# Dissecting I/O Burstiness in Machine Learning Cloud Platform: A Case Study on Alibaba's MLaaS

Qiang Zou, *Guangxi Minzu University, China*

Yuhui Deng, *Jinan University, China*
**Yifeng Zhu**, *University of Maine, USA*
Yi Zhou, *Columbus State University, USA*
Jianghe Cai, *Jinan University, China*
Shuibing He, *Zhejiang University, China*
yifeng.zhu@maine.edu
https://web.eece.maine.edu/~zhu/

# Outline

◆ Background

◆ Motivation

◆ Burstiness & Heavy-tailed Property

◆ Auto-correlation & Self-similarity

◆ Synthesis

◆ Conclusion

# Background

- ➢ Why <u>Alibaba's MLaaS</u> (Machine-Learning-as-a-Service)?
  - ✓ Alibaba Cloud launched **PAI** – the ML **P**latform for **A**rtificial **I**ntelligence
  - ✓ <u>Representative</u> - one of the leading MLaaS platforms in China
- ➢ For PAI:
  - ✓ Over 6500 GPUs across 1800 machines
  - ✓ See Fig. 1 for the architecture overview
- ➢ The <u>scalable storage solutions</u> rely on:
  - ✓ <u>Four key components</u> – Alibaba Cloud's object storage, distributed file system, database solutions, and elastic block storage
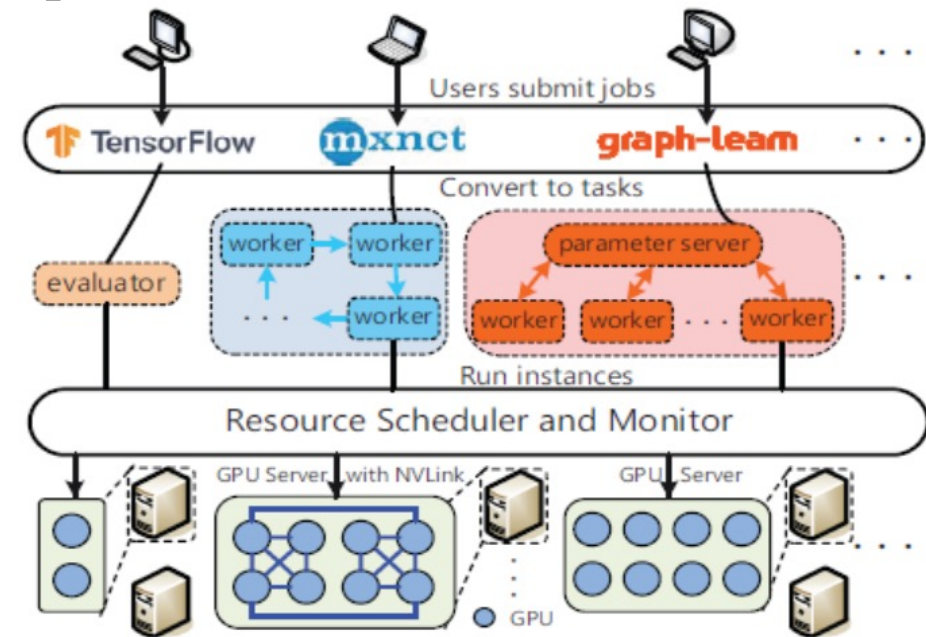


Fig. 1. PAI architecture overview.

# Background

➢ PAI traces collected on machines of GPU clusters

  ✓ The PAI traces at <u>the job, task, and instance levels</u> provide launch information including status, start_time, etc.

  ✓ The <u>machine-level</u> PAI trace contains information, such as timestamps, I/O waiting times (iowait), execution times in user and kernel modes, etc.

