# **Practical Root Cause Localization** for Microservice Systems via Trace Analysis

**Zeyan Li**<sup>1</sup>, Junjie Chen<sup>2</sup>, Rui Jiao<sup>1</sup>, Nengwen Zhao<sup>1</sup>, Zhijun Wang<sup>3</sup>, Shuwei Zhang<sup>3</sup>, Yanjun Wu<sup>3</sup>, Long Jiang<sup>3</sup>, Leiqin Yan<sup>3</sup>, Zikai Wang<sup>4</sup>, Zhekang Chen<sup>4</sup>, Wenchi Zhang<sup>4</sup>, Xiaohui Nie<sup>1</sup>, Kaixin Sui<sup>4</sup>, Dan Pei<sup>1</sup>





June 26 2021. IWQoS 2021





## Outline



Background

Approach





Summary

Microservice architecture: the lasted trend to build systems

- Microservice architecture: the lasted trend to build systems
- What is microservice architecture?
  - Loosely coupled and independently deployed 'micro' services
  - Communicating with other services to realize a user request



- Microservice architecture: the lasted trend to build systems
- What is microservice architecture?
  - Loosely coupled and independently deployed 'micro' services
  - Communicating with other services to realize a user request
- Benefits:
  - Fast delivery
  - Better scalability
  - Greater autonomy











#### Many related services also behave abnormally due to the dependency among services





#### Many related services also behave abnormally due to the dependency among services

#### Different services are developed and maintained by different teams







Many related services also behave abnormally due to the dependency among services

Different services are developed and maintained by different teams

When a fault happens, we need to localize the rootcause service at first before the corresponding team's further investigation







## Background **Challenges in Microservice Systems**

1. Complex dependencies among services





2. Dynamic runtime environments

#### 2. Dynamic runtime environments

| COMPANY               | DEPLOY<br>FREQUENCY    | DEPLOY<br>LEAD TIME | RELIABILITY | CUSTOMER<br>FEEDBACK |  |
|-----------------------|------------------------|---------------------|-------------|----------------------|--|
| AMAZON                | 23.000/day             | minutes             | high        | high                 |  |
| GOOGLE                | 5.500/day              | minutes             | high        | high                 |  |
| NETFLIX               | 500/day                | minutes             | high        | high                 |  |
| FACEBOOK              | 1/day                  | hours               | high        | high                 |  |
| TWITTER               | 3/week                 | hours               | high        | high                 |  |
| TYPICAL<br>ENTERPRISE | once every<br>9 months | months or more      | low/medium  | low/medium           |  |
|                       |                        |                     |             |                      |  |

From "The Phoenix Project"

#### 2. Dynamic runtime environments

| COMPANY               | DEPLOY<br>FREQUENCY    | DEPLOY<br>LEAD TIME | RELIABILITY | CUSTOMER<br>FEEDBACK |  |
|-----------------------|------------------------|---------------------|-------------|----------------------|--|
| AMAZON                | 23.000/day             | minutes             | high        | high                 |  |
| GOOGLE                | 5.500/day              | minutes             | high        | high                 |  |
| NETFLIX               | 500/day                | minutes             | high        | high                 |  |
| FACEBOOK              | 1/day                  | hours               | high        | high                 |  |
| TWITTER               | 3/week                 | hours               | high        | high                 |  |
| TYPICAL<br>ENTERPRISE | once every<br>9 months | months or more      | low/medium  | low/medium           |  |
|                       |                        |                     |             |                      |  |

From "The Phoenix Project"



3. Various types and huge volumes of monitoring data



### Background **Metrics or Logs Based Localization Approaches**



Metrics lack detailed contextual information and fine-grained dependencies of invocations

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-15 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition logs-0 as the leader reported an error: NOT\_LEADER\_ FOR\_PARTITION [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Error sending fetch request (sessionId=839052068, epoch=517118) to node 2: java.nio.channels.C losedSelectorException. [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-47 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-11 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-41 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-5 as the leader reported an erro r: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-35 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-17 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Error sending fetch request (sessionId=1303574239, epoch=127483) to node 0: java.nio.channels. 01----

