



A survey of root cause analysis

霍茵桐

2021/06/25 Workshop



香港中文大學
The Chinese University of Hong Kong

Problem Formulation

GOAL

Target to identify the underlying causes leading to a failure that has affected end users, which is often closely related to the failure diagnosis.

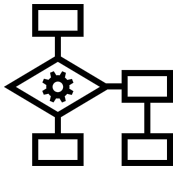
SOURCE

In the broad domain of root cause analysis, log-based failure diagnosis is now a standard practice.

CHALLENGES

- The complexity of modern software systems grows rapidly, making too complex to efficiently disclose relations of fault, failure and human-observed symptoms.
- As software systems become more mature, failures are becoming more and more hard to detect and diagnose.

In the very beginning...



Traditional rule-based methods are adopted.

- which heavily rely on a set of predefined rules (if-then logic) from the expert knowledge to diagnose failures.



Limitations:

Rule-based methods cannot be well generalized to unseen failures that are not included in the rules.

Outline

#1: Execution Replay Methods

#2: Model-based Methods

#3: Statistics-based Methods

#4: Retrieval-based Methods

Outline



#1: Execution Replay Methods

- Infer execution flow
- Trace back to software failure

#2: Model-based Methods

#3: Statistics-based Methods

#4: Retrieval-based Methods

Execution Replay Methods

Approach

Automatically infer the execution flow from logs and trace back to the software system failure.

Input:

- Log recorded during the failed execution
- Source code of the target program

Output:

- Derive the potential code paths and execution path

Execution Replay Methods (1)

