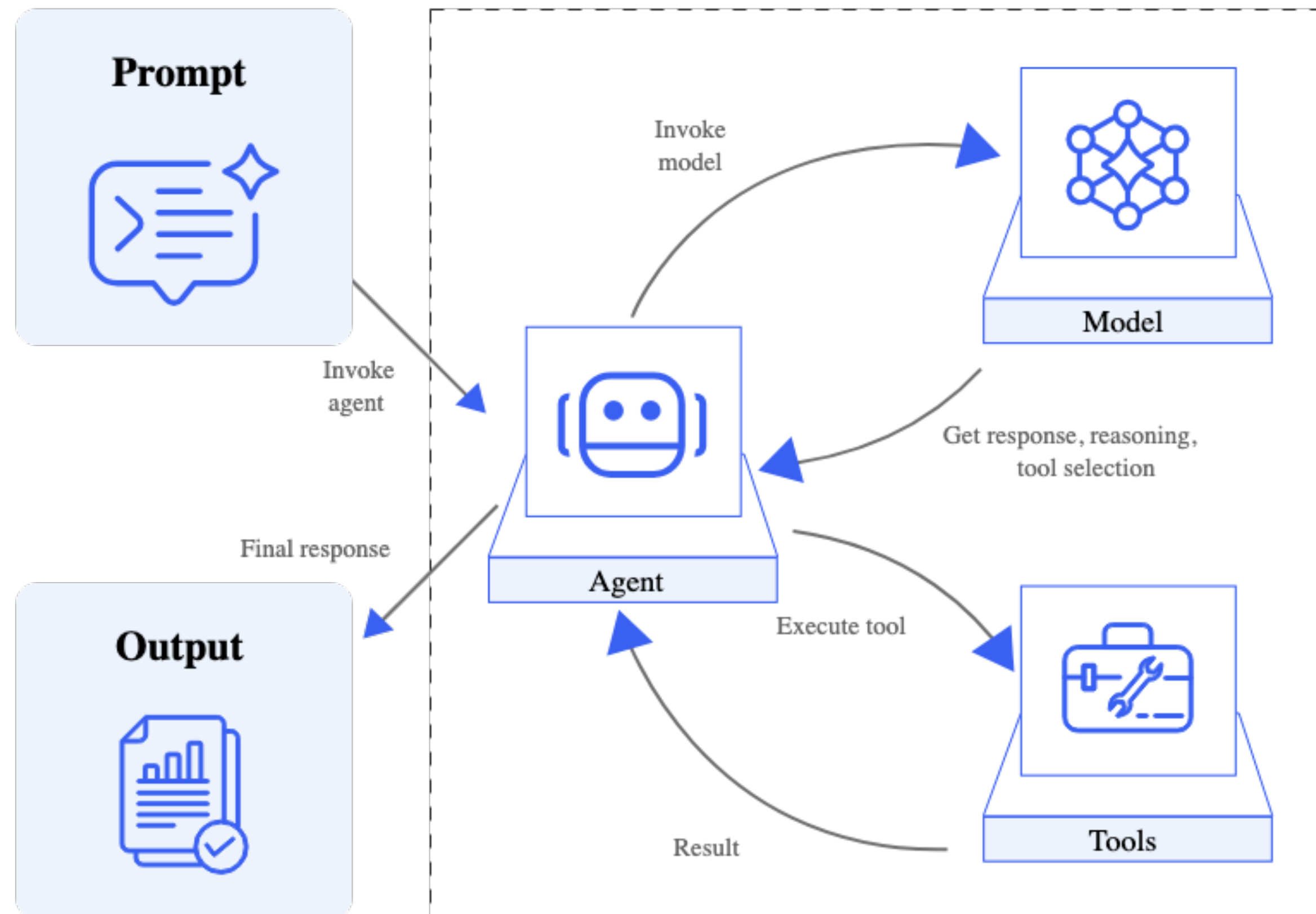
tom's **HARDWARE** [Follow](#)

Claude-powered AI coding agent deletes entire company database in 9 seconds — backups zapped, after Cursor tool powered by Anthropic's Claude goes rogue

TECH

20,000 job cuts at Meta, Microsoft raise concern that AI-driven labor crisis is here

Agentic AI explained



ZBrain

- Agentic AI consists of
 - LLM backend (Claude, GPT, Gemini ...)
 - a set of ***tools*** (bash, python, grep, text editor, ...)
 - a ***harness/scaffold*** (traditional code infrastructure managing model responses and tool use)
- Examples include Claude Code, Codex, Antigravity, ...

How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

- pheno tasks [ArgoLOOM, HEPTAPOD, MadAgents, CoLLM, FERMIACC, ColliderAgent/Magnus...]

The FERMIACC: Agents for Particle Theory

- experimental data analysis

AI Agents Can Already Autonomously Perform Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

- pheno tasks [ArgoLOOM, HEPTAPOD, MadAgents, CoLLM, FERMIACC, ColliderAgent/Magnus...]

- experimental data analysis

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

- Potential pitfalls of this line of work:

How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

- pheno tasks [ArgoLOOM, HEPTAPOD, MadAgents, CoLLM, FERMIACC, ColliderAgent/Magnus...]

- experimental data analysis

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

- Potential pitfalls of this line of work:

- LLM models improving rapidly — method papers quickly outdated?

How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

- pheno tasks [ArgoLOOM, HEPTAPOD, MadAgents, CoLLM, FERMIACC, ColliderAgent/Magnus...]

- experimental data analysis

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

- Potential pitfalls of this line of work:

- LLM models improving rapidly — method papers quickly outdated?
- custom harnesses — hard to compete with frontier companies

How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection
- pheno tasks [ArgoLOOM, HEPTAPOD, MadAgents, CoLLM, FERMIACC, ColliderAgent/Magnus...]
- experimental data analysis

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

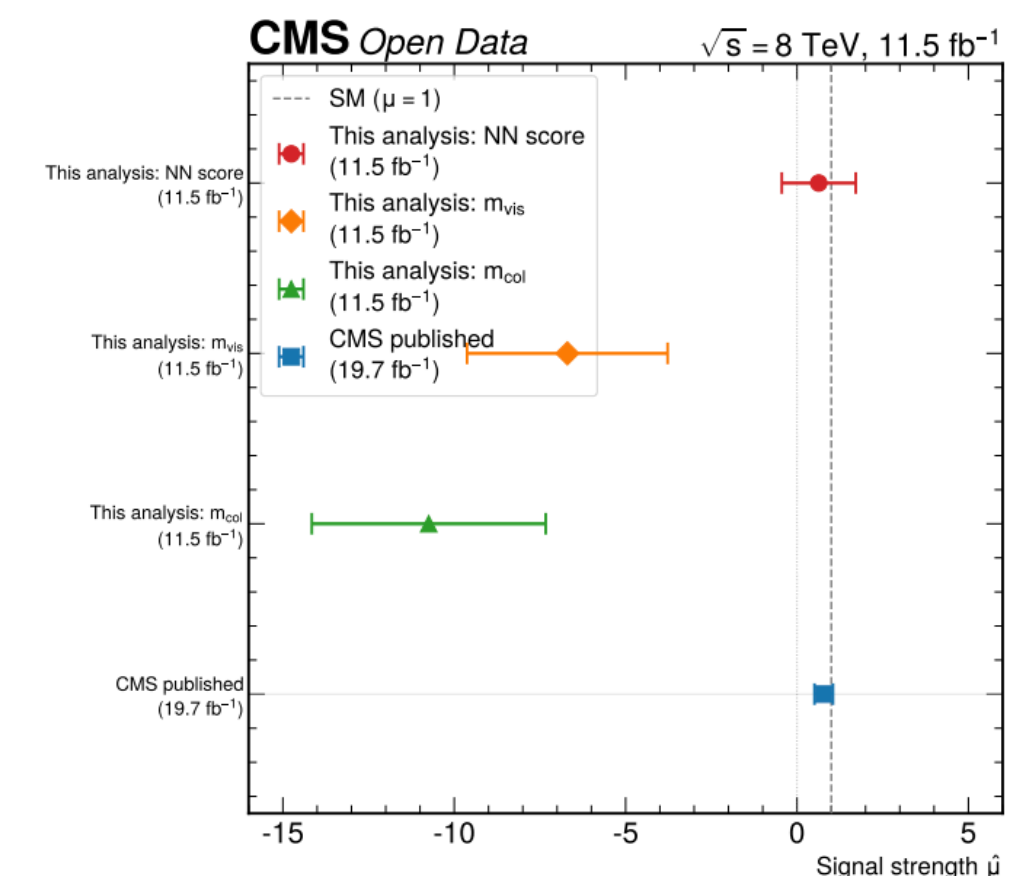
^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

2603.20179v2.pdf
Page 1 of 541

- Potential pitfalls of this line of work:

- LLM models improving rapidly — method papers quickly outdated?
- custom harnesses — hard to compete with frontier companies
- automation vs quality control — risk of AI slop



How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

- pheno task
CoLLM, FE

- experimental data analysis

Agents of Discovery

Increasing interest in rigorous benchmarking
— potentially more lasting value?

ukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform
 Experimental High Energy Physics

Eric A. Moreno^{*†1,2}, Samuel Bright-Thonney^{*†1,2}, Andrzej Novak^{*§1,2}, Dolores Garcia^{#3},
 and Philip Harris^{b1,2}

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

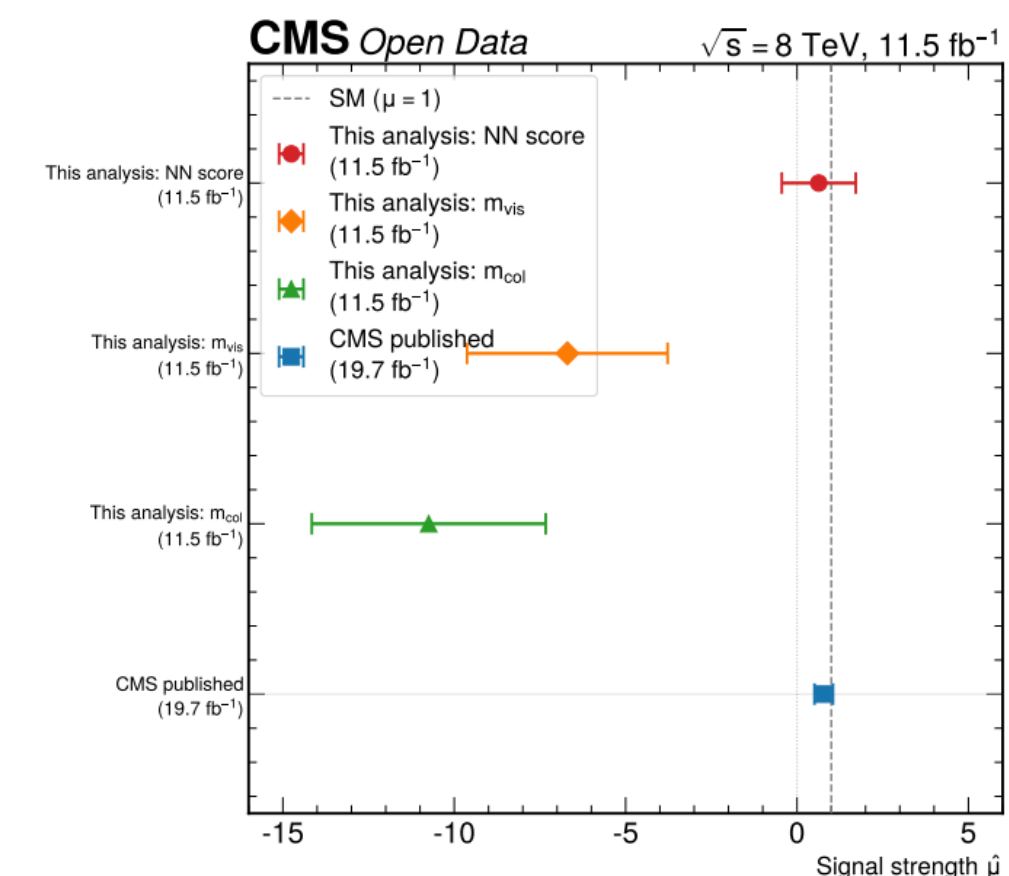
^aDepartment of Physics, University of California, Santa Barbara, CA 93106, USA

^bKavli Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

2603.20179v2.pdf
 Page 1 of 541

- Potential pitfalls of this line of work:

- LLM models improving rapidly — method papers quickly outdated?
- custom harnesses — hard to compete with frontier companies
- automation vs quality control — risk of AI slop



How are people using agentic AI?

- Particle physicists are exploring building custom harnesses for automating

- anomaly detection

- pheno task
CoLLM, FE

- experimental data analysis

Agents of Discovery

Increasing interest in rigorous benchmarking
— potentially more lasting value?

ukas²,

The FERMIACC: Agents for Particle Theory

AI Agents Can Already Autonomously Perform
Experimental High Energy Physics

Prateek Agrawal,^a Nathaniel Craig,^{a,b} Amalia Madden,^b and Iñigo Valenzuela Lomber^a

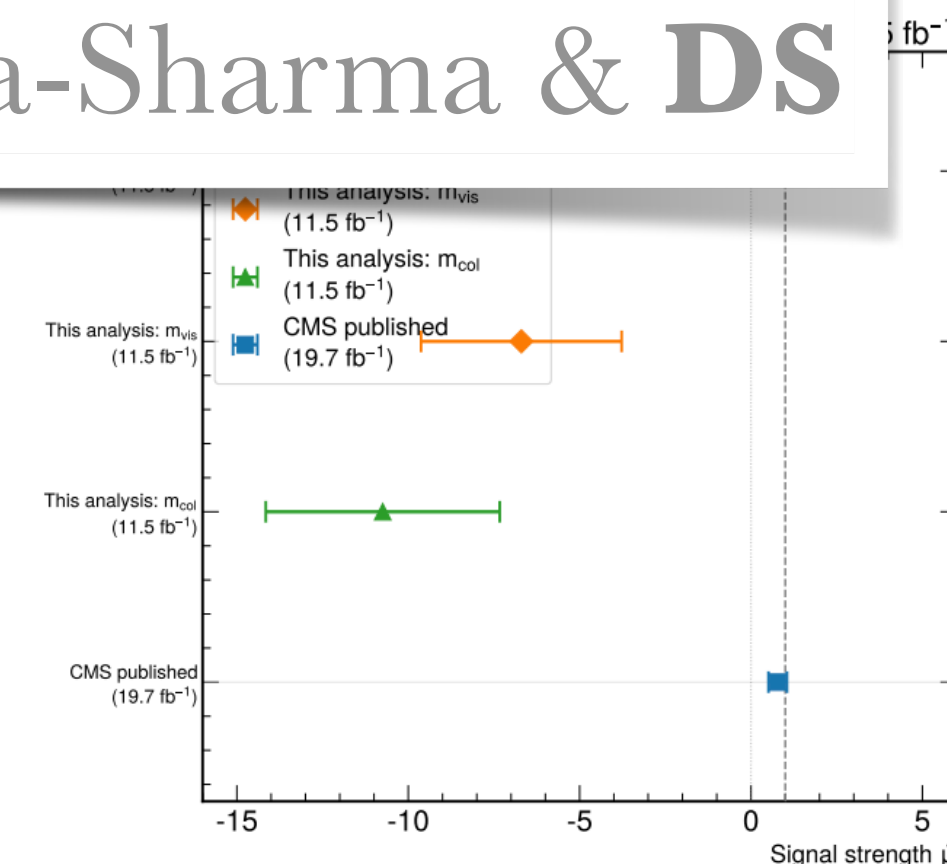
^aUniversity of California, Santa Barbara, CA 93106, USA
^bTheoretical Physics, University of California, Santa Barbara, CA 93106, USA

See: ColliderBench [2605.13950]

Faroughy, Palacios Schweitzer, Pang, Mishra-Sharma & DS

- Potential pitfalls of this line of work.

- LLM models improving rapidly — method papers quickly outdated?
- custom harnesses — hard to compete with frontier companies
- automation vs quality control — risk of AI slop



How are people using agentic AI?