Logs are usually generated in an arbitrary manner and most log messages contain little information for root cause localization



### Background **Metrics or Logs Based Localization Approaches**



Metrics lack detailed contextual information and fine-grained dependencies of invocations

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-15 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition logs-0 as the leader reported an error: NOT\_LEADER\_ FOR\_PARTITION [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Error sending fetch request (sessionId=839052068, epoch=517118) to node 2: java.nio.channels.C losedSelectorException. [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-47 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-11 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-41 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-5 as the leader reported an erro r: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-35 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-17 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Error sending fetch request (sessionId=1303574239, epoch=127483) to node 0: java.nio.channels. 01----

Logs are usually generated in an arbitrary manner and most log messages contain little information for root cause localization



### Background **Metrics or Logs Based Localization Approaches**



Metrics lack detailed contextual information and fine-grained dependencies of invocations

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-15 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition logs-0 as the leader reported an error: NOT\_LEADER\_ FOR\_PARTITION [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Opening socket connection to server kafka-zookeeper/10.47.244.48:2181. Will not attempt to aut henticate using SASL (unknown error) [kafka.log][INFO] Error sending fetch request (sessionId=839052068, epoch=517118) to node 2: java.nio.channels.C losedSelectorException. [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-47 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-11 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION [kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-41 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-5 as the leader reported an erro r: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-35 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Retrying leaderEpoch request for partition \_\_consumer\_offsets-17 as the leader reported an err or: NOT\_LEADER\_FOR\_PARTITION

[kafka.log][INFO] Error sending fetch request (sessionId=1303574239, epoch=127483) to node 0: java.nio.channels. 01----

Logs are usually generated in an arbitrary manner and most log messages contain little information for root cause localization

















Service Dependency Graph

**contextual information** for root-cause service localization

abnormal invocation of an abnormal trace as the root cause.





abnormal invocation of an abnormal trace as the root cause.





abnormal invocation of an abnormal trace as the root cause.



abnormal invocation of an abnormal trace as the root cause.



abnormal invocation of an abnormal trace as the root cause.





 MEPFL[FSE'19] uses a supervised machine learning model to "classify" the root-cause service for an abnormal trace.



- Frontend
- 1 Product
- 0 Cart

. . .

O Ad-Service
### Motivation **Existing Trace-Based Localization Approaches**

MEPFL[FSE'19] uses a supervised machine learning model to "classify" the root-cause service for an abnormal trace.





- Frontend
- Product
- Cart 0
- Ad-Service

Supervised approaches rely on training data with high coverage on all fault types and microservices



### Motivation **Existing Trace-Based Localization Approaches**

MEPFL[FSE'19] uses a supervised machine learning model to "classify" the root-cause service for an abnormal trace.





- Frontend
- Product
- Cart
- Ad-Service

Supervised approaches rely on training data with high coverage on all fault types and microservices

> It is hard for supervised approaches to handle new faults





The existing trace-based approaches utilize only *single* traces and *abnormal* traces.

The existing trace-based approaches utilize only single traces and abnormal traces.

less normal traces is more likely to be the root cause.

**Insight**: Intuitively, a microservice covered by more abnormal traces but



The existing trace-based approaches utilize only single traces and abnormal traces.

less normal traces is more likely to be the root cause.

Our insight utilizes the statistical information based on both normal and abnormal traces

**Insight**: Intuitively, a microservice covered by more abnormal traces but



is more likely to be the root cause.

### A microservice covered by more abnormal traces but less normal traces

13



is more likely to be the root cause.



is more likely to be the root cause.



is more likely to be the root cause.

#abnormal traces passing through s support(s) =#abnormal traces  $confidence(s) = \frac{\#abnormal traces passing through s}{}$ #traces passing through *s* 



is more likely to be the root cause.

#abnormal traces passing through s support(s) =support  $\uparrow \Rightarrow$  "more abnormal traces" #abnormal traces  $confidence(s) = \frac{\#abnormal traces passing through s}{}$ #traces passing through *s* 



is more likely to be the root cause.

