

# Perseus: A Fail-Slow Detection Framework for Cloud Storage Systems

**Ruiming Lu**, Erci Xu, Yiming Zhang,  
Fengyi Zhu, Zhaosheng Zhu, Mengtian Wang,  
Zongpeng Zhu, Guangtao Xue, Jiwu Shu, Minglu Li, Jiesheng Wu



# Data Center Instability

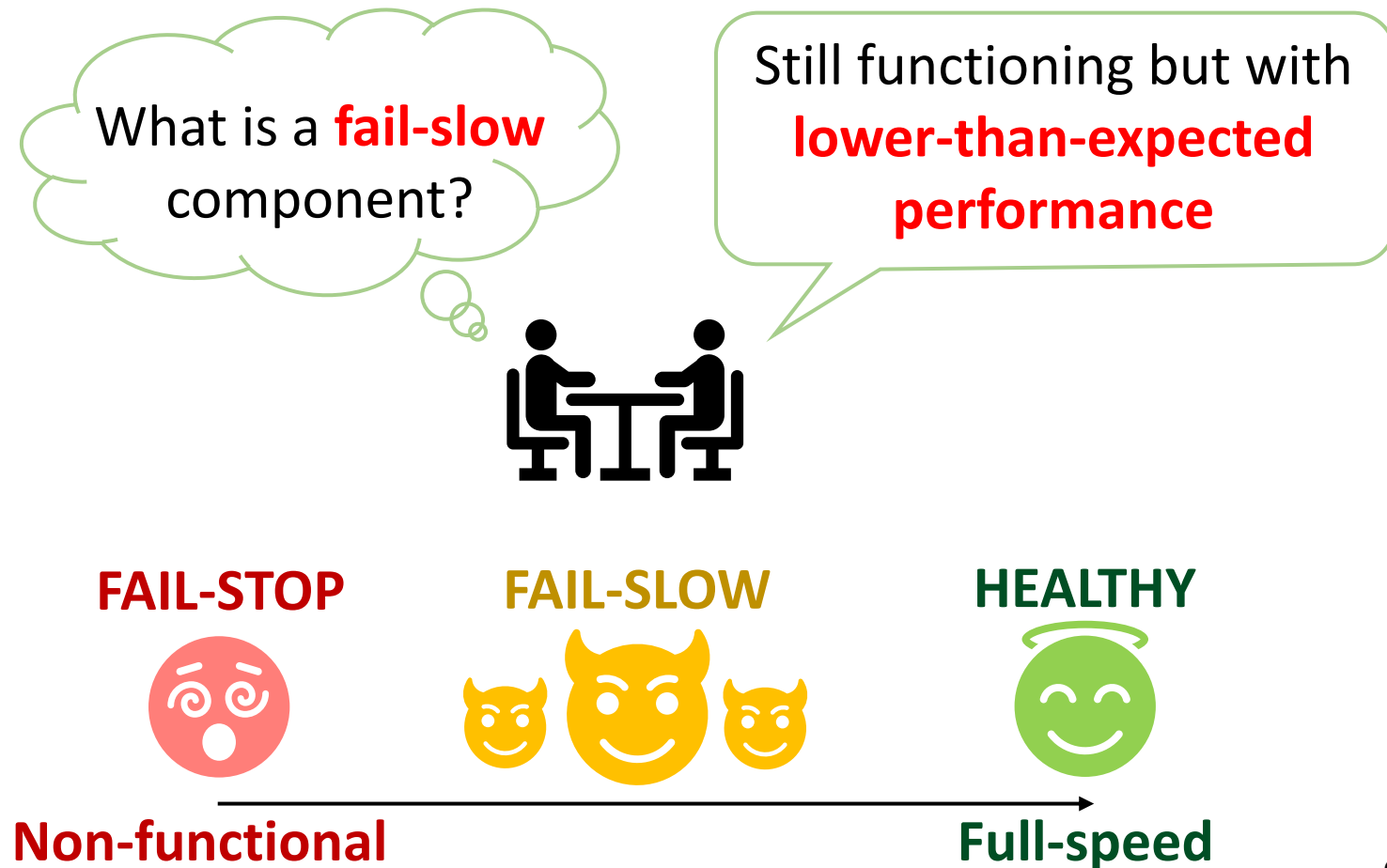
- Failures ⚡ in The Wild

- **Fail-Slow** 🔍

- Fail-Stop

- Byzantine

- ...



FAIL-SLOW



Severe

Not  
Uncommon

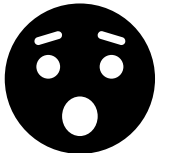
Confusing

“Fail-slow NVMe SSDs can **degrade to SATA SSD or HDD-level performance**<sup>[1]</sup>”

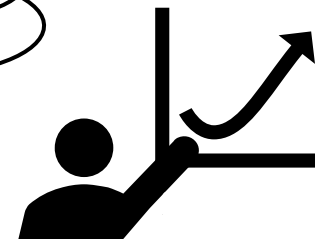


Annual fail-slow failure rate is **1-2%**<sup>[2]</sup>!

**As frequent as fail-stop incidents!**



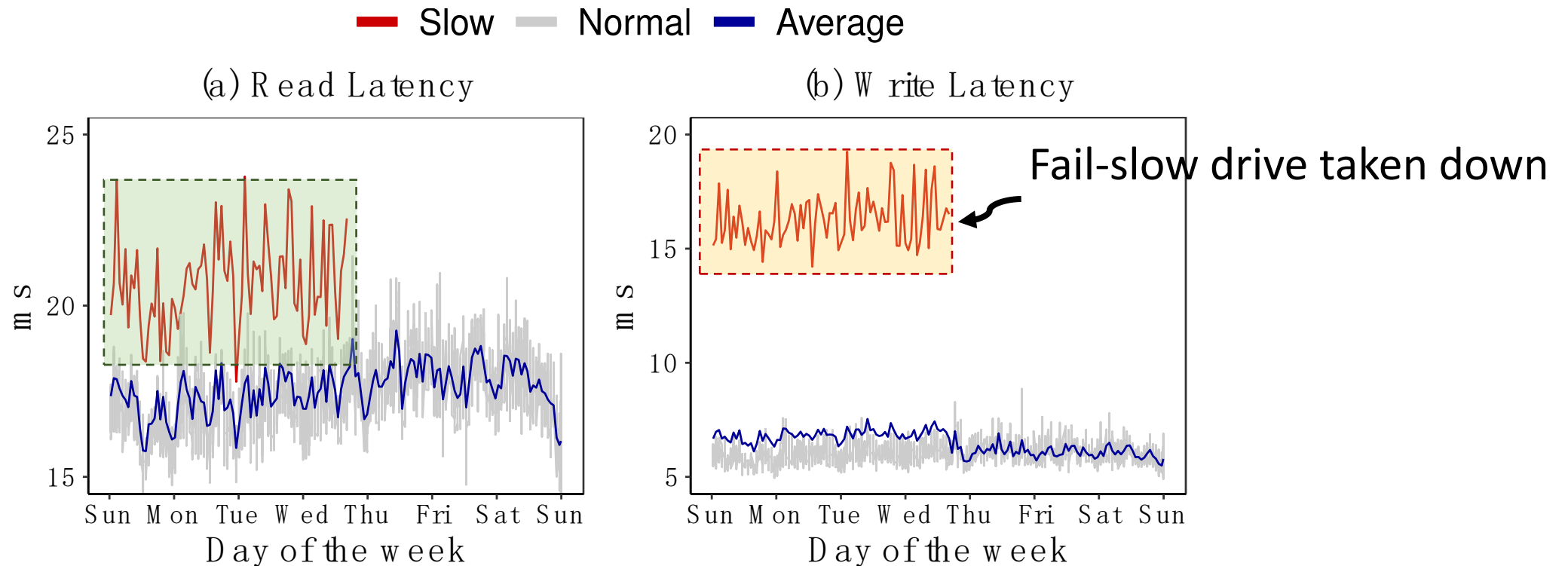
Fail-slow or just **normal variations**?



[1] NVMe SSD Failures in the Field: the Fail-Stop and the Fail-Slow, Lu et al.

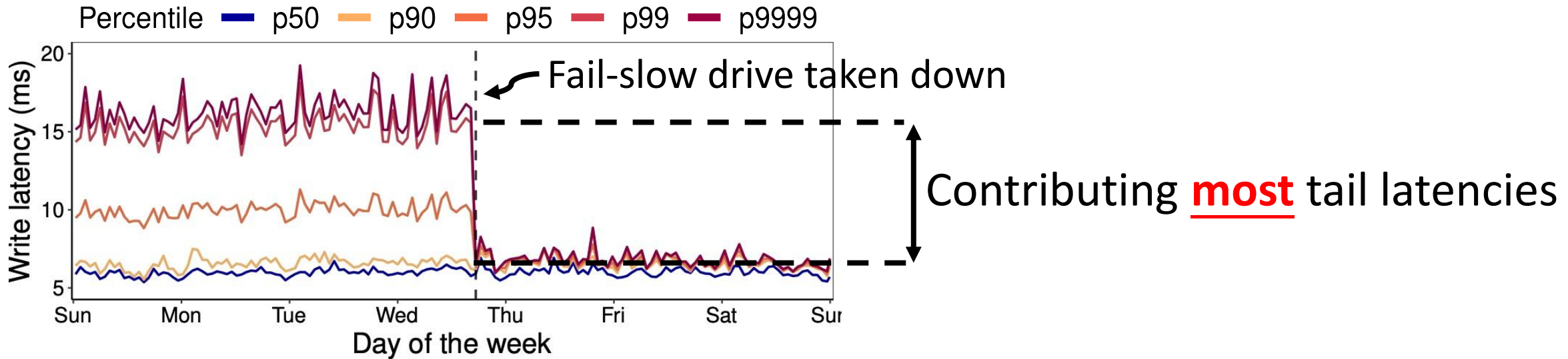
[2] IASO: A Fail-Slow Detection and Mitigation Framework for Distributed Storage Services, Panda et al.

## Fail-Slow in The Field:



**1.01-1.50X higher for read** **2.06-3.65X higher for write**

## Fail-Slow in The Field:



FAIL-SLOW



**Slow! Slow! Slow! Silent performance degradation!**

**Let's "dig" them out!**



- No Ground Truth in Identifying Fail-Slow



How slow is a drive to be considered fail-slow?

>100us?



>500us?



>1ms?



A “thousand” ways to define!

...



- Previous FSD Studies Are
  - Intrusive
    - Source Code Accessing/Altering
  - Coarse-grained
    - Node-Level Detection

## Capturing and Enhancing *In Situ* System Observability for Failure Detection

Peng Huang  
*Johns Hopkins University*

Chuanxiong Guo  
*ByteDance Inc.*

Jacob R. Lorch    Lidong Zhou  
*Microsoft Research*

Yingnong Dang  
*Microsoft*

## IASO: A Fail-Slow Detection and Mitigation Framework for Distributed Storage Services

Biswaranjan Panda, Deepthi Srinivasan, Huan Ke\*,  
Karan Gupta, Vinayak Khot, and Haryadi S. Gunawi\*

*Nutanix Inc.*

*University of Chicago\**

### Abstract

We address the problem of “fail-slow” fault, a fault where a hardware or software component can still function (does not fail-stop) but in much lower performance than expected.

absolute failure of sub-components but can also gracefully handle the occurrence of performance faults.

In this context, our work in this paper makes the two following contributions:

(1) Design and implementation of a fail-slow mitigation

- **Our Work Shares**

- Years of Experiences in FSD
- A Practical FSD Framework named **Perseus**
- Root Cause Analysis





INTRODUCTION



**DATASET**



FAILED  
ATTEMPTS



PERSEUS



EVALUATION &  
CONCLUSION

- 248K+ drives
  - 55% NVMe SSD + 45% SATA HDD
  - 4 manufacturers
  - 9 major drive models
  - Diverse cloud services:
    - Log service, big data, E-commerce, table storage, stream processing, database, object storage, data warehouse, block storage

- 248K+ drives
- 10-month performance logs (iostat)
  - Latency/throughput time series
- Test dataset released
  - <https://tianchi.aliyun.com/dataset/144479>