As research assistants (~graduate student) — “vibe physics”

arXiv > hep-th > arXiv:2603.11164 Search...
Help | Adv

High Energy Physics – Theory

[Submitted on 11 Mar 2026 (v1), last revised 11 Apr 2026 (this version, v2)]

Learning to Unscramble: Simplifying Symbolic Expressions via Self-Supervised Oracle Trajectories

David Shih

We present a new self-supervised machine learning approach for symbolic simplification of complex mathematical expressions. Training data is generated by scrambling simple expressions and recording the inverse operations, creating oracle trajectories that provide both goal states and explicit paths to reach them. A permutation-equivariant, transformer-based policy network is then trained on this data step-wise to predict the oracle action given the input expression. We demonstrate this approach on two problems in high-energy physics: dilogarithm reduction and spinor-helicity scattering amplitude simplification. In both cases, our trained policy network achieves near perfect solve rates across a wide range of difficulty levels, substantially outperforming prior approaches based on reinforcement learning and end-to-end regression. When combined with contrastive grouping and beam search, our model achieves a 100% full simplification rate on a representative selection of 5-point gluon tree-level amplitudes in Yang-Mills theory, including expressions with over 200 initial terms.

Comments: 14 pages, 6 figures, 2 tables; work done in collaboration with Claude Code; v2: refs added

arXiv > hep-ph > arXiv:2604.05034 Search...
Help | Adv

High Energy Physics – Phenomenology

[Submitted on 6 Apr 2026]

Learning to Unscramble Feynman Loop Integrals with SAILIR

David Shih

Integration-by-parts (IBP) reduction of Feynman integrals to master integrals is a key computational bottleneck in precision calculations in high-energy physics. Traditional approaches based on the Laporta algorithm require solving large systems of equations, leading to memory consumption that grows rapidly with integral complexity. We present SAILIR (Self-supervised AI for Loop Integral Reduction), a new machine learning approach in which a transformer-based classifier guides the reduction of integrals one step at a time in a fully online fashion. The classifier is trained in an entirely self-supervised manner on synthetic data generated by a scramble/unscramble procedure: known reduction identities are applied in reverse to build expressions of increasing complexity, and the classifier learns to undo these steps. When combined with beam search and a highly parallelized, asynchronous, single-episode reduction strategy, SAILIR can reduce integrals of arbitrarily high weight with bounded memory. We benchmark SAILIR on the two-loop triangle-box topology, comparing against the state-of-the-art IBP reduction code Kira across 16 integrals of varying complexity. While SAILIR is slower in wall-clock time, its per-worker memory consumption remains approximately flat regardless of integral complexity, in contrast to Kira whose memory grows rapidly with complexity. For the most complex integrals considered here, SAILIR uses only 40% of the memory of Kira while achieving comparable reduction times. This demonstrates a fundamentally new paradigm for IBP reduction in which the memory bottleneck of Laporta-based approaches could be entirely overcome, potentially opening the door to precision calculations that are currently intractable.

Comments: 16 pages, 3 figures, 5 tables, work done in collaboration with Claude Code

How are people using agentic AI?

As research assistants (~graduate student) — “vibe physics”



arXiv > hep-th > arXiv:2603.11164

Search...
Help | Adv

High Energy Physics – Theory

[Submitted on 11 Mar 2026 (v1), last revised 11 Apr 2026 (this version, v2)]

Learning to Unscramble: Simplifying Symbolic Expressions via Self-Supervised Oracle Trajectories

David Shih

We present a new self-supervised machine learning approach for symbolic simplification of complex mathematical expressions. Training data is generated by scrambling simple expressions and recording the inverse operations, creating oracle trajectories that provide both goal states and explicit paths to reach them. A permutation-equivariant, transformer-based policy network is then trained on this data step-wise to predict the oracle action given the input expression. We demonstrate this approach on two problems in high-energy physics: dilogarithm reduction and spinor-helicity scattering amplitude simplification. In both cases, our trained policy network achieves near perfect solve rates across a wide range of difficulty levels, substantially outperforming prior approaches based on reinforcement learning and end-to-end regression. When combined with contrastive grouping and beam search, our model achieves a 100% full simplification rate on a representative selection of 5-point gluon tree-level amplitudes in Yang-Mills theory, including expressions with over 200 initial terms.

Comments: 14 pages, 6 figures, 2 tables; work done in collaboration with Claude Code; v2: refs added

Risky projects
— branching out into new subfields

Approx 2 months of work
— compared to 6-12 months for a student



arXiv > hep-ph > arXiv:2604.05034

Search...
Help | Ad

High Energy Physics – Phenomenology

[Submitted on 6 Apr 2026]

Learning to Unscramble Feynman Loop Integrals with SAILIR

David Shih

Integration-by-parts (IBP) reduction of Feynman integrals to master integrals is a key computational bottleneck in precision calculations in high-energy physics. Traditional approaches based on the Laporta algorithm require solving large systems of equations, leading to memory consumption that grows rapidly with integral complexity. We present SAILIR (Self-supervised AI for Loop Integral Reduction), a new machine learning approach in which a transformer-based classifier guides the reduction of integrals one step at a time in a fully online fashion. The classifier is trained in an entirely self-supervised manner on synthetic data generated by a scramble/unscramble procedure: known reduction identities are applied in reverse to build expressions of increasing complexity, and the classifier learns to undo these steps. When combined with beam search and a highly parallelized, asynchronous, single-episode reduction strategy, SAILIR can reduce integrals of arbitrarily high weight with bounded memory. We benchmark SAILIR on the two-loop triangle-box topology, comparing against the state-of-the-art IBP reduction code Kira across 16 integrals of varying complexity. While SAILIR is slower in wall-clock time, its per-worker memory consumption remains approximately flat regardless of integral complexity, in contrast to Kira whose memory grows rapidly with complexity. For the most complex integrals considered here, SAILIR uses only 40% of the memory of Kira while achieving comparable reduction times. This demonstrates a fundamentally new paradigm for IBP reduction in which the memory bottleneck of Laporta-based approaches could be entirely overcome, potentially opening the door to precision calculations that are currently intractable.

Comments: 16 pages, 3 figures, 5 tables, work done in collaboration with Claude Code

How are people using agentic AI?

As research assistants (~graduate student) — “vibe physics”



arXiv > hep-th > arXiv:2603.11164

Search...
Help | Adv

High Energy Physics - Theory

[Submitted on 11 Mar 2026 (v1), last revised 11 Apr 2026 (this version, v2)]

Learning to Unscramble: Simplifying Symbolic Expressions via Self-Supervised Oracle Trajectories

David Shih

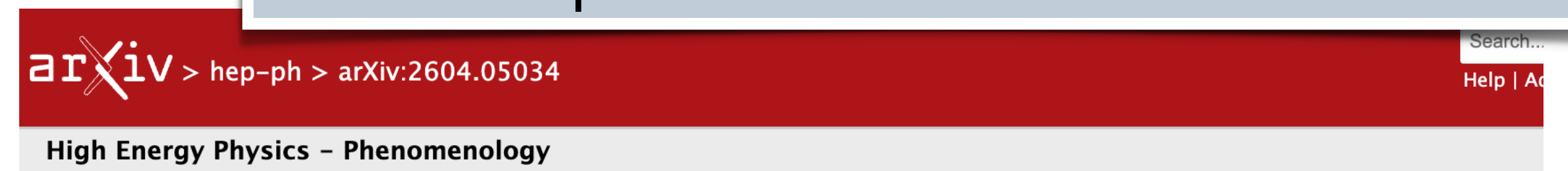
We present a new self-supervised machine learning approach for symbolic simplification of complex mathematical expressions. Training data is generated by scrambling simple expressions and recording the inverse operations, creating oracle trajectories that provide both goal states and explicit paths to reach them. A permutation-equivariant, transformer-based policy network is then trained on this data step-wise to predict the oracle action given the input expression. We demonstrate this approach on two problems in high-energy physics: dilogarithm reduction and spinor-helicity scattering amplitude simplification. In both cases, our trained policy network achieves near perfect solve rates across a wide range of difficulty levels, substantially outperforming prior approaches based on reinforcement learning and end-to-end regression. When combined with contrastive grouping and beam search, our model achieves a 100% full simplification rate on a representative selection of 5-point gluon tree-level amplitudes in Yang-Mills theory, including expressions with over 200 initial terms.

Comments: 14 pages, 6 figures, 2 tables; work done in collaboration with Claude Code; v2: refs added

see also Schwartz 2601.02484, Guevara et al 2602.12176, Zhang 2604.27050

Risky projects
— branching out into new subfields

Approx 2 months of work
— compared to 6-12 months for a student



arXiv > hep-ph > arXiv:2604.05034

Search...
Help | Ad

High Energy Physics - Phenomenology

[Submitted on 6 Apr 2026]

Learning to Unscramble Feynman Loop Integrals with SAILIR

David Shih

Integration-by-parts (IBP) reduction of Feynman integrals to master integrals is a key computational bottleneck in precision calculations in high-energy physics. Traditional approaches based on the Laporta algorithm require solving large systems of equations, leading to memory consumption that grows rapidly with integral complexity. We present SAILIR (Self-supervised AI for Loop Integral Reduction), a new machine learning approach in which a transformer-based classifier guides the reduction of integrals one step at a time in a fully online fashion. The classifier is trained in an entirely self-supervised manner on synthetic data generated by a scramble/unscramble procedure: known reduction identities are applied in reverse to build expressions of increasing complexity, and the classifier learns to undo these steps. When combined with beam search and a highly parallelized, asynchronous, single-episode reduction strategy, SAILIR can reduce integrals of arbitrarily high weight with bounded memory. We benchmark SAILIR on the two-loop triangle-box topology, comparing against the state-of-the-art IBP reduction code Kira across 16 integrals of varying complexity. While SAILIR is slower in wall-clock time, its per-worker memory consumption remains approximately flat regardless of integral complexity, in contrast to Kira whose memory grows rapidly with complexity. For the most complex integrals considered here, SAILIR uses only 40% of the memory of Kira while achieving comparable reduction times. This demonstrates a fundamentally new paradigm for IBP reduction in which the memory bottleneck of Laporta-based approaches could be entirely overcome, potentially opening the door to precision calculations that are currently intractable.

Comments: 16 pages, 3 figures, 5 tables, work done in collaboration with Claude Code

Outline

- Learning to unscramble: general framework
- **Application 1:** Dilog identities [skip for lack of time]
- **Application 2:** Tree-level YM scattering amplitudes
- **Application 3:** IBP reduction of Feynman loop integrals
- Lessons learned
- Bonus material: **ColliderBench** — a new LHC benchmark for LLM agents

Learning to Unscramble: general framework

- Goal: ML methods for mathematical (symbolic) simplification

$$E_{\text{complicated}} = X_1 + X_2 + X_3 + X_4 + \dots$$



mathematical identities

$$E_{\text{simple}} = Y_1$$

- Test beds: dilog sums, tree-level YM scattering amplitudes

(Dersy, Schwartz, Zhang [2206.04115](#); Cheung, Dersy, Schwartz [2408.04720](#))

Training data



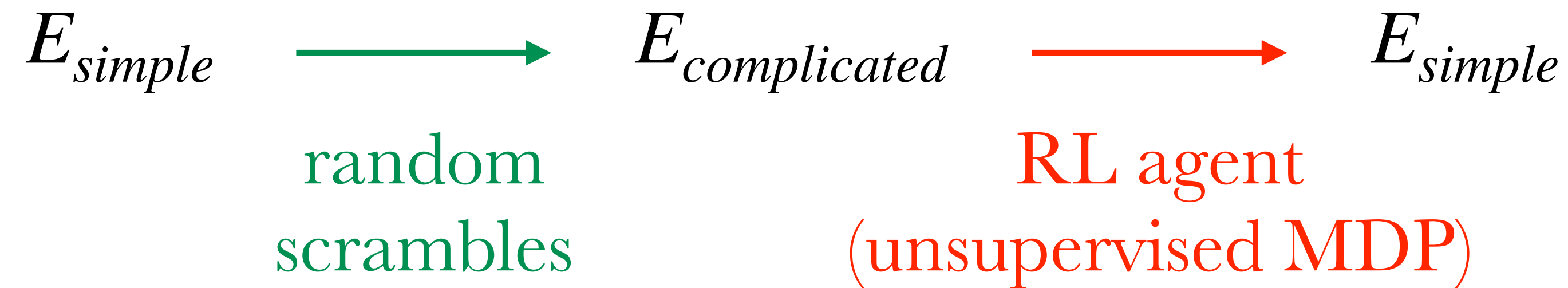
- ML methods need training data
- **Key insight:** can think of complicated expressions as *scrambles* of simple ones by mathematical identities.
- Process of simplification is then *unscrambling*. Dersy, Schwartz, Zhang [2206.04115](#)
- *Strong analogy with Rubik's cubes!*
- Can cheaply create training data by applying random identities to simple expressions



Prior approaches: RL

Dersy, Schwartz, Zhang [2206.04115](#)

MDP: Markov Decision Process
action classifier $\pi(a | s)$,
depends only on previous state



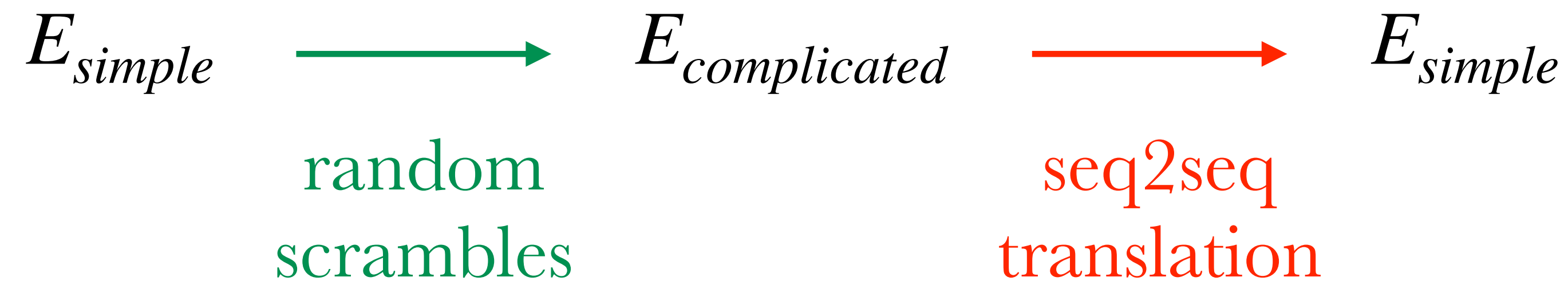
RL (PPO, TRPO) could not beat classical algorithm on even very simple simplification tasks (small dilog sums simplifying to zero).

Challenges: very sparse reward signal, non-monotonic reduction paths

Successful RL is very compute intensive and requires a lot of fine tuning.

Prior approaches: end-to-end regression

Cheung, Dersy, Schwartz [2408.04720](#)



End-to-end regression with a seq2seq transformer achieved decent performance (80-90%) for dilog sums and scattering amplitudes

Drawback: mathematical correctness not guaranteed, transformer can hallucinate simpler but incorrect expressions!

Our approach: “Learning to Unscramble”

Key insight: random scrambles can be reversed, step-by-step.

—> A wealth of simplification training data!!!

$$-\frac{\langle 13 \rangle [14][23]}{\langle 12 \rangle \langle 14 \rangle [12]^2}$$

$$\Downarrow \text{mom}^2: \langle 13 \rangle [31] + \langle 14 \rangle [41] + \langle 34 \rangle [43] = 0$$

$$\frac{\langle 14 \rangle [14]^2 [23] + \langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{mom}^2: \langle 12 \rangle [21] - \langle 34 \rangle [43] = 0$$

Idea from *ML for Rubik’s cube literature*:
Takano 2106.03157

$$\frac{-\langle 12 \rangle [12][14][23] + \langle 14 \rangle [14]^2 [23] + 2\langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{Schouten: } [24][13] + [21][34] + [23][41] = 0$$



$$\frac{1}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13][23]} \left(\langle 14 \rangle [12]^2 [34]^2 + \langle 14 \rangle [13]^2 [24]^2 + \langle 12 \rangle [12]^2 [23][34] \right. \\ \left. - 2\langle 34 \rangle [12][23][34]^2 - \langle 12 \rangle [12][13][23][24] - 2\langle 14 \rangle [12][13][24][34] + 2\langle 34 \rangle [13][23][24][34] \right)$$

Our approach: “Learning to Unscramble”

Key insight: random scrambles can be reversed, step-by-step.

—> A wealth of simplification training data!!!