#abnormal traces passing throug support(s) =#abnormal traces  $confidence(s) = \frac{\#abnormal\ traces\ passing\ throws the set of the set of$ #traces passing through s

| gh <u>s</u> | support $\uparrow \Rightarrow$ "more abnormal traces"  |
|-------------|--|
| ough s      | confidence $\uparrow \Rightarrow$ "less normal traces" |





| gh <u>s</u> | support $\uparrow \Rightarrow$ "more abnormal traces"  |
|-------------|--|
| ough s      | confidence $\uparrow \Rightarrow$ "less normal traces" |





| gh <u>s</u> | support $\uparrow \Rightarrow$ "more abnormal traces"  |
|-------------|--|
| ough s      | confidence $\uparrow \Rightarrow$ "less normal traces" |





| ugh s           | support ↑  | $\Rightarrow$ | "more abnormal traces"             |
|-----------------|------------|---------------|------------------------------------|
| nrough <u>s</u> | confidence | 1             | $\Rightarrow$ "less normal traces" |

|                   | support | confidence |
|-------------------|---------|------------|
| Frontend          |         |            |
| Product           |         |            |
| Recommendation    |         |            |
| <b>Ad-Service</b> |         |            |





| ugh s           | support ↑  | $\Rightarrow$ | "more abnormal traces"             |
|-----------------|------------|---------------|------------------------------------|
| nrough <u>s</u> | confidence | 1             | $\Rightarrow$ "less normal traces" |

|                   | support | confidence |
|-------------------|---------|------------|
| Frontend          |         |            |
| Product           |         |            |
| Recommendation    |         |            |
| <b>Ad-Service</b> |         |            |





| ugh <i>s</i> | support $\uparrow \Rightarrow$ "more abnormal traces  | 5"   |
|--------------|---|------|
| h s          | confidence $\uparrow \Rightarrow$ "less normal traces | , 77 |

|                | support | confidence |
|----------------|---------|------------|
| Frontend       | 3/3     |            |
| Product        |         |            |
| Recommendation |         |            |
| Ad-Service     |         |            |





| ugh <i>s</i> | support $\uparrow \Rightarrow$ "more abnormal traces  | 5"   |
|--------------|---|------|
| h s          | confidence $\uparrow \Rightarrow$ "less normal traces | , 77 |

|                | support | confidence |
|----------------|---------|------------|
| Frontend       | 3/3     |            |
| Product        |         |            |
| Recommendation |         |            |
| Ad-Service     |         |            |





| ugh <u>s</u> | support ↑  | $\Rightarrow$ | "more abnormal traces"             |
|--------------|------------|---------------|------------------------------------|
| h s          | confidence | 1             | $\Rightarrow$ "less normal traces" |

|                | support | confidence |
|----------------|---------|------------|
| Frontend       | 3/3     | 3/4        |
| Product        |         |            |
| Recommendation |         |            |
| Ad-Service     |         |            |





| ugh <u>s</u> | support $\uparrow \Rightarrow$ "more abnormal traces"  |
|--------------|--|
| rough s      | confidence $\uparrow \Rightarrow$ "less normal traces" |

|                | support | confidence |
|----------------|---------|------------|
| Frontend       | 3/3     | 3/4        |
| Product        | 3/3     | 3/3        |
| Recommendation | 3/3     | 3/4        |
| Ad-Service     | 1/3     | 1/2        |



We need to detect abnormal traces



We need to detect abnormal traces



We need to detect abnormal traces

 $S_2$  $S_1$ 





We need to detect abnormal traces

 $S_1$  $S_2$ 

resource exhaustion



We need to detect abnormal traces







We need to detect abnormal traces



response time ↑ accept all new connections



We need to detect abnormal traces



response time ↑ accept all new connections

response rate  $\downarrow$ 

refuse new connections after the connection pool is exhausted



We need to detect abnormal traces



response time ↑ accept all new connections

response rate  $\downarrow$ 

refuse new connections after the connection pool is exhausted



We need to detect abnormal traces



Only some features are related to the concerned fault, but there could be also anomalies in other irrelevant features





We need to detect abnormal traces



Only some features are related to the concerned fault, but there could be also anomalies in other irrelevant features







