



ICLR

International Conference On
Learning Representations



ICLR

International Conference On
Learning Representations

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal “World Models”

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal "World Models"

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal “World Models”

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal “World Models”

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal “World Models”

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Observability Foundation Models

From Scaling Time Series Foundation Models to Multimodal “World Models”

ICLR 2026

Othmane Abou-Amal – Director, Datadog AI Research



DATADOG

Agenda

01 Observability 101

02 Why Foundation Models for observability

03 Scaling Toto - From the Bert moment to the GPT-2 Moment

04 From TSFMs to Multimodal "World Models"

05 Questions?

Observability 101



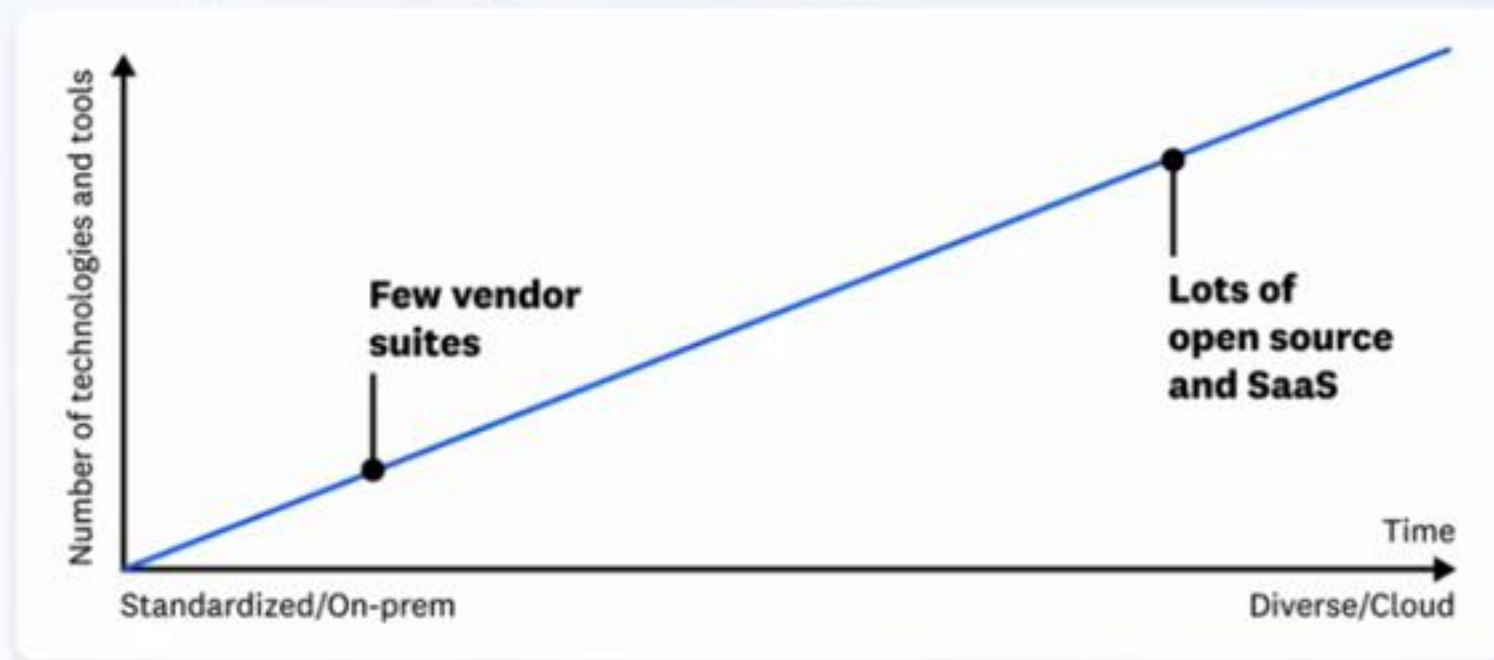
The **3am** page problem!



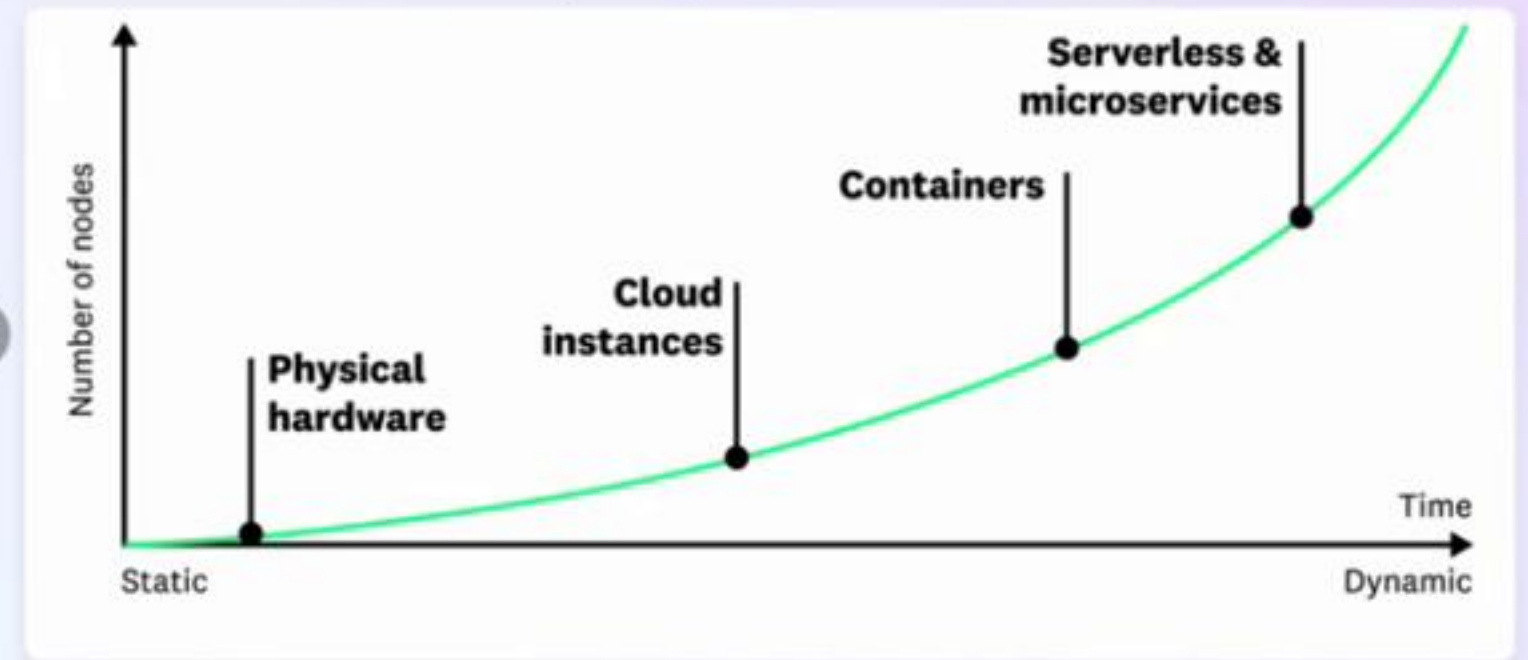
The **3am** page problem!

