# 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

- **Failures ⚠ in The Wild**
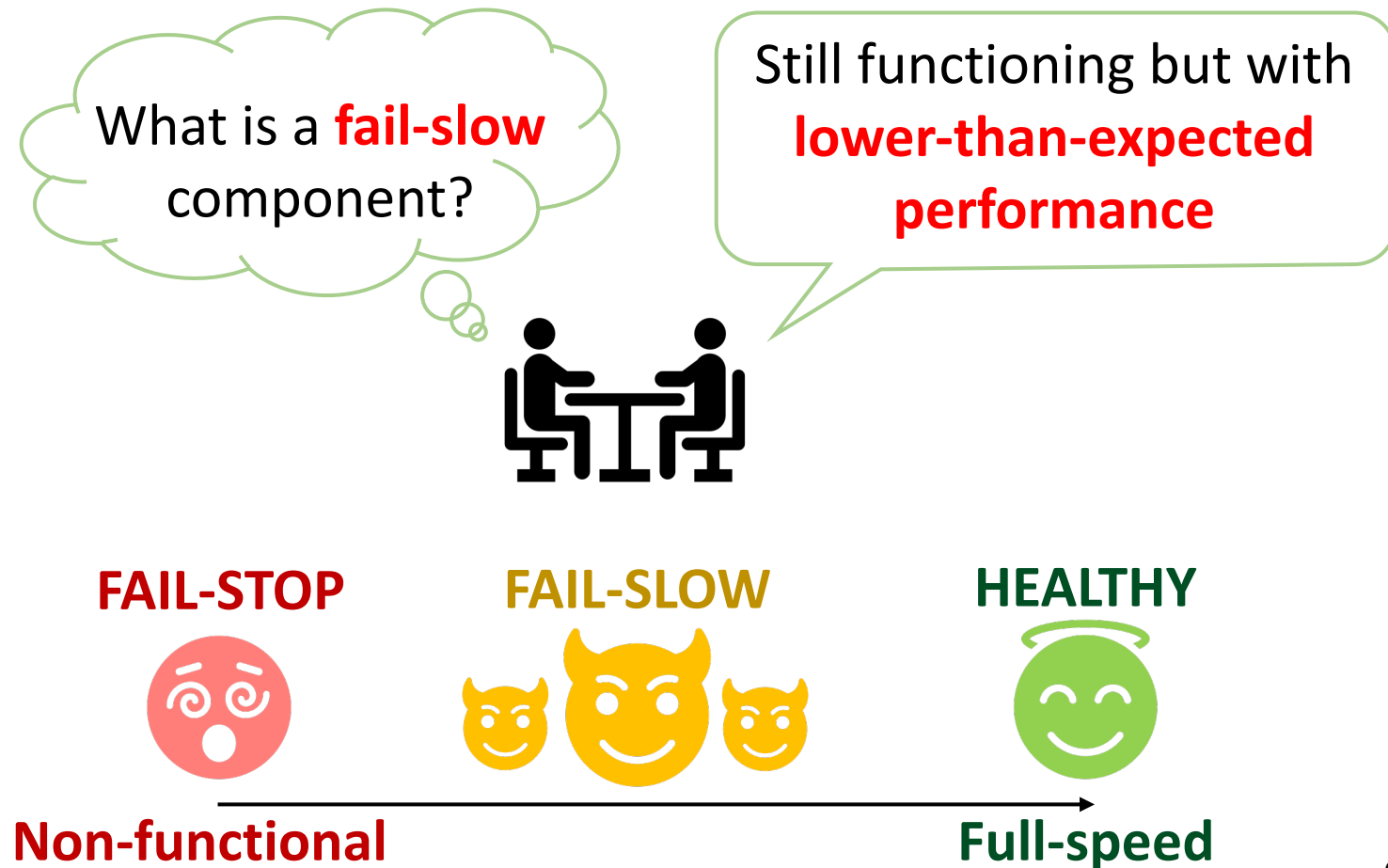
  - **Fail-Slow** 🔍
  - Fail-Stop
  - Byzantine
  - …

What is a **fail-slow** component?

