

# Invisible Cities

Imagining the next era of AI-enabled fundamental physics research

Mariel Pettee

February 10th, 2026



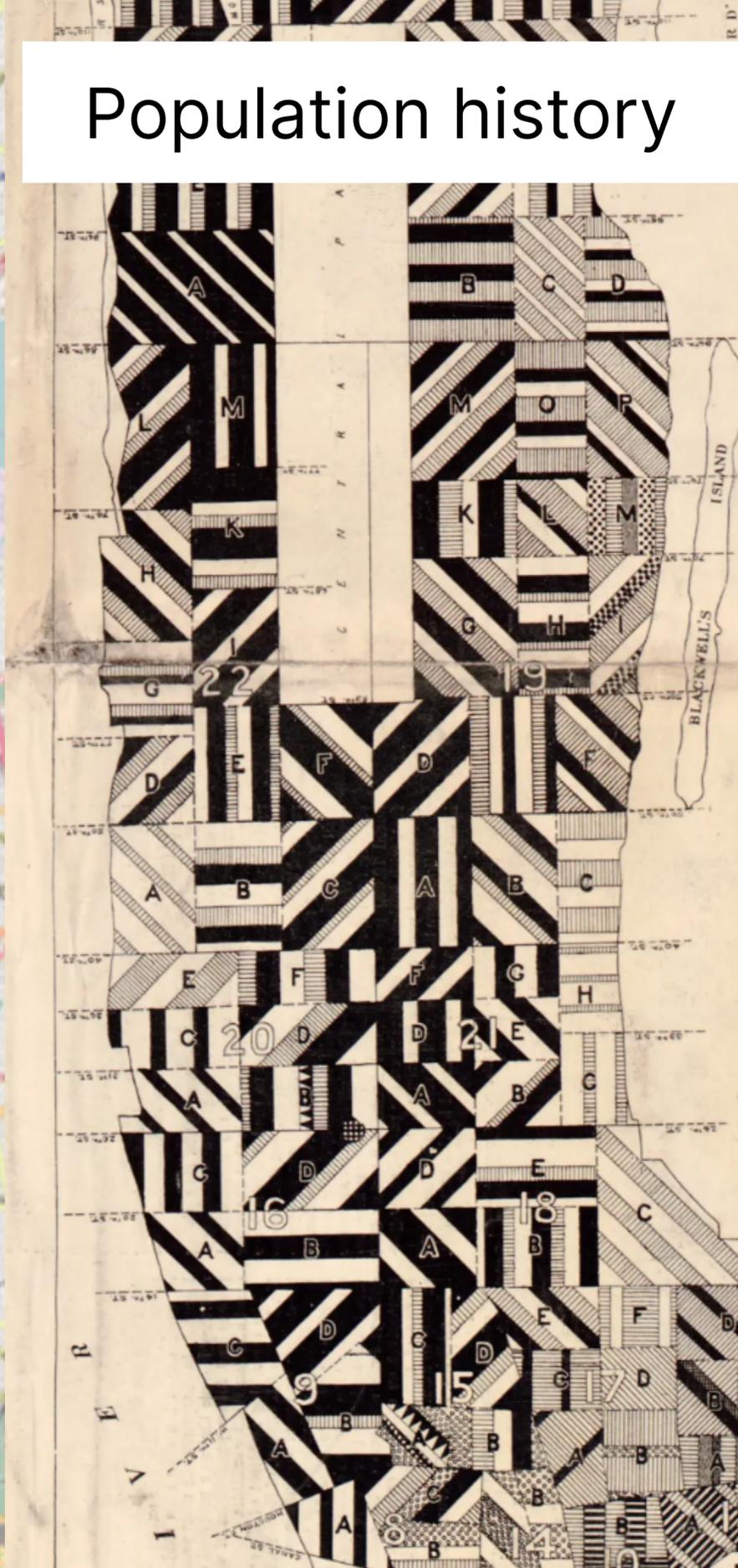




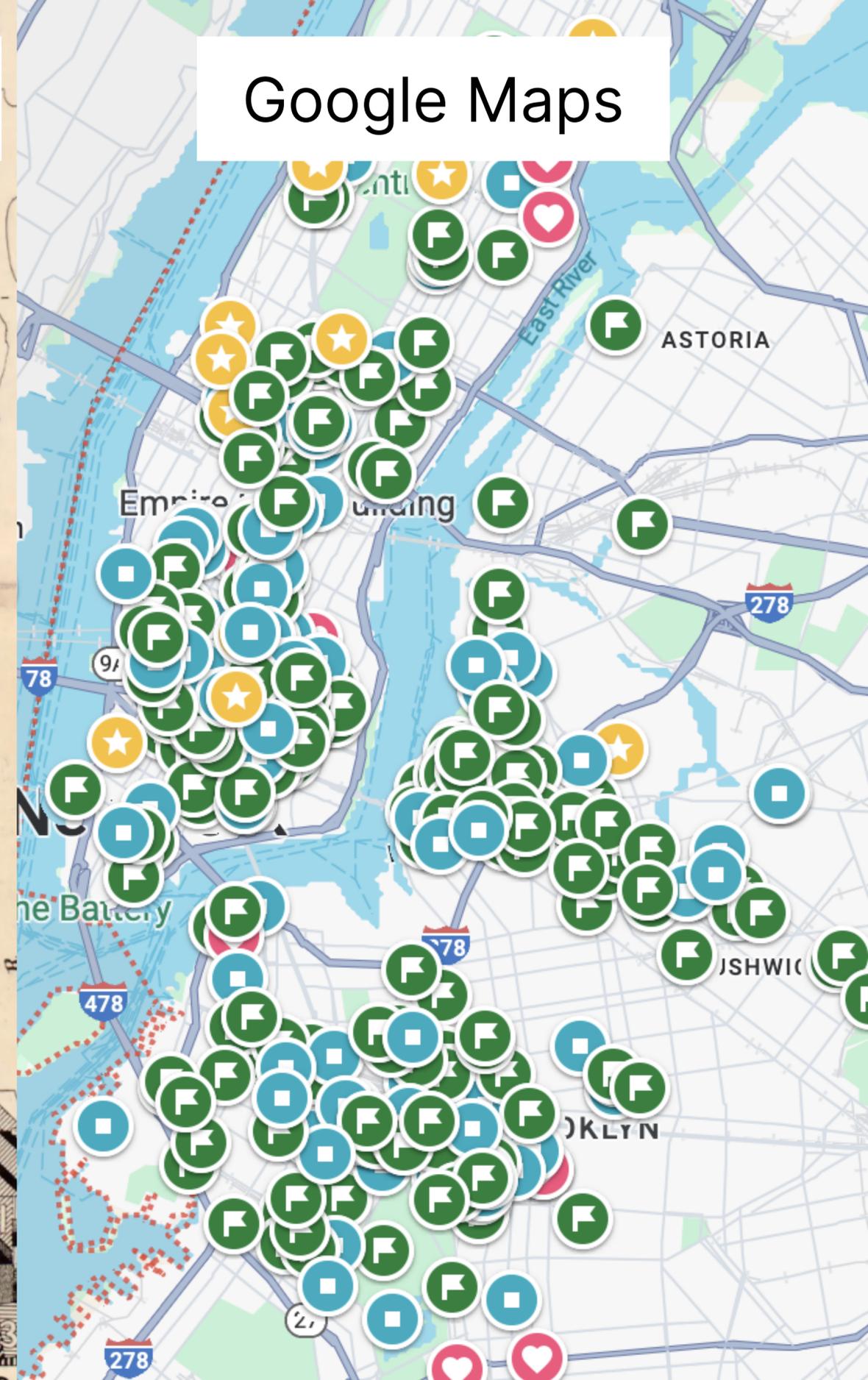
Highest-rated pizza



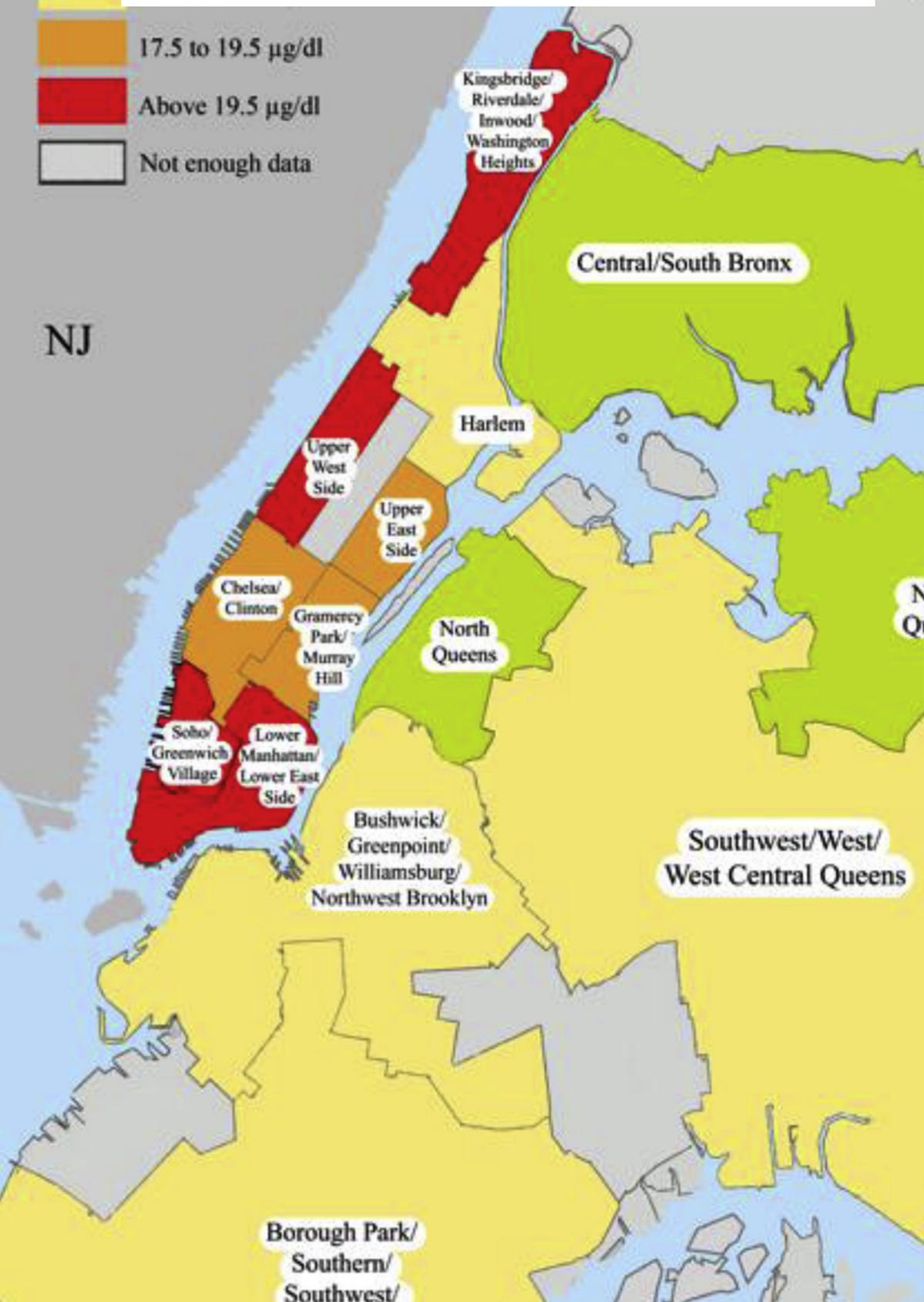
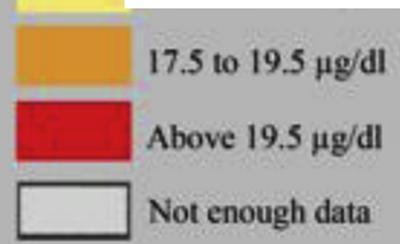
Population history



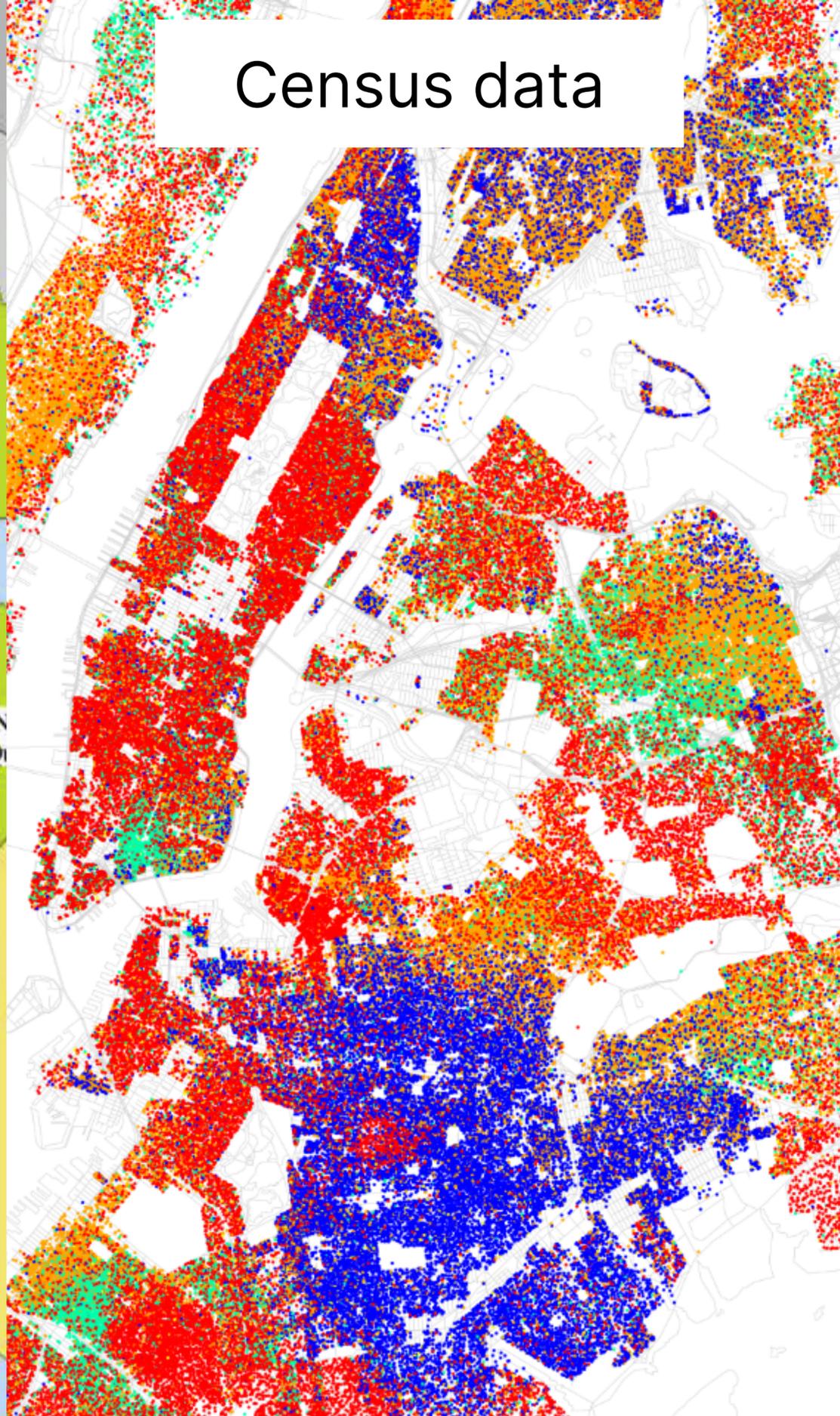
Google Maps



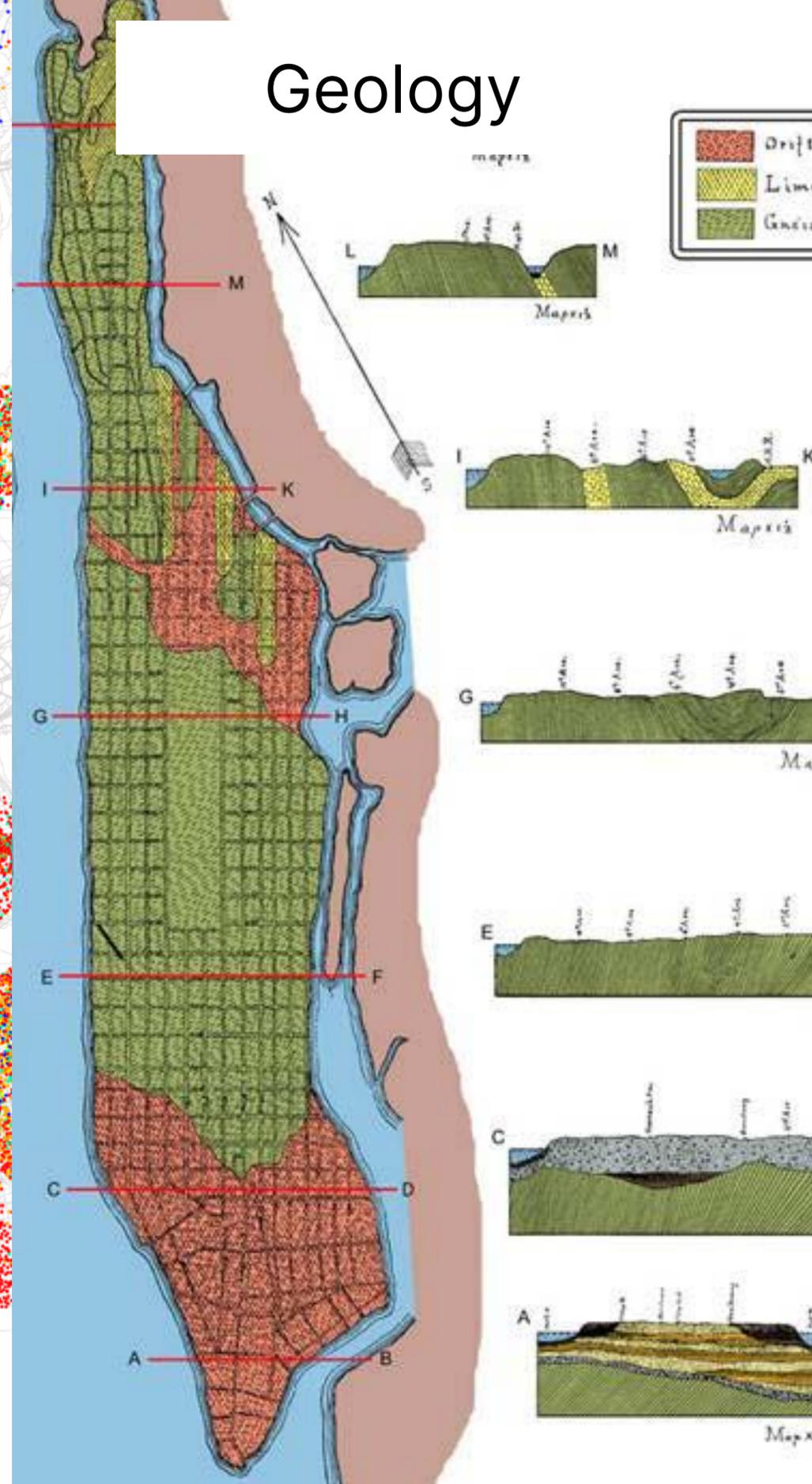
# Lead levels in pigeons



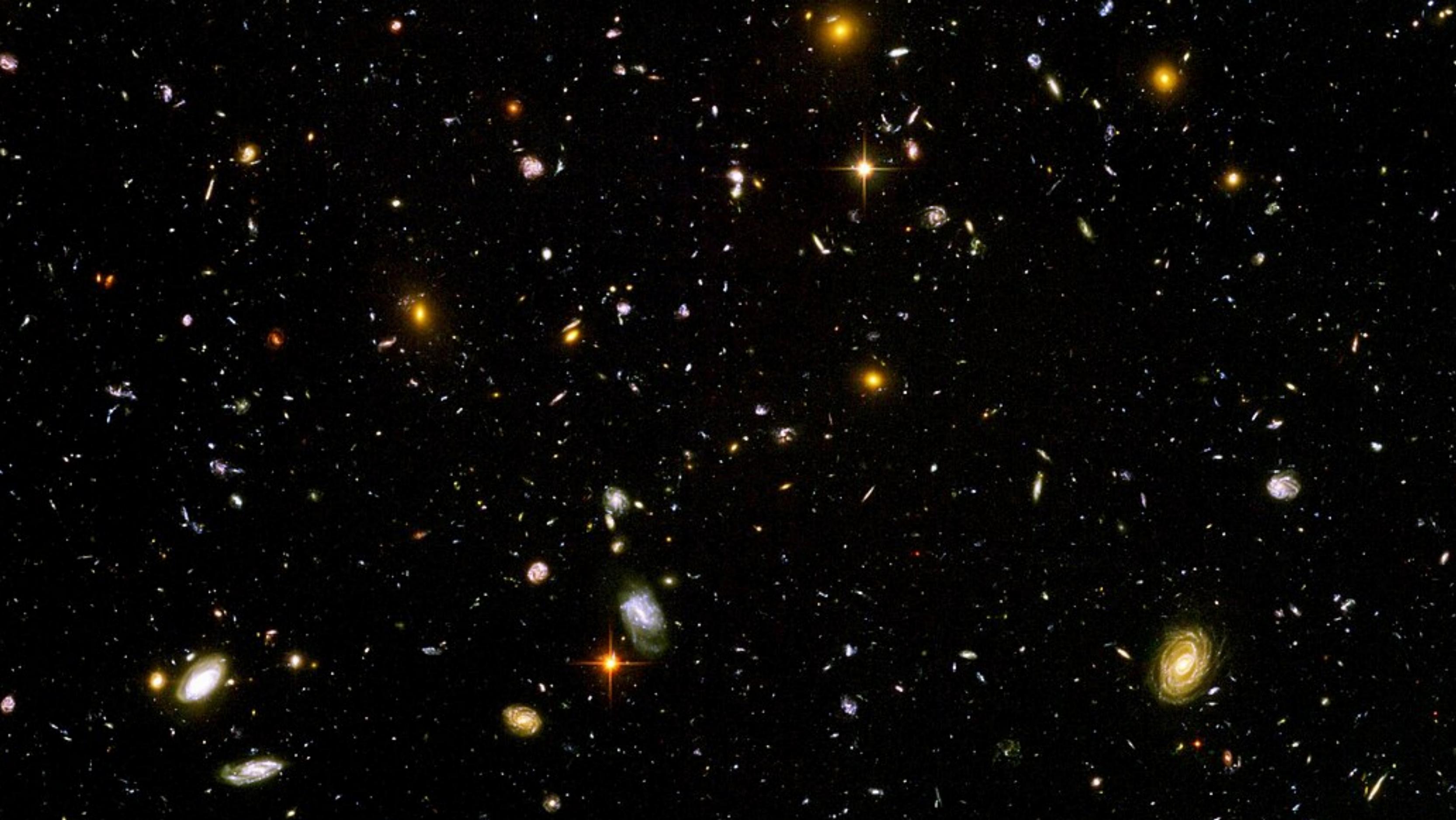
# Census data



# Geology



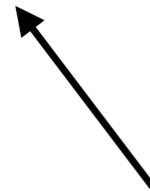




(multi-detector &  
multi-disciplinary)

(multimodal)

(self-supervised)



I want to construct a **wide-ranging, multi-layered, data-driven** representation of our Universe that can both complement and deepen our existing scientific understanding.



||

Illustration by Vivek Thakker

ANNALS OF ARTIFICIAL INTELLIGENCE

# CHATGPT IS A BLURRY JPEG OF THE WEB

The AI industry has already established some effective methods for integrating text, images, video, speech, etc. across domains.

Will physics data present new challenges for these methods?

Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

**Highly multi-modal** → unify inputs across an entire detector

**Informed by universal physical laws** → multi-detector/disciplinary analyses

Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

Diverse data formats → custom vs. generic embeddings

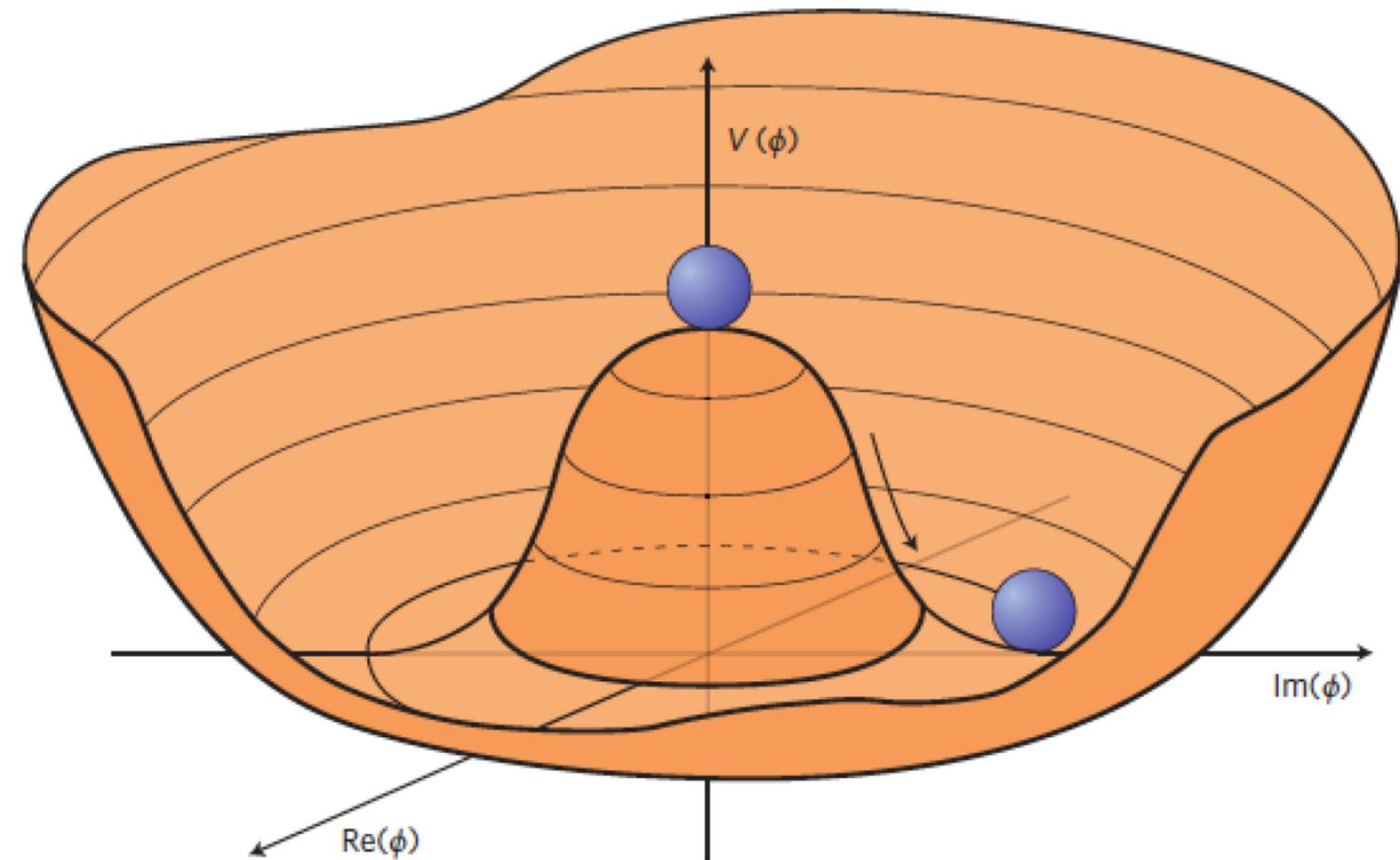
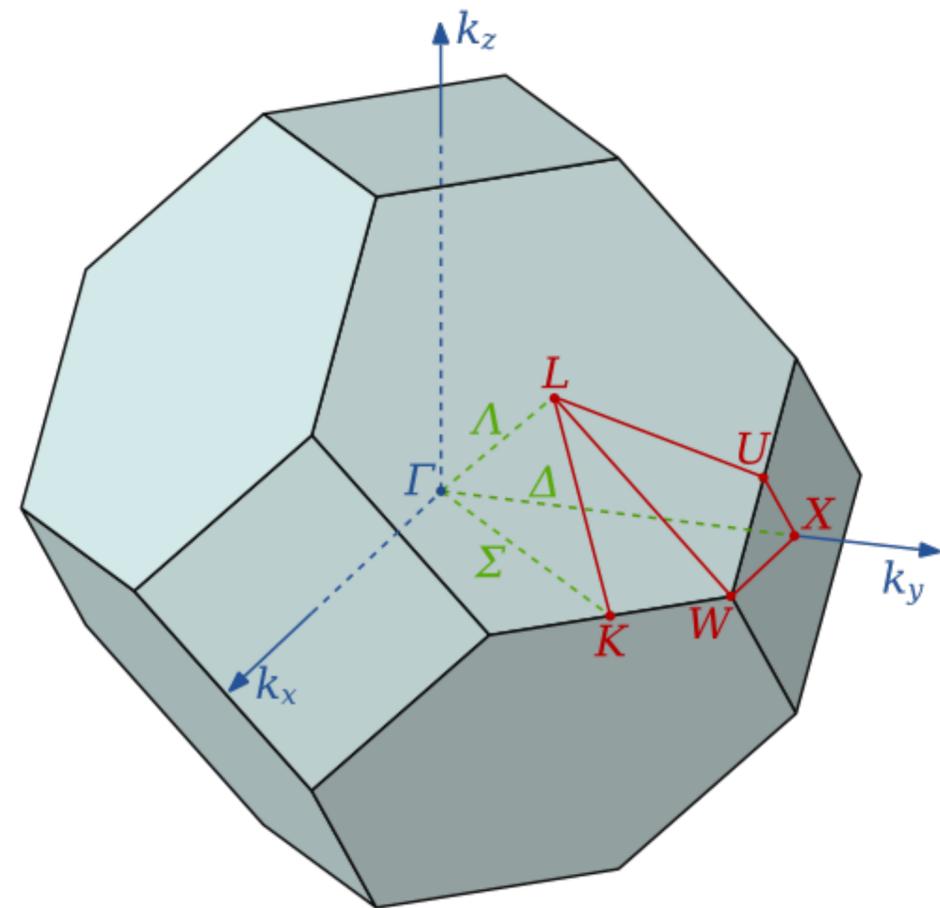
Detector-specific biases → publish data without detector effects

Group anomalies → statistical anomaly detection

Highly multi-modal → unify inputs across an entire detector

Informed by universal physical laws → multi-detector/disciplinary analyses

Unlike most fields, physics is grounded in deep & well-known mathematical frameworks.



