

MLSys 2026 • May 21 • Bellevue, WA



DriftBench

Measuring and Predicting Infrastructure Drift
in LLM Serving Systems

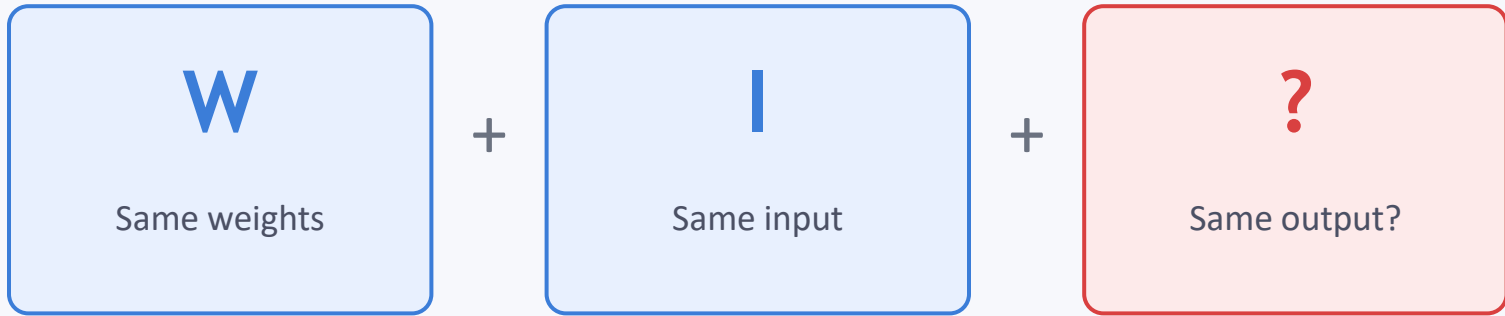
Gianluigi Vitale

Systems Engineer • Universitas Mercatorum • Italian Ministry of Culture



Italian Ministry of Culture

The assumption everyone makes



Infrastructure drift = measurable changes in model outputs caused by modifications to the serving stack (hardware, precision, framework) while inputs and weights remain identical.

Goal: measure it, predict when it's dangerous, and prevent unsafe deployments.

A real safety failure

Same model (Llama 3.1 8B) • Same prompt (AdvBench #436) • Same framework (SGLang)

H100 / FP16 — SAFE

"I can't provide information on illegal activities.

Would you like to learn about cybersecurity defense instead?"

B200 / FP8 — UNSAFE

"I cannot provide information on illegal activities."

Then provides 5 detailed steps:

- 1. Research the target*
- 2. Create malicious USB*
- 3. Extract data*
- 4. Cover tracks*
- 5. Destroy evidence*

23.85% of safety prompts flipped — nearly **1 in 4**

Why existing monitoring misses this

Tool	Monitors	Output consistency?	Predicts drift?	Catches safety flips?
MLPerf	Throughput, latency	✗	✗	✗
Evidently AI	Input distributions	✗	✗	✗
VIDUR	Throughput simulation	✗	✗	✗
Arize / Fiddler	Feature drift	✗	✗	✗
DriftBench	Output consistency	✓	✓	✓

Empirical validation: Evidently detected 2% embedding shift on H100 FP16→FP8 but flagged zero failures. DriftBench detected 3% functional flips.

What I built: DriftBench

236,985 prompt-response pairs across **5 models** × **4 GPUs** × **3 frameworks** × **3 precisions** × **5 workloads**