Still functioning but with **lower-than-expected performance**

**FAIL-STOP**  **FAIL-SLOW**  **HEALTHY**

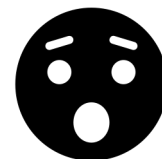**Non-functional**  **Full-speed**

**FAIL-SLOW**

Severe

Not Uncommon

Confusing

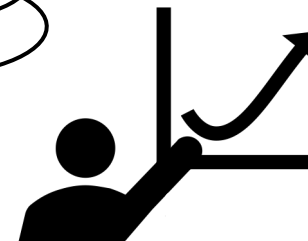"Fail-slow NVMe SSDs can **degrade to SATA SSD or HDD-level performance**[1]"

Annual fail-slow failure rate is **1-2%**[2]!
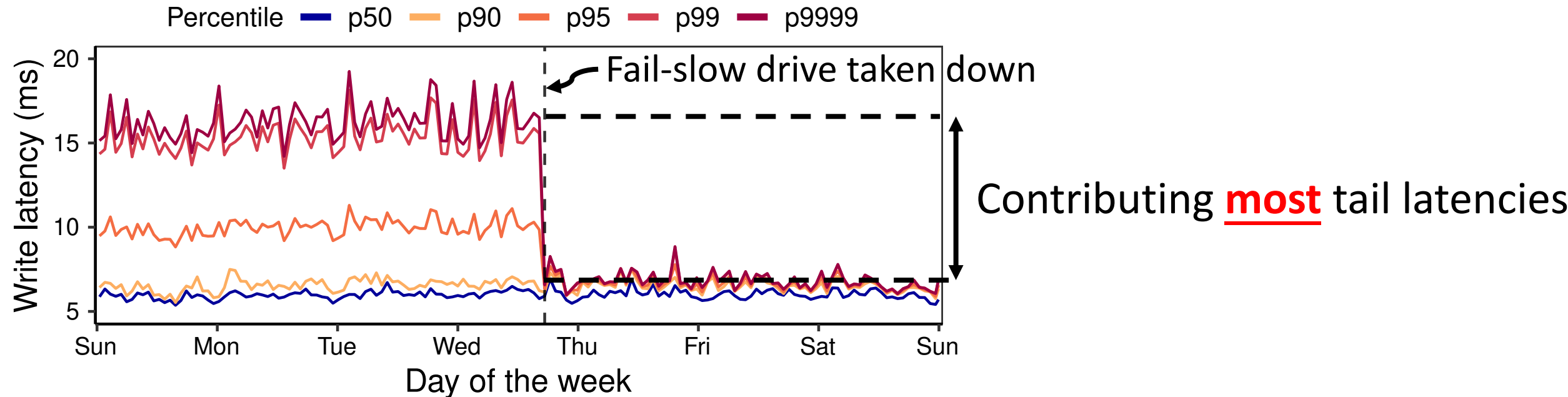
**As frequent as fail-stop** incidents!

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

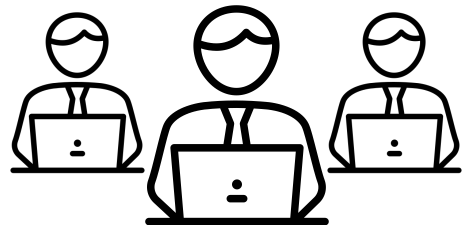[1] NVMe SSD Failures in the Field: the Fail-Stop and the Fail-Slow, Lu et al.
[2] IASO: A Fail-Slow Detection and Mitigation Framework for Distributed Storage Services, Panda et al.

3

- **No Ground Truth in <u>Identifying Fail-Slow</u>**

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

>100us?

>500us?

>1ms?

⋮

A "thousand" ways to define!