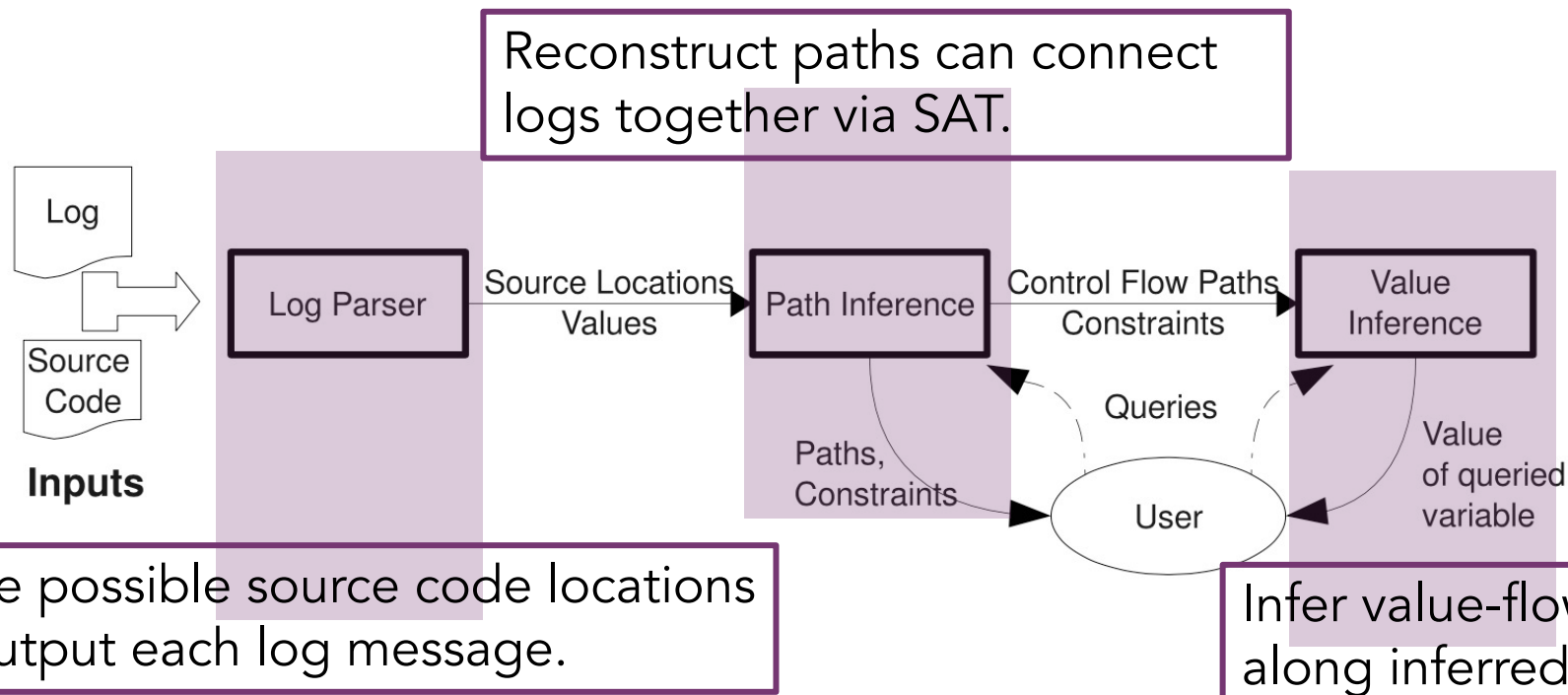
LogMap Approach

1. LogMap *retrieves log messages* from bug reports.
2. Apply static analysis technique to *identify corresponding logging lines* in source code.
3. Traverse through logging lines to *derive the potential code paths*, which help reconstruct the execution path and assist debugging process.

Execution Replay Methods (2)

Goal: Find *information* from printed log

- *Must-have*: Partial execution paths that were definitely executed.
- *May-Have*: Partial execution paths that may have been executed
- *Must-Not-Have*: Execution paths that were definitely NOT executed.



Execution Replay Methods (2)

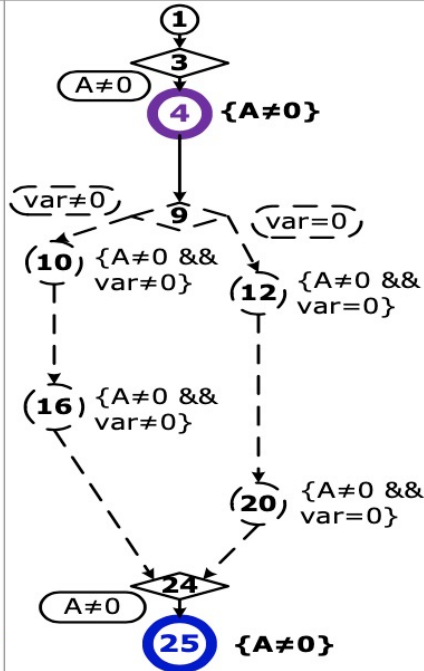
A log:

msg1
msg2

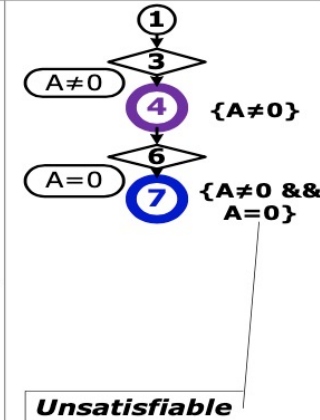
```

1 main() {
2   ...
3   if (A)
4     log(msg1);
5
6   if (!A)
7     log(msg2);
8
9   if (var)
10    b1();
11  else
12    b2();
13 }
14
15 b1() {
16   c();
17 }
18
19 b2() {
20   c();
21 }
22
23 c() {
24   if (A)
25     log(msg2);
26 }
    
```

Must-Path & May-Path:



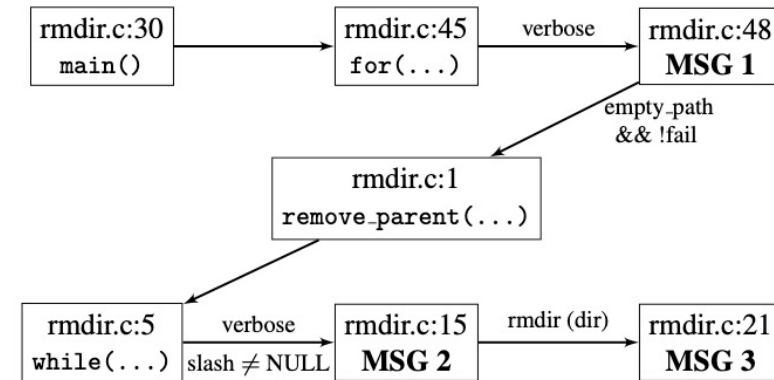
Pruned-Path:



SherLog Report for rmdir

Paths

Path 1:



Path 2:

...