Idea from *ML for Rubik’s cube literature*:
Takano 2106.03157



$$-\frac{\langle 13 \rangle [14][23]}{\langle 12 \rangle \langle 14 \rangle [12]^2}$$

$$\Downarrow \text{mom}^2: \langle 13 \rangle [31] + \langle 14 \rangle [41] + \langle 34 \rangle [43] = 0$$

$$\frac{\langle 14 \rangle [14]^2 [23] + \langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{mom}^2: \langle 12 \rangle [21] - \langle 34 \rangle [43] = 0$$

$$\frac{-\langle 12 \rangle [12][14][23] + \langle 14 \rangle [14]^2 [23] + 2\langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{Schouten: } [24][13] + [21][34] + [23][41] = 0$$

$$\frac{1}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13][23]} \left(\langle 14 \rangle [12]^2 [34]^2 + \langle 14 \rangle [13]^2 [24]^2 + \langle 12 \rangle [12]^2 [23][34] \right. \\ \left. - 2\langle 34 \rangle [12][23][34]^2 - \langle 12 \rangle [12][13][23][24] - 2\langle 14 \rangle [12][13][24][34] + 2\langle 34 \rangle [13][23][24][34] \right)$$

Our approach: “Learning to Unscramble”

Key insight: random scrambles can be reversed, step-by-step.

—> A wealth of simplification training data!!!

Idea from *ML for Rubik’s cube literature*:
Takano 2106.03157



$$-\frac{\langle 13 \rangle [14][23]}{\langle 12 \rangle \langle 14 \rangle [12]^2}$$

$$\Downarrow \text{mom}^2: \langle 13 \rangle [31] + \langle 14 \rangle [41] + \langle 34 \rangle [43] = 0$$

$$\frac{\langle 14 \rangle [14]^2 [23] + \langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{mom}^2: \langle 12 \rangle [21] - \langle 34 \rangle [43] = 0$$

$$\frac{-\langle 12 \rangle [12][14][23] + \langle 14 \rangle [14]^2 [23] + 2\langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{Schouten: } [24][13] + [21][34] + [23][41] = 0$$

$$\frac{1}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13][23]} \left(\langle 14 \rangle [12]^2 [34]^2 + \langle 14 \rangle [13]^2 [24]^2 + \langle 12 \rangle [12]^2 [23][34] \right. \\ \left. - 2\langle 34 \rangle [12][23][34]^2 - \langle 12 \rangle [12][13][23][24] - 2\langle 14 \rangle [12][13][24][34] + 2\langle 34 \rangle [13][23][24][34] \right)$$

Our approach: “Learning to Unscramble”

Key insight: random scrambles can be reversed, step-by-step.

—> A wealth of simplification training data!!!

Idea from *ML for Rubik’s cube literature*:
Takano 2106.03157



$$-\frac{\langle 13 \rangle [14][23]}{\langle 12 \rangle \langle 14 \rangle [12]^2}$$

$$\Downarrow \text{mom}^2: \langle 13 \rangle [31] + \langle 14 \rangle [41] + \langle 34 \rangle [43] = 0$$

$$\frac{\langle 14 \rangle [14]^2 [23] + \langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

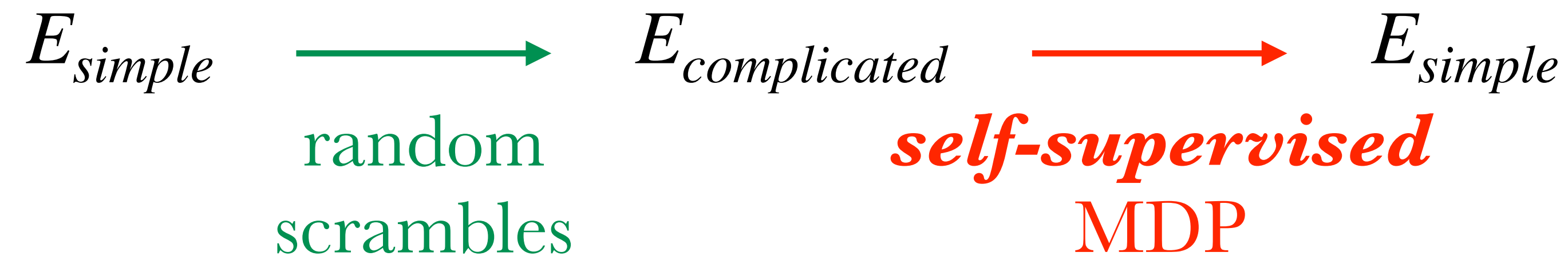
$$\Downarrow \text{mom}^2: \langle 12 \rangle [21] - \langle 34 \rangle [43] = 0$$

$$\frac{-\langle 12 \rangle [12][14][23] + \langle 14 \rangle [14]^2 [23] + 2\langle 34 \rangle [14][23][34]}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13]}$$

$$\Downarrow \text{Schouten: } [24][13] + [21][34] + [23][41] = 0$$

$$\frac{1}{\langle 12 \rangle \langle 14 \rangle [12]^2 [13][23]} \left(\langle 14 \rangle [12]^2 [34]^2 + \langle 14 \rangle [13]^2 [24]^2 + \langle 12 \rangle [12]^2 [23][34] \right. \\ \left. - 2\langle 34 \rangle [12][23][34]^2 - \langle 12 \rangle [12][13][23][24] - 2\langle 14 \rangle [12][13][24][34] + 2\langle 34 \rangle [13][23][24][34] \right)$$

Our approach: “Learning to Unscramble”



Instead of unsupervised MDP of RL, train a **self-supervised MDP** on **reversed random scramble sequences**

Learns to predict the best action from unscrambling steps

**Addresses limitations of both previous approaches:
highly dense reward signal and no hallucinations!**

Application 2: scattering amplitudes

[DS 2603.11164]

Application 2: scattering amplitudes

[DS 2603.11164]

- Example of 5pt YM tree amplitude computed with Feynman diagrams (from Cheung, Dersy, Schwartz)

- Known to simplify to **Parke-Taylor amplitude**

- through repeated application of identities:

Schouten identity :
$$\begin{cases} \langle ij \rangle \langle kl \rangle = \langle il \rangle \langle kj \rangle + \langle ik \rangle \langle jl \rangle \\ [ij][kl] = [il][kj] + [ik][jl] \end{cases}$$

momentum squared :
$$\sum_{\substack{i < j \\ (i,j) \in S_1^n}} \langle ij \rangle [ji] = \sum_{\substack{k < l \\ (k,l) \in S_2^n}} \langle kl \rangle [lk]$$

momentum conservation :
$$\sum_{j=1}^n \langle ij \rangle [jk] = 0 \quad \forall i, k$$

$$A_5 = \frac{\langle 12 \rangle^3}{\langle 13 \rangle \langle 23 \rangle \langle 34 \rangle \langle 45 \rangle \langle 51 \rangle} + \dots$$

$\frac{5}{[45]}$
 $\frac{1}{[45]}$

Scattering amplitudes: setup

- Start from 1-3 terms (shared denominator, same mass dimension and little-group scaling)
- Scramble with identities; also multiplication by unity and addition by zero
- Record correct unscrambling action at every step

$$\overline{\mathcal{M}} \rightarrow \overline{\mathcal{M}} \frac{\langle 12 \rangle}{\langle 12 \rangle} \rightarrow \overline{\mathcal{M}} \frac{\langle 13 \rangle \langle 25 \rangle - \langle 15 \rangle \langle 23 \rangle}{\langle 12 \rangle \langle 35 \rangle}$$

$$\overline{\mathcal{M}} + \frac{[34] - [34]_{\mathcal{F}}}{\mathcal{D}} \rightarrow \overline{\mathcal{M}} + \frac{[14][35] - [13][45] - [15][34]_{\mathcal{F}}}{\mathcal{D}[15]}$$

Training data

500k scramble trajectories.
—> ~1M single steps
(per n-pt).

Test data

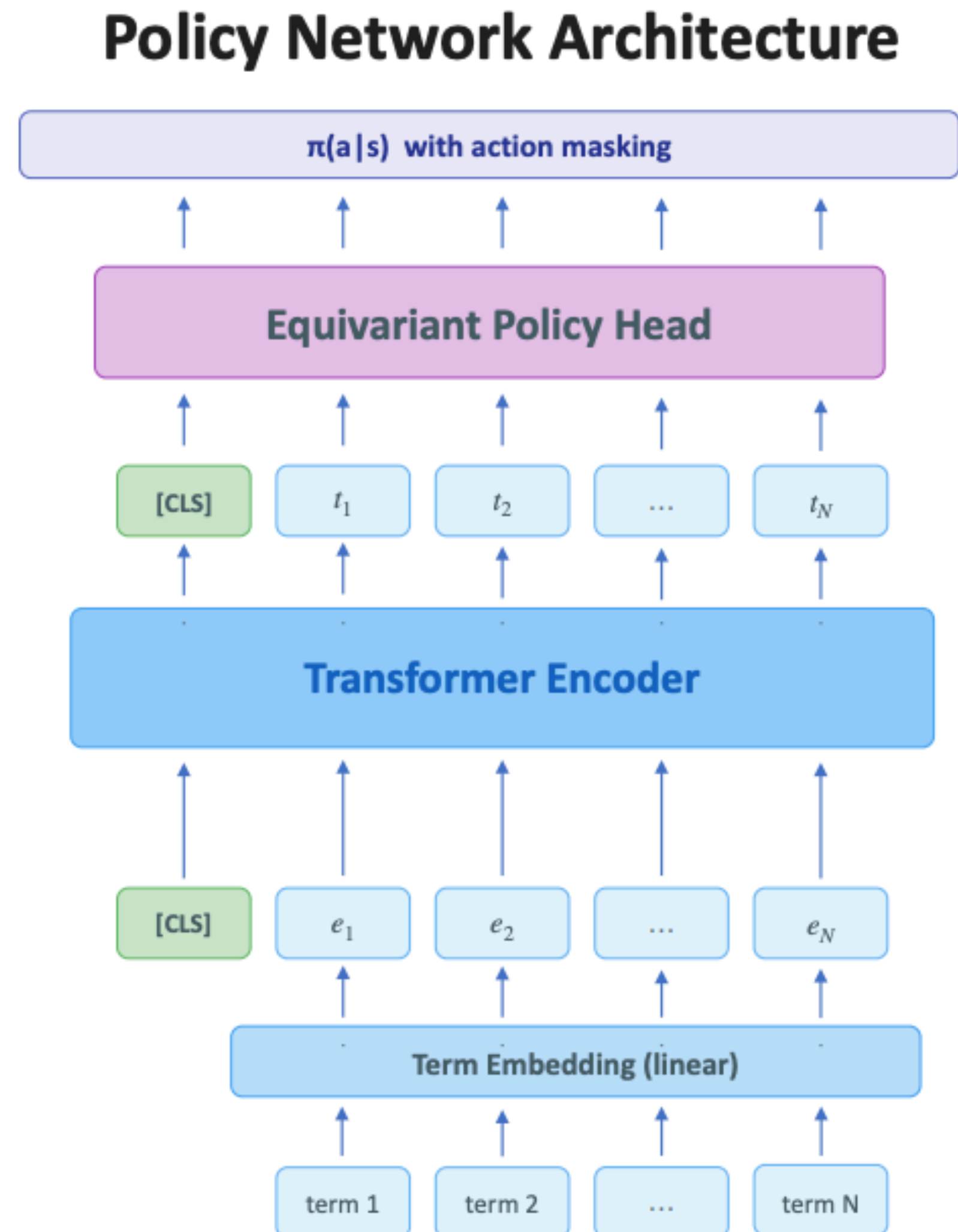
10k test expressions from
CDS enabling head-to-head
comparison.

Additional test data

100 actual 5pt YM
amplitudes
(all simplifying down
to PT form)

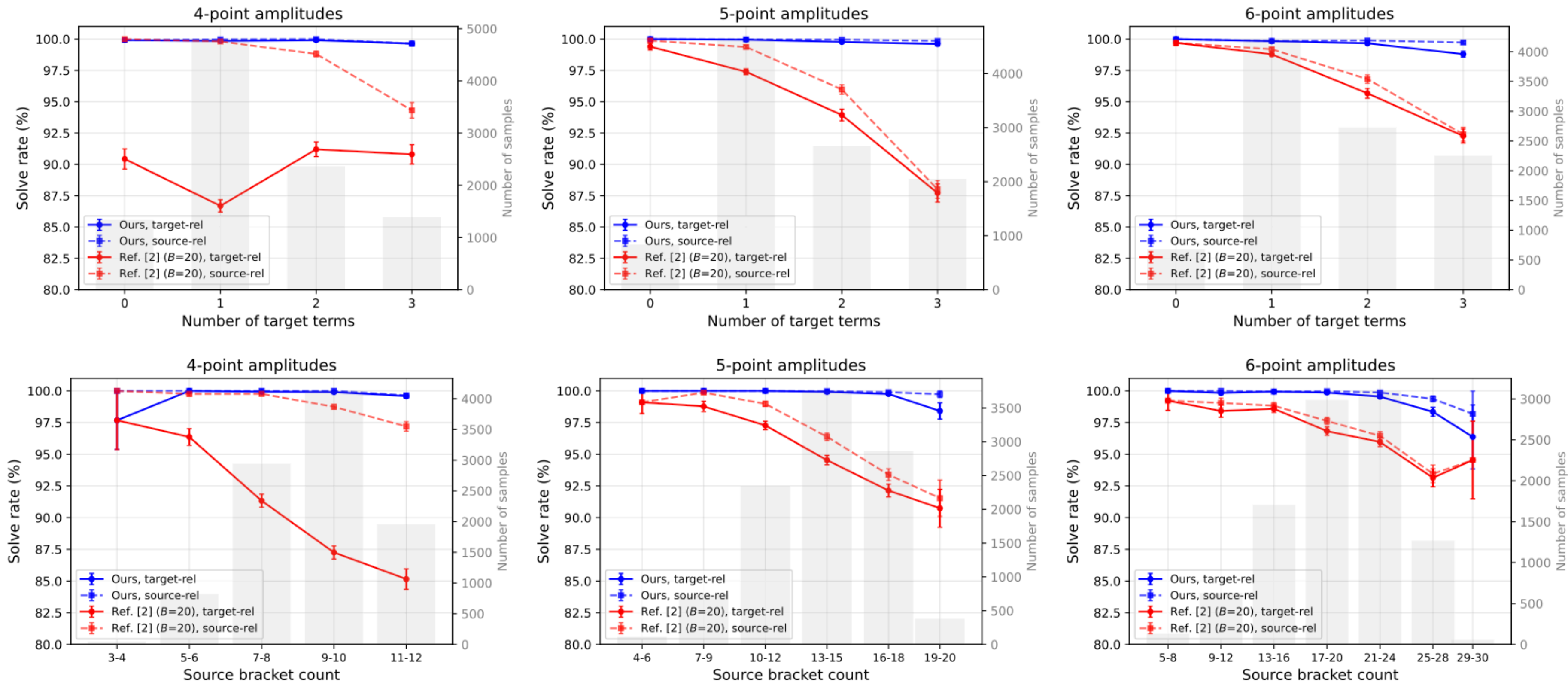
Scattering amplitudes: setup

- **Permutation-equivariant action classifier:** interchanging order of terms should change action classifier output correspondingly
- **Action space** [$\sim 1.4k$ for 4 pts, $\sim 4k$ for 5 pts and $\sim 30k$ for 6 pts]
 - choice of:
 - term (up to 10 for 4-5pts, 30 for 6 pts)
 - bracket type
 - momentum
- **Term encoding** — for each factor (up to 8):
 - coefficient
 - bracket type
 - momenta
 - power