# Fail-Slow Detection (FSD)

- **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:

# Fail-Slow Detection (FSD)

- **Our Work Shares**

  - Years of Experiences in FSD

  - A Practical FSD Framework named **Perseus**

  - Root Cause Analysis

~~INTRODUCTION~~        **DATASET**        FAILED ATTEMPTS        PERSEUS        EVALUATION & CONCLUSION

# Our Dataset

- **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

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

**Threshold are hard to decide**

(d)

(c)

Latency

(b)

(a)

$t_0$      $t_1$      $t_2$

Time

**Dilemma** { Relaxed => More False Positives

Strict   => Many Fail-Slow Undiscovered

**Sliding window**

$Event_3$

$Event_1$

$Event_2$

$Event_4$

$Event_5$

Latency

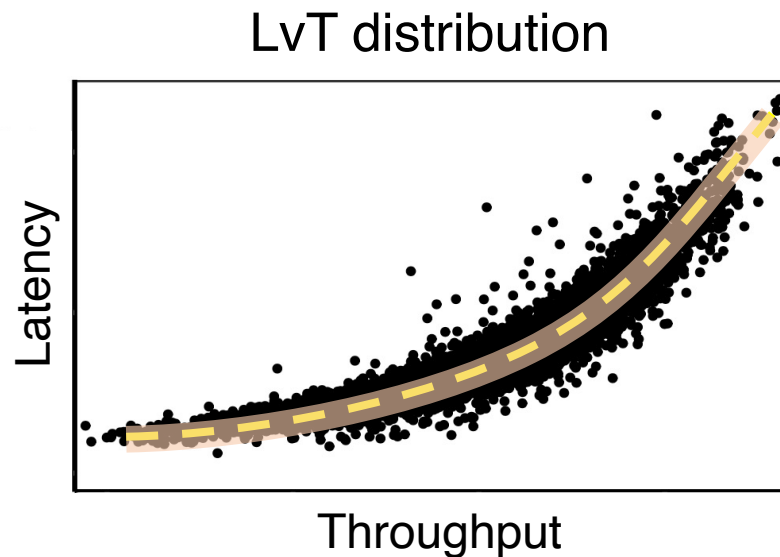**Most peers are normal and stable => Let their peers decide!**

Time

**Time-consuming to fine-tune**

**Insight: "Workload pressure can affect latency variations"**

- How to model such a <u>positive correlation</u>?

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

LvT distribution



**Linear regression?**
=> Define a *statistically normal* drive

**Insight: "No golden standards to identify fail-slow "**

FAIL-STOP       FAIL-SLOW       HEALTHY

Non-functional                     Full-speed

**Guideline: Non-binary output**

- Model the likelihood of fail-slow

LvT distribution of one storage node

**From one fail-slow drive**

Latency

Throughput

**From 10+ normal peer drives**

23

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

Fail-slow event

Latency

Time

Time

**Revisit the temporal dimension**

**2** Minor → Extreme

**Risk levels**

**3** Risk scores

| Drive | Score |
|-------|-------|
| D1 | 70 |
| D2 | 59 |
| D3 | 20 |
| ... | |