Value Inference

Query results of the value of path in remove_parents():
 path = "dir1/dir2/" on Line 3
 path = "dir1/dir2" on Line 11

Outline

#1: Execution Replay Methods



#2: Model-based Methods

- Build reference model
- Check violate

#3: Statistics-based Methods

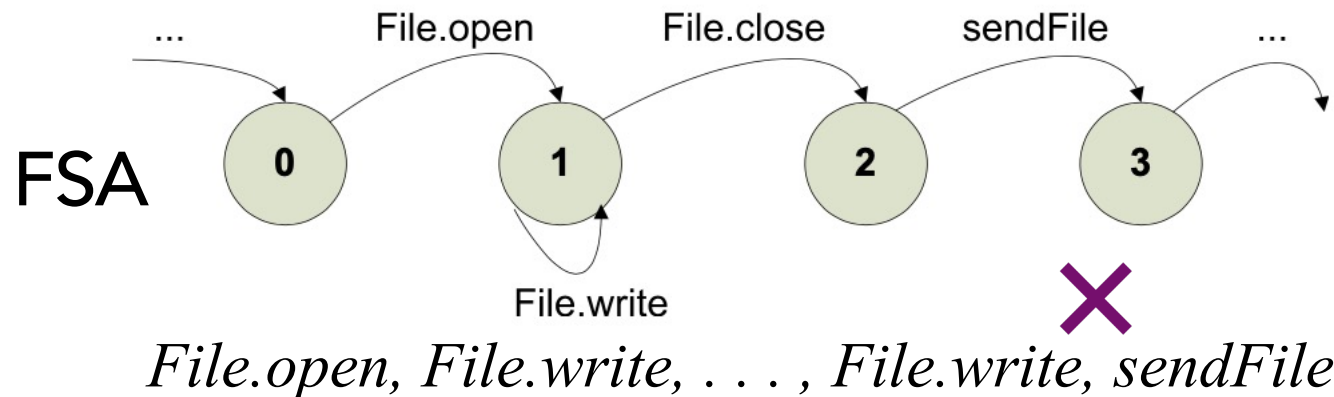
#4: Retrieval-based Methods

Model-based Methods

Approach

Utilize logs to build the reference model for a software system, then check which log events violate the reference model.

- 1 *Generate a model* of the legal behavior of a target system by tracing successful executions at testing-time.
- 2 *Compare* the set of events detected during failing executions with the generated models.



Model-based Methods (1)

How to generate reference model from log?

How to compare?

Raw Logs:

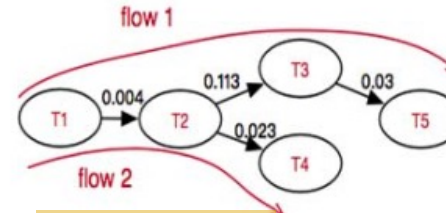
```
(1)2016-10-10T01:00:18.866Z parse_params: {"inline-relations-depth"=>"1"} flow 1
(2)2016-10-10T01:00:18.867Z parse_params: {"inline-relations-depth"=>"1", "q"=>"name:test-app-1476011110087"} flow 2
(3)2016-10-10T01:00:18.870Z cc.dispatch host: 9e65a679-f602-4366-bfe8-24b3efa73230 flow 1
(4)2016-10-10T01:00:18.957Z Request failed: 400: {"code"=>"170002", "description"=>"App has not finished staging"} parameters flow 1
(5)2016-10-10T01:00:18.964Z cc.dispatch host: 9e65a679-f602-4366-bfe8-24b3efa73230 flow 2
(6)2016-10-10T01:00:18.987Z Completed 200 vcap-request-id: d89bb92b-bcc5-467b-4120-54e82eac7873::2faebb11-7343-4d95-bb75-a7e7e59eedea flow 2
(7)2016-10-10T01:00:18.987Z Completed 400 vcap-request-id: f010f420-04f5-425a-4aab-a38bbcc18a9f::bf38597a-acc8-4180-8bb6-943c356b5d40 flow 1
```

Log Templates:

```
1)parse_params: {"*"}
2)parse_params: {"*"}
3)cc.dispatch host: *
4)Request failed: 400:{"code"=>*, "description"=>""}
5)cc.dispatch host: *
6)Completed 200 vcap-request-id: *
7)Completed 400 vcap-request-id: *
```

T1
T1
T2
T3
T2
T4
T5

Ground Truth time-weighted CFG:



Mine time-weighted control flow graphs (TCFG) from interleaved logs of a system.

