# Vertical AI Meets Observability

Ameet Talwalkar

12/7/25

Carnegie Mellon University

DATADOG

# AI has taken over the world

Frontpage news

Startup valuations

Big tech spending

AI research pace



THE WALL STREET JOURNAL.

DOW JONES | News Corp ★ ★ ★ ★ ★ ★  TUESDAY, JANUARY 28, 2025 ~ VOL. CCLXXXV NO. 22  WSJ.com  ★ ★ ★ ★ $5.00

DJIA 44713.58 ▲ 289.33 0.65% | NASDAQ 19341.83 ▼ 3.1% | STOXX 600 529.69 ▼ 0.1% | 10-YR. TREAS. ▲ 24/32, yield 4.529% | OIL $73.17 ▼ $1.49 | GOLD $2,737.50 ▼ $39.80 | EURO $1.0493 | YEN 154.51

## What's News

**Business & Finance**

◆ **Financial markets** swooned at the emergence of a dark-horse power in artificial intelligence, which sent shares of Nvidia down 17% and posed a fresh threat to the multitrillion-dollar boom in the U.S. tech sector. The S&P 500 and Nasdaq slid 1.5% and 3.1%, respectively, while the Dow rose 0.7%. **A1, A4**

◆ **The Senate confirmed** Scott Bessent as treasury secretary, putting the longtime investor at the center of Trump's efforts to cut taxes, fight inflation and im-

## DeepSeek Flips Script on AI

Chinese dark horse emerges, threatening a market darling and other big tech stocks

*By Gunjan Banerji, Asa Fitch and Karen Langley*

For two years, markets' belief that the rise of artificial intelligence would usher in a new era of productivity growth has fueled trillions of dollars in stock-market gains.

Nvidia, the maker of the computer chips at the heart of

the AI boom, has been in the vanguard of this advance. Wall Street has perceived the company to have an almost unreachable defense against competition with its offerings of high-tech chips. The company's rapid growth and windfall profits have helped push other technology firms and the Nasdaq Composite Index to record after record, with giddy investors expecting more of the same down the road.
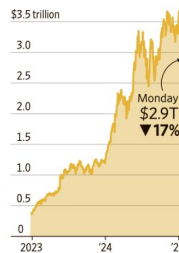
On Monday, the mood turned sour. DeepSeek, a dark-horse power in artificial intelligence, emerged from China. That rattled big tech stocks, led by a plunge of almost $600 billion in

Nvidia, which only last week was the world's most valuable company. Nvidia's fall marked the largest one-day loss in market value for any public company.

DeepSeek released last week an AI model that appeared to perform on par with a cutting-edge counterpart from OpenAI, the U.S. firm at the heart of the AI craze. The twist: Creative engineering tricks meant DeepSeek needed far less computing power. The upshot is that the AI models of the future might not require as many high-end Nvidia chips as investors have counted on.

*Please turn to page A4*

**Nvidia market value**

$3.5 trillion

Monday
$2.9T
▼17%

2023  '24  '25

Source: FactSet

## Market Plunges As China Firm Stirs Worries

Fresh threat to AI in the U.S. wipes out about a trillion dollars from stock market

Financial markets swooned on Monday at the emergence of
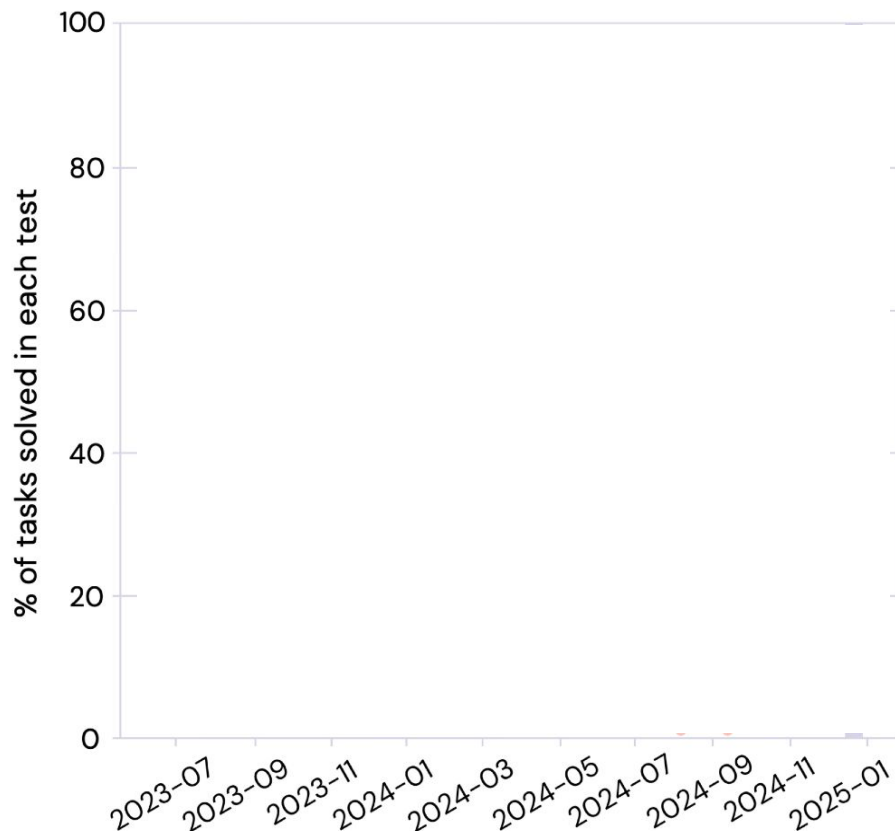
# AI has taken over the world?



**Hype**

Credit: GPT-4o

*'AI eating the world'*

!=

**Impact**

Search

Coding

Why?