**1** Daily slow-down **duration**  Daily slow-down **degree**

SSD

**4** Manual inspection

**Quantify the slowness of drives**
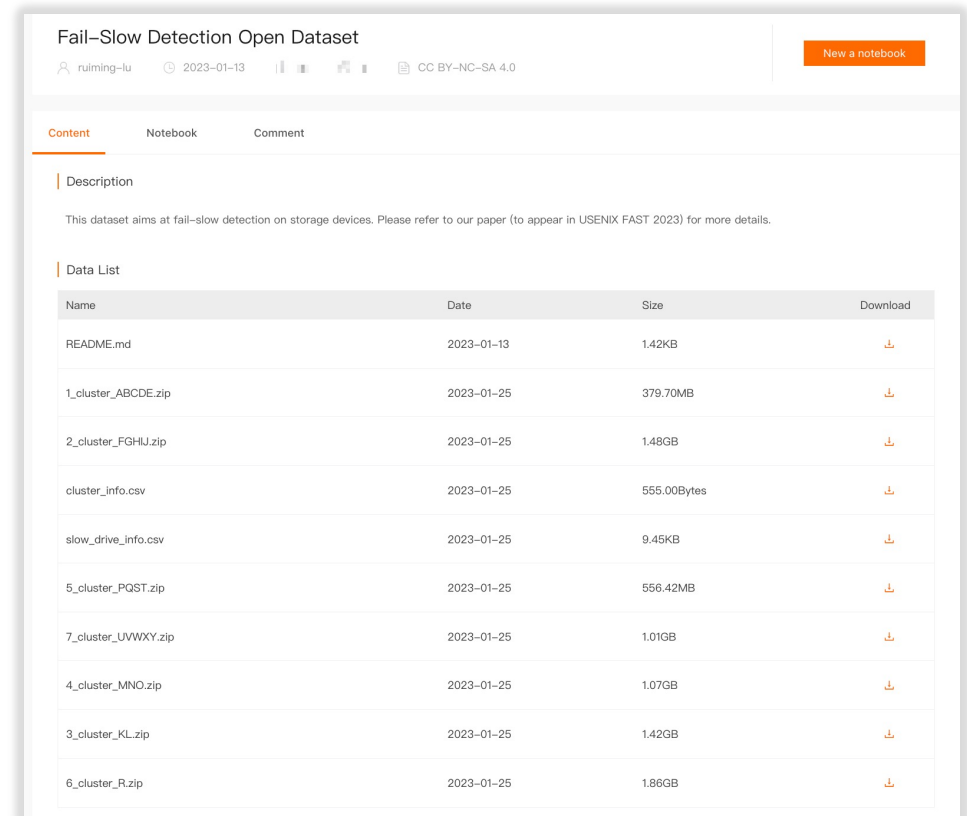
INTRODUCTION

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



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

# Evaluations

- **Perseus outperforms all previous attempts (§5.4)**

- **Effectiveness of Perseus's Design Choices (§5.5)**

- **Reduce Tail Latency By 31-48% (§5.6)**
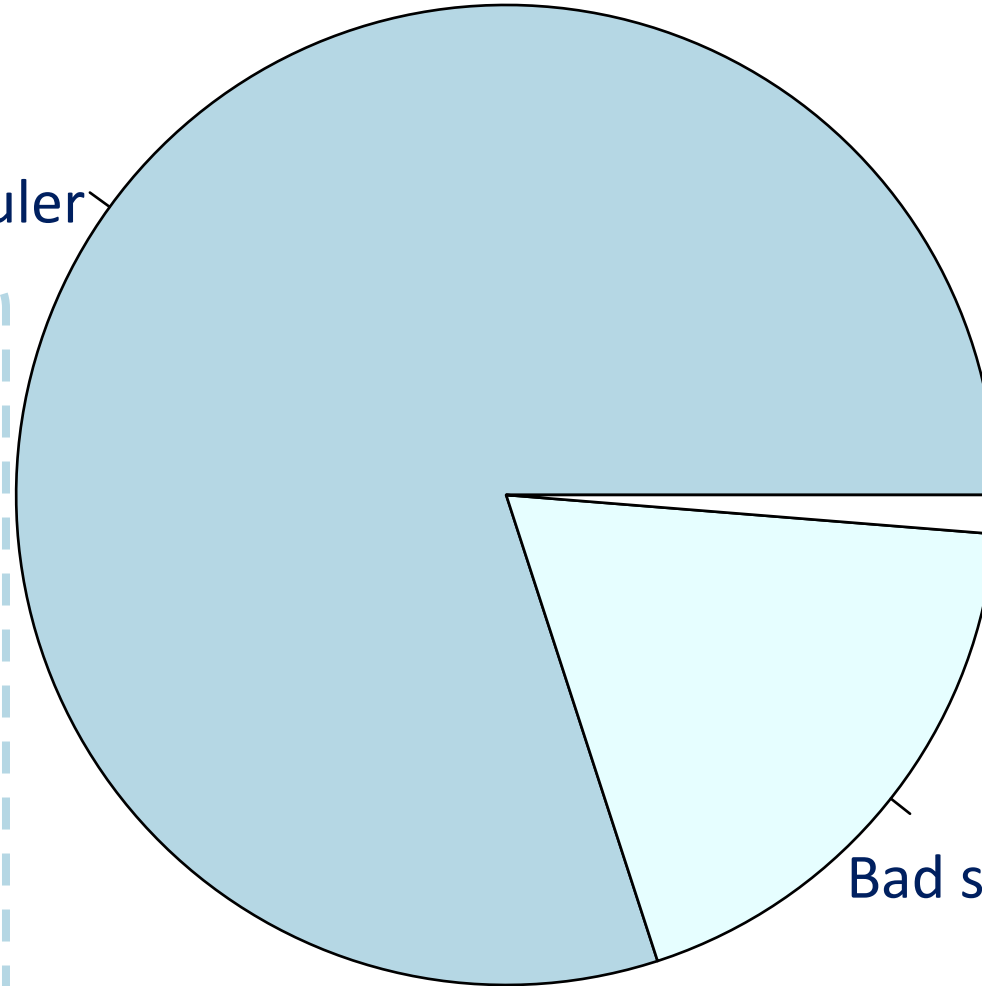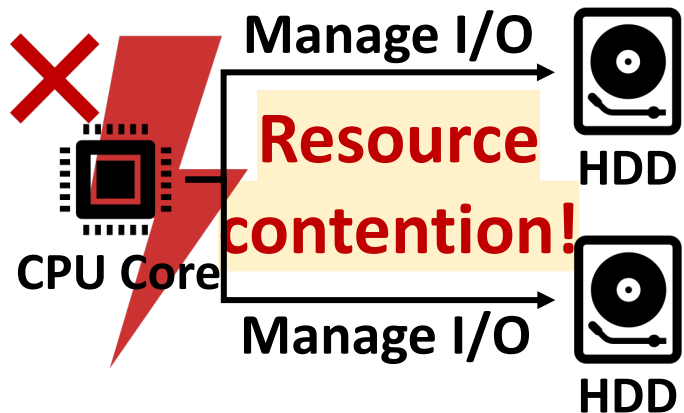
- **Root Cause Analysis (§6)**

**More details in the paper!**
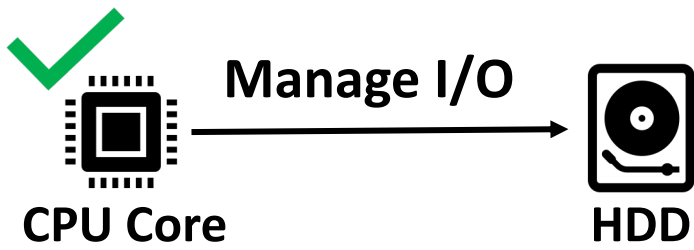
**Among 315 Confirmed Fail-Slow Drives:**

**Software (252):**
Ill-implemented scheduler



CPU Core → Manage I/O → HDD

Resource contention!

**Environment (4):**
Power & temperature

**Hardware (59):**
Bad sector, bad capacitors, etc.

# Perseus

---

**Detection Framework**

**Efficient**
**Fail-Slow**

🚩 **Non-intrusive**
(Performance) log-based
No source code altering

👆 **Accurate**
Recall/precision rate > 0.99

**Detection**

➤ **Fine-grained**
Device-level detection

✓ **General**
One set of parameters fits all scenarios

---

**Storage Devices**

· · · **Adaptable to Other Problem Domains**

33

# 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: lrm318@sjtu.edu.cn