

#### Actionable and Interpretable Fault Localization for Recurring Failures in Online Service Systems

**Zeyan Li**, Nengwen Zhao, Mingjie Li, Xianglin Lu, Lixin Wang, Dongdong Chang, Xiaohui Nie, Li Cao, Wenzhi Zhang, Kaixin Sui, Yanhua Wang, Xu Du, Guoqiang Duan, Dan Pei

ESEC/FSE 2022

# Table of Contents

- Background
- Design
- Evaluation
- Limitations and Future Works

# Background



# Background

Span3 User Span2 Span1 Online service systems comprise Spanrequests many different components (e.g., services, containers and servers). Service B Service C Services ServiceD Service E Service A Container Container Container Containers The large scale and complex Process Process Process dependencies among components make system failures inevitable and hard to diagnose. Infrastructure/ Server Storage Network

# **Monitoring Metrics**



Various metrics are closely monitored on a 24×7 basis. They serve as the most direct signals to the underlying failures





**Recurring Failures** 

How to determine recurring failures? Recurring failures should be discovered, mitigated and fixed in similar ways in engineers' perspective.

# **Recurring Failures**

We analyzed 576 failure tickets, spanning one year, of a production banking information system.

Typical recurring failures:

- Failed disks
- Unavailable third-party services
- Missing database indexes
- Slow SQLs



Typical non-recurring failures:

- Design flaw
- Data inconsistency
- Code defects

Recurring failures are largely prevalent in practice.

# Core Ideas





# **Core Ideas**





Various metrics are monitored on different components.

Failure propagates from faulty components due to the complex dependencies among components.

Generalizing to previously unseen failures (no failures of the same kinds have occurred at the same locations).

It is hard for engineers to trust the localization results from black-box models without interpretation.

# DéjàVu Overview



The name DéjàVu comes from a French phrase (Déjà vu) which means "already seen".



#### The vertices of an FDG is all the candidate failure units of an online service system

 Service
 English

 DB
 CPU Exhaustion : CPU utilization, I/O wait...

 Docker
 Packet loss: Recv/Sent packets, Recv/Sent bytes.....

 Server
 Out of memry: free bytes, cached bytes, total bytes.....

Engineers are supposed to summarize the actionable fault types of historical failures and the corresponding input metrics

We can generate all the candidate failure units by combining all components and the corresponding fault types.



We connect the candidate failure units on an FDG by the call or deployment relationships of the components

FDGs are automatically constructed for different failures since FDGs are dynamic due to frequent changes

#### Model Architecture



# **Interpretation Methods**

Local interpretation: finding similar historical failures based on the aggregated features





Figure 15: A decision tree trained for the failure class OS Net work of  $\mathcal{A}$ . Note that the metrics values are normalized.

Global interpretation: using simpler but interpretable models to approximate (rather than outperform) it as accurately as possible.

### **Evaluation: Datasets**

Dataset	#Failures	<b>#Metrics</b>	<b>#Failure Units</b>	<b>#Failure Classes</b>	System	Failures
$\mathcal{A}$	188	710	102	18	A production system	Injected
${\mathcal B}$	158	2419	189	20	A production system	Injected
${\mathcal C}$	99	2594	41	41	A production system	Real-world failures
${\cal D}$	156	5724	1044	23	An open-source benchmark	Injected

Our experiments are conducted on four systems and 601 failures, including three production systems and 99 real-world failures.

The typical failure injection types include CPU/memory stress, database session limit, and network packet delay/loss.

# **Evaluation: Overall Performance**

Category	Approach	Mean Average Rank (MAR)	Top-3 Accuracy	p-value (MAR, t-test)
Supervised	DéjàVu	2.82	87.45%	
Similar failure	JSS'20	103.56	30.63%	1.1e-16
matching	iSUQAD	80.29	34.34%	1.8e-14
Traditional	Decision Tree	49.25	63.90%	7.3e-17
supervised	Gradient Boosting	12.83	71.59%	1.5e-10
	Random Forest	6.20	88.78%	2.1e-08
	SVM	33.34	26.48%	3.1e-11
Unsupervised	RandomWalk@Metric	36.19	19.13%	7.4e-17
heuristic	RandomWalk@FI	135.44	20.49%	8.7e-18

The localization performance of DéjàVu significantly outperforms the baselines.

# Limitations and Future Work



Non-recurring failures: the localization models should be able to identify potential non-recurring failures and notify engineers.



Frequent changes: the localization models should be updated after every change without relying on new failures.

# Thank you

https://github.com/NetManAIOps/DejaVu