Node: Template

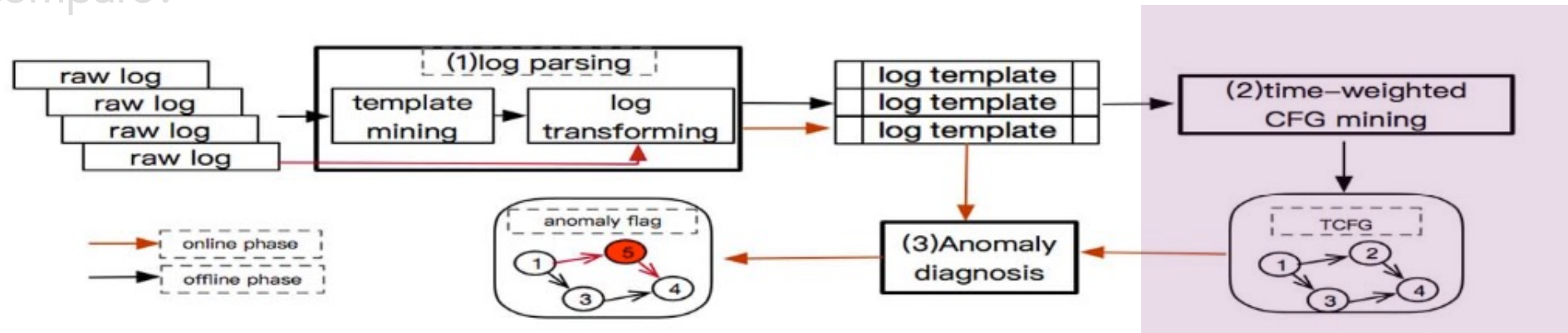
Edge: Transition from one template to another

Weighted: Execution time between templates

Model-based Methods (1)

● How to generate reference model from log?

● How to compare?



Challenges: Without a certain transaction ID.

Motivation: Logs of correct immediately succeeding template will appear in a short period.

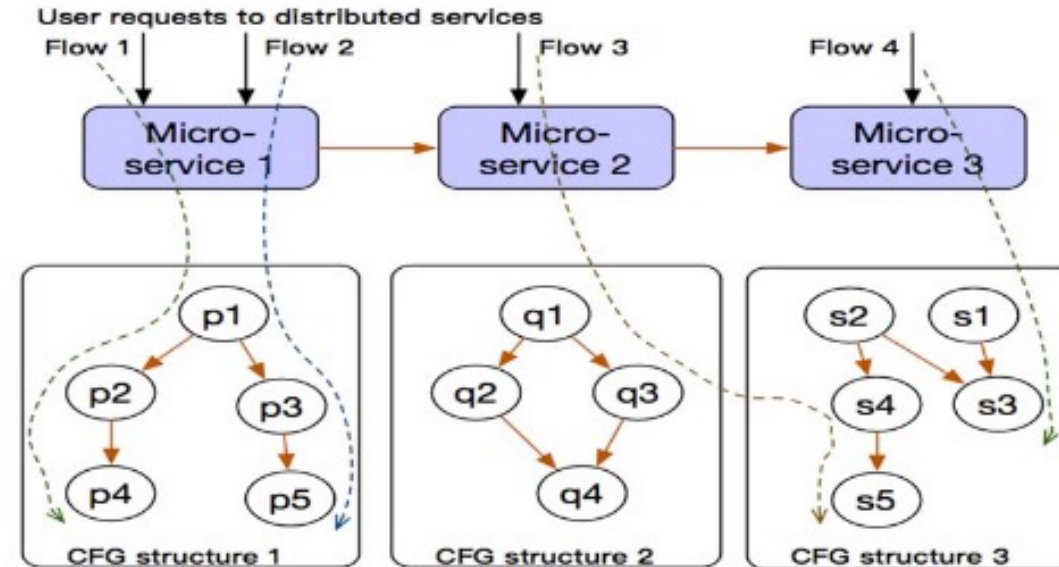
Two Stages:

- (1) Compute frequent successor groups called FS groups of each template.
- (2) Leverage this FS group to mine only the immediate successors of each template.

Model-based Methods (2)

● How to generate reference model from log?

● How to compare?



Model the causal relationships between services as a **service topology**.

***Node:** Service*

***Edge:** Causal relationship between two services*

Mining approach: PC-algorithm

- How to generate reference model from log?
- How to compare?