  ✓ See the referenced literature [1] for the details
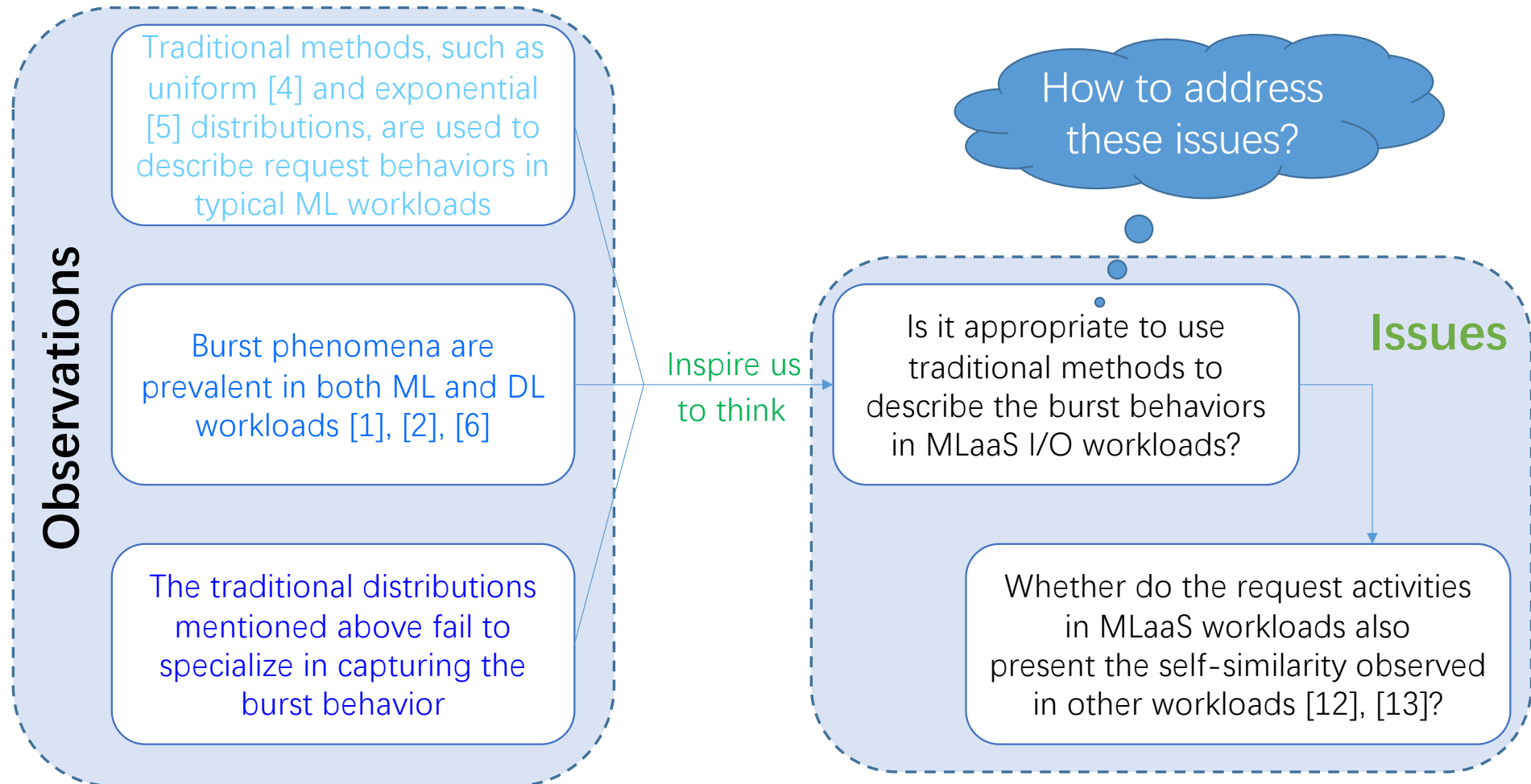
➢ <u>This study aims at:</u>

  ✓ The <u>machine-level</u> trace

  ✓ <u>Timestamp</u> information (in seconds for about two months)

TABLE I
SUMMARY OF PAI TRACE AND MACHINE SPECS OF GPU CLUSTERS [1].

| #Machines | 1800 | | Duration | 2 months | |
|---|---|---|---|---|---|
| Memory (GiB) | 512 | 512 | 512 | 384 | 512/384 |
| GPU type | P100 | T4 | Misc. | V100M32 | V100 |
| #GPUs | 2 | 2 | 8 | 8 | 8 |
| #Nodes | 798 | 497 | 280 | 135 | 104 |

# Motivation

**Observations**

Traditional methods, such as uniform [4] and exponential [5] distributions, are used to describe request behaviors in typical ML workloads

Burst phenomena are prevalent in both ML and DL workloads [1], [2], [6]

The traditional distributions mentioned above fail to specialize in capturing the burst behavior

Inspire us to think

How to address these issues?

