

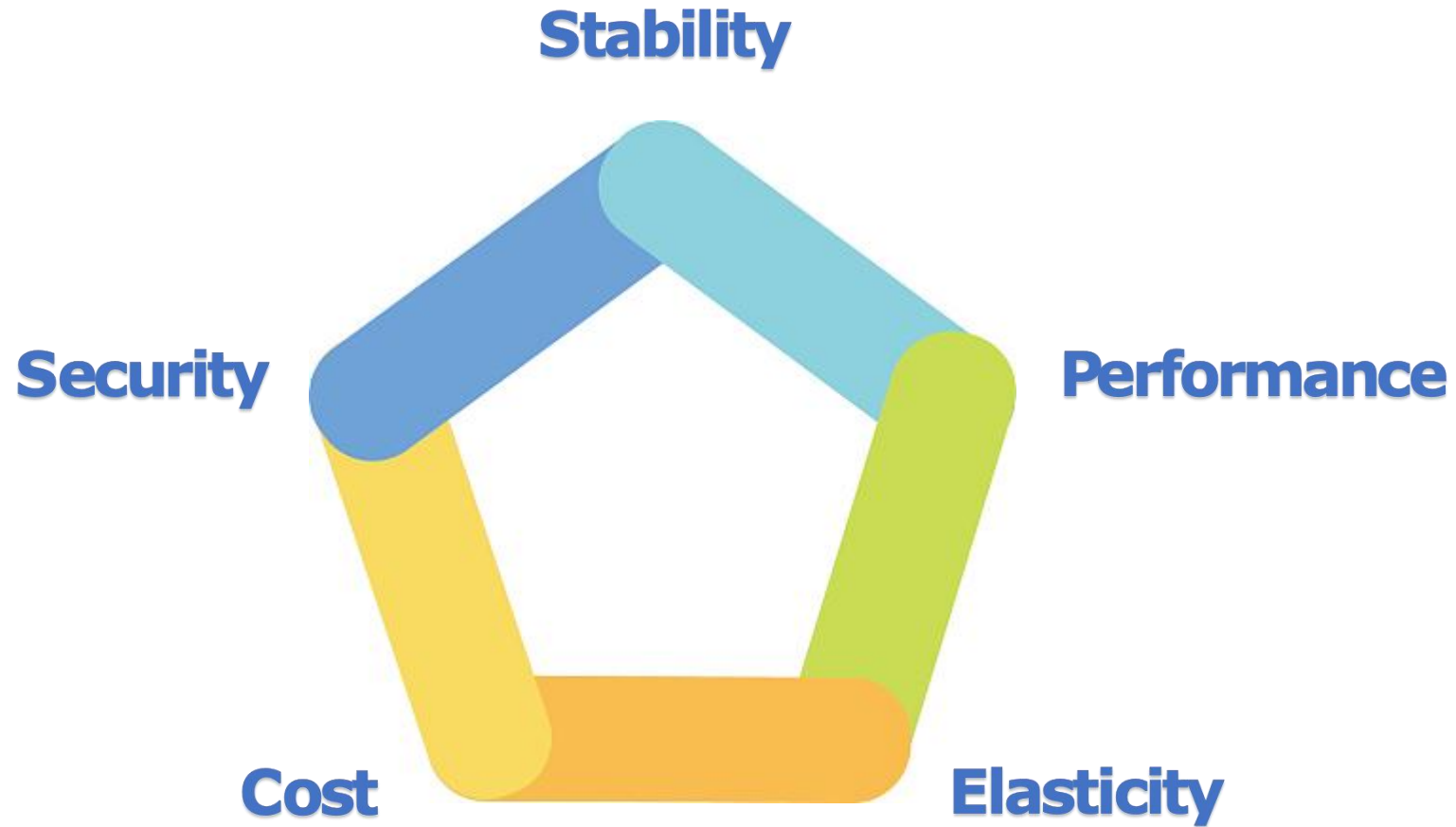
Stability is not Downtime: Comprehensive Stability Evaluation for Large-Scale Cloud Servers in Alibaba Cloud

Haoyu Wang, Zhicheng Liu, Yeliang Qiu, Haozhe Li, Hongke Guo,
Zhaoliang Zhu, You Zhang, Yu Zhou, Xudong Zheng

Alibaba Cloud Computing

ICDE 2025

Why Cloud Service?