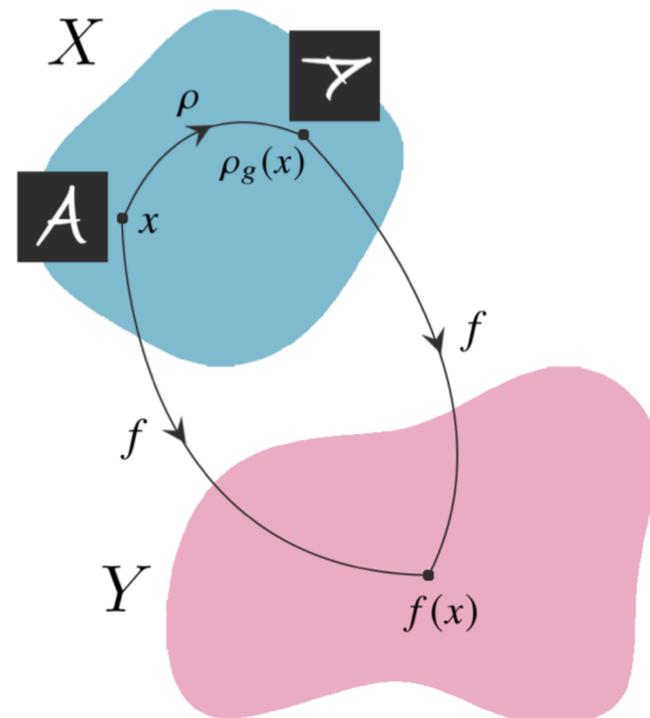
(particularly useful: symmetries!)

Unlike most fields, physics is grounded in deep & well-known mathematical frameworks.

Symmetries can be implicitly or explicitly encoded into neural networks:

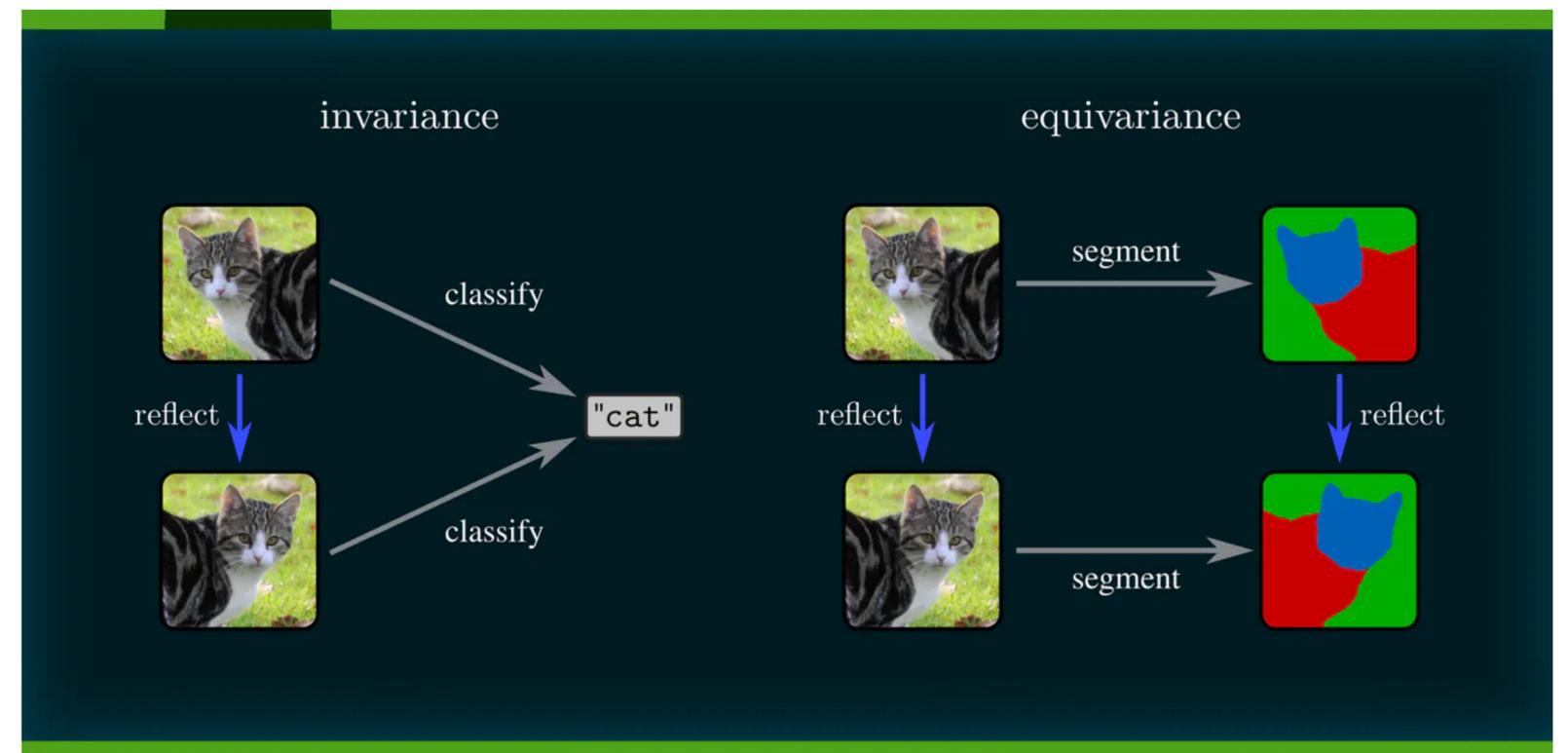
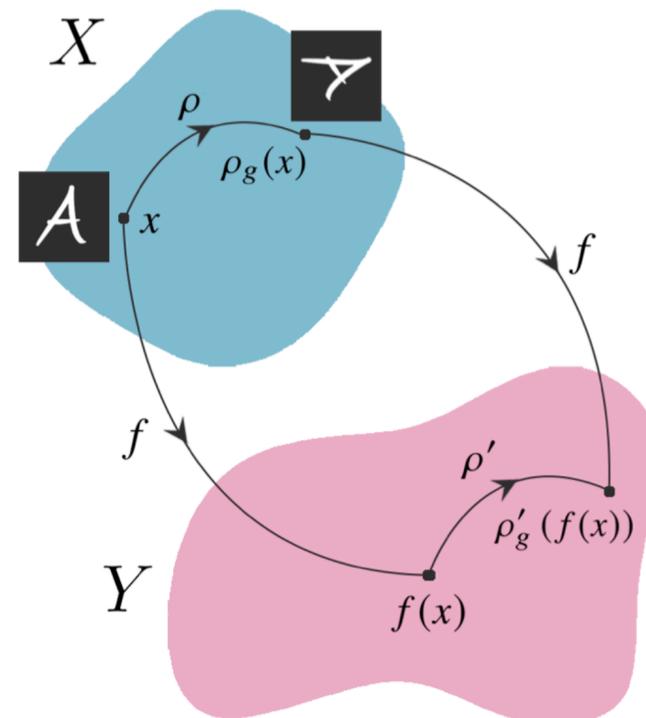
### Invariance

$$f(\rho_g(x)) = f(x)$$



### Equivariance

$$f(\rho_g(x)) = \rho'_g(f(x))$$



Unlike most fields, physics is grounded in deep & well-known mathematical frameworks.

We can do this even for much more complicated-looking symmetries:

$$SO^+(1, 3) \rtimes \mathbb{R}^{(1,3)} \times SU(3)_C \times SU(2)_L \times U(1)_Y$$



---

**Group Equivariant Convolutional Networks**

---

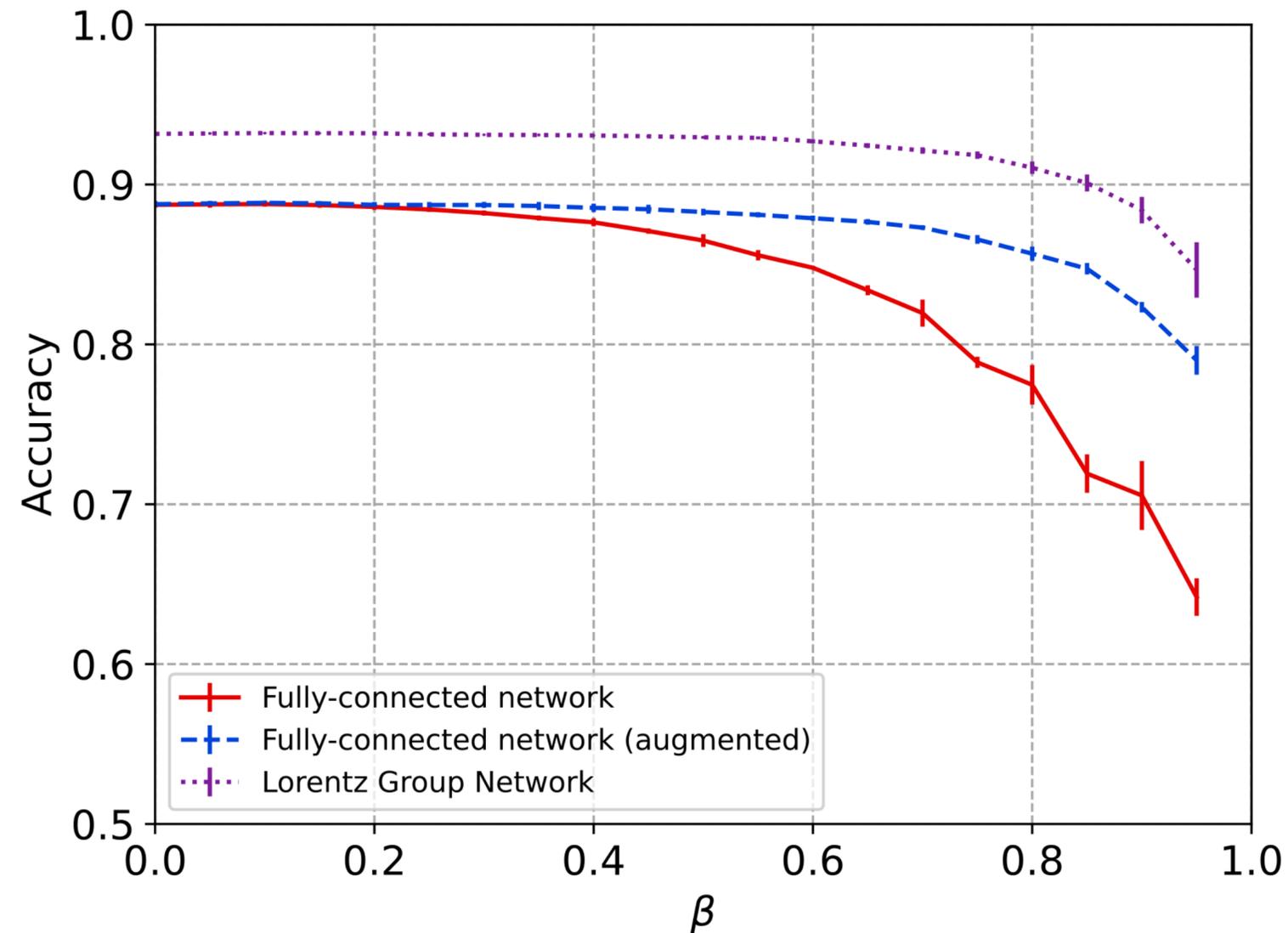
**Taco S. Cohen**  
University of Amsterdam

T.S.COHEN@UVA.NL

**Max Welling**  
University of Amsterdam  
University of California Irvine  
Canadian Institute for Advanced Research

M.WELLING@UVA.NL

Unlike most fields, physics is grounded in deep & well-known mathematical frameworks.



That said, the benefits of equivariant networks are most striking when data is highly limited — in a data-rich regime, the benefits are less clear.

Training %	Model	Accuracy	AUC	$1/\epsilon_B$
0.5%	ParticleNet	0.913	0.9687	$199 \pm 14$
	LorentzNet	<b>0.932</b>	<b>0.9793</b>	<b><math>562 \pm 72</math></b>
1%	ParticleNet	0.919	0.9734	$287 \pm 19$
	LorentzNet	<b>0.932</b>	<b>0.9812</b>	<b><math>697 \pm 58</math></b>
5%	ParticleNet	0.931	0.9807	$609 \pm 35$
	LorentzNet	<b>0.937</b>	<b>0.9839</b>	<b><math>1108 \pm 84</math></b>

# What level of structure is best for representing our data?

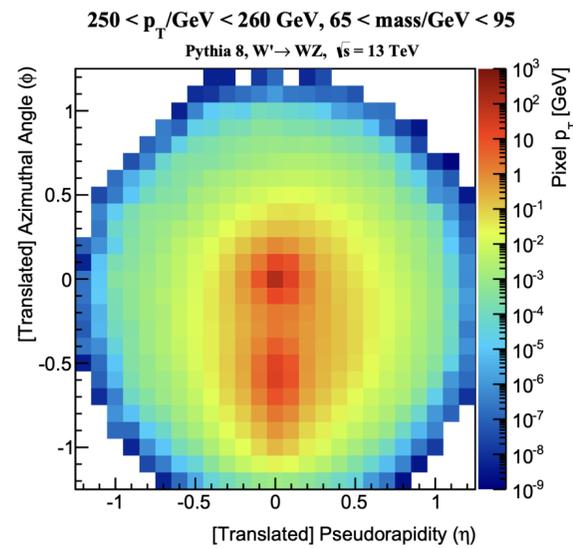
## Prescribed structure

## Learned structure

### Industry-standard formats

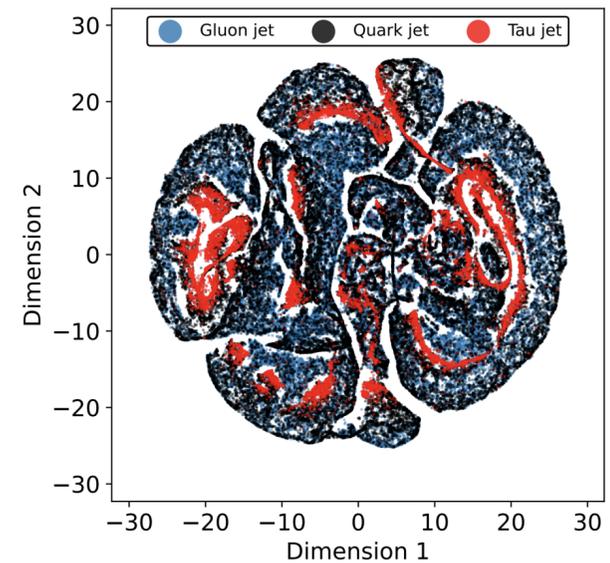
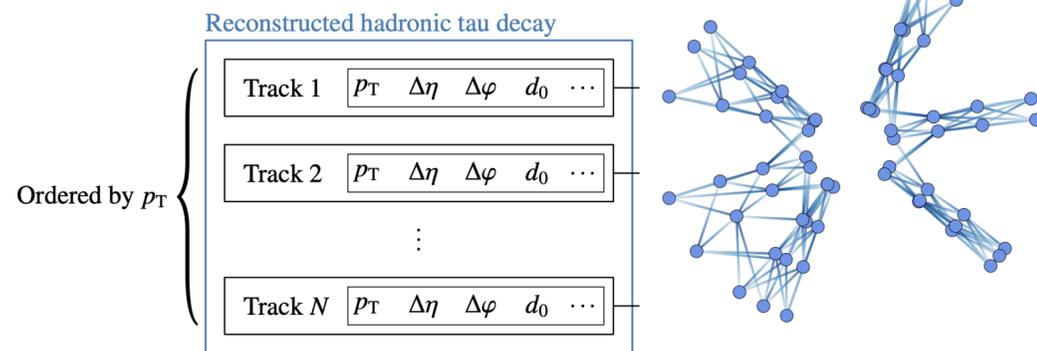
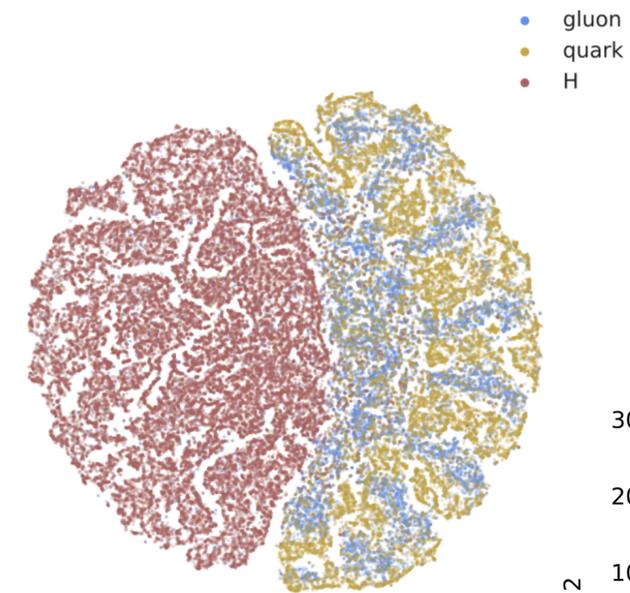
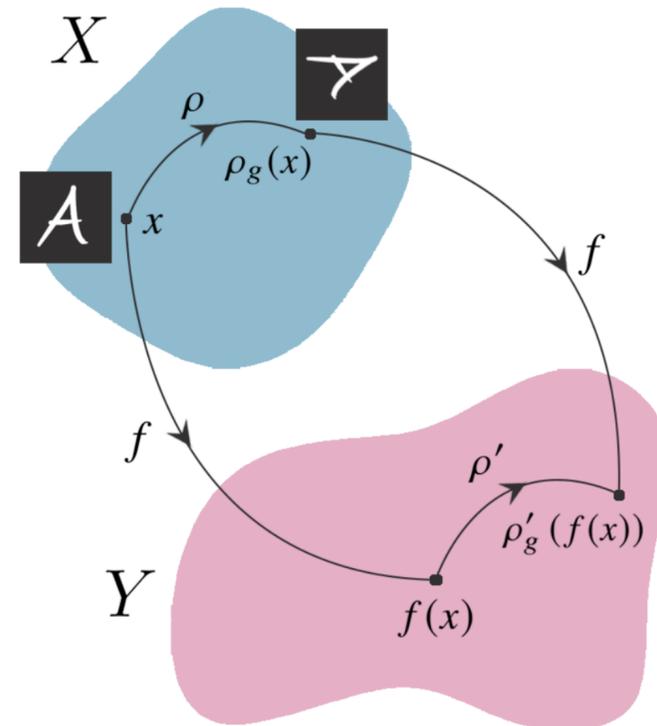
### Physics-informed

### Self-supervised



### Equivariance

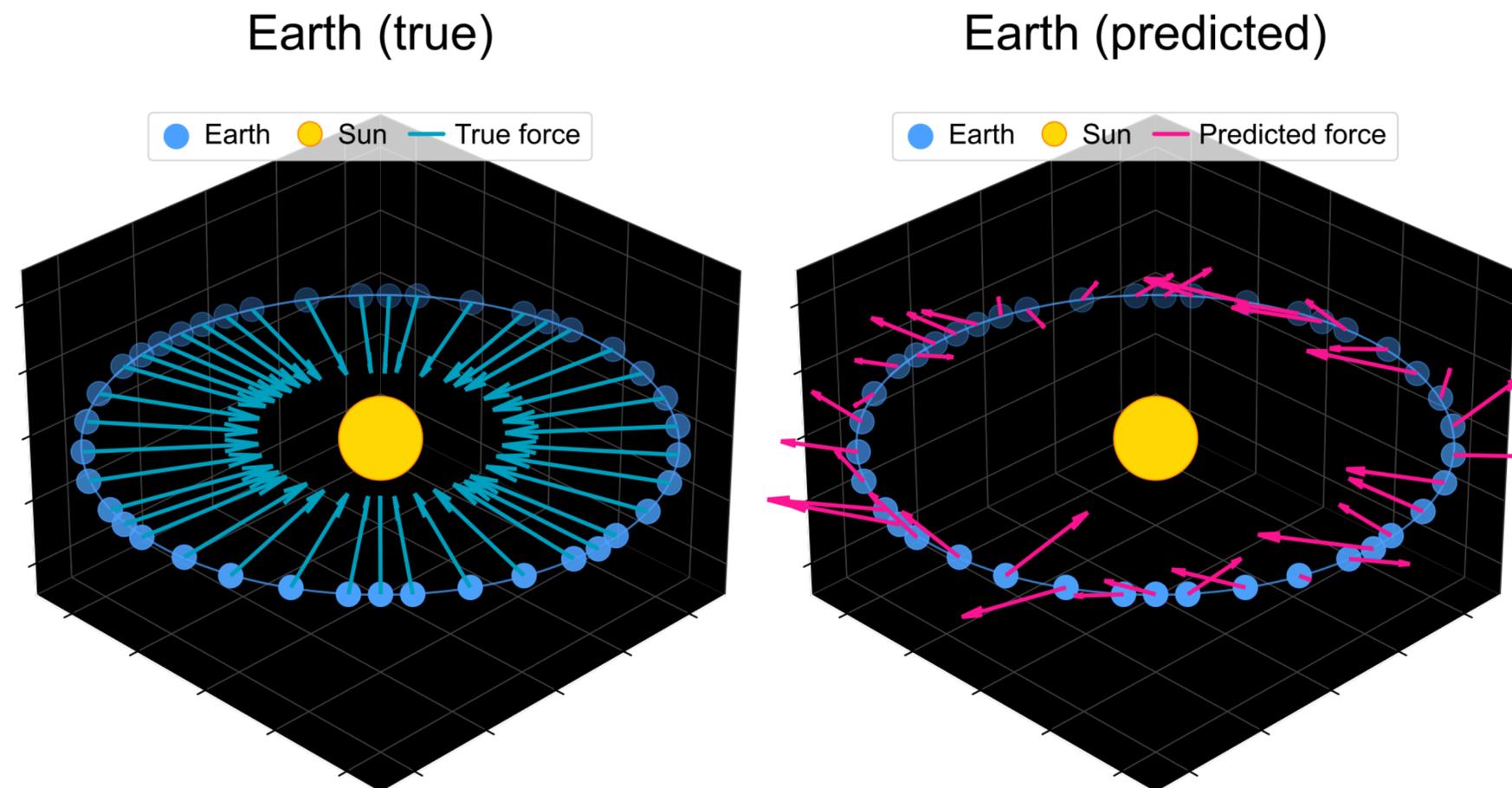
$$f(\rho_g(x)) = \rho'_g(f(x))$$



Unlike most fields, physics is grounded in deep & well-known mathematical frameworks.

Another benefit: **we often have high-quality simulations!**

But also: **a rigorous standard for interpretability** (maybe)



Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

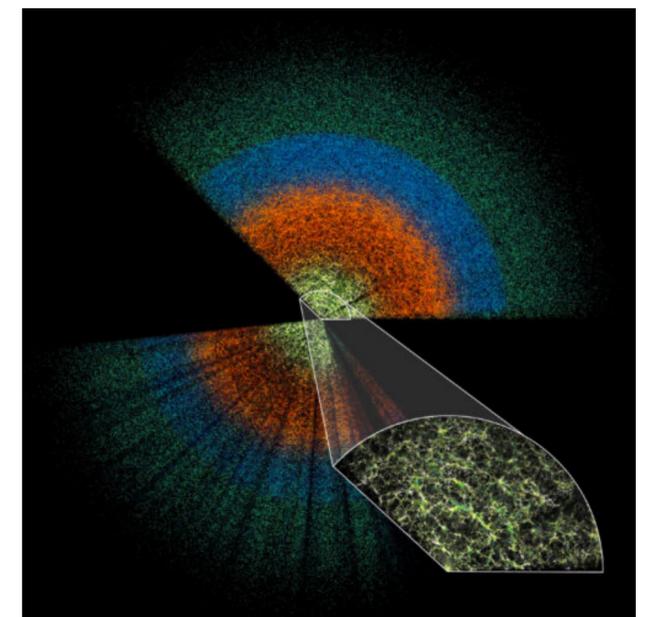
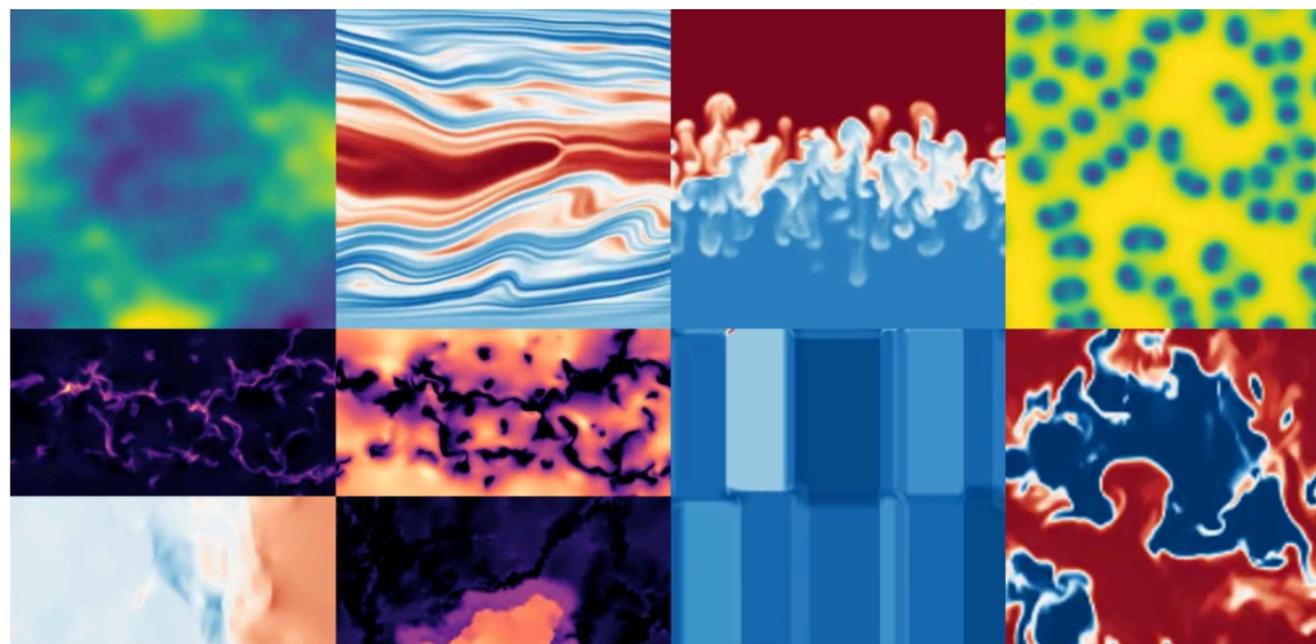
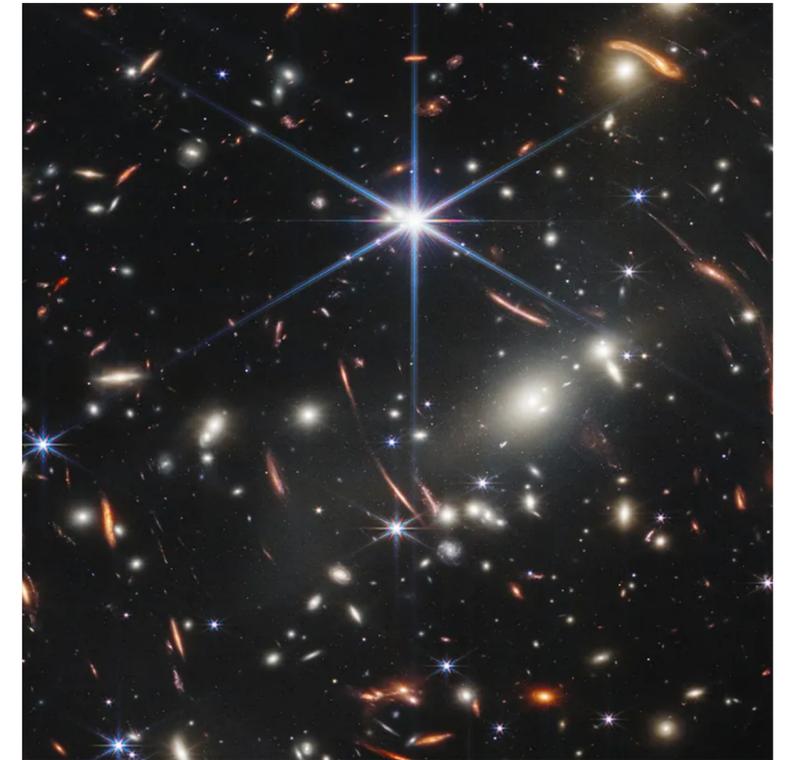
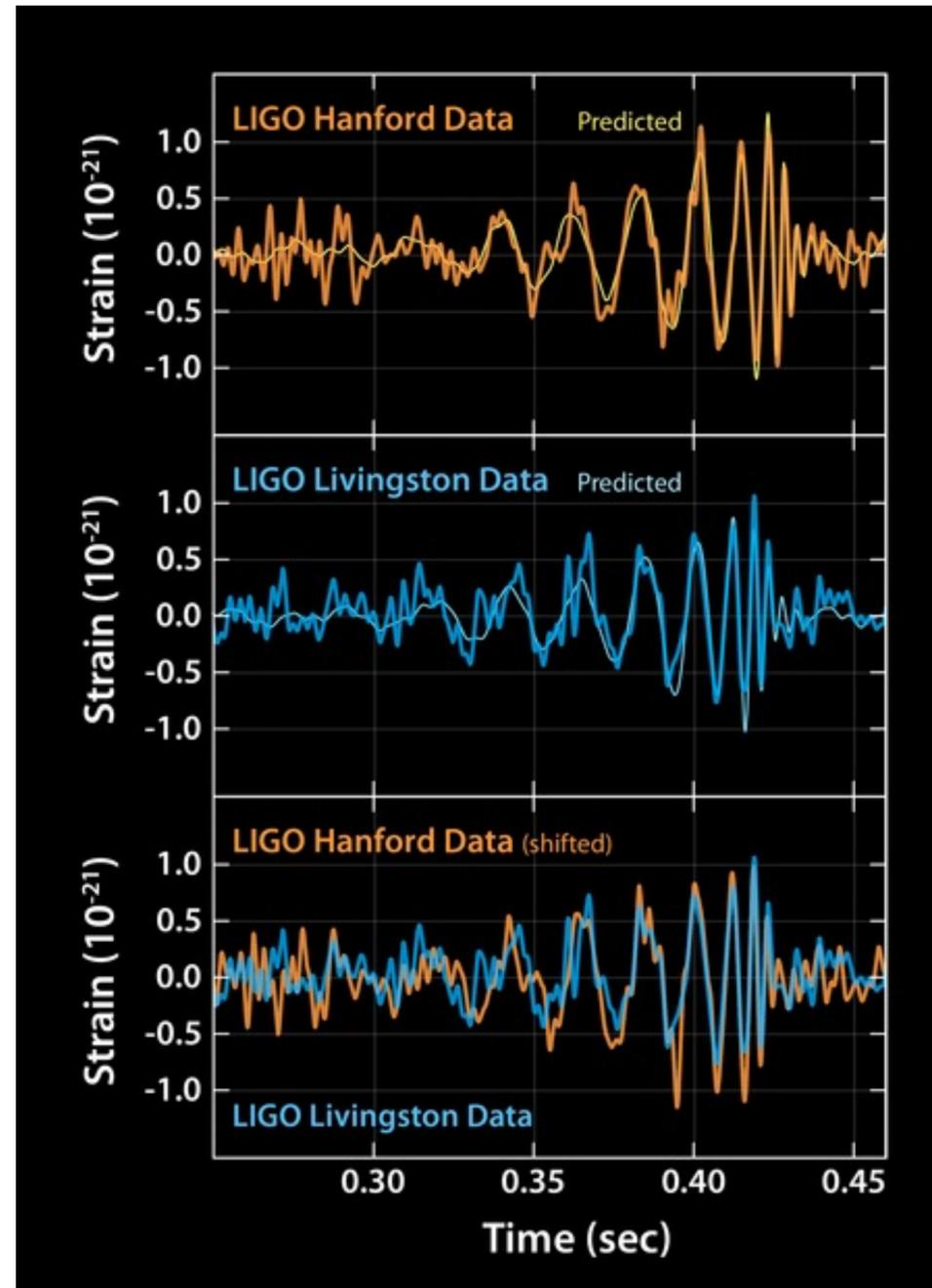
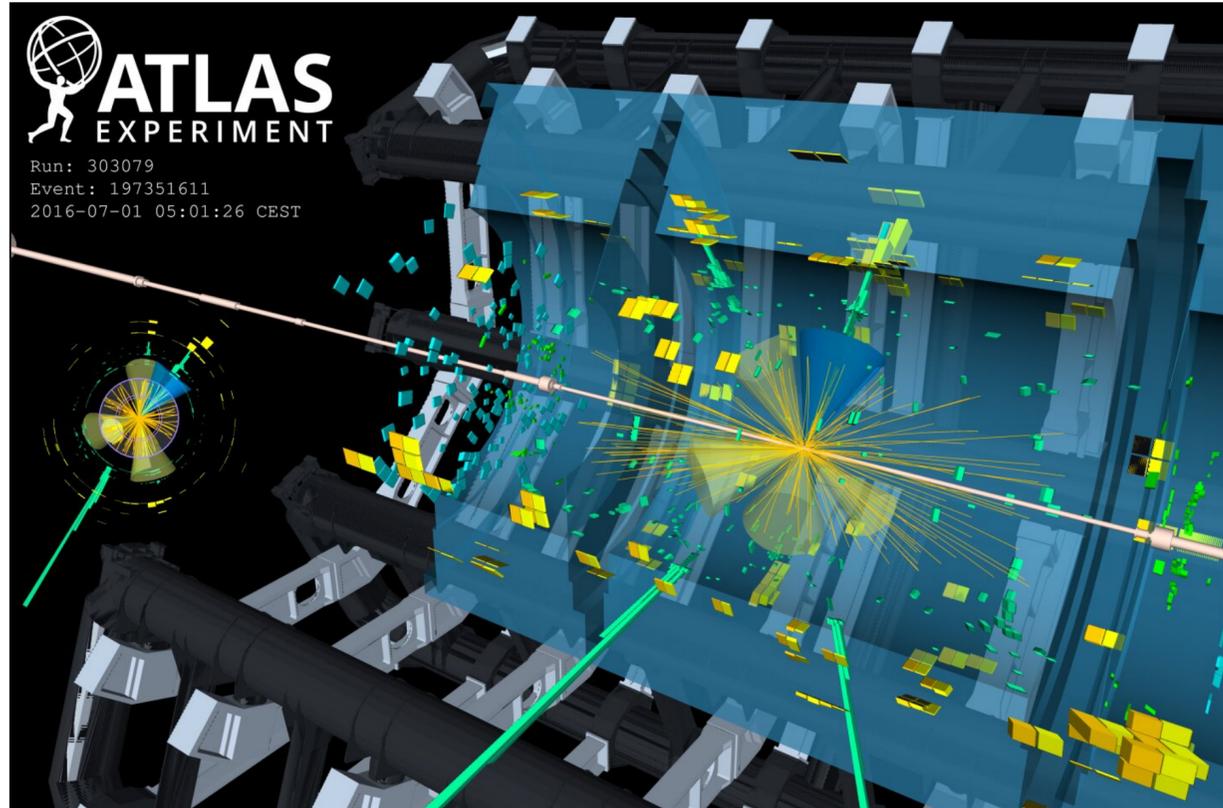
**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