Scattering amplitudes: results

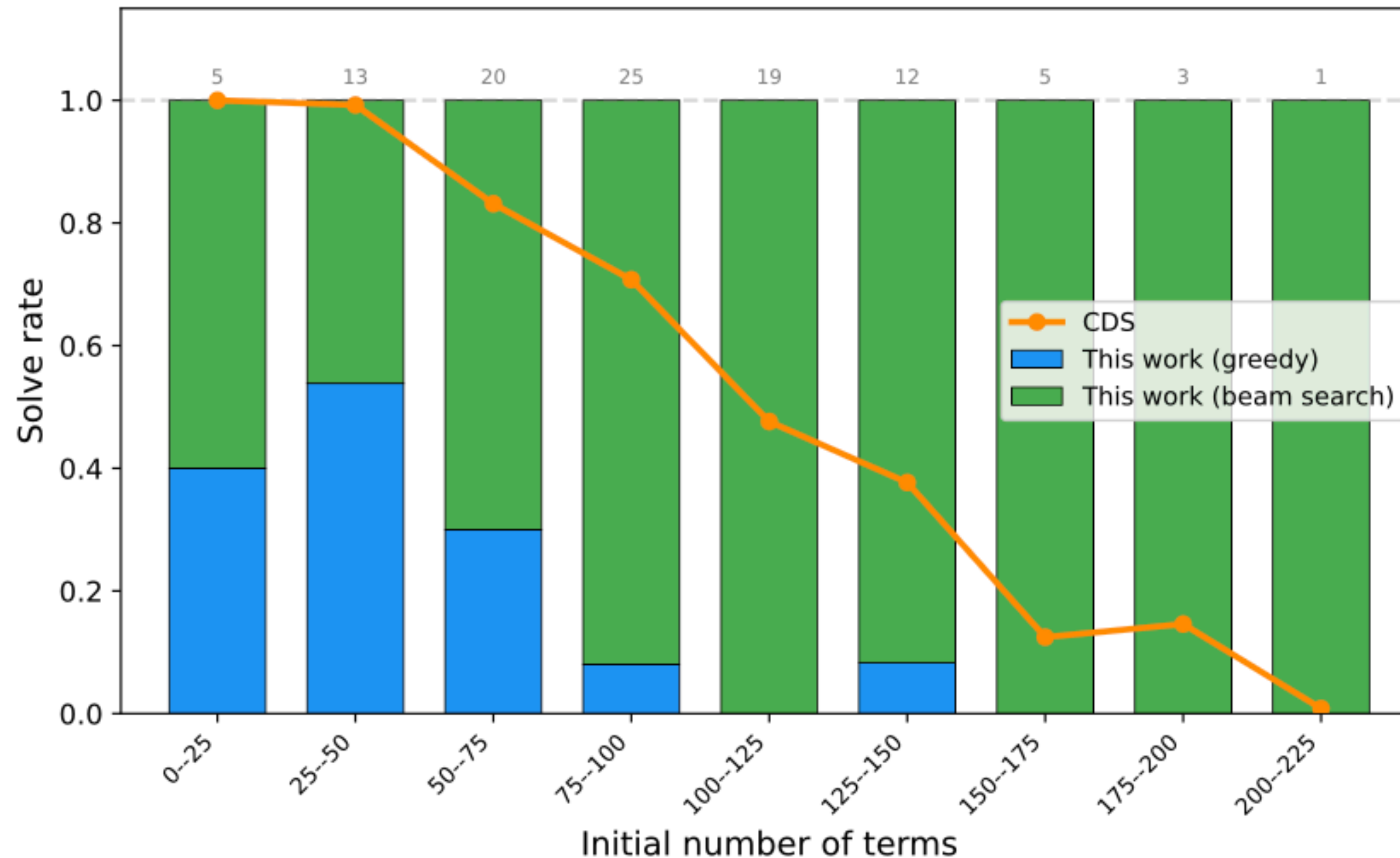
Target relative: simplifies completely to target
Source relative: simplifies source any amount



- We achieve nearly perfect simplification rate across all target and source complexities
- Big improvement vs the end-to-end transformer of CDS

Scattering amplitudes: results

- We also achieved 100% simplification of actual 5pt YM amplitudes (up to ~200 initial terms), another big improvement over CDS!



Application 3: IBP reduction

[DS 2604.05034]

- At the precision frontier, calculation of SM processes with Feynman diagrams can generate hundreds of loop integrals. These need to be reduced before numerically evaluating.
- Integration-by-parts identities are linear relations among loop integrals, parametrized by *seed* \mathbf{a} and *operation* (choice of loop momentum k_i and additional momentum v^μ):

$$I[\mathbf{a}] = \int \prod_{j=1}^L d^d k_j \prod_{i=1}^{N_p+N_s} \frac{1}{D_i^{a_i}} \quad \int \prod_{j=1}^L d^d k_j \frac{\partial}{\partial k_i^\mu} (v^\mu \cdot \text{integrand}) = 0,$$

- Can use IBP identities to reduce any loop integral to a basis of master integrals

Laporta algorithm

r: total denominator
(propagator) weight

$$w(I[\mathbf{a}]) = \left(\sum_i \max(a_i, 0), \sum_i |\min(a_i, 0)|, \mathbf{w}_3 \right)$$

s: total numerator
weight

- Basis for general purpose IBP reduction programs like **Kira**

- Make a giant linear system of IBPs
- Use Gaussian elimination to target highest weight IBP
- Continue until fully reduced to master integrals

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1k} \\ a_{21} & a_{22} & \cdots & a_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kk} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix},$$

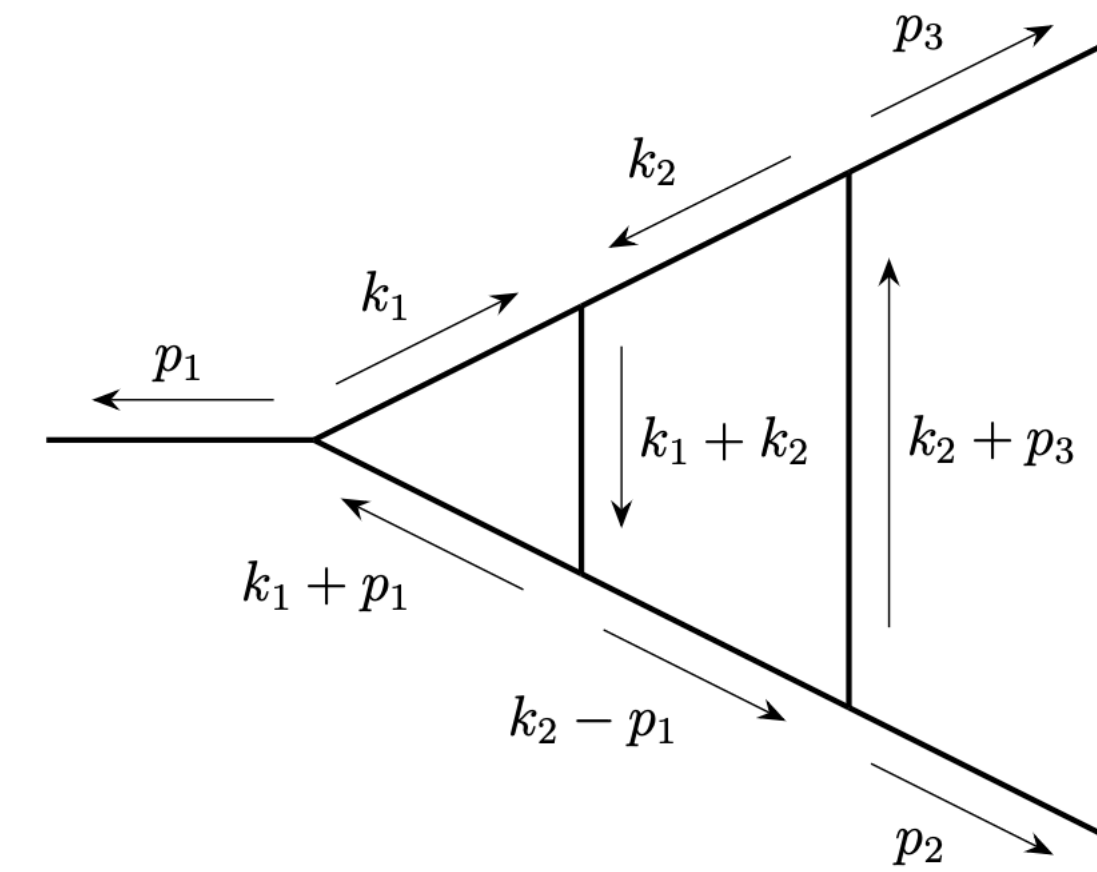
$$\left[\begin{array}{cccc|c} a'_{11} & a'_{12} & \cdots & a'_{1k} & b'_1 \\ 0 & a'_{22} & \cdots & a'_{2k} & b'_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & a'_{kk} & b'_k \end{array} \right].$$

- Guaranteed to work: number of IBP identities grows faster than the number of seeds
- But very computationally inefficient — **memory wall**

Self-supervised AI for Loop Integral Reduction (SAILIR)

- Learning to Unscramble applied to IBP reduction
- Some challenges unique to IBP reduction:
 - weight reduction instead of term reduction
 - much bigger action space (number of IBPs combinatorially large)
 - but set of valid actions much smaller (10-100 per expression) but variable (depends on expression)
 - need special classifier architecture to handle variable action space

Triangle-box topology



from von Hippel & Wilhelm (2025)

- 16 master integrals

$I[0, 0, 1, 1, 1, 0, 0],$	$I[0, 1, 1, 1, 0, 0, 0],$
$I[0, 1, 1, 1, 1, 0, 0],$	$I[-1, 1, 1, 1, 1, 0, 0],$
$I[1, 0, 0, 1, 1, 1, 0],$	$I[1, 0, 1, 0, 0, 1, 0],$
$I[1, 0, 1, 0, 1, 0, 0],$	$I[1, 0, 1, 0, 1, 1, 0],$
$I[1, -1, 1, 0, 1, 1, 0],$	$I[1, 0, 1, 1, 1, 0, 0],$
$I[1, -1, 1, 1, 1, 0, 0],$	$I[1, 0, 1, 1, 1, 1, 0],$
$I[1, 1, 0, 1, 0, 1, 0],$	$I[1, 1, 0, 1, 1, 0, 0],$
$I[1, 1, 0, 1, 1, 1, 0],$	$I[1, 1, 1, 1, 1, 0, 0].$

$$I[a_0, a_1, a_2, a_3, a_4, a_5, a_6] = \int \frac{d^d k_1 d^d k_2}{D_1^{a_0} D_2^{a_1} \dots D_7^{a_6}}$$

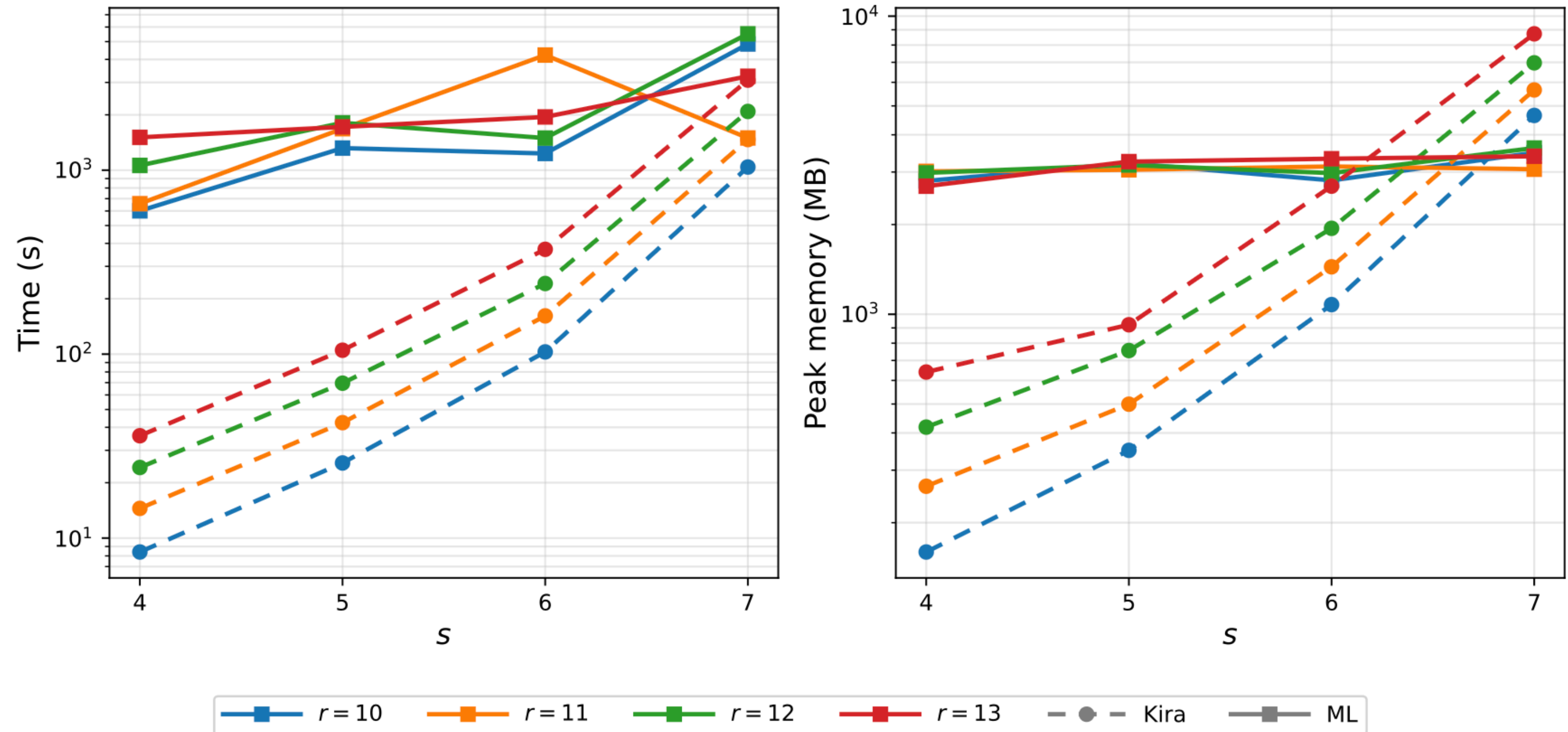
$D_1 = k_1^2,$	$D_2 = k_2^2,$
$D_3 = (k_1 + k_2)^2,$	$D_4 = (k_1 + p_1)^2$
$D_5 = (k_2 + p_3)^2,$	$D_6 = (k_2 - p_1)^2$
$D_7 = (k_1 + p_3)^2$ (ISP),	

- **Training data:** 1M training samples from randomly applying IBP identities to corner integrals

- 3 massive external momenta, 2 massless loop momenta

IBP reduction: results

r	s	Integral	r	s	Integral
10	4	$I[2, 1, 2, 1, 2, 2, -4]$	10	6	$I[1, 1, 3, 2, 2, 1, -6]$
11	4	$I[2, 2, 2, 1, 1, 3, -4]$	11	6	$I[1, 4, 2, 1, 2, 1, -6]$
12	4	$I[2, 3, 1, 3, 1, 2, -4]$	12	6	$I[3, 2, 3, 2, 1, 1, -6]$
13	4	$I[2, 3, 3, 3, 1, 1, -4]$	13	6	$I[3, 2, 1, 3, 2, 2, -6]$
10	5	$I[1, 1, 2, 2, 1, 3, -5]$	10	7	$I[2, 3, 1, 1, 2, 1, -7]$
11	5	$I[1, 1, 2, 3, 2, 2, -5]$	11	7	$I[2, 1, 1, 2, 3, 2, -7]$
12	5	$I[1, 2, 2, 2, 1, 4, -5]$	12	7	$I[3, 1, 1, 1, 1, 5, -7]$
13	5	$I[2, 2, 3, 3, 2, 1, -5]$	13	7	$I[2, 2, 3, 3, 1, 2, -7]$



- Tested on 16 triangle-box integrals with increasing weight

- SAILIR is memory stable while **Kira** memory grows rapidly with weight!