INTRODUCTION



DATASET



**FAILED  
ATTEMPTS**



PERSEUS



EVALUATION &  
CONCLUSION

# Ideal Fail-Slow Detection Should Be ...?

## Efficient Fail-Slow Detection

---



### Non-intrusive

No source code altering  
External performance log-based



### Accurate

High precision/recall rate



### Fine-grained

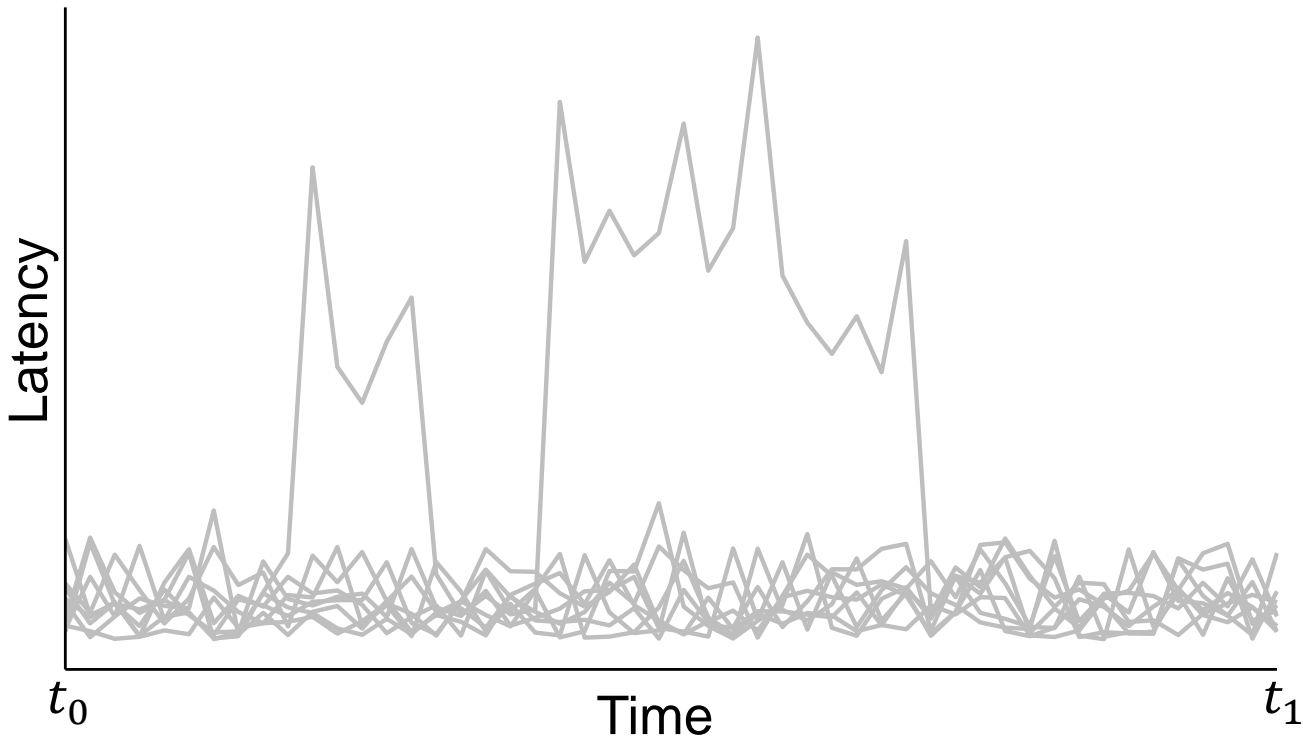
Device-level detection



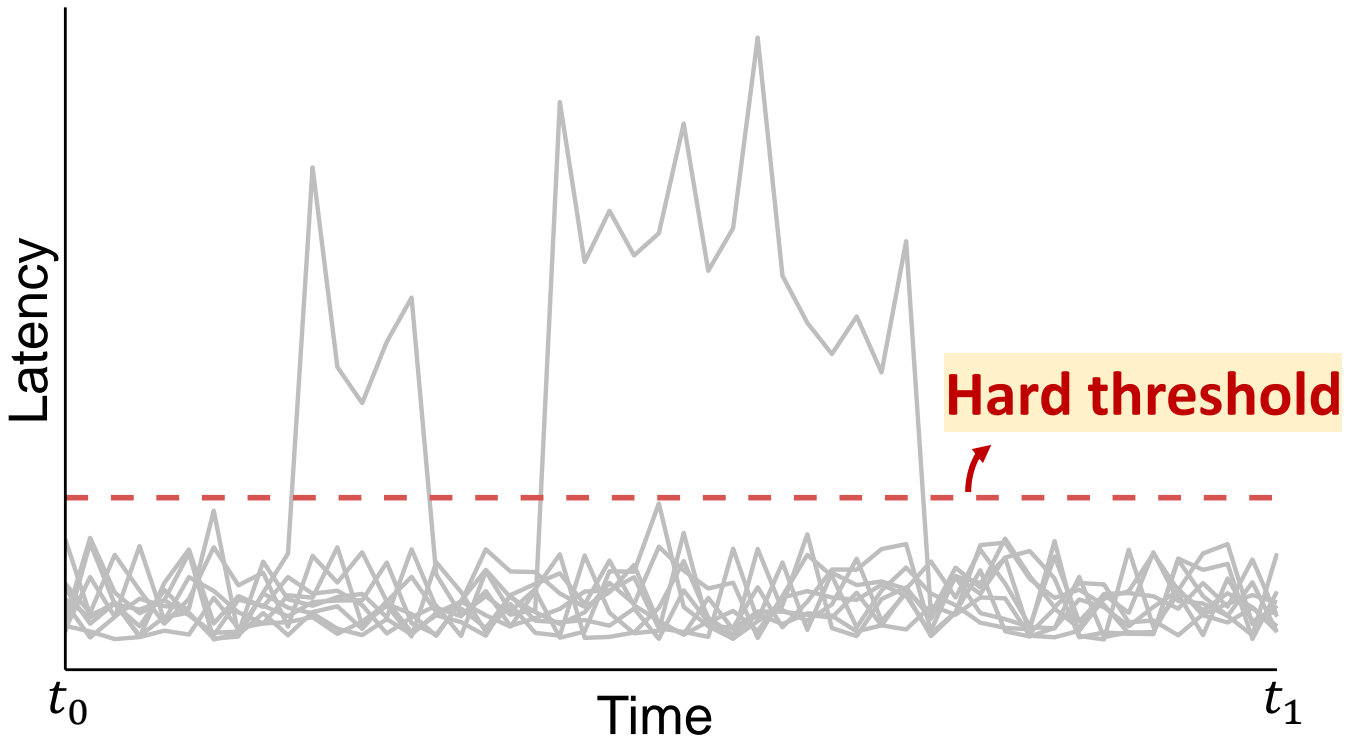
### General

Adaptable to {  
SSD/HDD clusters  
Various cloud services

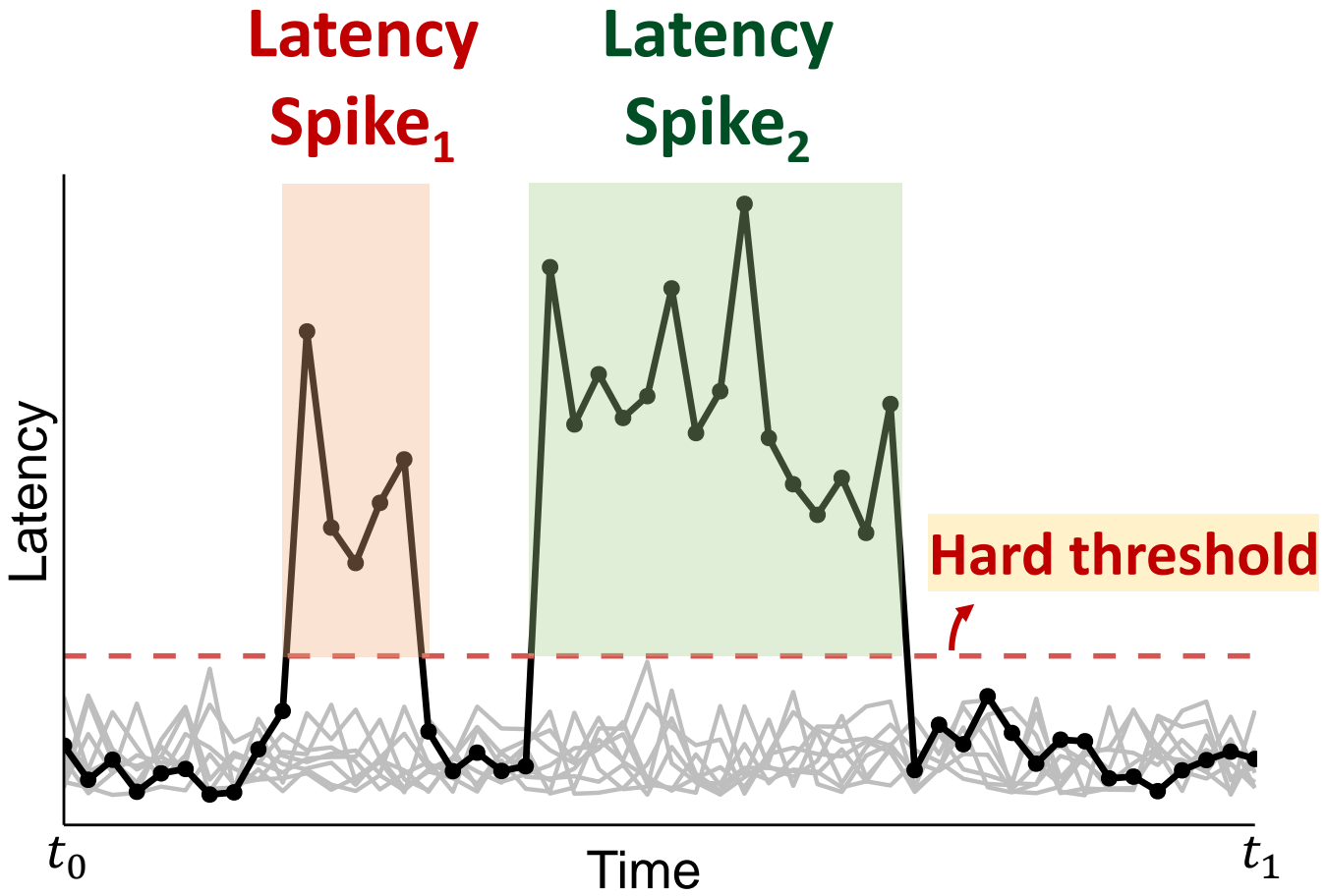
# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering



# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering

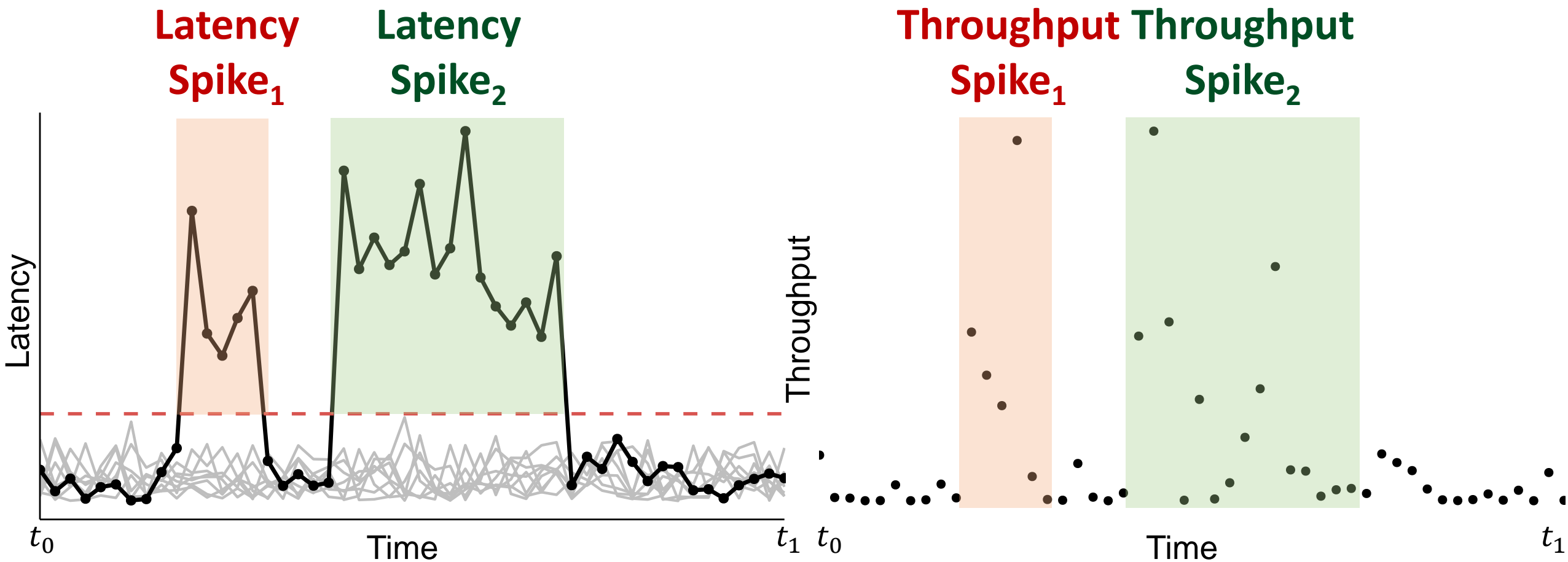


# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering



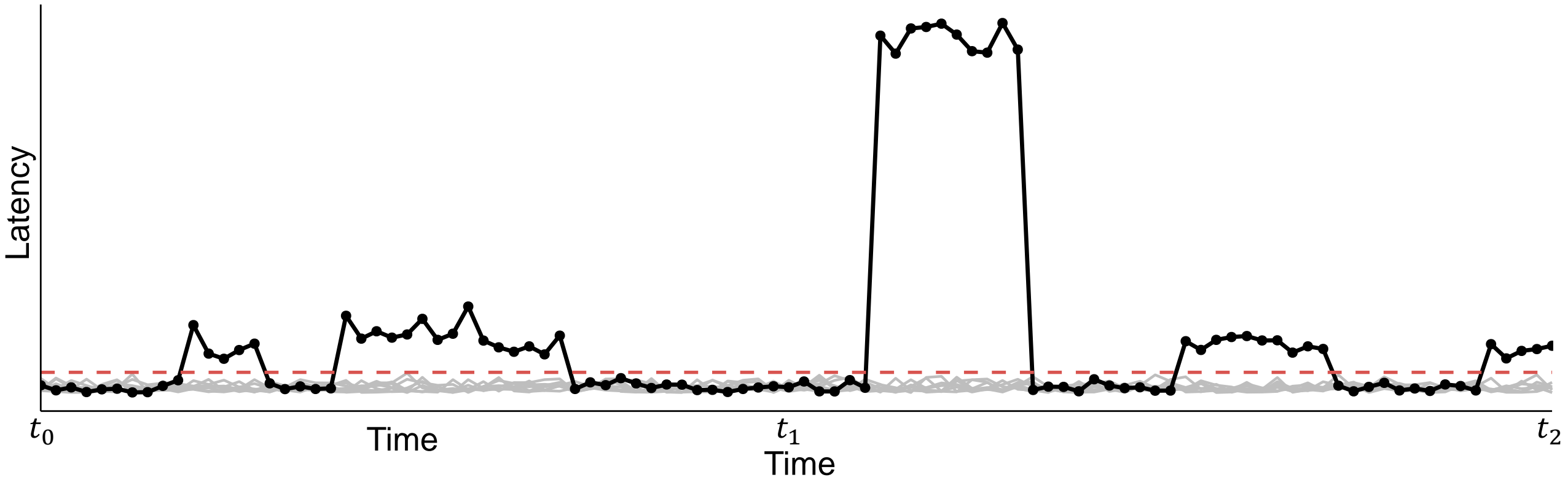


# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering



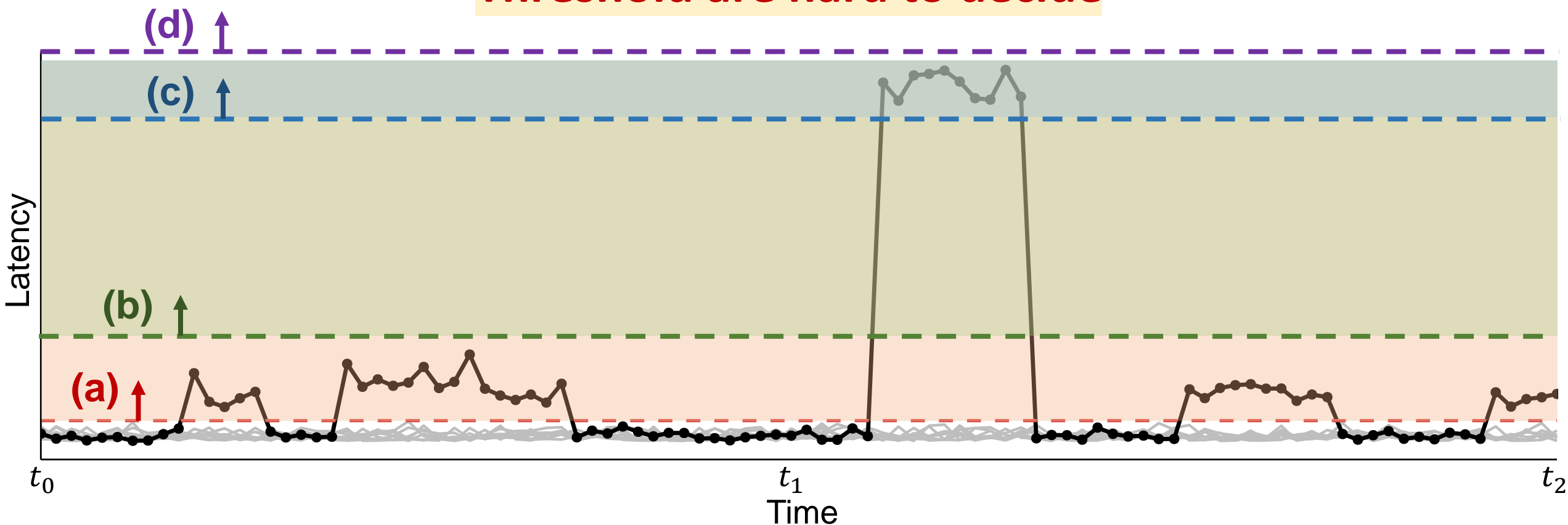
**Workload bursts are common causes of latency variations**

# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering



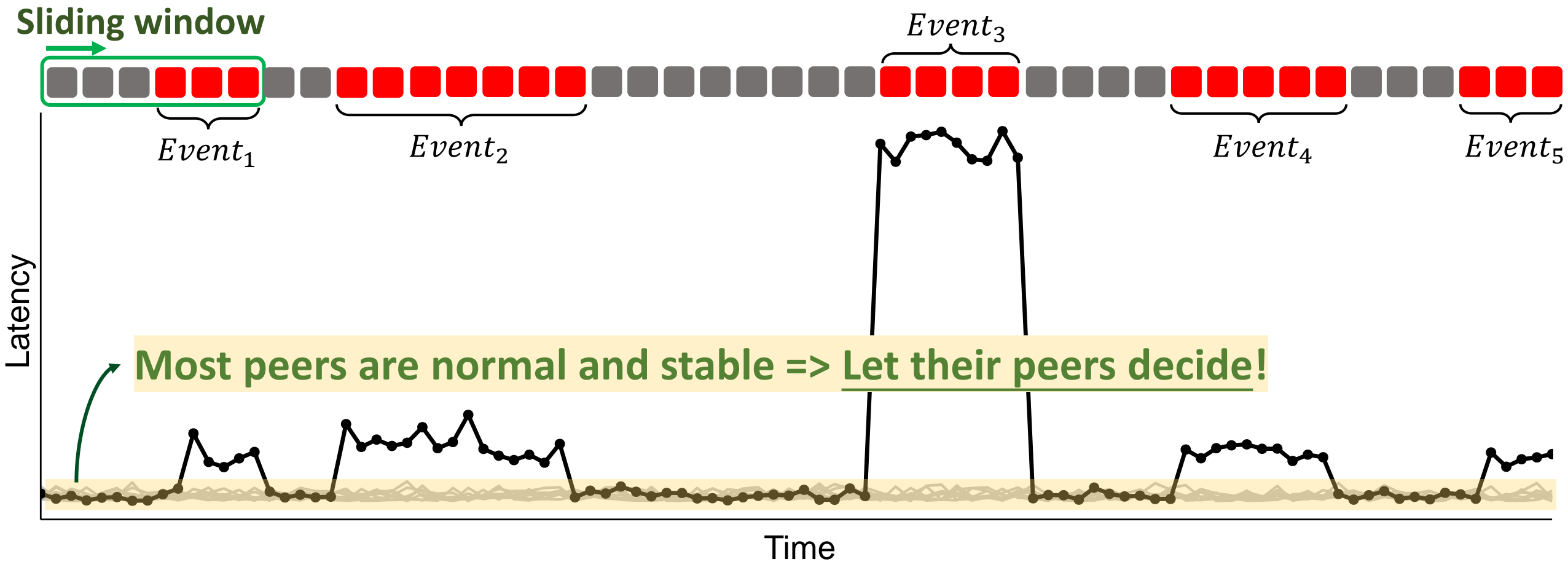
# FAST<sup>1</sup><sub>23</sub> Failed Attempt: Threshold Filtering

Threshold are hard to decide



**Dilemma** { Relaxed => More False Positives  
Strict => Many Fail-Slow Undiscovered

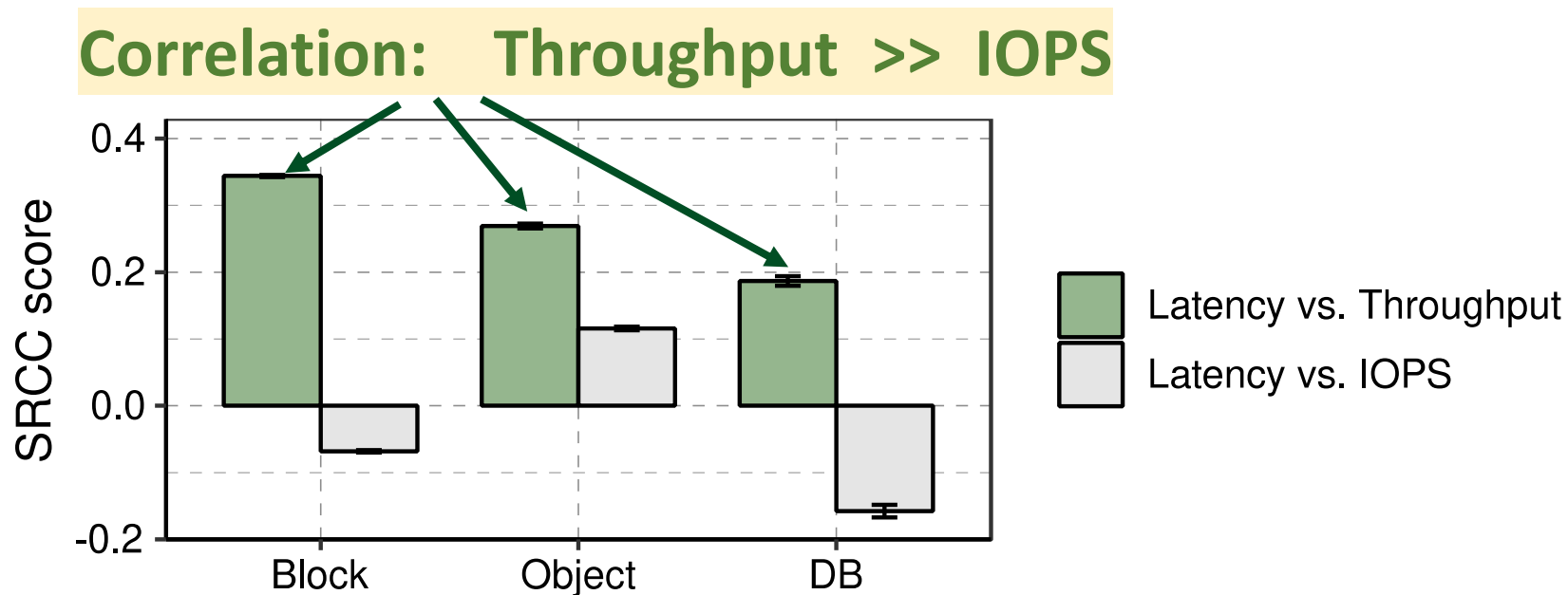
# Failed Attempt: Peer Evaluation



**Time-consuming to fine-tune**

Insight: “**Workload pressure** can affect latency variations”

- Throughput or IOPS?

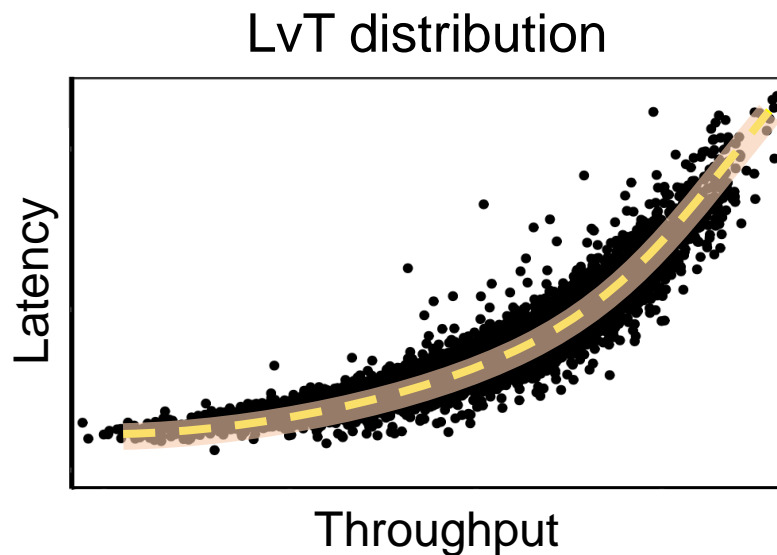


**Guideline 1: Use throughput to model the workload pressure**

Insight: “Workload pressure can affect latency variations”

- How to model such a positive correlation?

**Guideline 2: Model the latency-vs-throughput (LvT) distribution**

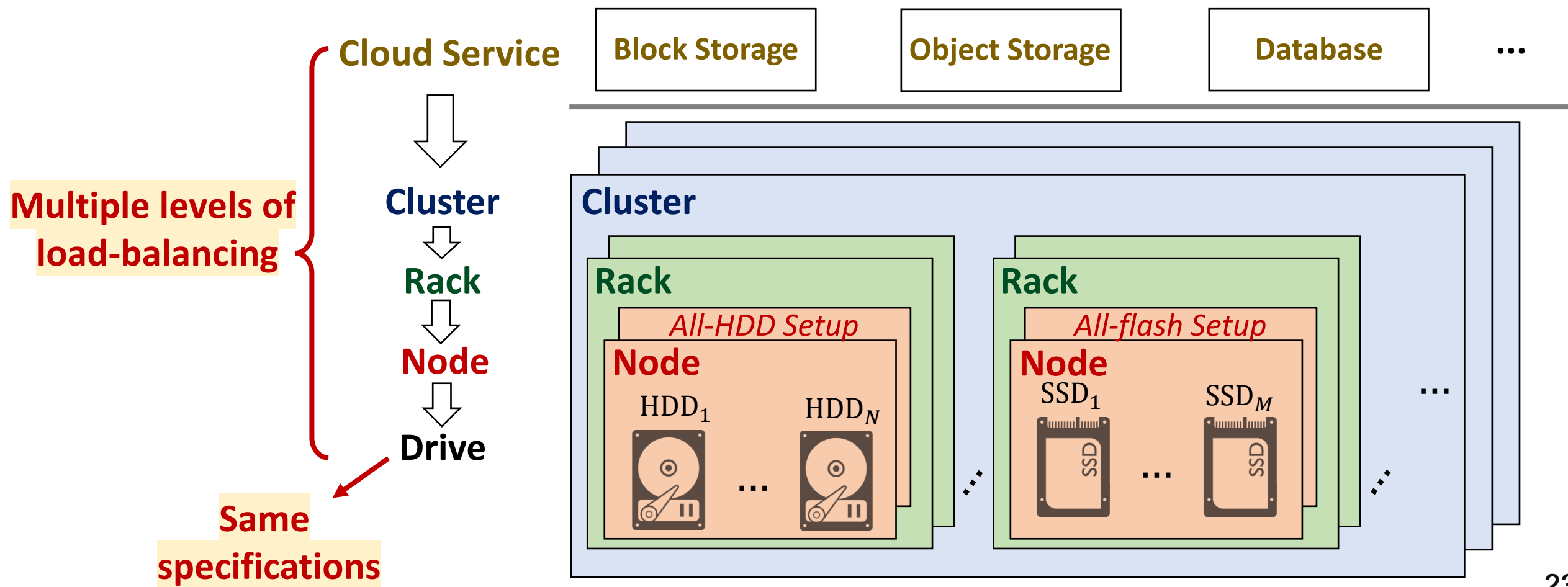


**Linear regression?**

**=> Define a statistically normal drive**

# Design Guidelines (III)

Insight: “Homogeneous setup: similar workload + same drive spec”



Insight: “Homogeneous setup: similar workload + same drive spec”

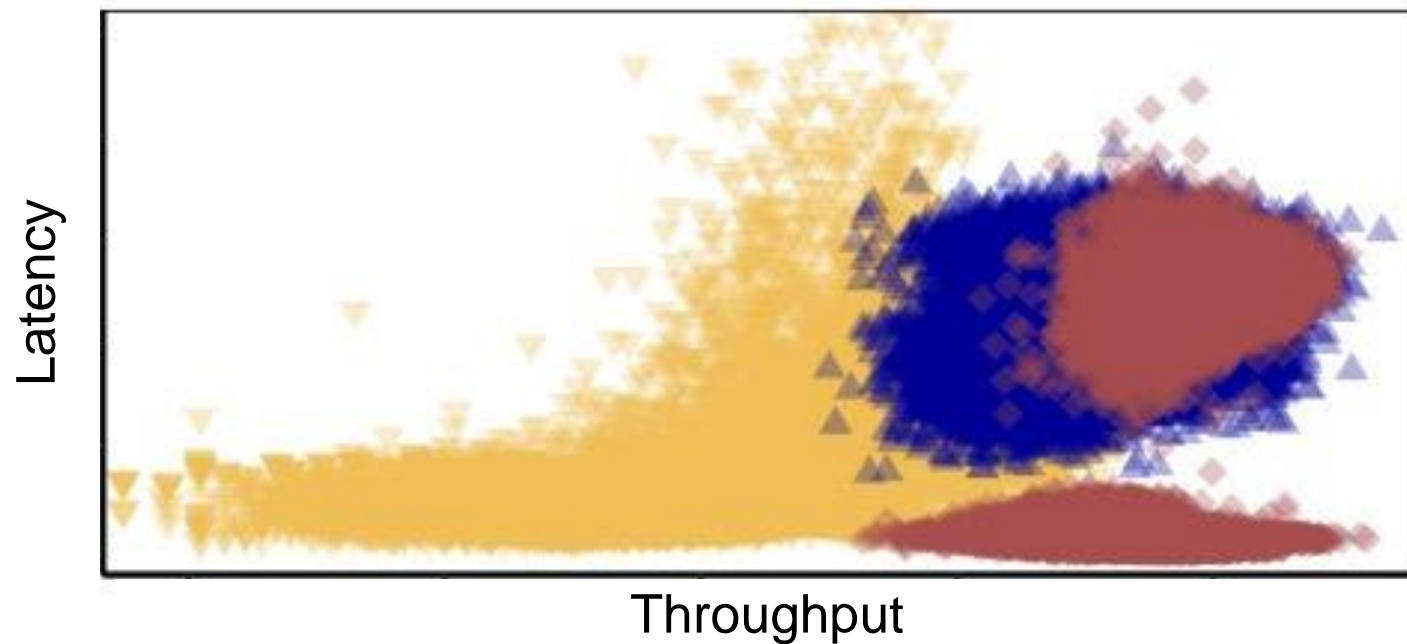
- Need to determine the scope of drives to model
  - Drives from the same service?
  - Drives from the same cluster?
  - Drives from the same node?