**Highly multi-modal** → unify inputs across an entire detector

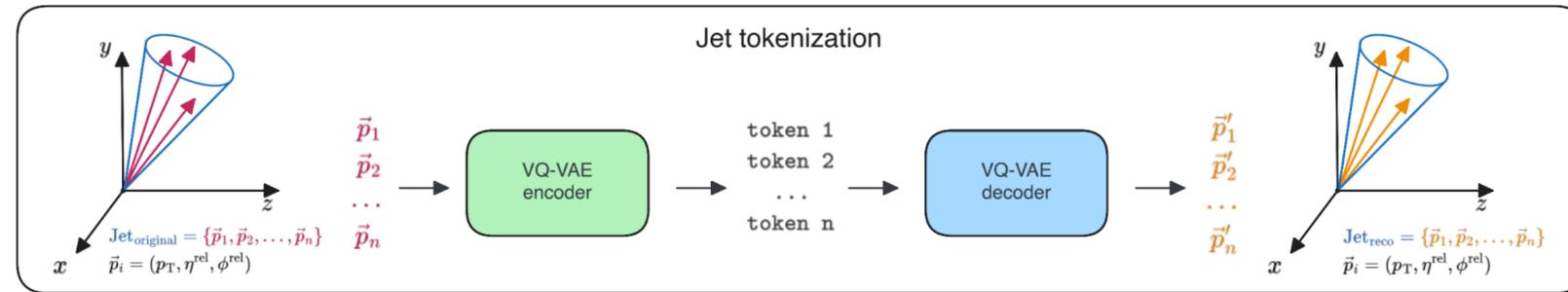
**Informed by universal physical laws** → multi-detector/disciplinary analyses

Physics data comes in all kinds of custom formats that are difficult to natively combine.

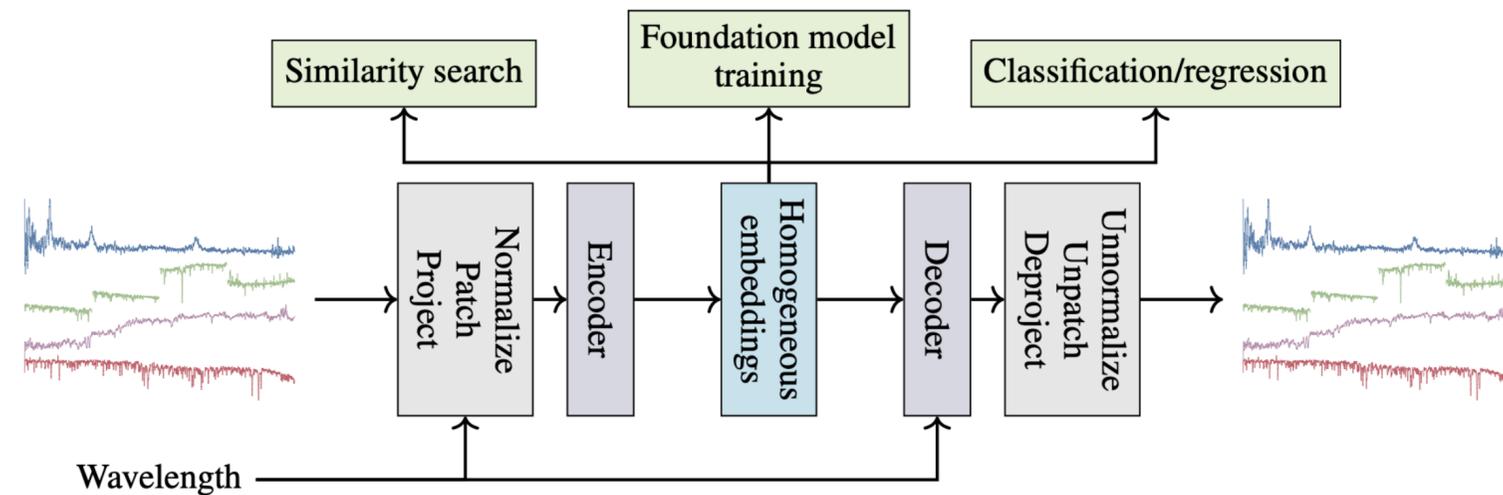


Physics data comes in all kinds of custom formats that are difficult to natively combine.

**Custom embeddings:**  
precise, but potentially tedious



[arXiv:2403.05618](https://arxiv.org/abs/2403.05618)



[arXiv:2510.17959](https://arxiv.org/abs/2510.17959)

**Generic embedding?**

```
{'description':{
  'planet0':{
    'm':4.842,'a':2.679,'e':0.923
  },
  'planet1':{
    'm':4.813,'a':1.0,'e':0.191
  },
  'stepsize': 0.8},
'data':[
  [[2.525,-0.397],[0.963,-0.232]],
  [[2.557,0.112],[0.830,0.542]],
  [[2.483,0.618],[0.1832,0.981]],
  [[2.318,1.100],[-0.577,0.823]],
  [[2.066,1.541],[-0.994,0.173]]
]}
```

Physics data comes in all kinds of custom formats that are difficult to natively combine.

Rendering data as text is easy to combine, but most language tokenizers hinder mathematical abilities.

e.g. GPT-4 fails at 4-digit multiplication over 95% of the time without tool use:

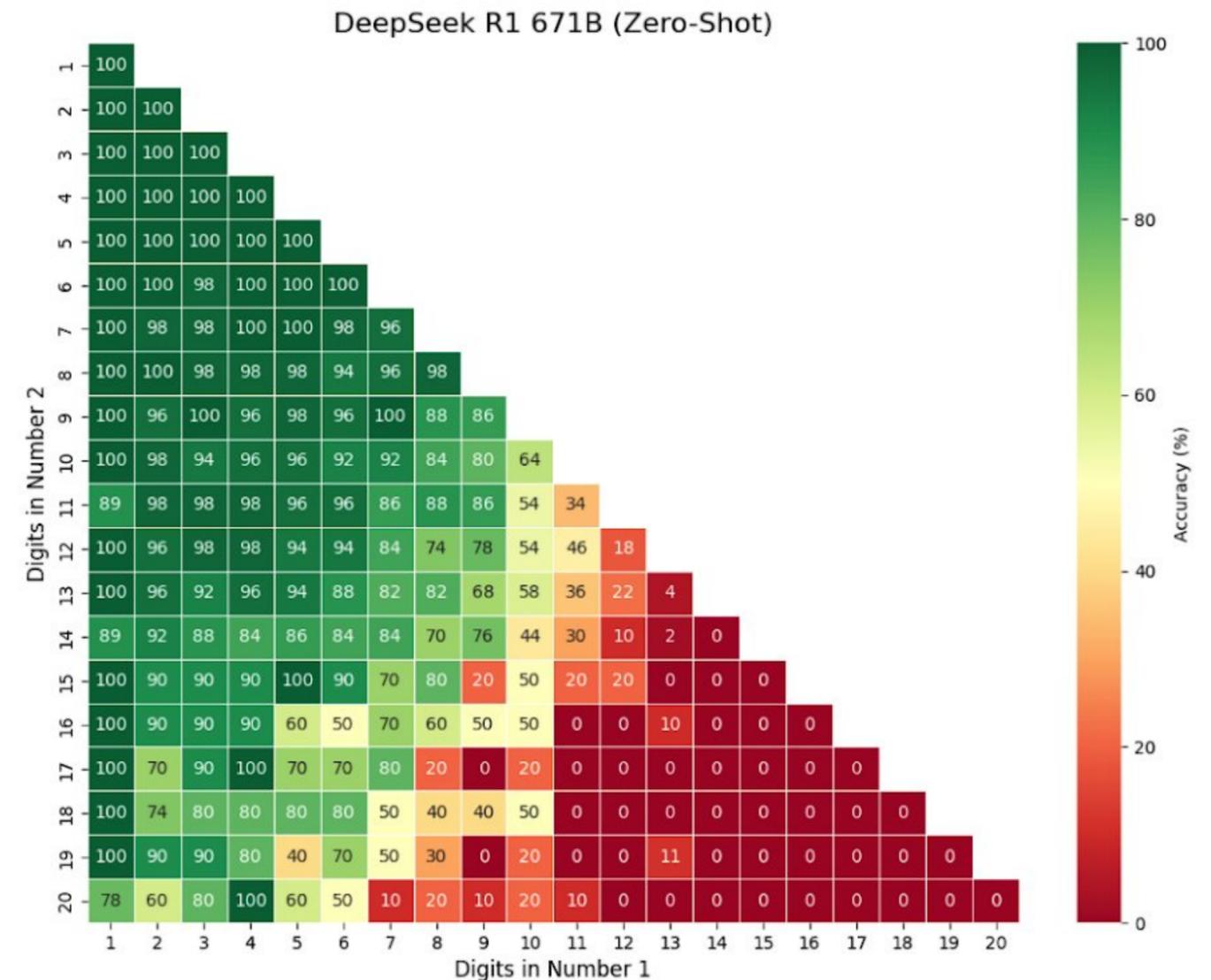
**You**  
What is  $3582 \times 2911$ ?

**Wrong! It's 10,427,202.**

**ChatGPT**  
3582 multiplied by 2911 equals 10,425,302.



[arXiv:2305.18654](https://arxiv.org/abs/2305.18654)



Physics data comes in all kinds of custom formats that are difficult to natively combine.

LLMs struggle to understand what makes numbers distinct from other language-based tokens.

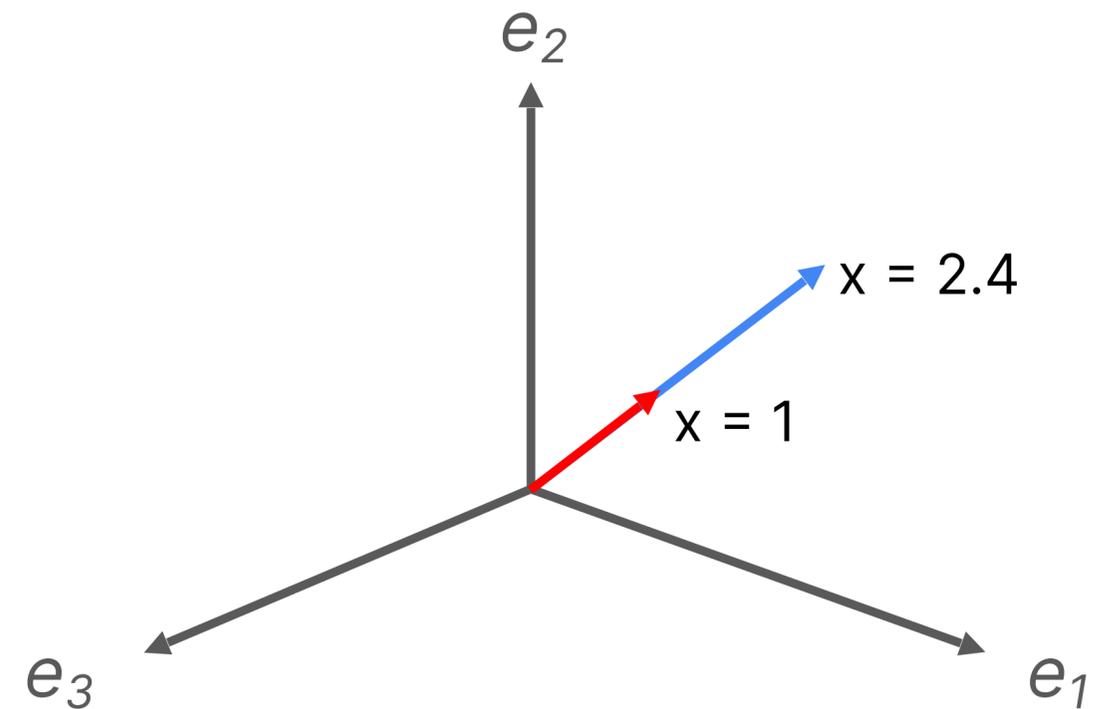
Moreover, they struggle to represent numbers with a finite dictionary of tokens:



Encoding	3.14	$-6.02 \cdot 10^{23}$	Tokens / coefficient	Size of vocabulary
P10	[+, 3, 1, 4, E-2]	[-, 6, 0, 2, E21]	5	210
P1000	[+, 314, E-2]	[-, 602, E21]	3	1100
B1999	[314, E-2]	[-602, E21]	2	2000
FP15	[FP314/-2]	[FP-602/21]	1	30000

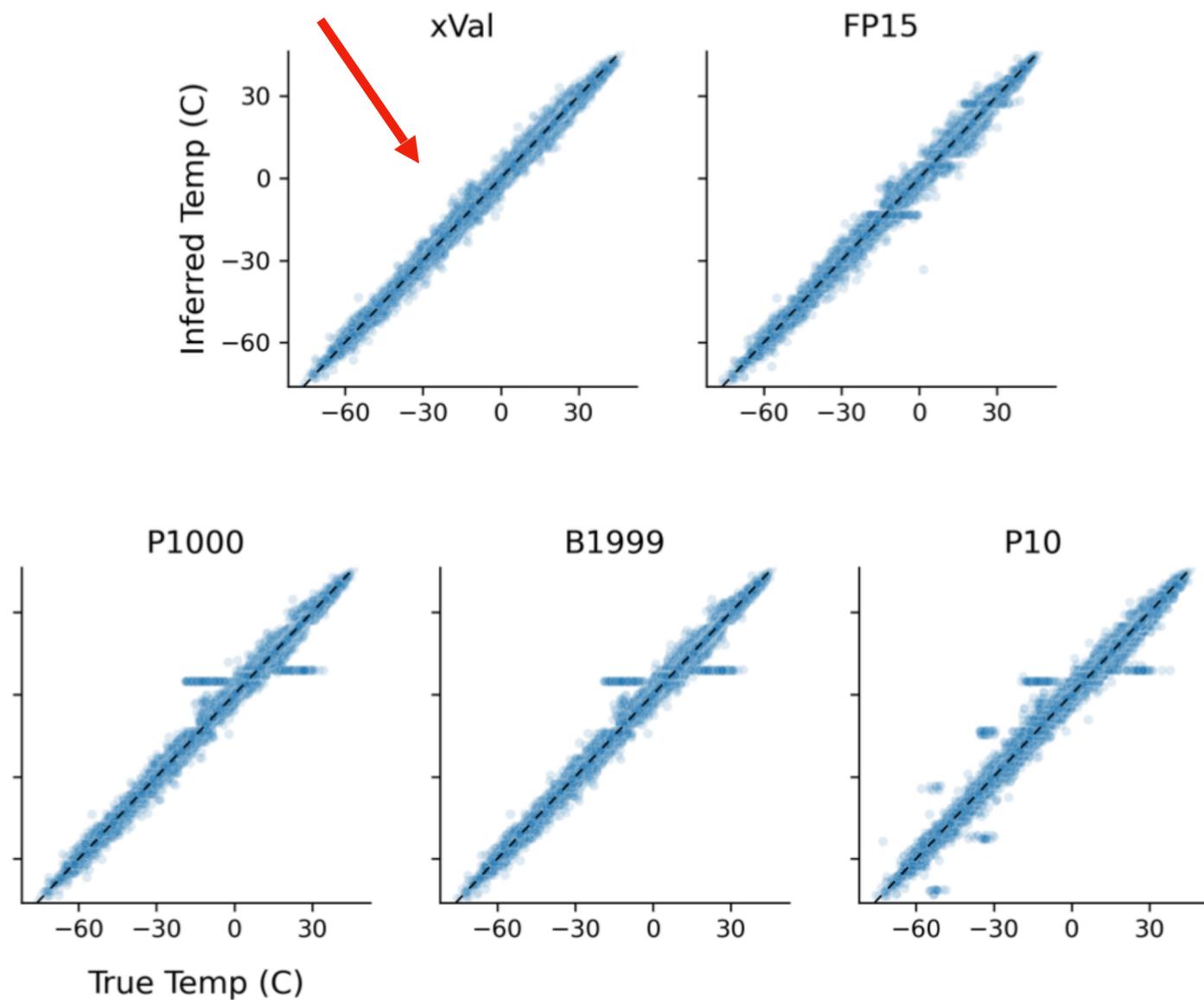
Physics data comes in all kinds of custom formats that are difficult to natively combine.

One possibility is to add a dedicated numerical token with a magnitude that scales with the number itself:

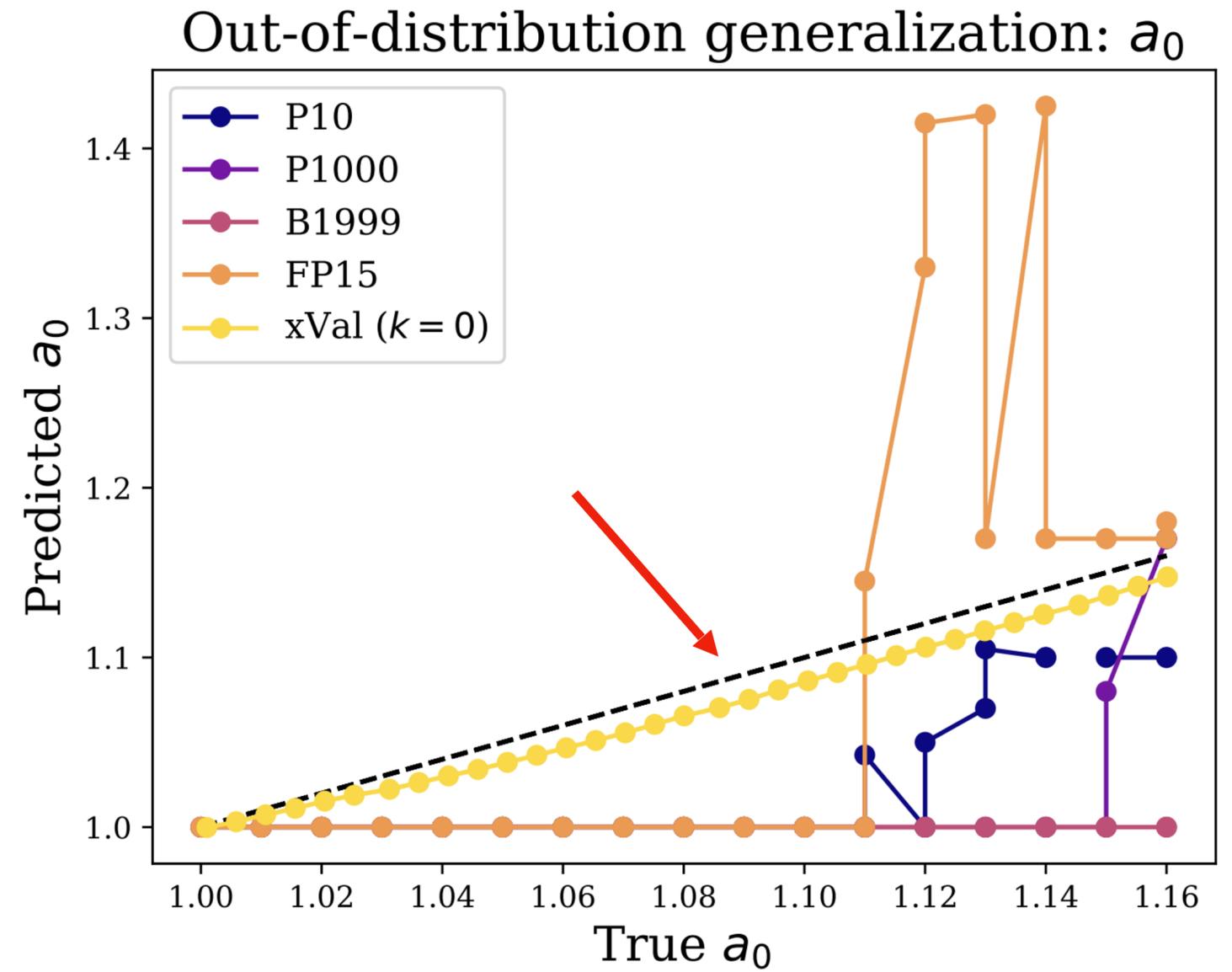


Dedicated numerical encodings can help enable rendering physics data as language.

## In-distribution

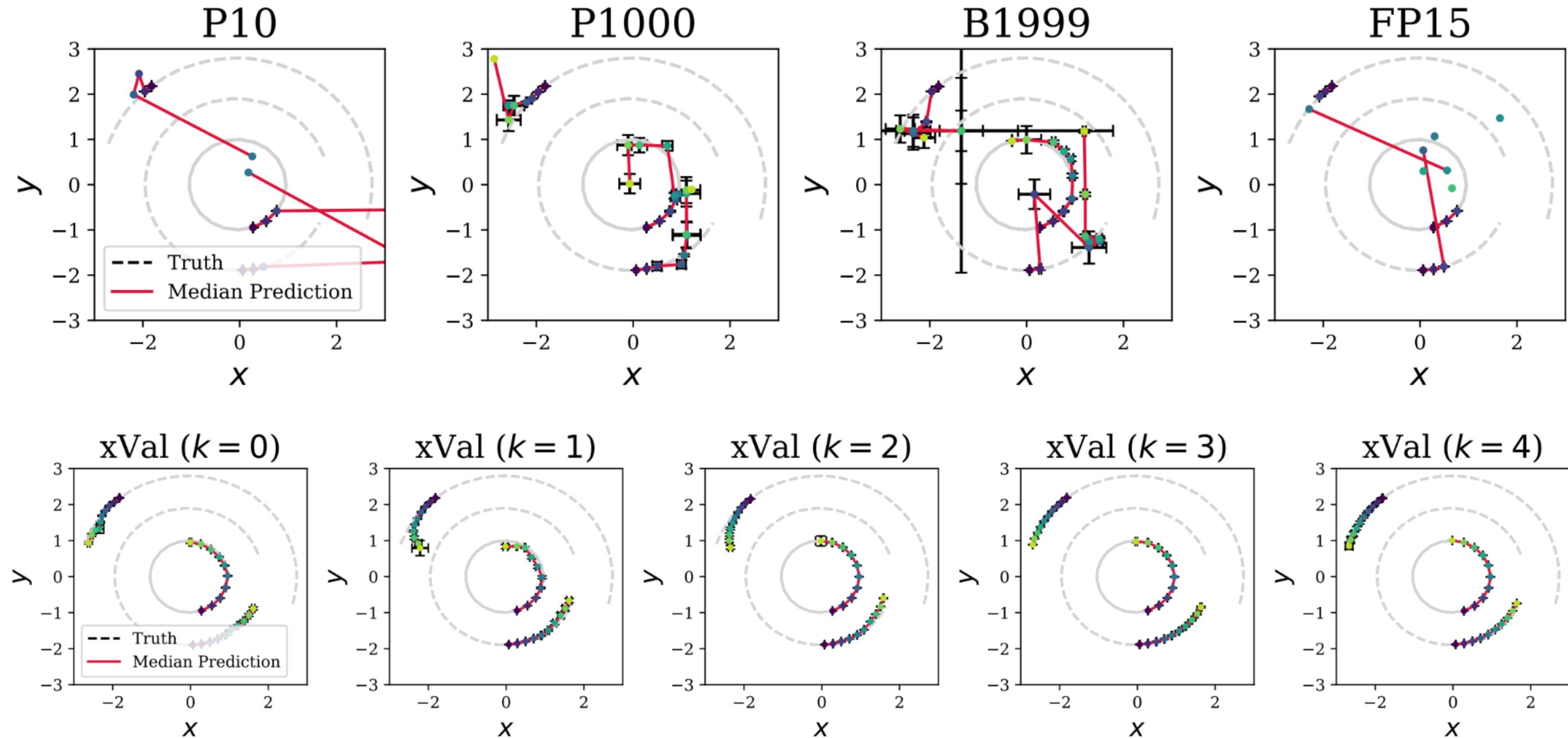


## Out-of-distribution

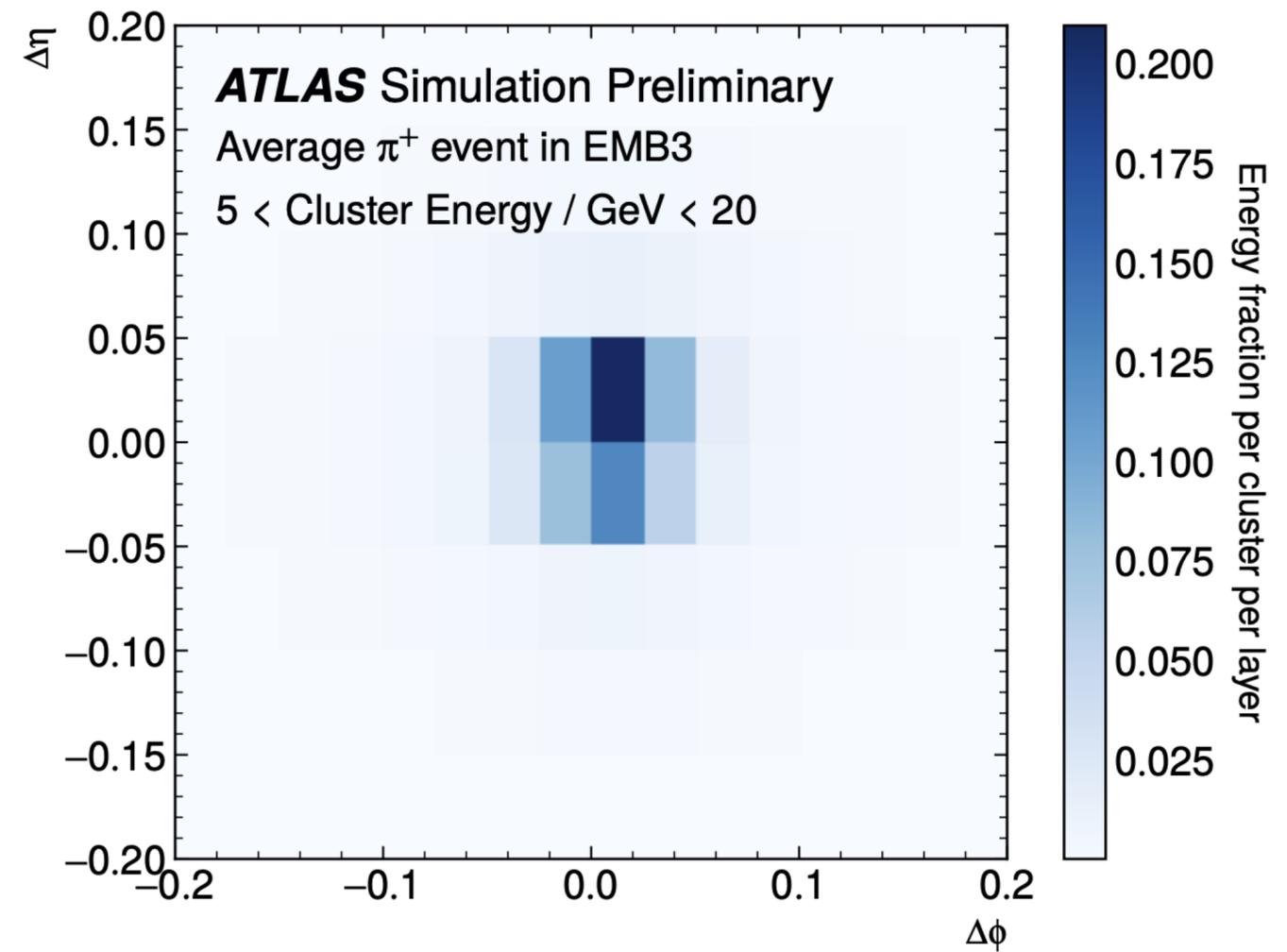
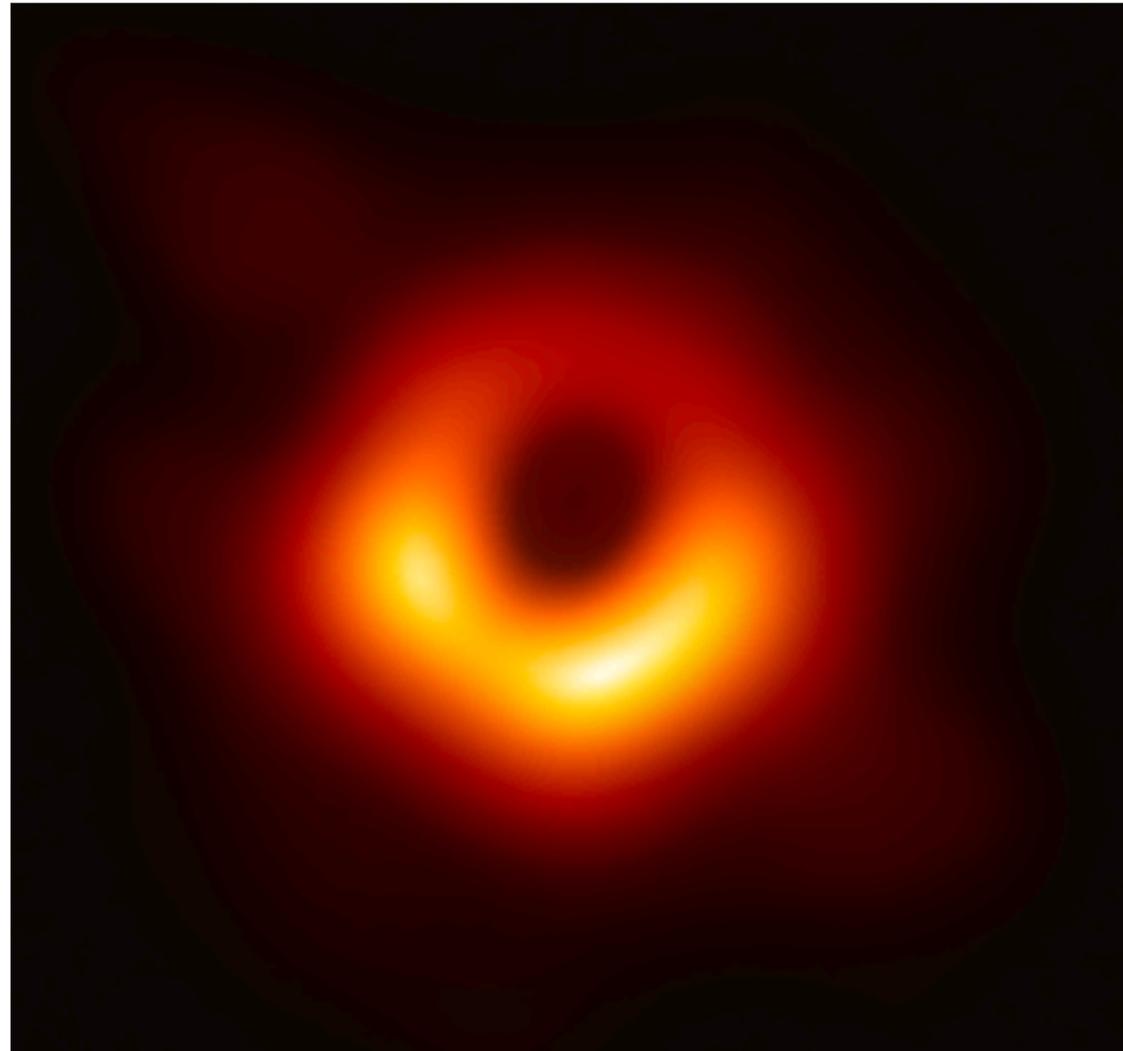


Dedicated numerical encodings can help enable rendering physics data as language.