The problem: an explosion of complexity

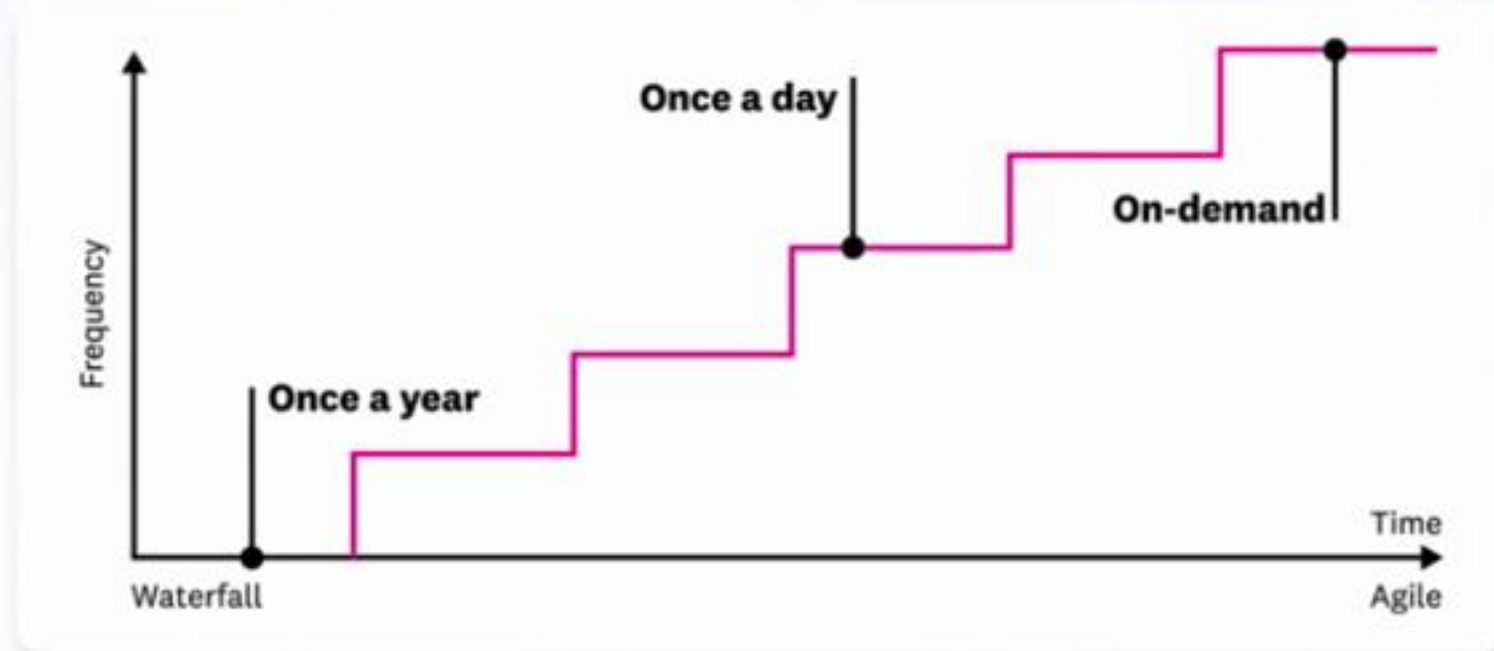
Diversity of technologies in use



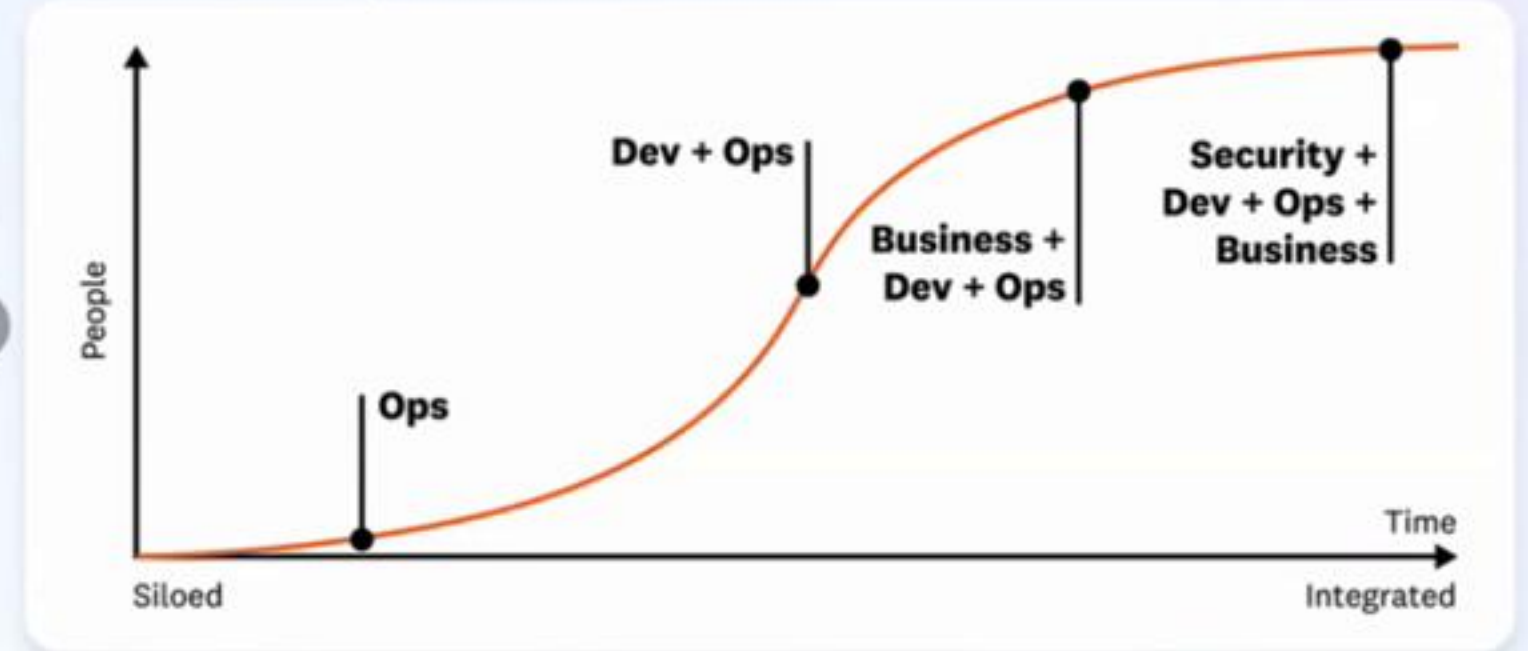
Scale in number of computing units



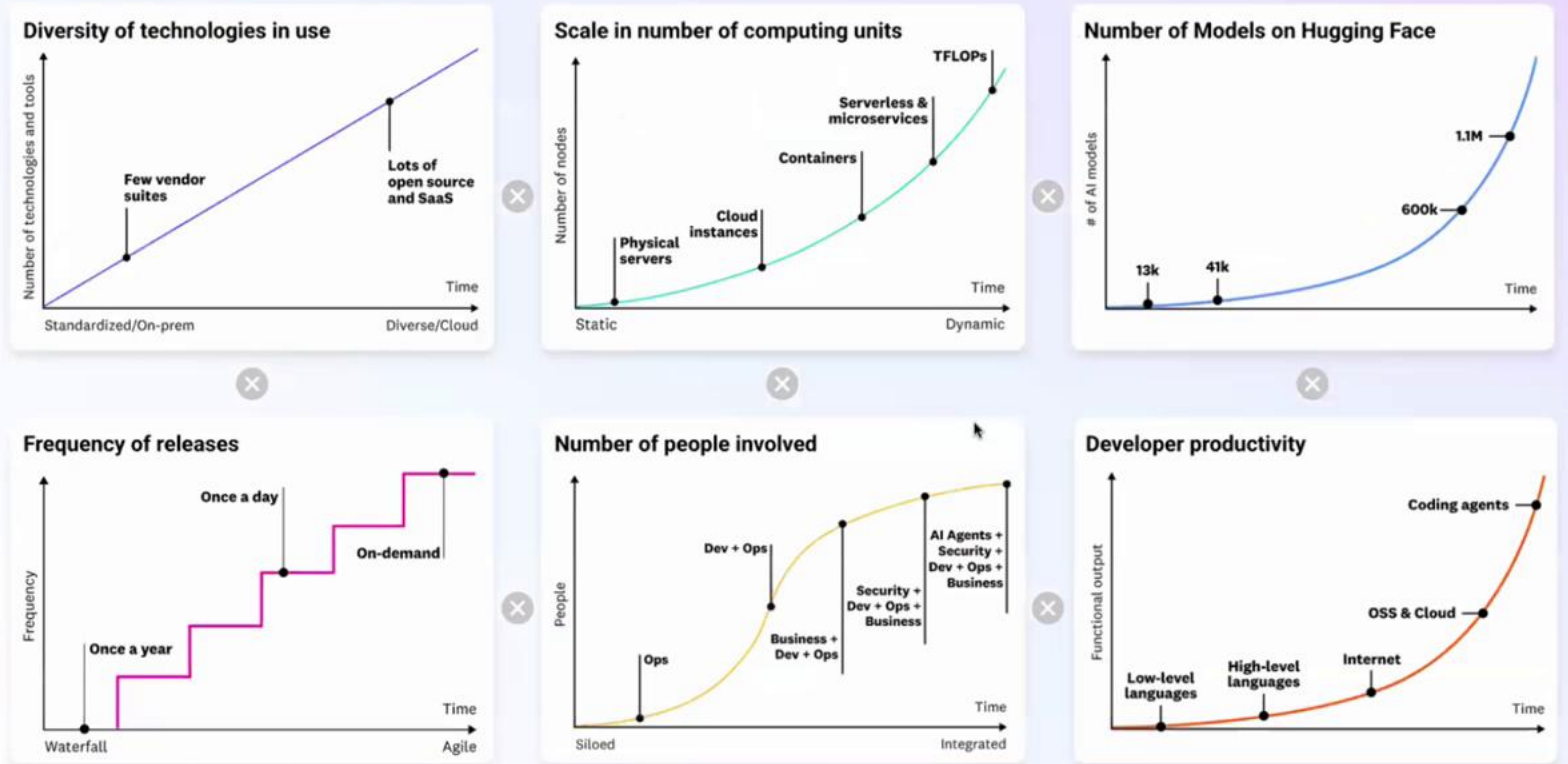
Frequency of release



Number of people involved

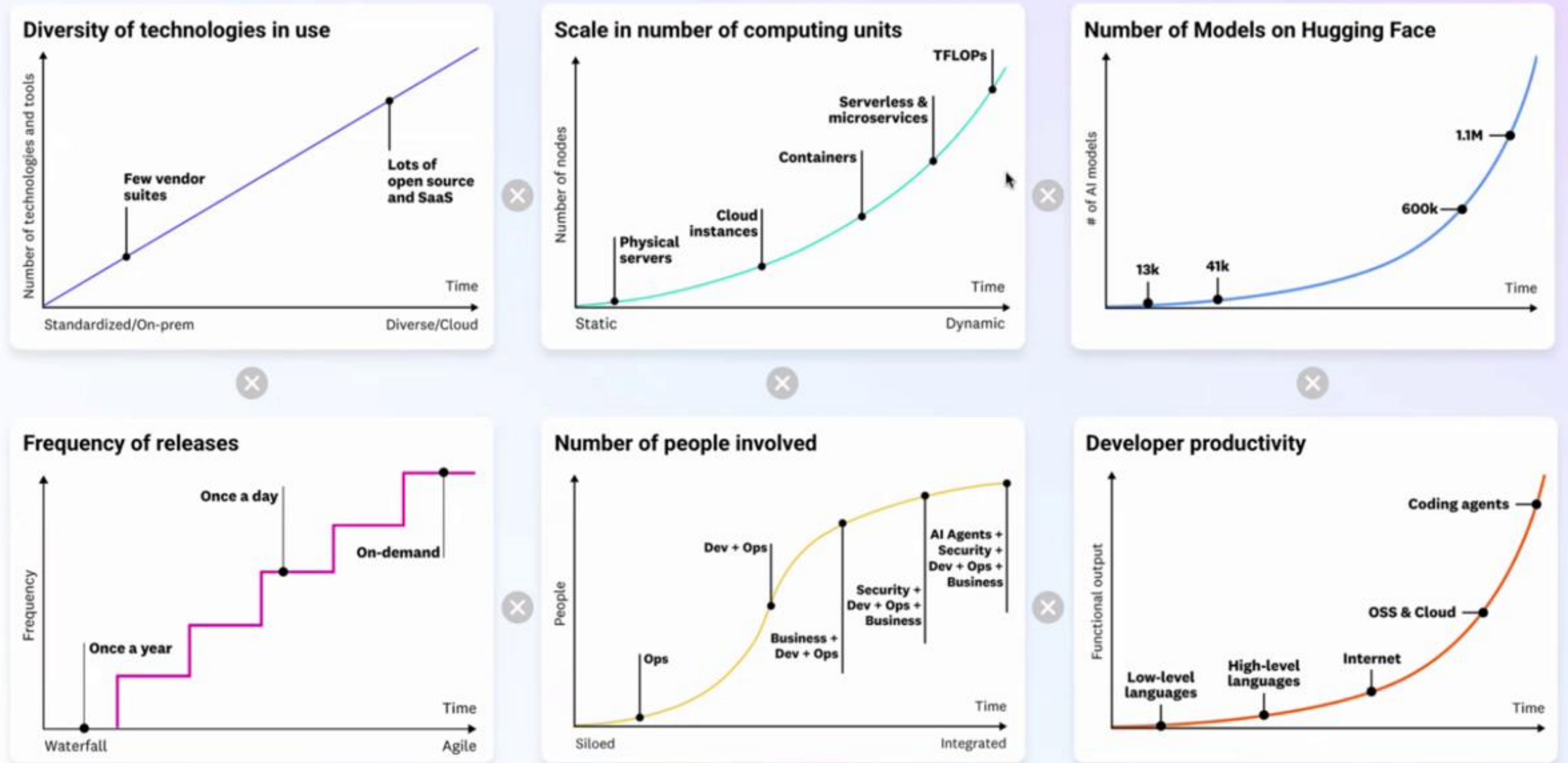


AI compounds complexity



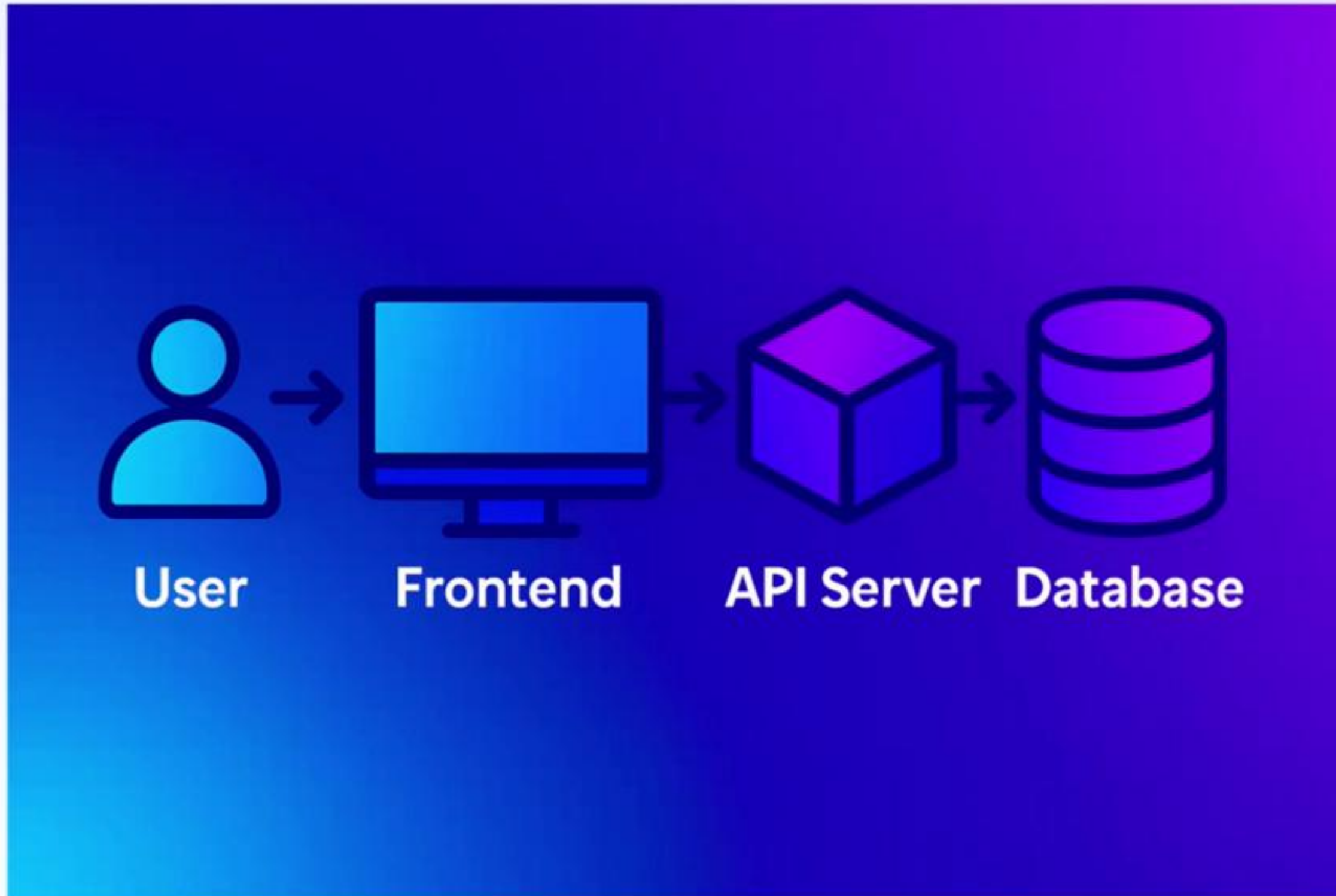
Source for number of models: Hugging Face Hub Stats Dashboard, cfahlgren1, 2025.

AI compounds complexity

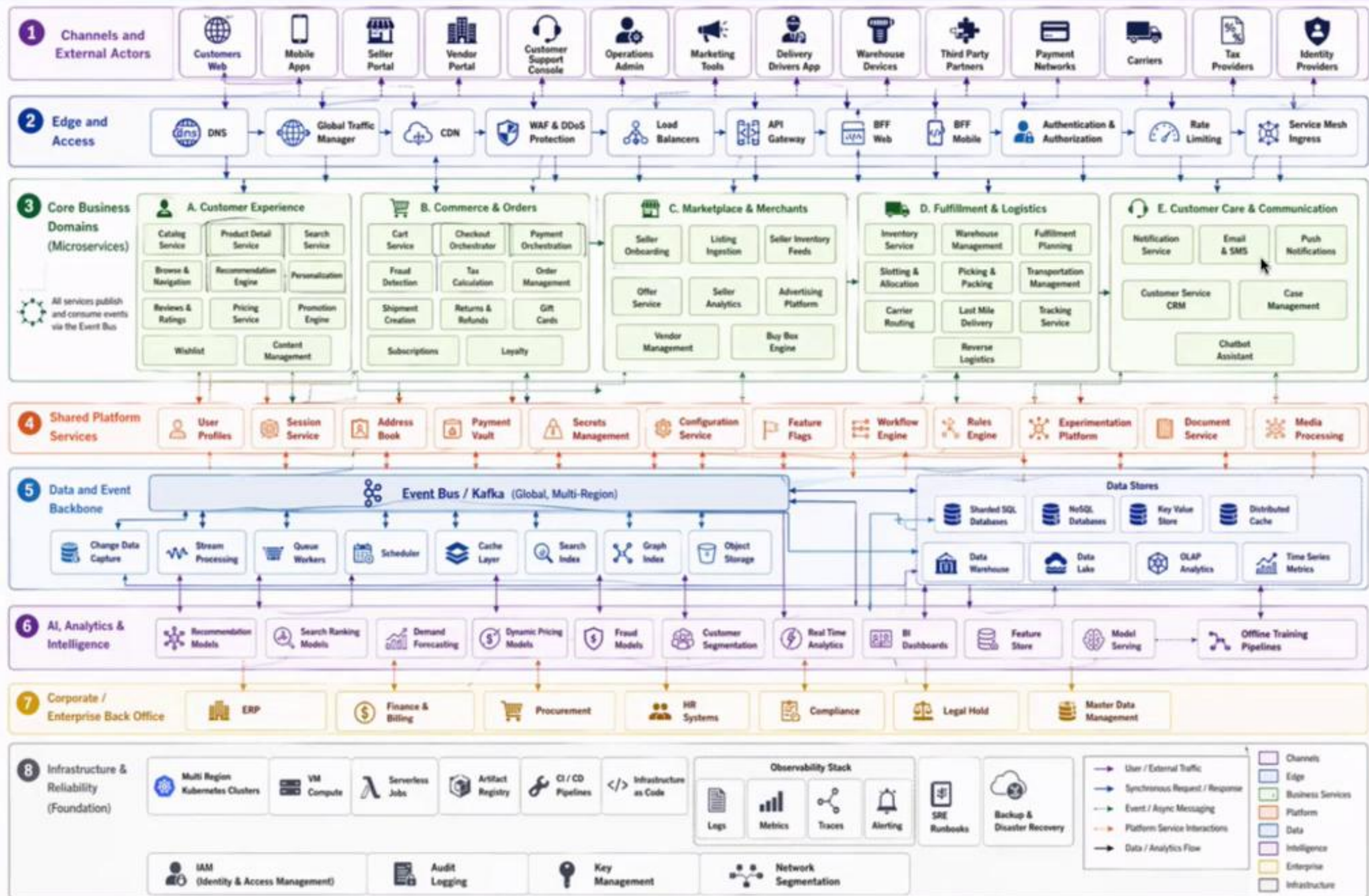


Source for number of models: Hugging Face Hub Stats Dashboard, cfahlgren1, 2025.

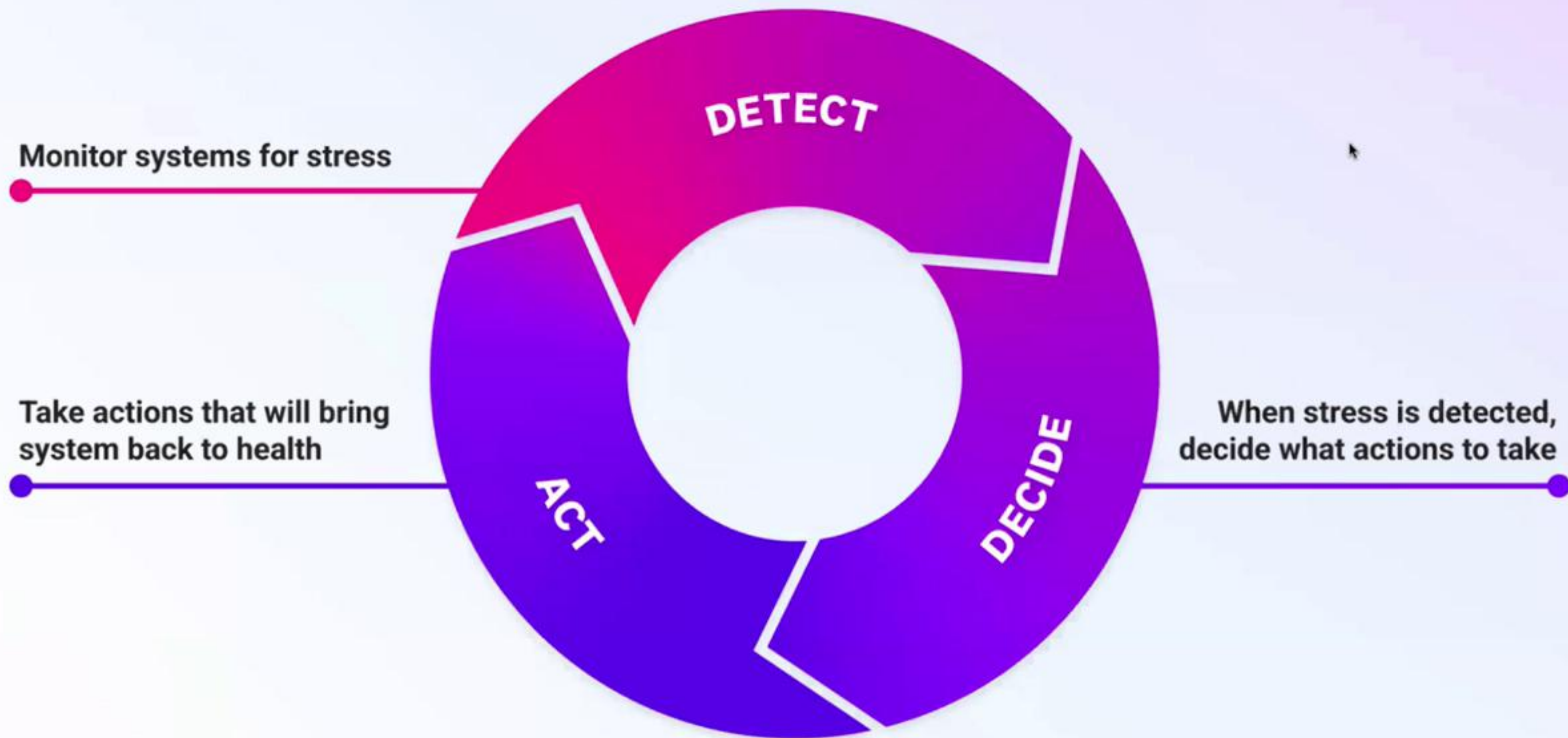
Example: A “simple” e-commerce site



Exa



Close the DevOps loop



Why Foundation Models for observability

This not a solved problem.