Importance of Stability

- Serious faults in cloud computing may bring a very significant impact on society.



Importance of Stability

- A short period of unavailability may have a catastrophic impact on the businesses of customers.
- A customer experienced an unexpected memory hardware failure.

Instance type: ecs.g7.2xlarge

Resource: 8 cores + 32 GB Mem

Price: ~1000 CNY / month

VS

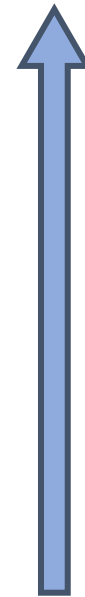
Service unavailability: ~10 min

Business loss: 70,000+ CNY

Complaint: 20+

Stability Challenge

- Elastic Compute Service (ECS) in Alibaba Cloud has an enormous scale, facing significant pressures in stability.
 - We built an AIOps system, CloudBot, to maintain stability.



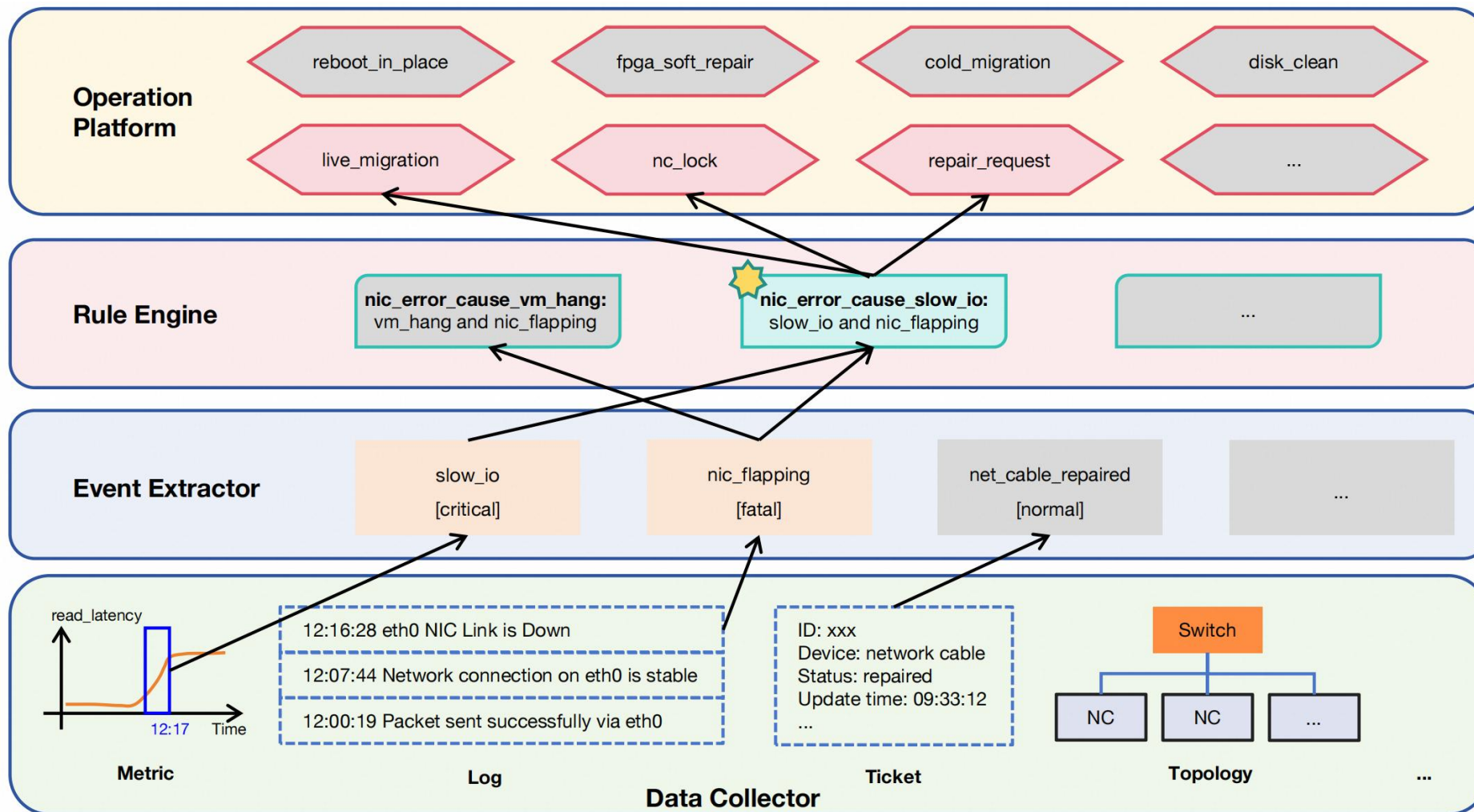
10,000,000+ VMs

1,000,000+ Servers

5,000+ Clusters

29 Regions

CloudBot



Quantitative Stability Evaluation

- A quantitative stability evaluation is required to drive the evolution of CloudBot.
- **Interpretability:** an interpretable evaluation is more convincing than a black-box one especially in cross-team collaborations.
- **Cost-Effectiveness:** heavyweight evaluation would result in significant extra resource usage due to the vast scale.
- **Non-Invasiveness:** it is not allowed to operate inside the VM.
- **Comprehensiveness:** The cloud server is a complex product including computing, storage and networking.