# Hypothesis 1: Just wait, it's coming



% of tasks solved in each test

- FrontierMath: Advanced mathematics
- ARC-AGI: Abstract reasoning (semi-secret evaluation)
- SWE-bench: Real-world software engineering
- GPQA: Graduate-level science
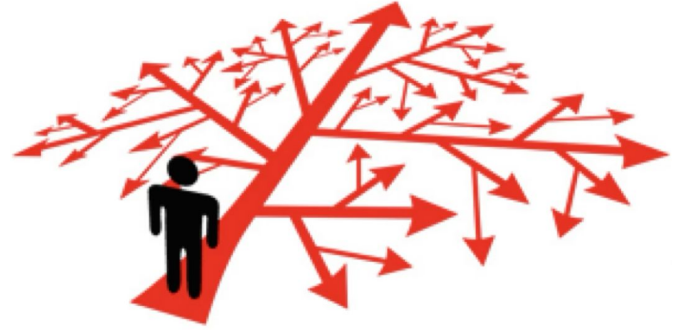- AIME 2024: Mathematics competition for elite students

Source: International AI Safety Report, Jan 2025

# Hypothesis 2: Specialization matters!

## People specialize to become experts
- E.g., scientists and athletes
- E.g., the human brain itself

## AI will need to specialize
- Accuracy
- Efficiency

# It's not either/or

Vertical

General

# How good are today's specialized FMs?

## Toto: An Observability TSFM

**Before**: Train directly on supervised data

**BERT Moment**: Pre-train on massive corpora, then fine-tune

**After**: Nobody uses supervised learning alone

2012

2018

**NLP Timeline**

# How good are today's Specialized FMs?

[Gupta*-Xu*-Cheng-Shen-Shen-**T**-Khodak, ICLR25]

**Genomics**

**Satellite Imaging**

**Time Series**

NLP (2018 - 2019)

# How good are today's Specialized FMs?

[Gupta*-Xu*-Cheng-Shen-Shen-**T**-Khodak, ICLR25]

**Genomics**

**Satellite Imaging**

**Time Series**

Don't assume specialized FMs work!

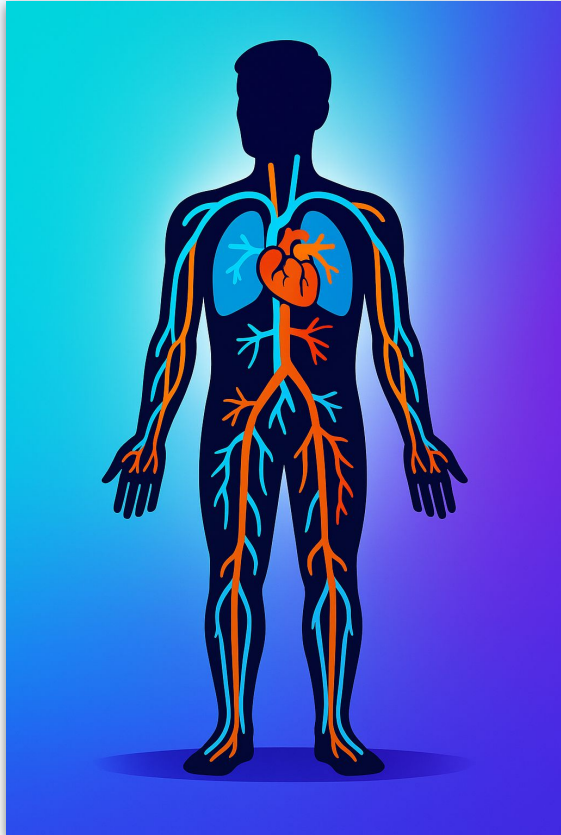Baselines & benchmarks are important

# How good are today's specialized FMs?
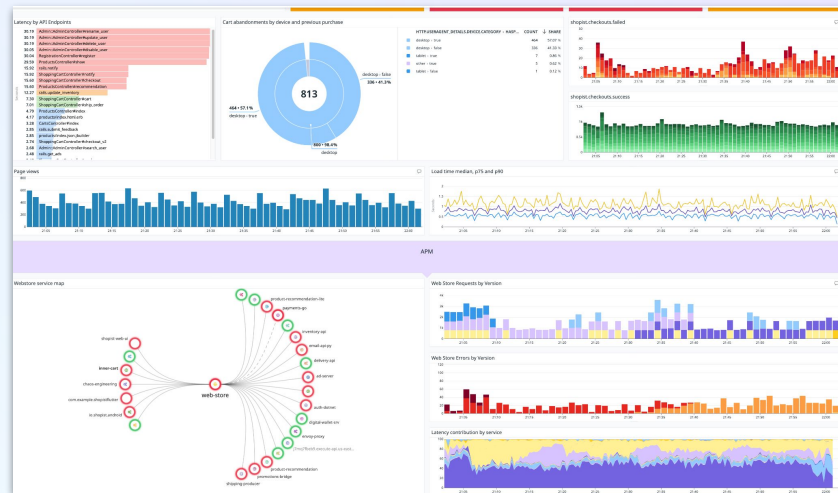
## Toto: An Observability TSFM

Datadog AI Research

# What is
# **Observability?**

# *Observing/monitoring:* The Human Body



*AI-generated Illustrations*

DATADOG

# *Observing/monitoring:* Computer Systems



*AI-generated Illustrations*

# 1,000s of hosts, pods, containers, etc.

DATADOG

# Trillions

of data points/hour

19

DATADOG

# Types of Observability Data

**Telemetry Data**
- Metrics
- Logs
- Traces
- Network Flows
- Source Code
- Cloud Cost
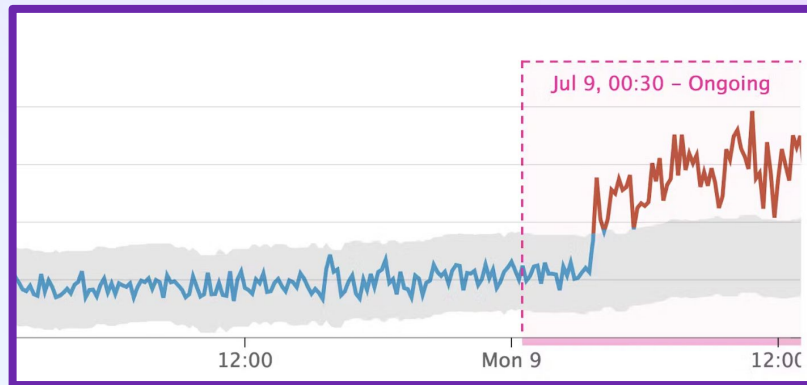- CI / CD Pipelines
- Security Signals
- …

**Human Interaction Data**
- Monitors configuration
- Dashboards configuration
- Notebooks configuration
- Interactive usage during an investigation
- …