# Design Guidelines (III)

Insight: “Homogeneous setup: similar workload + same drive spec”

(a) Same service, different clusters

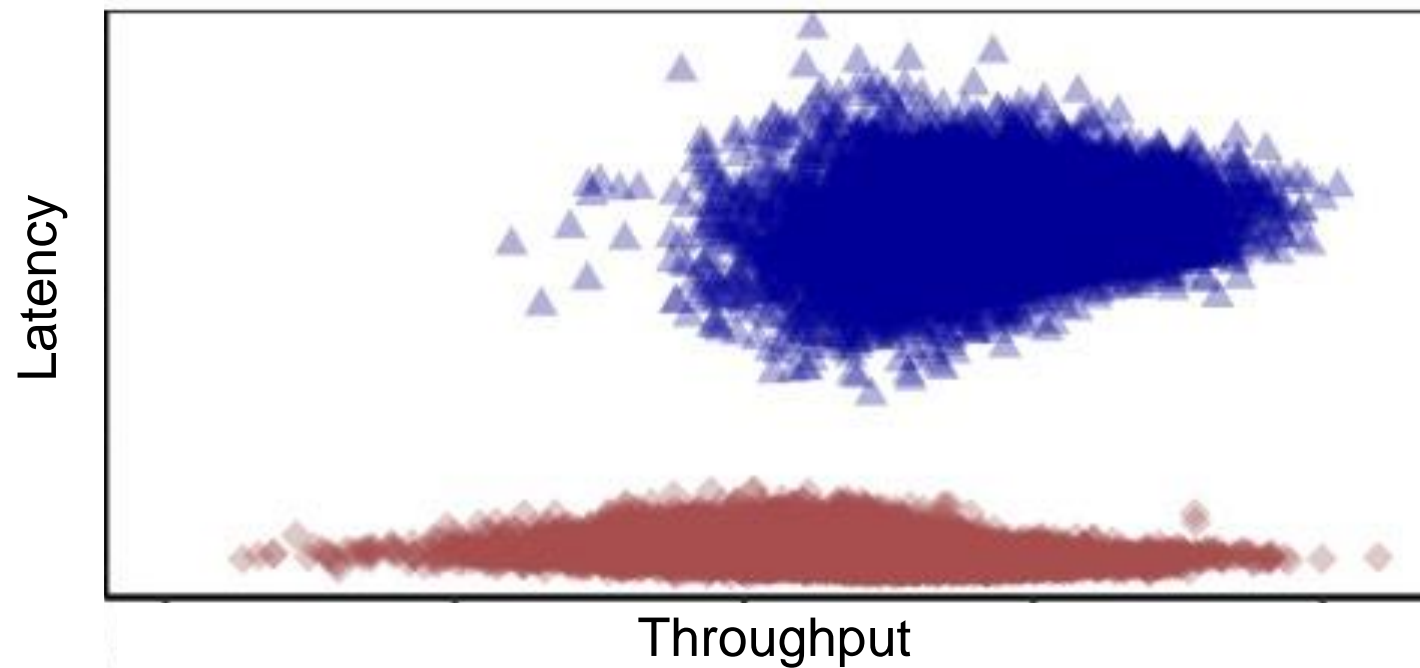


**✗ Distinct LvT distributions**

# Design Guidelines (III)

Insight: “Homogeneous setup: similar workload + same drive spec”

(b) Same cluster, different nodes

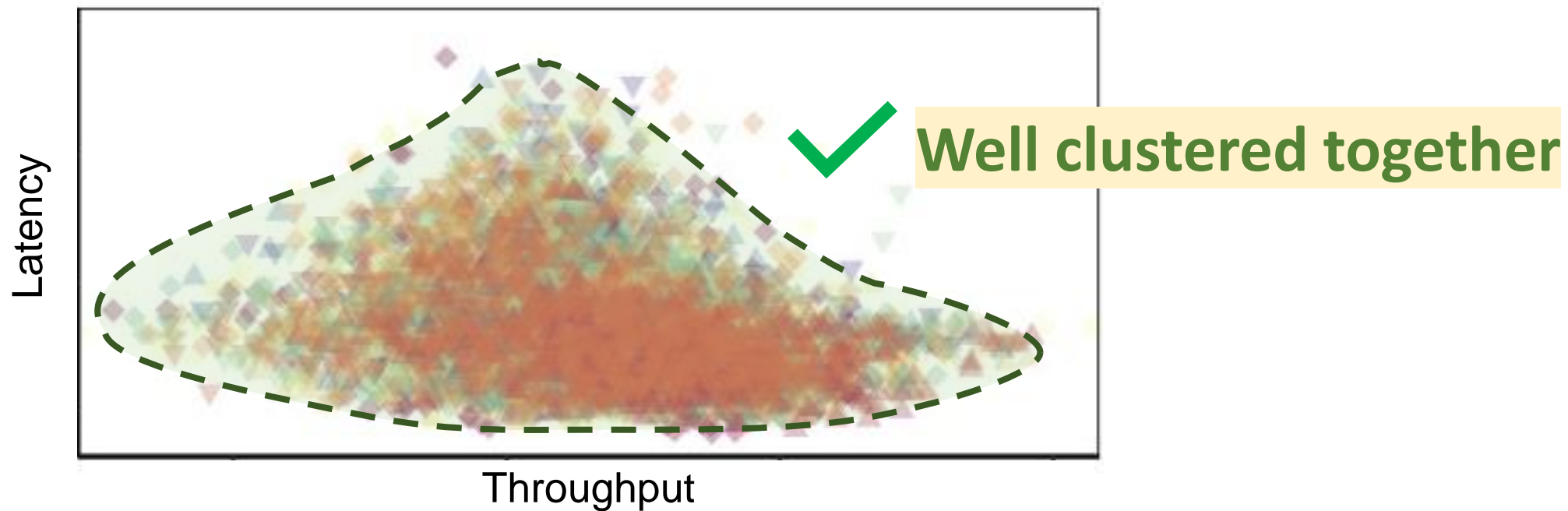


**✗ Distinct LvT distributions**

# Design Guidelines (III)

Insight: “Homogeneous setup: similar workload + same drive spec”

(c) Same node, different drives



**Guideline 3: Use node-wise samples to model LvT distribution**

Insight: “No golden standards to identify fail-slow ”