Related Work Comparison

- Benchmark-based evaluation: SuperBench...
- Downtime-based evaluation: Downtime Percentage, Annual Interruption Rate...

Dimension	Benchmark-based	Downtime-based
Interpretability	Yes	Yes
Cost-Effectiveness	No	Yes
Non-Invasiveness	No	Yes
Comprehensiveness	Yes	No

Stability \neq Downtime

- Complaints about VM performance
 - The performance is impacted by ECS underlying changes during its release.
- Incident on Nov. 12, 2023
 - VMs are not controllable due to faulty logic within the AccessKey system.

[Exception (recovered)] Abnormalities in Alibaba Cloud product console access and API calls

•

Scope of impact

- Some products such as Object Storage Service (OSS), Tablestore, Log Service, and Message Service (MNS) went through partial service disruption, while service provision of most products remained intact, such as Elastic Compute Service (ECS), ApsaraDB RDS, and all networking products.
- The consoles and management APIs of cloud products were affected.

Time

November 12, 2023 17:39 – 19:20 (UTC+8)

Definition & Category

■ Definition

- The stability of cloud servers is defined as its capacity to deliver and **manage** computational resources in a **continuous** and **consistent** manner.

■ Category

- Unavailability
- Performance
- Control-Plane

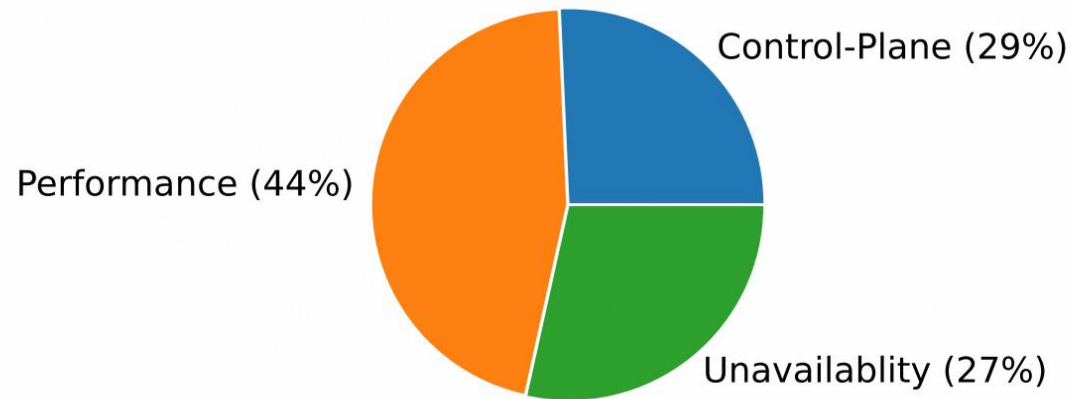
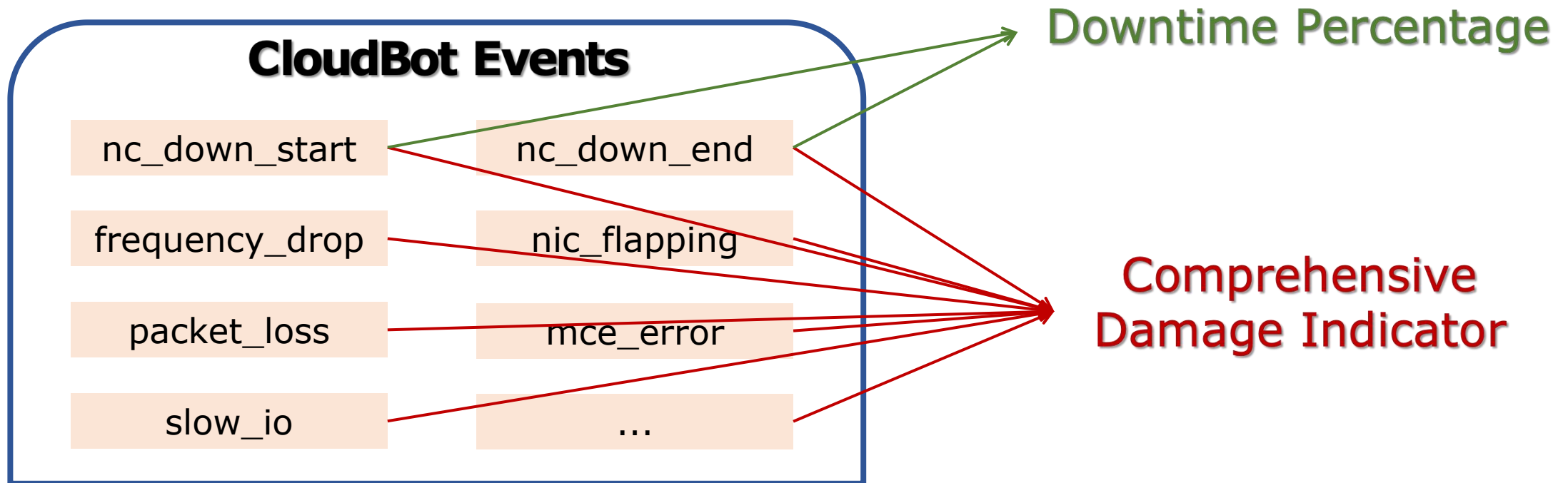