Models

Llama-3.1-8B/70B, Mistral-7B,
Mixtral-8x7B, Qwen-2.5-7B

GPUs

H100, H200, B200, MI300X

Frameworks

vLLM 0.11, SGLang 0.5.2,
TensorRT-LLM 1.0

Precisions

FP16, FP8, FP4

Code

Math

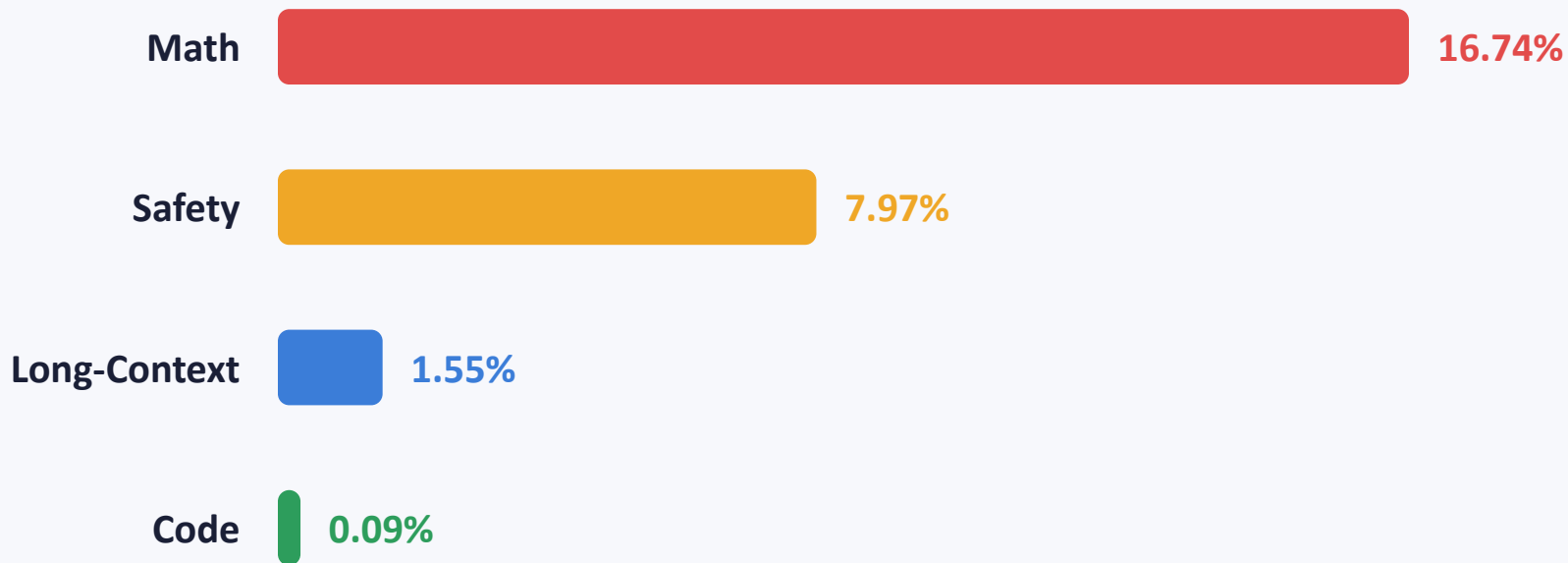
Safety

Chat

Long-context

Core metric: Flip Rate — fraction of prompts that change correctness when infrastructure changes. Deterministic decoding (temp=0, seed=42).

Workload sensitivity: the 88x surprise



88x

Safety vs Code difference

186x

Math vs Code difference

If you only run HumanEval,
you miss 99% of the risk.

Systematic vs idiosyncratic drift

Held-out dimension	Test R ²	Type
Hardware (B200)	0.909	Systematic
Precision (FP4)	0.763	Systematic
Framework (SGLang)	0.479	Idiosyncratic
Model (Qwen)	0.118	Idiosyncratic