We need to detect abnormal traces



Only some features are related to the concerned fault, but there could be also anomalies in other irrelevant features









The number of anomalies with respect to useful features and irrelevant features are different





The number of anomalies with respect to useful features and irrelevant features are different

Highlight 1: We select useful features by testing whether the distribution of abnormal and normal invocations with respect to a feature changes after the fault happens


# **Core Idea 1. Multi-Feature Anomaly Detection**



| CPU usage | Before | After |
|-----------|--------|-------|
| Normal    | 99%    | 90%   |
| Abnormal  | 1%     | 10%   |



The number of anomalies with respect to useful features and irrelevant features are different

Highlight 1: We select useful features by testing whether the distribution of abnormal and normal invocations with respect to a feature changes after the fault happens



# **Core Idea 1. Multi-Feature Anomaly Detection**



### Highlight 1: We select useful features by testing whether the distribution of abnormal and normal invocations with respect to a feature changes after the fault happens

| CPU usage | Before | After |
|-----------|--------|-------|
| Normal    | 99%    | 90%   |
| Abnormal  | 1%     | 10%   |



#### The number of anomalies with respect to useful features and irrelevant features are different

| Throughput | Before | After |
|------------|--------|-------|
| Normal     | 98%    | 98%   |
| Abnormal   | 2%     | 2%    |





### Highlight 2: (Extending Our Insight) We mine microservice sets with high supports and confidences rather than single microservices



### Highlight 2: (Extending Our Insight) We mine microservice sets with high supports and confidences rather than single microservices





### Highlight 2: (Extending Our Insight) We mine microservice sets with high supports and confidences rather than single microservices





### Highlight 2: (Extending Our Insight) We mine **microservice sets** with high supports and confidences rather than **single microservices**



|   | support | confidence |
|---|---------|------------|
| Α | 4/4     | 4/6        |
| B | 4/4     | 4/6        |
| С | 4/4     | 4/6        |
| D | 2/4     | 2/6        |



### Highlight 2: (Extending Our Insight) We mine **microservice sets** with high supports and confidences rather than **single microservices**



| Only the | ose trad |
|----------|----------|
|----------|----------|

|   | support | confidence |   |
|---|---------|------------|---|
| A | 4/4     | 4/6        | × |
| В | 4/4     | 4/6        |   |
| С | 4/4     | 4/6        |   |
| D | 2/4     | 2/6        |   |



### Highlight 2: (Extending Our Insight) We mine microservice sets with high supports and confidences rather than single microservices



|        | support | confidence |
|--------|---------|------------|
| {A, B} | 4/4     | 4/6        |
| {A, C} | 2/4     | 2/4        |
| {B, C} | 2/4     | 2/2        |
| {B, D} | 2/4     | 2/2        |



### Highlight 2: (Extending Our Insight) We mine microservice sets with high supports and confidences rather than single microservices



|        | support | confidence |  |
|--------|---------|------------|--|
| {A, B} | 4/4     | 4/6        |  |
| {A, C} | 2/4     | 2/4        |  |
| {B, C} | 2/4     | 2/2        |  |
| {B, D} | 2/4     | 2/2        |  |





Our insight and the two key metrics are validated on the two systems





Our insight and the two key metrics are validated on the two systems



Target Output: A ordered list of microservices

Target Output: A ordered list of microservices





Target Output: A ordered list of microservices







Target Output: A ordered list of microservices



### We have mined **{A, B}**, but which one should be further investigated by operators first?

#### Existing approaches presume the most upstream service is the root cause (page. 10)

С

Target Output: A ordered list of microservices



### We have mined **{A, B}**, but which one should be further investigated by operators first?

#### Existing approaches presume the most upstream service is the root cause (page. 10)

С

Target Output: A ordered list of microservices



Existing approaches presume the most

### We have mined **{A, B}**, but which one should be further investigated by operators first?

#### Existing approaches presume the most upstream service is the root cause (page. 10)

С













Highlight 2: For any service on an abnormal trace, if it has both in-coming and outcoming abnormal invocations, it is probably just propagating the anomaly.