But

**Data is our
differentiated
advantage**

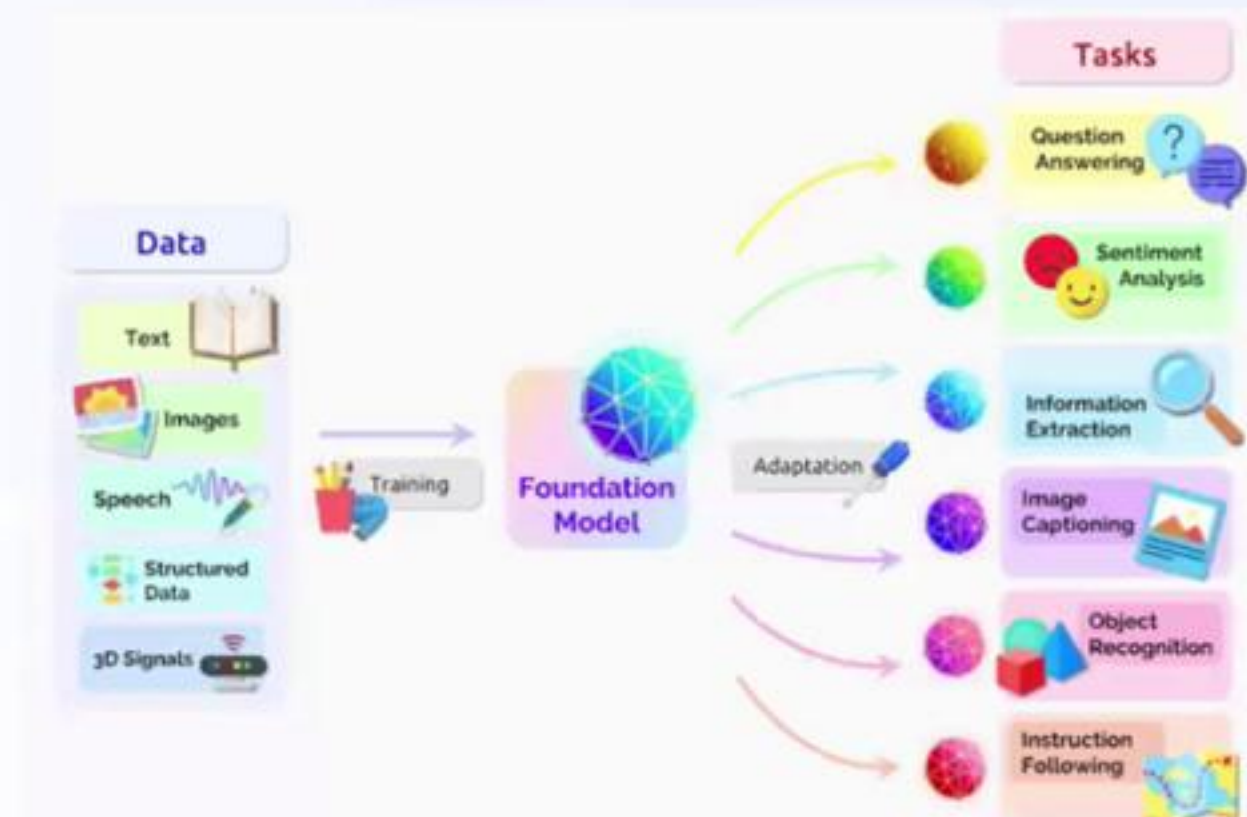
Trillions of datapoints per hour

Trillions of
metrics

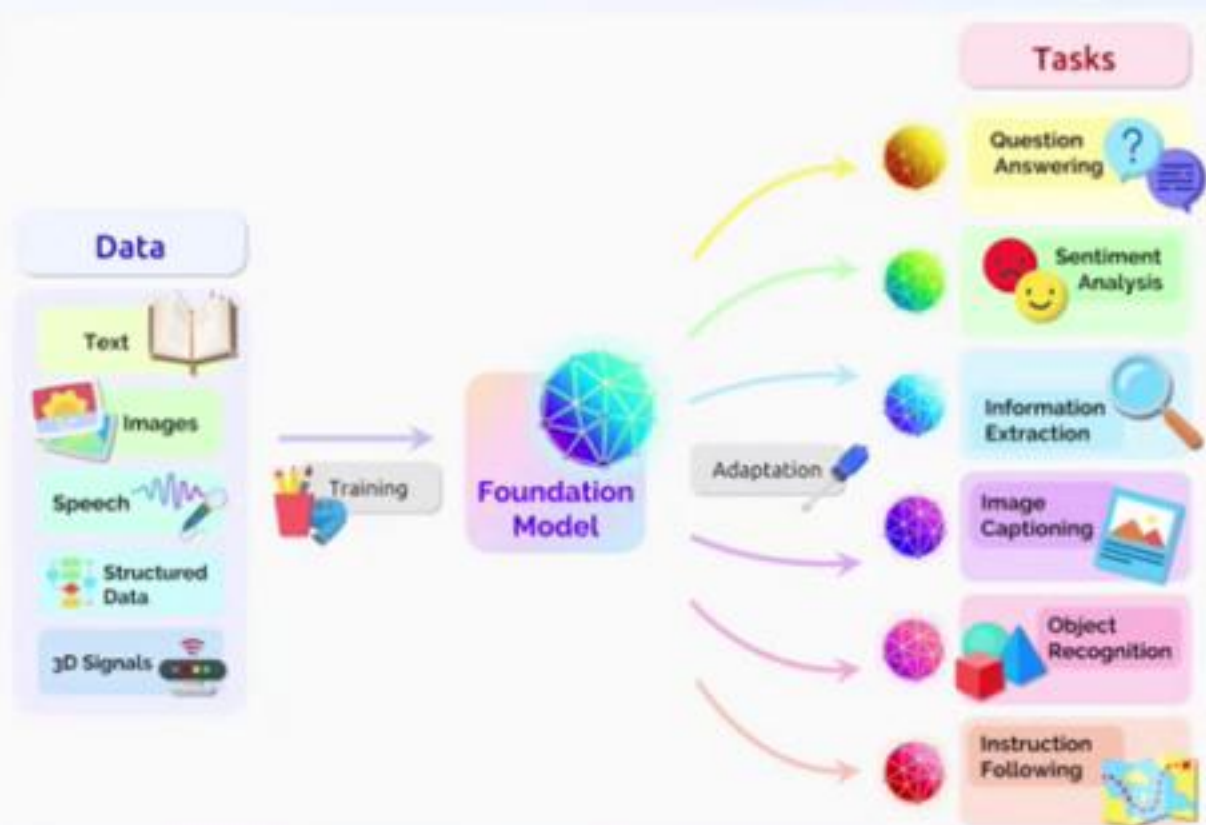
Billions of
application traces

Exabytes
of logs

Foundation Models unlock scale and performance across modalities



Foundation Models unlock scale and performance across modalities



**Observability
is next?**

TOTO

Time Series Foundation Model (TSFM)

SOTA zero-shot forecasts

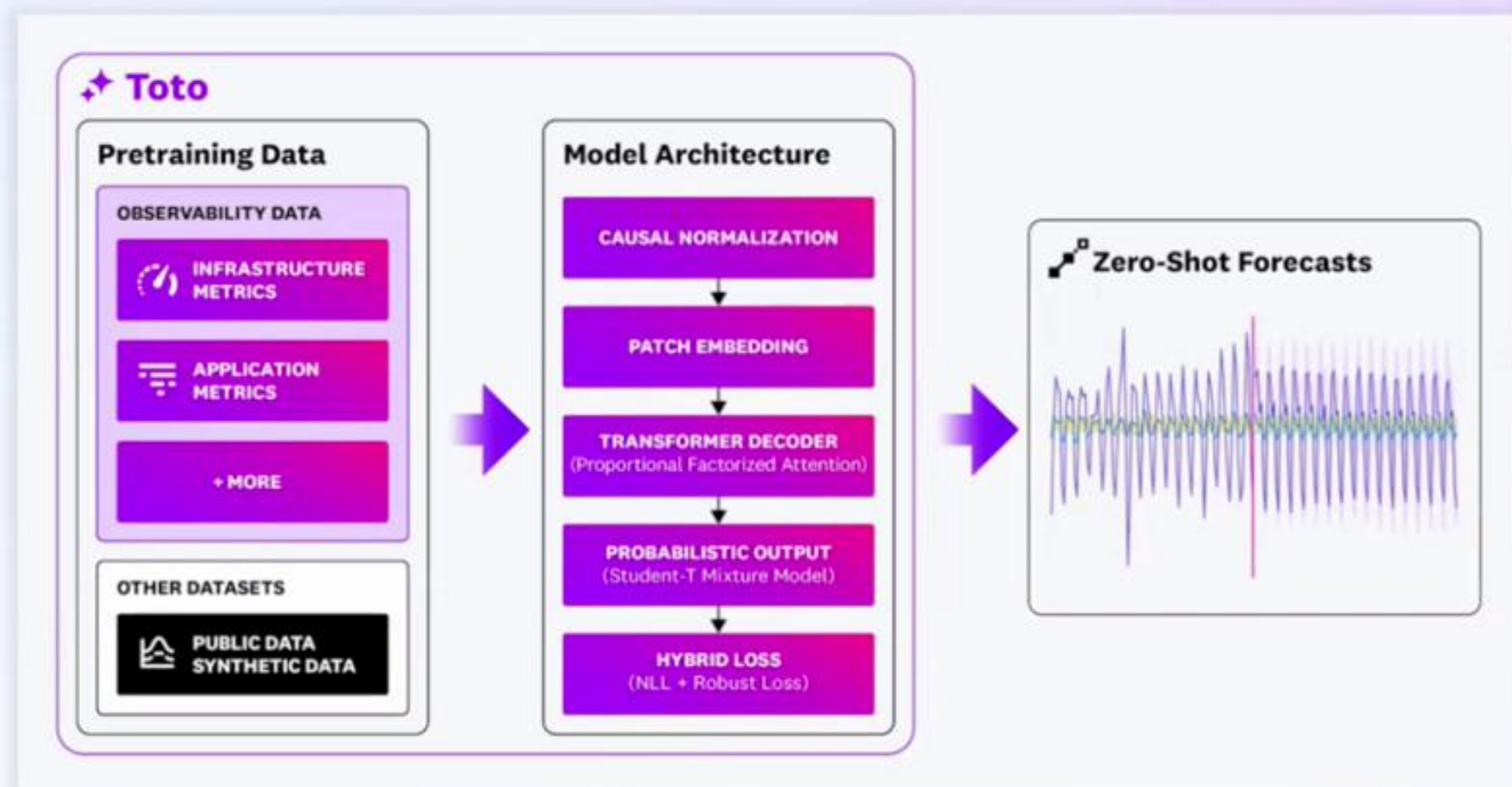
Optimized for Observability

And also SOTA on general-purpose time series forecasting

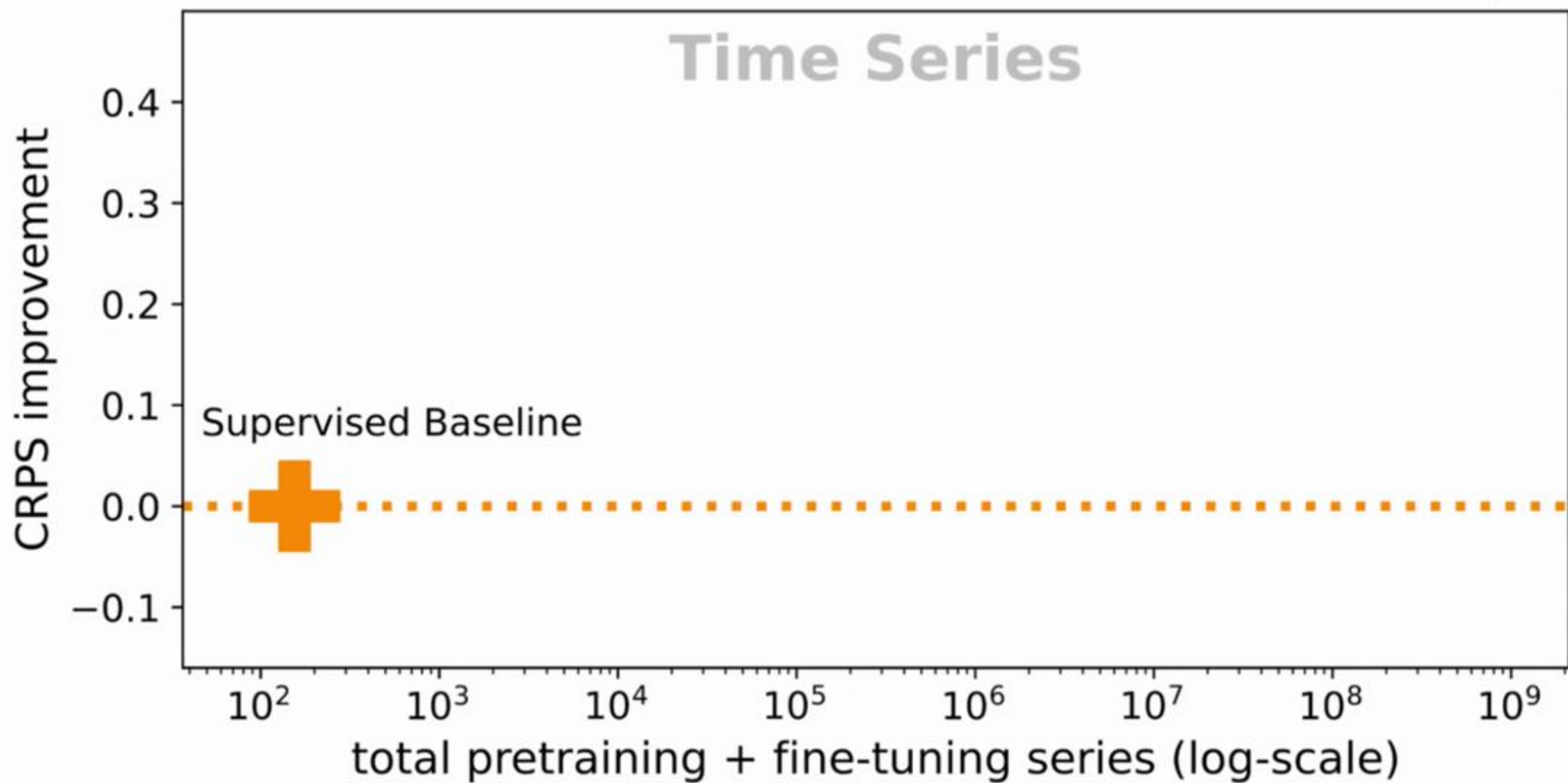
Open Source/Weights

Apache 2.0

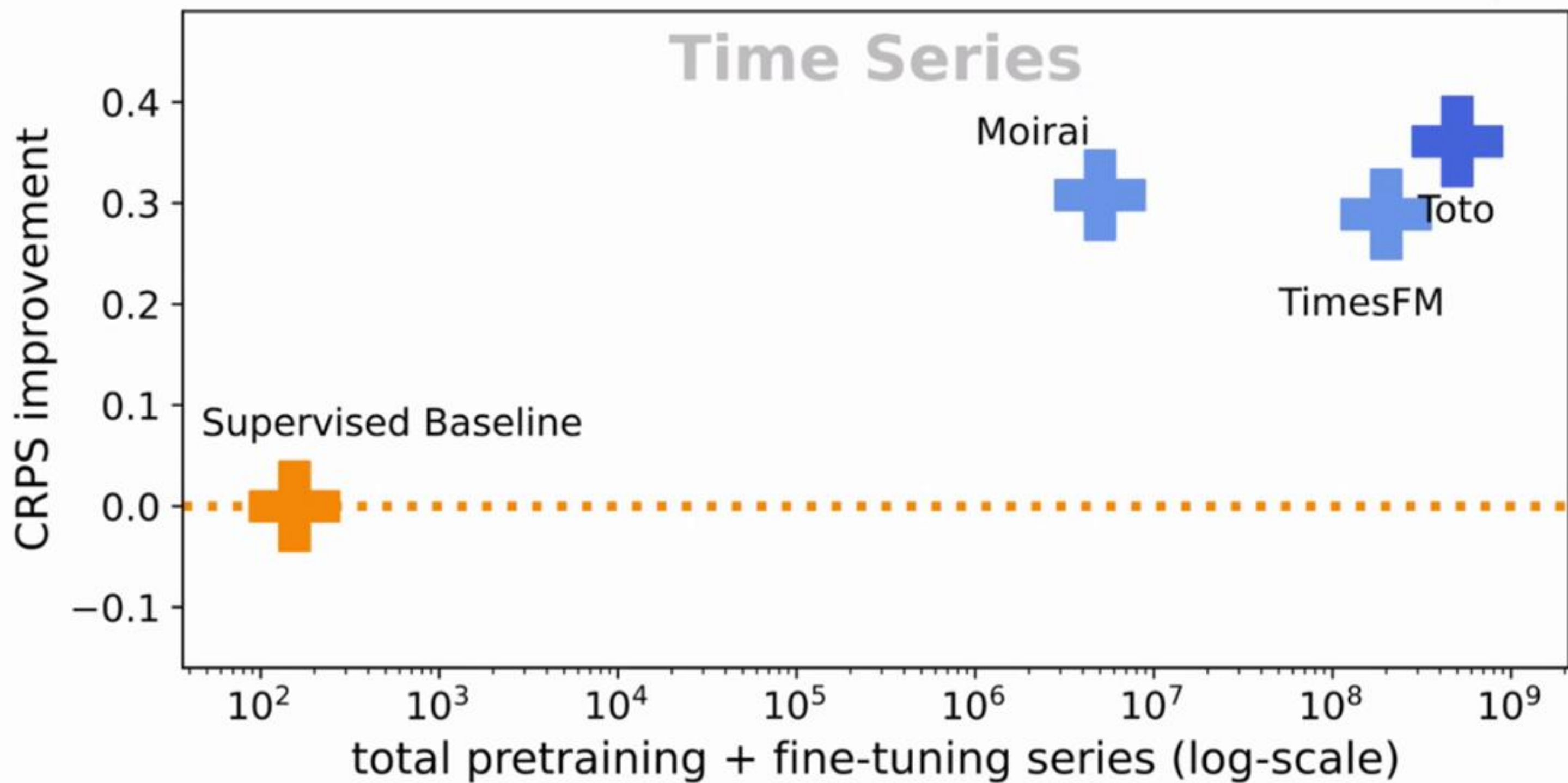
Downloaded >9M times from HF



BOOM Results



BOOM Results



Scaling Toto – From the BERT moment to the GPT-2 Moment

Scaling Toto – From the BERT moment to the GPT-2 Moment

Toto-2.0



Toto-2.0



Toto-2.0



Toto-2.0




Toto-2.0






Toto-2.0

Toto-2.0

updated 3 days ago

 Add a description




Datadog/Toto-2.0-4m

 Time Series Forecasting · Updated 3 days ago ·  47 ·  1




Datadog/Toto-2.0-22m

 Time Series Forecasting · Updated 3 days ago ·  151 ·  1




Datadog/Toto-2.0-313m

 Time Series Forecasting · Updated 3 days ago ·  42 ·  1

Datadog/Toto-2.0-1B

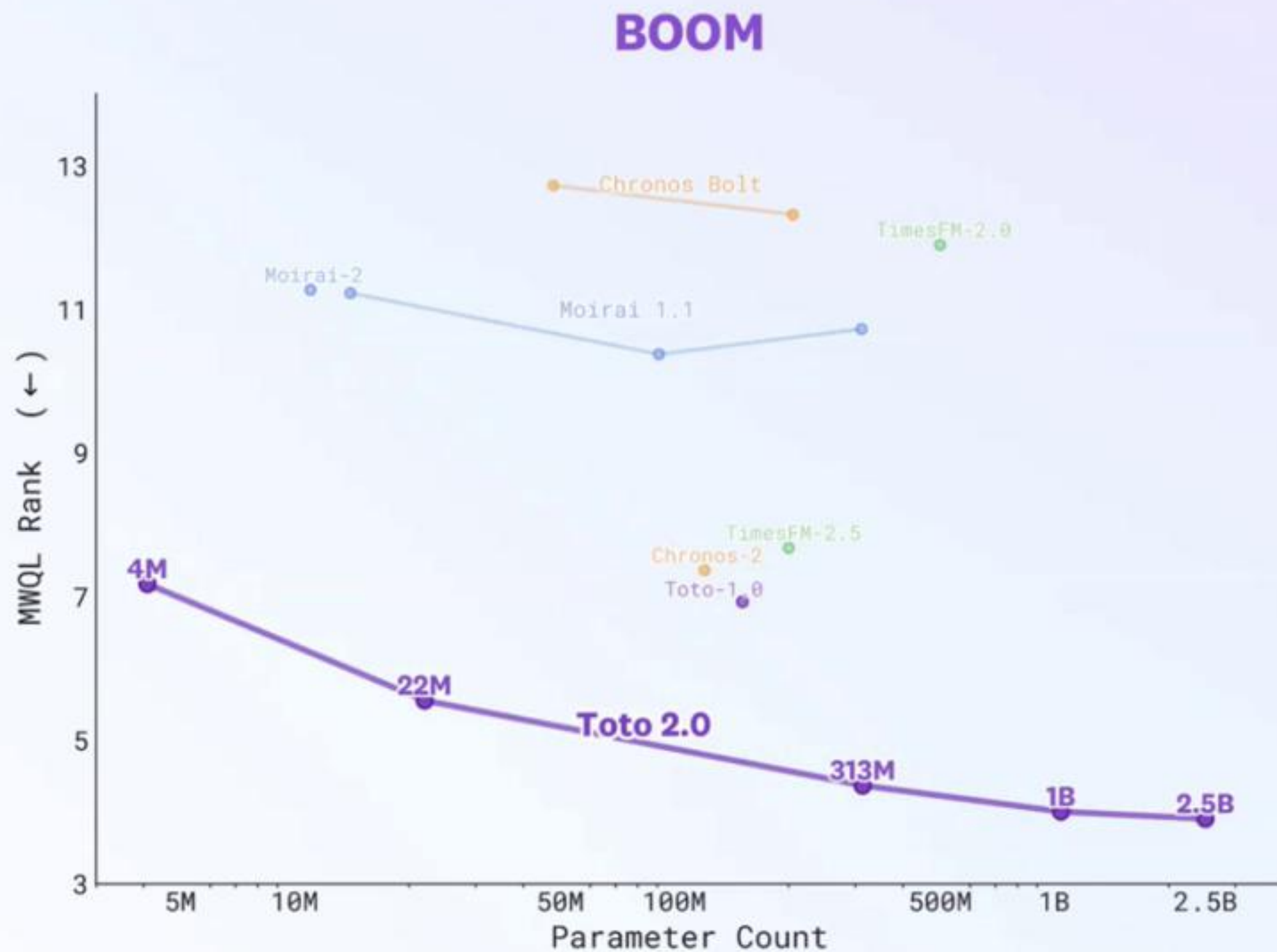
 Time Series Forecasting · Updated 3 days ago ·  50 ·  1