## Guideline 4: Non-binary output

- Model the likelihood of fail-slow



~~INTRODUCTION~~



~~DATASET~~



~~FAILED  
ATTEMPTS~~

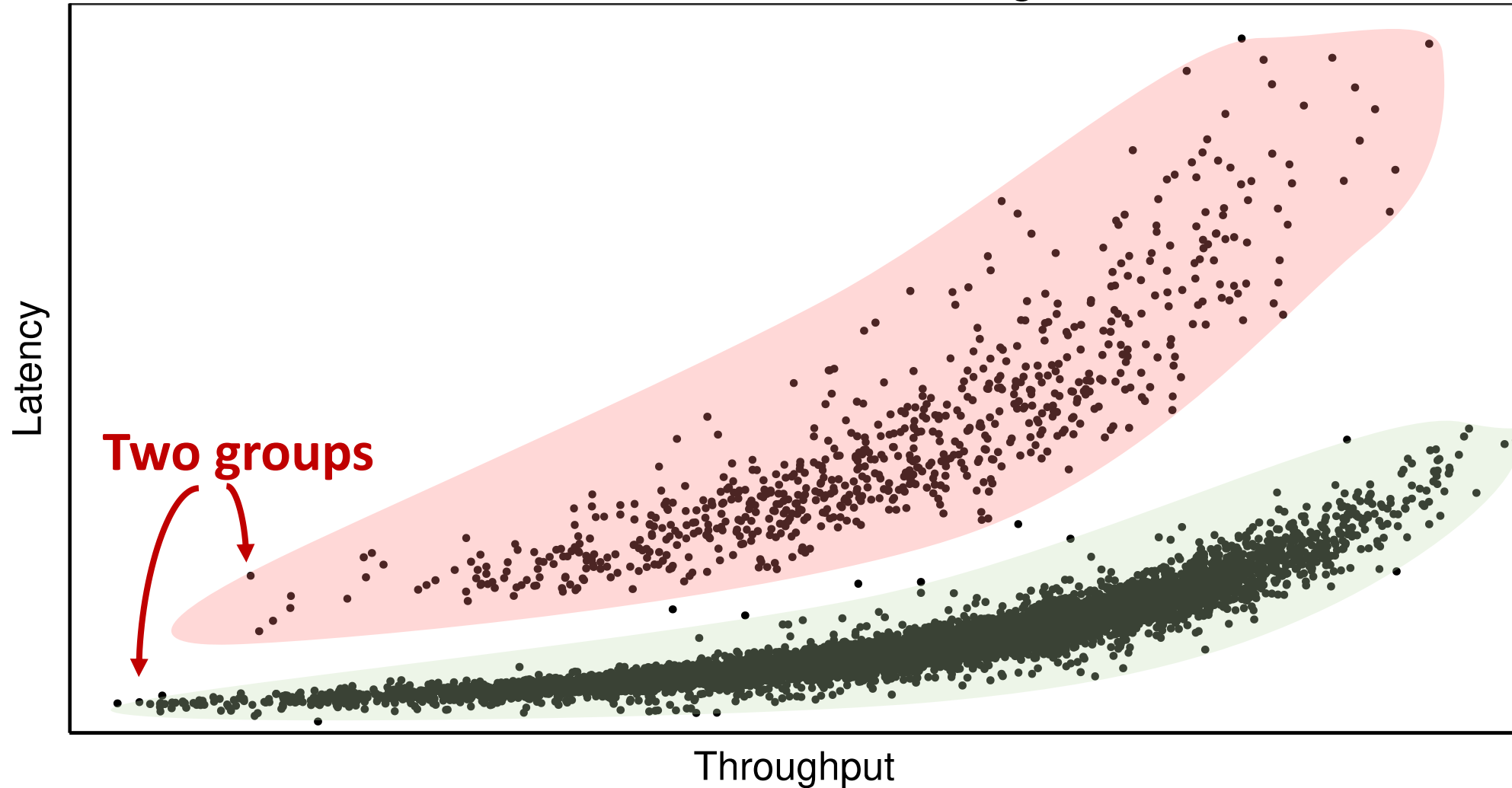


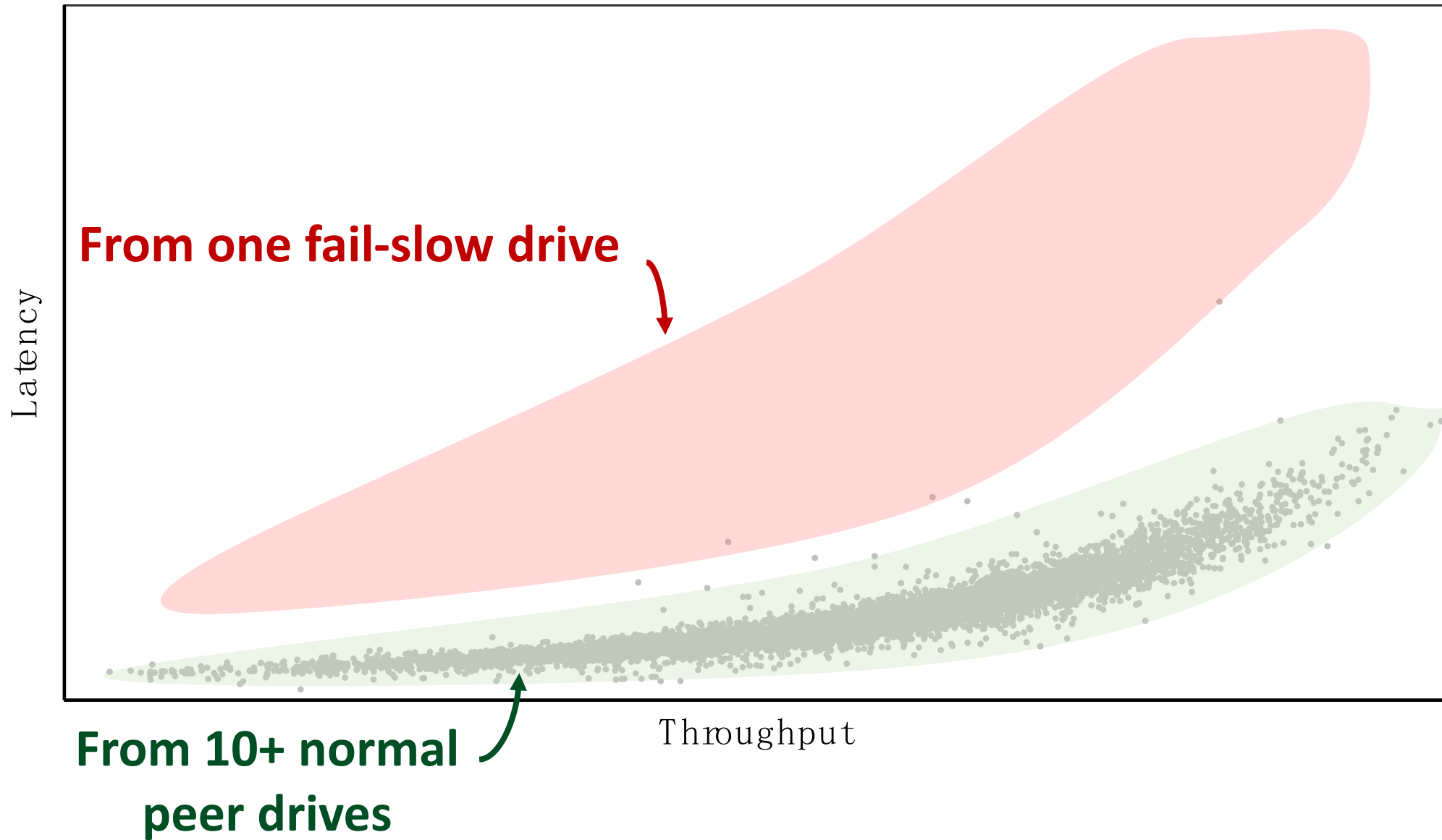
**PERSEUS**



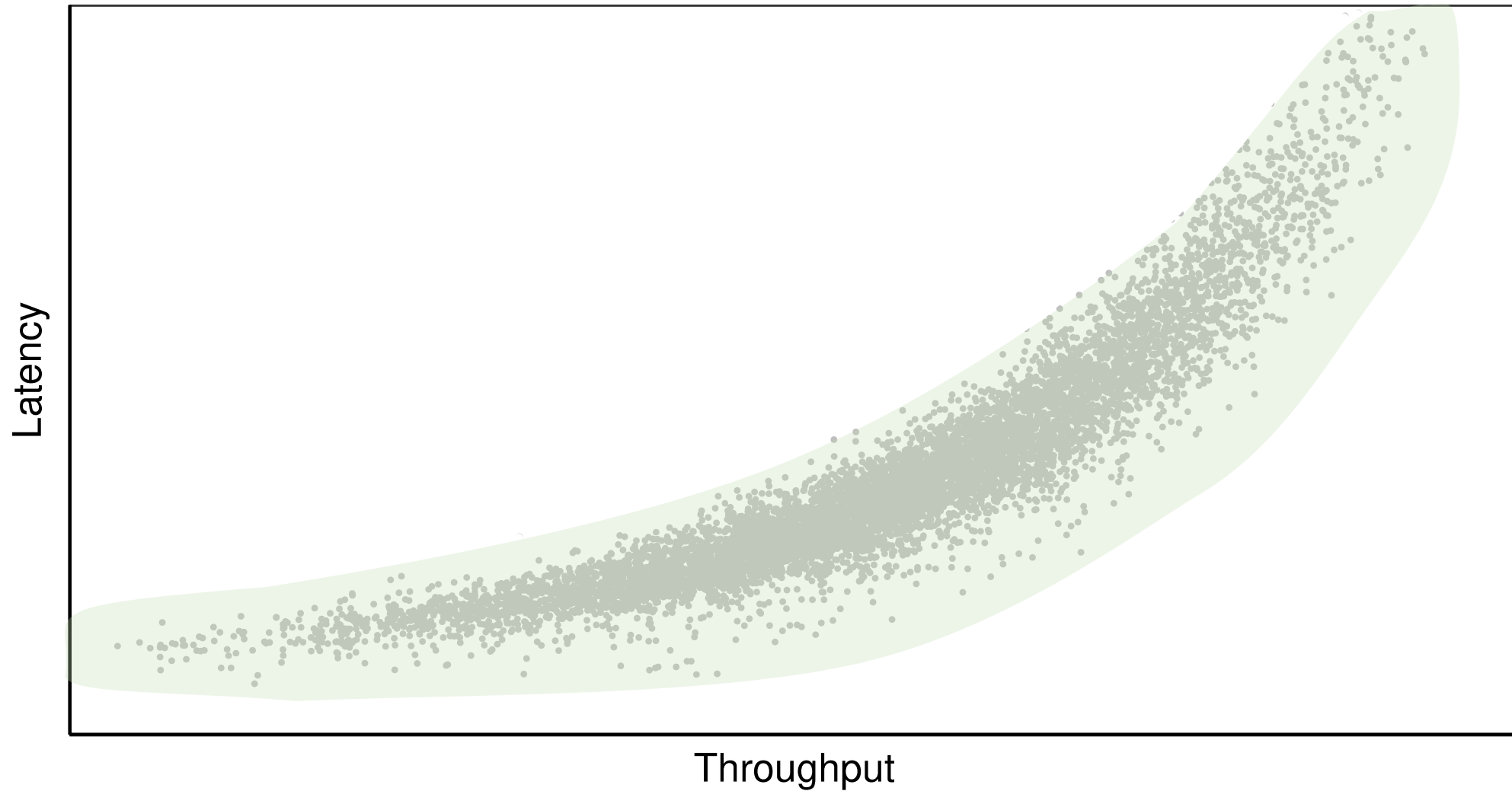
**EVALUATION &  
CONCLUSION**

LvT distribution of one storage node



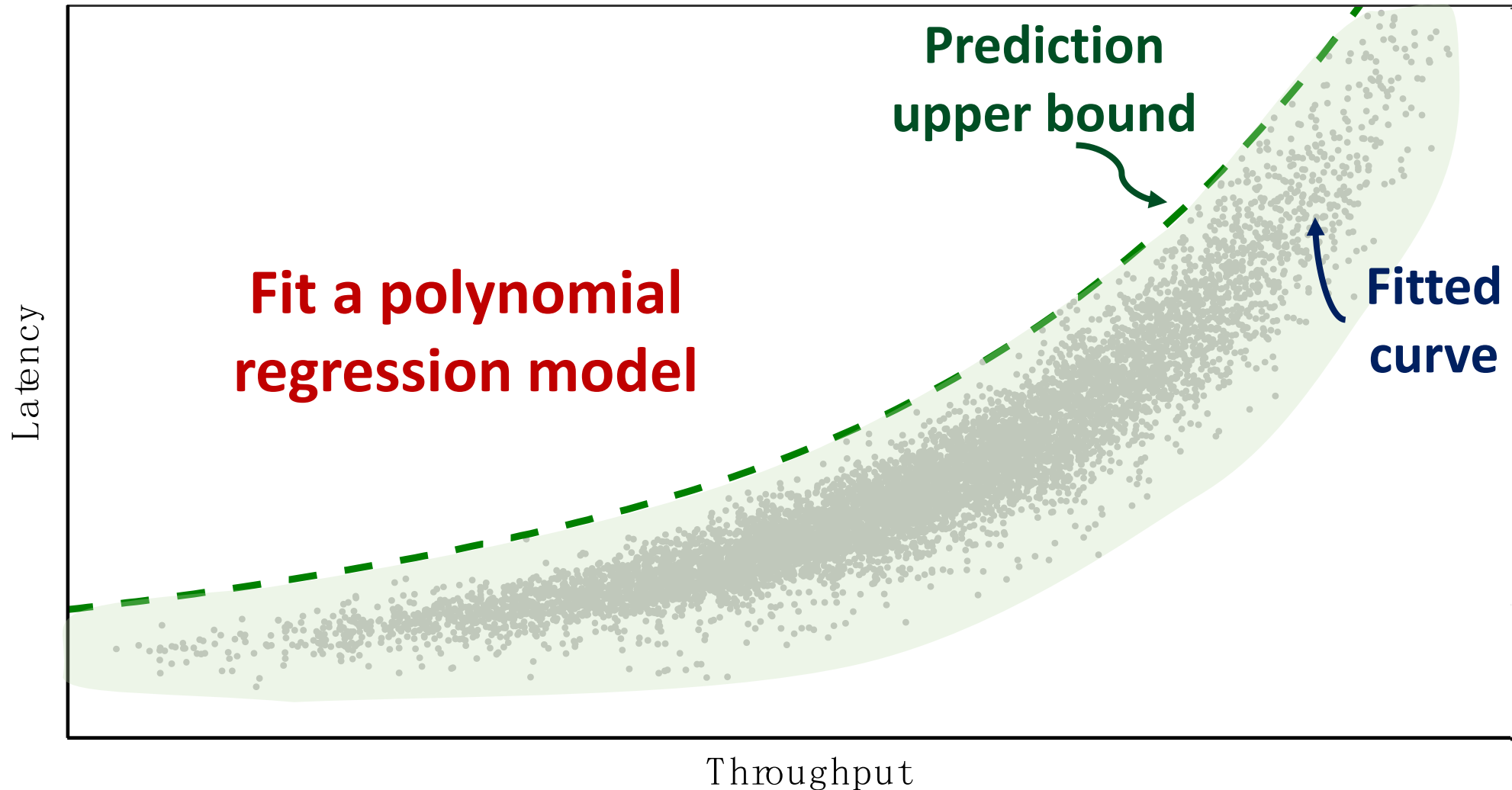


# Step 1: Outlier Detection

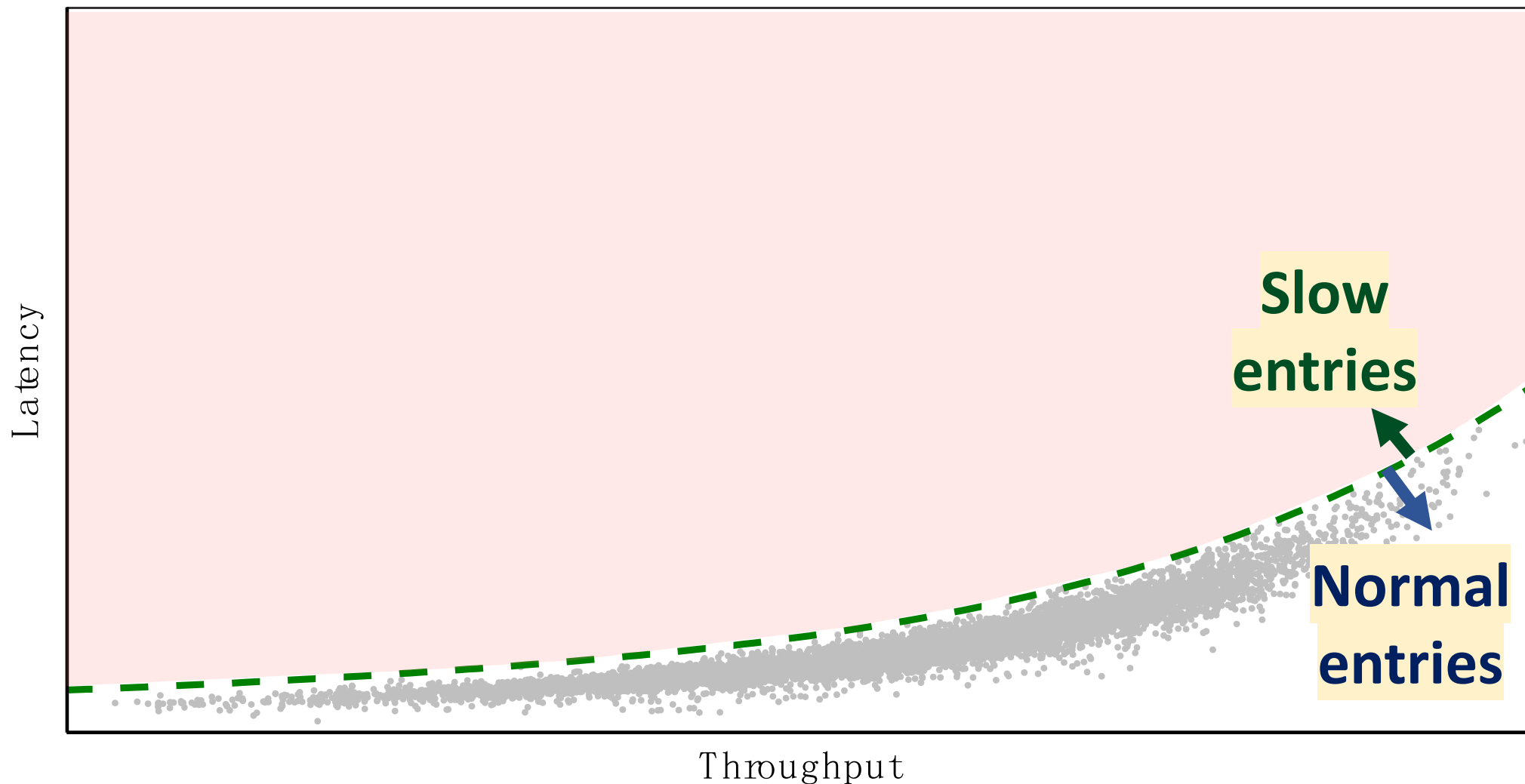




# FAST<sup>1</sup><sub>23</sub> Step 2: Building Regression Model



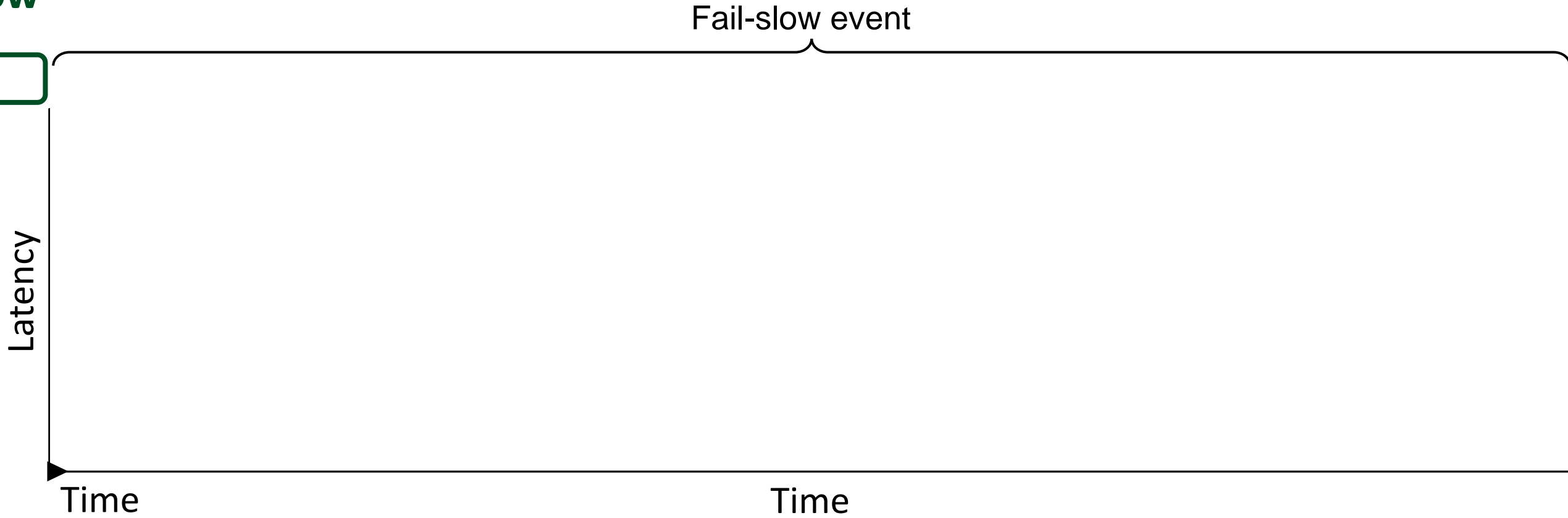
# FAST<sup>1</sup><sub>23</sub> Step 2: Building Regression Model