The diagram illustrates the mapping of user requests to distributed services and their corresponding CFG structures. At the top, four user request flows (Flow 1, Flow 2, Flow 3, Flow 4) are shown as arrows pointing to three micro-services: Micro-service 1, Micro-service 2, and Micro-service 3. The flows are distributed as follows: Flow 1 and Flow 2 point to Micro-service 1; Flow 3 points to Micro-service 2; and Flow 4 points to Micro-service 3. Below each micro-service, a CFG structure is shown, representing the internal state transitions of that service. The CFG structures are labeled CFG structure 1, CFG structure 2, and CFG structure 3. CFG structure 1 (under Micro-service 1) has nodes p1, p2, p3, p4, and p5. CFG structure 2 (under Micro-service 2) has nodes q1, q2, q3, and q4. CFG structure 3 (under Micro-service 3) has nodes s1, s2, s3, s4, and s5. Dashed arrows indicate the mapping from the user request flows to the specific nodes in the CFG structures: Flow 1 maps to p1 and p4; Flow 2 maps to p1 and p5; Flow 3 maps to q1 and q4; and Flow 4 maps to s1 and s5. Solid orange arrows show the sequence of micro-services: Micro-service 1 to Micro-service 2 to Micro-service 3.

Examine the TCFGs of its predecessor services in the service topology.

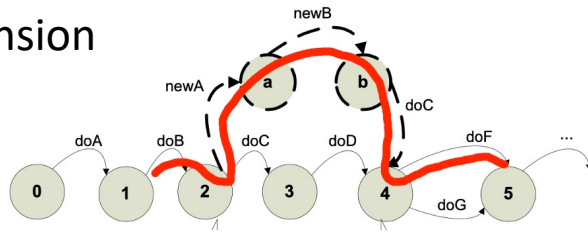
Model-based Methods

How to generate reference model from log?

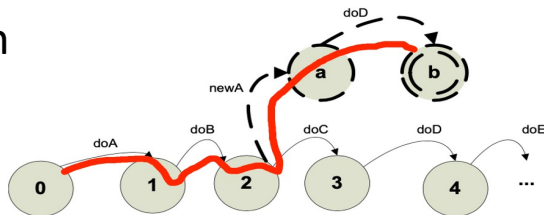
How to compare?

1 Select the expected behaviors to be considered for local analysis.

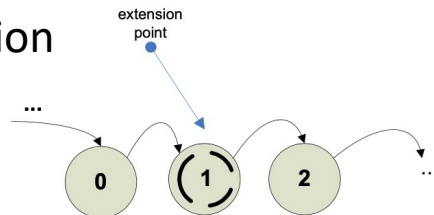
Branch extension



Tail extension



Final state extension



(delete, insert, replace and final state)

2 Compare expected behavior with the observed failed sequence.

Example insert interpretation:

expected sequence: doA doB - - - doC

observed sequence: doA doE doF doG doH doC

Example delete interpretation:

expected sequence: doL doL doF doG - doM

observed sequence: doL - - - doN doM

Sets specific to anticipation

$I.newEvents = doE doF doG doH$

$D.removedEvents = doL doF doG$

Aligned sequences:

doE doF doG doH

doL doF doG -

$align = doF doG$

$notAlign.n = doL$

$notAlign.g = doH$

General sets

$I.m = doA doC$

$D.m = doL doM$

$I.d = /$

$D.a = doN$

$$score = \frac{2*2-2*1-1+2+2-0-1}{6+6} = 0.33$$

confidence value = $(0.33+1)/2=0.665$
(normalized score)