# Can we just apply existing Time Series FMs (TSFMs)?

**Forecasting & Anomaly Detection**



**Promise:**

- Several models in recent years
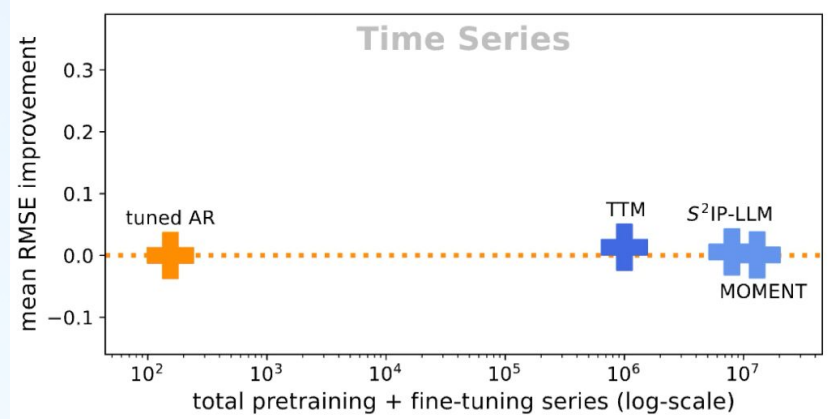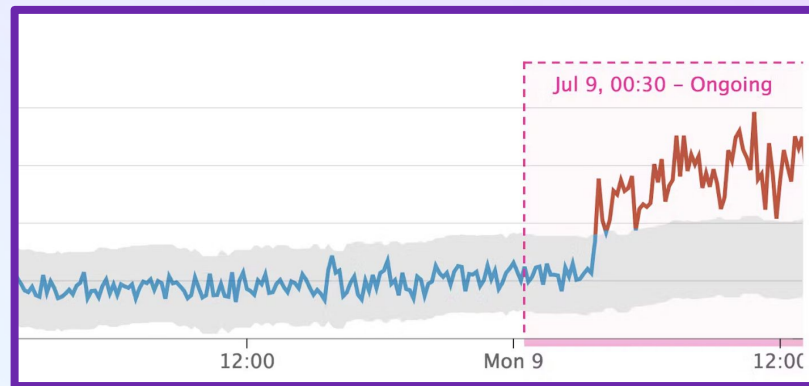- Zero shot capabilities

DATADOG

# Can we just apply existing Time Series FMs (TSFMs)?

## Issues:
- Don't beat supervised baselines
- Not tailored to observability

**Forecasting & Anomaly Detection**



Jul 9, 00:30 – Ongoing

12:00    Mon 9    12:00



Time Series

mean RMSE improvement

0.3
0.2
0.1
0.0
-0.1

tuned AR    TTM    $S^2$IP-LLM

MOMENT

$10^2$  $10^3$  $10^4$  $10^5$  $10^6$  $10^7$

total pretraining + fine-tuning series (log-scale)

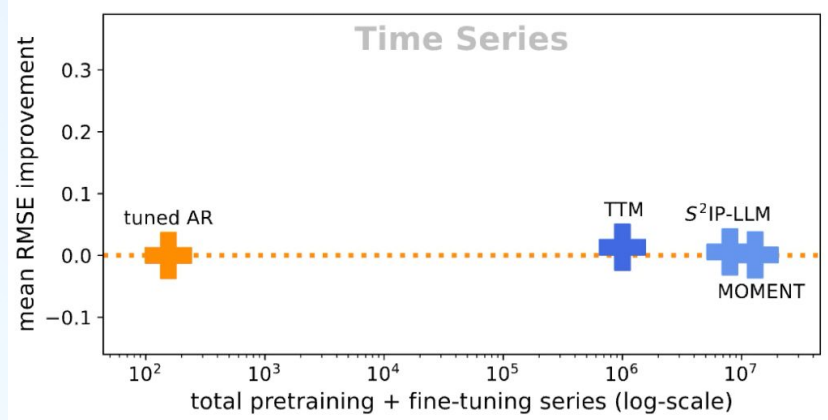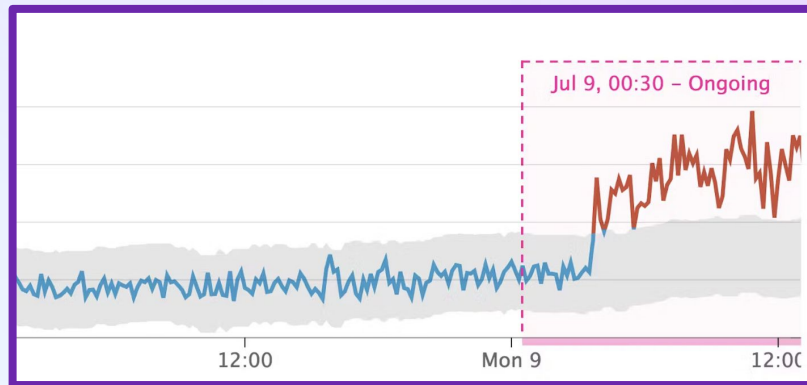DATADOG

# Can we just apply existing Time Series FMs (TSFMs)?

**Our work:** specialize eval, data & modeling for Observability!

## Forecasting & Anomaly Detection



Jul 9, 00:30 – Ongoing



Time Series

tuned AR

TTM

$S^2$IP-LLM

MOMENT

mean RMSE improvement

total pretraining + fine-tuning series (log-scale)