Datadog/Toto-2.0-2.5B

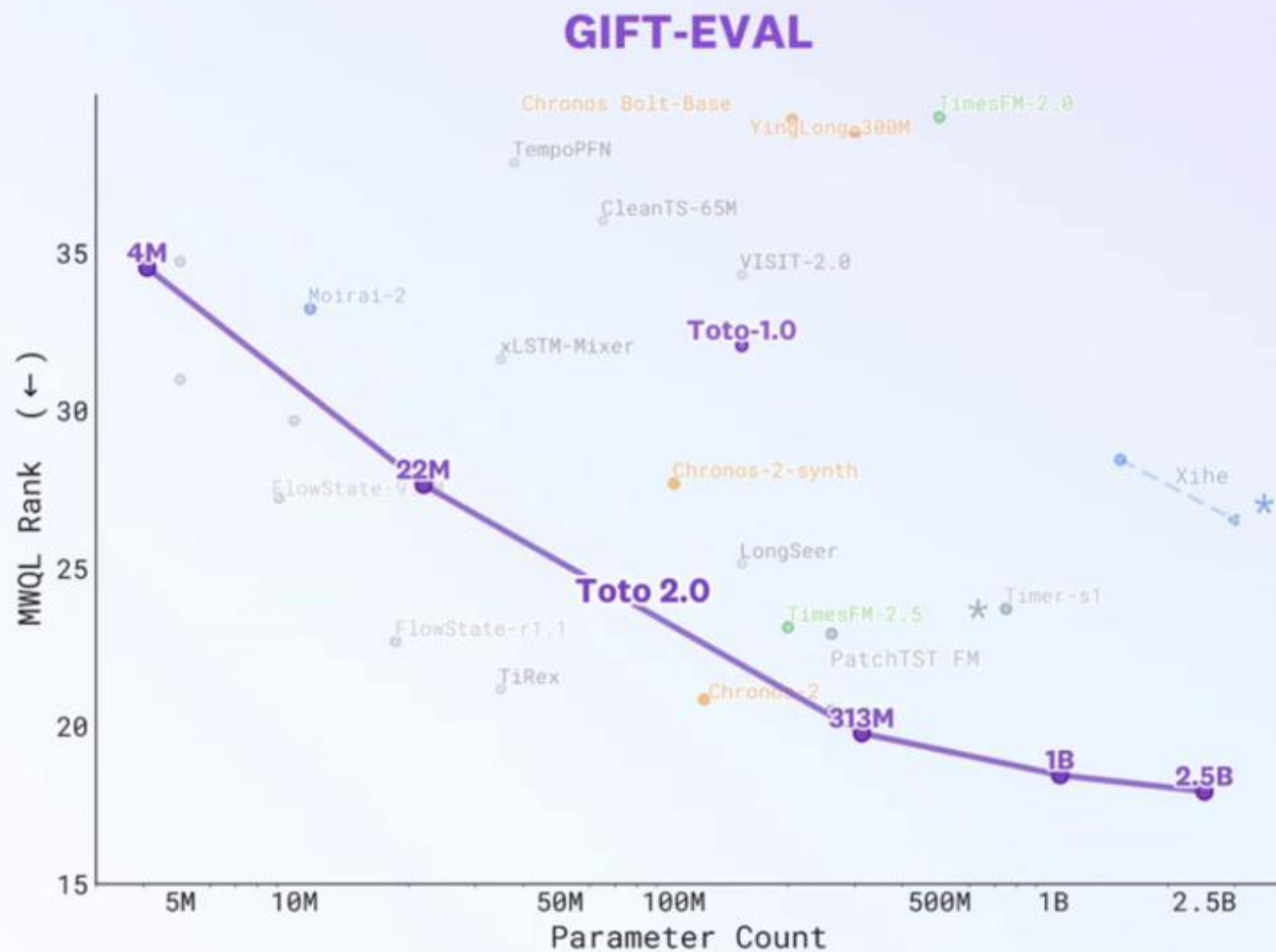
 Time Series Forecasting · Updated 3 days ago ·  11 ·  3



tl;dr We can scale to a SOTA model (the GPT-2 Moment)



tl;dr We can scale to a SOTA model (the GPT-2 Moment)



Emaad Khwaja, Chris Lettieri, Gerald Woo

Eden Belouadah, Marc Cenac, Guillaume Jarry, Enguerrand Paquin, Xunyi Zhao, Viktoriya Zhukova

Othmane Abou-Amal, Ameet Talwalkar, David Asker

Emaad Khwaja, Chris Lettieri, Gerald Woo

Eden Belouadah, Marc Cenac, Guillaume Jarry, Enguerrand Paquin, Xunyi Zhao, Viktoriya Zhukova

Othmane Abou-Amal, Ameet Talwalkar, David Asker

What changed from Toto 1.0?

Architecture

Contiguous Patch Masking (CPM)

Training



Inference



Inputs



Toto 2.0

Robust Causal Scaler

Patch Embedding + CPM

Input Residual MLP

Variate-Time Transformer Decoder

Output Residual MLP

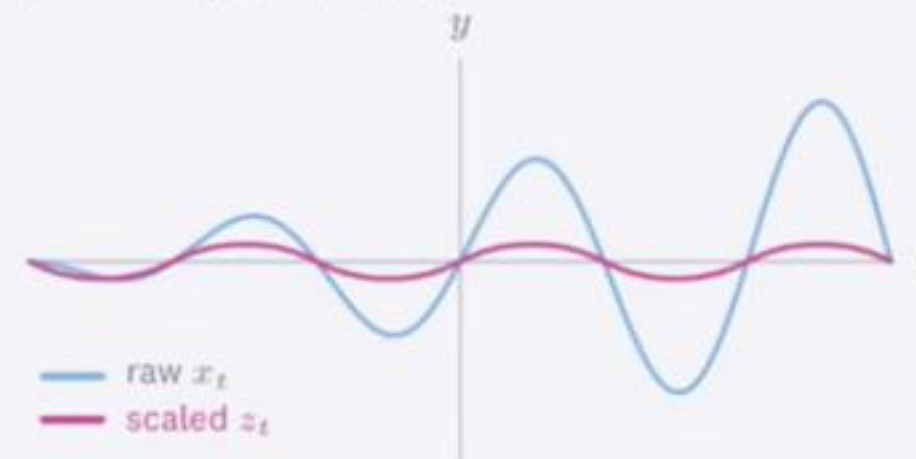
Quantile Knots Output Head

Forecasts

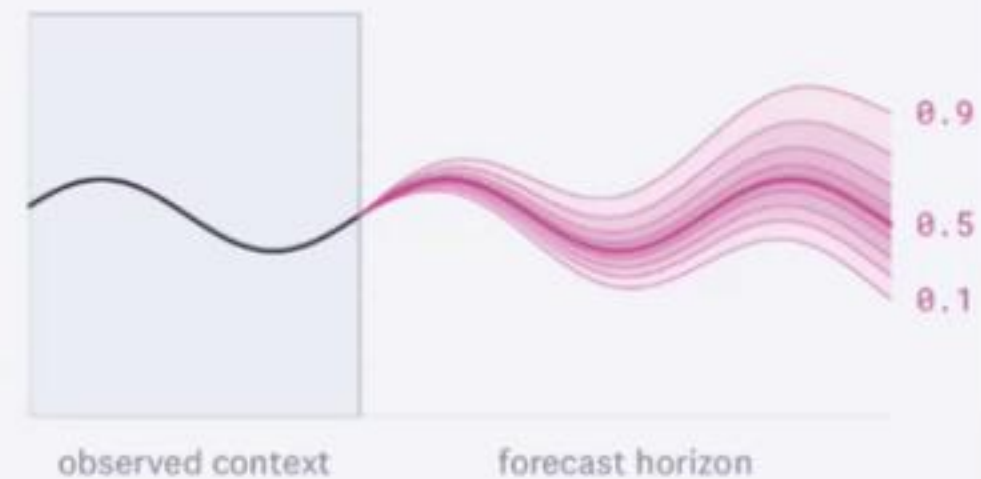


Robust Causal Scaler

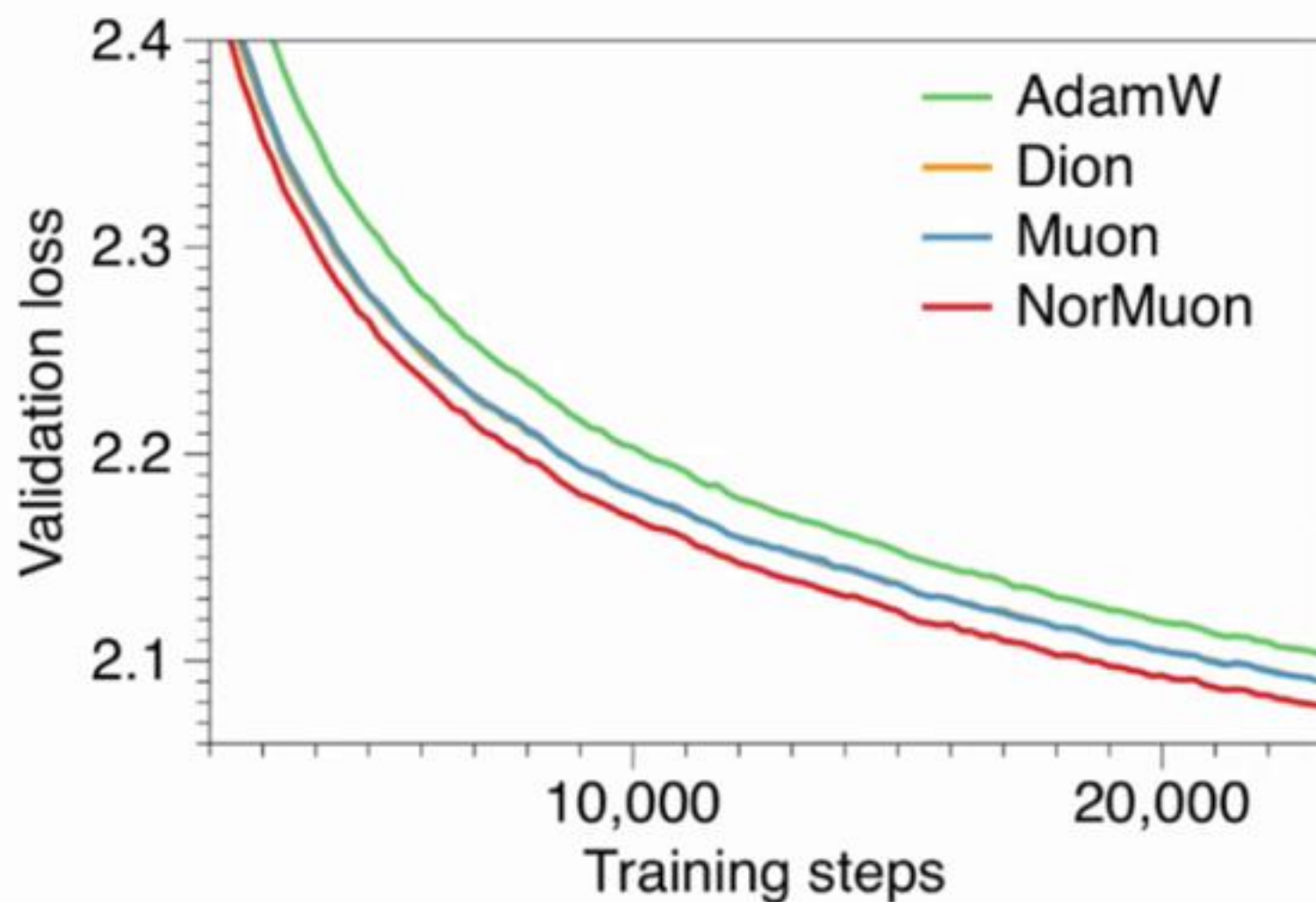
$$z_t = \text{asinh}((x_t - \mu_t) / \sigma_t)$$



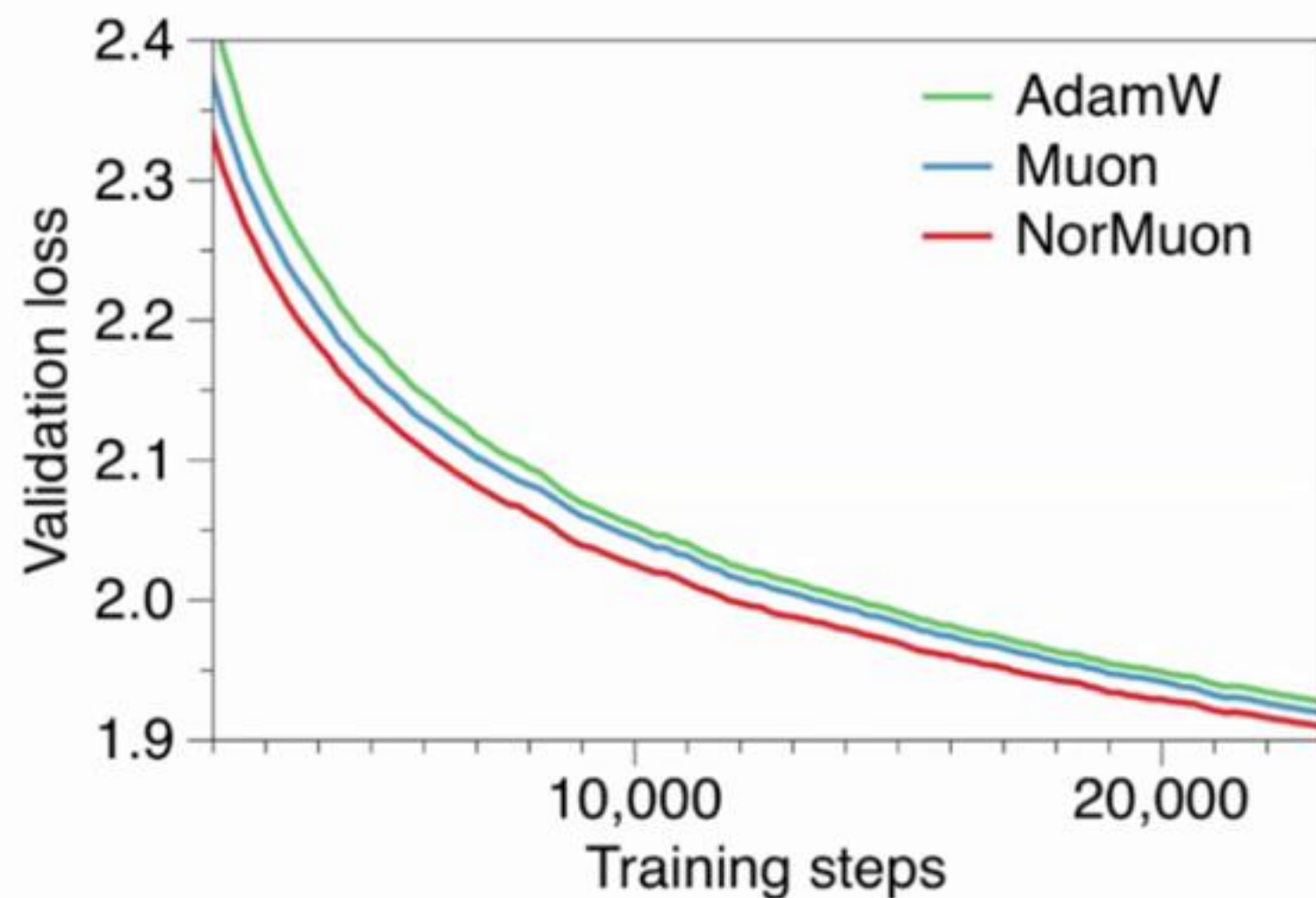
Quantile Knots Head



Optimizer: Normuon + AdamW



(a) Pretraining results of 1.1B model.



(b) Pretraining results of 5.4B model.

Beating GPT-2 for <<\$100: the nanochat journey #481

karpathy announced in Announcements



karpathy

on Jan 31

Maintainer

edited



When OpenAI released GPT-2 in February 2019, training the largest model (1.5B parameters) required serious compute:

- **Hardware:** 32 TPU v3 chips (256 TPU v3 cores, 8 cores per chip)
- **Training time:** "A bit over a week" (~168 hours)
- **Cloud cost:** At \$8/hour per TPU v3, that's $32 \times 168 \times \$8 = \$43,000$

Sources: [Reddit thread from 2019](#), [HuggingFace model card](#).