by abs(#incoming abnormal invocations - #outcoming abnormal invocations)



















see details in our paper






























| Dataset | Microservice<br>Benchmark | Fault Type                         | Root-Cause<br>Component Level | # Faults |
|---------|---------------------------|------------------------------------|-------------------------------|----------|
| A       |                           | <b>Application Bug</b>             | Microservice                  | 58       |
|         | Train-Ticket              | CPU exhausted                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Container                     | 10       |
|         |                           | Network Delay                      | API                           | 14       |
| B       | Production<br>System      | 5 types provided by the developers |                               | 22       |



| ype      | Root-Cause<br>Component Level | # Faults |
|----------|-------------------------------|----------|
| on Bug   | Microservice                  | 58       |
| austed   | Microservice                  | 59       |
| Delay    | Microservice                  | 59       |
| Delay    | Container                     | 10       |
| Delay    | API                           | 14       |
| vided by |                               | つつ       |
| opers    |                               |          |

| D | ataset | Microservice<br>Benchmark | Fault Type                         | Root-Cause<br>Component Level | # Faults |
|---|--------|---------------------------|------------------------------------|-------------------------------|----------|
|   | A      |                           | Application Bug                    | Microservice                  | 58       |
|   |        | Troip Tiolrot             | CPU exhausted                      | Microservice                  | 59       |
|   |        | TCE'12                    | Network Delay                      | Microservice                  | 59       |
|   |        |                           | Network Delay                      | Container                     | 10       |
|   |        |                           | Network Delay                      | API                           | 14       |
|   | B      | Production<br>System      | 5 types provided by the developers |                               | 22       |

| Dataset | Microservice<br>Benchmark | Fault Type                            | Root-Cause<br>Component Level | # Faults |
|---------|---------------------------|---------------------------------------|-------------------------------|----------|
|         |                           | <b>Application Bug</b>                | Microservice                  | 58       |
| A       | Train Tialcat             | CPU exhausted                         | Microservice                  | 59       |
|         | TCE'121                   | Network Delay                         | Microservice                  | 59       |
|         |                           | Network Delay                         | Container                     | 10       |
|         |                           | Network Delay                         | API                           | 14       |
| B       | Production<br>System      | 5 types provided by<br>the developers | У                             | 22       |

| Dataset | Microservice<br>Benchmark | Fault Type                         | Root-Cause<br>Component Level | # Faults |
|---------|---------------------------|------------------------------------|-------------------------------|----------|
|         |                           | Application Bug                    | Microservice                  | 58       |
| A       | Troip Tiolrot             | CPU exhausted                      | Microservice                  | 59       |
|         | TCE'12                    | Network Delay                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Container                     | 10       |
|         |                           | Network Delay                      | API                           | 14       |
| B       | Production<br>System      | 5 types provided by the developers |                               | 22       |

| Dataset | Microservice<br>Benchmark | Fault Type                         | Root-Cause<br>Component Level | # Faults |
|---------|---------------------------|------------------------------------|-------------------------------|----------|
| A       |                           | <b>Application Bug</b>             | Microservice                  | 58       |
|         | Train-Ticket<br>[TSE'18]  | CPU exhausted                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Container                     | 10       |
|         |                           | Network Delay                      | API                           | 14       |
| B       | Production<br>System      | 5 types provided by the developers |                               | 22       |
|         |                           |                                    |                               |          |

| Dataset | Microservice<br>Benchmark | Fault Type                         | Root-Cause<br>Component Level | # Faults |
|---------|---------------------------|------------------------------------|-------------------------------|----------|
|         |                           | Application Bug                    | Microservice                  | 58       |
| A       | Train-Ticket              | CPU exhausted                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Microservice                  | 59       |
|         |                           | Network Delay                      | Container                     | 10       |
|         |                           | Network Delay                      | API                           | 14       |
| B       | Production<br>System      | 5 types provided by the developers | _                             | 22       |