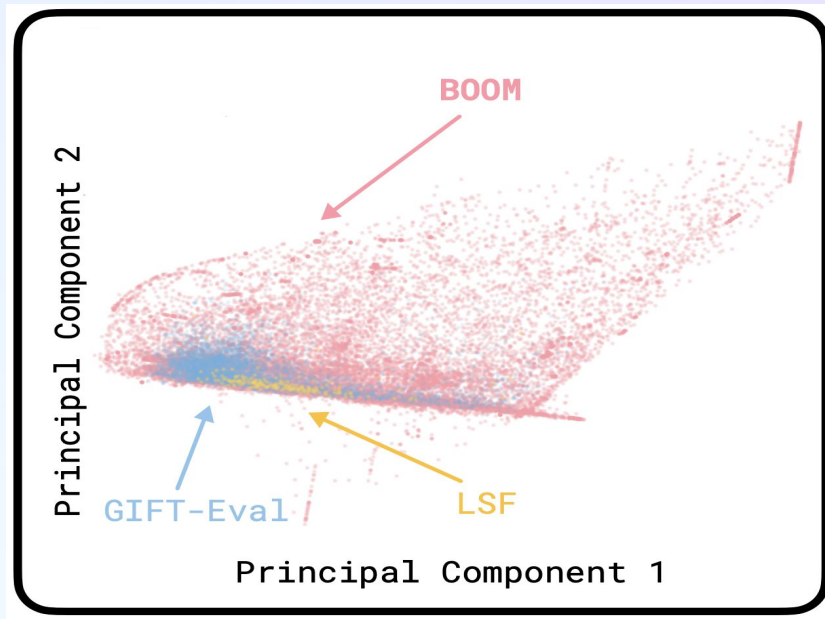DATADOG

# BOOM

## New Observability Benchmark
Largest time series benchmark
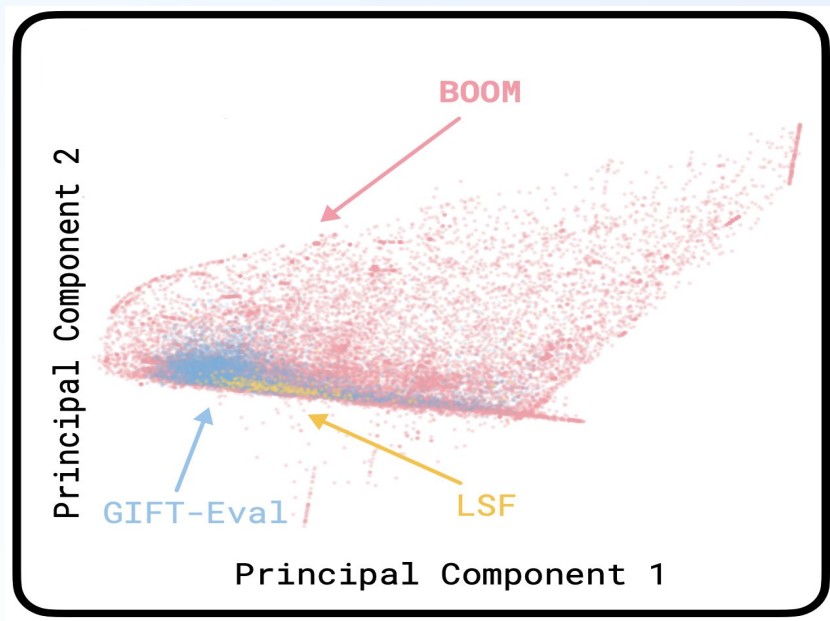
---

## Comprised of Real Data
Internal observability data from Datadog

---

## Open Source
Apache 2.0
27K HF downloads

# Captures challenge of real-world observability data



| Dataset | # Series | # Variates | # Points |
|---|---|---|---|
| BOOM | 2,807 | 32,887 | 350 M |
| BOOMLET | 32 | 1,627 | 23M |
| GIFT-Eval | 144,246 | 147,688 | 158 M |
| LSF | 6 | 370 | 11 M |

# TOTO

**Time Series Foundation Model**

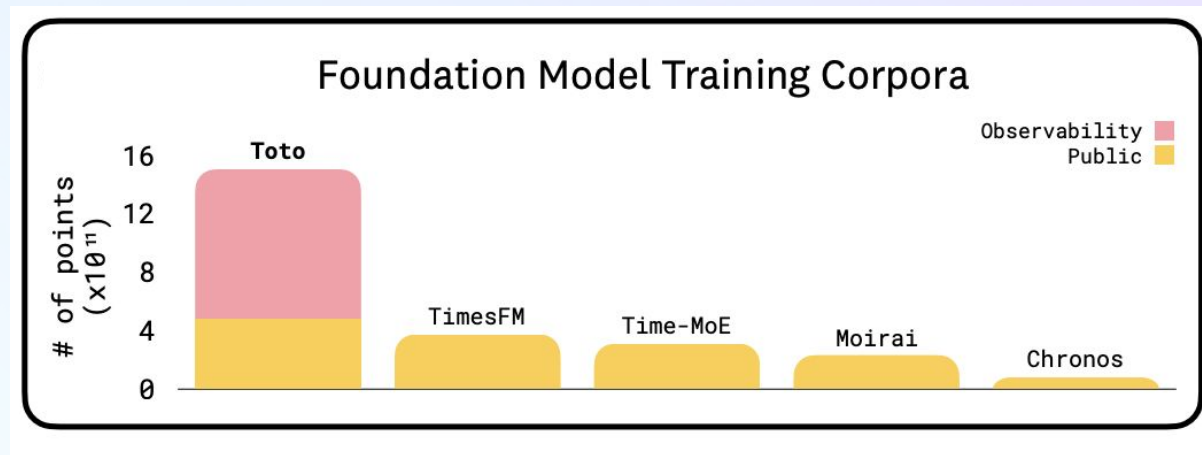150M param decoder-only architecture

---

**Optimized for Observability**

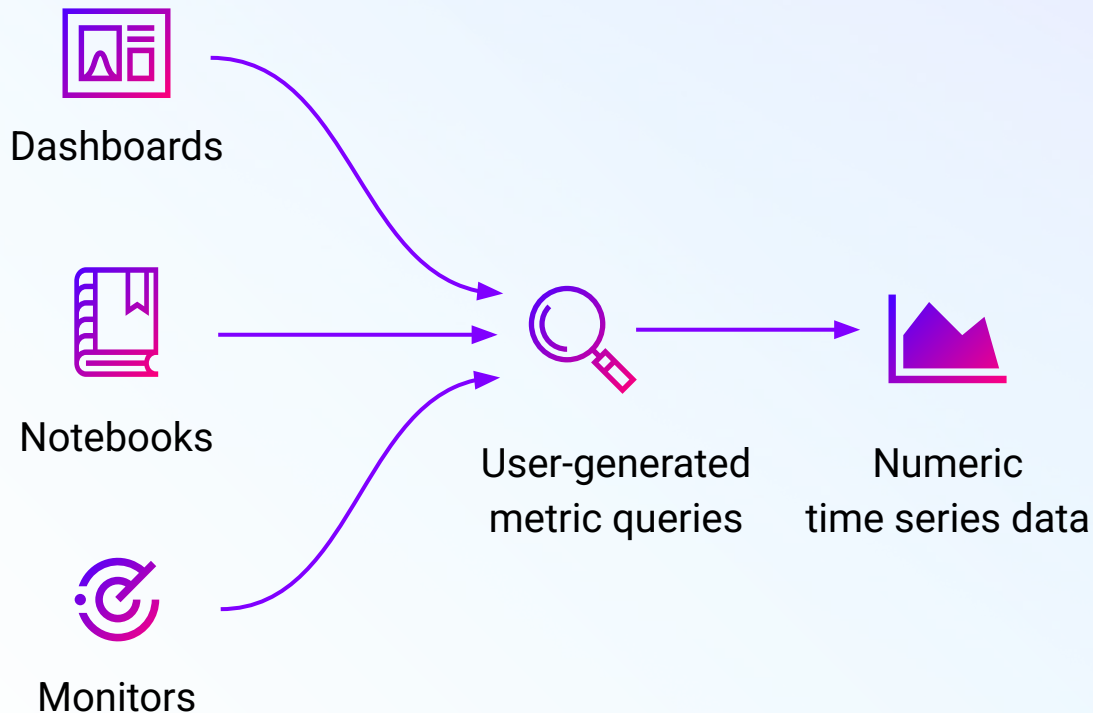And also SOTA on general-purpose
time series forecasting