Beating GPT-2 for <\$100 from scratch has been a bit of an odd obsession for me but finally here we are. Seven years later, we can beat GPT-2's performance in nanochat ~1000 lines of code running on a single 8XH100 GPU node for ~3 hours. At ~\$24/hour for an 8xH100 node, that's **\$73**, i.e. **~600x cost reduction**. That is, each year the cost to train GPT-2 is falling to approximately 40% of the previous year. (I think this is an underestimate and that further improvements are still quite possible). The gains come from everywhere: better hardware (H100 vs TPU v3), better software (Flash Attention 3, torch.compile), better algorithms (Muon optimizer, architectural improvements), and better data (FineWeb-edu).

val/bpb

val/bpb

val/bpb

core_metric

0.26

0.14

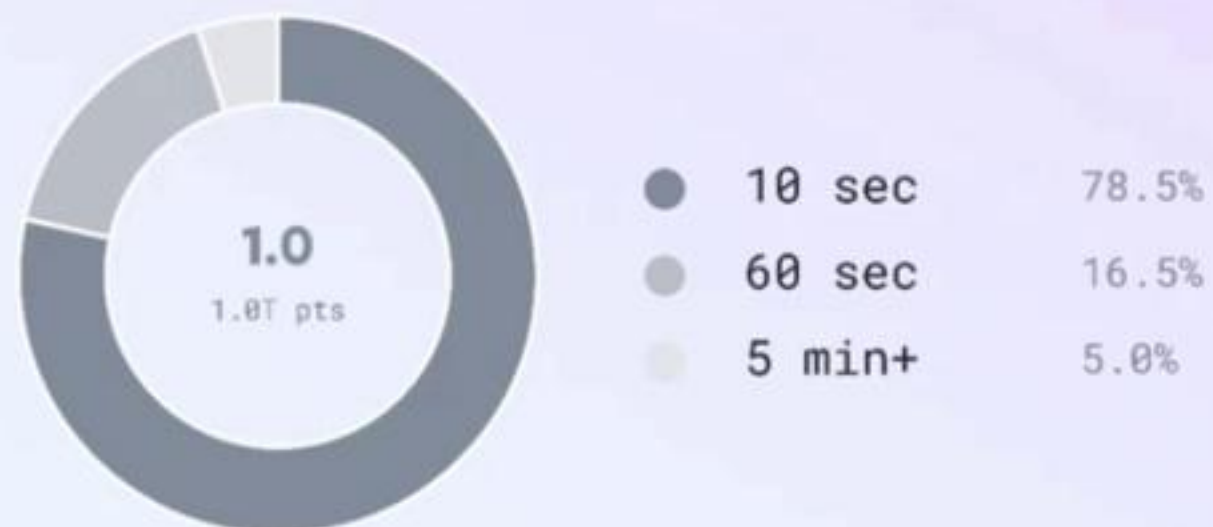
Toto 2.0 Training Data

Composition

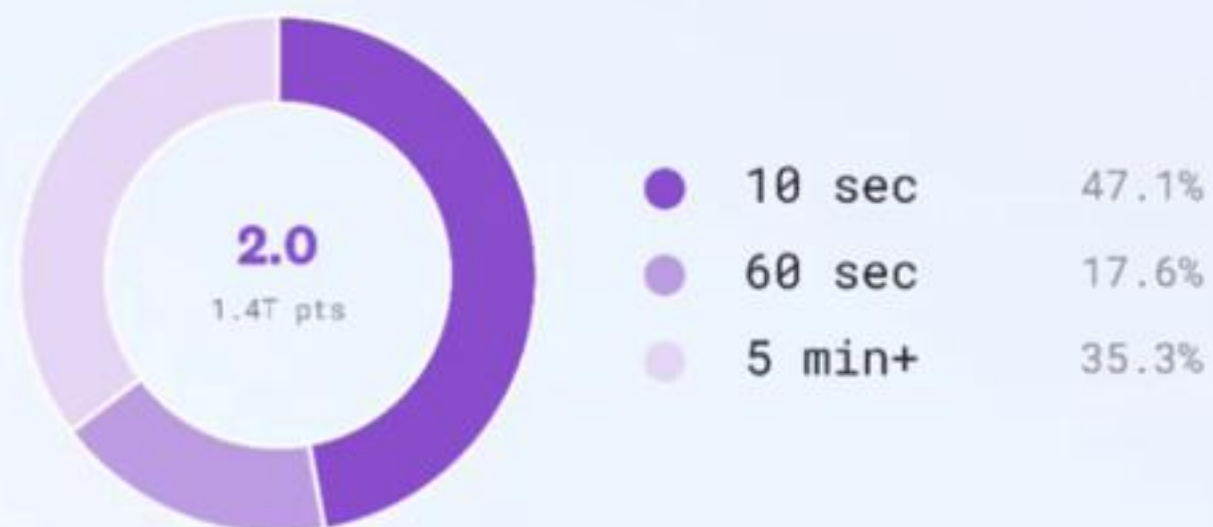
Toto 1.0



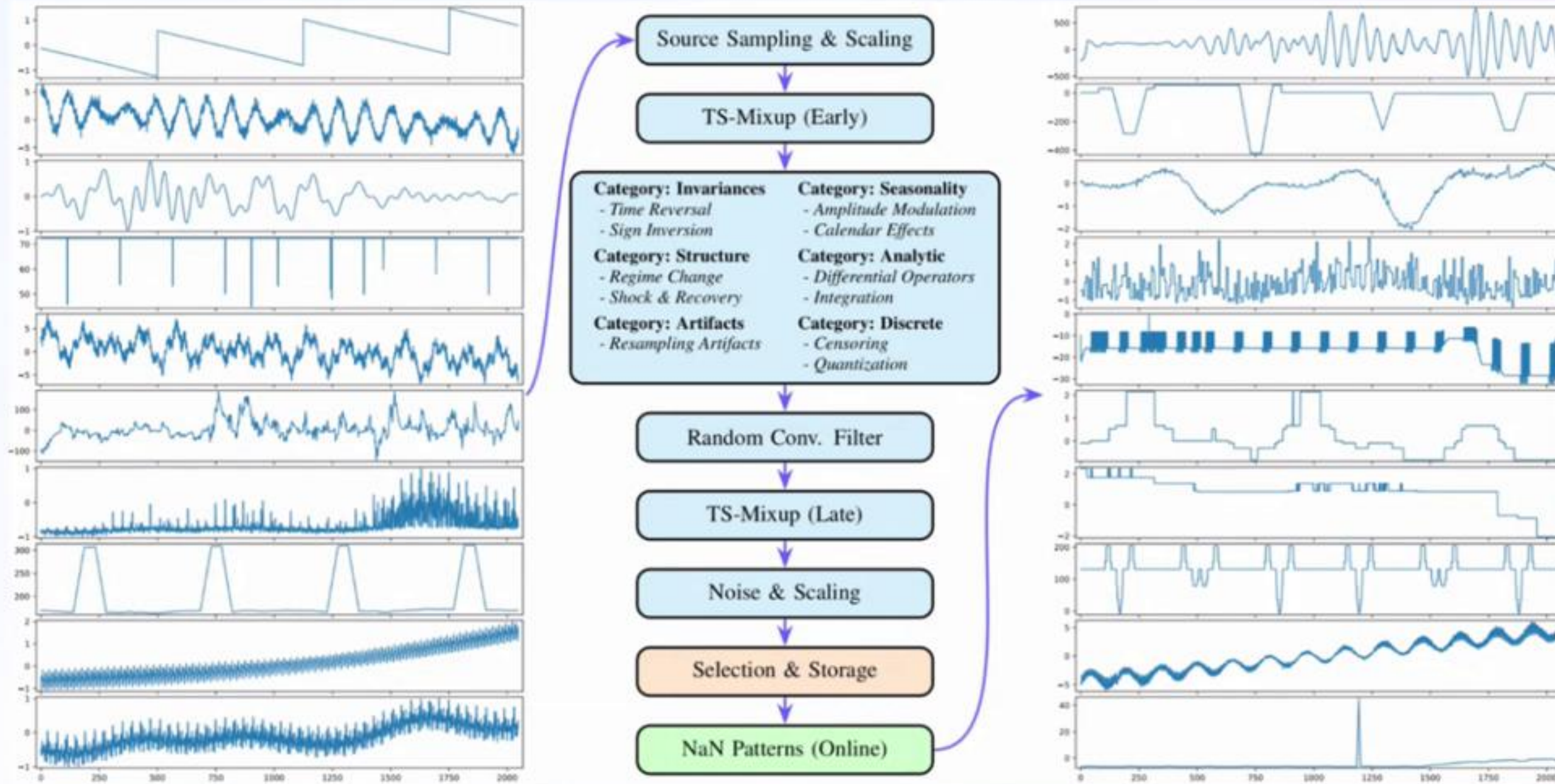
Interval



Toto 2.0



TempoPFN Synthetic Data



Scaling Efficiently





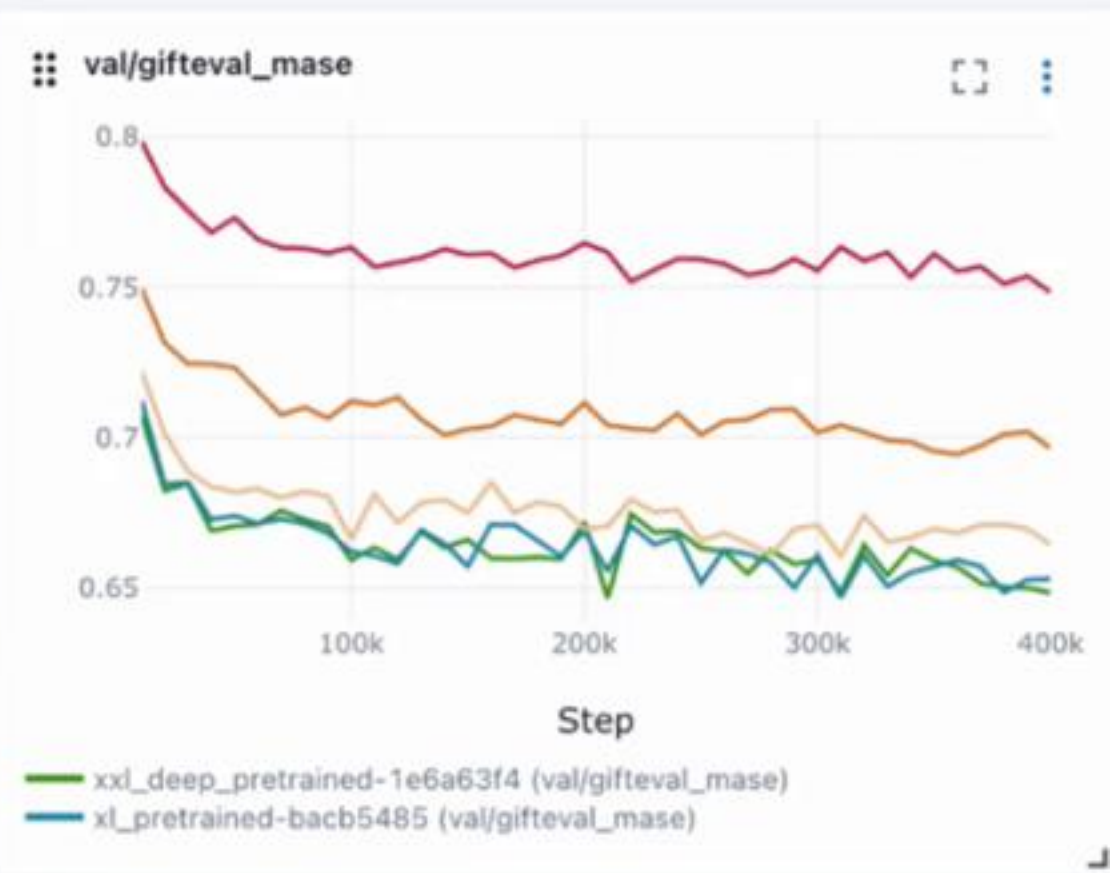
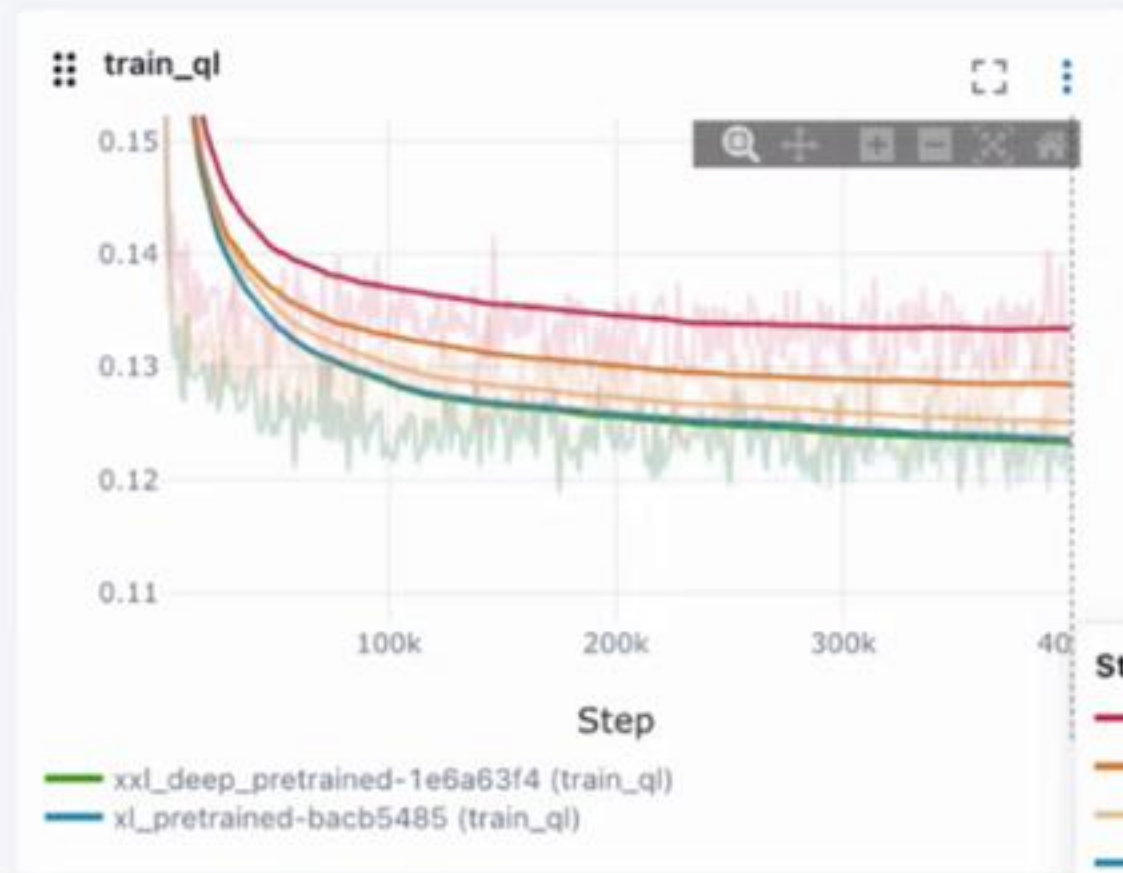
Searching for Model Parameters with REx



Zero shot scale up to 2.5B!

Model metrics (3)

+ Add chart



Step 398799

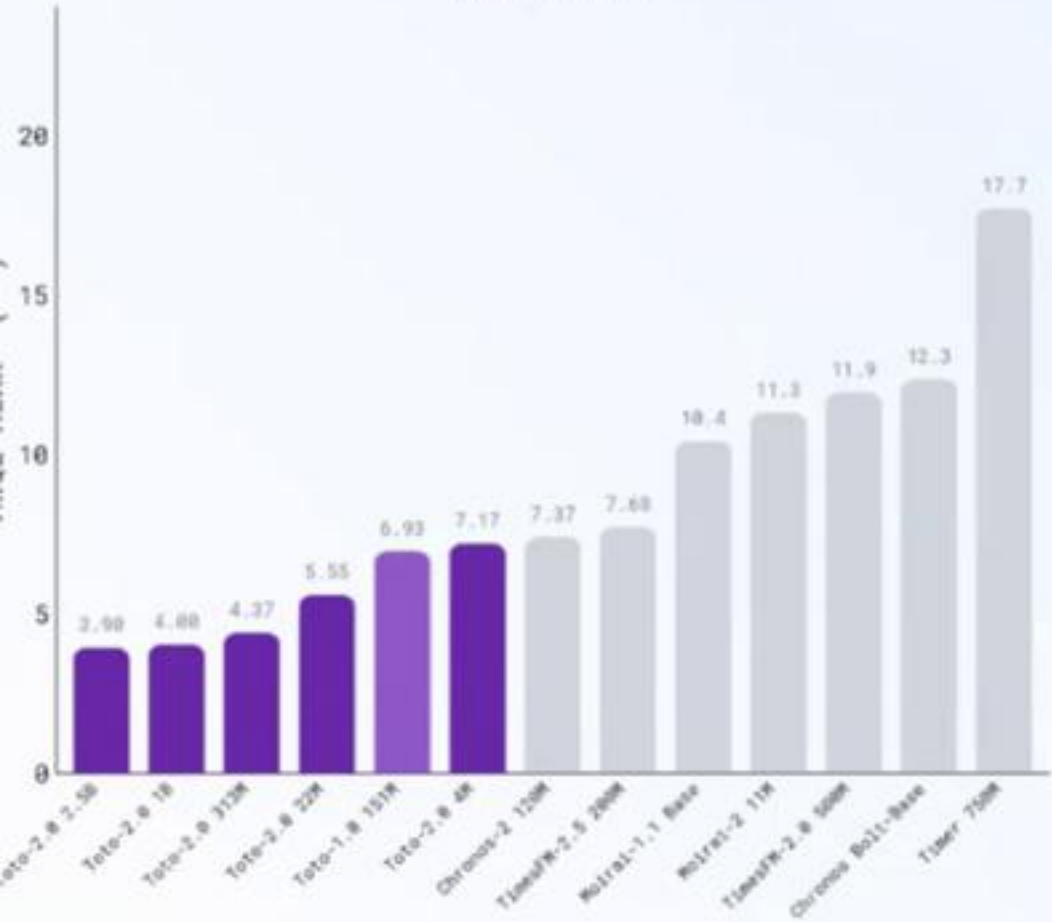
tiny_pretrained_ps32-b3448c57	0.13336994280526543
small_pretrained_ps32-e309c867	0.12845698537834563
large_pretrained-95fb2190	0.1250865281509384
xl_pretrained-bacb5485	0.12364319184133438
xxl_deep_pretrained-1e6a63f4	0.12329482017193934

System metrics (0)