**Issues**

Is it appropriate to use traditional methods to describe the burst behaviors in MLaaS I/O workloads?

Whether do the request activities in MLaaS workloads also present the self-similarity observed in other workloads [12], [13]?

# **Burstiness** & Heavy-



Fig. 2. Empirical CDF of request arrival intervals in the PAI workload.

➤ To show the burstiness quantitatively,

  ✓ By concept: ***non-stationary***, a large ***variance***;

  ✓ Approach – empirical study

➤ Non-stationary:

  ✓ 83% of I/O requests arrive within an interval of no more than 1 second

  ✓ up to 72% of requests arrive simultaneously at certain moments (in seconds)

➤ Variance:  as high as 8892

➤ To measure the strength of burstiness,

  ✓ Using the ***index of dispersion for intervals*** (IDI) [20]

  ✓ A larger value of the index of dispersion indicates stronger burstiness

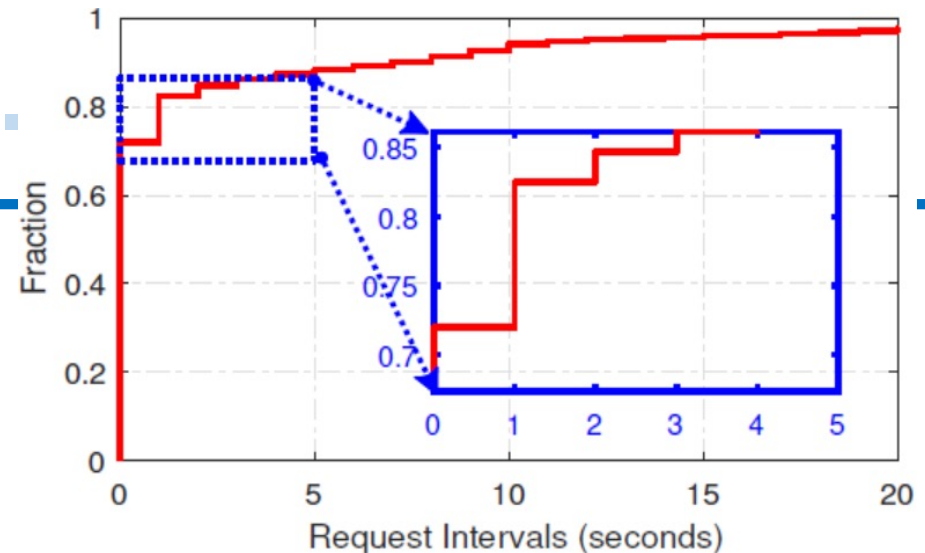  ✓ We calculate the IDI for I/O arrivals in the PAI workload as 1519 (significant bursty)

# Burstiness & Heavy-tailed Property

➢ Gaussianity Test:

   ✓ helps accurately describe <u>the tail trend</u> in the distribution of access characteristics

   ✓ can be conducted using a quantile-quantile (**QQ**) **plot**

➢ **For PAI** (see Figure 3):

   ✓ The corresponding scatter points clearly <u>do not fall on a straight line</u>

   ✓ Instead, the curve is concave upward, indicating *a heavy-tailed trend*

   ✓ Suggesting that the I/O behaviors in the PAI workload are <u>non-Gaussian</u>
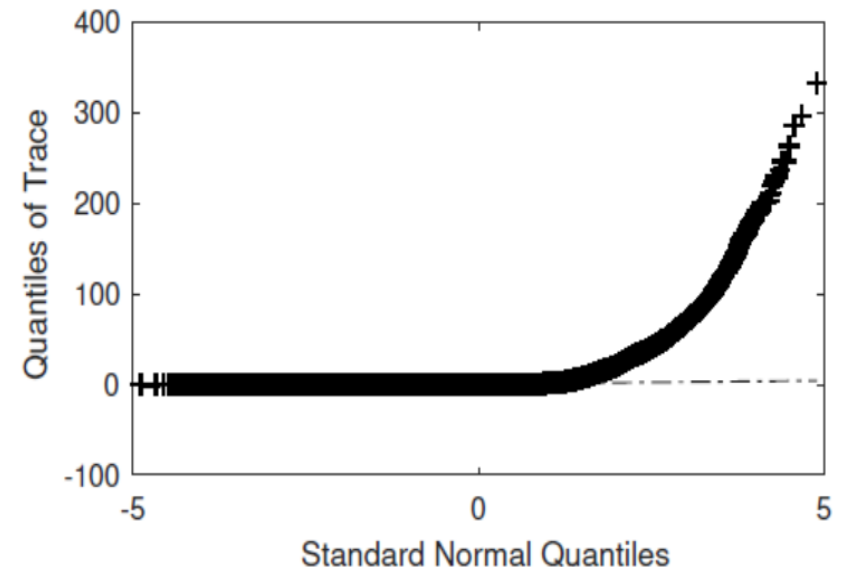