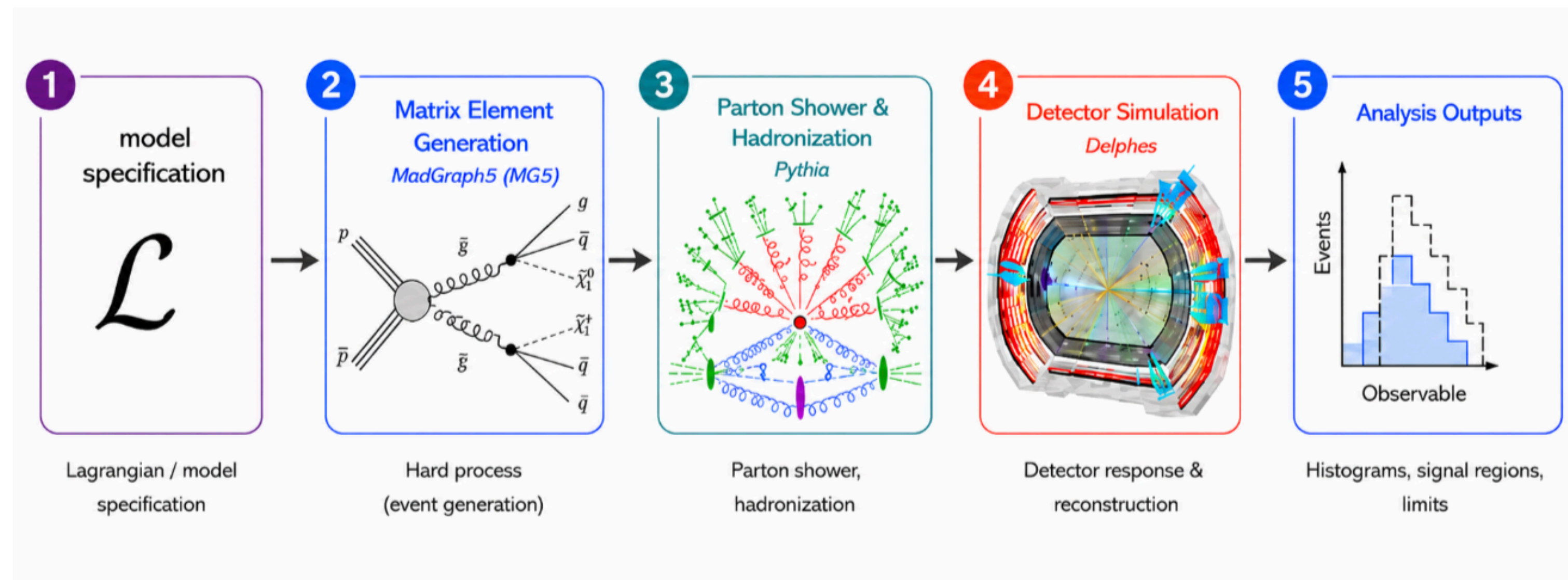
Lessons learned

- These projects were both intensely exhilarating and extremely frustrating.
- In both, Claude Code played the role of a graduate student that I interacted with entirely through text chat. Under my supervision, it wrote and ran all of the code, did all of the analysis, made all of the plots, and iterated with me on drafts of the paper.
- It was insanely fast but made tons of mistakes and kept forgetting things. Without constant close supervision and thorough validation and detailed cross checking, these projects would not have succeeded.
- Despite their flaws, I believe these agentic AI assistants will significantly lower the barrier to exploring and implementing new ideas. **Huge implications for our field!**
 - What big new questions can we now tackle that we couldn't before?
 - How do we ensure quality control? How do we handle the coming flood of papers?
 - How do we protect the future of our graduate students and postdocs? Urgently need to train them to use this amazing but flawed tool!

Bonus material: ColliderBench

[Faroughy, Palacios Schweitzer, Pang, Mishra Sharma, DS 2605.13950]

- A new benchmark for evaluating LLM agents on automated, long-horizon, real-world scientific tasks
- Setting: reproducing published LHC analyses (new physics searches) using public simulation tools — “recasting”



Challenges

- Unique challenges:
 - experimental papers from CMS and ATLAS often missing crucial details, inaccuracies, errors.
 - Gap between public simulation tools and proprietary CMS/ATLAS tools, correction factors, etc
 - Agent cannot get a perfect reproduction. Needs combination of physical intuition, domain knowledge, guesswork

ColliderBench: inputs

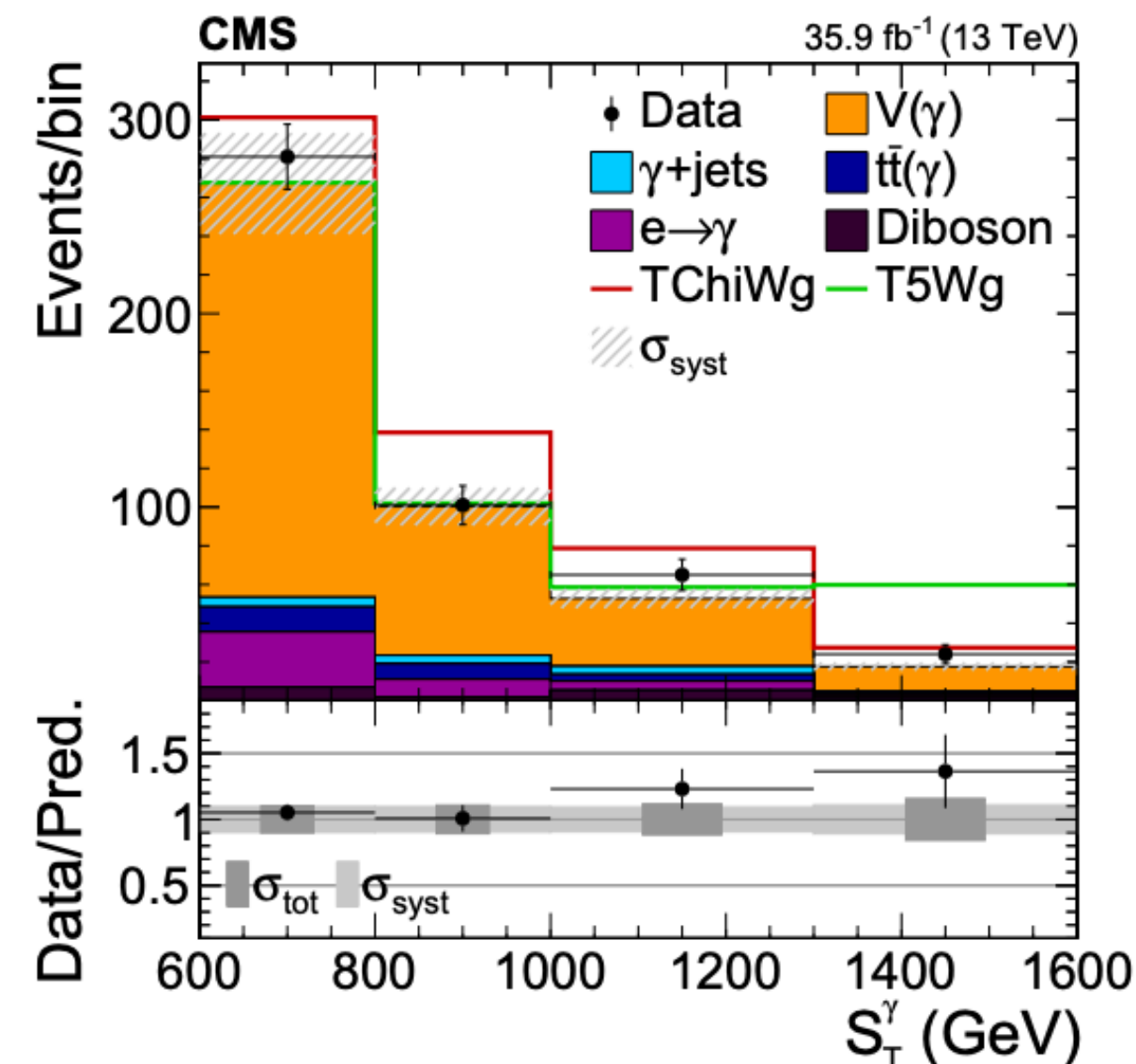
Search for gauge-mediated supersymmetry in events with at least one photon and missing transverse momentum in pp collisions at $\sqrt{s} = 13$ TeV

- Paper (PDF) from CMS/ATLAS

The CMS Collaboration*

Table 1: Summary of the event selection criteria required for the control, validation, and signal regions.

Region	Selection
Preselection	p_T^{miss} filters
	At least one reconstructed vertex
	At least one photon with $p_T > 180$ GeV
	$\Delta R(\gamma, \text{jet}) > 0.5$ $\Delta\phi(\vec{p}_T^{\text{miss}}, \text{jet}) > 0.3$ rad, if $p_T(\text{jet}) > 100$ GeV
Control region	Preselection
	$p_T^{\text{miss}} > 100$ GeV
	$M_T(\gamma, \vec{p}_T^{\text{miss}}) > 100$ GeV
	$p_T^{\text{miss}} < 300$ GeV or $M_T(\gamma, \vec{p}_T^{\text{miss}}) < 300$ GeV
Validation region	Preselection
	$p_T^{\text{miss}} > 300$ GeV
	$M_T(\gamma, \vec{p}_T^{\text{miss}}) > 300$ GeV
	$S_T^\gamma < 600$ GeV
Signal region	Preselection
	$p_T^{\text{miss}} > 300$ GeV
	$M_T(\gamma, \vec{p}_T^{\text{miss}}) > 300$ GeV
	$S_T^\gamma > 600$ GeV



Abstract

A search for gauge-mediated supersymmetry (SUSY) in final states with photons and large missing transverse momentum is presented. The data sample of pp collisions at $\sqrt{s} = 13$ TeV was collected with the CMS detector at the CERN LHC and corresponds to an integrated luminosity of 35.9 fb^{-1} . Data are compared with models in which the lightest neutralino has bino- or wino-like components, resulting in decays to photons and gravitinos, where the gravitinos escape detection. The event selection is optimized for both electroweak (EWK) and strong production SUSY scenarios. The observed data are consistent with standard model predictions, and limits are set in the context of a general gauge mediation model in which gaugino masses up to 980 GeV are excluded at 95% confidence level. Gaugino masses below 780 and 950 GeV are excluded in two simplified models with EWK production of mass-degenerate charginos and neutralinos. Stringent limits are set on simplified models based on gluino and squark pair production, excluding gluino (squark) masses up to 2100 (1750) GeV depending on the assumptions made for the decay modes and intermediate particle masses. This analysis sets the highest mass limits to date in the studied EWK models, and in the considered strong production models when the mass difference between the gauginos and the squarks or gluinos is small.

ColliderBench: inputs

- Task specification prompt
 - new physics signal to simulate
 - histogram from paper to reproduce
 - other output requirements

TASK.md

Paper: CMS-SUS-16-046

Centre-of-mass energy: 13 TeV

Luminosity: 35.9 fb^{-1}

Task type: simulation (shape-only)

Signal benchmark: TChiWg_700

Observable: S_T^γ

Task

Implement the search analysis described in **CMS-SUS-16-046** and use it to predict the normalized event distribution (shape) of S_T^γ for the signal benchmark point TChiWg_700, in the analysis's signal region.

The agent should:

1. Generate TChiWg_700 events using a matrix-element generator, parton shower, and detector simulation chain of its choice.
2. Read the paper to determine the object identification, event-selection requirements, and signal-region cuts that define the analysis, then apply them to the generated events.
3. Histogram the surviving events in S_T^γ using the bin edges already present in the `results/*.yaml` template. The bin edges must not be modified.

Definitions

- S_T^γ is the scalar sum of p_T^{miss} and the transverse momenta of all photons in the event.
- TChiWg_700 denotes electroweak associated production of mass-degenerate winos at $m(\tilde{W}) = 700 \text{ GeV}$, decaying to a massless gravitino LSP and Standard Model gauge bosons, with at least one photon in the final state through the TChiWg simplified-model topology.

Output requirements

Artifact	Purpose
<code>results/*.yaml</code>	Fill null bin values with predicted relative signal bin contents.
<code>analysis/*.py</code>	Event-selection code, runnable on the generated sample.
<code>data/*.root, sims/*.dat</code>	Selected-event files and generator/detector cards.
<code>report.md</code>	Methodological choices and deviations from the paper.

Important

The goal is to predict the signal shape from the agent's own simulation and analysis pipeline. The agent must not extract or digitize the target bin values from the paper's figures, tables, or HEPData record.

ColliderBench: inputs

- HEP toolbox
 - md files that teach agents how to use HEP software tools

```
1  √ # `feynrules` – browse and download UFO models from the FeynRules database
2
3  **Purpose.** If the bundled MadGraph UFO models (`sm`, `MSSM_SLHA2`,
4  `SMEFTsim`) don't cover your paper's BSM scenario, fetch an additional
5  model from the FeynRules wiki.
6
7  **When to use.** Only when `bin/simulate info` doesn't already list a
8  suitable UFO. For standard SUSY / SMEFT / top-EFT the bundled models are
9  usually enough.
10
11 √ ## Invocation
12
13 ```bash
14 bin/feynrules categories                # 7 top-level categories
15 bin/feynrules list --category SusyModels # browse one category
16 bin/feynrules list --search "vector-like quark" # substring search over slug/title/description
17 bin/feynrules info <model-slug>        # show all attachments for a model
18 bin/feynrules fetch <model-slug> [--file <name>] [--all] [--dest <dir>] [--extract]
19 bin/feynrules refresh-catalog           # re-scrape the wiki (catalog is cached)
20 ```
21
```

```
1  √ # Simulation stack – MadGraph5 / Pythia8 / Delphes
2
3  **Purpose.** Generate signal events end-to-end (parton-level → parton
4  shower → detector simulation) when no matching MC is available in CMS
5  Open Data. The CLIs are simple enough that there's no wrapper; you call
6  them directly. `bin/simulate info` lists installed models and cards;
7  `bin/simulate --doc` shows this page.
8
9  **When to use.** For BSM signals the paper defines but nobody has
10 produced for Open Data (SUSY benchmark points, EFT signals, etc.). For
11 SM backgrounds you usually want `cms-opendata` instead – save compute.
12
13 ---
14
15 √ ## Quick discovery
16
17 ```bash
18 bin/simulate info          # installed UFO models + Delphes cards + tool paths
19 bin/simulate --doc        # this page
20 ```
21
22 The relevant env vars (set inside the bench image; fall back to the
23 in-repo paths otherwise):
24
```

ColliderBench: metrics

$$d(\hat{y}, y^*) = \sqrt{\frac{\sum_{k=1}^K (\hat{y}_k - y_k^*)^2}{\sum_{k=1}^K y_k^{*2}}}.$$

L2 histogram distance

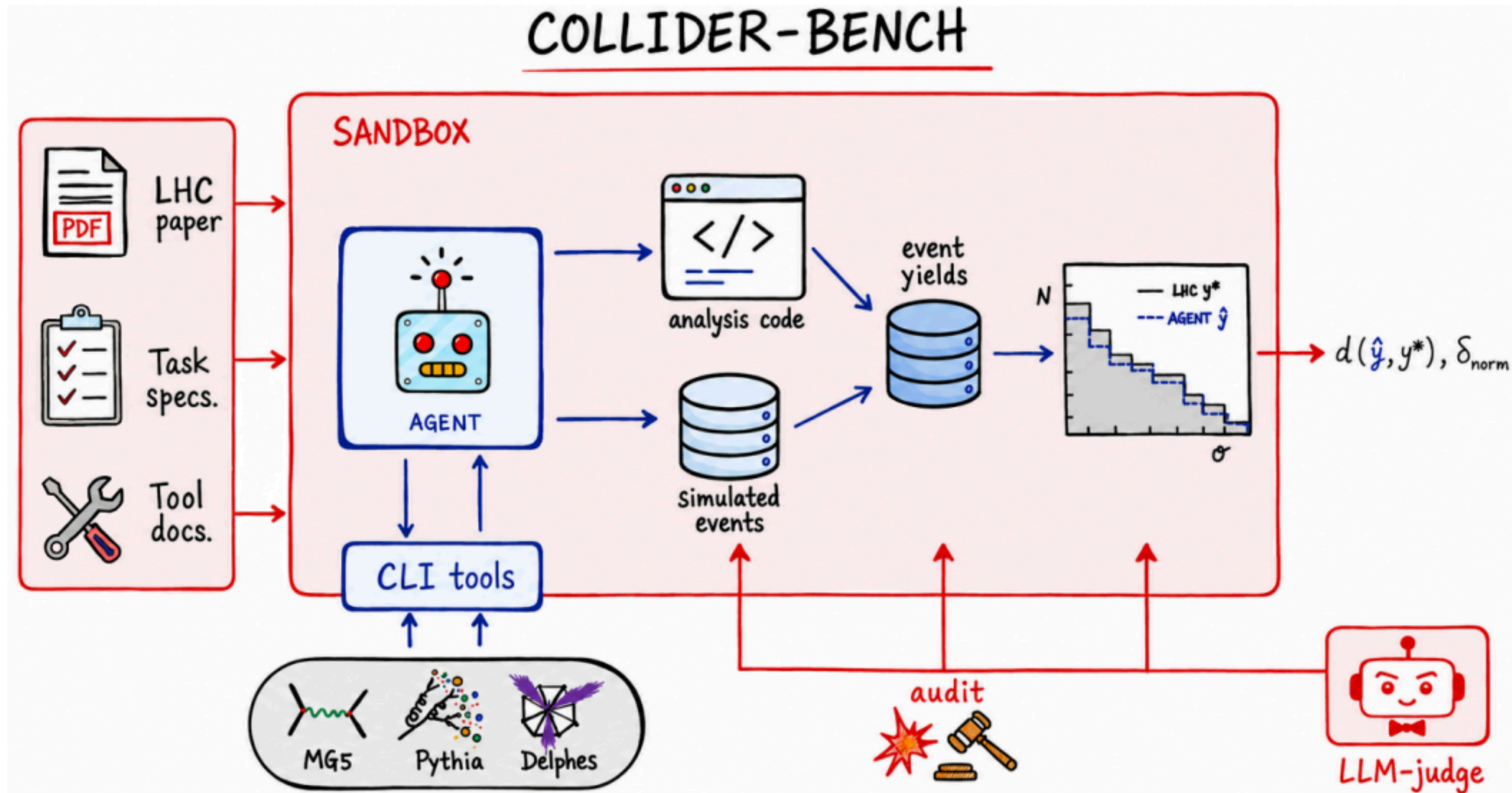
$$\delta_{\text{norm}} = \frac{|\hat{Y} - Y^*|}{Y^*}.$$