Fig. 2: Distribution of tickets related to ECS stability

From Downtime to CDI

- Comprehensive Damage Indicator (CDI)
 - Evaluate large-scale cloud server stability based on **interpretable intermediate results**, CloudBot events.



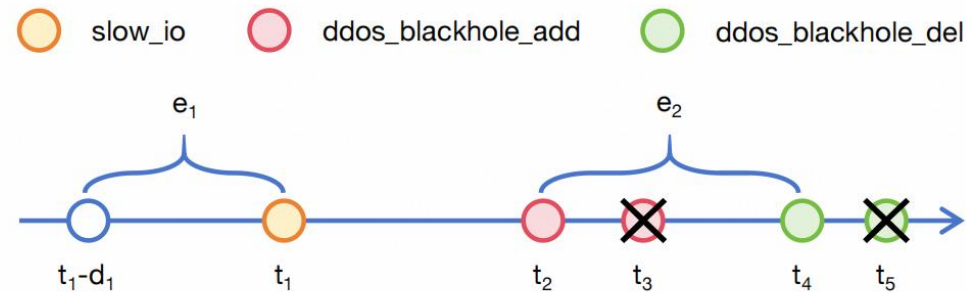
Comprehensive Damage Indicator

■ Sub-metrics

- Unavailability & Performance & Control-Plane Indicator
- Only difference lies in the specific events they rely on

■ Event Period

- Stateless events -> single
- Stateful events -> paired



■ Event Weight

- Expert weight: according to experience
- Customer weight: according to ticket number
- Analytic Hierarchy Process (AHP) to integrate the weights

Comprehensive Damage Indicator

■ Example

VM	Service Time	Event	Period	Weight	CDI
1	60min	packet_loss	10:08-10:10	0.3	0.020
		packet_loss	10:10-10:12	0.3	
2	1440min	vcpu_high	13:25-13:30	0.6	0.002
3	1000min	slow_io	08:08-08:10	0.5	0.004
		slow_io	08:10-08:12	0.5	
		vcpu_high	08:10-08:15	0.6	
All	2500min	-	-	-	0.003