**Prediction upper bounds** as adaptive latency thresholds without fine-tuning

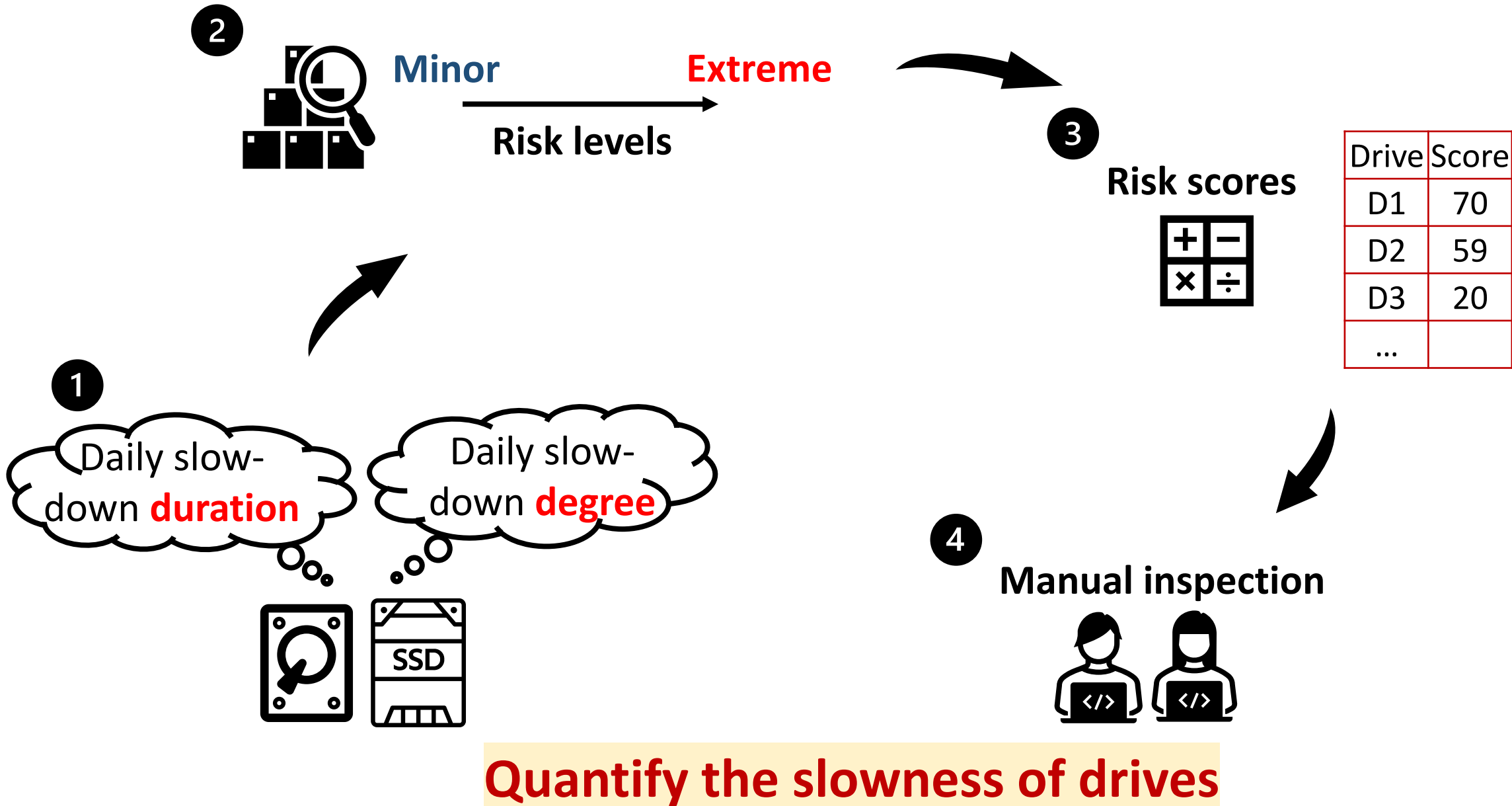
# FAST<sup>1</sup><sub>23</sub> Step 3: Identifying Fail-Slow Event

ow



**Revisit the temporal dimension**

# Step 4: Evaluating Risk





~~INTRODUCTION~~



~~DATASET~~



~~FAILED  
ATTEMPTS~~

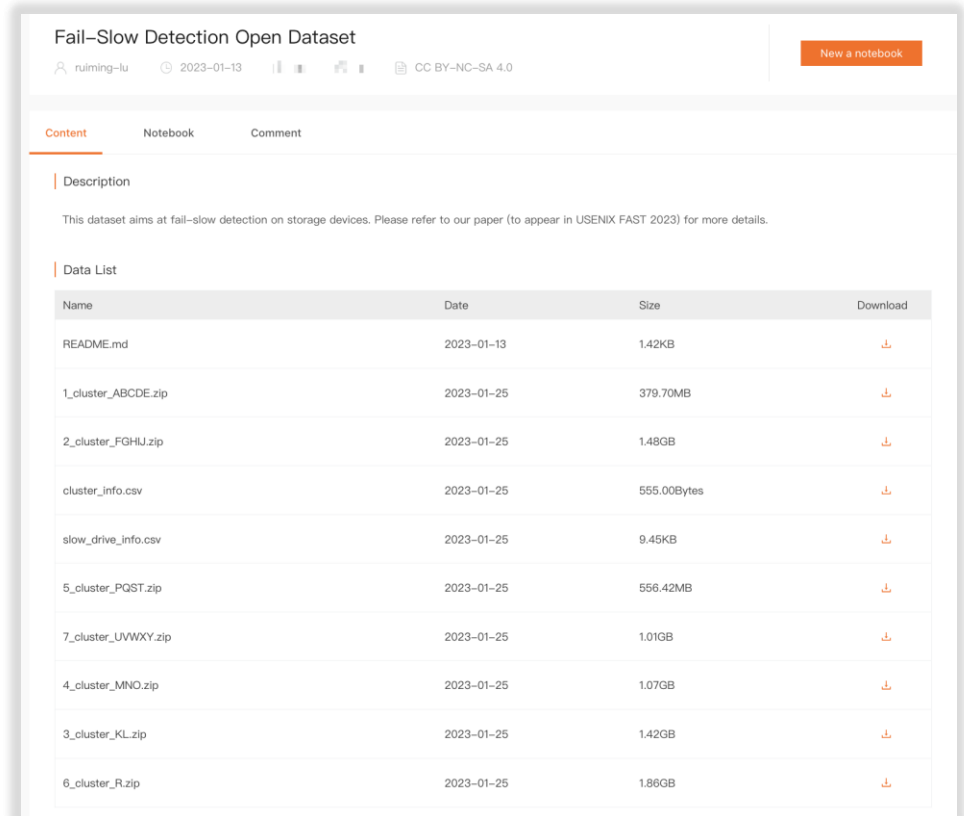


~~PERSEUS~~



**EVALUATION &  
CONCLUSION**

- Built and released our self-assembled test dataset
  - Clear labels (fail-slow or not)
  - 15 days of operational traces
  - 41K drives
  - ~300 fail-slow drives



Fail-Slow Detection Open Dataset

ruiming-lu 2023-01-13 CC BY-NC-SA 4.0

New a notebook

Content Notebook Comment

Description

This dataset aims at fail-slow detection on storage devices. Please refer to our paper (to appear in USENIX FAST 2023) for more details.

Data List

Name	Date	Size	Download
README.md	2023-01-13	1.42KB	<a href="#">Download</a>
1_cluster_ABCDE.zip	2023-01-25	379.70MB	<a href="#">Download</a>
2_cluster_FGHIJ.zip	2023-01-25	1.48GB	<a href="#">Download</a>
cluster_info.csv	2023-01-25	555.00Bytes	<a href="#">Download</a>
slow_drive_info.csv	2023-01-25	9.45KB	<a href="#">Download</a>
5_cluster_PQRST.zip	2023-01-25	556.42MB	<a href="#">Download</a>
7_cluster_UVWXY.zip	2023-01-25	1.01GB	<a href="#">Download</a>
4_cluster_MNO.zip	2023-01-25	1.07GB	<a href="#">Download</a>
3_cluster_KL.zip	2023-01-25	1.42GB	<a href="#">Download</a>
6_cluster_R.zip	2023-01-25	1.86GB	<a href="#">Download</a>

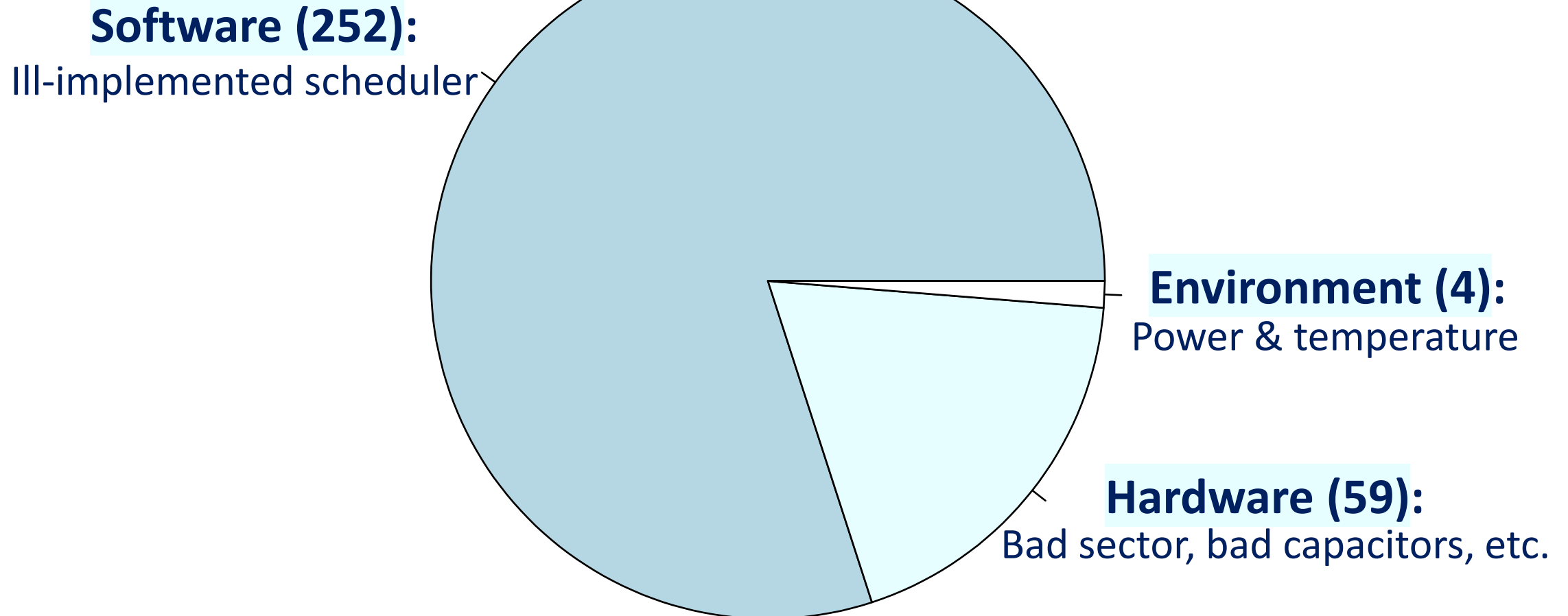
<https://tianchi.aliyun.com/dataset/144479>

- Perseus outperforms all previous attempts (§5.4)
- Effectiveness of Perseus's Design Choices (§5.5)
- Reduce Tail Latency By 31-48% (§5.6)
- Root Cause Analysis (§6)

**More details in the paper!**

# Root Cause Distribution

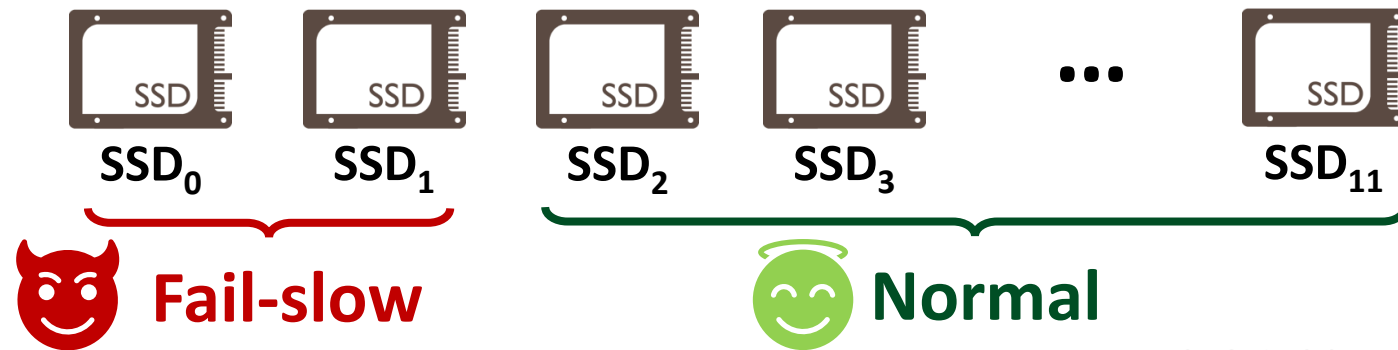
Among 315 Confirmed Fail-Slow Drives:





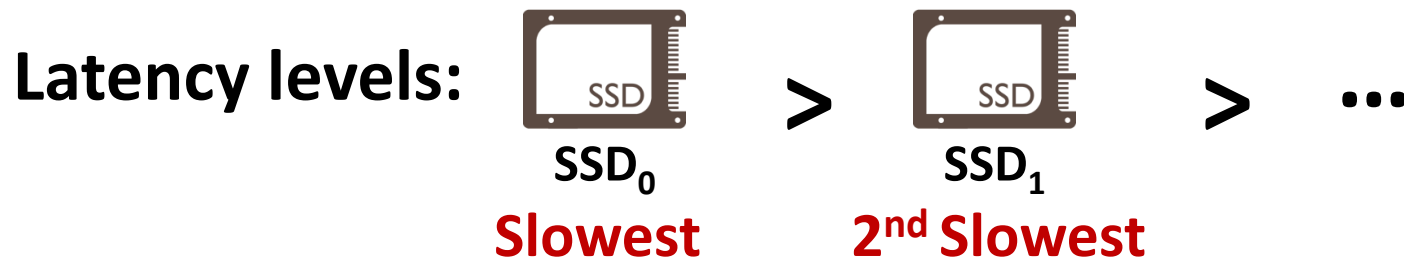
## Case I: In Open-Channel SSD Cluster

1. Every node always has 2 fail-slow drives

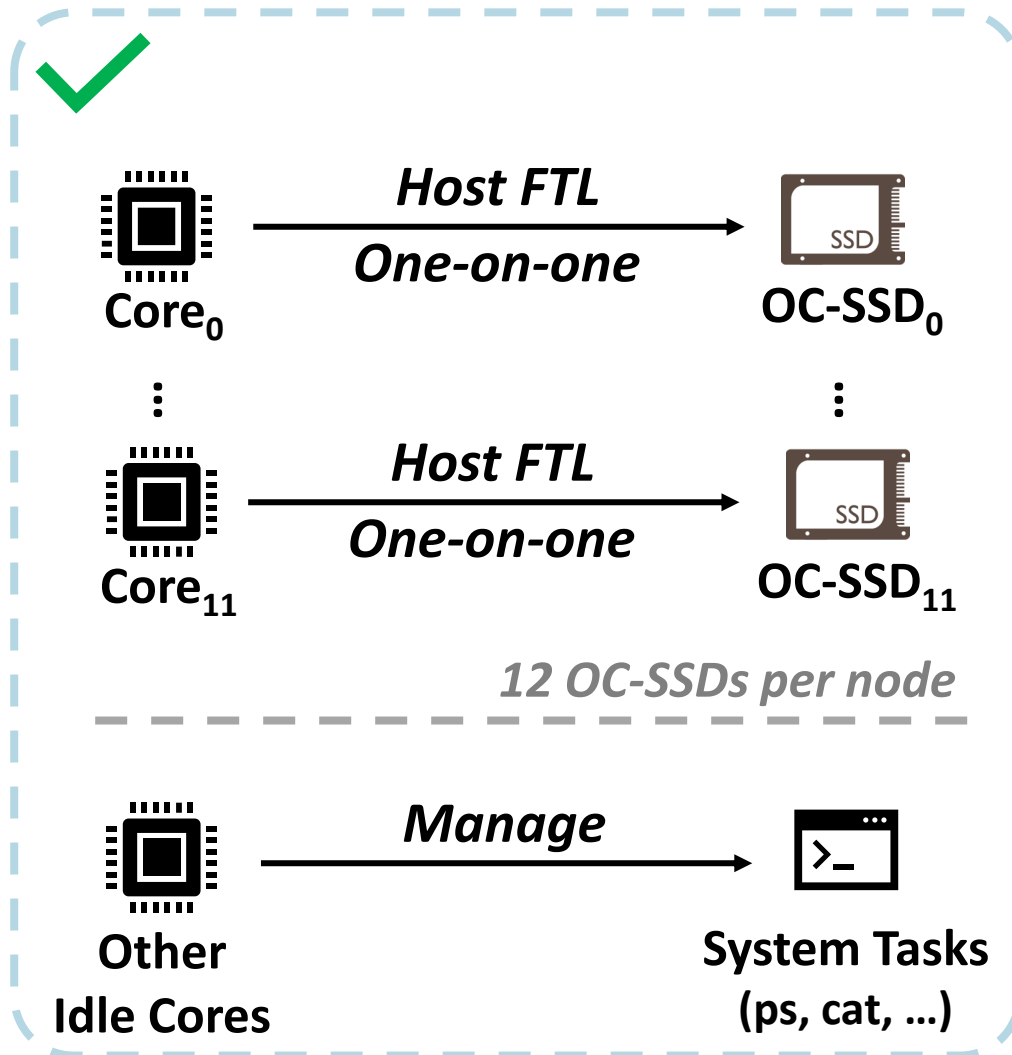


*12 OC-SSDs per node*

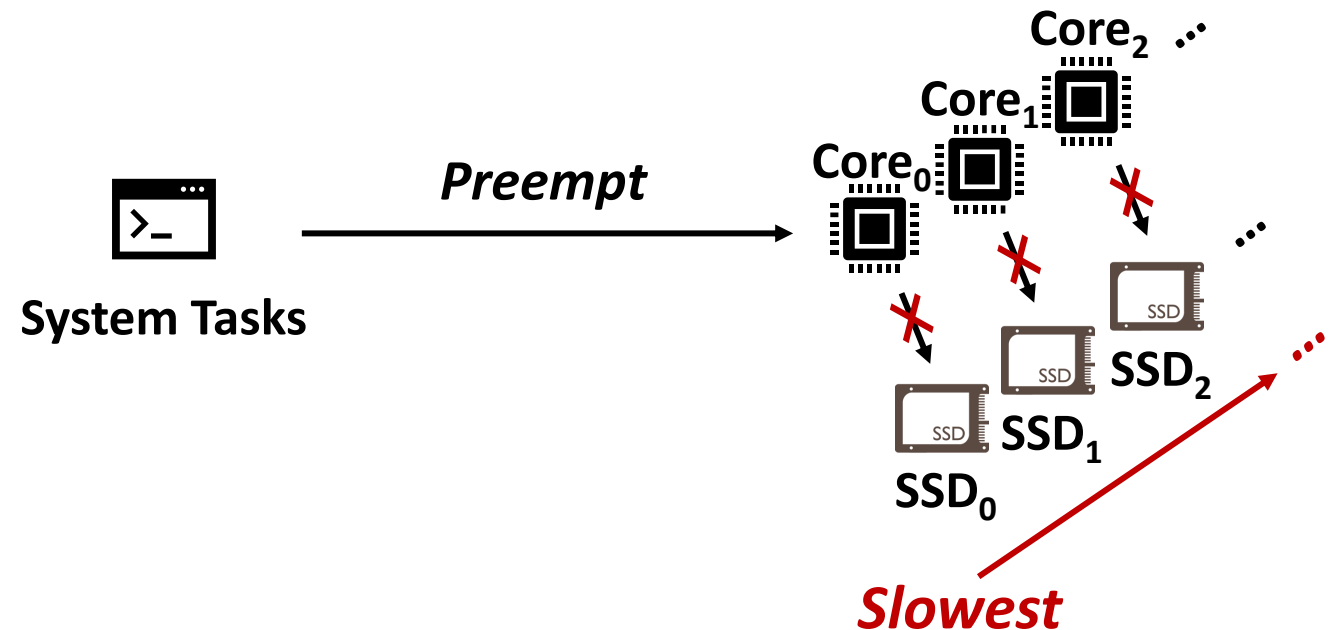
2. Latency levels follow ascending order of logical IDs



## Case I: In Open-Channel SSD Cluster



✗ *If No Idles Cores Available...*



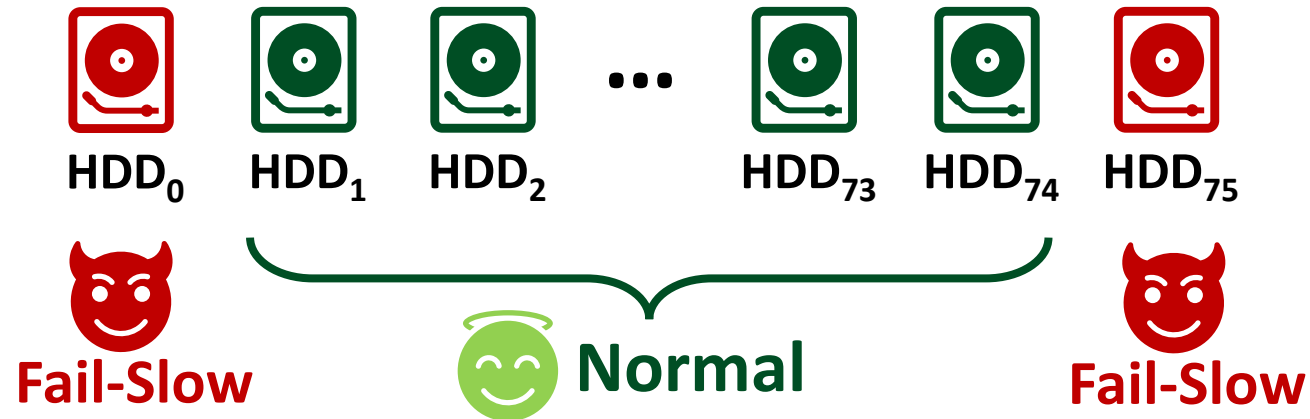
**=> Fail-Slow OC-SSDs in Ascending Order** ⚡

## Case II: In All-HDD Cluster

### 1. Fail-slow drives always appear in fixed pairs

----- 76 HDDs per node

*Two fail-slow disks in a node:*

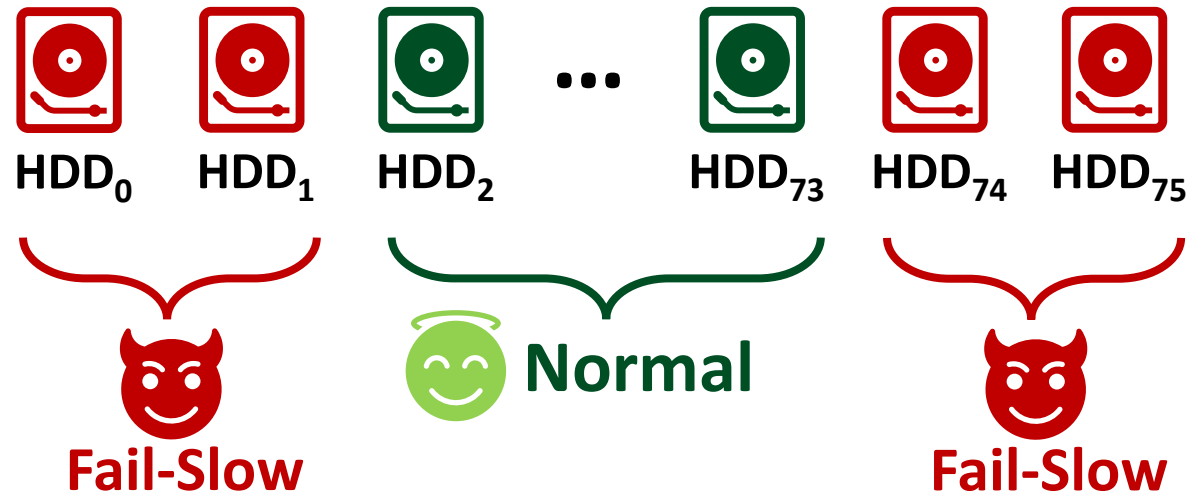


## Case II: In All-HDD Cluster

### 1. Fail-slow drives always appear in fixed pairs

----- 76 HDDs per node

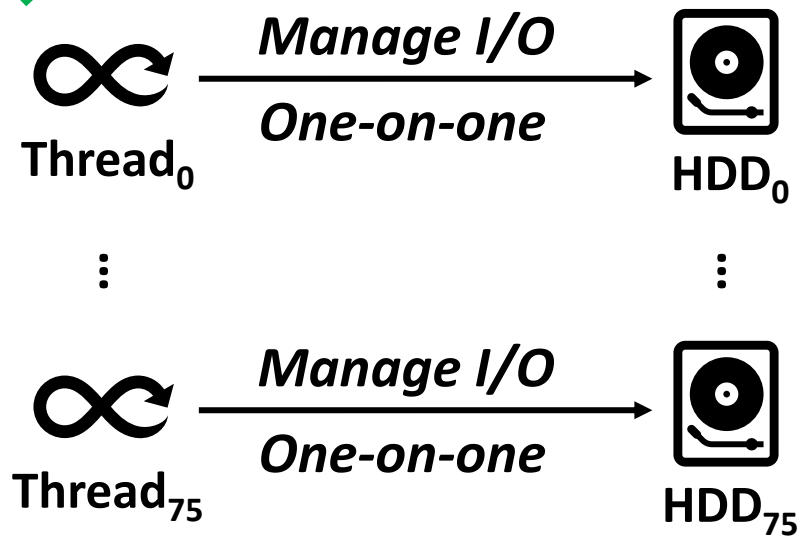
*Four fail-slow disks in a node:*



## Case II: In All-HDD Cluster

1. Fail-slow drives always appear in fixed pairs
2. All fail-slow drives are experiencing similar slowdown
3.  $\#Fail\text{-}slow = 2 \times \#Offline$

## Case II: In All-HDD Cluster



76 HDDs per node

$$\text{Thread}_{ID} = \text{Disk}_{ID} \bmod \#Drives$$



**When Disks Taken Down For Repair...**

e.g., HDD<sub>20</sub> taken down #Drives = 75



HDD<sub>0</sub>

$$0 \bmod 75 = 0$$



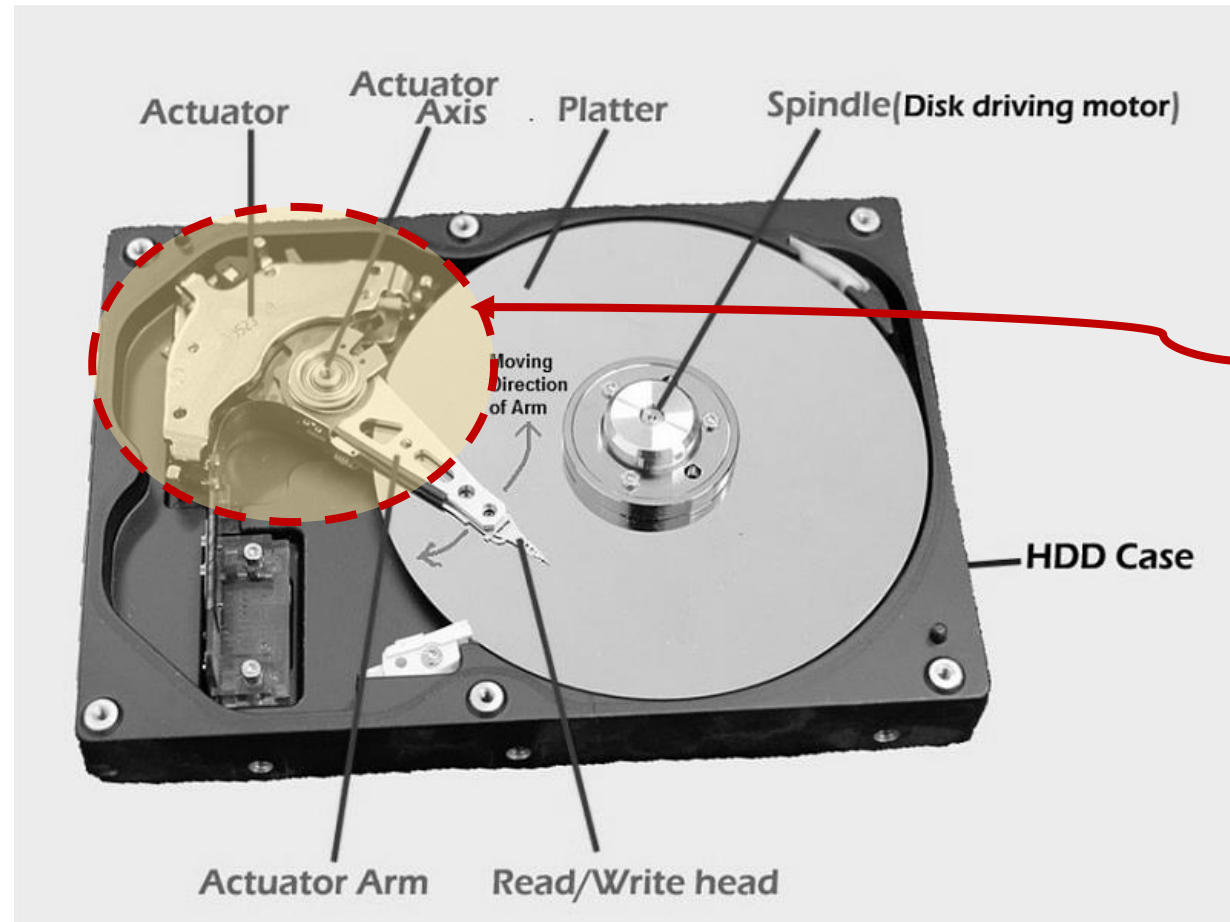
HDD<sub>75</sub>

$$75 \bmod 75 = 0$$

Thread<sub>0</sub> (represented by an infinity symbol)