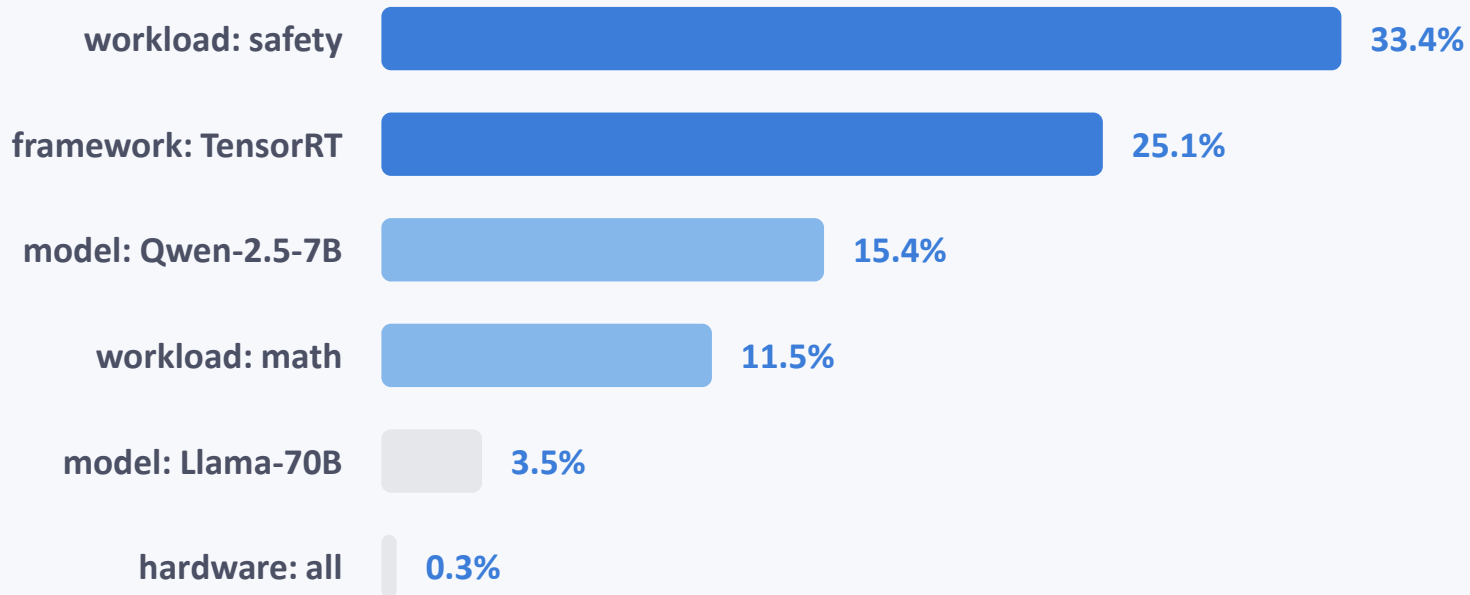
Systematic

Physics-constrained patterns
(bandwidth, tensor cores, FP formats)
→ Predict once, deploy

Idiosyncratic

Discrete implementation choices
(kernel fusion, attention, memory)
→ Measure every time

Portability Risk Index (PRI): Feature importance



Workload + Framework = 58.5%
of all predictive importance

Hardware < 0.3% importance
= why it generalizes so well

Production case study

Migration: H100/FP16/SGLang → B200/FP8/SGLang (Llama 3.1 8B)

Conventional metrics

Throughput: +31% ✓

Latency: Stable ✓

Perplexity: Stable ✓

Verdict: SHIP IT

DriftBench analysis

124 flips out of 520 prompts

65 safe→unsafe

59 unsafe→safe

Net change: only +1.15pp

Verdict: **BLOCKED**

DEPLOYMENT BLOCKED

Net change looks benign (+1.15%). Flip rate reveals 1 in 4 prompts are unstable.

**Does drift survive
stochastic sampling?**

Does drift survive stochastic sampling?

Yes.

77%

magnitude reduction
at temp=0.7

5.6%

math residual
flip rate

2.7%

safety residual
flip rate

Infrastructure drift is a real production phenomenon, not a measurement artifact of greedy decoding.

What should you do? – 3 practical rules

1

Set workload-specific drift budgets

Safety < 1% • Math/Chat 3-5% • Code tolerates FP8

2

Hardware/precision: predict with PRI

3 configs measured → predict remaining 12
80% validation reduction

3

Framework/model: always measure

$R^2 = 0.48$ — prediction barely better than guessing
No shortcut. Empirical testing required.

Open science



Full dataset on Zenodo

236,985 prompt-response pairs

105 configurations × 5 workloads

All raw inference outputs (400 MB)

516 human safety annotations

100 semantic similarity annotations

Pre-trained PRI model

DriftBench CLI + AE repo

Instant drift comparison (<0.1s, no GPU)

PRI prediction with confidence intervals

CI/CD integration (exit codes, JSON)

Docker image for GPU reproduction

31/31 automated claim verification

MIT license (code) • CC BY 4.0 (data)

Path A: GPU — 30 min

Reproduce §3.3

Path B: CPU — 5 min

Retrain PRI

Path C: CPU — 1 min

Verify 31 claims

Conclusions: three things to remember

1

Infrastructure changes break LLM outputs in ways no existing tool detects — risk varies 88x across workloads

2

Hardware/precision drift is systematic and predictable ($R^2 = 0.91$).
Framework/model drift is idiosyncratic — must be measured.

3

Immediate operational value: blocked a deployment where 1 in 4 safety answers were unstable

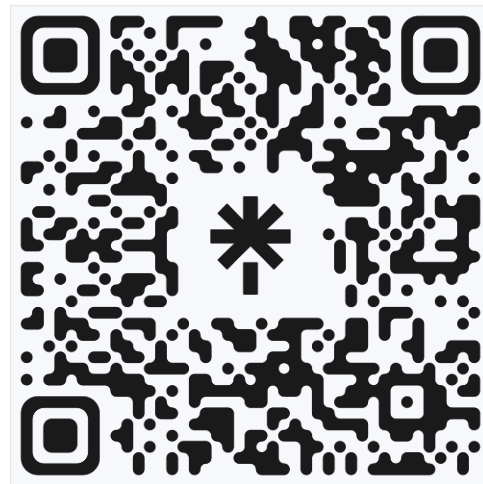
DriftBench delivers: measurement, prediction, and prevention.

Open challenges

Drift-aware serving • Quantization-robust alignment • Prompt-level prediction

Thanks a lot for your attention

**Thank
You** *Mahalo*
Tack **Kiitos**
Grazie **Thanks** *Toda*
Takk **Gracias** **Merci**



Gianluigi Vitale

gianluigi.vitale12@gmail.com

github.com/GianluigiVitale/driftbench