■ Solution

$$\square Q_1 = (2*0.3+2*0.3) / 60 = 1.2 / 60 = 0.020$$

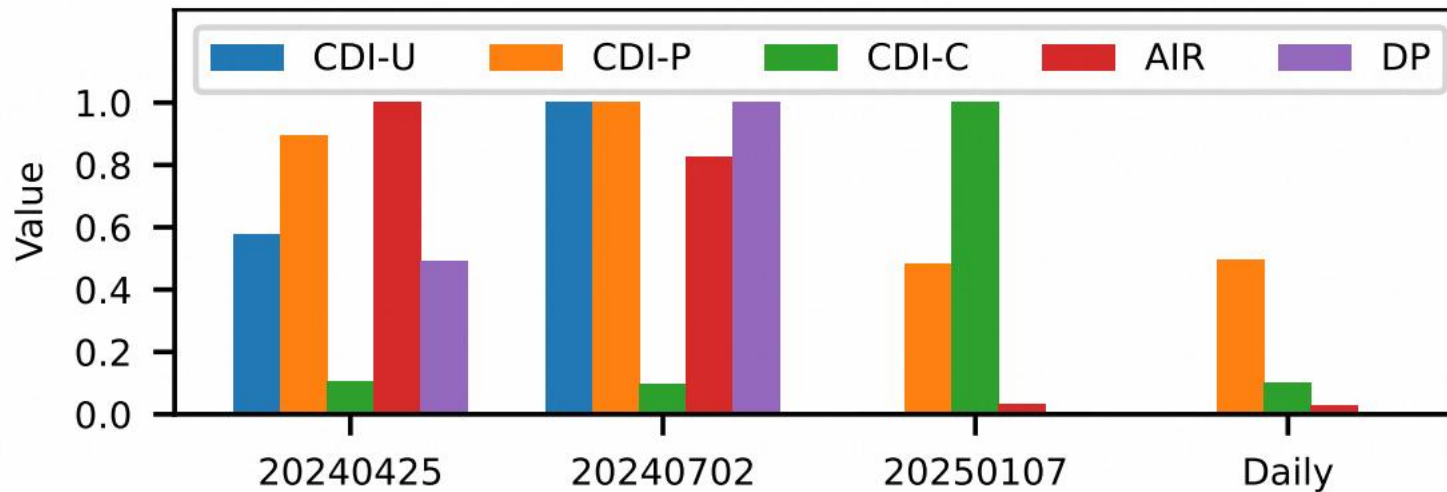
$$\square Q_2 = (5*0.6) / 1440 = 3 / 1440 = 0.002$$

$$\square Q_3 = (2*0.5+5*0.6) / 1000 = 4 / 1000 = 0.004$$

$$\square Q_{all} = (0.020*60+0.002*1440+0.004*1000) / (60+1440+1000) = 0.003$$

Application: Stability Evaluation

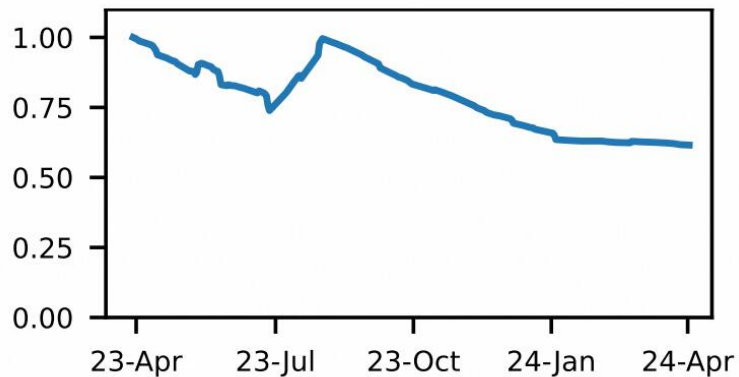
- Stability Evaluation on Incidents
 - 20240425: ECS service exception
 - 20240702: network access abnormalities
 - 20250107: abnormalities in purchase and modify



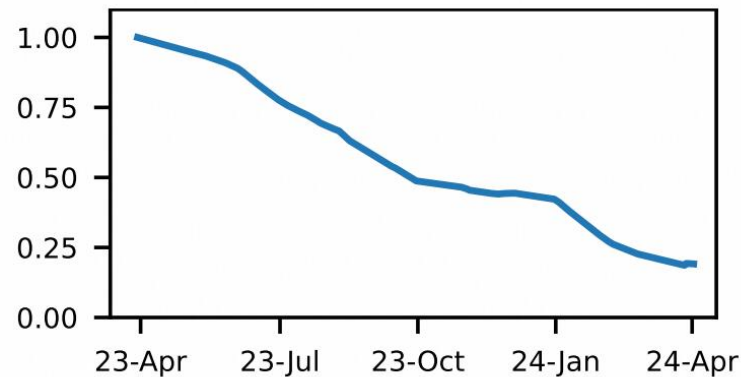
CDI-U: Unavailability Indicator
CDI-P: Performance Indicator
CDI-C: Control-Plane Indicator
DP: Downtime Percentage
AIR: Annual Interruption Rate