---

**Open Weights**

Apache 2.0

~8M HF downloads



(Datadog internal data only)

DATADOG

# Data Collection (Datadog internal data only)

Dashboards

Notebooks

Monitors

User-generated
metric queries

Numeric
time series data

Collect each query over:
- Multiple time slices
- Different time intervals

**High-cardinality multivariate data**

**Proportional Attention**: judiciously attend across covariates
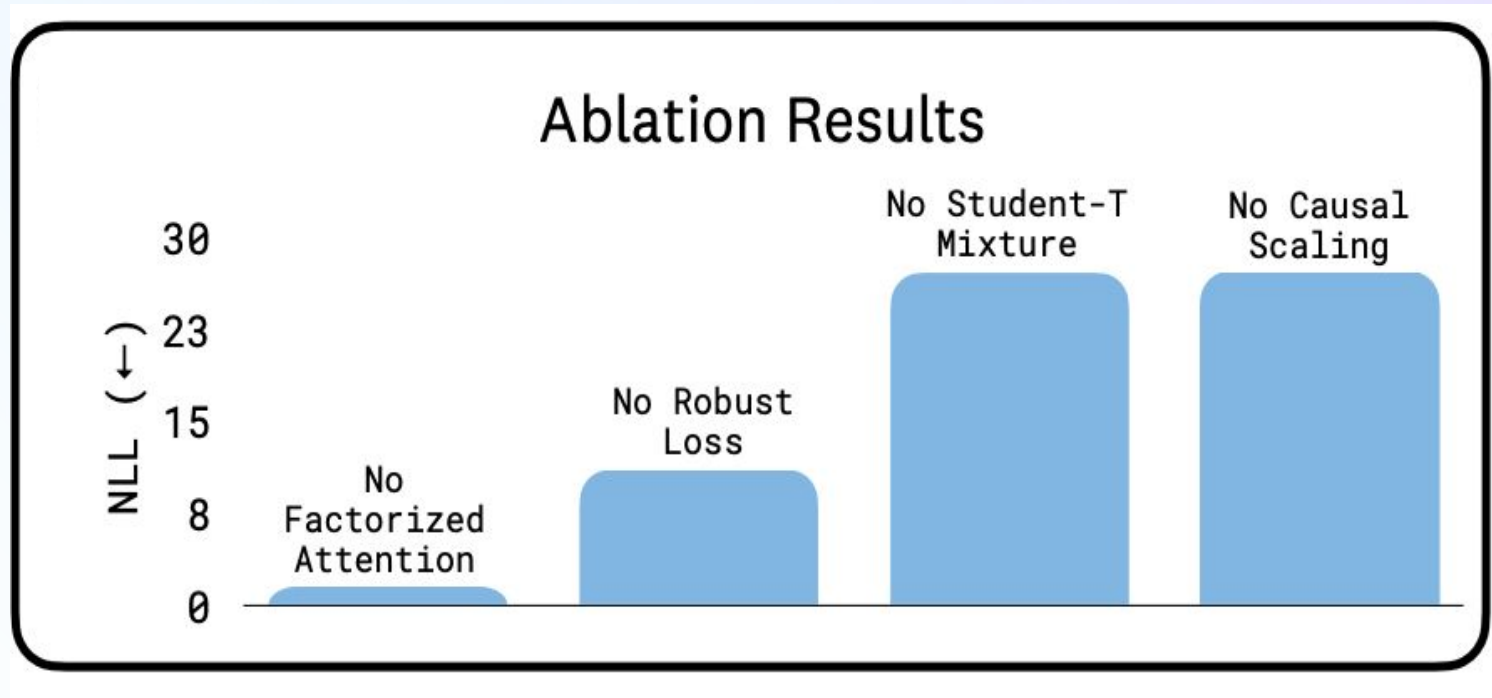
**Skewed, heavy tailed distributions**