This can even predict time series data, despite not being trained autoregressively:



Moreover, physics data often needs context — perhaps we should treat text as a critical modality:



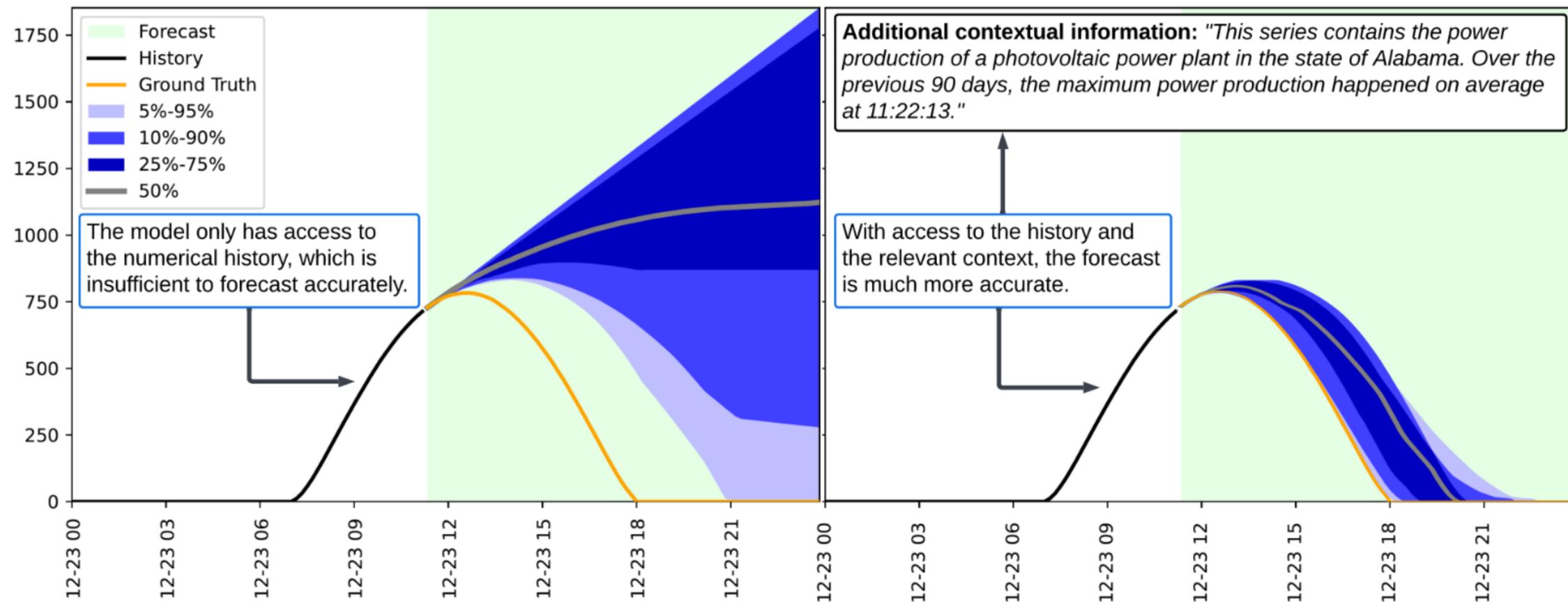
Moreover, physics data often needs context — perhaps we should treat text as a critical modality:

---

## Context is Key: A Benchmark for Forecasting with Essential Textual Information

---

\*Andrew R. Williams<sup>1 2 3</sup> \*Arjun Ashok<sup>1 2 3</sup>  
†Étienne Marcotte<sup>1</sup> †Valentina Zantedeschi<sup>1 4</sup> Jithendaraa Subramanian<sup>1 2 5</sup> Roland Riachi<sup>2</sup>  
James Requeima<sup>6</sup> Alexandre Lacoste<sup>1</sup> Irina Rish<sup>2 3</sup> †Nicolas Chapados<sup>1 2 7</sup> †Alexandre Drouin<sup>1 2 4</sup>



Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

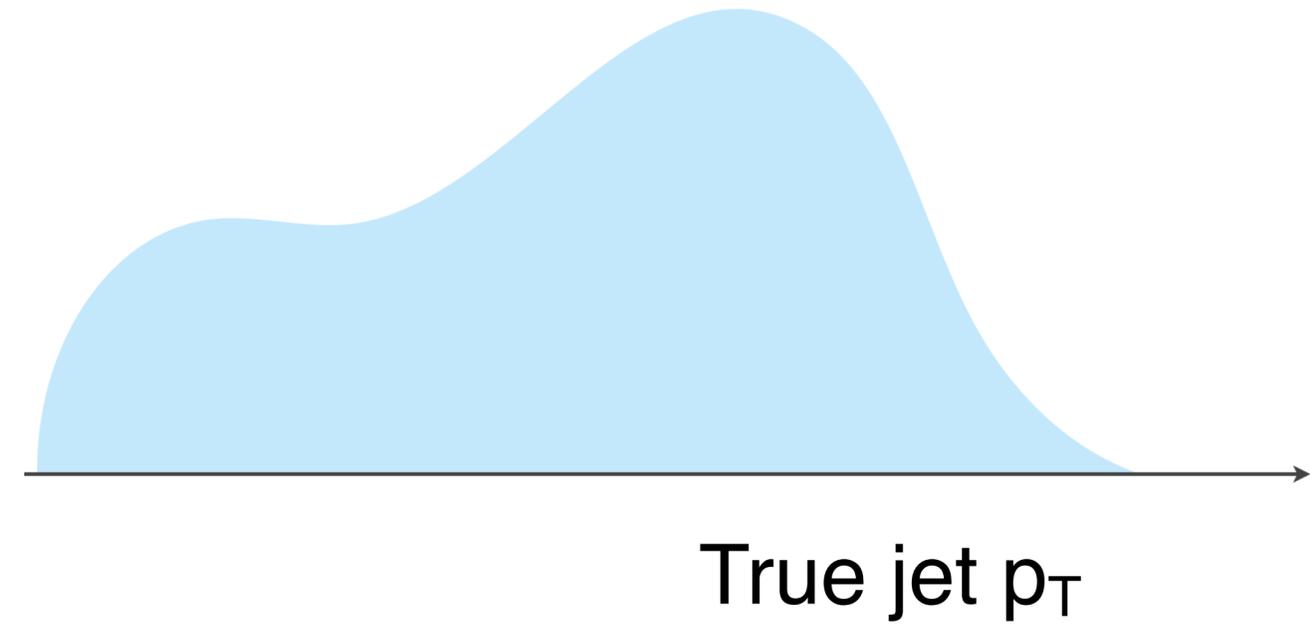
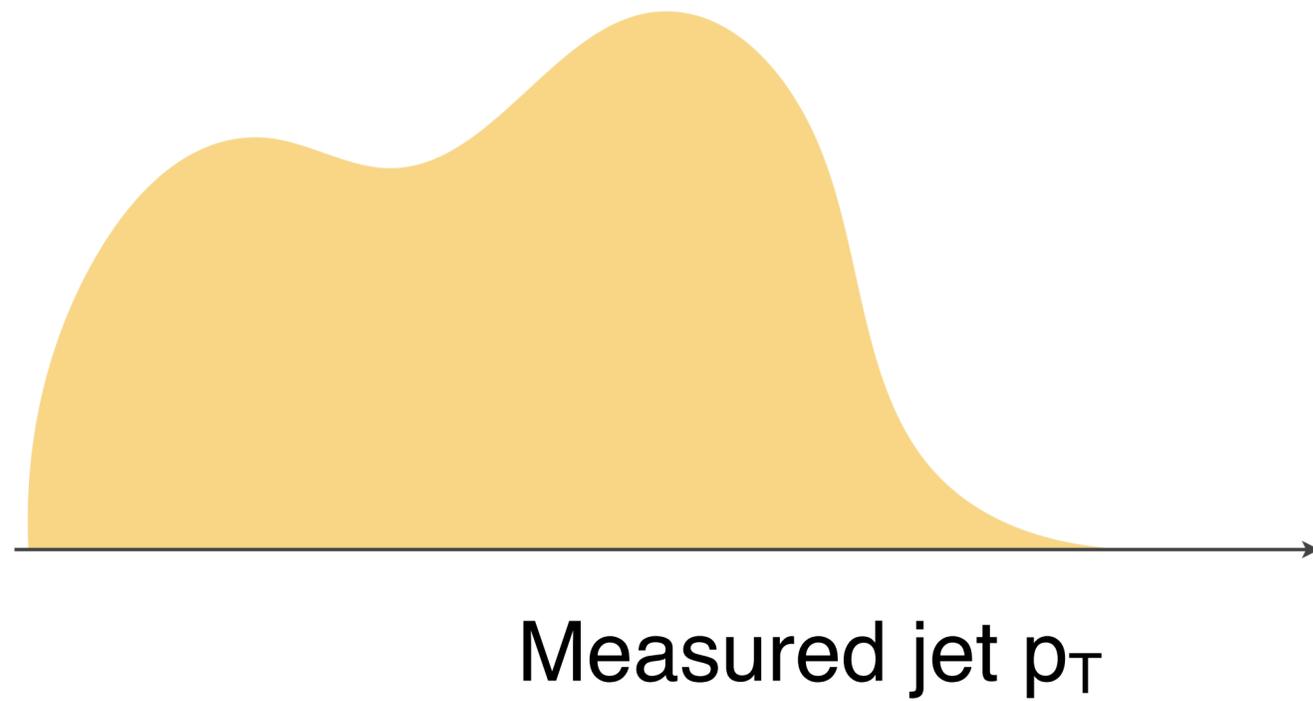
**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

**Highly multi-modal** → unify inputs across an entire detector

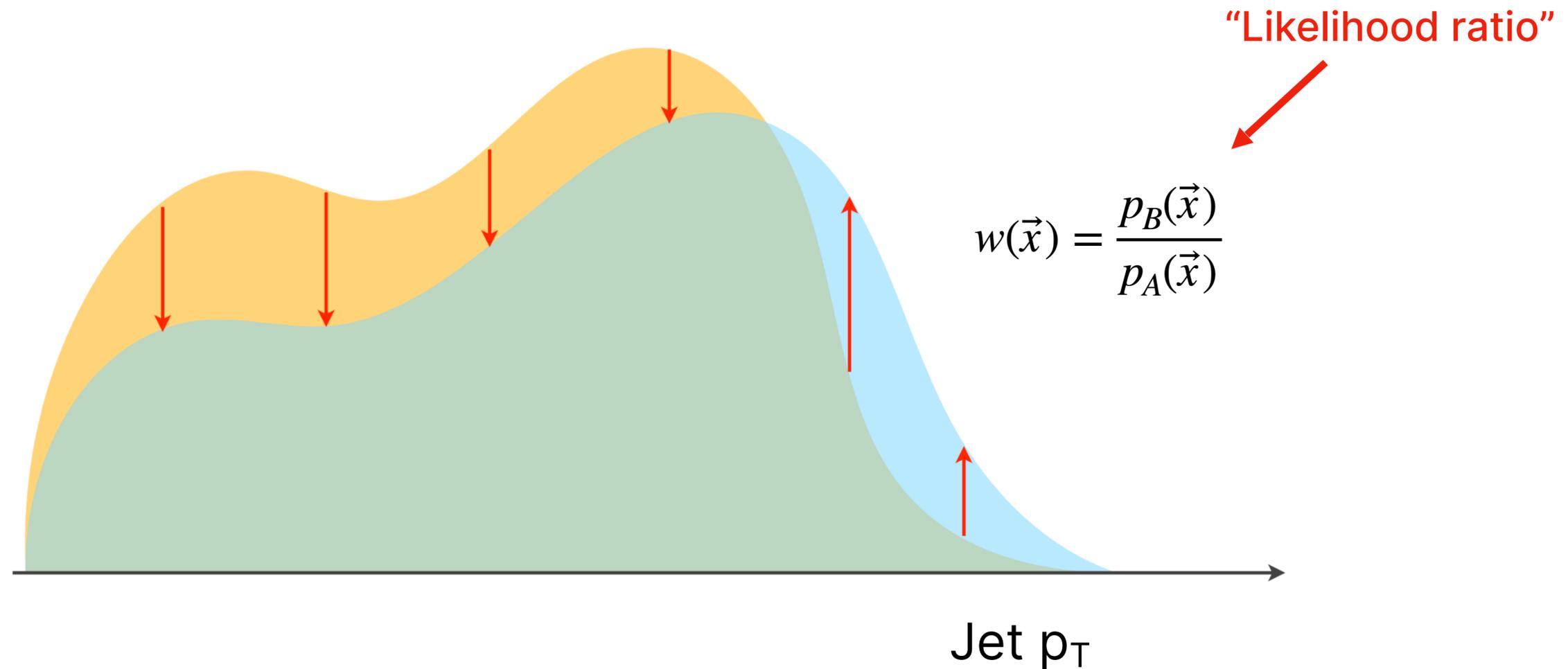
**Informed by universal physical laws** → multi-detector/disciplinary analyses

High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

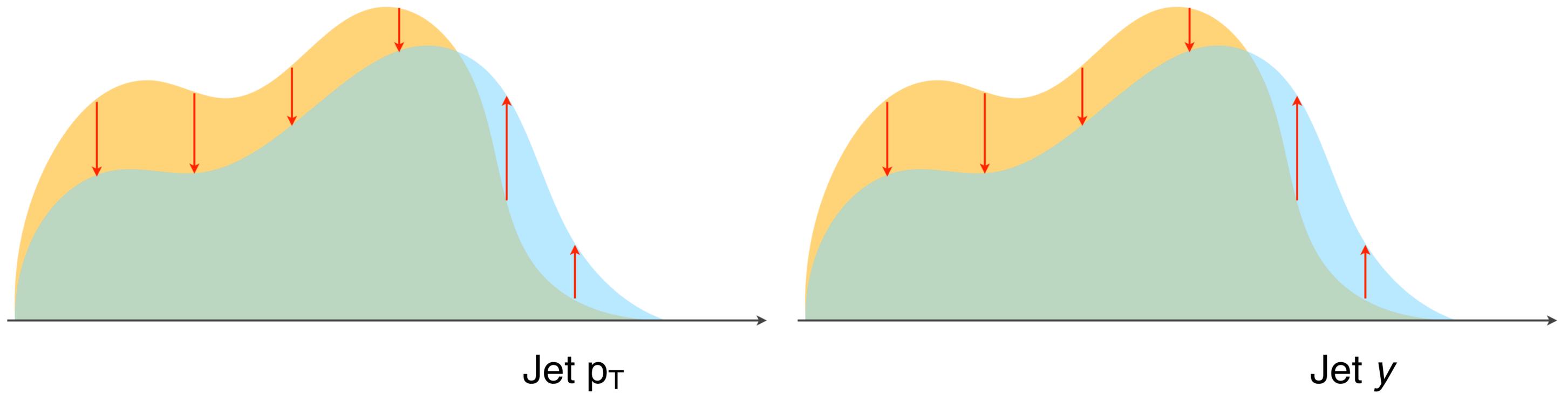


High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

Learns a high-dimensional reweighting function based on the ratio of the probability densities:



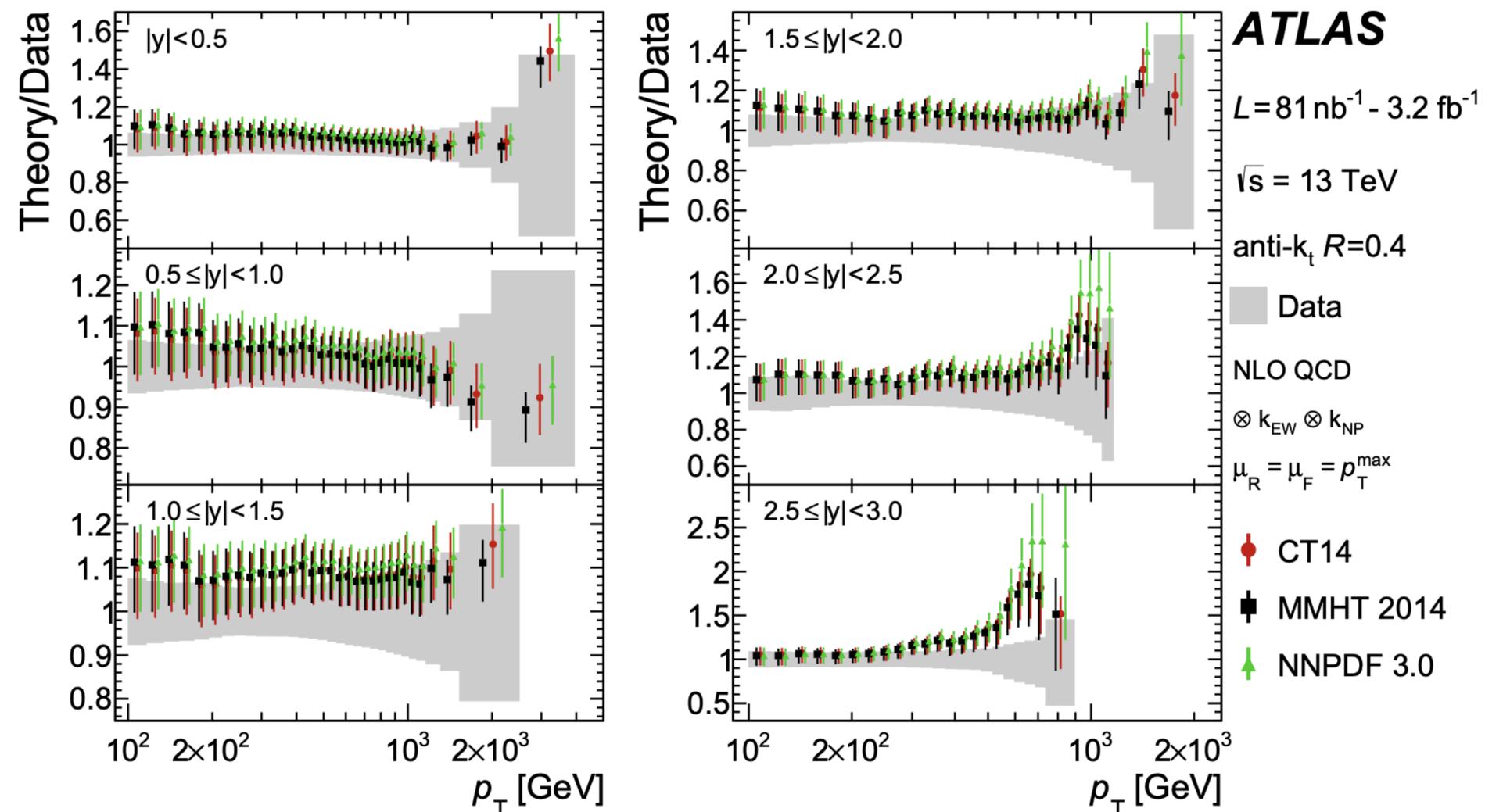
High-dimensional & unpinned unfolding allows us to publish more flexible datasets.



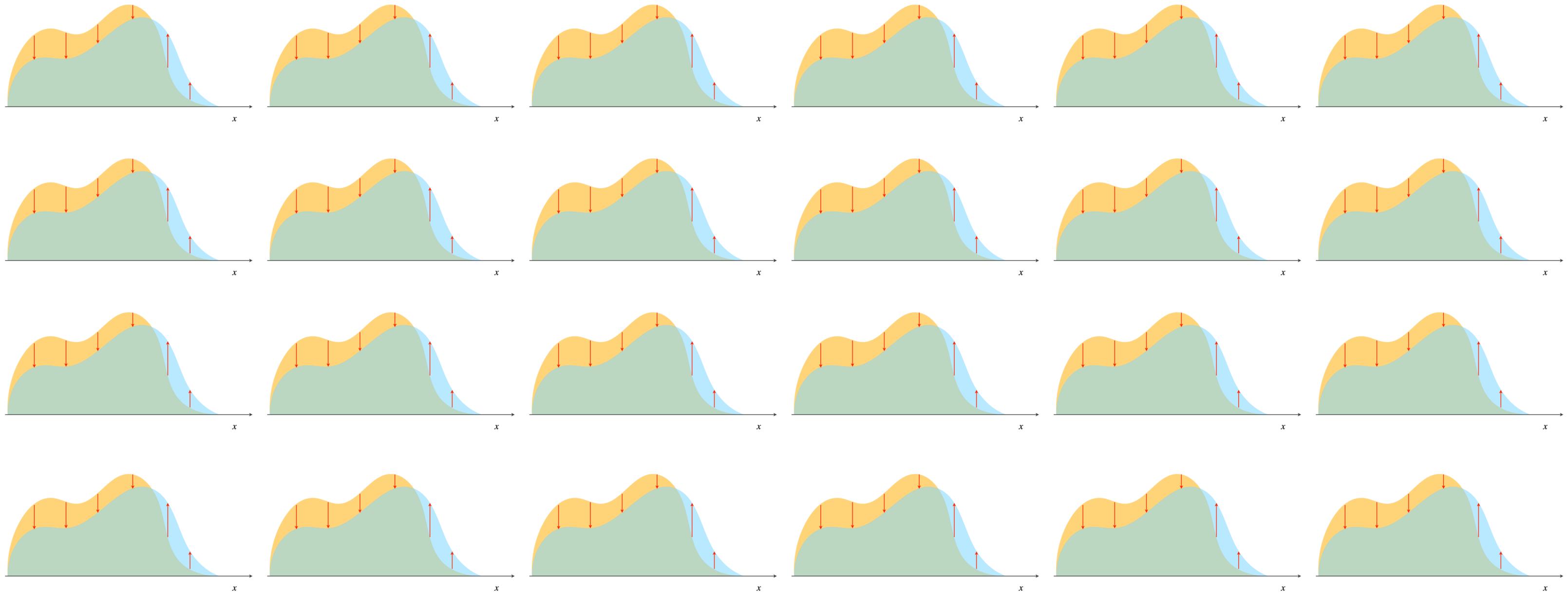
High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

Many LHC measurements are published as binned differential cross-sections at the particle level.

(So far from the LHC, for only one or two observables at a time).



High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

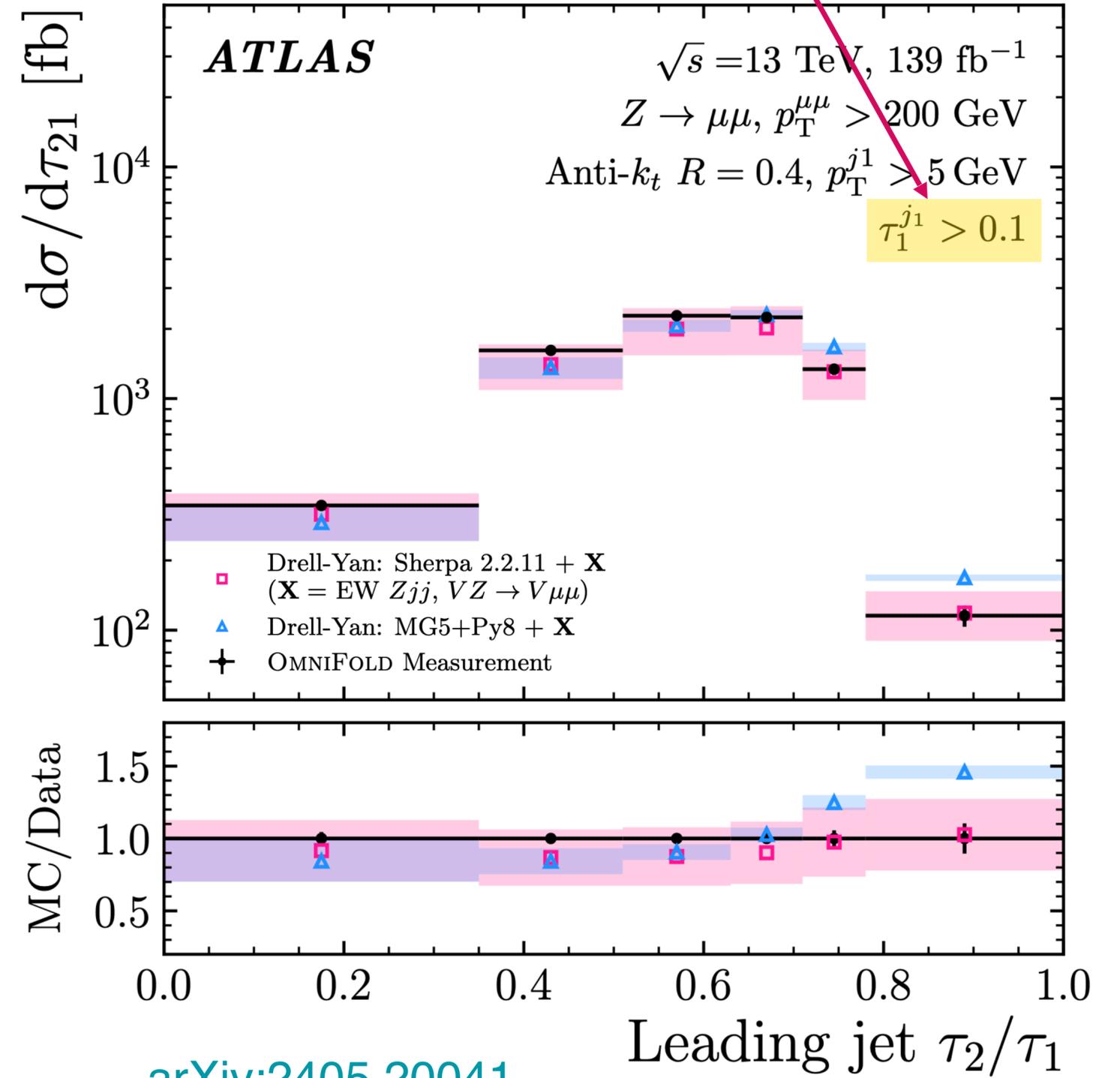
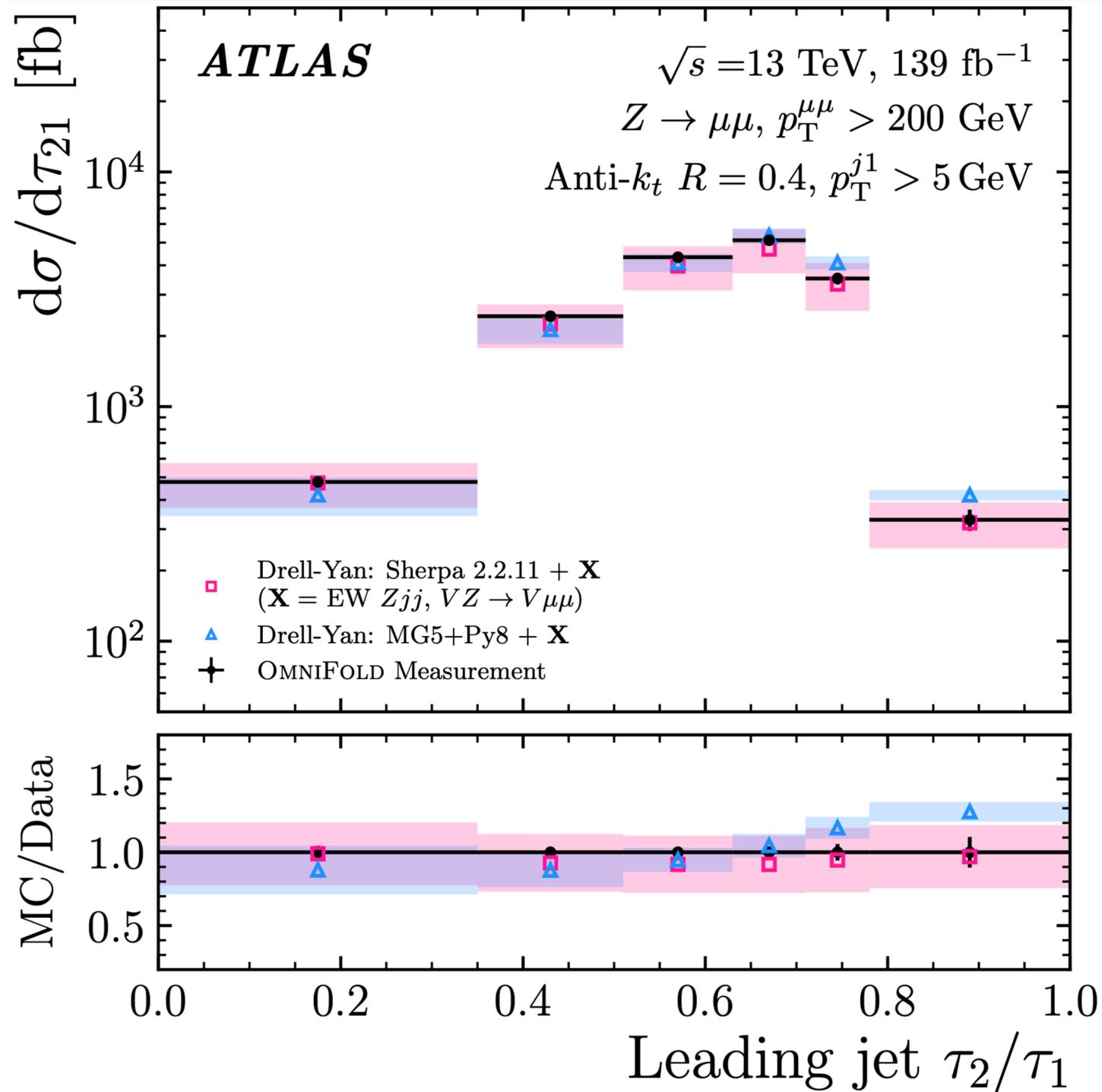




Let's construct some new observables...