Application: Stability Evaluation

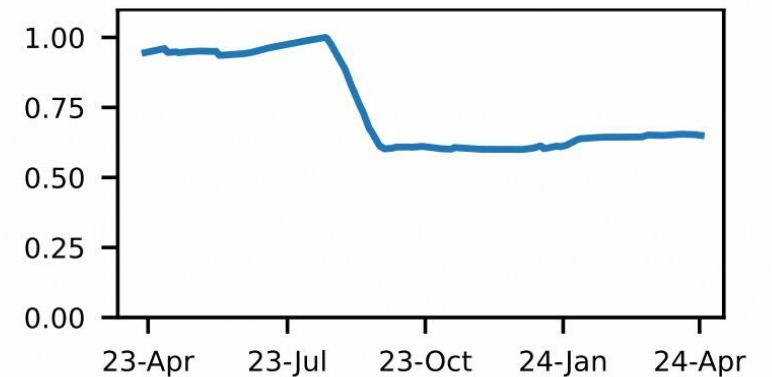
- Daily stability evaluation for FY24 (April 2023 to March 2024)
 - ▣ Unavailability Indicator: **40%** decrease
 - ▣ Performance Indicator: **80%** decrease
 - ▣ Control-plane Indicator: **35%** decrease



(a) Unavailability Indicator



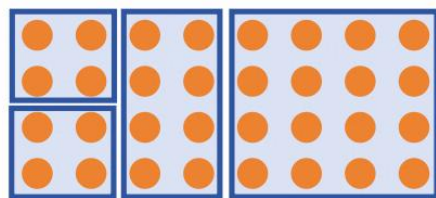
(b) Performance Indicator



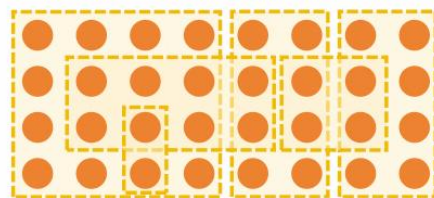
(c) Control-plane Indicator

Application: Architecture Comparison

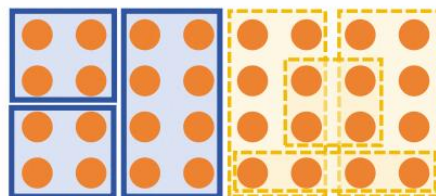
- Homogeneous-deployment to hybrid-deployment
- Monitor CDI to prevent stability loss from architecture evolution.