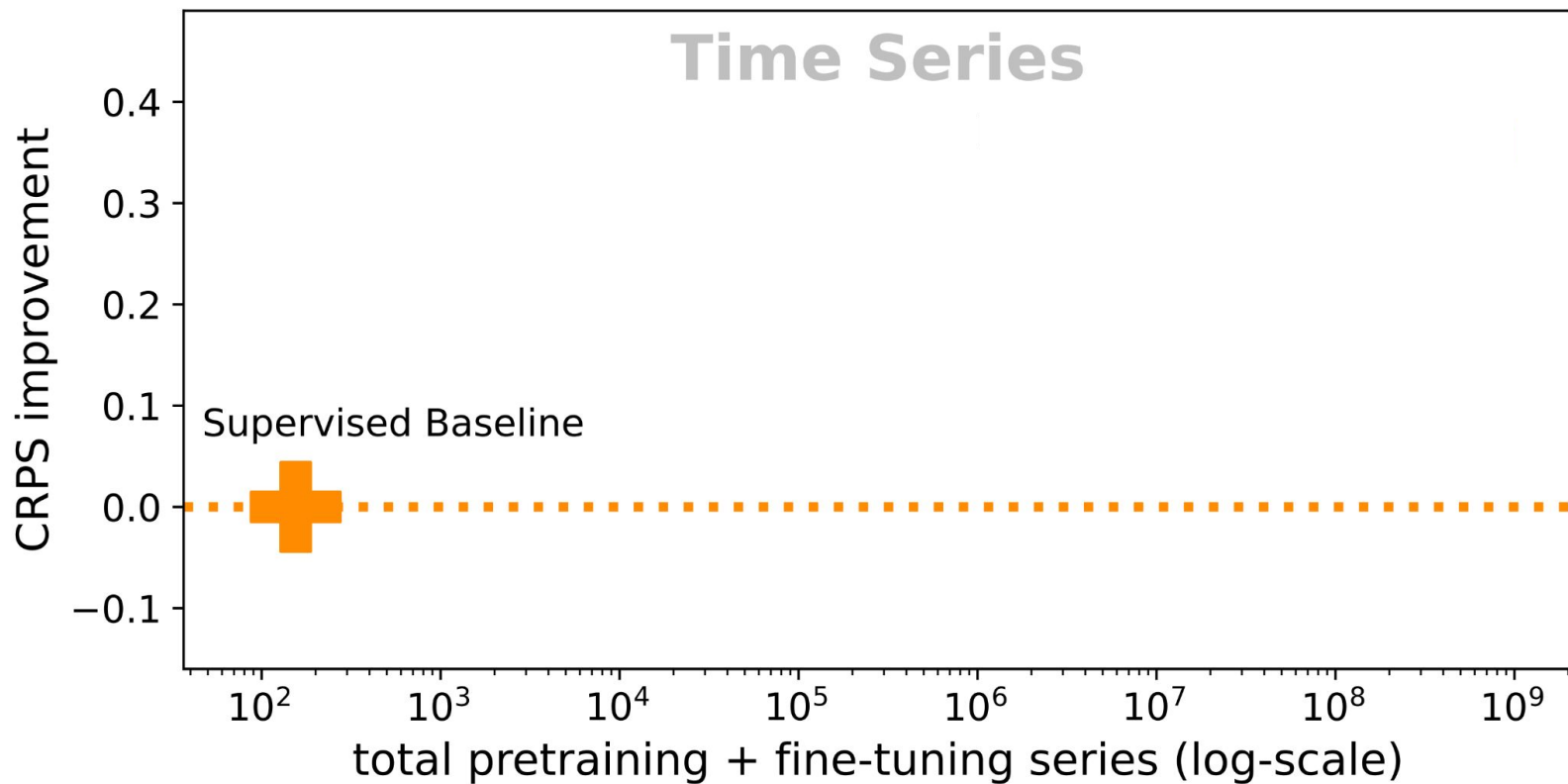**Student-T mixture & robust loss:** for improved modeling and learning

**Extreme dynamic range, nonstationarity**

**Patch-based causal scaling**: address highly non-stationary TS
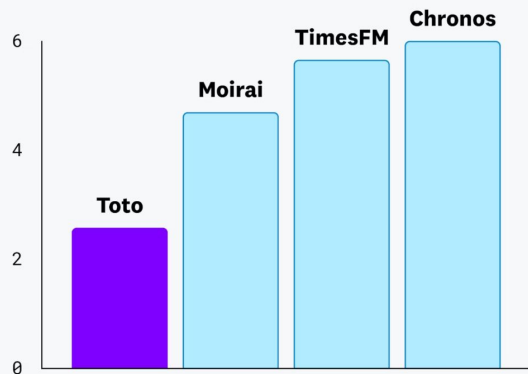
# These modifications make a big difference

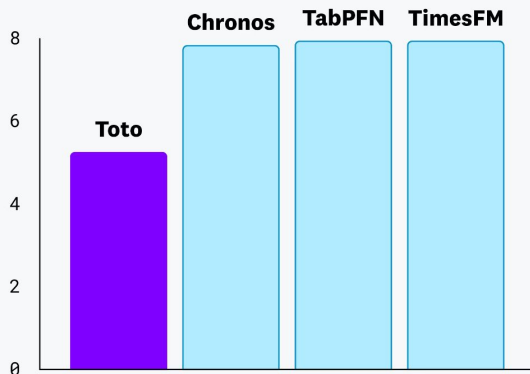# BOOM Results

# Specialized Observability FMs

Datadog AI Research

[Cohen*-Khwaja*-et al.]



**BOOM Results (Rank ↓)**

Toto, Moirai, TimesFM, Chronos

**GIFT-Eval Results (Rank ↓)**

Toto, Chronos, TabPFN, TimesFM

TSFMs have achieved their "BERT moment"!