normalization error

$$\text{Acc}_\tau = \mathbb{I}[d_{\text{task}} < \tau],$$

pass/fail metric

ColliderBench: overview



LLM judge to catch cheating, hallucinations

ColliderBench V1

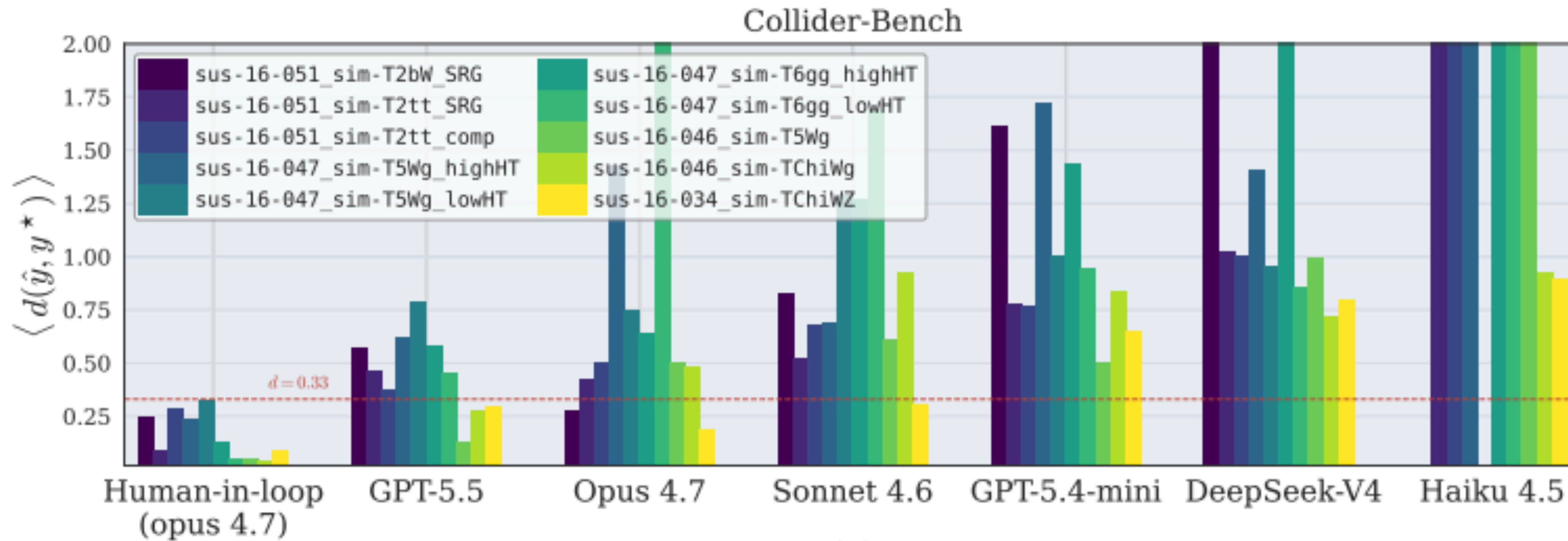
- 10 tasks over 4 CMS supersymmetry searches from ~2016 (13 TeV, 36/fb)

Table 1: Source analyses and simulation tasks in COLLIDER-BENCH. Each task fixes a paper, signal benchmark, target observable or signal region, and output template. Difficulty is graded from the relative- L^2 residual of our physicist-in-the-loop reproduction of each task: ★ (easy), ★★ (medium), ★★★ (hard).

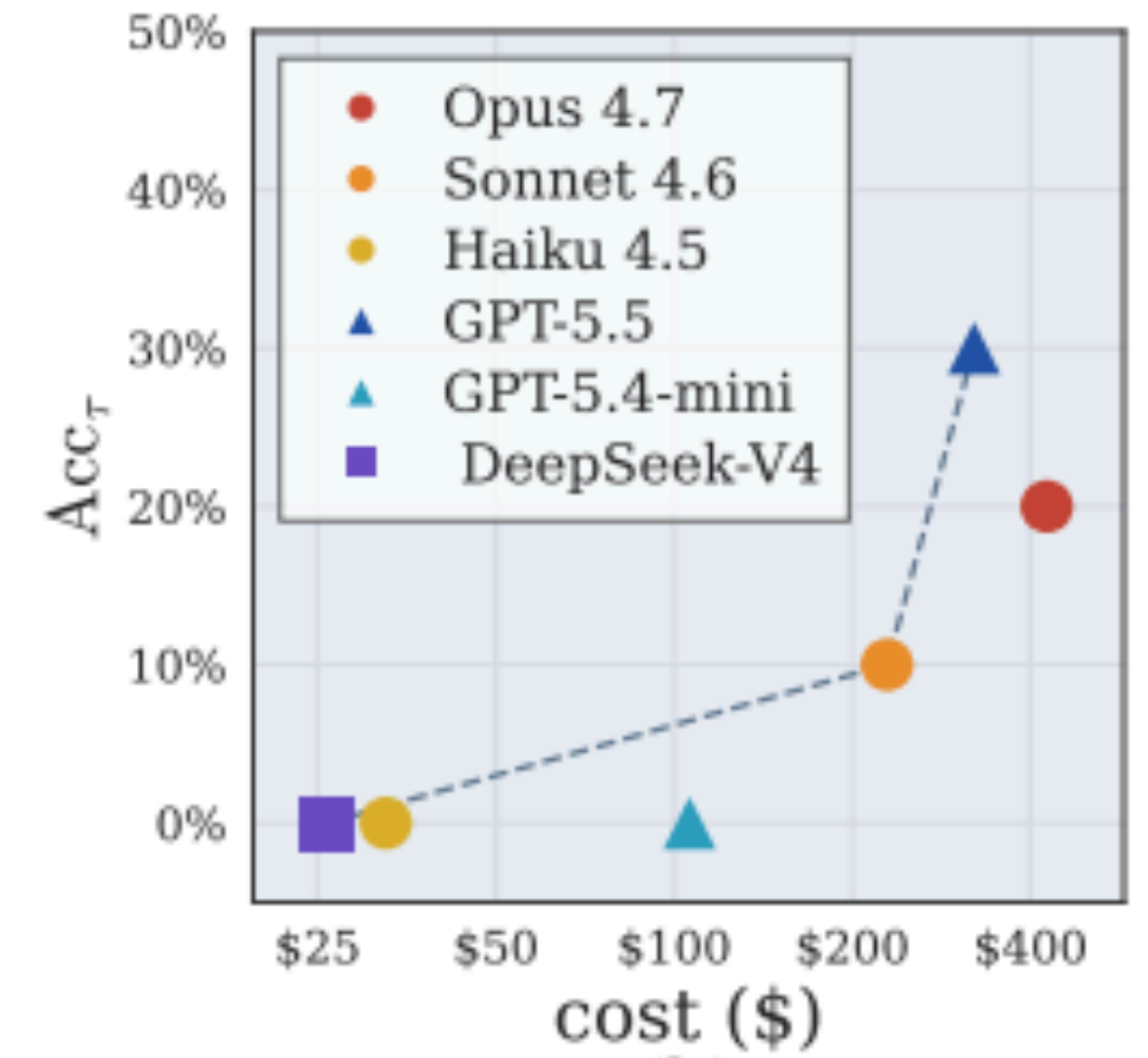
Task	Analysis target	Signal s	obs. \mathcal{O}	LHC search	Reference	Diff.
sus-16-034_sim-TChiWZ	leptons + jets	TChiWZ	E_T^{miss}	CMS-SUS-16-034	Sirunyan et al. [2018b]	★
sus-16-046_sim-T5Wg	photons	T5Wg	S_T^γ	CMS-SUS-16-046	Sirunyan et al. [2018a]	★
sus-16-046_sim-TChiWg	photons	TChiWg	S_T^γ	CMS-SUS-16-046	Sirunyan et al. [2018a]	★
sus-16-047_sim-T5Wg_highHT	photons	T5Wg, high- H_T	p_T^{miss}	CMS-SUS-16-047	Sirunyan et al. [2017b]	★★
sus-16-047_sim-T5Wg_lowHT	photons	T5Wg, low- H_T	p_T^{miss}	CMS-SUS-16-047	Sirunyan et al. [2017b]	★★★
sus-16-047_sim-T6gg_highHT	photons	T6gg, high- H_T	p_T^{miss}	CMS-SUS-16-047	Sirunyan et al. [2017b]	★★
sus-16-047_sim-T6gg_lowHT	photons	T6gg, low- H_T	p_T^{miss}	CMS-SUS-16-047	Sirunyan et al. [2017b]	★
sus-16-051_sim-T2tt	single lepton	T2tt	E_T^{miss}	CMS-SUS-16-051	Sirunyan et al. [2017a]	★
sus-16-051_sim-T2tt_comp	single lepton	T2tt, compressed	E_T^{miss}	CMS-SUS-16-051	Sirunyan et al. [2017a]	★★★
sus-16-051_sim-T2bW	single lepton	T2bW	E_T^{miss}	CMS-SUS-16-051	Sirunyan et al. [2017a]	★★

- Hundreds of more tasks and searches available, will be added for V2!

ColliderBench: results

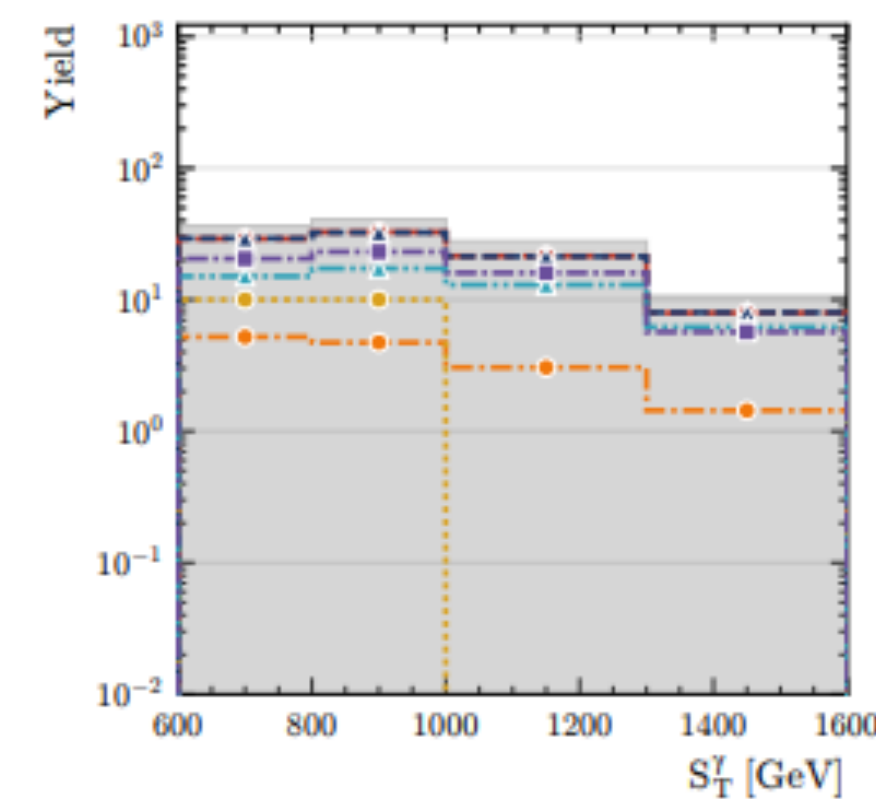


(a)

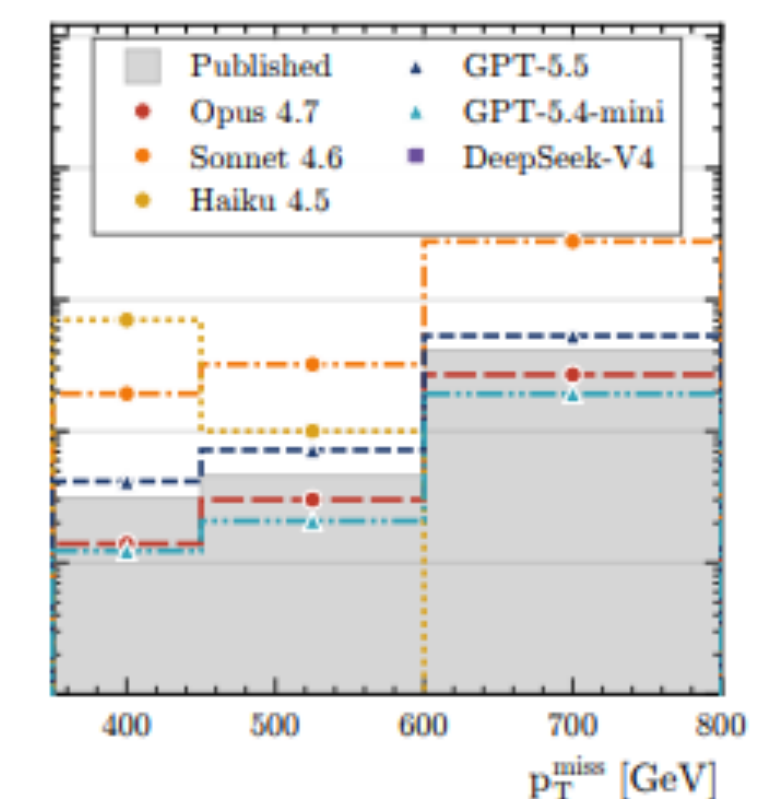


(b)

- Clear capability ladder
(Opus>Sonnet>Haiku; GPT5.5>GPT-5.4-mini)
- Clear Pareto front in pass/fail vs cost
- No LLM agent beats human-in-the-loop



(a)



(b)

Summary

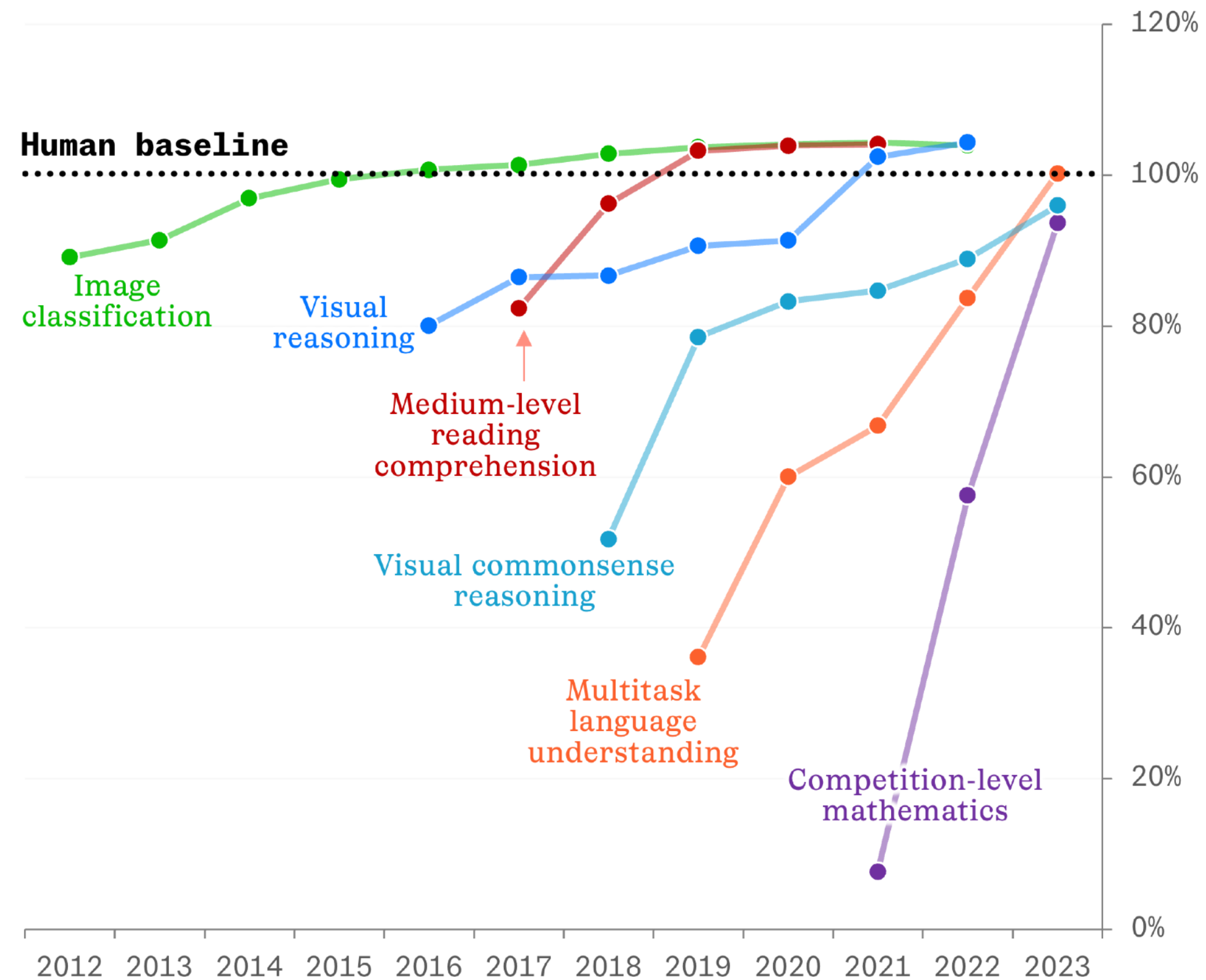
- Agentic AI is going to be a game changer for our field (and well well beyond)
- In this talk I presented two uses of agentic AI for research
 - As a hands-on research assistant -> two papers on “Learning to Unscramble”, applications to scattering amplitudes and Feynman loop integral reduction
 - As a fully automated system for reproducing LHC analyses (an important and common task in pheno research)

Where are we headed?