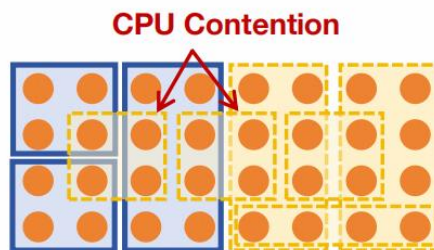
(a) Homogeneous dedicated VM



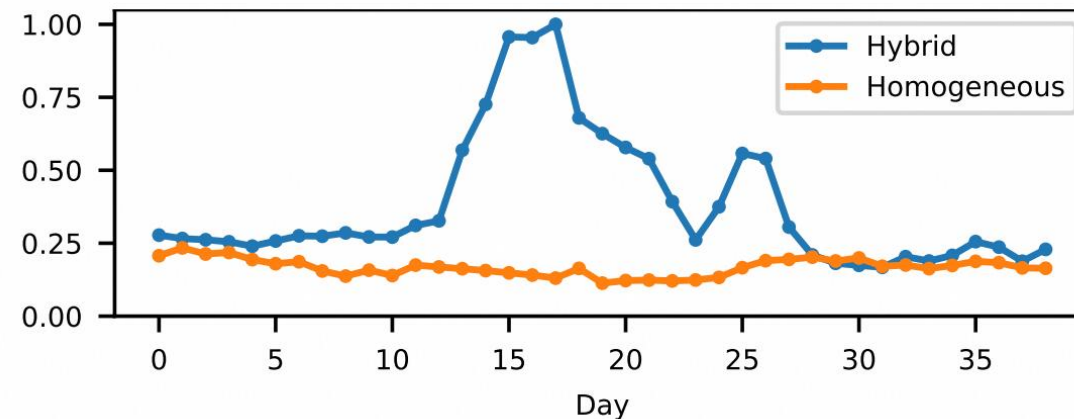
(b) Homogeneous shared VM



(c) Hybrid VM

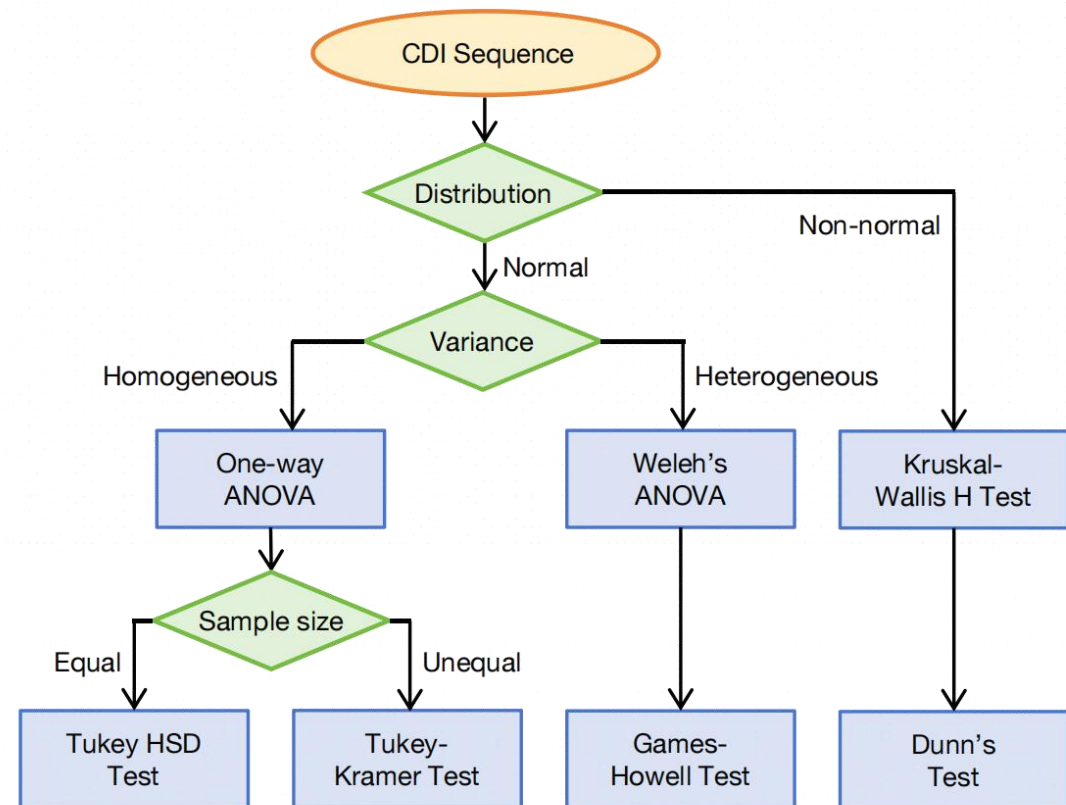
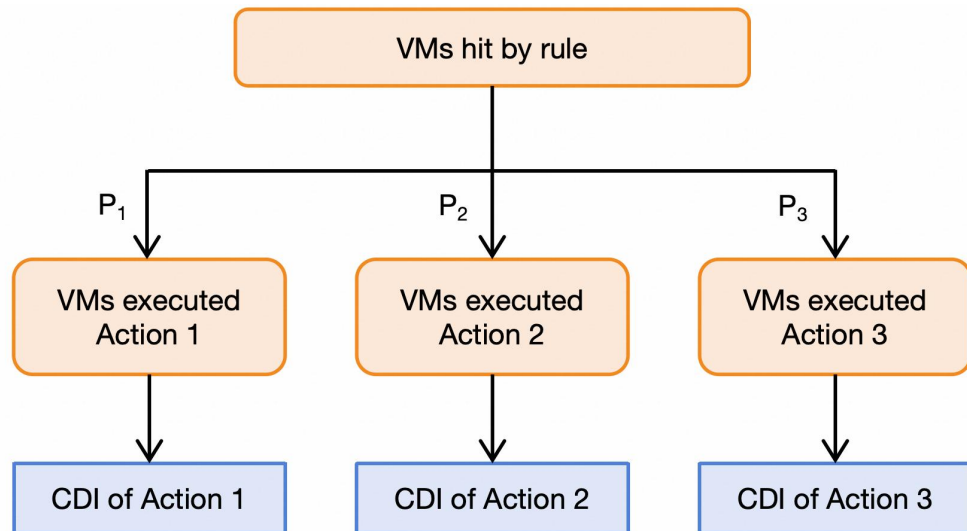