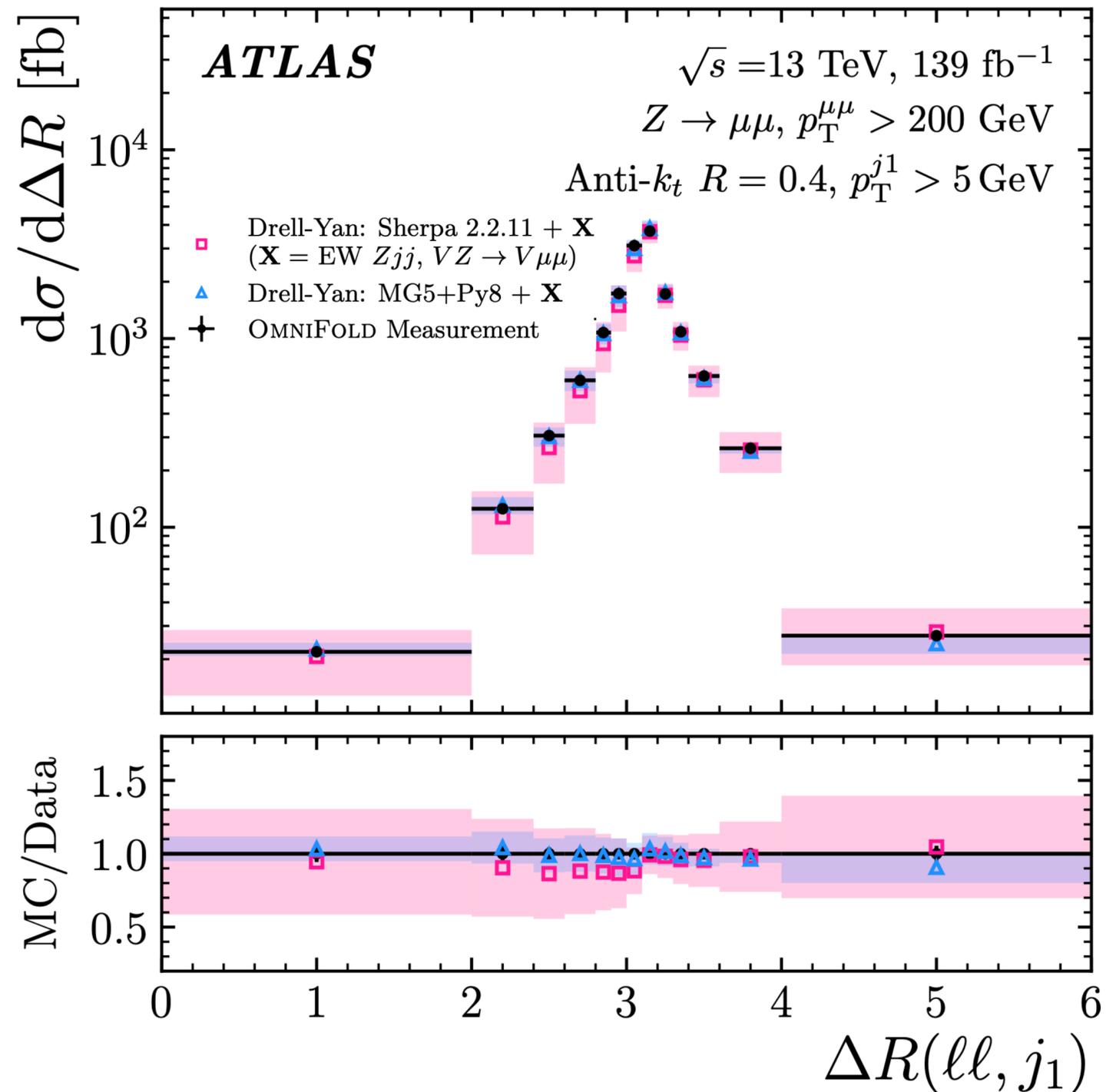
$$\tau_{21} = \tau_2 / \tau_1$$

IRC-safe phase space!

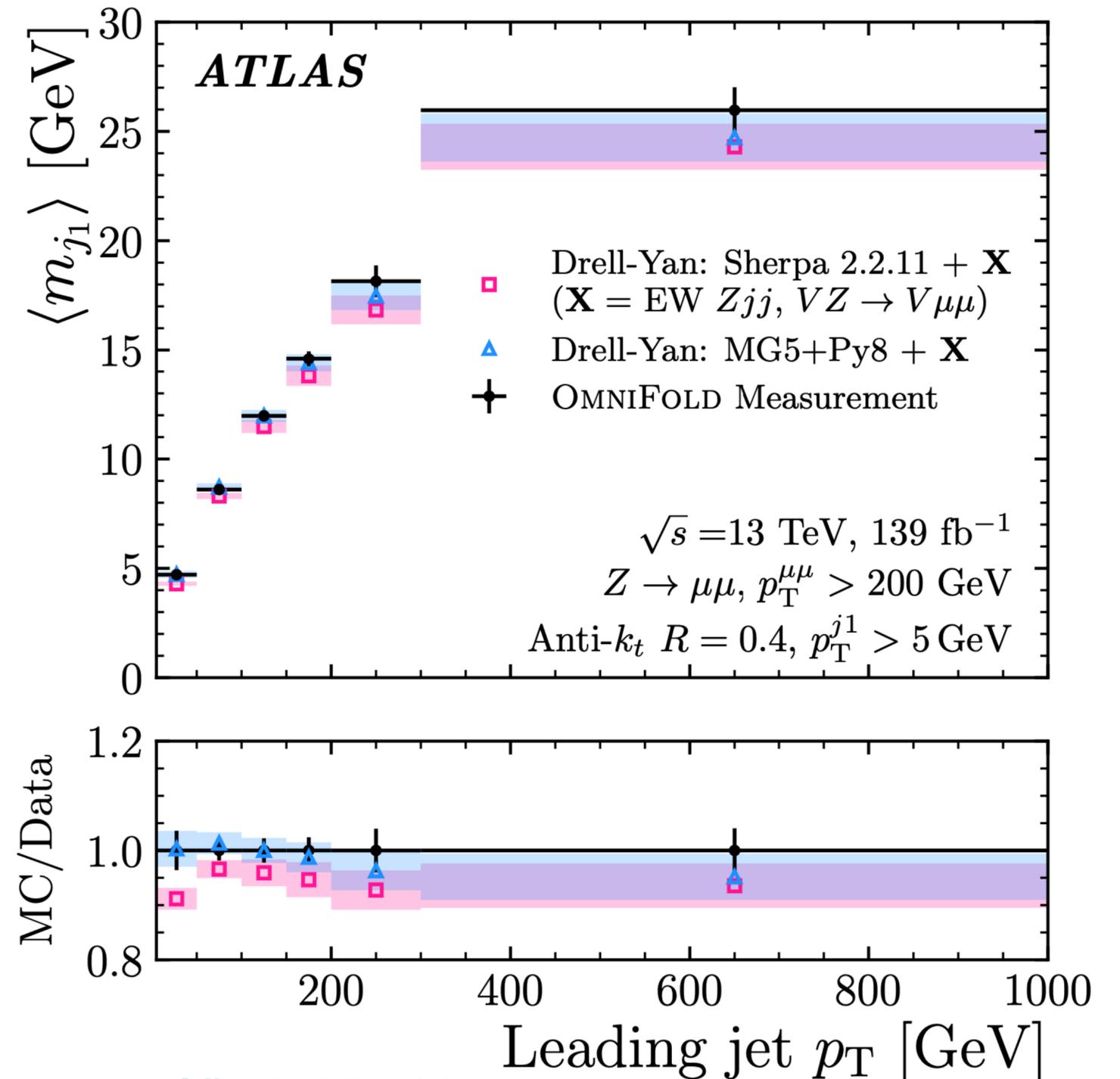


Let's construct some new observables...

$$\Delta R(\ell\ell, j_1)$$



Average  $m_{j_1}$  in bins of leading jet  $p_T$



High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

	pT_ll	pT_l1	pT_l2	eta_l1	...	weights_trackPtScale	weights_theoryPSjet	weights_theoryPSsoft
<b>0</b>	479.442780	288.466919	198.183929	-0.117443	...	0.003174	0.002844	0.003195
<b>1</b>	274.524994	166.120789	125.378044	0.313321	...	0.008168	0.008563	0.008236
<b>2</b>	462.713226	335.697479	133.157684	0.766387	...	0.001638	0.001724	0.001890
<b>3</b>	215.157608	189.518021	25.711994	1.083798	...	0.004669	0.004622	0.004648
<b>4</b>	222.458313	128.850159	108.589226	-0.635713	...	0.002102	0.002417	0.002129
...	...	...	...	...	...	...	...	...
<b>418009</b>	934.971924	738.464722	196.525192	0.102944	...	0.000069	0.000070	0.000061
<b>418010</b>	245.813461	166.847061	93.757919	1.308837	...	0.000193	0.000189	0.000203
<b>418011</b>	478.670349	378.737518	108.016479	-0.328871	...	0.001969	0.001813	0.001825
<b>418012</b>	278.586029	249.255356	43.581135	0.632484	...	0.003238	0.003101	0.003090
<b>418013</b>	244.505249	219.796280	40.357105	1.833223	...	0.000947	0.000968	0.000957

High-dimensional & unpinned unfolding allows us to publish more flexible datasets.

## A Practical Guide to Unbinned Unfolding

Florencia Canelli,<sup>1,\*</sup> Kyle Cormier,<sup>1,\*</sup> Andrew Cudd,<sup>2,†</sup> Dag Gillberg,<sup>3,‡</sup> Roger G. Huang,<sup>4,†</sup> Weijie Jin,<sup>1,\*</sup> Sookhyun Lee,<sup>5,§</sup> Vinicius Mikuni,<sup>4,6,¶</sup> Laura Miller,<sup>7</sup> Benjamin Nachman,<sup>4,8,9,\*\*</sup> Jingjing Pan,<sup>4,10,\*\*</sup> Tanmay Pani,<sup>11,††</sup> Mariel Pettee,<sup>4,‡</sup> Youqi Song,<sup>10,††</sup> and Fernando Torales Acosta<sup>4,¶</sup>

Experiment	Paper Link	Dimensions	Final State	Momentum Selection
ATLAS [23]	<a href="#">2405.20041</a>	24	Z+jets	$p_T^{\ell\ell} > 200 \text{ GeV}$
ATLAS [24]	<a href="#">2502.02062</a>	6	Dijets	$p_T^{j1} > 240 \text{ GeV} \ \& \ p_T^{j1} < 1.5 p_T^{j2}$
CMS [25]	<a href="#">2505.17850</a>	8	Minimum bias	$> 2$ charged particles with $p_T > 0.5 \text{ GeV}$
H1 [26]	<a href="#">2108.12376</a>	8*	High $Q^2$ DIS	$Q^2 > 150 \text{ GeV}^2$
H1 [27]	<a href="#">2303.13620</a>	10	High $Q^2$ DIS	$Q^2 > 150 \text{ GeV}^2$
H1 [28]	<a href="#">2412.14092</a>	8*	High $Q^2$ DIS	$Q^2 > 150 \text{ GeV}^2$
H1 [29]	<a href="#">H1prelim-25-031</a>	Variable	High $Q^2$ DIS	$Q^2 > 150 \text{ GeV}^2$
LHCb [30]	<a href="#">2208.11691</a>	4	Z+hadrons in jets	$20 < p_T^j < 100 \text{ GeV}$ and $p_T^h > 0.25 \text{ GeV}$
STAR [31]	<a href="#">2307.07718</a>	6	Jets	$20 < p_T^j < 50 \text{ GeV}$
STAR [32]	<a href="#">2403.13921</a>	7	Jets in heavy ions	$20 < p_T^j < 45 \text{ GeV}$
T2K [33]	<a href="#">2504.06857</a>	6	Muon + Proton	$p^p > 450 \text{ MeV}$ for single transverse variables

Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

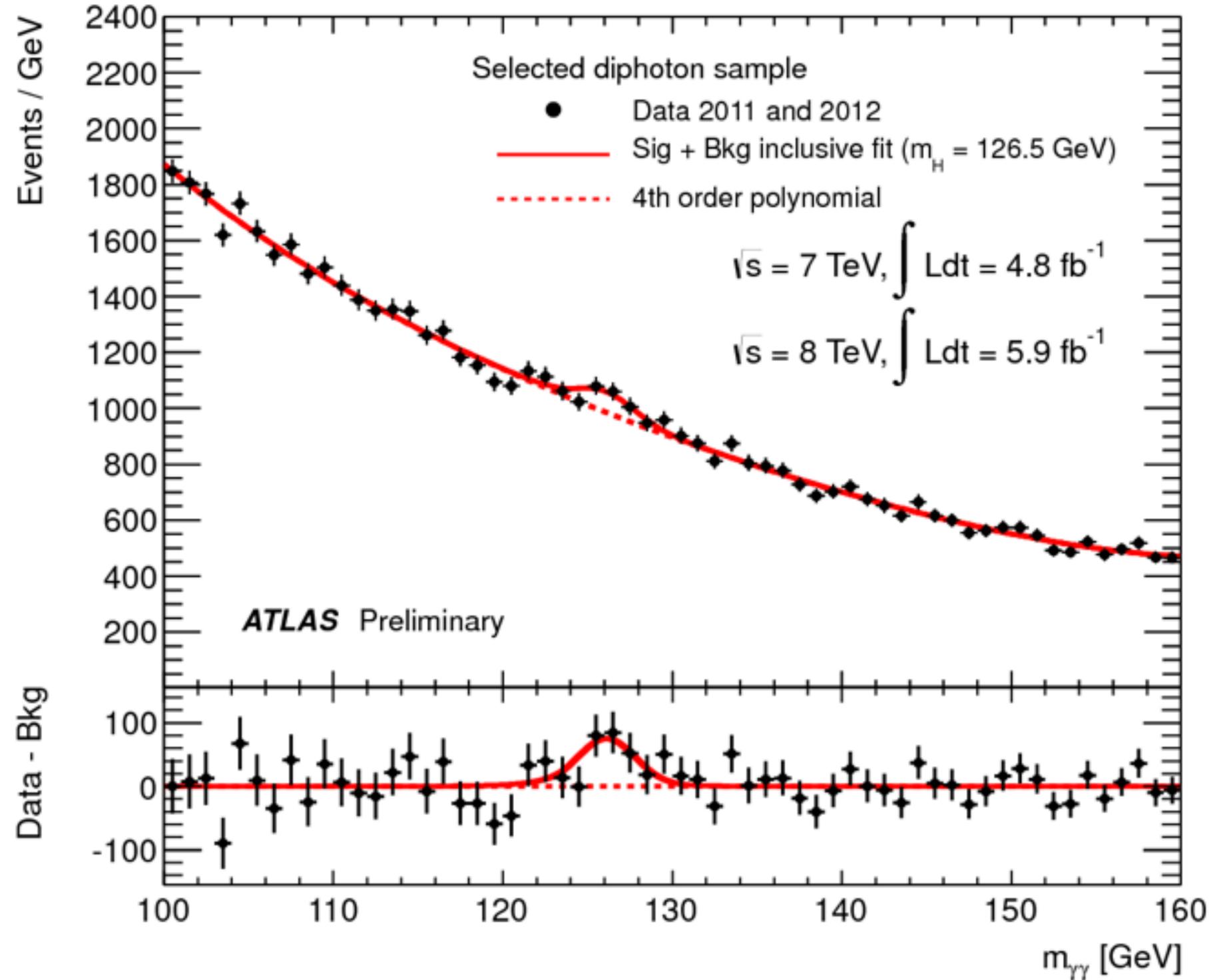
**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

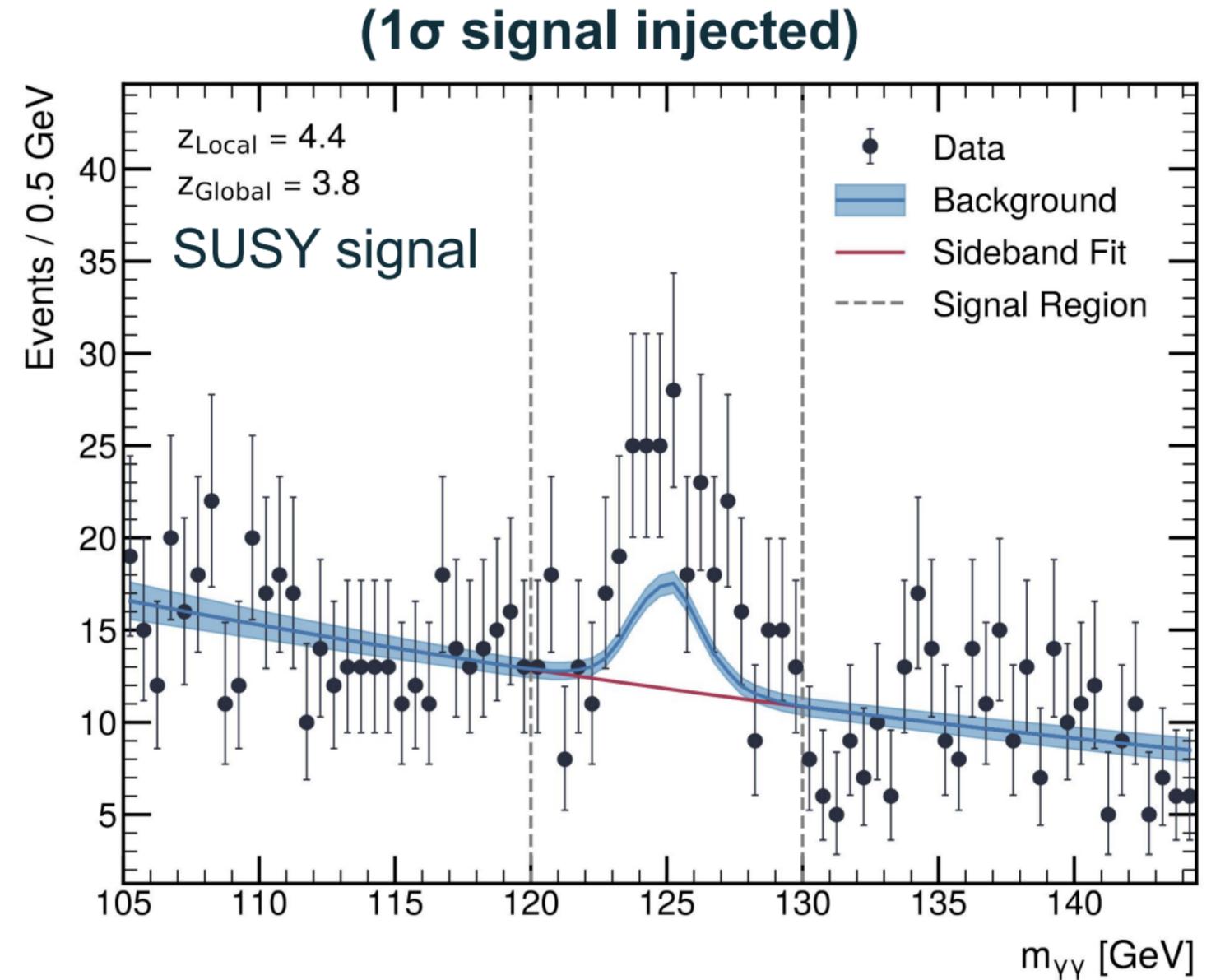
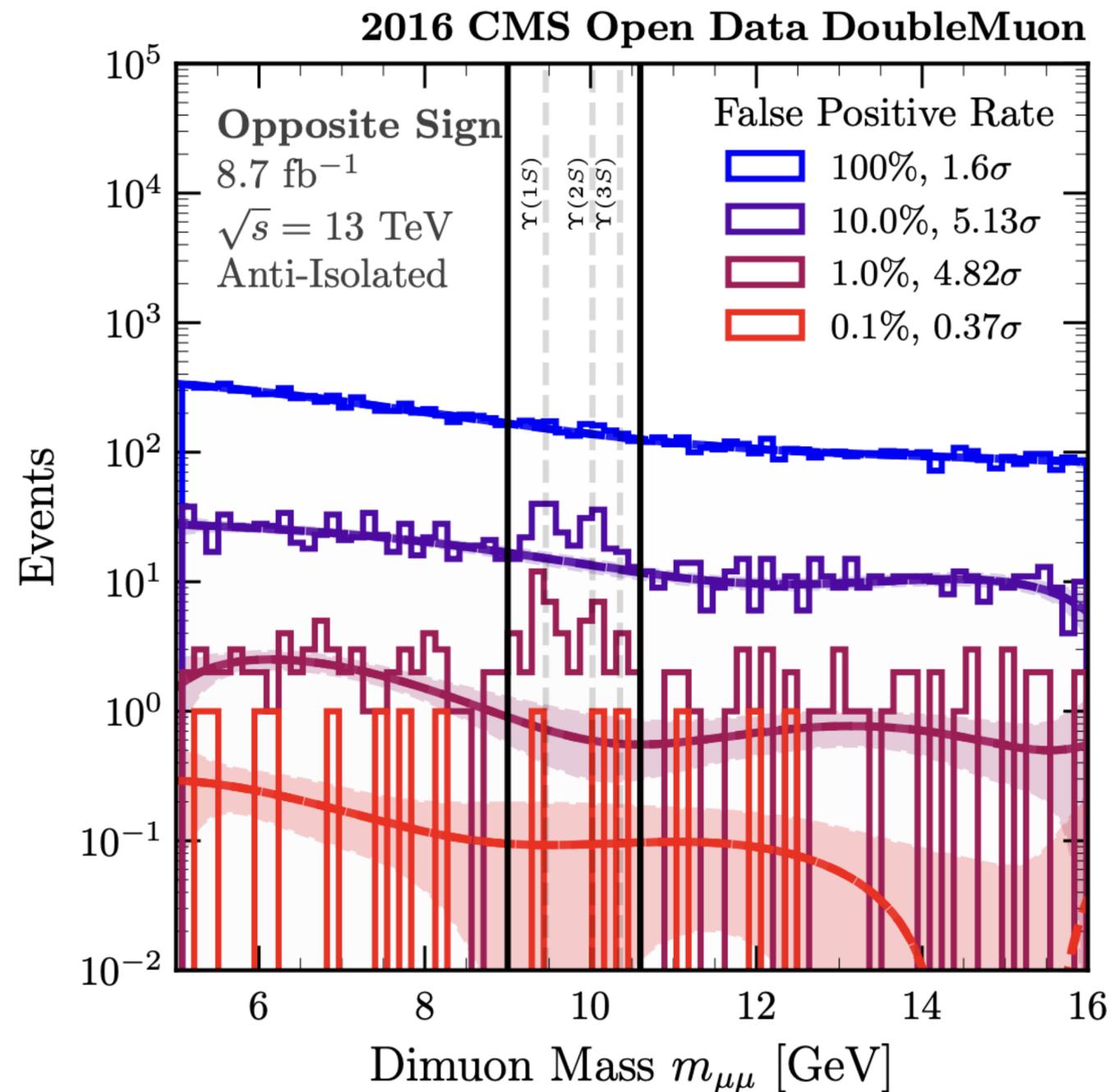
**Highly multi-modal** → unify inputs across an entire detector

**Informed by universal physical laws** → multi-detector/disciplinary analyses

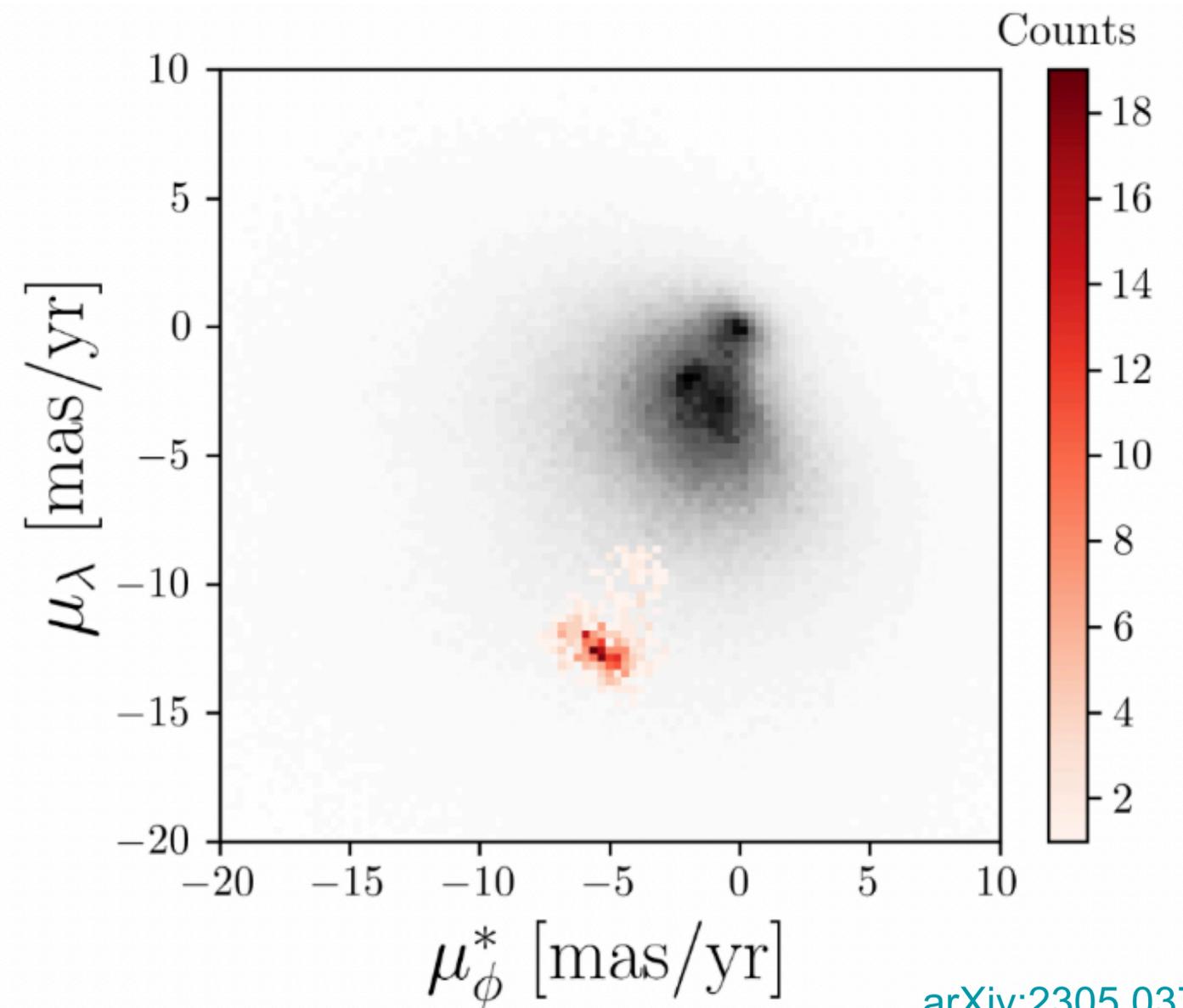
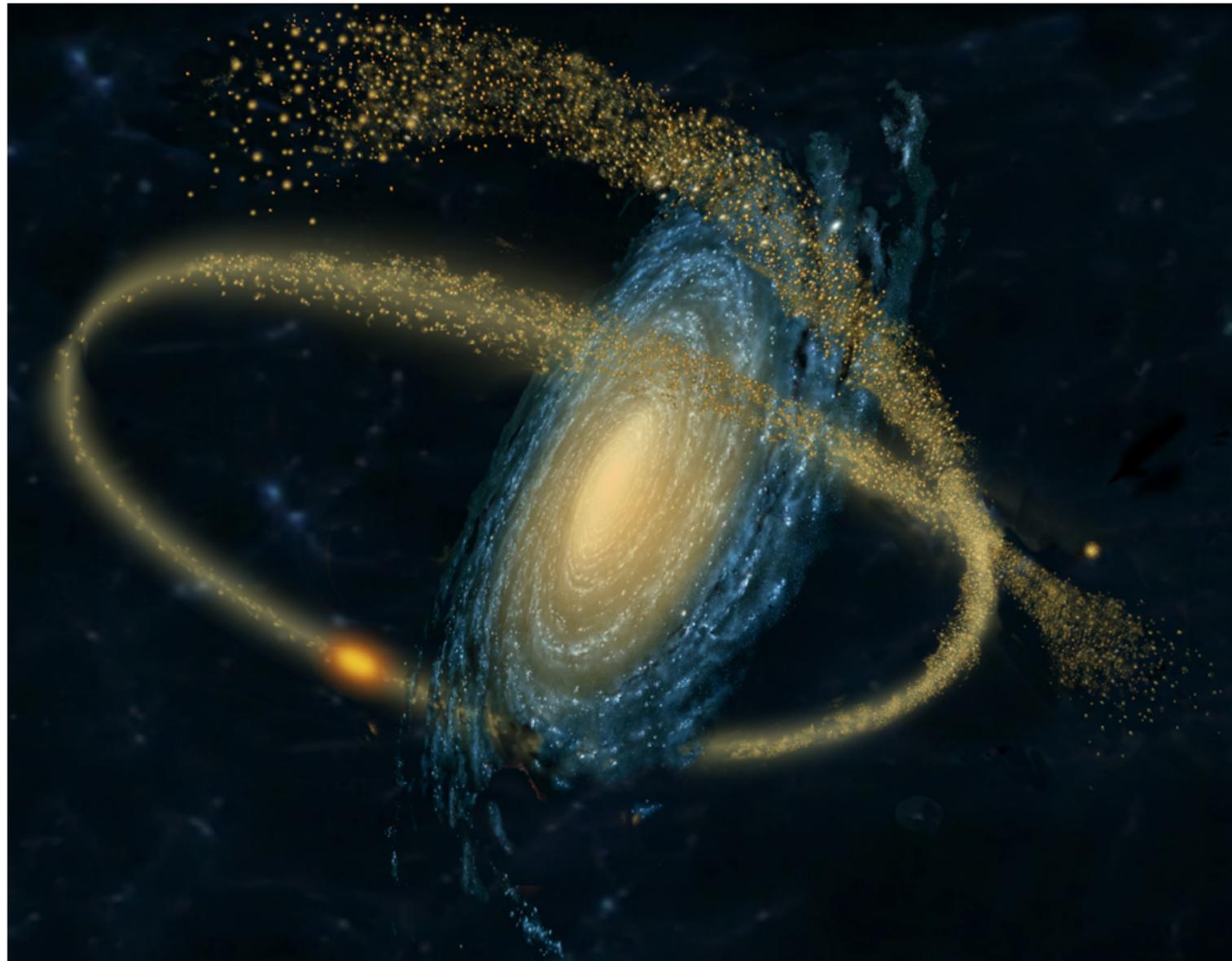
Anomaly detection in particle physics targets **group anomalies**, not single-event outliers.