#### **Evaluation** Overall Performance

#### **Evaluation Overall Performance**

TABLE III: Comparison of root cause localization on faults of different levels on

|              |              | -    |             |      |               |      |             |      |               |      |          |
|--------------|--------------|------|-------------|------|---------------|------|-------------|------|---------------|------|----------|
| Subject      | Algorithm    | A@1  | <b>↑A@1</b> | A@2  | <b>↑A@2</b>   | A@3  | <b>↑A@3</b> | MAR  | ↑МАR          | MFR  | ↑MFR     |
|              | TraceRCA     | 0.83 |             | 0.93 |               | 0.97 |             | 1.39 |               | 1.34 |          |
|              | MicroScope   | 0.56 | 46.67%      | 0.62 | <b>49.49%</b> | 0.70 | 37.33%      | 3.64 | 61.77%        | 3.47 | 61.26%   |
| Mionocomico  | MEPFL (RF)   | 0.94 | -12.00%     | 0.97 | -4.82%        | 0.97 | -0.64%      | 1.42 | 1.98%         | 1.38 | 2.71%    |
| MICIOSEIVICE | Random Walk  | 0.51 | 61.96%      | 0.86 | 7.25%         | 0.94 | 3.00%       | 1.97 | <b>29.37%</b> | 1.91 | 29.51%   |
|              | RCSF         | 0.52 | 60.00%      | 0.86 | 7.64%         | 0.93 | 3.69%       | 1.68 | <b>16.98%</b> | 1.60 | 16.02%   |
|              | TraceAnomaly | 0.49 | 70.85%      | 0.59 | 58.14%        | 0.63 | 53.11%      | 4.42 | <b>68.54%</b> | 4.34 | 69.13%   |
|              | TraceRCA     | 0.80 |             | 0.80 |               | 0.80 |             | 3.80 | 0%            | 3.80 |          |
|              | MicroScope   | 0.20 | 300.00%     | 0.40 | 100.00%       | 0.40 | 100.00%     | 7.20 | 47.22%        | 7.20 | 47.22%   |
| Containan    | MEPFL (RF)   | 0.80 | 0.00%       | 0.80 | 0.00%         | 1.00 | -20.00%     | 1.40 | -171.43%      | 1.40 | -171.43% |
| Container    | Random Walk  | 0.40 | 100.00%     | 0.60 | 33.33%        | 0.60 | 33.33%      | 8.40 | 54.76%        | 8.40 | 54.76%   |
|              | RCSF         | 0.40 | 100.00%     | 0.60 | 33.33%        | 0.60 | 33.33%      | 3.60 | -5.56%        | 3.60 | -5.56%   |
|              | TraceAnomaly | 0.20 | 300.00%     | 0.30 | 166.67%       | 0.30 | 166.67%     | 7.10 | 46.48%        | 7.10 | 46.48%   |
|              | TraceRCA     | 0.83 |             | 0.83 |               | 0.83 |             | 1.75 |               | 1.75 |          |
|              | MicroScope   | 0.58 | 42.86%      | 0.67 | 64.29%        | 0.92 | -9.09%      | 2.00 | 12.50%        | 2.00 | 12.50%   |
| A DI         | MEPFL (RF)   | 0.83 | 0.00%       | 1.00 | -16.67%       | 1.00 | -16.67%     | 1.17 | -50.00%       | 1.17 | -50.00%  |
| AFI          | Random Walk  | 0.58 | 42.86%      | 0.75 | 11.11%        | 0.64 | 29.63%      | 2.33 | 25.00%        | 2.33 | 25.00%   |
|              | RCSF         | 0.42 | 100.00%     | 0.58 | 42.86%        | 0.67 | 25.00%      | 2.58 | 32.26%        | 2.58 | 32.26%   |
|              | TraceAnomaly | 0.21 | 287.33%     | 0.36 | 132.40%       | 0.50 | 66.00%      | 5.86 | 70.12%        | 5.86 | 70.12%   |