(d) Hybrid VM with issue



Application: Operation Action Optimization

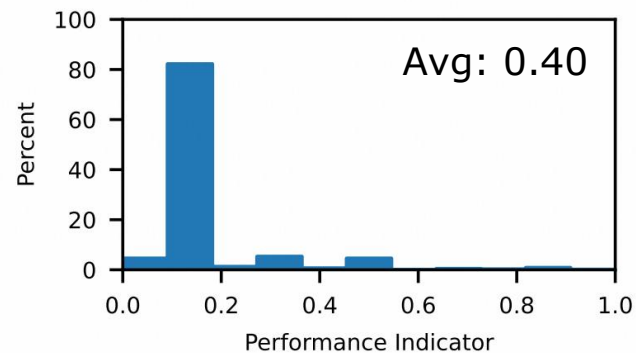
- Select optimal action from several candidates
 - A/B Test: randomly execute action and calculate its CDI for the next two days
 - Hypothesis testing: omnibus test + post-hoc analysis



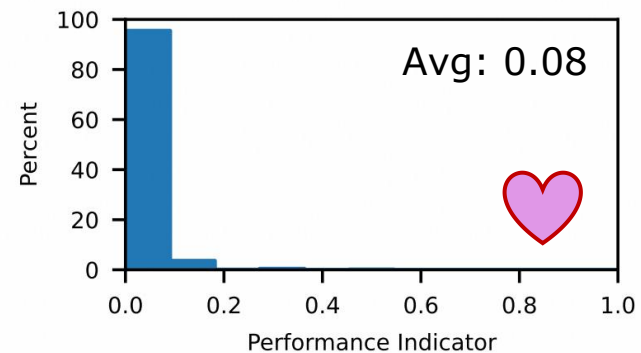
Application: Operation Action Optimization

- Select optimal action from several candidates
 - A/B Test: randomly execute action and calculate its CDI for the next two days
 - Hypothesis testing: omnibus test + post-hoc analysis

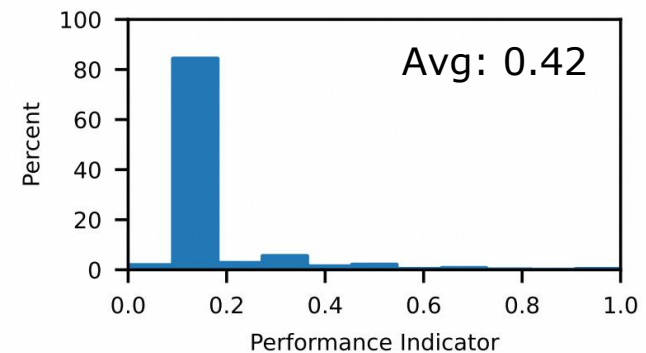
Sub-metric	Omnibus Test		Post-hoc Analysis		
	P-value	Sign.	Pair	P-value	Sign.
Unavailability	0.47	False			
Control-plane	0.89	False			
Performance	0	True	A-B	0	True
			A-C	0.03	True
			B-C	0	True



(a) Action A



(b) Action B



(c) Action C

Conclusion

- In this paper, we propose Comprehensive Damage Indicator (CDI).
- The **first** comprehensive quantitative stability evaluation metric for large-scale cloud servers, to the best of our knowledge.
- Extending stability evaluation beyond system unavailability to include **performance** and **control-plane** issues.
- Deployed in various domains for over two years within Alibaba Cloud ECS, providing support for the evolution of stability.

THANKS FOR YOUR LISTENING

Slides are shown on the personal homepage <https://wanghy.pages.dev>