MAN VS. MACHINE

AI Models Are Improving Every Year

AI Technical Performance [Selected measures, 100% = human baseline]



CHARTR

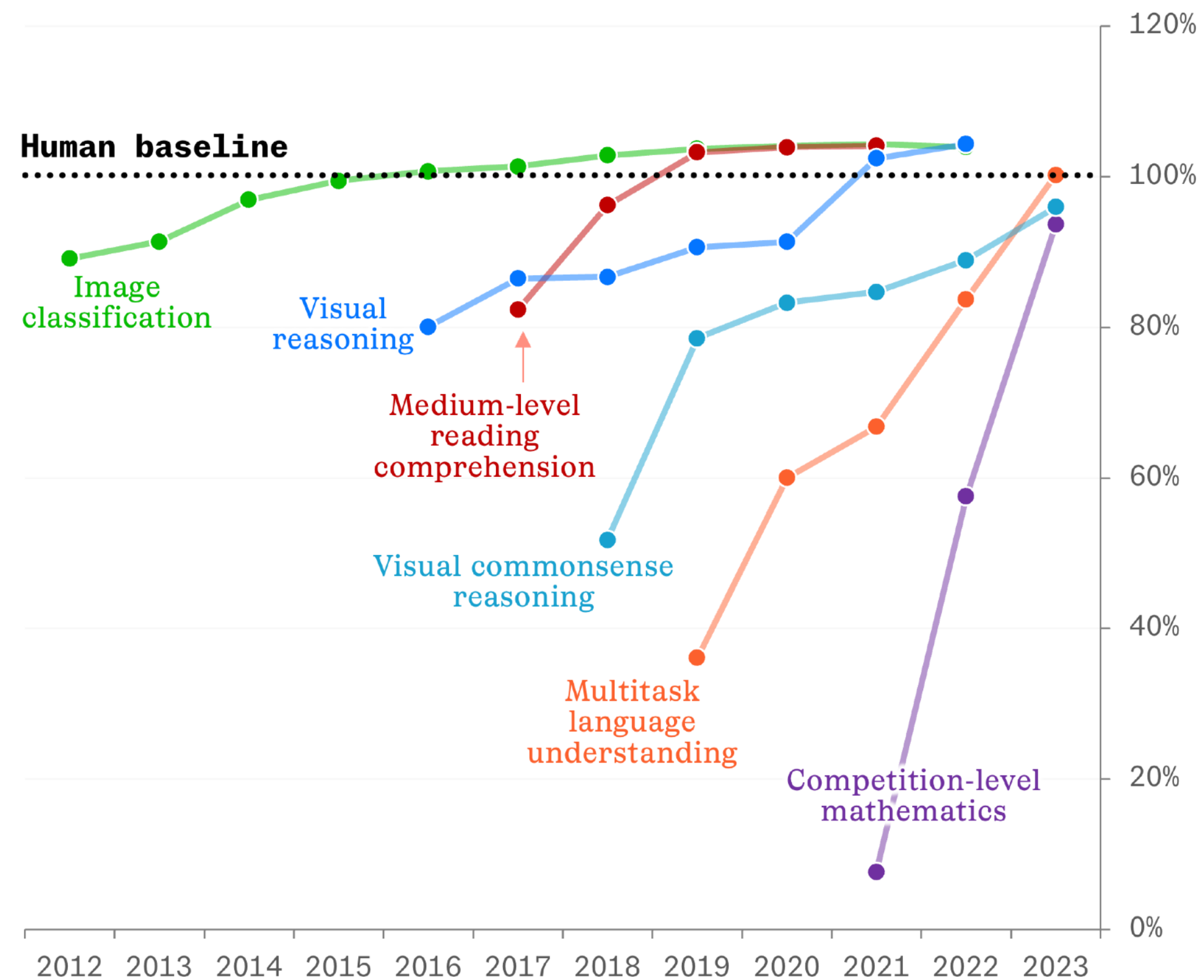
Source: Stanford University AI Index Report 2024

Where are we headed?

MAN VS. MACHINE

AI Models Are Improving Every Year

AI Technical Performance [Selected measures, 100% = human baseline]



CHARTR

Source: Stanford University AI Index Report 2024

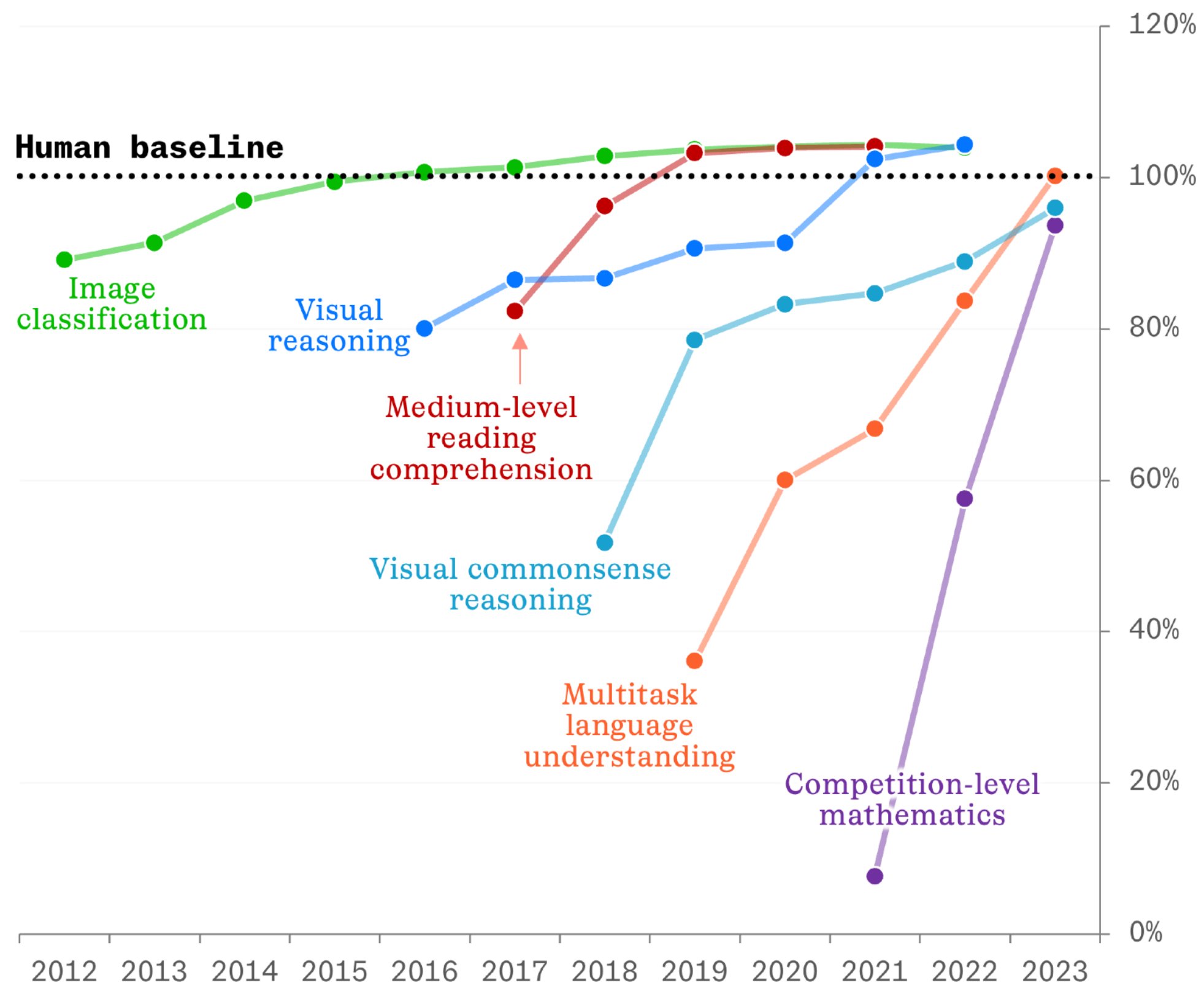


Matt Schwartz

Where are we headed?

MAN VS. MACHINE AI Models Are Improving Every Year

AI Technical Performance [Selected measures, 100% = human baseline]

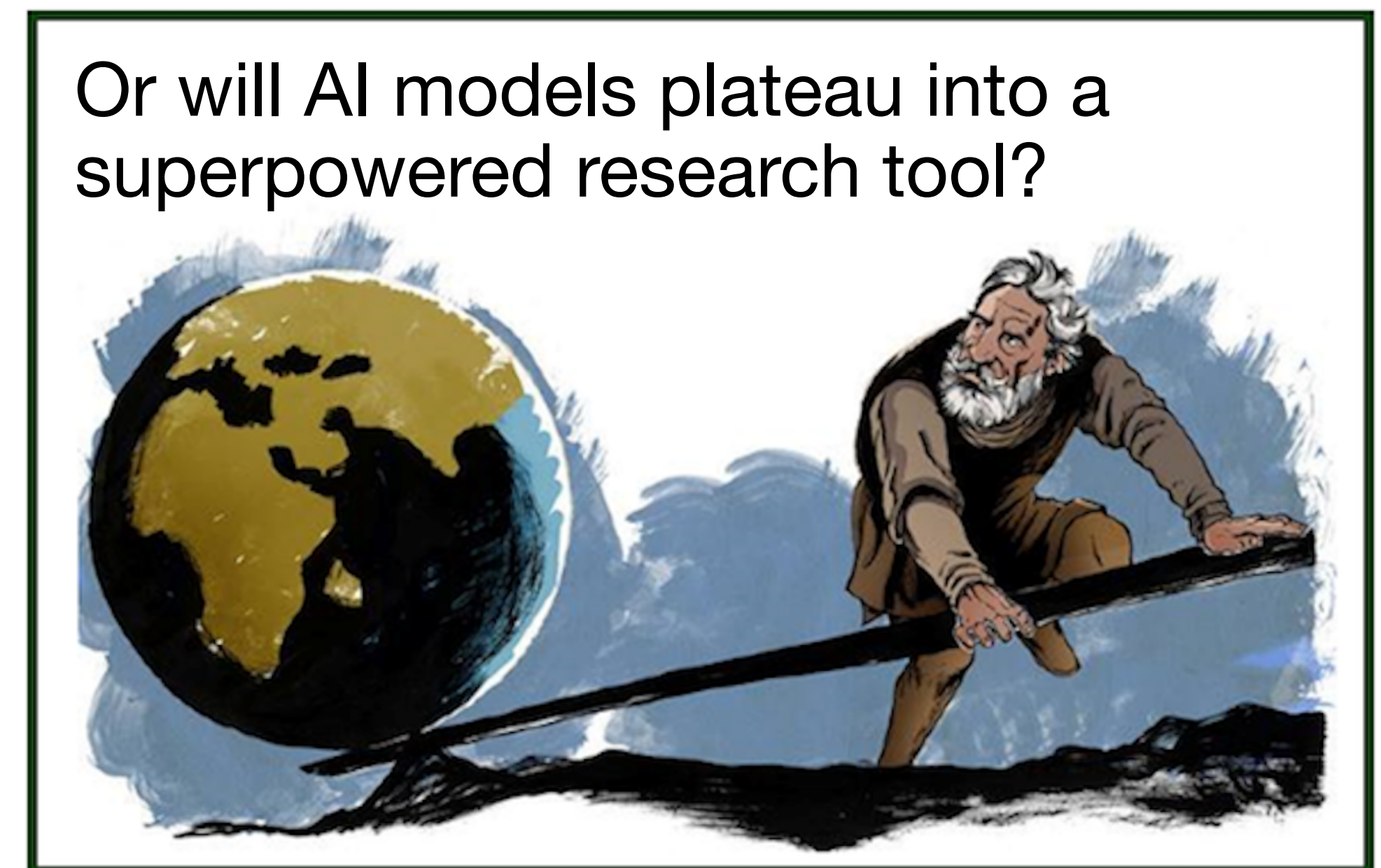


CHARTR

Source: Stanford University AI Index Report 2024



Matt Schwartz



Thanks!