Anomaly detection in particle physics targets **group anomalies**, not single-event outliers.

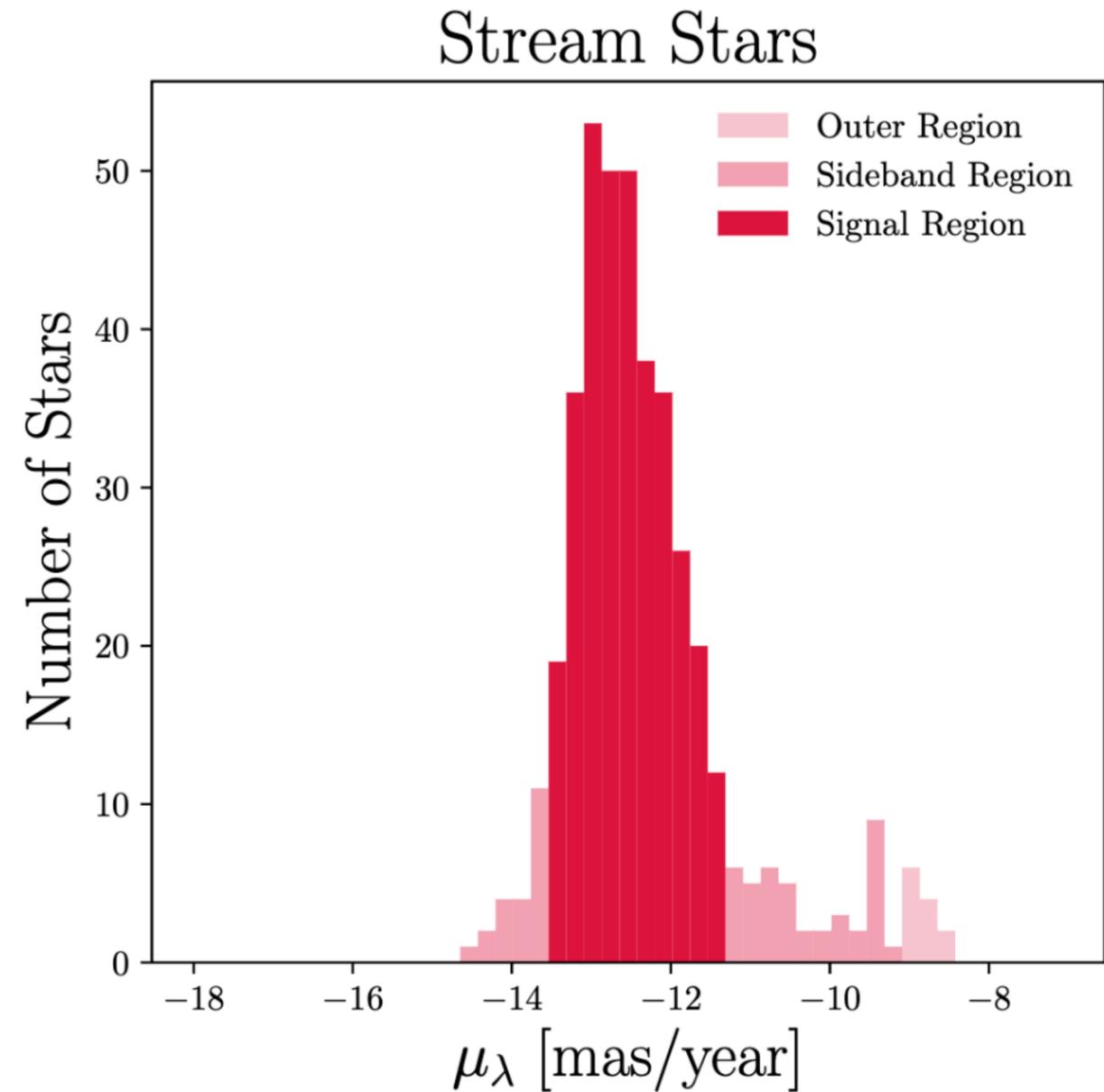
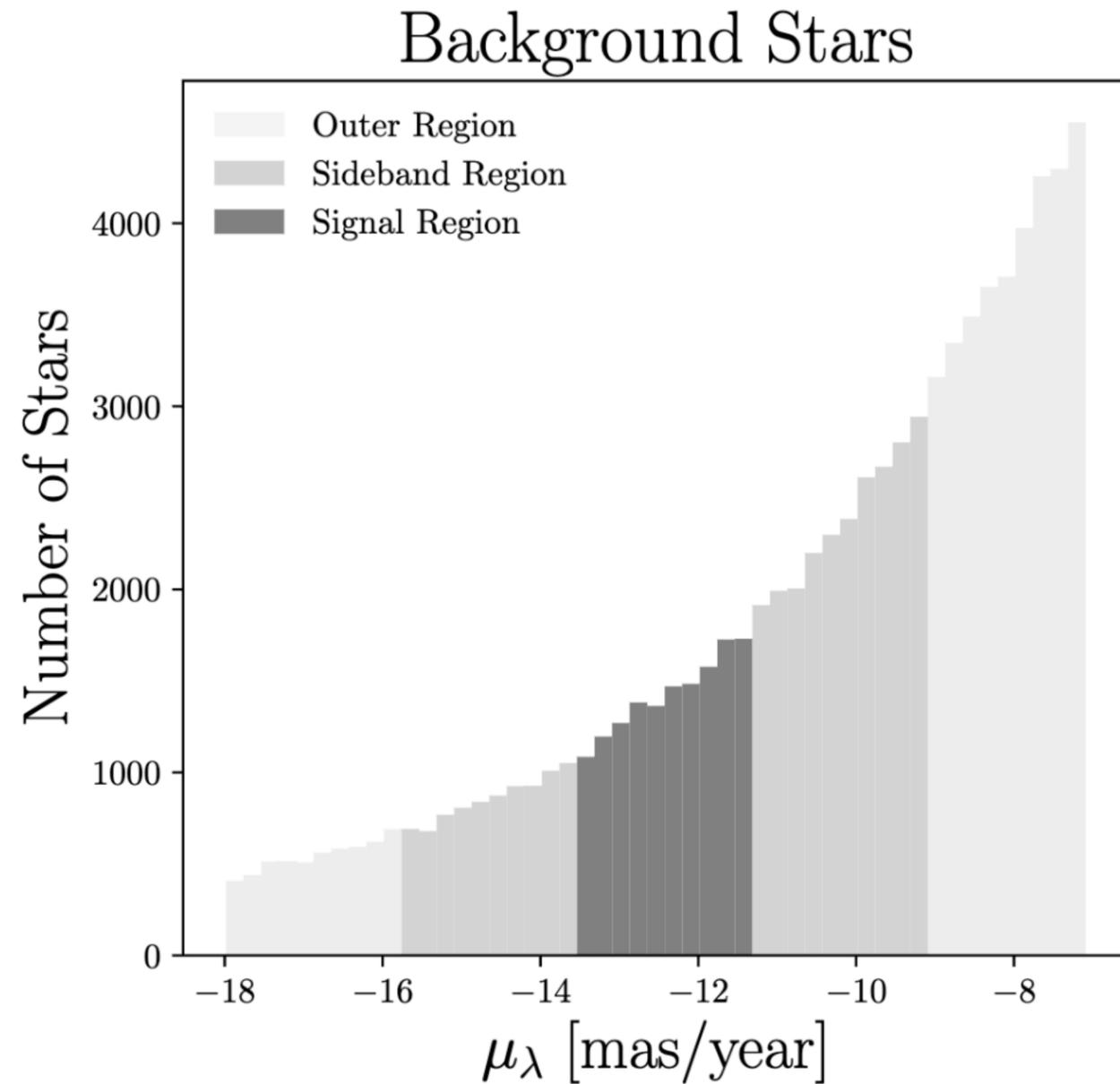


Anomaly detection in particle physics targets **group anomalies**, not single-event outliers.  
(and group anomalies show up in astrophysics, too!)

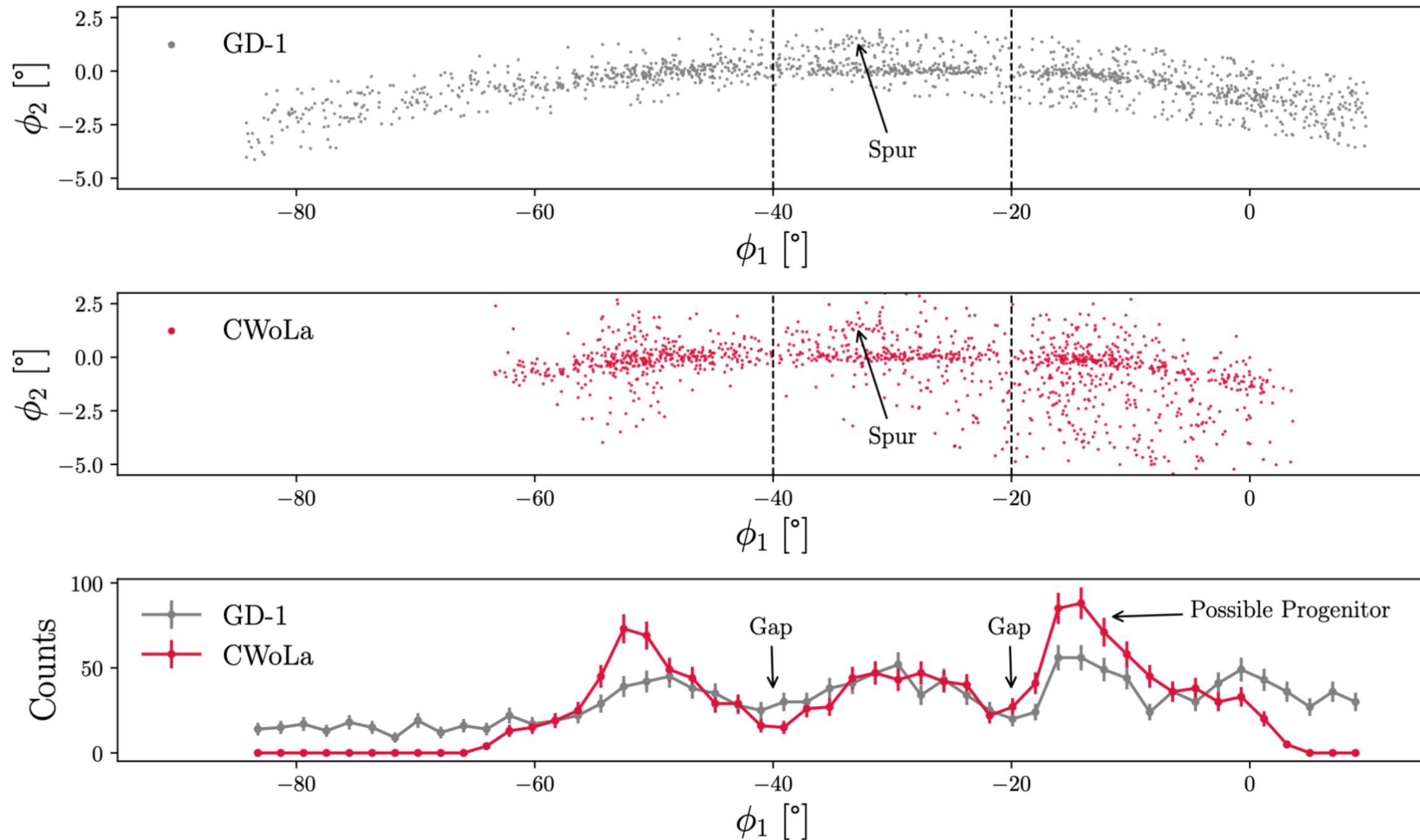


[arXiv:2305.03761](https://arxiv.org/abs/2305.03761)

Anomaly detection in particle physics targets **group anomalies**, not single-event outliers.  
(and group anomalies show up in astrophysics, too!)



Anomaly detection in particle physics targets **group anomalies**, not single-event outliers.  
(and group anomalies show up in astrophysics, too!)



Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

**Highly multi-modal** → unify inputs across an entire detector

**Informed by universal physical laws** → multi-detector/disciplinary analyses

A **modality** is a specific way of perceiving the world – it's a concept deeply rooted in our own five senses.

(In ML, common modalities include text, audio, video, etc.)

AI industry has multimodal models, but the modalities are relatively few and well-defined:



# Multimodal Intelligence: Next Token Prediction & Beyond

**ICLR 2026 Workshop**

**Location:** Rio de Janeiro, Brazil

**Date:** 26 or 27 April 2026

[SUBMIT YOUR PAPER](#)

[CALL FOR REVIEWERS](#)

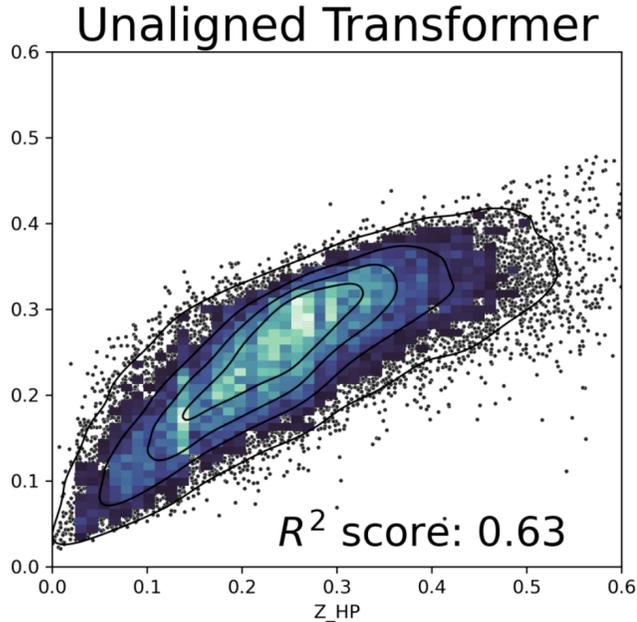
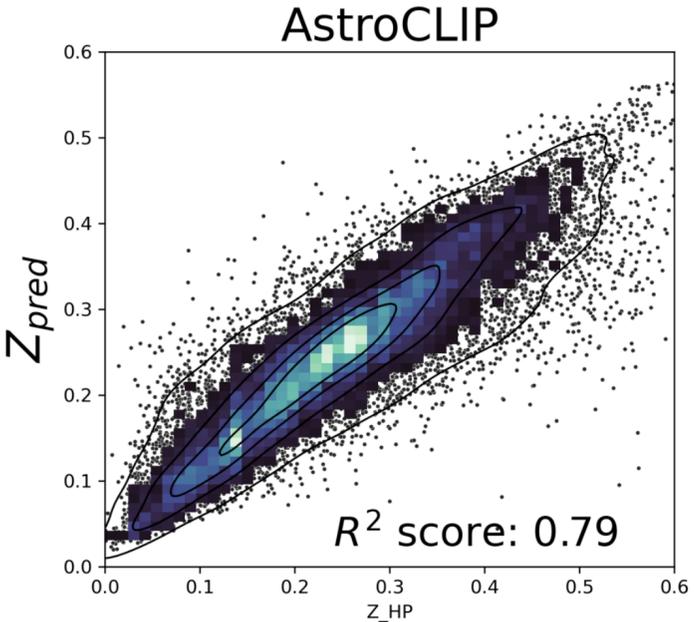
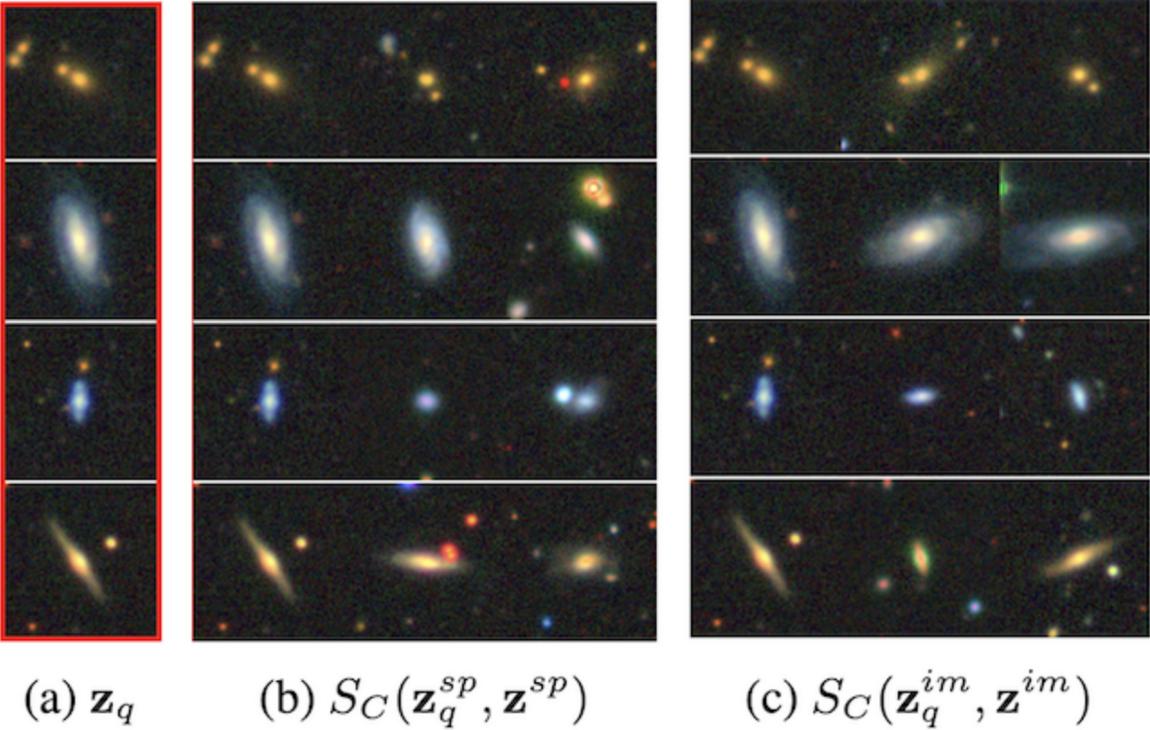
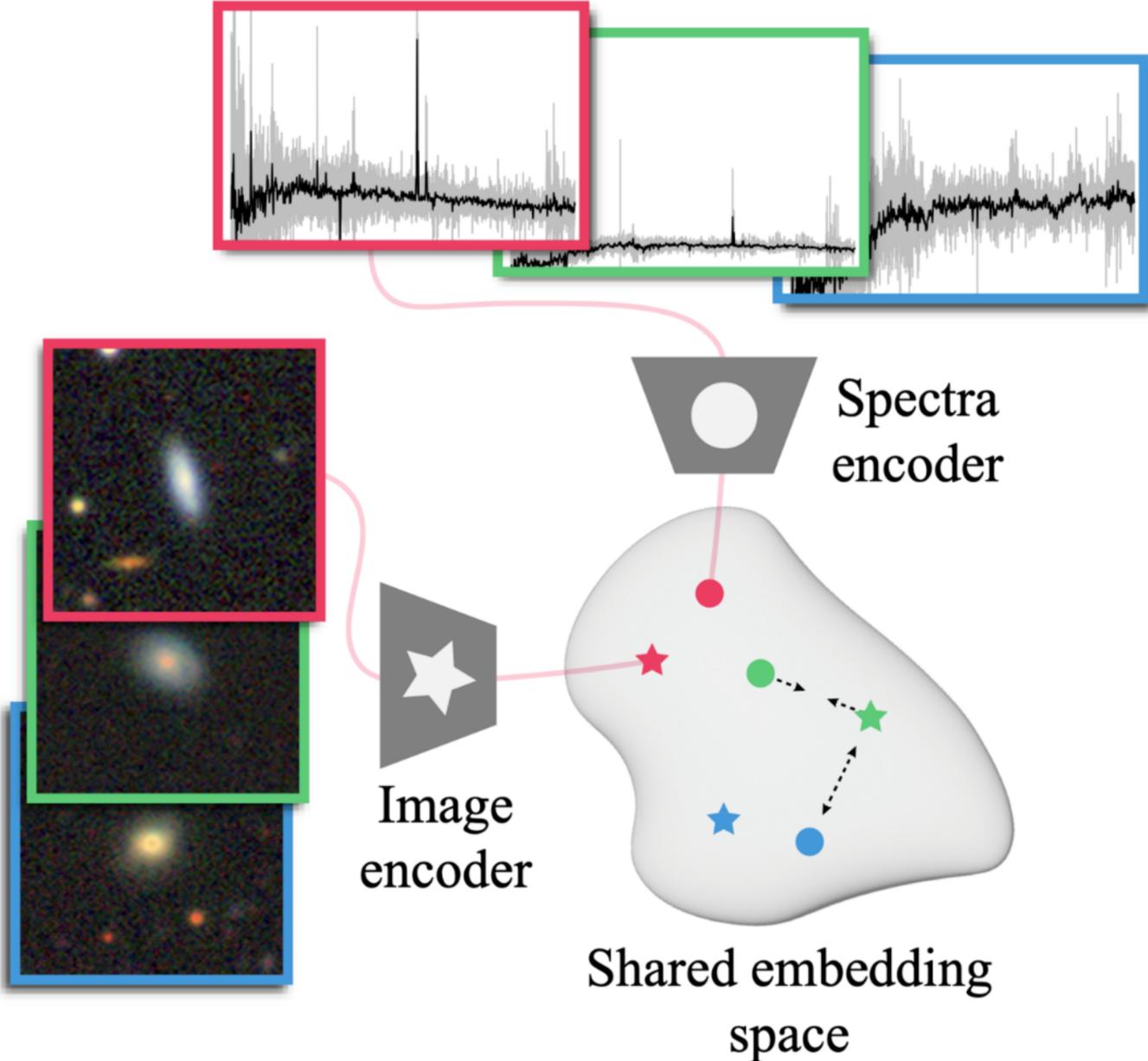
# Multimodal models see promising performance in industry and in science.

Significant evidence suggests that including multiple views/ modalities in a single model can improve performance on downstream tasks relative to unimodal models:

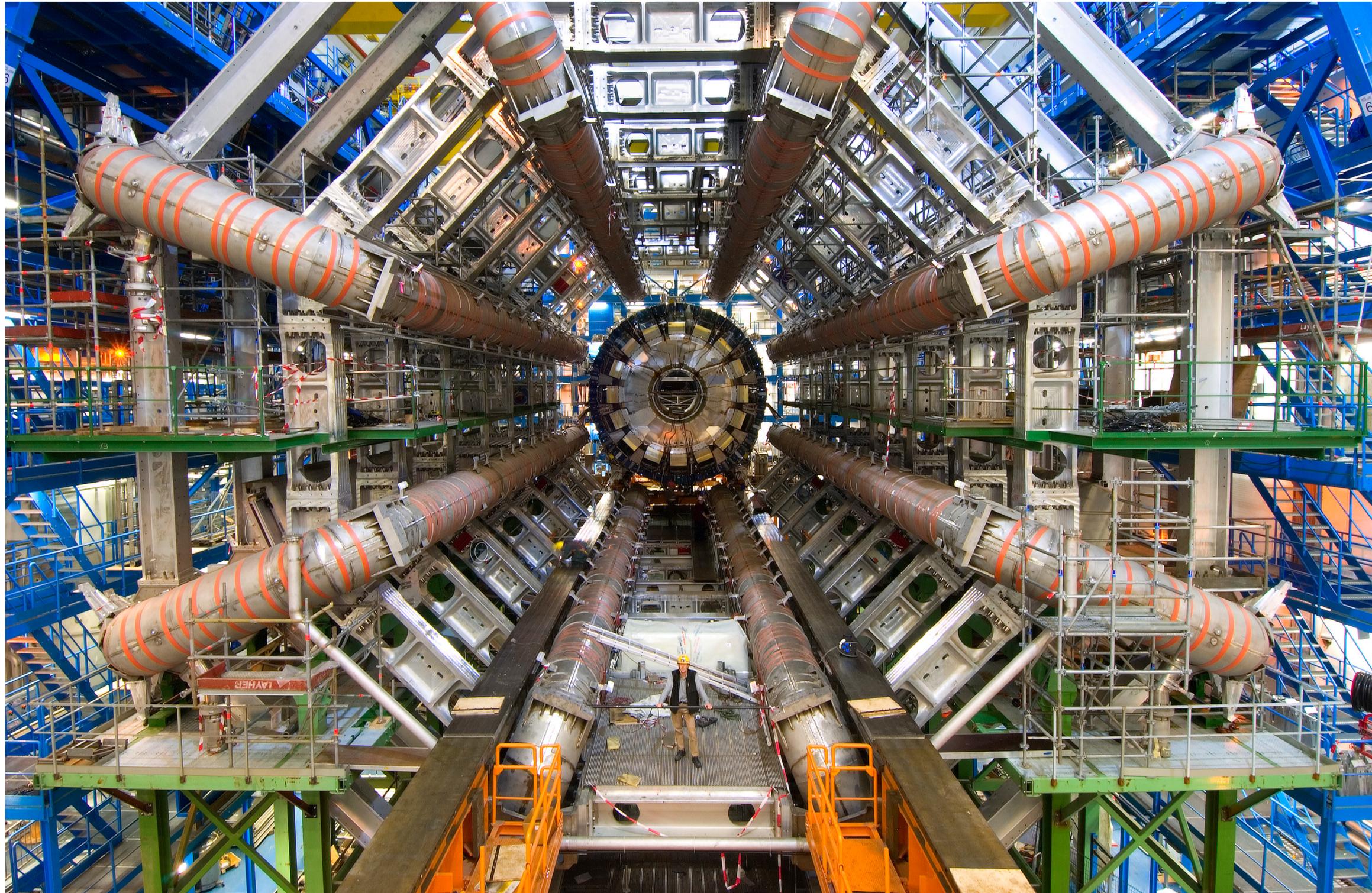
	Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	$\Delta$ Score
ImageNet		76.2	76.2	0%
ImageNetV2		64.3	70.1	+5.8%
ImageNet-R		37.7	88.9	+51.2%
ObjectNet		32.6	72.3	+39.7%
ImageNet Sketch		25.2	60.2	+35.0%
ImageNet-A		2.7	77.1	+74.4%

(CLIP, 2021)

Multimodal models see promising performance in industry and in science.



But what happens when we try to encode 10, 20, 50+ modalities at once?

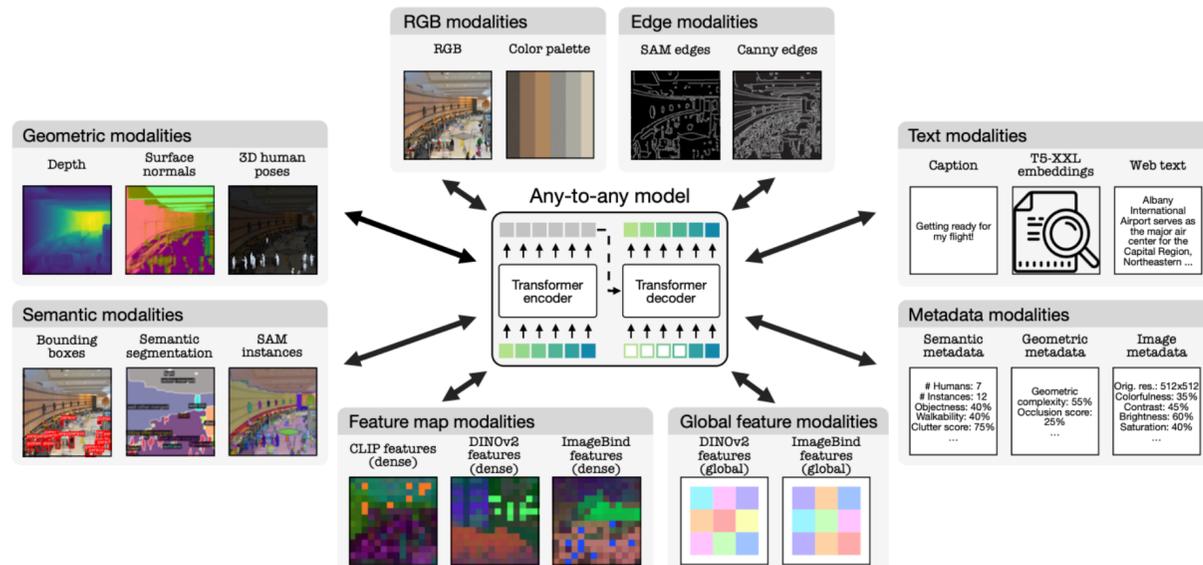


But what happens when we try to encode 10, 20, 50+ modalities at once?