Results

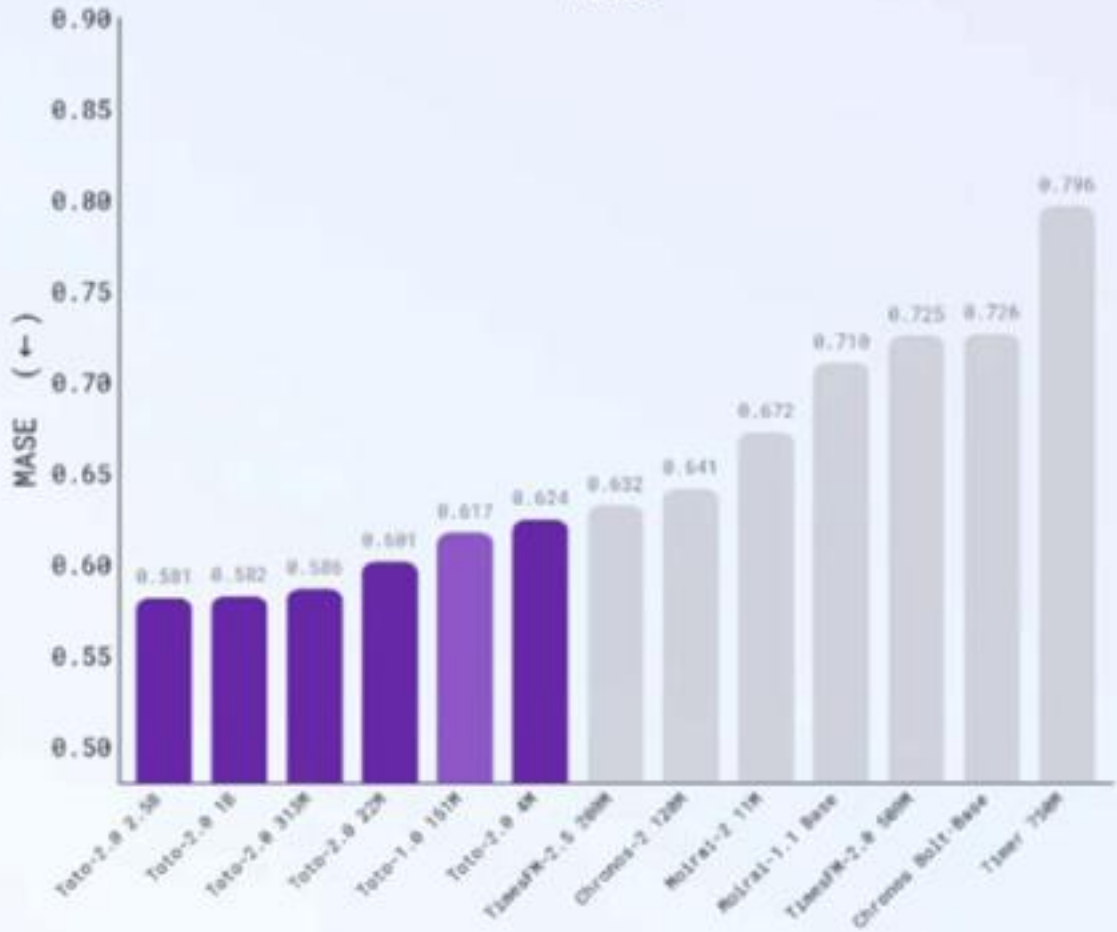
BOOM

MWQL RANK



BOOM

MASE



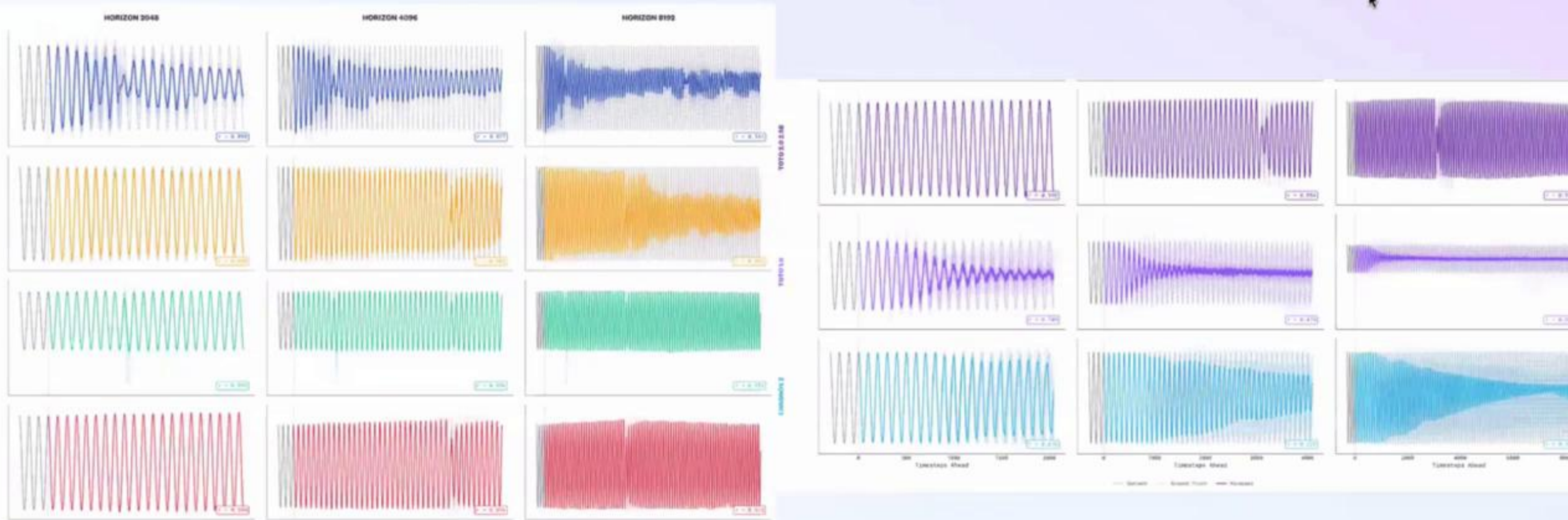
MWQL



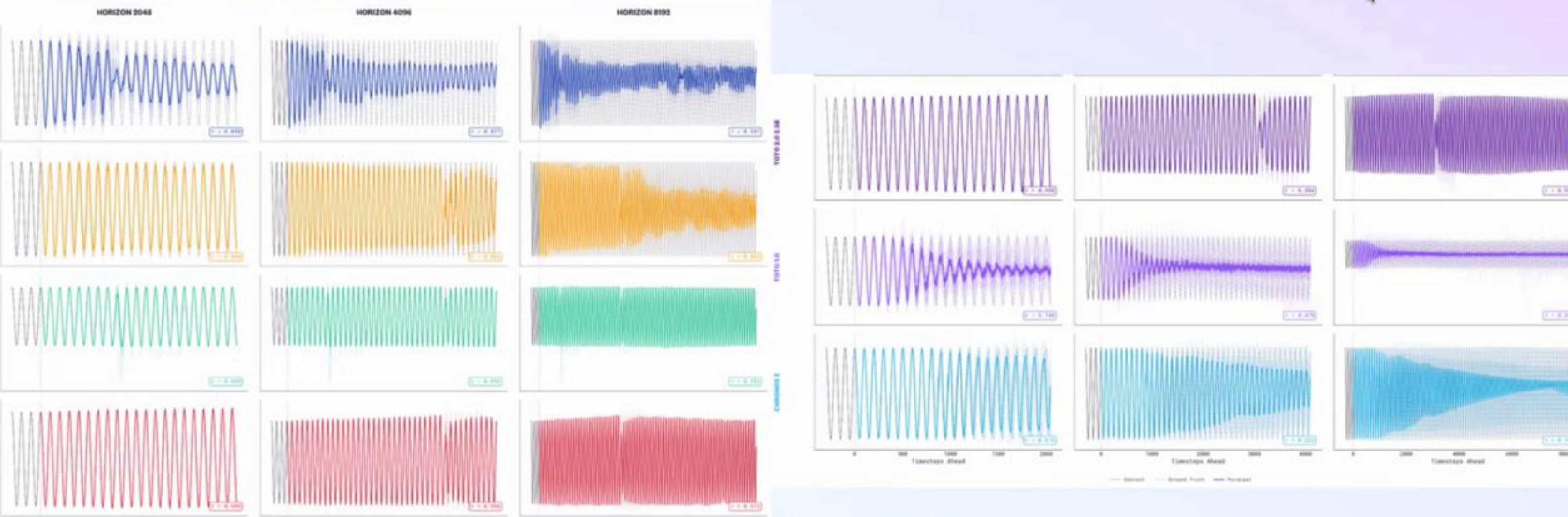
GIFT-Eval (update in progress)

T	model ▲	Organization ▲	Test Leak. ▲	Replication Code	MASE ▲	MASE_Rank ▲	CRPS ▲	CRPS_Rank ▲
●	Toto-2.0-2.5B (code)	Datadog	No	Yes	0.698	19.041	0.481	18.031
●	Toto-2.0-1B (code)	Datadog	No	Yes	0.700	19.701	0.48	18.546
●	Chronos-2 (code)	AWS	No	Yes	0.698	19.928	0.485	20.093
●	Toto-2.0-313m (code)	Datadog	No	Yes	0.701	21.474	0.486	19.918
●	Timer-S1 (code)	Tsinghua & Bytel	No	Yes	0.693	24.052	0.485	24.01
●	PatchTST-FM-r1 (code)	IBM TSFM & Renss	No	Yes	0.707	24.175	0.483	20.68
●	FlowState-r1.1 (code)	IBM TSFM	No	Yes	0.701	24.227	0.487	22.918
●	TimesFM-2.5 (code)	Google Research	No	Yes	0.705	24.567	0.49	23.32
●	TiRex (code)	NXAI	No	Yes	0.716	26.742	0.488	21.412
●	Granite-PatchTST-FM-r1	IBM TSFM & Renss	No	Yes	0.717	27.866	0.488	23.124
●	Toto-2.0-22m (code)	Datadog	No	Yes	0.719	28.258	0.497	24.526
●	Chronos-2-Synth (code)	AWS	No	Yes	0.720	29.01	0.496	27.918

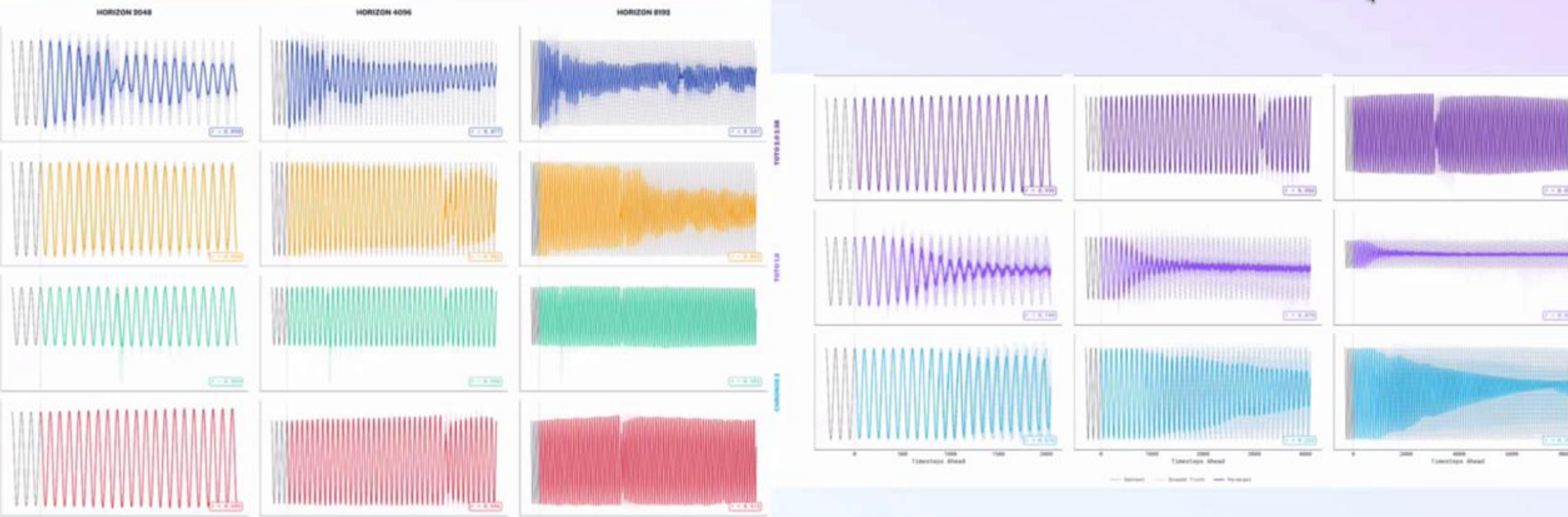
Discovery: Long Horizon Stability Scales with Size



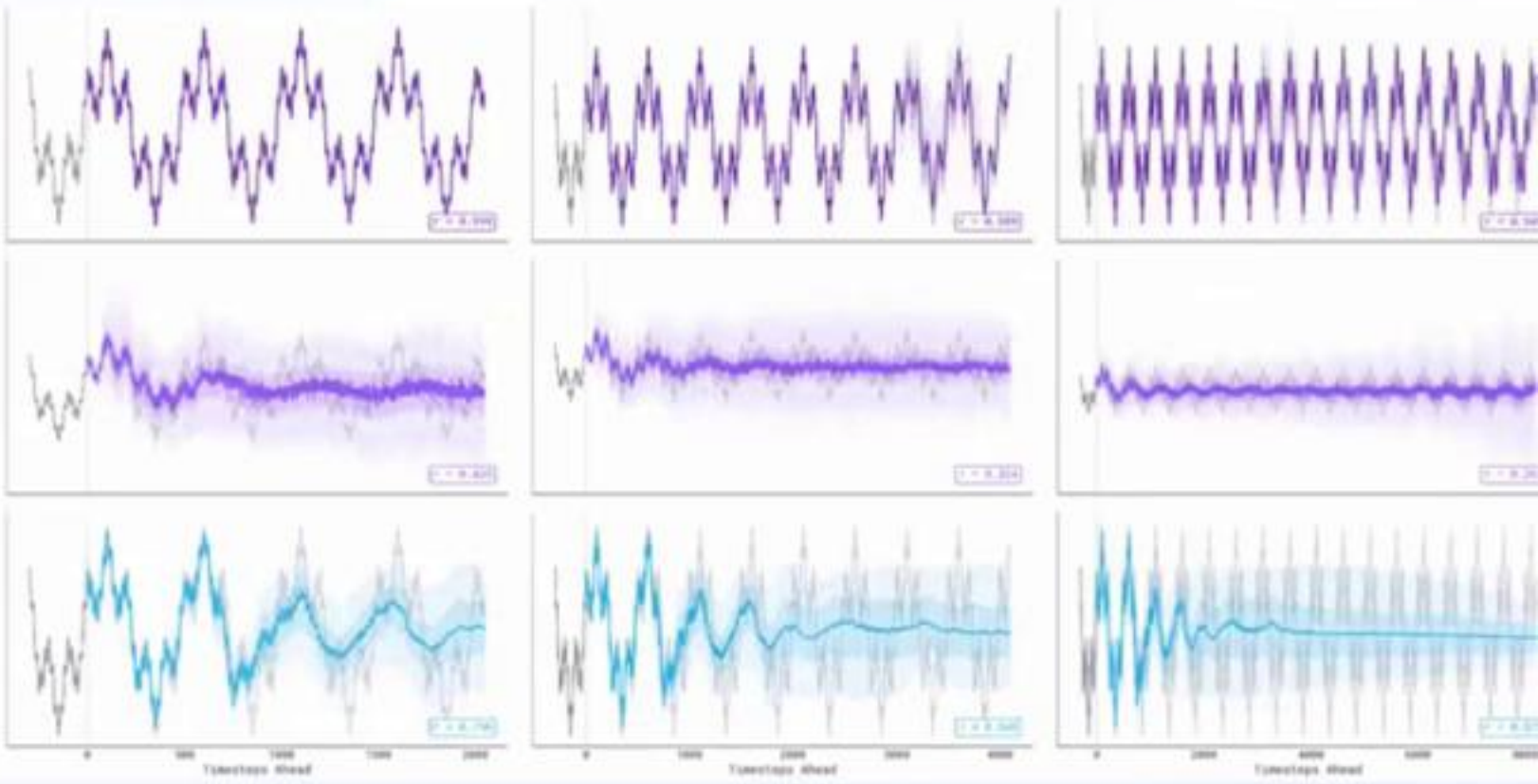
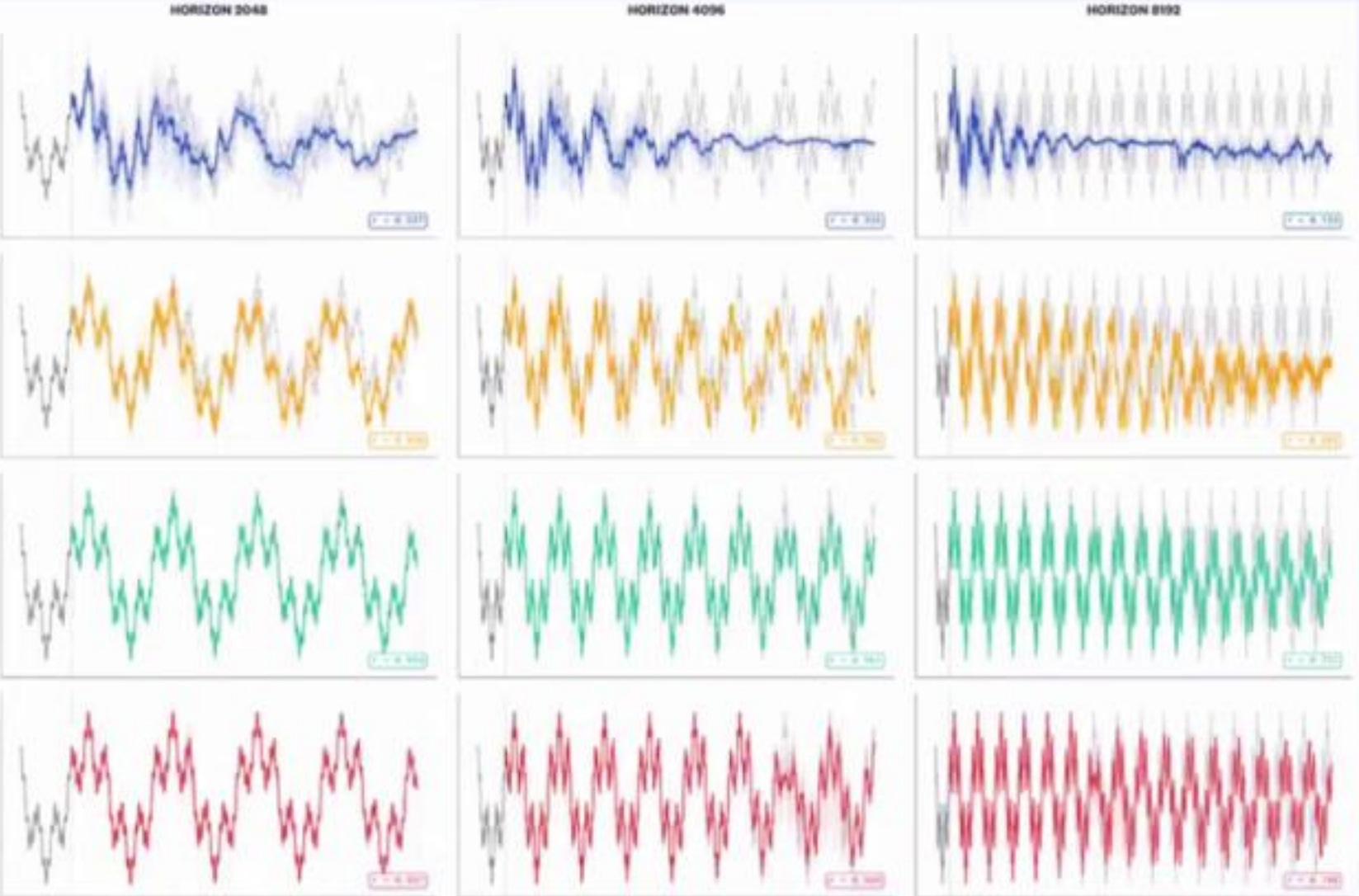
Discovery: Long Horizon Stability Scales with Size