Application 1: Dilog sums

[DS 2603.11164]

$$\text{Li}_2(x) = - \int_0^x dz \frac{\log(1-z)}{z}$$

$$(inversion) \quad \text{Li}_2(x) = -\text{Li}_2\left(\frac{1}{x}\right) - \frac{\pi^2}{6} - \frac{\ln^2(-x)}{2}$$

$$(reflection) \quad \text{Li}_2(x) = -\text{Li}_2(1-x) + \frac{\pi^2}{6} - \ln(x) \ln(1-x)$$

$$(duplication) \quad \text{Li}_2(x) = -\text{Li}_2(-x) + \frac{1}{2}\text{Li}_2(x^2)$$

- Toy version of Feynman loop integral simplification
- Class of expressions that Mathematica struggles to simplify

$$-4 \text{Li}_2\left(\frac{4}{x^2-8x+16}\right) + 8 \text{Li}_2\left(\frac{2}{x-4}\right) - 8 \text{Li}_2\left(\frac{2}{x-2}\right)$$

$$-6 \text{Li}_2\left(\frac{1}{x^2+2x}\right) - 6 \text{Li}_2(x^2+2x) + \frac{7}{2} \text{Li}_2\left(\frac{4}{x^2-4x+4}\right)$$

$$-7 \text{Li}_2\left(\frac{x}{2}\right) - 7 \text{Li}_2\left(\frac{2}{x-2}\right)$$

$$-8 \text{Li}_2\left(\frac{x^2}{2x-2}\right) - 8 \text{Li}_2\left(-\frac{x^2}{2x-2}\right) - \frac{7}{2} \text{Li}_2\left(\frac{4}{x^2+2x+1}\right)$$

$$+4 \text{Li}_2(x^2+4x+4) + 4 \text{Li}_2\left(\frac{x^4}{4x^2-8x+4}\right) - 7 \text{Li}_2\left(-\frac{x}{2} - \frac{1}{2}\right)$$

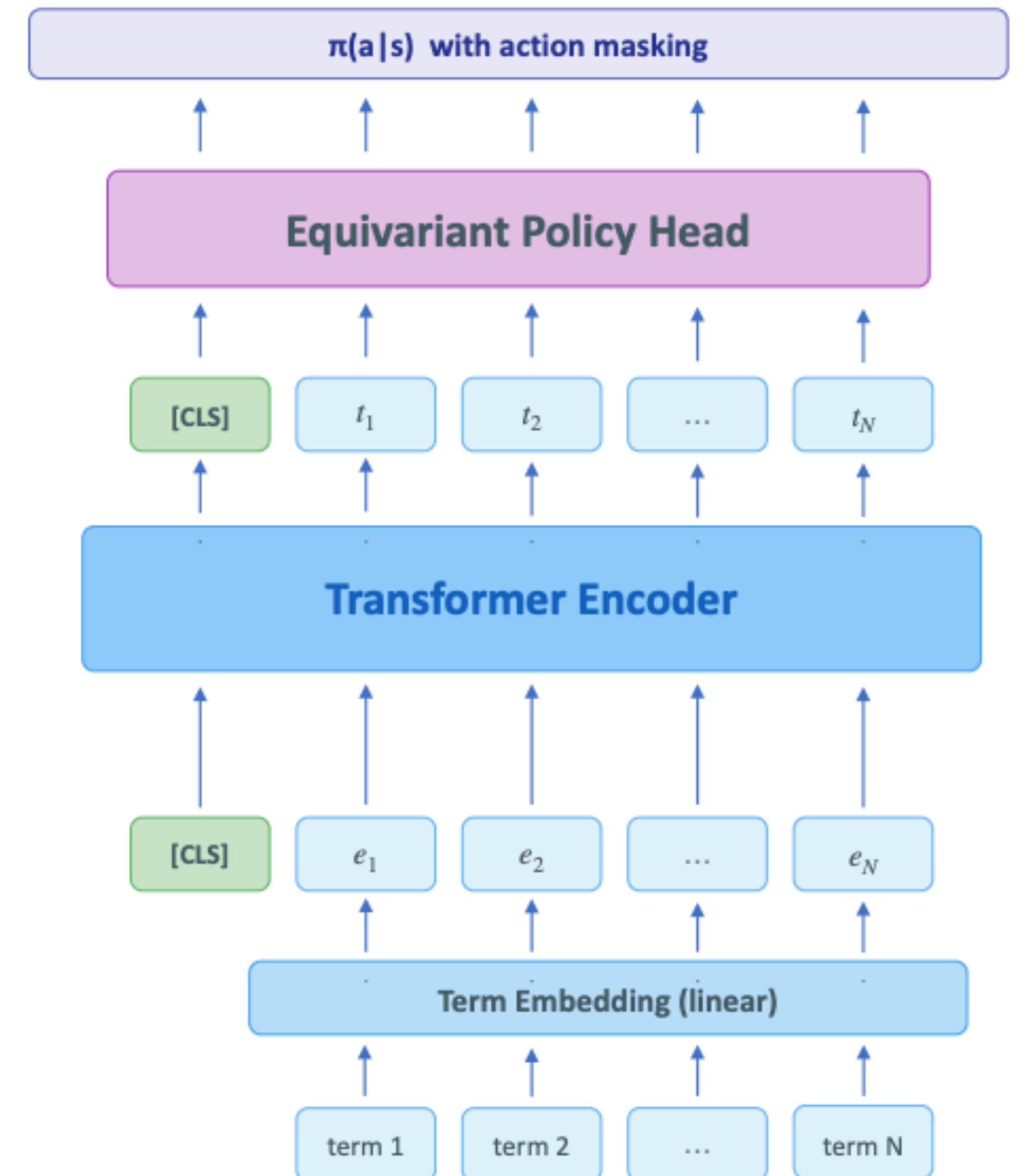
$$-2 \text{Li}_2\left(\frac{1}{x+2}\right) + 8 \text{Li}_2\left(-\frac{1}{x+2}\right) - 10 \text{Li}_2(x+2)$$

examples from Dersy,
Schwartz, Zhang

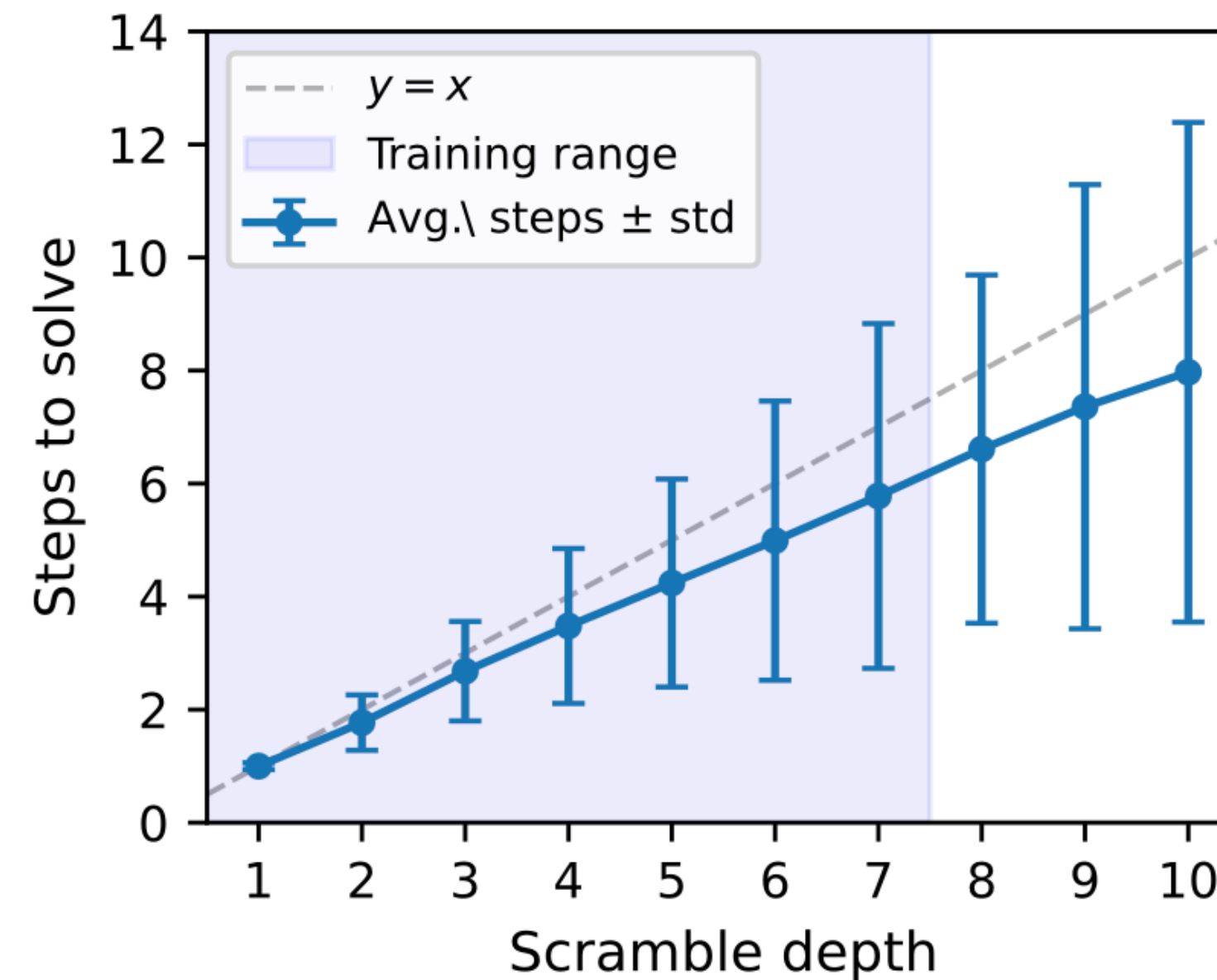
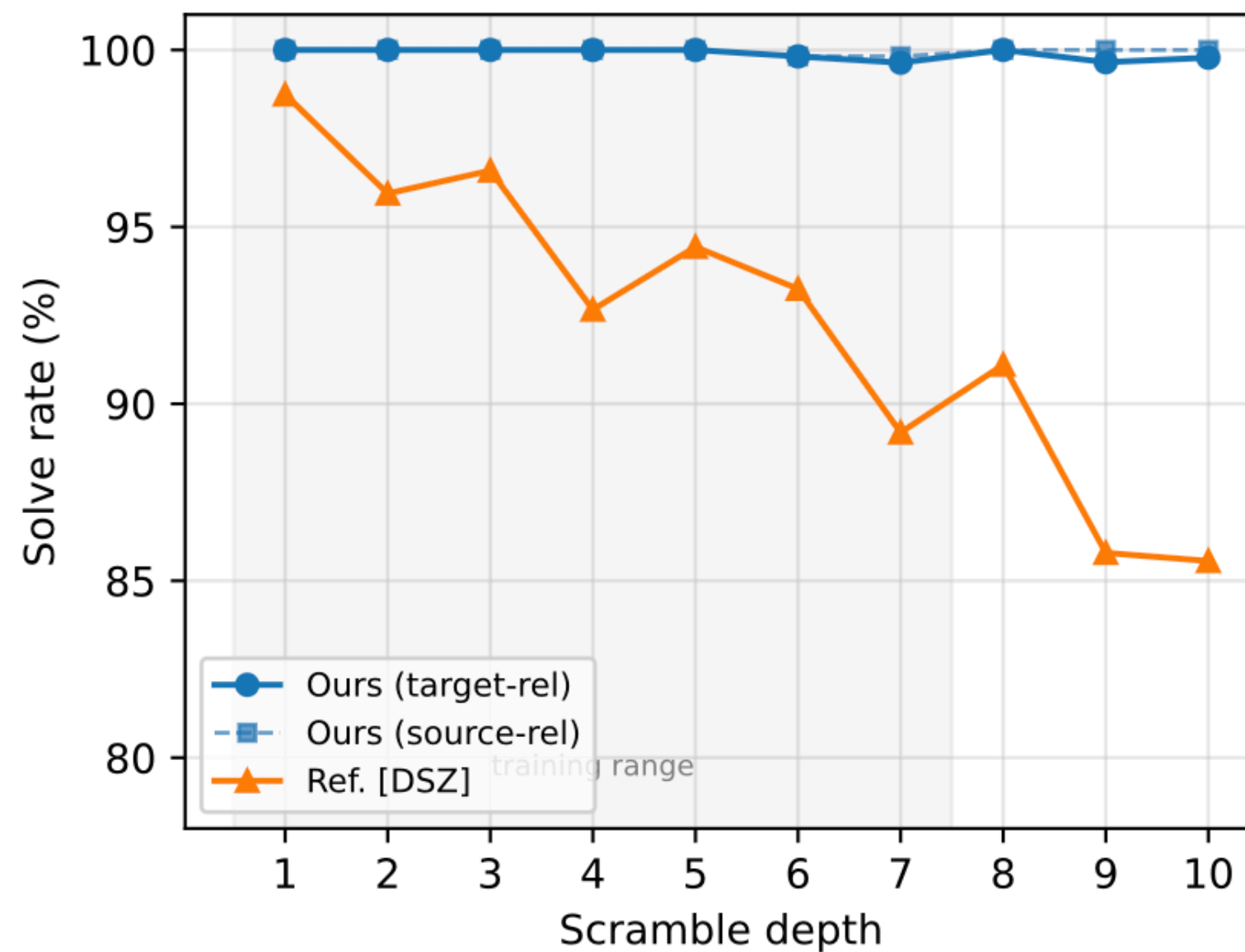
Learning to unscramble dialogs: setup

- Action space: (3 actions) x (up to 15 terms)
- Permutation-equivariant action classifier: 45 dimensional softmax (irrelevant actions masked)
- **Training data:** 100k samples scrambled from simple expressions -> 500k individual unscramble steps
- **Test data:** 5k expressions from DSZ

Policy Network Architecture



Learning to unscramble dialogs



- Nearly 100% simplification rate! Much better than end-to-end transformer (DSZ)
- Only 6 failures out of 5k test set. (Only 2 if we demand some simplification)
- Excellent performance generalizes beyond training set!
- Finds *shorter* simplification paths than training data on average

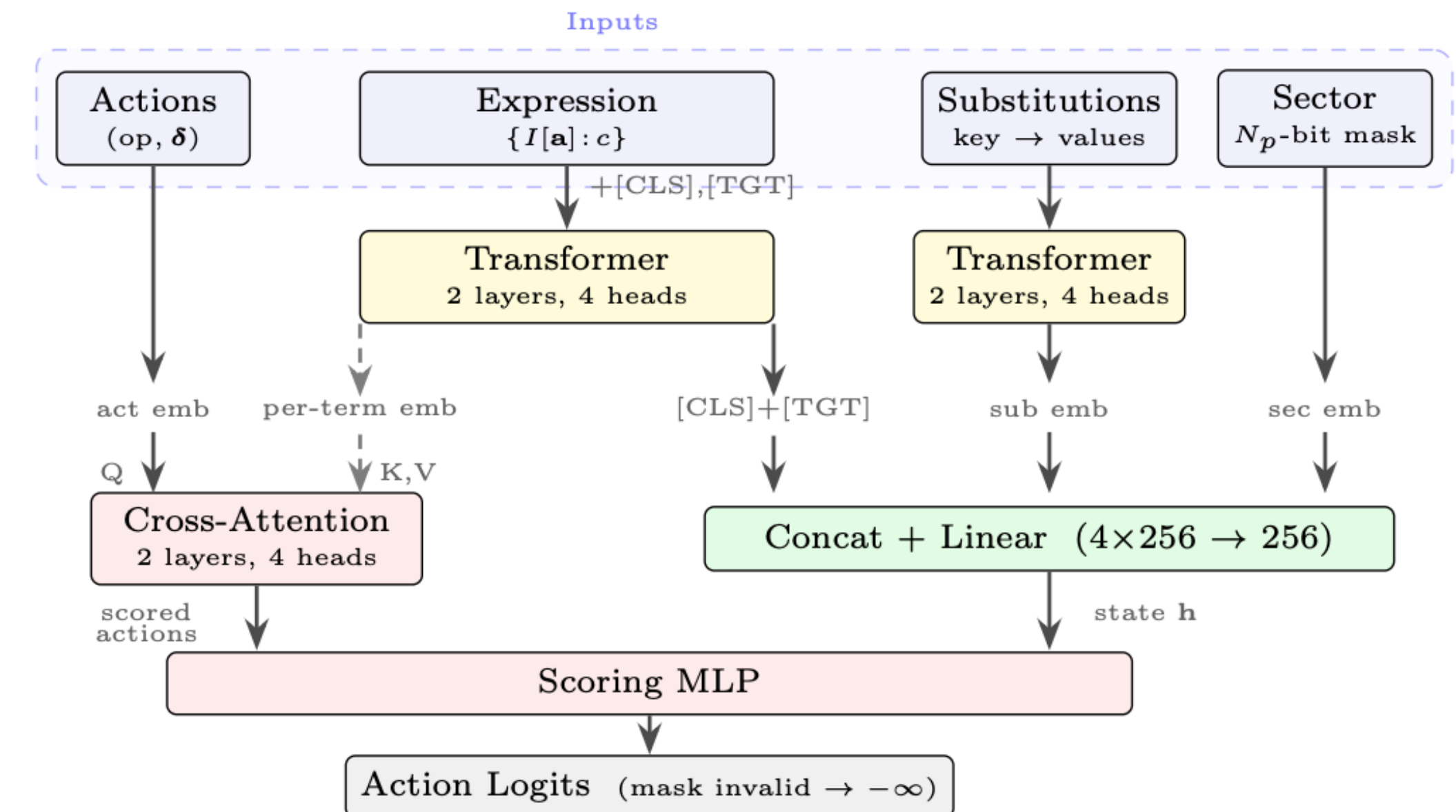
Weight vs term reduction

- Weight reordering:
 - to train the MDP to reduce by weight, we reorder the unscrambles to monotonically reduce highest weight integral in the expression at each step (can prove that this can always be done)
- Sector-wise reduction:
 - Loop integrals labeled by indices $\mathbf{a} = (a_1, a_2, \dots)$ organize into sectors — which propagator factors are present
 - Notion of subsectors $(2, 1, -1, 1) \rightarrow (1, 1, 0, 1)$
 $(-1, 0, -1, 3) \rightarrow (0, 0, 0, 1)$
 - Laporta reduction can work sectorwise: reduce to **corner integral** $I[1, 1, 0, 1]$, $I[0, 0, 0, 1]$, etc — modulo subsectors

Variable action space: learning to rank

- Used in information retrieval (ML for search engines)

- Expression \leftrightarrow question
- actions \leftrightarrow potential answers that we want to rank



- Use cross attention between encoded actions and expression to produce a per-action score
- Softmax gives classifier over all valid actions