Polymathic

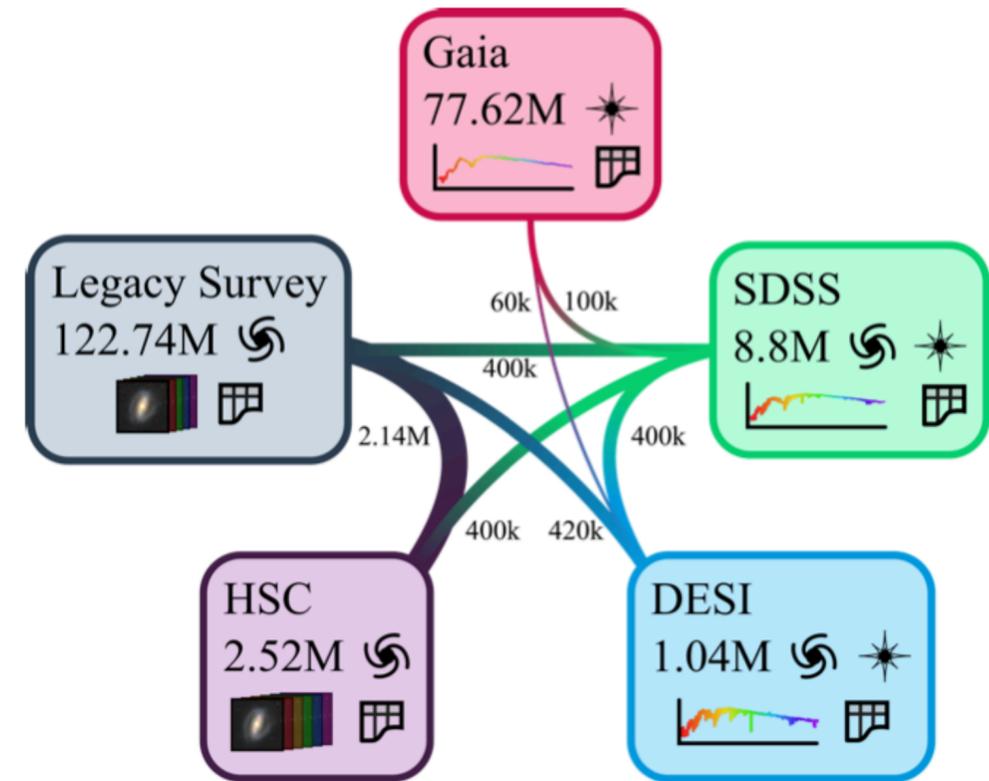
## 4M-21: An Any-to-Any Vision Model for Tens of Tasks and Modalities

Roman Bachmann<sup>1†\*</sup> Oğuzhan Fatih Kar<sup>1\*</sup> David Mizrahi<sup>2†\*</sup> Ali Garjani<sup>1</sup>  
 Mingfei Gao<sup>2</sup> David Griffiths<sup>2</sup> Jiaming Hu<sup>2</sup> Afshin Dehghan<sup>2</sup> Amir Zamir<sup>1</sup>  
<sup>1</sup>Swiss Federal Institute of Technology Lausanne (EPFL) <sup>2</sup>Apple  
<https://4m.epfl.ch>

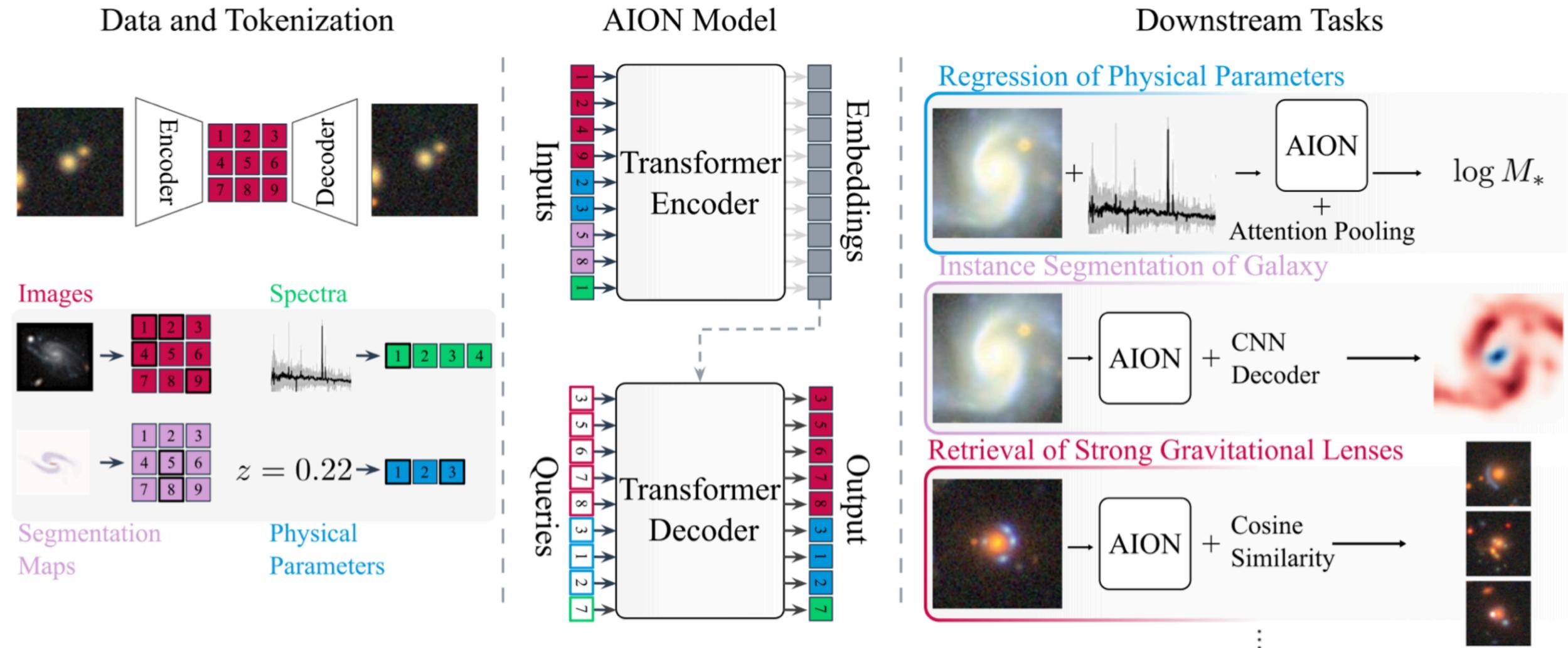


## AION-1: Omnimodal Foundation Model for Astronomical Sciences

Liam Parker<sup>\*,1,2,3,4</sup>, Francois Lanusse<sup>\*,5,2</sup>, Jeff Shen<sup>\*,6</sup>, Ollie Liu<sup>7</sup>, Tom Hehir<sup>8</sup>, Leopoldo Sarra<sup>3</sup>, Lucas Meyer<sup>3</sup>, Micah Bowles<sup>9</sup>, Sebastian Wagner-Carena<sup>2,3</sup>, Helen Qu<sup>2</sup>, Siavash Golkar<sup>2,3</sup>, Alberto Bietti<sup>2</sup>, Hatim Bourfoune<sup>10</sup>, Nathan Cassereau<sup>10</sup>, Pierre Cornette<sup>10</sup>, Keiya Hirashima<sup>2,11</sup>, Geraud Krawezik<sup>2</sup>, Ruben Ohana<sup>2</sup>, Nicholas Lourie<sup>3</sup>, Michael McCabe<sup>2,3</sup>, Rudy Morel<sup>2</sup>, Payel Mukhopadhyay<sup>1,8</sup>, Mariel Pettee<sup>12</sup>, Bruno Regaldo-Saint Blancard<sup>2</sup>, Kyunghyun Cho<sup>3</sup>, Miles Cranmer<sup>8</sup>, Shirley Ho<sup>2,3,6</sup>



But what happens when we try to encode 10, 20, 50+ modalities at once?



**Figure 1:** AION-1 integrates 39 different data modalities — multiband images, optical spectra, and various properties and measurements — into a single model usable for a wide range of downstream applications. It implements a two-step process: first, bespoke tokenization strategies that homogenize the diverse scientific data, followed by multimodal masked modeling that learns how different observations relate, inducing a deep understanding of the underlying physical objects. Astronomers can then leverage AION-1’s rich astrophysical understanding for a variety of downstream tasks.

Physics data has some qualities that are relatively under-explored in mainstream ML.

**Deep theoretical grounding** → exploit known symmetries & simulations

**Diverse data formats** → custom vs. generic embeddings

**Detector-specific biases** → publish data without detector effects

**Group anomalies** → statistical anomaly detection

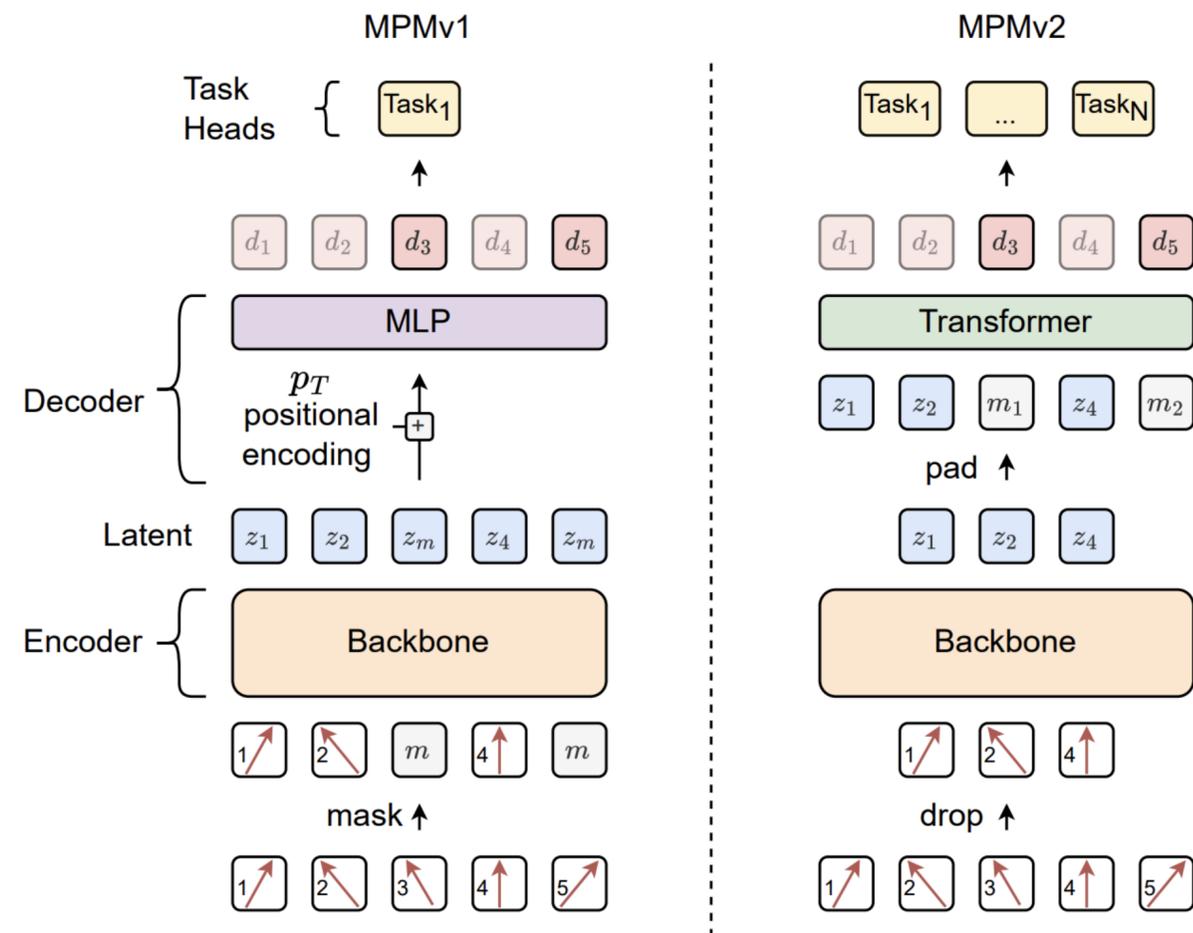
**Highly multi-modal** → unify inputs across an entire detector

**Informed by universal physical laws** → multi-detector/disciplinary analyses

We've started to see some domain-specific foundation models in physics...

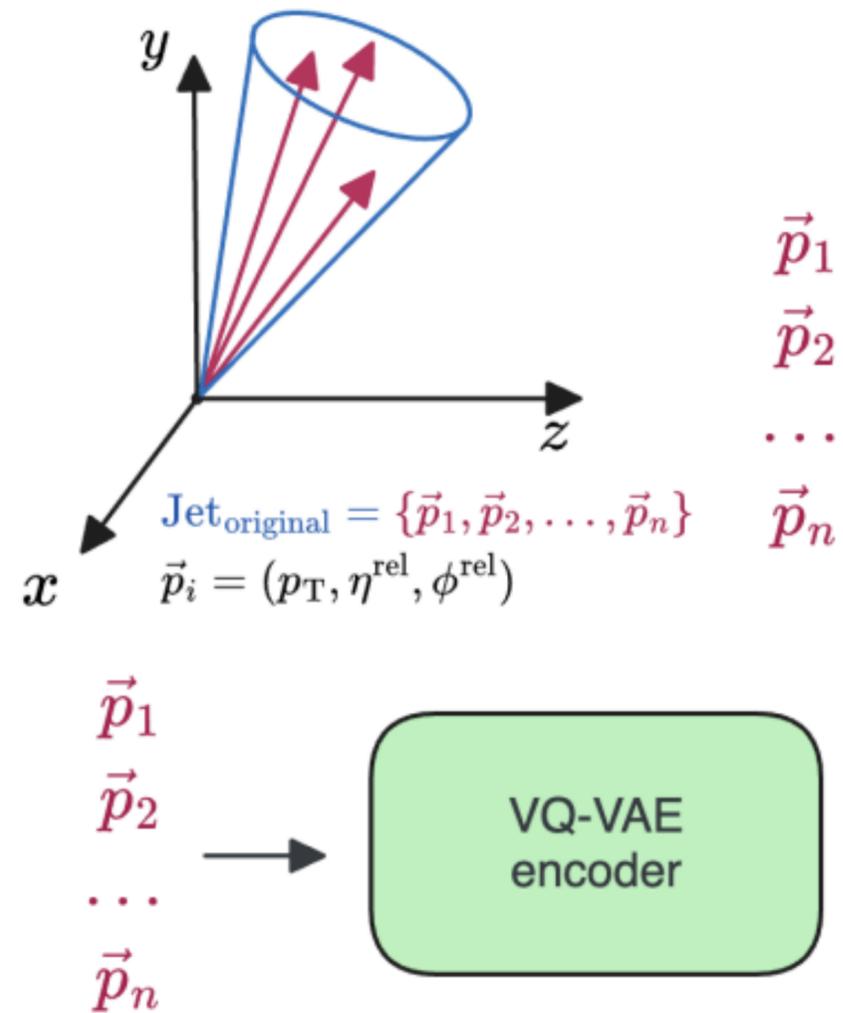
### Masked Particle Modeling (MPM)

[arXiv:2401.13537](https://arxiv.org/abs/2401.13537) & [arXiv:2409.12589](https://arxiv.org/abs/2409.12589)



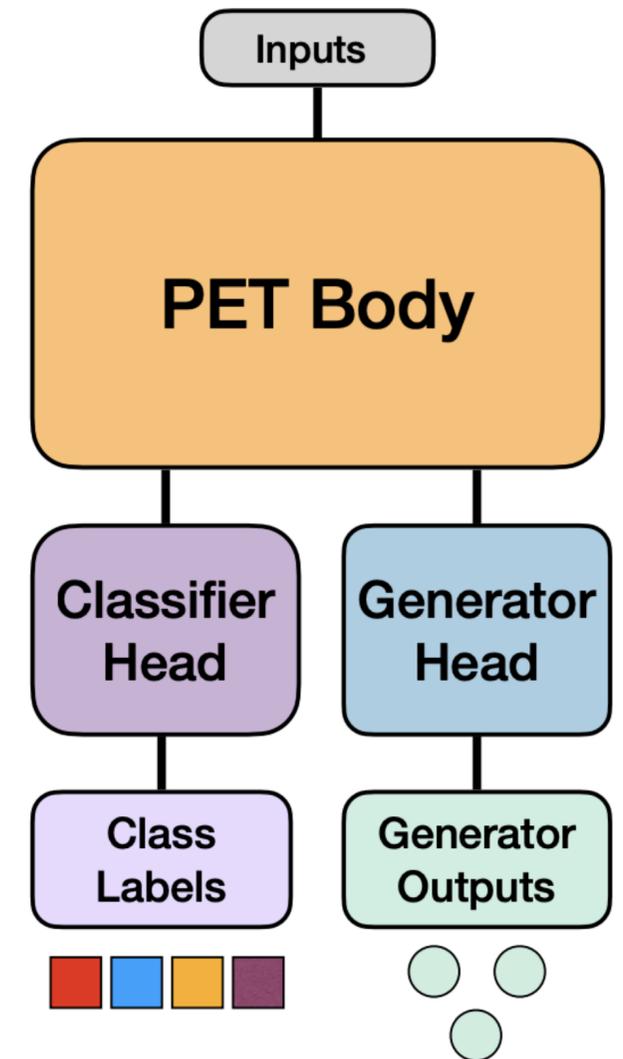
### Omni-Jet $\alpha$

[arXiv:2403.05618](https://arxiv.org/abs/2403.05618)



### OmniLearn

[arXiv:2404.16091](https://arxiv.org/abs/2404.16091)  
[arXiv:2502.14652](https://arxiv.org/abs/2502.14652)



...and it seems like they can be fine-tuned to transfer surprisingly well out-of-domain:

## OMNICOSMOS: Transferring Particle Physics Knowledge Across the Cosmos

Vinicius Mikuni,<sup>1,\*</sup> Ibrahim Elsharkawy,<sup>2,3,†</sup> and Benjamin Nachman<sup>4,5,‡</sup>

<sup>1</sup>*Nagoya University, Kobayashi-Maskawa Institute, Aichi 464-8602, Japan*

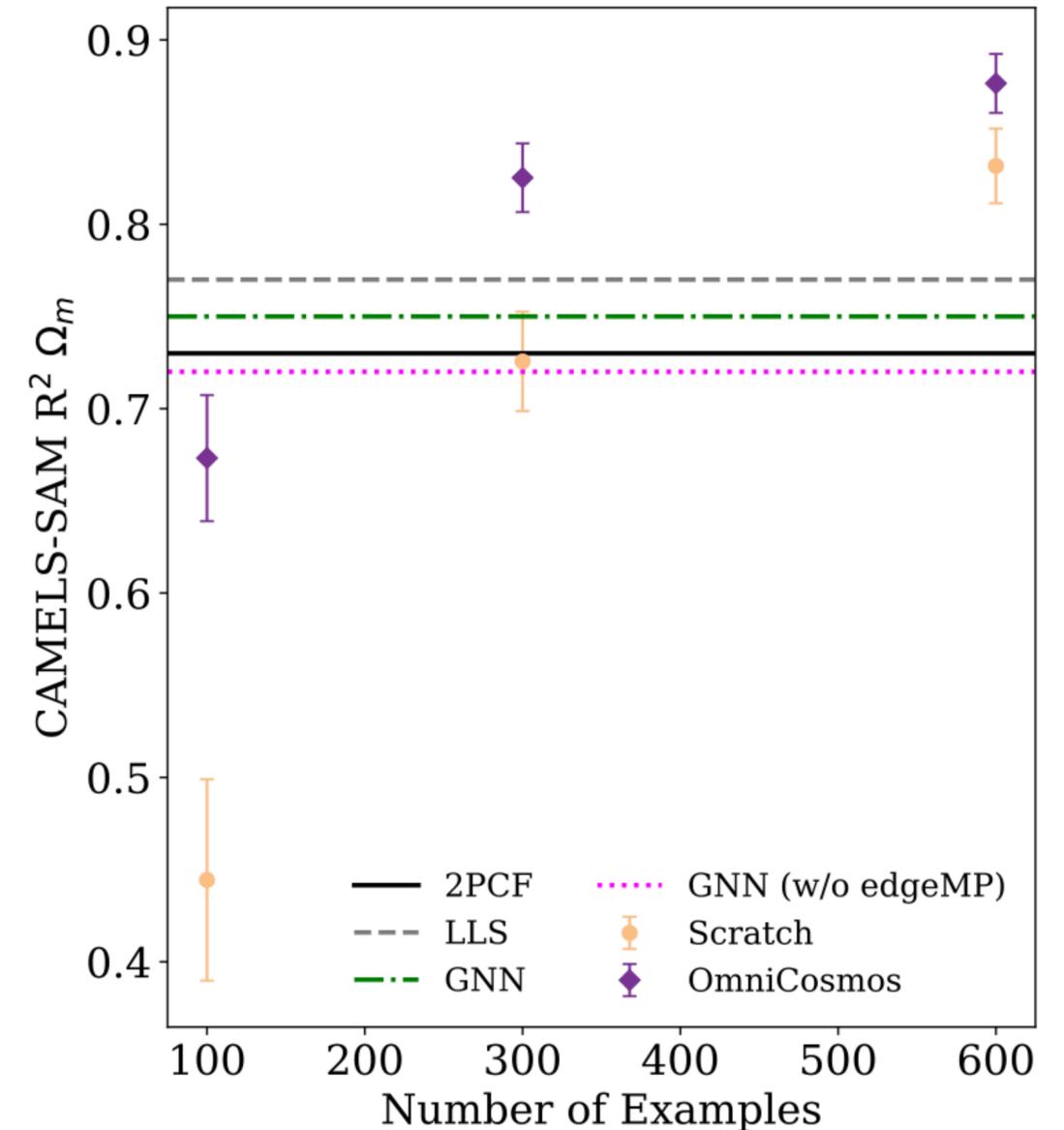
<sup>2</sup>*National Energy Research Scientific Computing Center (NERSC),  
Lawrence Berkeley National Laboratory, Berkeley, CA, USA*

<sup>3</sup>*Department of Physics, University of Toronto, Toronto, ON, Canada*

<sup>4</sup>*Department of Particle Physics and Astrophysics, Stanford University, Stanford, CA 94305, USA*

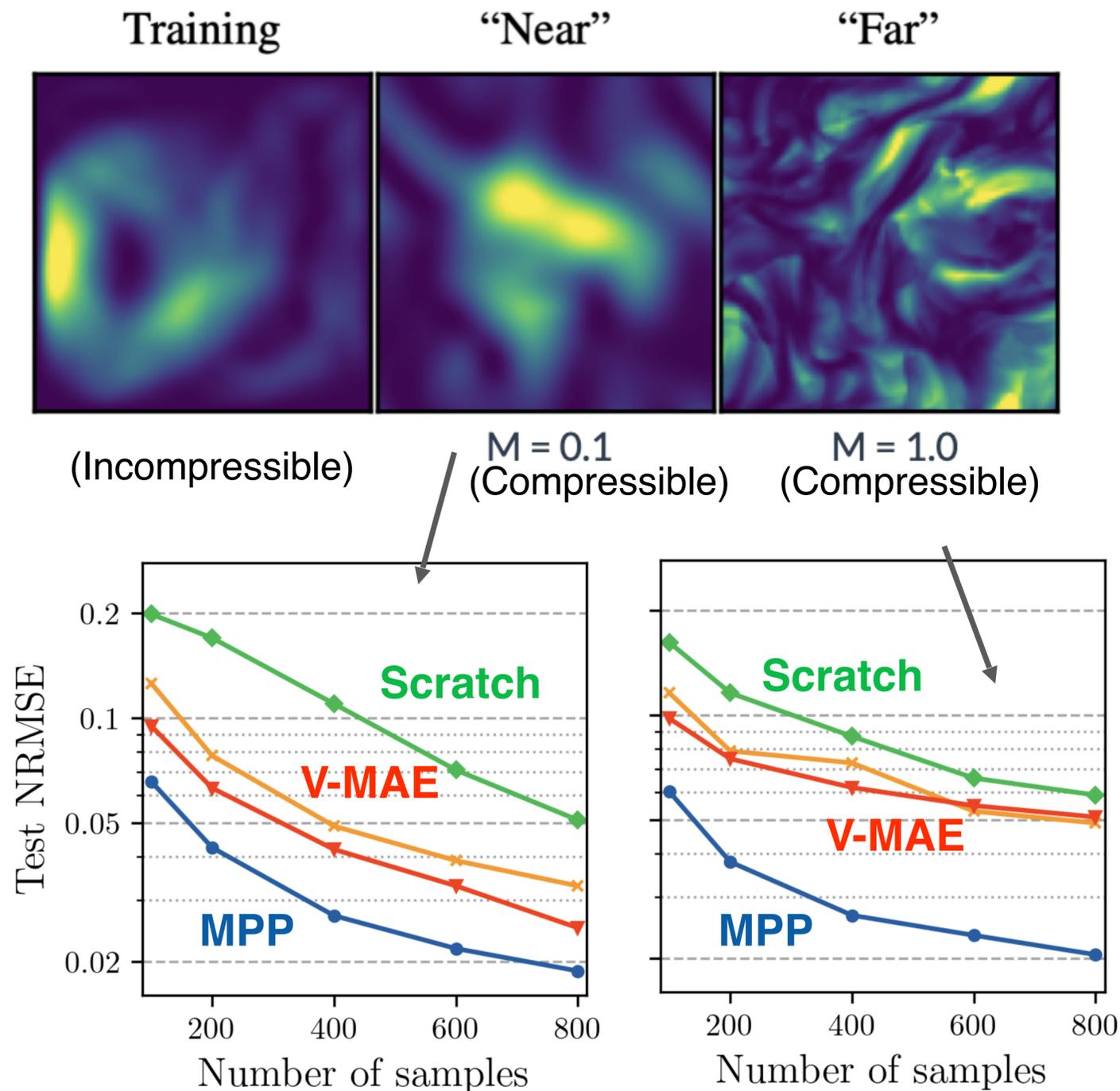
<sup>5</sup>*Fundamental Physics Directorate, SLAC National Accelerator Laboratory, Menlo Park, CA 94025, USA*

Foundation models build an effective representations of data that can be deployed on diverse downstream tasks. Previous research developed the OMNILEARNED foundation model for collider physics and showed that it could significantly advance discovery potential across collider experiments. In this paper we go beyond collider physics and show that Foundation Models trained on collider data can help improve the prediction of cosmological parameters and to predict halo and galaxy velocities in different datasets from COSMOBENCH. This is the first time a collider physics model is shown to generalize across scientific fields.





Will our analyses benefit from incorporating information from other datasets?



---

The Multimodal Universe:  
Enabling Large-Scale Machine Learning with  
100 TB of Astronomical Scientific Data

---



MULTIMODAL  
UNIVERSE

How much does the shared information content affect downstream performance?

# Multimodal Datasets with Controllable Mutual Information

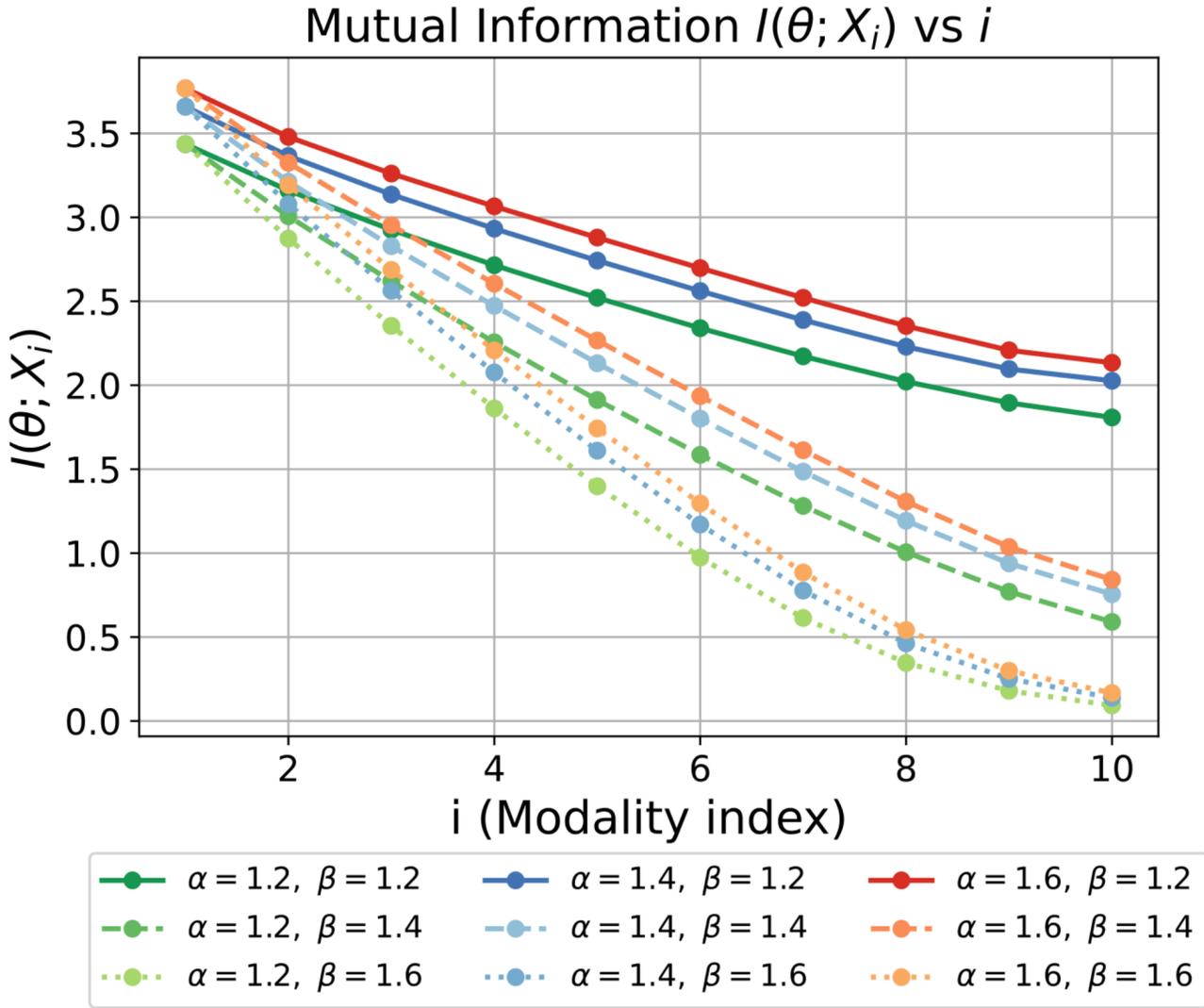
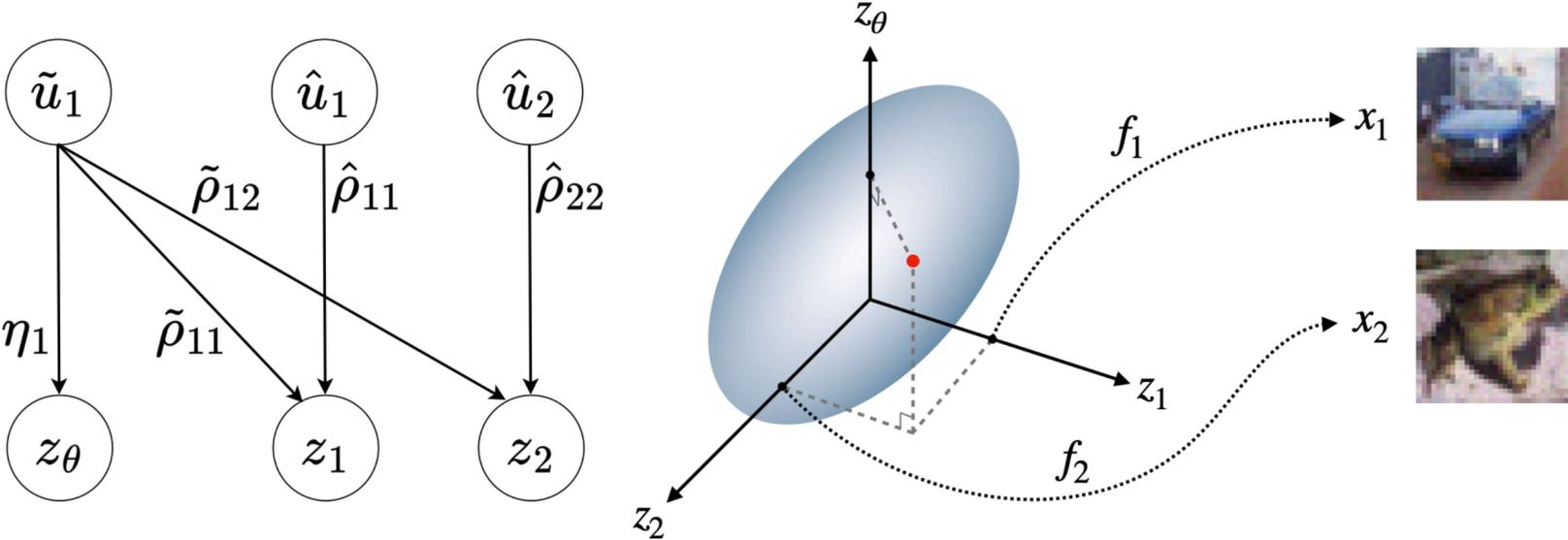
**Raheem Karim Hashmani**  
 University of Wisconsin–Madison  
 hashmani@wisc.edu