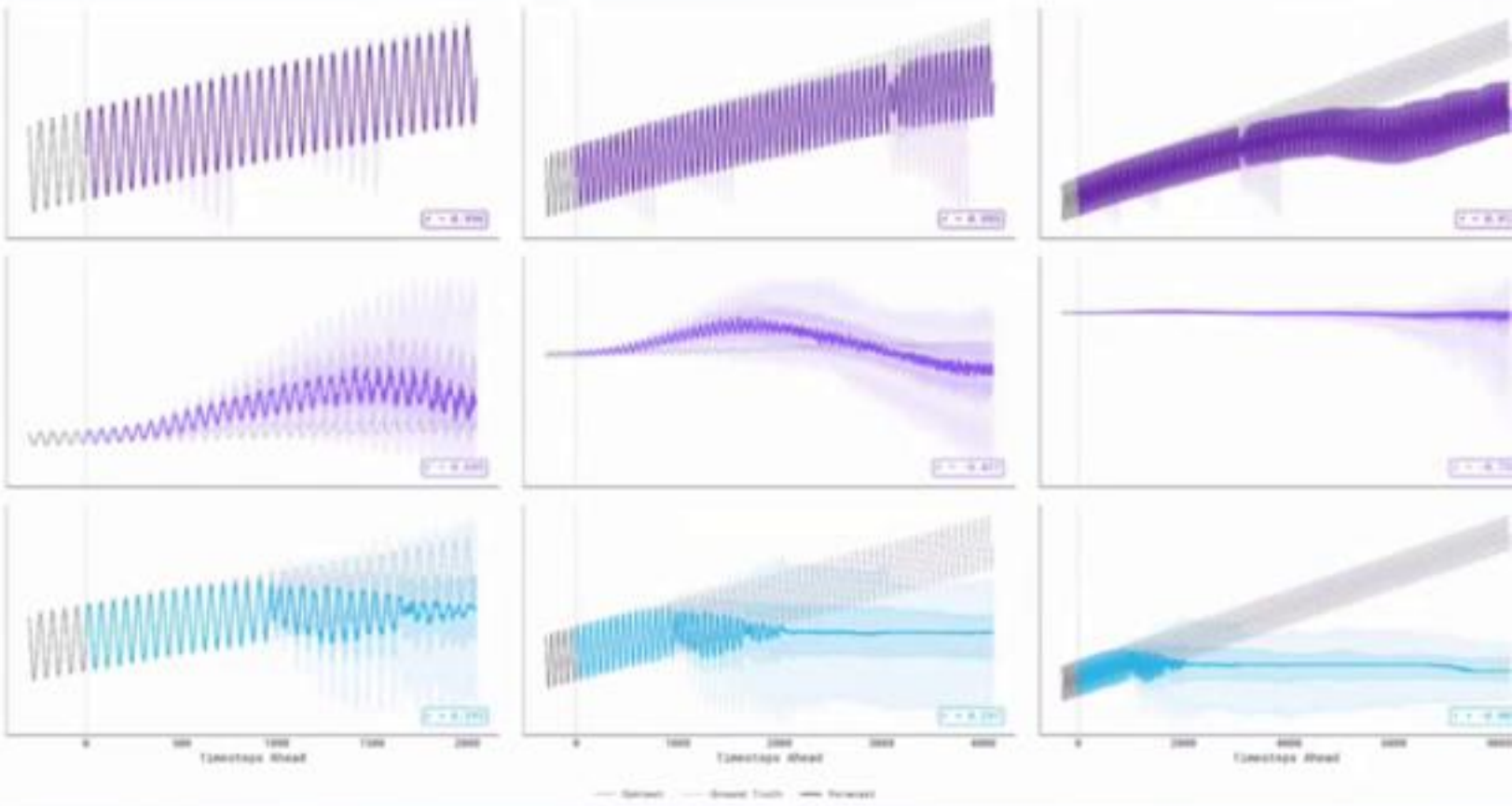
Discovery: Long Horizon Stability Scales with Size



Discovery: Long Horizon Stability Scales with Size

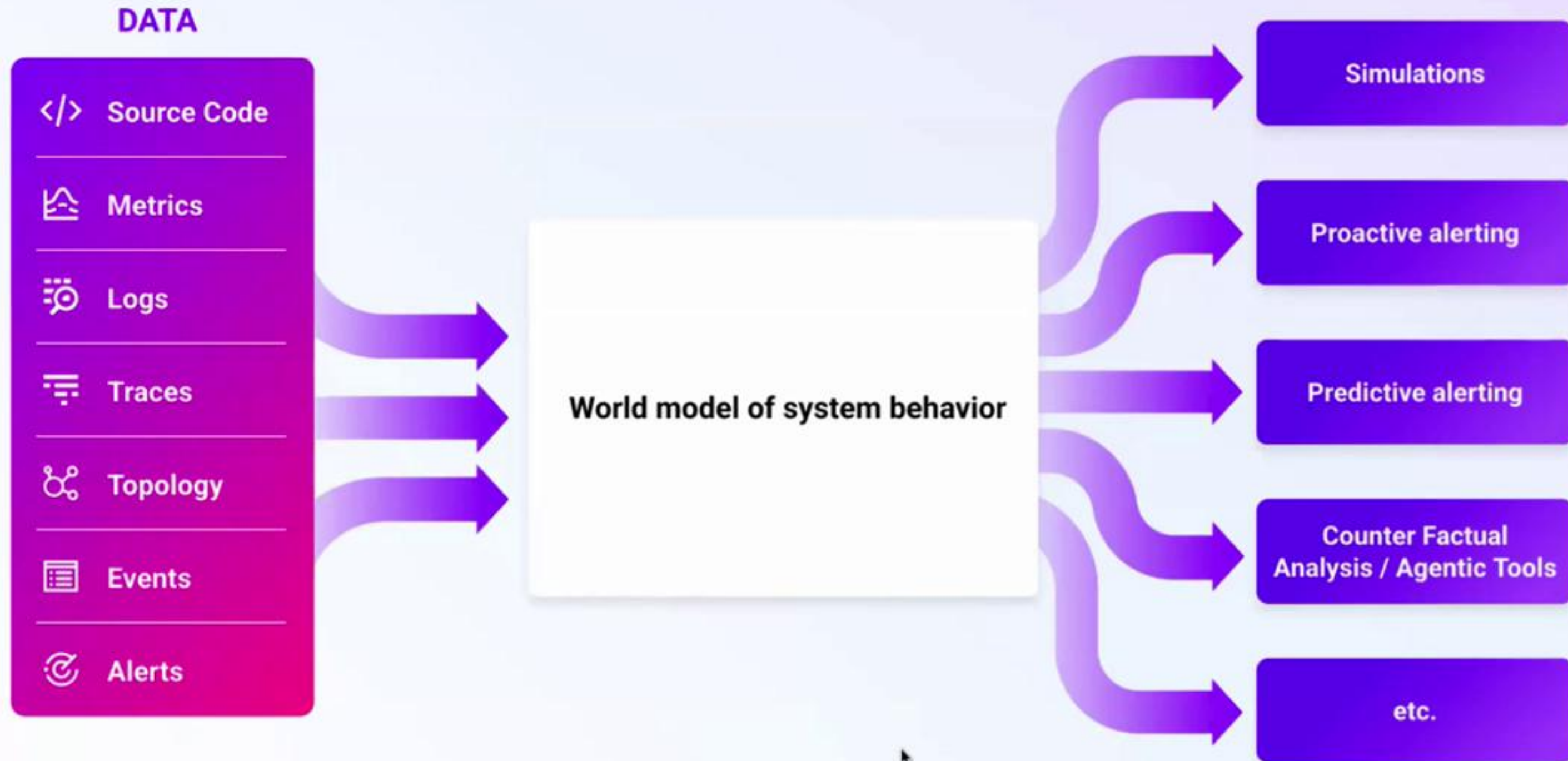


Discovery: Long Horizon Stability Scales with Size



From Time Series Foundation Models to Multi Modal “World Models”

Multi Modality & World Models



ARFBench: TSQA Benchmark based of internal incidents

ARFBench: Benchmarking Time Series Question Answering Ability for Software Incident Response

Stephan Xie^{1,2}, Ben Cohen², Mononito Goswami³, Junhong Shen¹,
Emaad Khwaja², Chenghao Liu², David Asker², Othmane Abou-Amal², Ameet Talwalkar^{1,2}

¹Machine Learning Department, Carnegie Mellon University, Pittsburgh, PA, USA

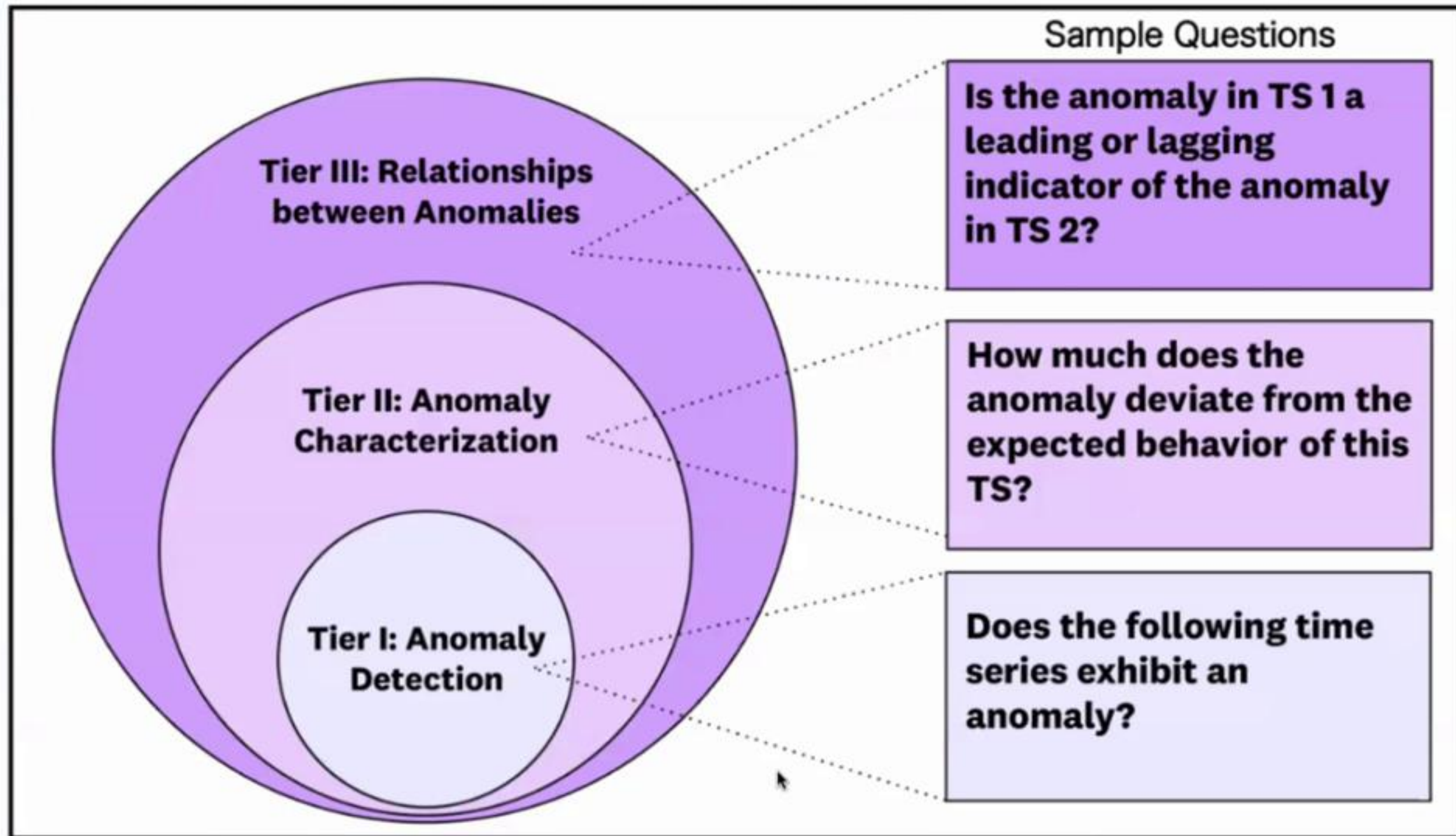
²Datadog AI Research, New York, NY, USA

³Amazon AI Research, Seattle, WA, USA

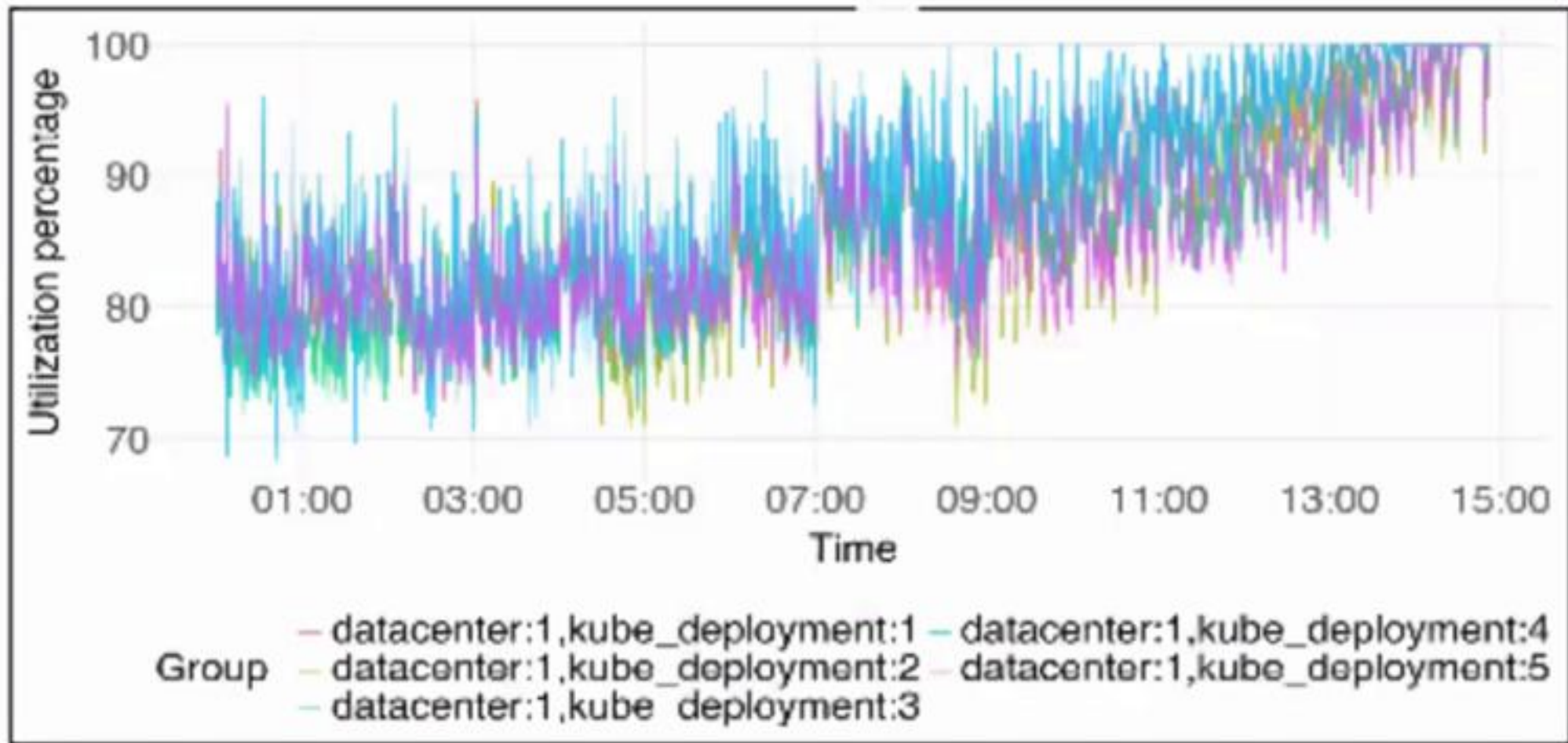
{stephan.xie, ameer.talwalkar}@datadoghq.com