Fig. 3. Examine the Gaussianity of I/O request activities in the PAI workload through QQ plot of the PAI trace data versus standard normal, respectively.

# **Auto-correlation** **& Self-similarity**

➢ Tool: Auto-Correlation Function (ACF)

  ✓ For a time series $Y = \{Y_t: t = 1, \ 2, \ldots, n\}$, $\theta = E[Y_t]$, $y_t = Y_t - \theta$,

  ✓ Correlation coefficients: $R(k) = \dfrac{E[y_t \cdot y_{t+k}]}{E[y_t{}^2]}$, for $k \geq 0$

  ✓ A correlation coefficient forms <u>a mapping relationship</u> with a time interval (also called *lag*) $k$

➢ How is the auto-correlation curve *related to request activities*?

  ✓ **If** the correlation coefficients of arrival intervals decrease rapidly with the increase of *lag* and approach 0, there is <u>almost no correlation.</u>

  ✓ **Otherwise**, there is <u>a certain degree of correlation</u> for requests

# Auto-correlation & Self-similarity

➢ For I/O requests in the PAI workload, see Figure 4,

  ✓ As the lag increases from 0 to 100, the correlation coefficients of I/O requests do not approach zero sharply; instead, they exhibit a gradual declining trend

  ✓ There is *a noticeable degree of correlation* between request arrivals in the PAI workload

➢ Therefore, exploring self-similarity in the PAI workload becomes essential to accurately understand request behaviors
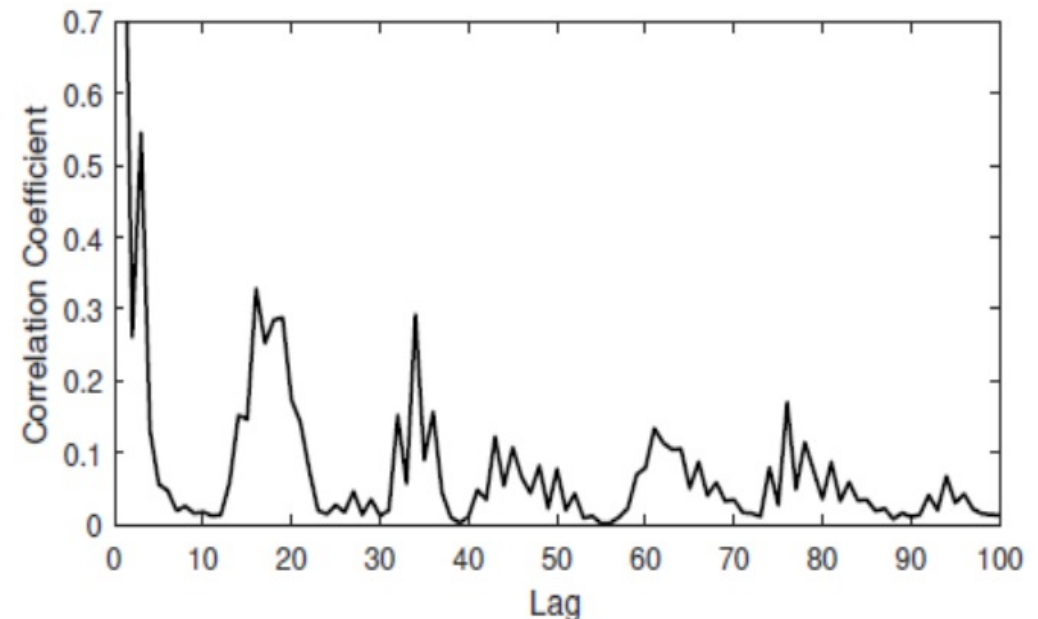


Fig. 4. Auto-correlation function of I/O request arrivals in the PAI workload.