**Garrett W. Merz**  
 University of Wisconsin–Madison  
 garrett.merz@wisc.edu

**Helen Qu**  
 Flatiron Institute  
 hqu@flatironinstitute.org

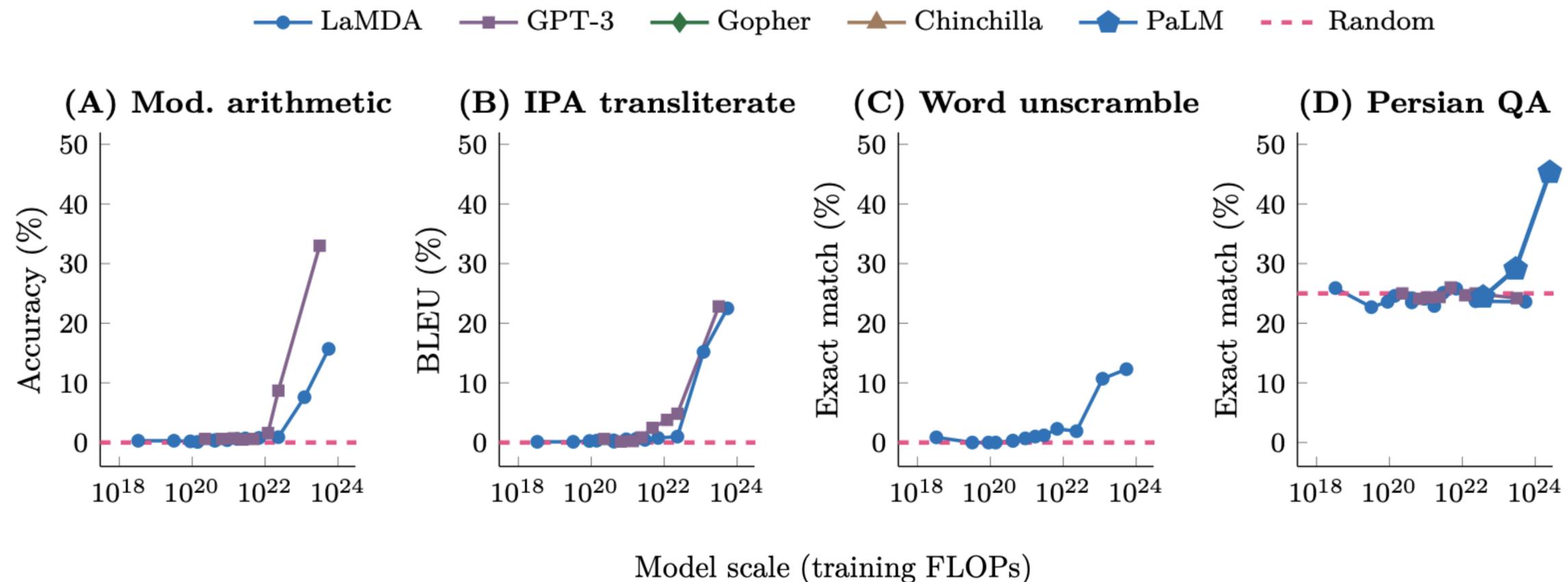
**Mariel Pettee**  
 University of Wisconsin–Madison  
 mpettee@wisc.edu

**Kyle Cranmer**  
 University of Wisconsin–Madison



# Will our analyses benefit from incorporating information from other datasets?

ChatGPT and other LLMs were trained on  $O(1\%)$  of the Internet. Their powerful generalization capabilities emerged at some point due to a combination of model capacity, heterogeneity, and training dataset size.

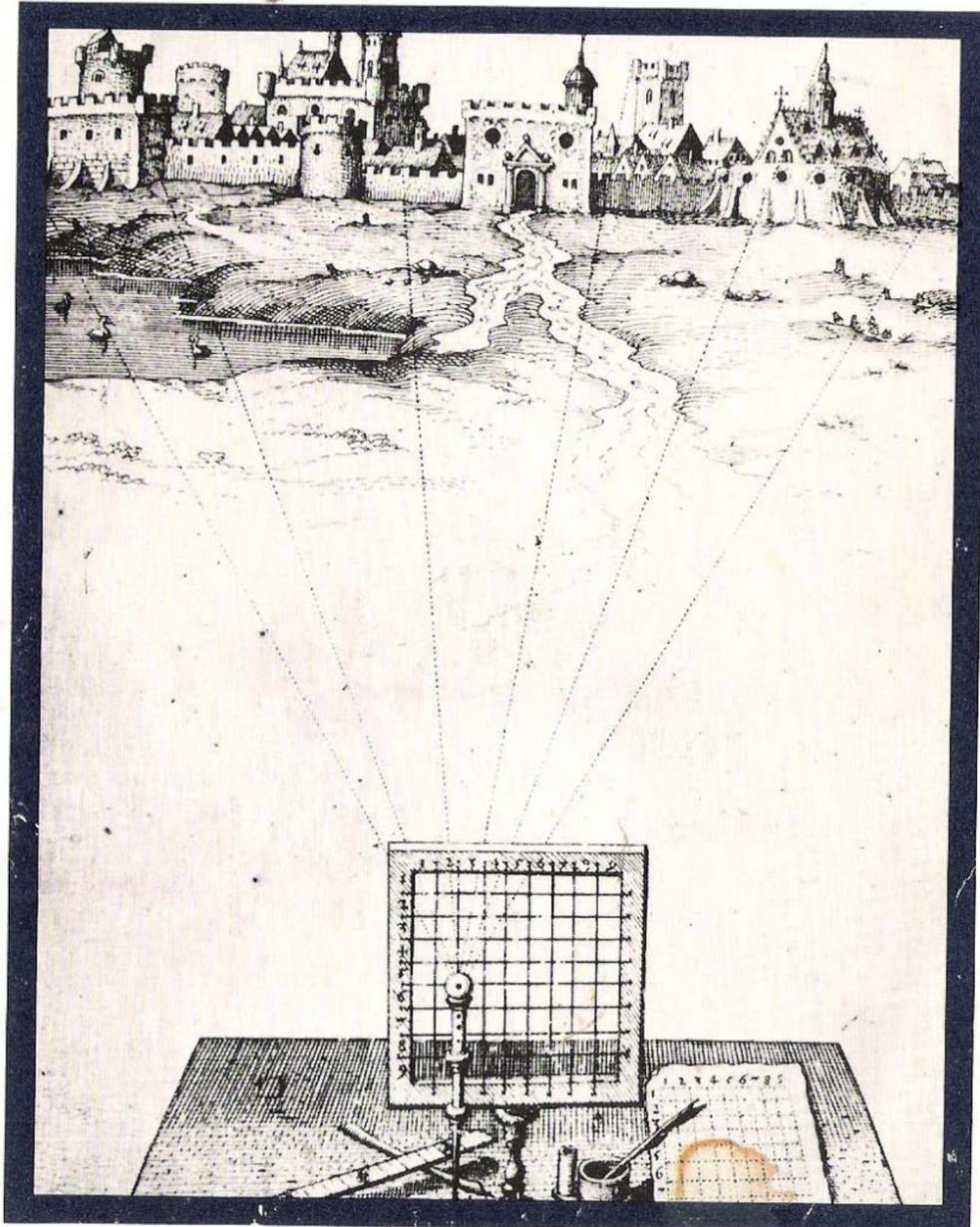


[arXiv:2206.07682](https://arxiv.org/abs/2206.07682)

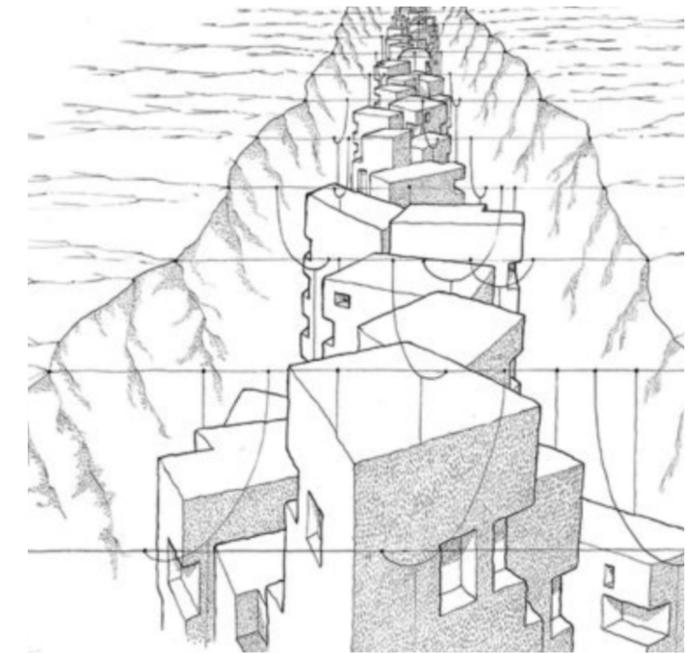
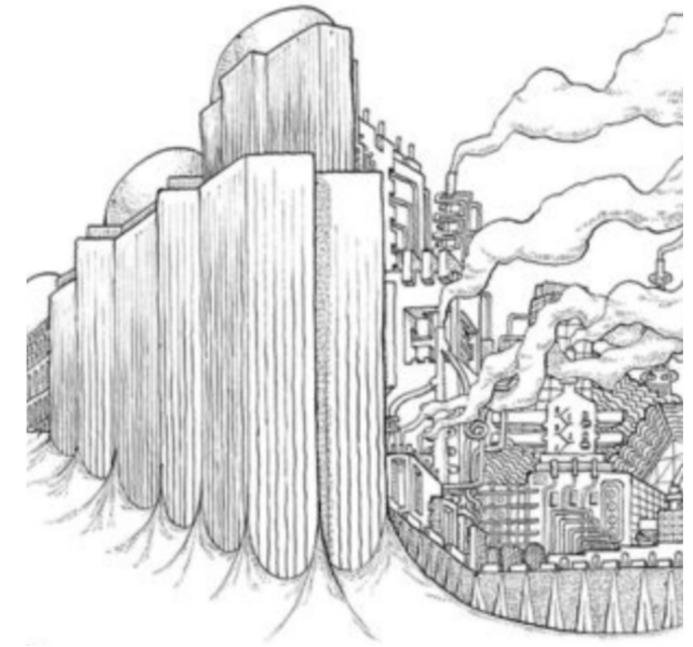
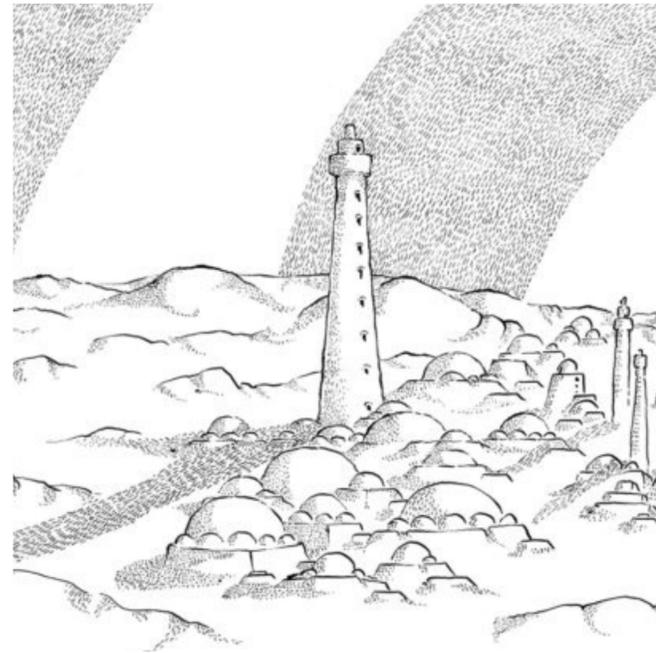
Will scientific foundation models trained on vast, heterogeneous datasets also yield emergent out-of-the-box generalization properties, as we see in many industry-scale foundation models?

A Harvest/HBJ Book

# INVISIBLE CITIES ITALO CALVINO

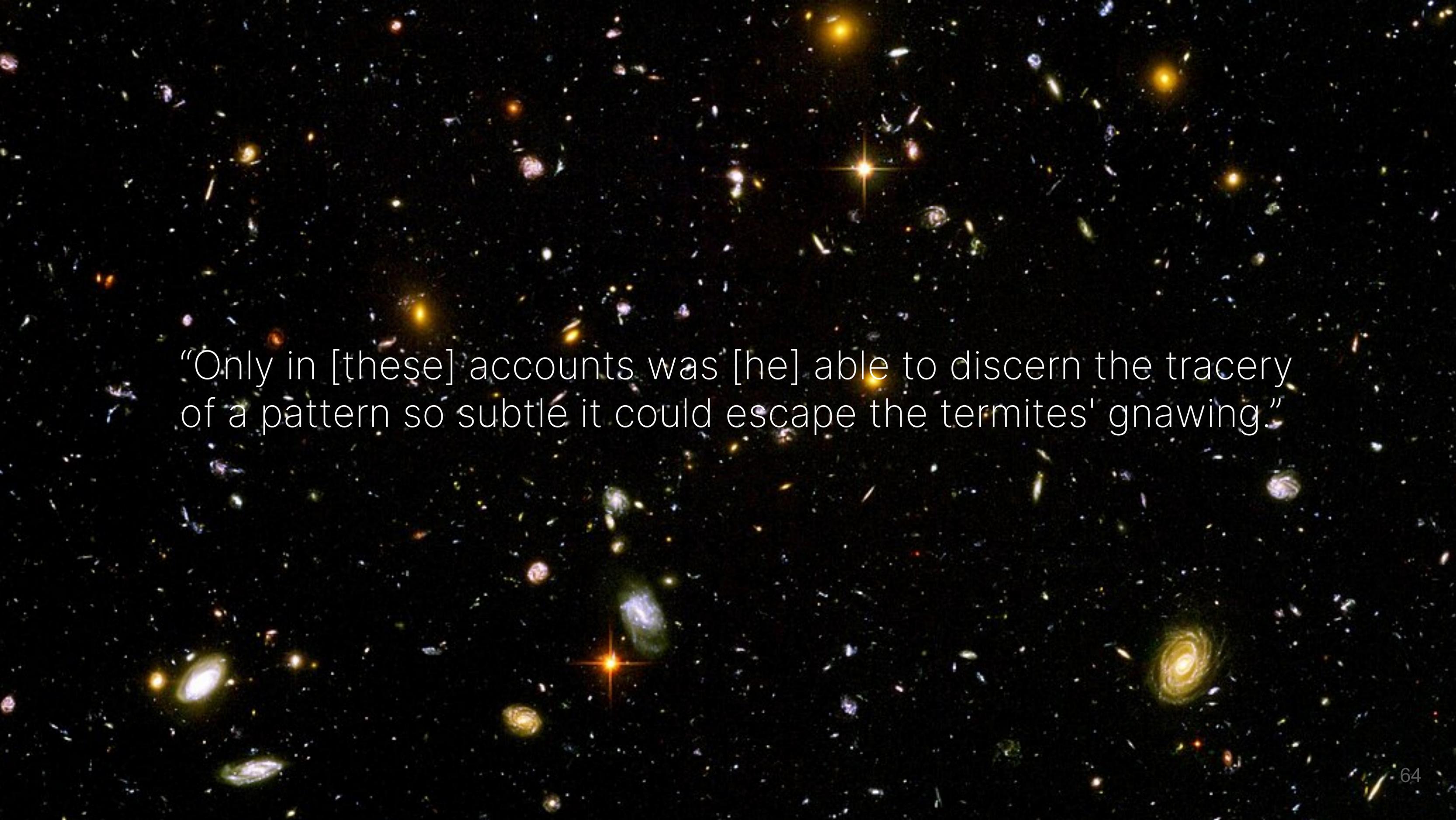


“Every time I describe a city I am saying something about Venice.”

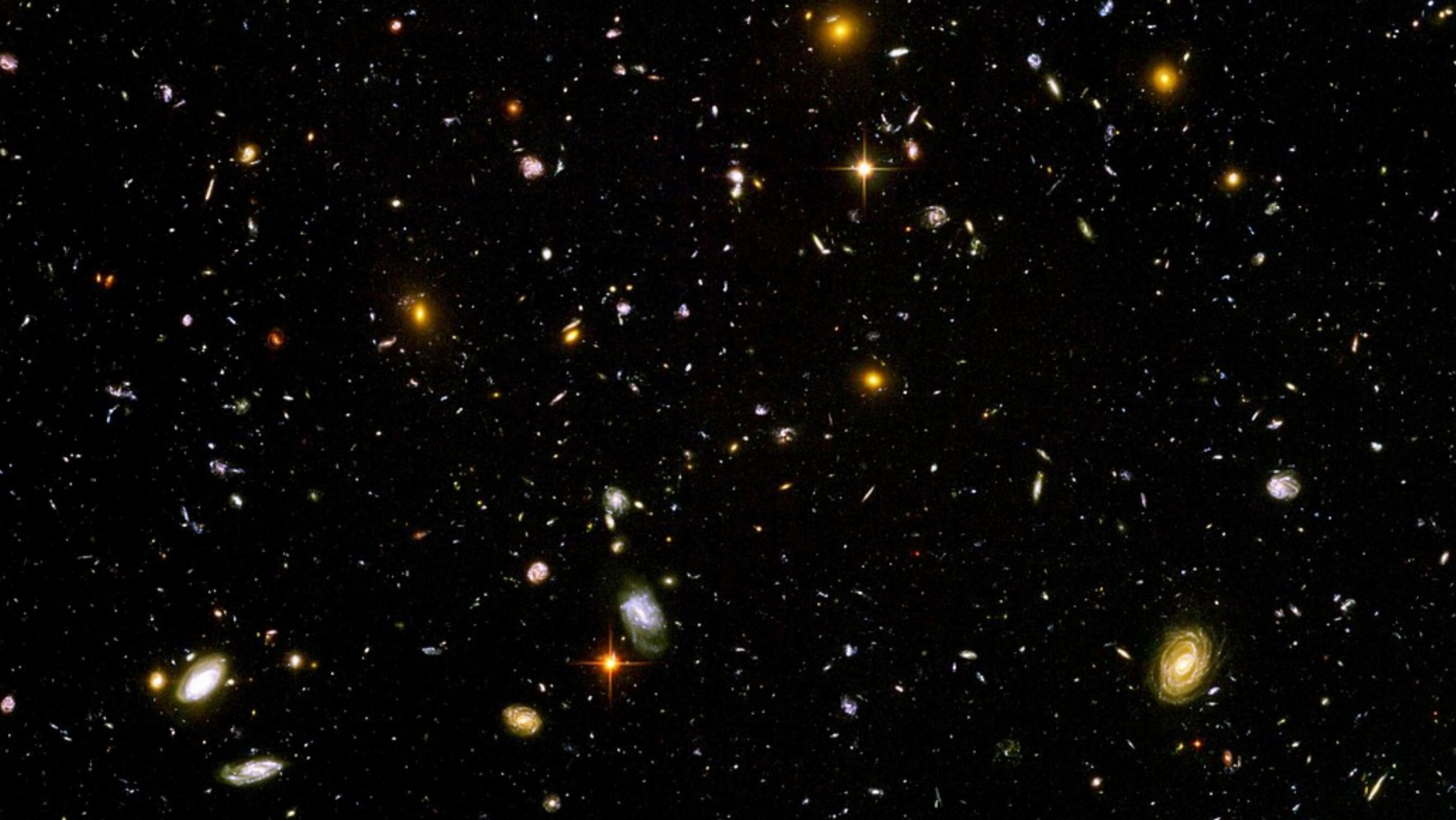


Art by Matteo Pericoli

“Only in [these] accounts was [he] able to discern the tracery of a pattern so subtle it could escape the termites' gnawing.”

A deep-field astronomical image showing a vast field of galaxies. The galaxies are scattered across the dark background, appearing in various colors (yellow, orange, blue, purple) and shapes (spiral, elliptical, irregular). Two prominent bright stars with diffraction patterns are visible in the foreground, one near the top center and one near the bottom center. The text is overlaid in the center of the image.

“Only in [these] accounts was [he] able to discern the tracery of a pattern so subtle it could escape the termites' gnawing.”





## Classifier functions can be re-used to directly approximate a likelihood ratio.

A vanilla NN classifying between two classes could be trained using **binary cross-entropy loss**:

$$\mathcal{L}_{\text{BCE}}[f] = - \int d\vec{x} (p_A(\vec{x}) \log(f(\vec{x})) + p_B(\vec{x}) \log(1 - f(\vec{x})))$$

The function that minimizes this expression will satisfy, for any small variatic  $\delta f(\vec{x})$  :

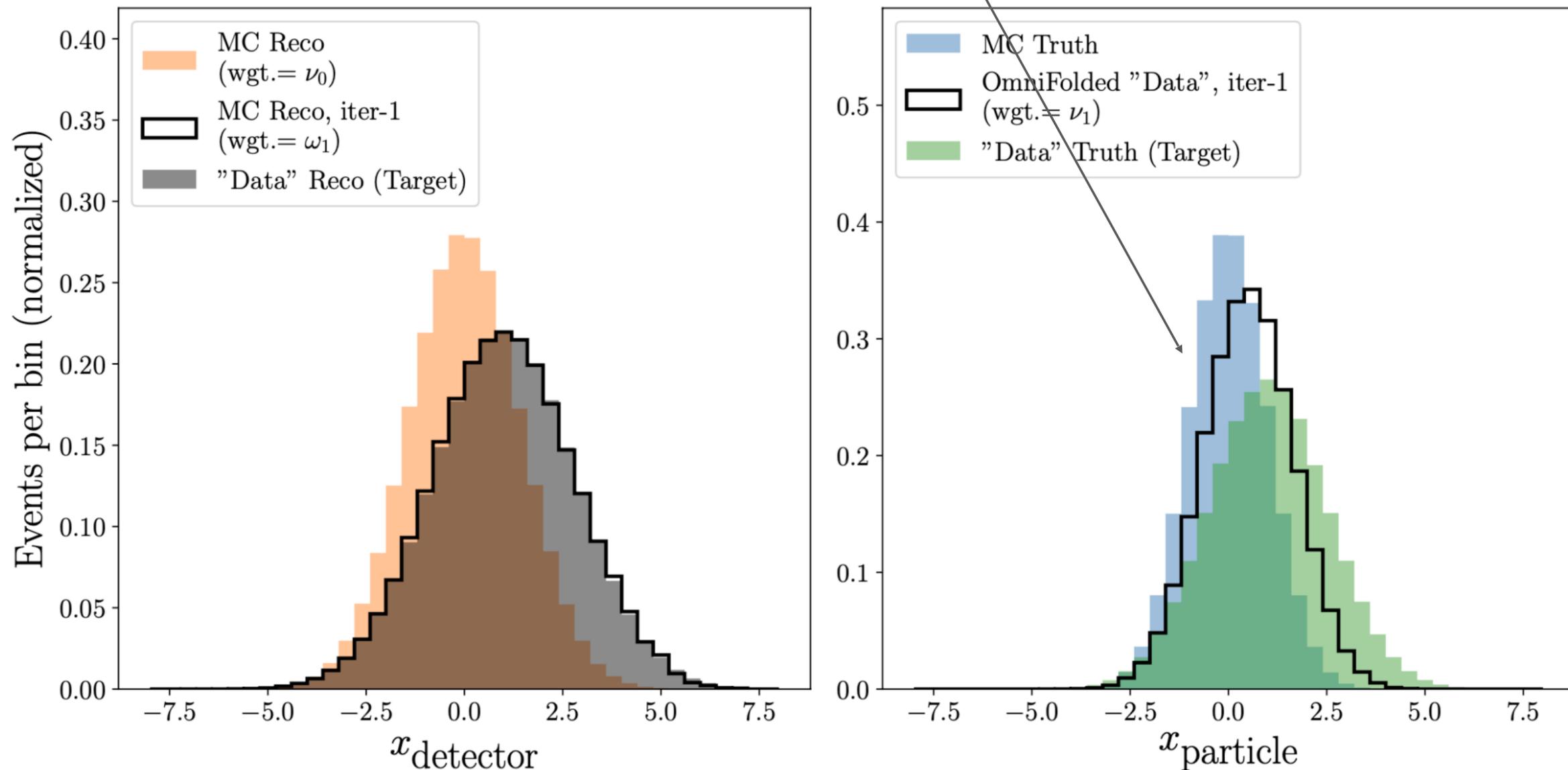
$$\delta \mathcal{L}_{\text{BCE}} = \int d\vec{x} \left( -\frac{p_A(\vec{x})}{f(\vec{x})} + \frac{p_B(\vec{x})}{1 - f(\vec{x})} \right) \delta f(\vec{x}) = 0.$$


$$-\frac{p_A(\vec{x})}{f(\vec{x})} + \frac{p_B(\vec{x})}{1 - f(\vec{x})} = 0$$


$$\text{Rescaling of classifier output} \quad \frac{f(\vec{x})}{1 - f(\vec{x})} = \frac{p_A(\vec{x})}{p_B(\vec{x})} \quad \text{Likelihood ratio}$$

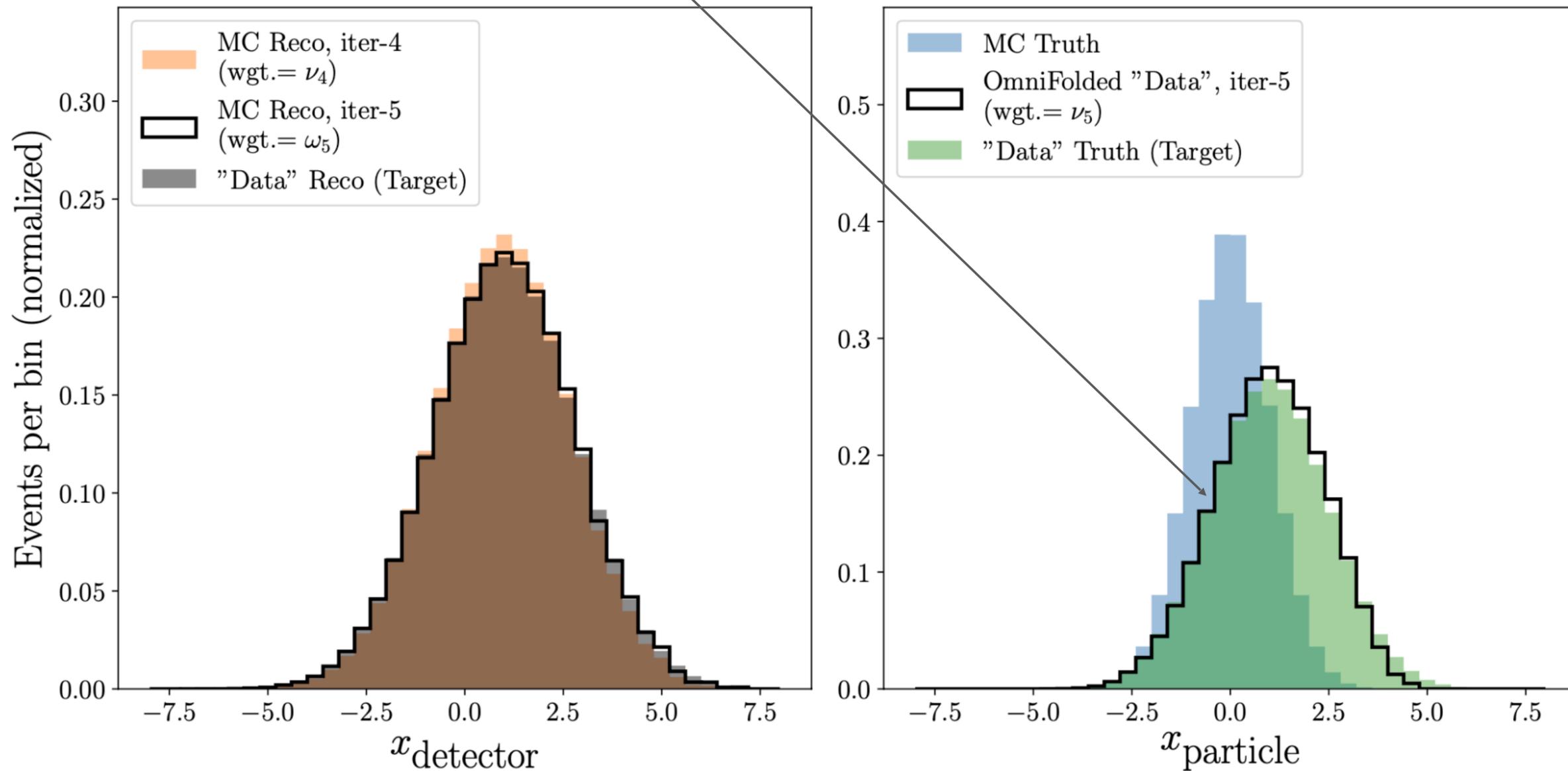
# OmniFold reweights truth-level MC to estimate “truth-level” data.

After one iteration of our method, the **truth-level MC** is **reweighted** to look more like **truth-level data**.



# OmniFold reweights truth-level MC to estimate “truth-level” data.

After five iterations of our method, the **reweighted truth-level MC** is closely aligned with **truth-level data**.



The OmniFold procedure requires two datasets\* as inputs:

(\*In practice, we use samples from different MC generators as well as systematically-shifted samples to determine uncertainties.)

- MC sample with events at both detector-level and particle-level
- Real data

In a multi-stage and iterative process, a series of neural networks are trained to learn a reweighting function that maps particle-level MC distributions to particle-level data distributions.