Specialization matters

***At time of release***

# Ongoing Work



Product Applications



Multimodality



Scaling

DATADOG

# Ongoing Work



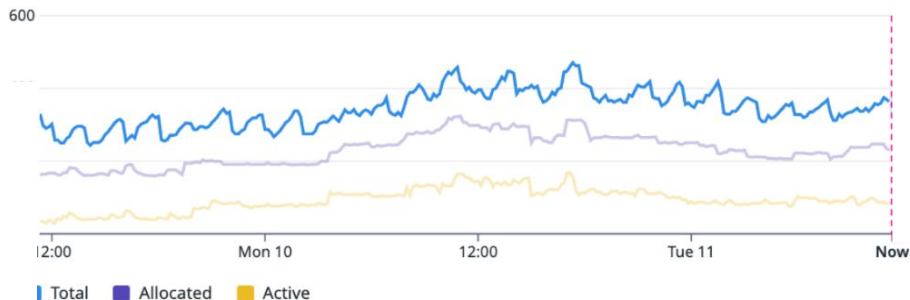**Product Applications**



Multimodality



Scaling

DATADOG

# GPU Monitoring: how many GPUs will I need?



**Device distribution across your fleet**

Visualize GPU allocation to optimize capacity planning and performance

DEVICE ALLOCATION OVER TIME

600

12:00    Mon 10    12:00    Tue 11    Now

Total    Allocated    Active

Important for budgeting/planning

Mature production tool already exists

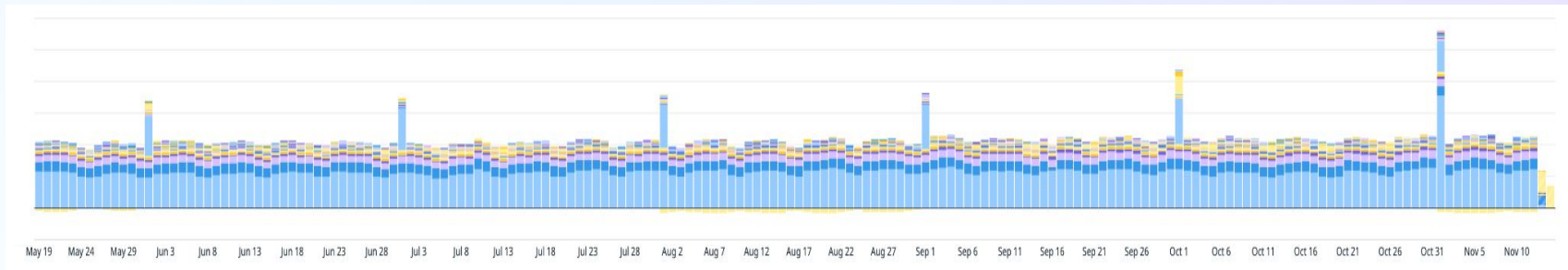Natural application of Zero Shot Toto

DATADOG

# How does Toto perform?

No clear winner between ZS Toto and mature production tool

ZS Toto preferred by humans in 71% of cases in blind evaluation

**In production as of last month!**
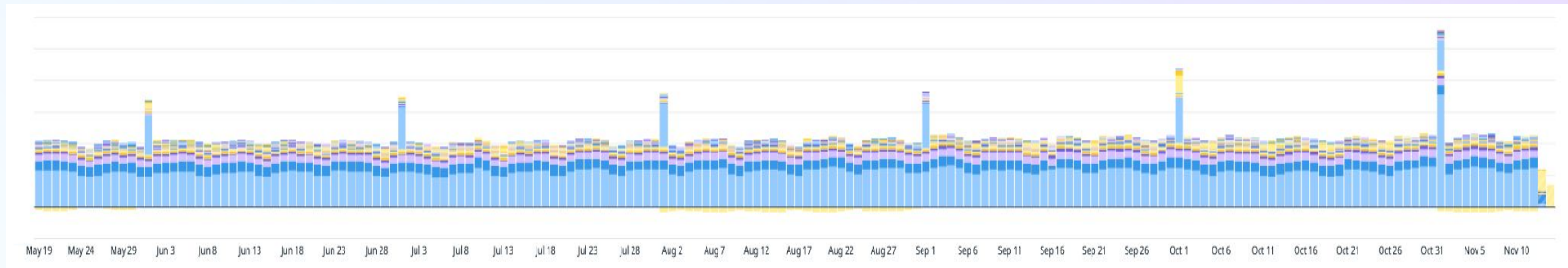
DATADOG

# Forecasting Cloud Costs: how much $ will I spend?



Important for budgeting/planning

Also a seemingly natural application of Zero Shot Toto

# Forecasting Cloud Costs: how much $ will I spend?
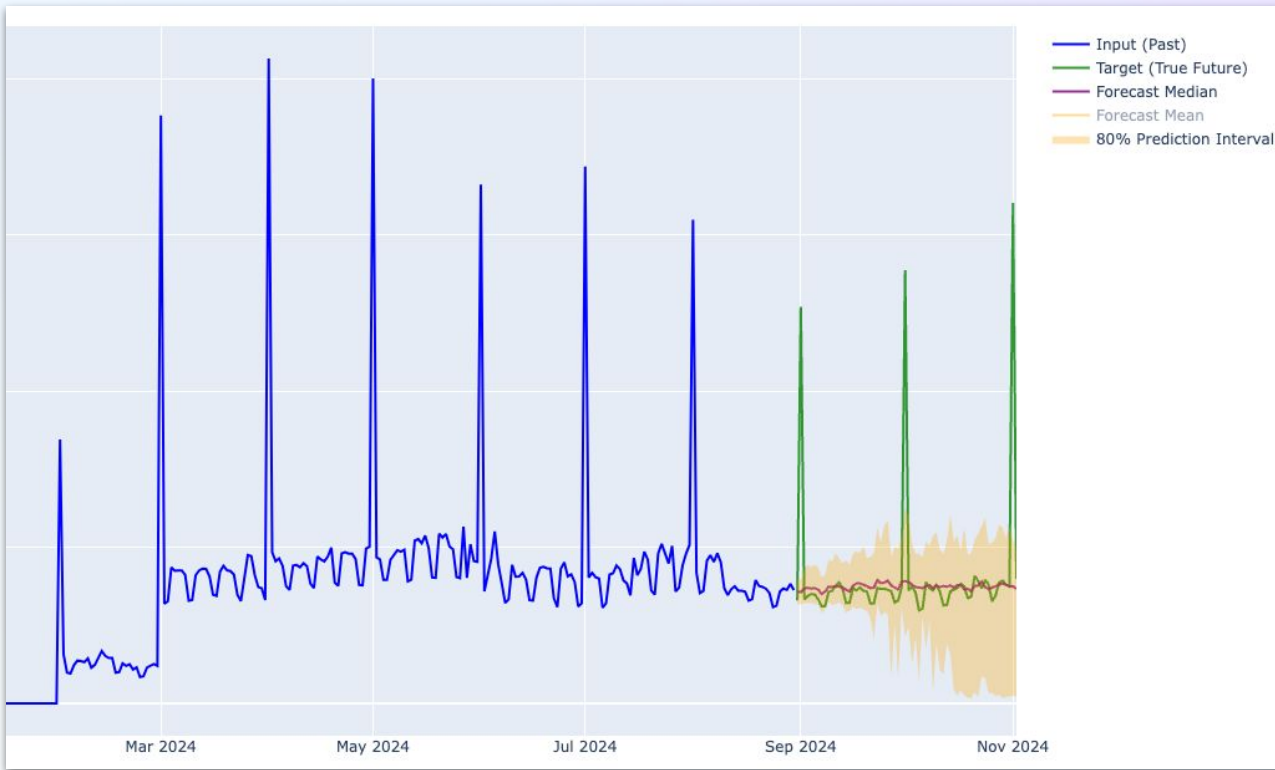


Important for budgeting/planning

Challenges:

- Large, irregular seasonality effects (e.g. day-of-month)
- Not much historical data available
- Product-specific eval differs from typical TSFM benchmarking
- Strong latency constraints

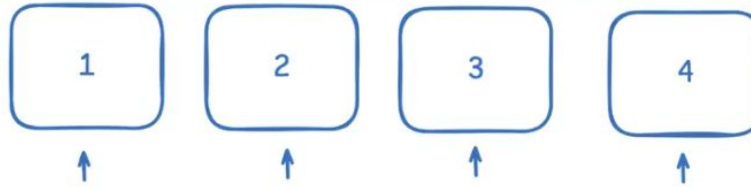DATADOG

# Zero shot Toto misses the spikes
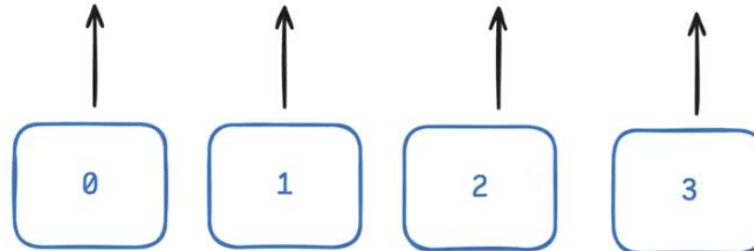


Especially tricky b/c of uneven period lengths

DATADOG

Provide Toto with 'exogenous variable', e.g., day of month!