9

# Auto-correlation & **Self-similarity**

➢ <u>What is</u> self-similarity?
  - ✓ In brief, the characteristics of a certain process are similar from different time scales

➢ <u>How to explore</u> the self-similarity in system workloads?
  - ✓ Showing the visualization
  - ✓ Providing theoretical evidence
  - ✓ Estimating the *Hurst* parameter ($0.5 < H < 1$)

➢ The well-known <u>tools</u> to estimate the *Hurst* Parameter:
  - ✓ Variance-time plot [12]
  - ✓ <u>R/S</u> (rescaled adjusted range) analysis (also called *Pox plot*) [26]

# Self-similarity (Visualization)

➢ Main trait: <u>the persistence of bursts and burst aggregations at various timescales</u>

➢ For PAI, see Figure 5:

  ✓ Three different timescales in subplots (a)-(c);

  ✓ each subsequent timescale being <u>ten times larger than</u> the previous one

  ✓ Each subplot is derived from a subinterval <u>randomly selected from</u> the time range depicted in the following subplot and it enhances the temporal resolution by a factor of 10

➢ **Finding**: <u>The time range</u> characterized by bursty requests <u>consists of nested subintervals, each is made of even smaller subintervals with similar burst behaviors.</u>
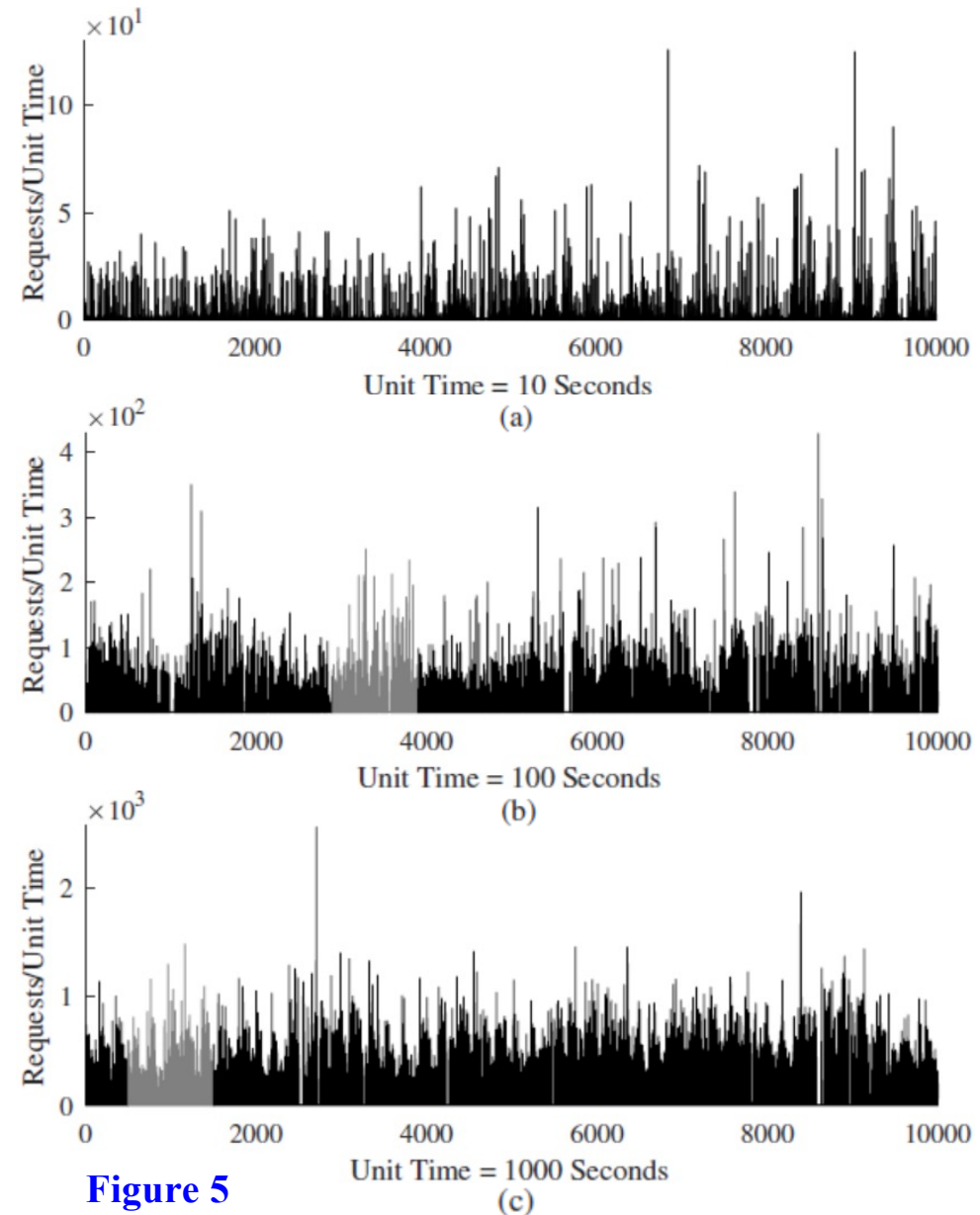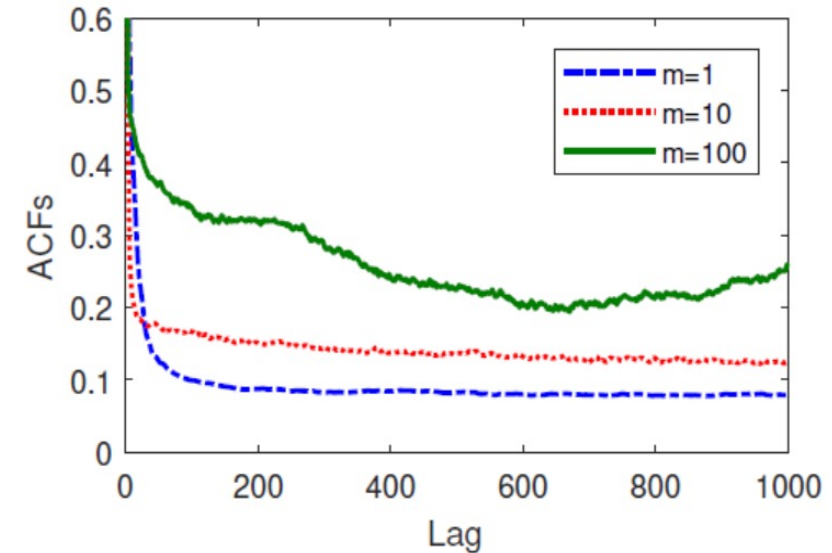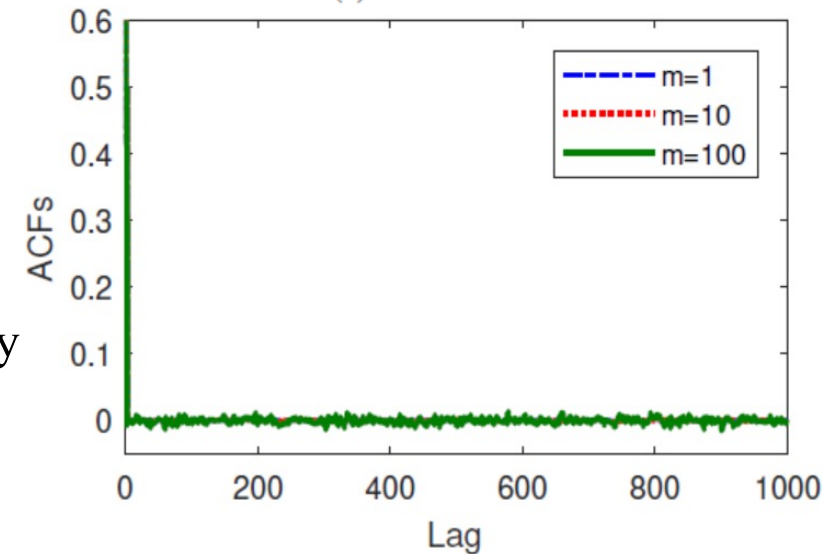


Figure 5

# Self-similarity **(Theoretical evidence)**


(a) PAI workload

➤ Theoretical basis: see the statements regarding the structure of $R^{(m)}(k)$ in Section IV-A