ARFBench: TSQA Benchmark based of internal incidents

ARFBench: TSQA Benchmark based of internal incidents

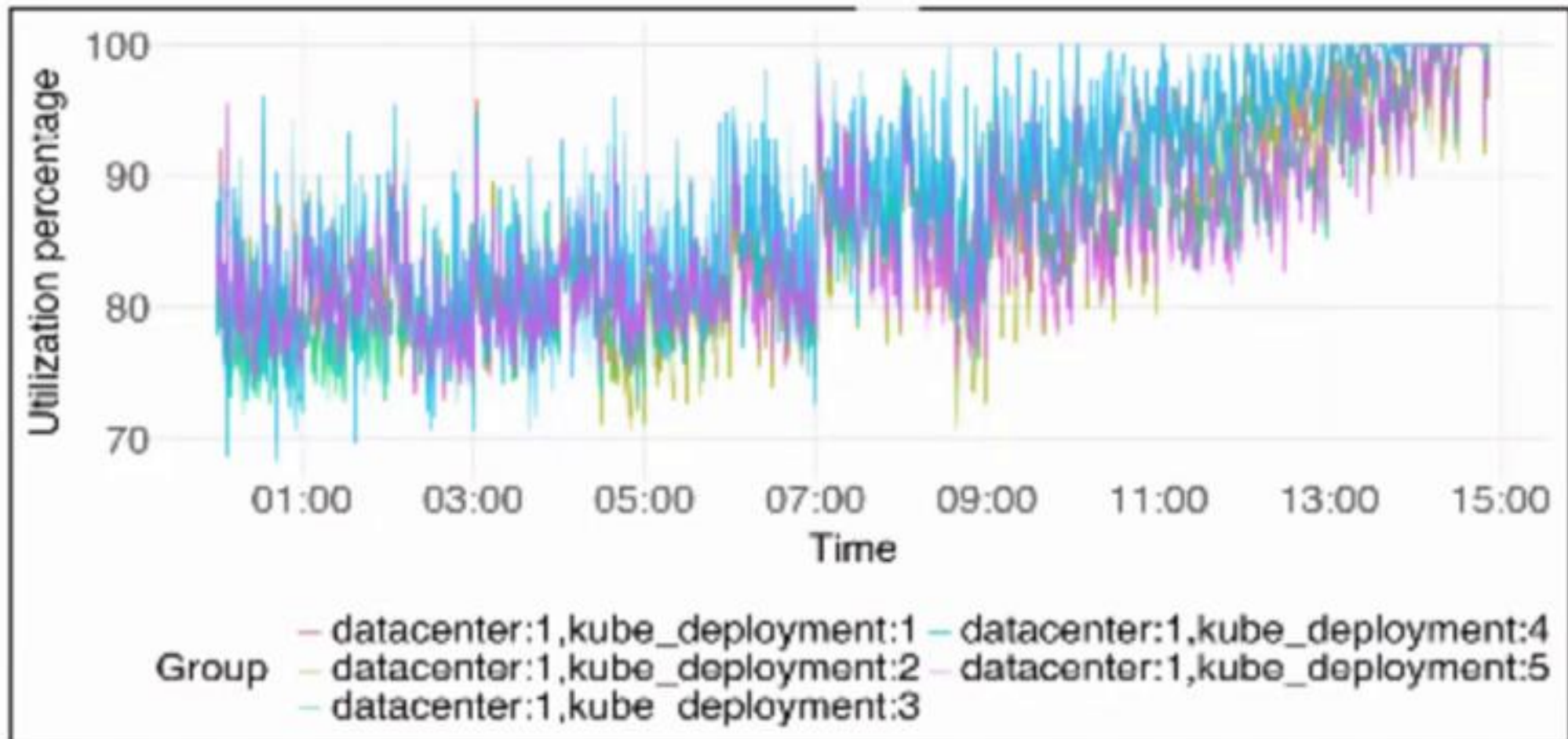


Time series question answering (TSQA)



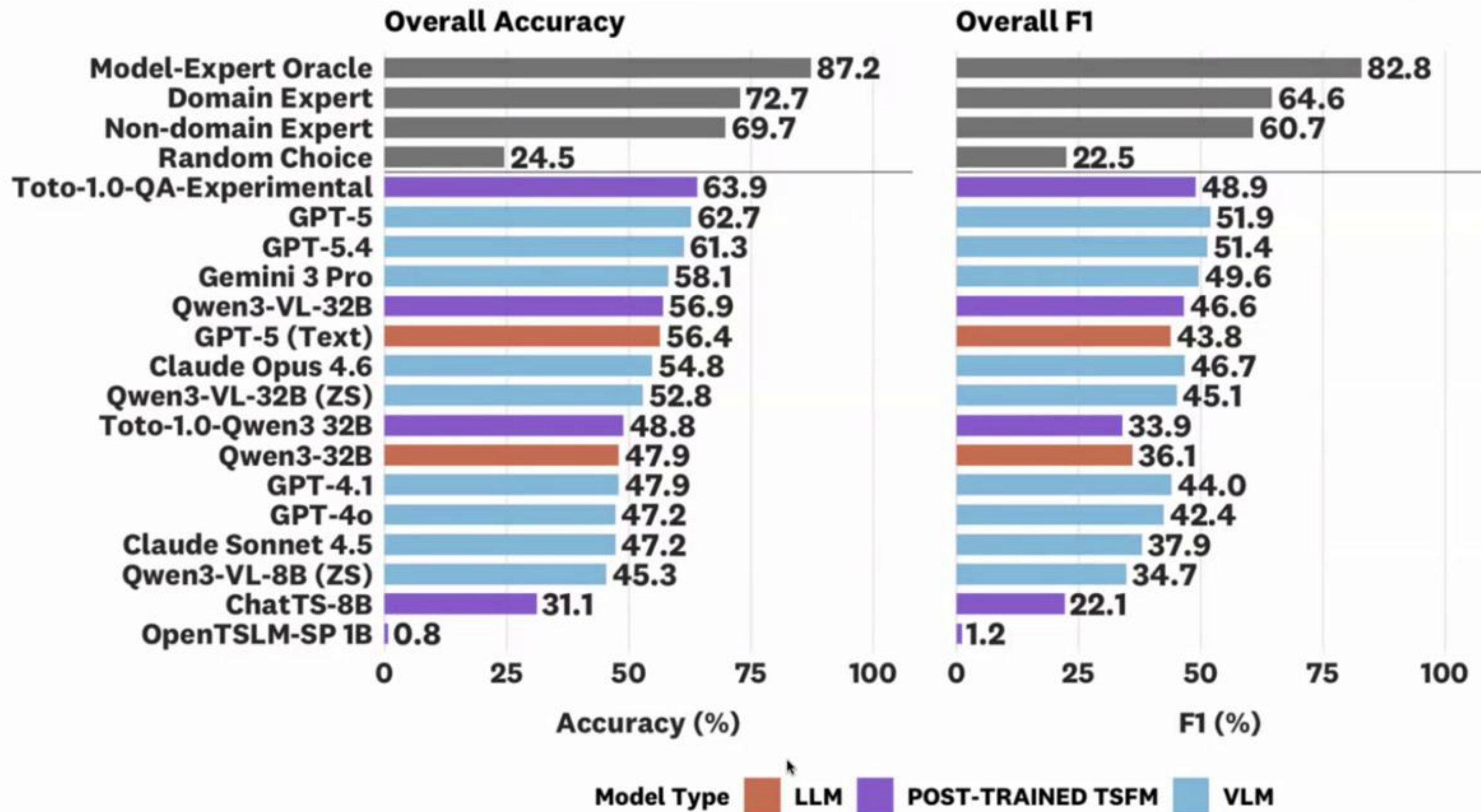
Avg utilization % for a metadata queue,
grouped by k8s deployment and datacenter

Time series question answering (TSQA)

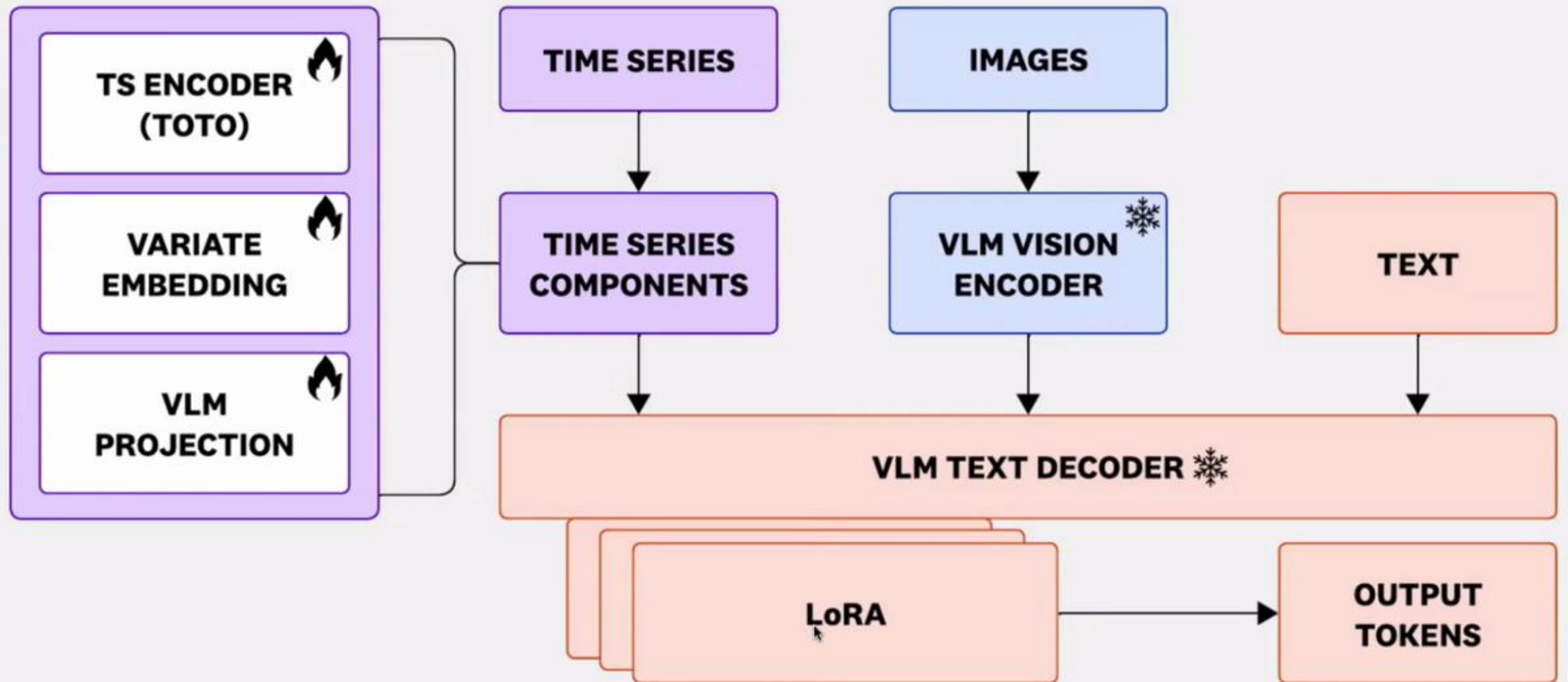


Is there an anomaly in this time series?

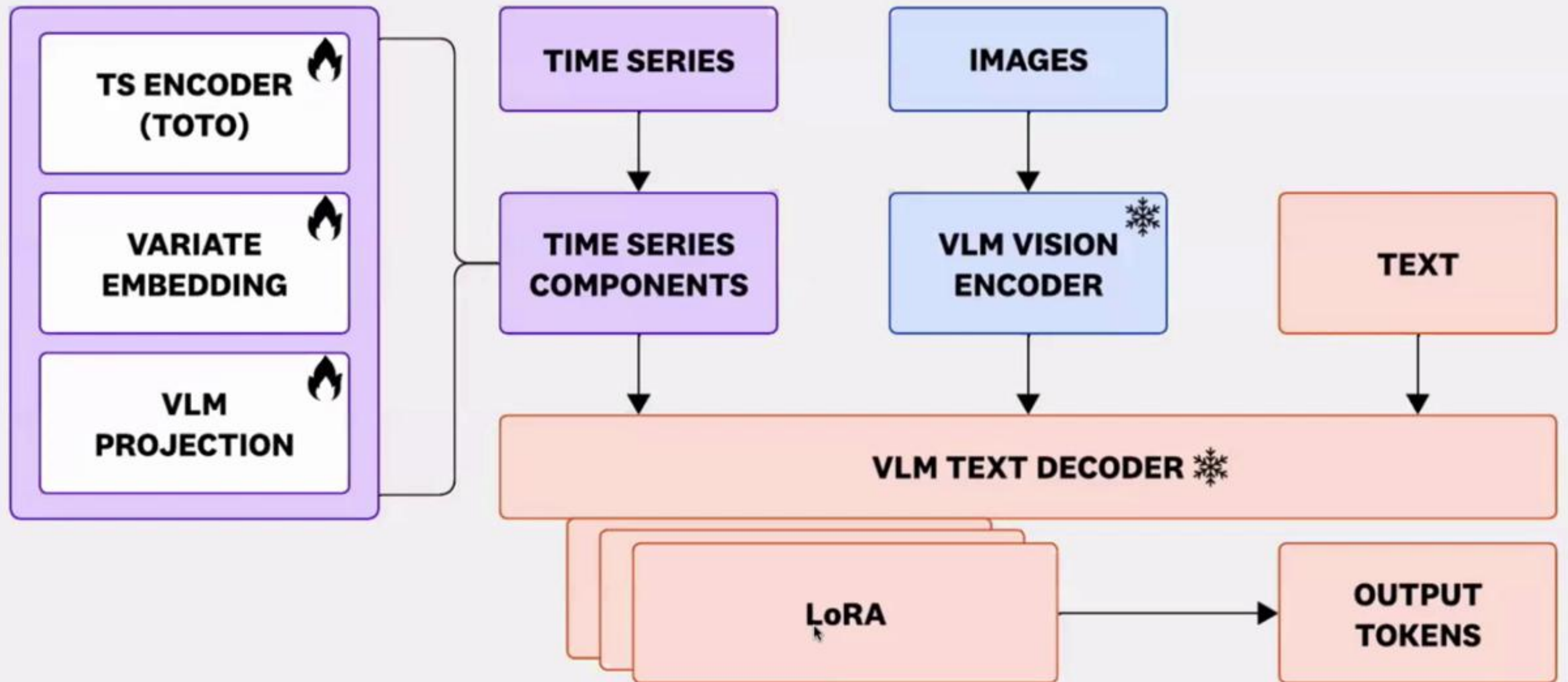
Avg utilization % for a metadata queue, grouped by k8s deployment and datacenter



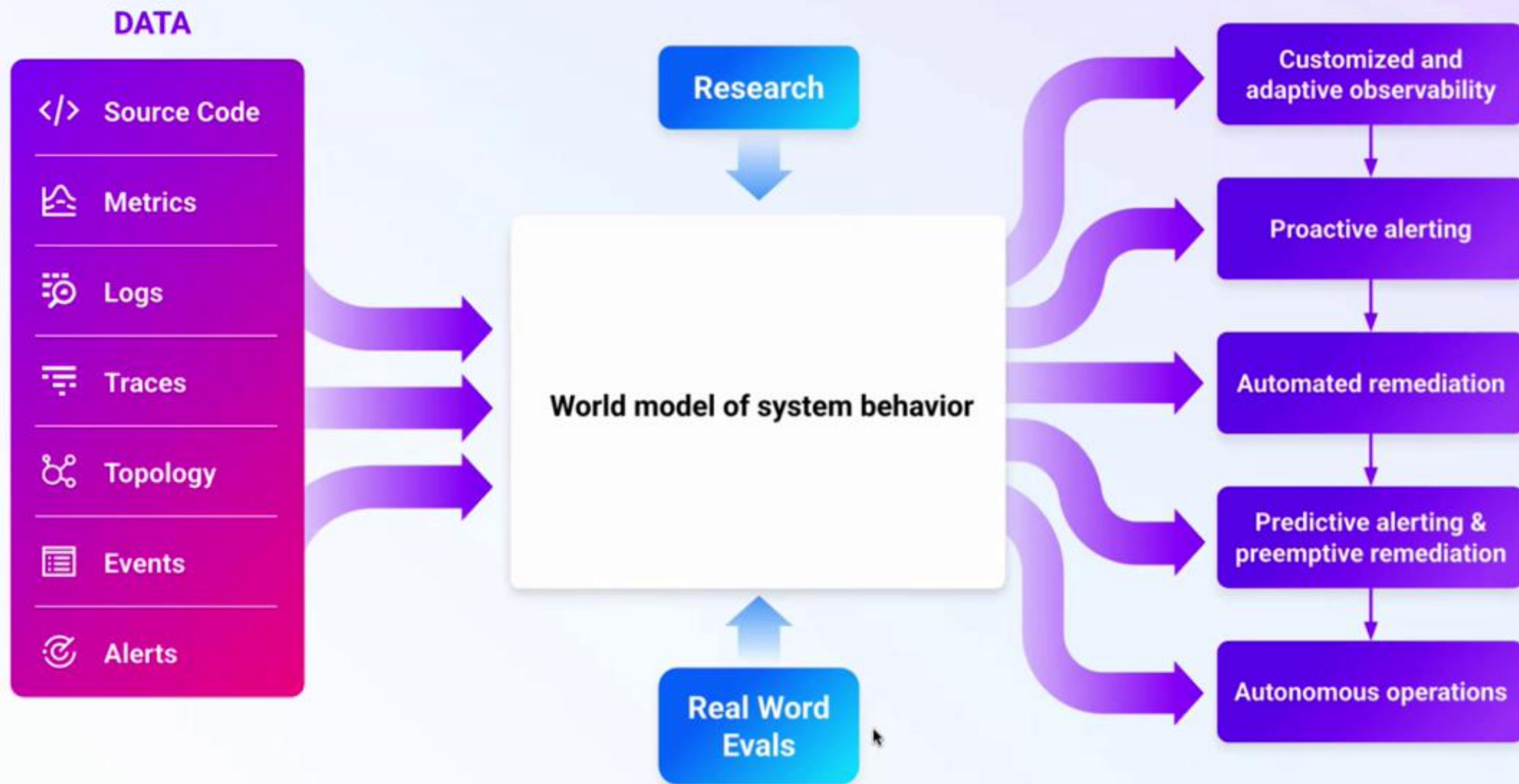
Toto-1.0-QA-Experimental architecture



Toto-1.0-QA-Experimental architecture



Towards autonomous operations



BERT Moment For Time Series

BERT Moment For
Time Series

GPT-2 Moment For
Time Series

BERT Moment For
Time Series

GPT-2 Moment For
Time Series

Multimodal World
Models for o11y

Datadog AI Research

We are hiring!



Sesh
(VP)



David
(EM II)



Afshin
(Research
Scientist)



Sam
(Research
Scientist)



Bogna
(Engineer)



Arun
(Engineer)



Rocco
(Recruiter)



ICLR

International Conference On
Learning Representations