# Exogenous variable fine-tuning
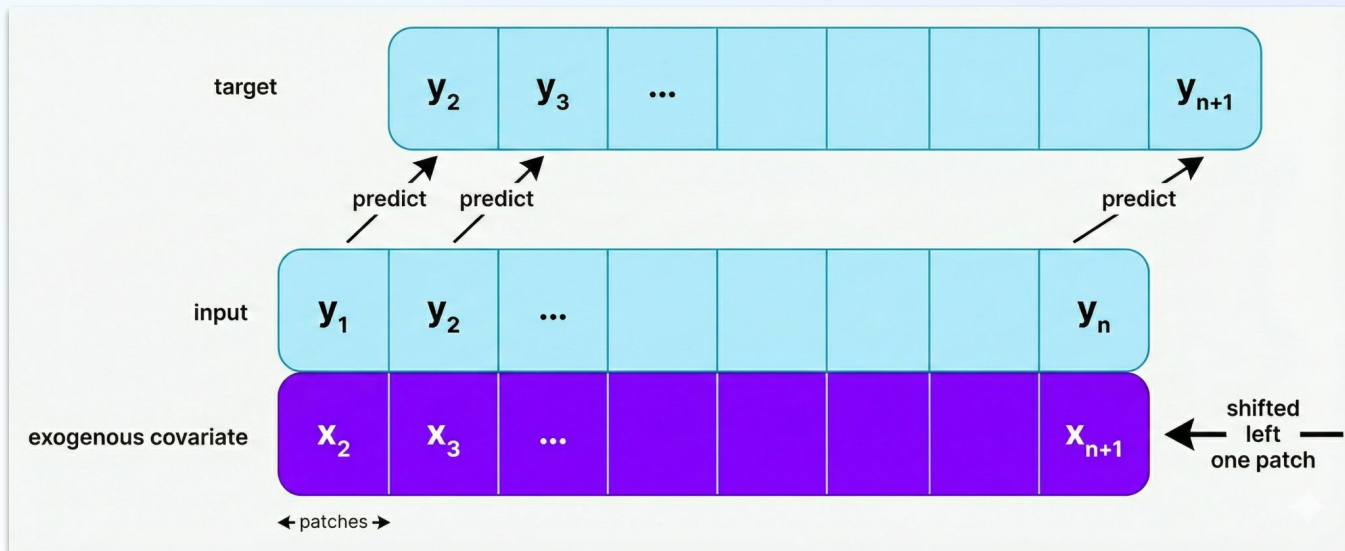
**Input preprocessing**
- Use dummy variable for first of the month
- Shifted exogenous variables one patch into the future
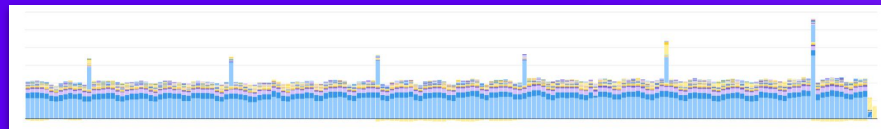- Stacked along variate dimension

**Training**
- Mask loss for exogenous variable

**Inference**
- Inject known future exogenous values during decoding

# Forecasting Cloud Costs



## CHALLENGES

Irregular seasonality effects

---

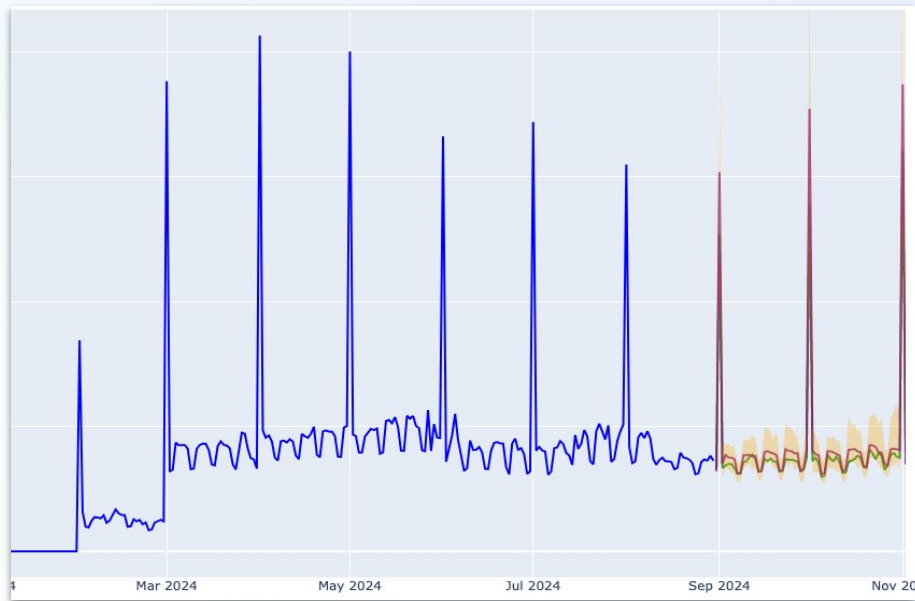Product-imposed latency constraints

---

Product-specific evaluation

## SOLUTIONS

Toto + FT + Exogenous variables

---

Toto inference meets latency requirements

---

New benchmark & metrics

DATADOG

# The end result...



17% improvement over baseline
41% over zero-shot
Satisfies latency reqs

**Coming soon:** Support for
fine-tuning and exogenous variables
https://github.com/DataDog/toto

# Toto applications, next steps

**Autoscaling** – Forecast demand so services can right-size

**Predictive alerting** – Forecast issues before they happen (and ideally fix them w/o needing to page an engineer)

# Datadog AI Research:
# Vertical AI for observability

# We are hiring!