➤ Examining the auto-correlation functions of the agg-regated time series of the request sequence at multiple aggregation levels

➤ For PAI, see Figure 6:
  ✓ Plot (a) depicts the ACFs of the aggregated time series of the request sequence at multiple aggregation levels, that appear to converge to a similar function structure;
  ✓ Plot (b) demonstrates that the auto-correlation coefficients of Poisson workload at each aggregation level are generally very small and almost equal to zero

➤ Quite different from Poisson, request activities in the PAI workload behave like a self-similar process

**Figure 6**
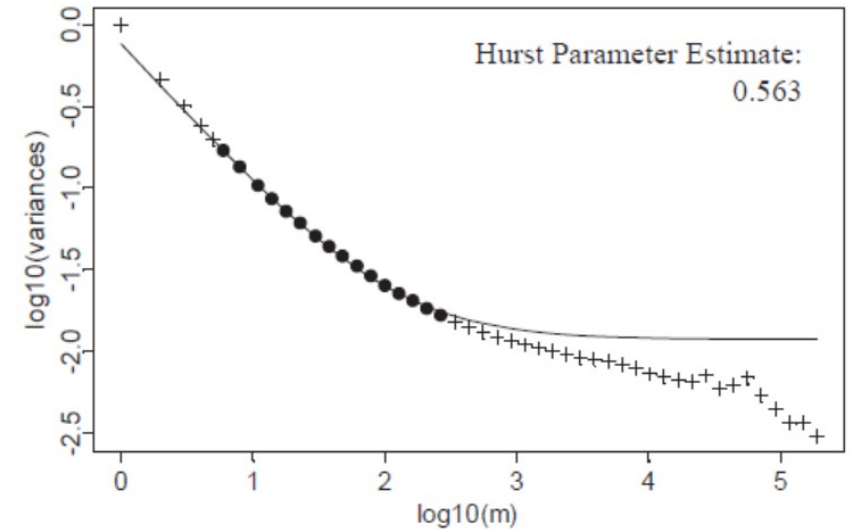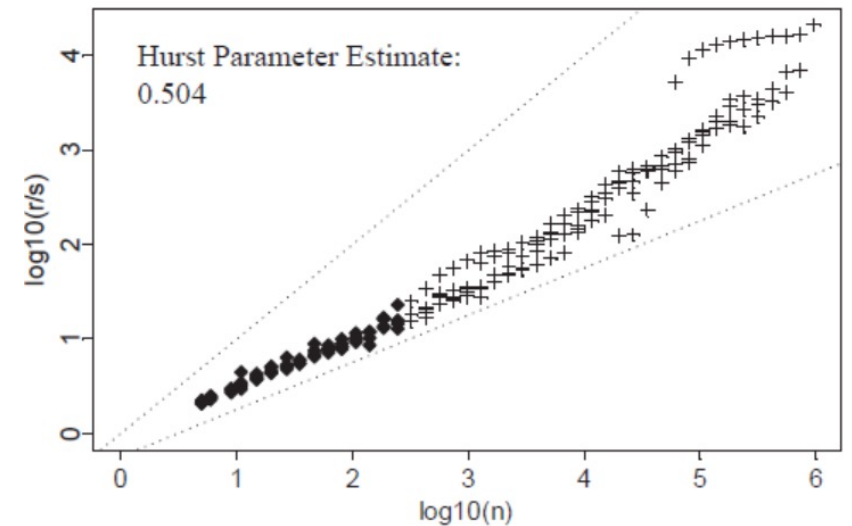

(b) Poisson workload

12

# Self-similarity (Hurst Parameter)

➢ For the I/O request sequence in the PAI workload, the Hurst parameter can be estimated by:

  ✓ The variance-time plot: see Fig. 7(a);

  ✓ The Pox plot: see Fig. 7(b)

➢ Finding:

  ✓ All Hurst parameter estimates are greater than 0.5;

  ✓ Quantitatively confirming the existence of self-similarity



Hurst Parameter Estimate: 0.563

(a) The variance-time plot



Hurst Parameter Estimate: 0.504

**Figure 7**

(b) The Pox plot

13

# Synthesis

➢ We have made the following findings:

  ✓ The arrival process of I/O requests is <u>highly bursty</u>
  ✓ Traditional methods struggle to accurately characterize the PAI work-
     load, as the I/O arrivals show <u>a certain degree of correlation</u>
  ✓ There seems to be <u>self-similarity</u> in the PAI workload
  ✓ The I/O request activities in PAI appear to be <u>non-Gaussian</u>

➢ These findings <u>inspire us</u> to use several methods to synthesize

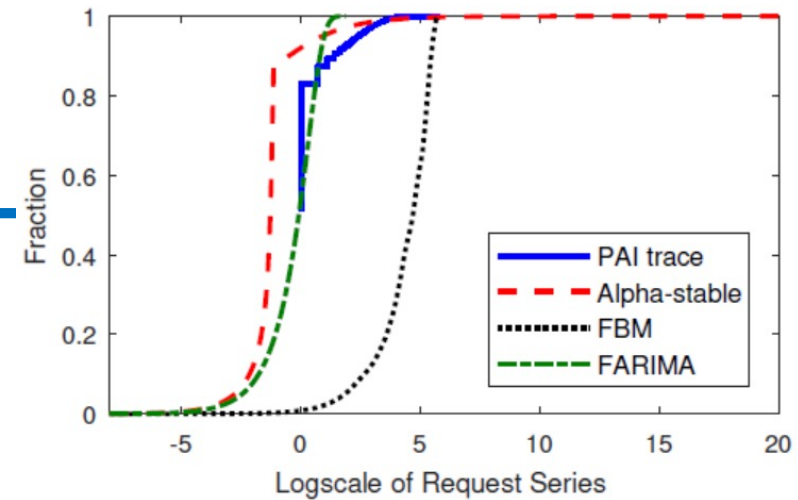  I/O request series for the **self-similar** PAI workload

# Synthesis

➢ Two typical <u>self-similar</u> workload models are chosen:

   ✓ Fractional Brownian motion (<u>FBM</u>) [27] – is adept at characterizing self-similarity under Gaussian conditions

   ✓ Fractional autoregressive integrated moving average (<u>FARIMA</u>) [28] – is well-known for its ability to describe both long-range and short-range dependences

➢ The versatile <u>alpha-stable</u> model [8] is also extended,

   ✓ **by** redefining its model parameters to synthesize request series for PAI

   ✓ **to** faithfully describe the bursts and heavy-tailed properties under non-Gaussian conditions

# **Synthesis**

➢ To evaluate the accuracy of these models, we adopt:

    ✓ The trimmed mean of errors [8];

    ✓ The cumulative distribution functions (CDFs)

➢ The <u>trimmed mean of errors</u> for FBM, FARIMA, and alpha-stable models:

    ✓ is 109.46, 1.26, and 0.78, respectively;

    ✓ the trimmed mean of the errors for the alpha-stable synthetic sequence (i.e., 0.78) is very <u>close to</u> that for the FARIMA synthetic one (i.e., 1.26)

➢ For the <u>CDFs</u> of the actual series and the synthetic ones, as shown in Figure 8:

    ✓ <u>Both</u> the FARIMA synthetic sequence and the alpha-stable synthetic one <u>exhibit</u> <u>convincing matching degrees</u>

    ✓ <u>One advantage of the latter over the former</u> is its ability to better capture the heavy-tailed feature

16

# Conclusion

➢ Characterizing the request behaviors in MLaaS workloads is crucial for scheduling and managing the I/O subsystem in GPU clusters.

➢ This paper studies the burstiness of the I/O requests in a representative and real-world MLaaS workload – the PAI workload, and shows the existence of self-similarity in the PAI  workload.

➢ Based on the inputs measured from real trace data, we deploy self-similar workload models to synthesize I/O request sequences for the PAI workload.

# Dissecting I/O Burstiness in Machine Learning Cloud Platform: A Case Study on Alibaba's MLaaS

Qiang Zou (qiangzou@hotmail.com), Yuhui Deng (tyhdeng@jnu.edu.cn),

**Yifeng Zhu (yifeng.zhu@maine.edu)**, Yi Zhou (zhou_yi@columbusstate.edu),

Jianghe Cai (761571151@qq.com), Shuibing He(heshuibing@zju.edu.cn)

# *Thank you!*

GUANGXI MINZU UNIVERSITY

1906 JINAN UNIVERSITY

1865

COLUMBUS STATE UNIVERSITY

ZHEJIANG UNIVERSITY