TABLE IV: Comparison of root cause localization on multi-root-cause faults of  $\mathcal{A}$ 

| Subject    | Algorithm    | A@1  | <b>↑A@1</b> | A@2  | <b>↑A@2</b> | A@3  | <b>↑A@3</b>   | MAR  | ↑MAR          | MFR  | ↑MFR   |
|------------|--------------|------|-------------|------|-------------|------|---------------|------|---------------|------|--------|
|            | TraceRCA     | 0.45 |             | 0.82 |             | 0.95 |               | 1.77 |               | 1.09 |        |
| multi-     | MicroScope   | 0.27 | 66.67%      | 0.27 | 200.00%     | 0.41 | 133.33%       | 5.18 | <b>65.79%</b> | 2.73 | 60.00% |
|            | MEPFL (RF)   | 0.45 | 0.00%       | 0.95 | -14.29%     | 0.95 | 0.00%         | 1.64 | -8.33%        | 1.09 | 0.00%  |
| root-cause | Random Walk  | 0.41 | 11.11%      | 0.64 | 28.57%      | 0.82 | <b>16.67%</b> | 2.27 | 22.00%        | 1.36 | 20.00% |
| cases on A | RCSF         | 0.23 | 100.00%     | 0.50 | 63.64%      | 0.73 | 31.25%        | 2.82 | 37.10%        | 1.73 | 36.84% |
|            | TraceAnomaly | 0.50 | -10.00%     | 0.73 | 12.75%      | 0.82 | 16.11%        | 2.50 | 29.20%        | 1.00 | -9.00% |

| Λ                  |  |
|--------------------|--|
| $\mathbf{\Lambda}$ |  |

#### **Evaluation Overall Performance**

TABLE III: Comparison of root cause localization on faults of different levels on

| Subject      | Algorithm    | A@1  | <b>↑A@1</b> | A@2  | <b>↑</b> A@2  | A@3  | <b>↑A@3</b> | MAR  | ↑MAR          | MFR  | ↑MFR          |  |
|--------------|--------------|------|-------------|------|---------------|------|-------------|------|---------------|------|---------------|--|
|              | TraceRCA     | 0.83 |             | 0.93 |               | 0.97 |             | 1.39 |               | 1.34 |               |  |
|              | MicroScope   | 0.56 | 46.67%      | 0.62 | <b>49.49%</b> | 0.70 | 37.33%      | 3.64 | 61.77%        | 3.47 | 61.26%        |  |
| Microservice | MEPFL (RF)   | 0.94 | -12.00%     | 0.97 | -4.82%        | 0.97 | -0.64%      | 1.42 | 1.98%         | 1.38 | 2.71%         |  |
|              | Random Walk  | 0.51 | 61.96%      | 0.86 | 7.25%         | 0.94 | 3.00%       | 1.97 | <b>29.37%</b> | 1.91 | <b>29.51%</b> |  |
|              | RCSF         | 0.52 | 60.00%      | 0.86 | 7.64%         | 0.93 | 3.69%       | 1.68 | <b>16.98%</b> | 1.60 | 16.02%        |  |
|              | TraceAnomaly | 0.49 | 70.85%      | 0.59 | 58.14%        | 0.63 | 53.11%      | 4.42 | 68.54%        | 4.34 | 69.13%        |  |
|              | TraceRCA     | 0.80 |             | 0.80 |               | 0.80 |             | 3.80 | 0%            | 3.80 |               |  |
|              | MicroScope   | 0.20 | 300.00%     | 0.40 | 100.00%       | 0.40 | 100.00%     | 7.20 | 47.22%        | 7.20 | 47.22%        |  |
| Container    | MEPFL (RF)   | 0.80 | 0.00%       | 0.80 | 0.00%         | 1.00 | -20.00%     | 1.40 | -171.43%      | 1.40 | -171.43%      |  |
| Container    | Random Walk  | 0.40 | 100.00%     | 0.60 | 33.33%        | 0.60 | 33.33%      | 8.40 | 54.76%        | 8.40 | 54.76%        | TraceRCA outperforms unsupervised base   |
|              | RCSF         | 0.40 | 100.00%     | 0.60 | 33.33%        | 0.60 | 33.33%      | 3.60 | -5.56%        | 3.60 | -5.56%        | nucerier europenenne uneuperviced suce   |
|              | TraceAnomaly | 0.20 | 300.00%     | 0.30 | 166.67%       | 0.30 | 166.67%     | 7.10 | 46.48%        | 7.10 | 46.48%        | and parform as well as auparvised approx |
|              | TraceRCA     | 0.83 |             | 0.83 |               | 0.83 |             | 1.75 |               | 1.75 |               | and perform as well as supervised approa |
|              | MicroScope   | 0.58 | 42.86%      | 0.67 | 64.29%        | 0.92 | -9.09%      | 2.00 | 12.50%        | 2.00 | 12.50%        |  |
| Δ            | MEPFL (RF)   | 0.83 | 0.00%       | 1.00 | -16.67%       | 1.00 | -16.67%     | 1.17 | -50.00%       | 1.17 | -50.00%       | under different situations               |
|              | Random Walk  | 0.58 | 42.86%      | 0.75 | 11.11%        | 0.64 | 29.63%      | 2.33 | 25.00%        | 2.33 | 25.00%        |  |
|              | RCSF         | 0.42 | 100.00%     | 0.58 | 42.86%        | 0.67 | 25.00%      | 2.58 | 32.26%        | 2.58 | 32.26%        |  |
|              | TraceAnomaly | 0.21 | 287.33%     | 0.36 | 132.40%       | 0.50 | 66.00%      | 5.86 | 70.12%        | 5.86 | 70.12%        |  |