**=> Resource Contention-Induced Fail-Slow**

- Rotor Eccentricity



**Frequently miss targeted positions**

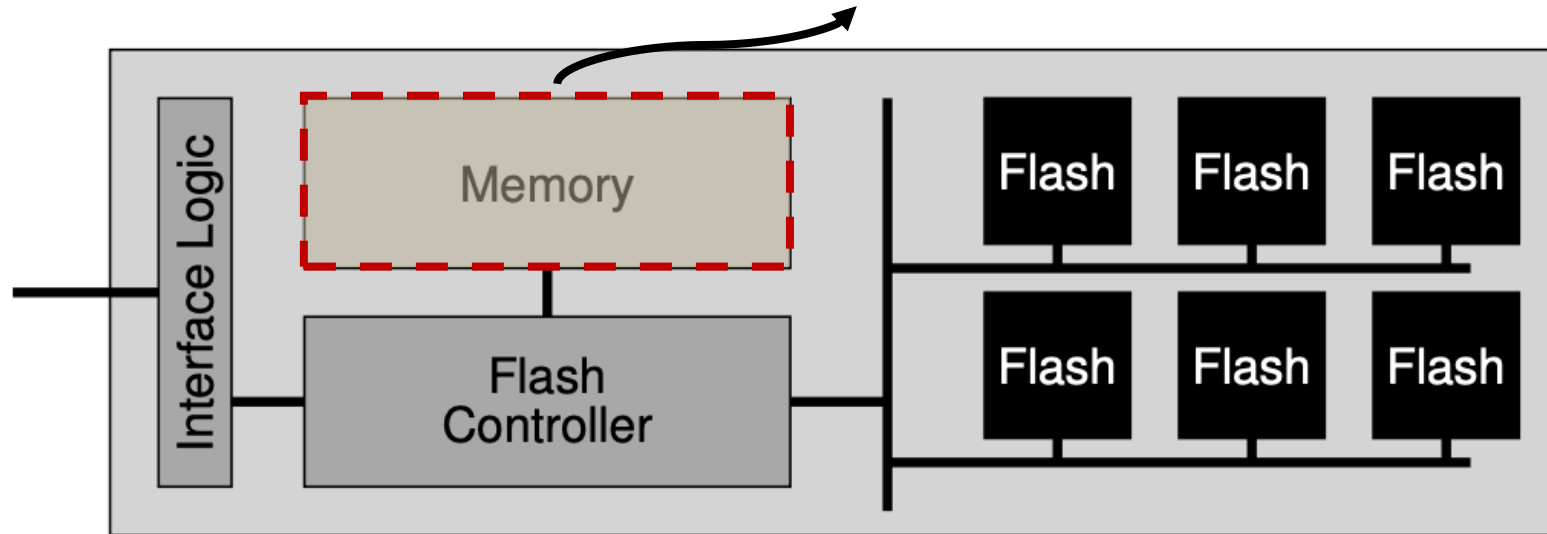


**I/O delay**

Source: <https://www.techintangent.com/hard-disk-description/>

- Rotor Eccentricity
- Bad Capacitors

## DRAM as an internal write-back cache



Source: Operating Systems: Three Easy Pieces

**DRAM capacitors failed => Delayed writes**



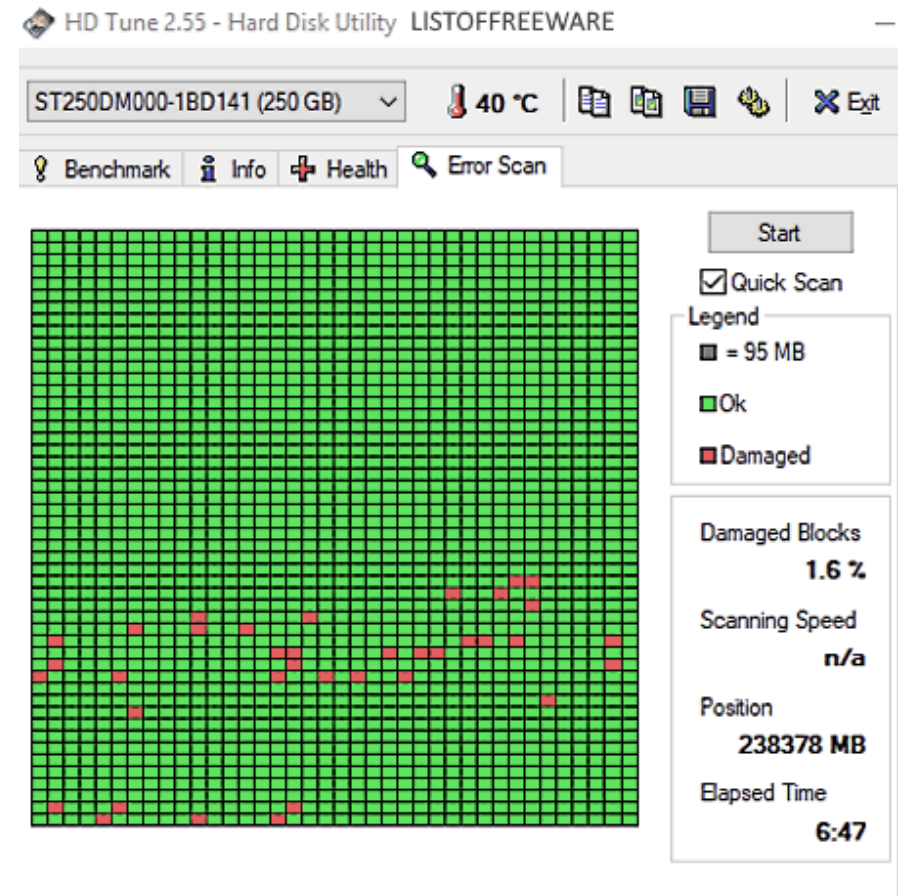
# Root Cause: Hardware

- Rotor Eccentricity
- Bad Capacitors
- Bad Sectors:

**Data reallocate to spare sectors**



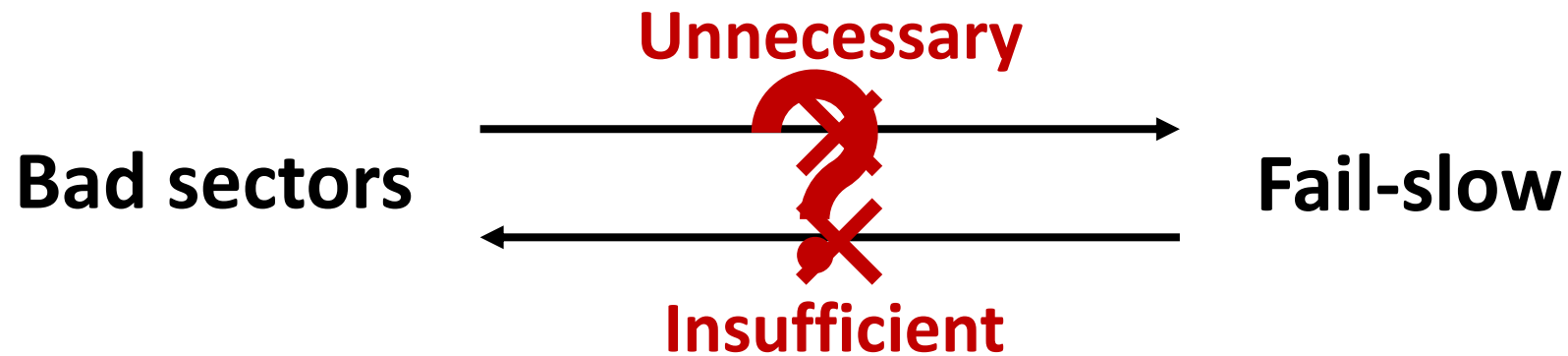
**Long seek time**



Source: <https://listoffreeware.com/best-free-software-to-check-hard-drive-bad-sectors-windows/>

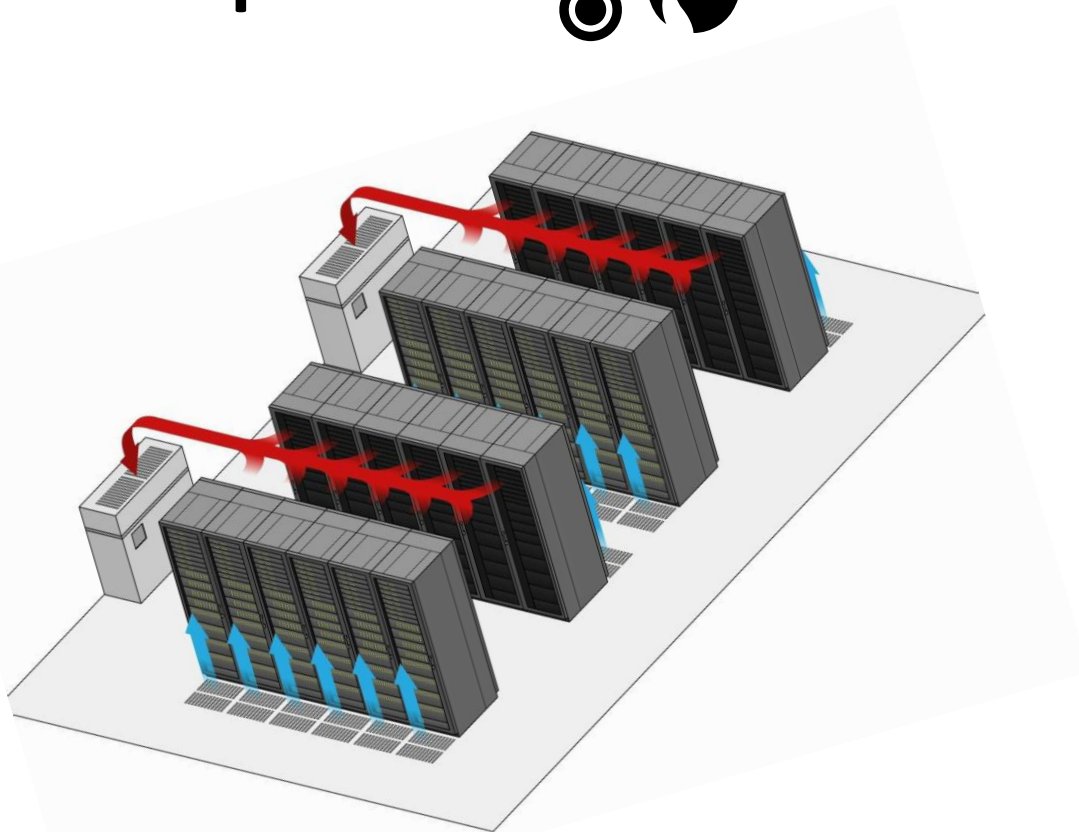
# Root Cause: Hardware

- Rotor Eccentricity
- Bad Capacitors
- Bad Sectors:



# Root Cause: Environment

- Temperature 🌡️ 🔥



Source: <https://www.upsite.com/blog/helping-your-data-center-breathe-easier-with-good-air-flow-management/>

- Power 🔌 ⚖️



Source: <https://www.ecsintl.com/how-data-centers-can-effectively-manage-power-surges/>

# Perseus

---

## Detection Framework

Efficient



**Non-intrusive**

(Performance) log-based

No source code altering

Fail-Slow

Detection



**Fine-grained**

Device-level detection



**Accurate**

Recall/precision rate > 0.99



**General**

One set of parameters fits all scenarios

---

## Storage Devices



...

**Adaptable to Other Problem Domains**

# Thank you!

Perseus: A Fail-Slow Detection Framework  
for Cloud Storage Systems

**Ruiming Lu**, Erci Xu, Yiming Zhang, Fengyi Zhu, Zhaosheng Zhu,  
Mengtian Wang, Zongpeng Zhu, Guangtao Xue, Jiwu Shu, Minglu Li, Jiesheng Wu

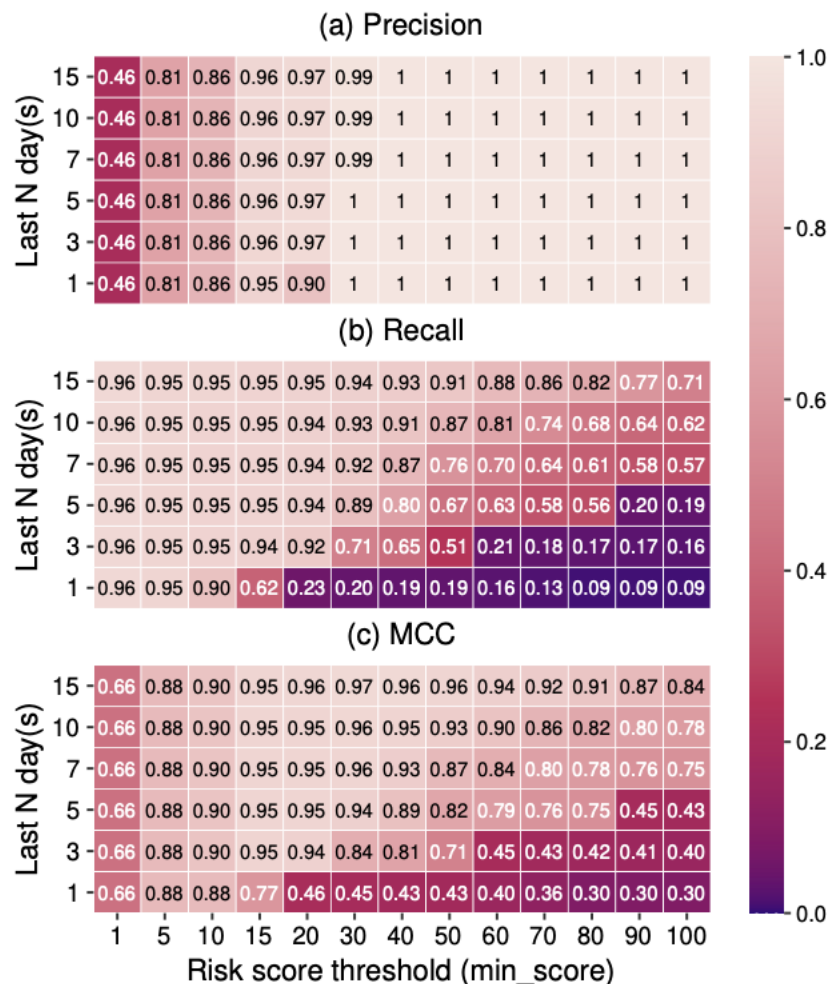
Contact email: [lrn318@sjtu.edu.cn](mailto:lrn318@sjtu.edu.cn)

## Effectiveness of Perseus's design choices (§5.5)

Parameter	Range	Description
S1: Outlier detection (§4.2)		
<i>PCA</i>	On/Off	Transform the coordinates w.r.t. the principal components.
<i>DBSCAN</i>	On/Off	Density-based outlier detection.
S3: Identifying fail-slow event (§4.4)		
<i>X</i>	95~99.9	Use the <i>X</i> % prediction upper bound as the latency upper bound.
S4: Evaluating risk (§4.5)		
<i>min_score</i>	1~100	Risk score threshold.
<i>N</i>	1~15	Evaluate the risk score of the most recent <i>N</i> days.

Metric	w/o Outlier	w/o PCA	p95	p99	p999	Deployed
Full-set						
Precision	0.98	0.55	0.99	1.00	1.00	0.99
Recall	0.51	0.43	0.99	0.93	0.93	1.00
MCC	0.71	0.49	0.99	0.96	0.96	0.99
Subset (excluding software-induced)						
Precision	0.95	0.36	0.94	0.98	1.00	0.94
Recall	0.82	0.91	0.95	0.92	0.95	1.00
MCC	0.88	0.57	0.95	0.95	0.98	0.97

## Perseus's design tradeoffs (§5.5)



Metric	Thresh-Stat	Thresh-Emp	Peer Eval	IASO-Based	PERSEUS-Deployed
Full-set					
Precision	1.00	1.00	0.98	0.48	0.99
Recall	0.52	0.02	0.57	0.24	1.00
MCC	0.72	0.14	0.74	0.32	0.99
Subset (excluding software-induced)					
Precision	1.00	1.00	1.00	0.45	0.94
Recall	0.71	0.09	0.65	0.61	1.00
MCC	0.84	0.30	0.80	0.52	0.97