TABLE IV: Comparison of root cause localization on multi-root-cause faults of  $\mathcal{A}$ 

| Subject  | Algorithm    | A@1  | <b>↑A@1</b> | A@2  | <b>↑A@2</b> | A@3  | <b>↑A@3</b>   | MAR  | <b>↑MAR</b> | MFR  | ↑MFR   |
|--|--------------|------|-------------|------|-------------|------|---------------|------|-------------|------|--------|
|  | TraceRCA     | 0.45 |             | 0.82 |             | 0.95 |               | 1.77 |             | 1.09 |        |
| multi-<br>root-cause<br>cases on $\mathcal{A}$ | MicroScope   | 0.27 | 66.67%      | 0.27 | 200.00%     | 0.41 | 133.33%       | 5.18 | 65.79%      | 2.73 | 60.00% |
|  | MEPFL (RF)   | 0.45 | 0.00%       | 0.95 | -14.29%     | 0.95 | 0.00%         | 1.64 | -8.33%      | 1.09 | 0.00%  |
|  | Random Walk  | 0.41 | 11.11%      | 0.64 | 28.57%      | 0.82 | <b>16.67%</b> | 2.27 | 22.00%      | 1.36 | 20.00% |
|  | RCSF         | 0.23 | 100.00%     | 0.50 | 63.64%      | 0.73 | 31.25%        | 2.82 | 37.10%      | 1.73 | 36.84% |
|  | TraceAnomaly | 0.50 | -10.00%     | 0.73 | 12.75%      | 0.82 | 16.11%        | 2.50 | 29.20%      | 1.00 | -9.00% |
|  |              |      |             |      |             |      |               |      |             |      |        |

| Λ                          |  |
|----------------------------|--|
| $\boldsymbol{\mathcal{A}}$ |  |





### TraceRCA handles 10,000 traces per second per core.

24









Fig. 14: Efficiency improvement and performance degradation of *TraceRCA* with trace sampling.

#### TraceRCA handles 10,000 traces per second per core.





Fig. 13: Comparison of efficiency



Fig. 14: Efficiency improvement and performance degradation of *TraceRCA* with trace sampling.

#### TraceRCA handles 10,000 traces per second per core.

#### The running time of TraceRCA is linear to #traces.







Fig. 13: Comparison of efficiency



Fig. 14: Efficiency improvement and performance degradation of *TraceRCA* with trace sampling.

#### TraceRCA handles 10,000 traces per second per core.

The running time of TraceRCA is linear to #traces.

**TraceRCA** achieves relative good performance with only 1/10 sampled traces.









Motivation: Trace-based root-cause service localization; our insight and two key metrics



Motivation: Trace-based root-cause service localization; our insight and two key metrics



Solution: TraceRCA, containing trace anomaly detection, microservice set mining, suspicious microservice ranking



Motivation: Trace-based root-cause service localization; our insight and two key metrics



Solution: TraceRCA, containing trace anomaly detection, microservice set mining, suspicious microservice ranking



Evaluation: Experiments based on an open-source benchmark and a production system