Outline

#1: Execution Replay Methods

#2: Model-based Methods



#3: Statistics-based Methods

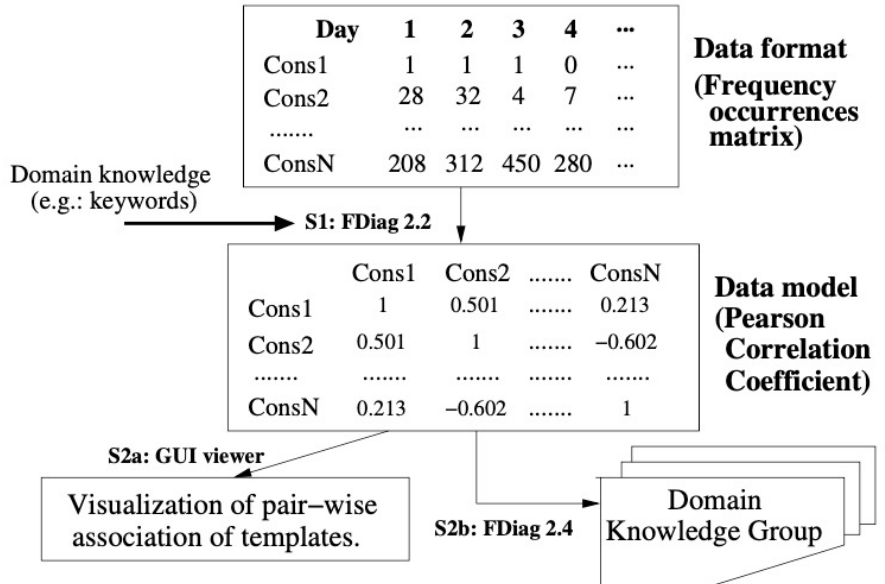
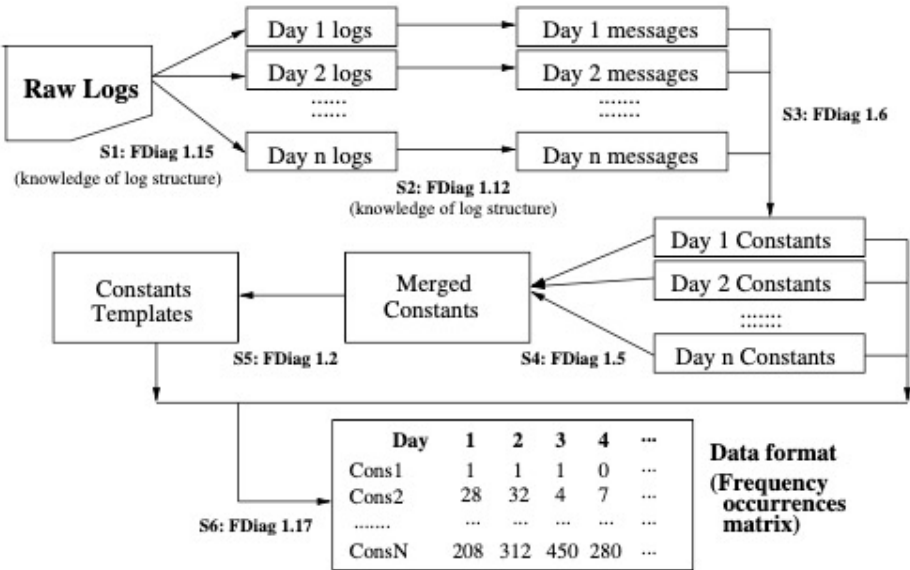
- Statistical analysis
-

#4: Retrieval-based Methods

Approach

Employ some statistical techniques to capture the relationships between logs and the consequent failures.

*Given knowledge of any system event, find all the system events that are **associated** with the given event, and **how strongly associated** these system events are to the given event.*



Outline

#1: Execution Replay Methods

#2: Model-based Methods

#3: Statistics-based Methods



#4: Retrieval-based Methods

- Information Retrieval

Retrieval-based Methods

Approach

Retrieve similar failures in a knowledge base composed of failures in history.

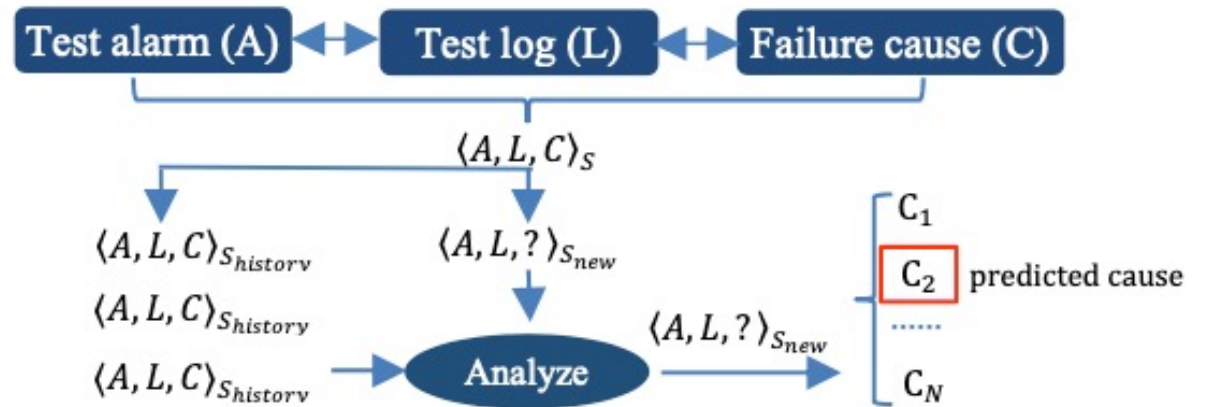
Motivation: In practice, failures that previously occurred are valuable since they can aid developers in better diagnosing newly-occurred failures

Retrieval-based Methods

Basic Idea: Detect the test logs of historical test alarms that may share the same causes with the new test logs.



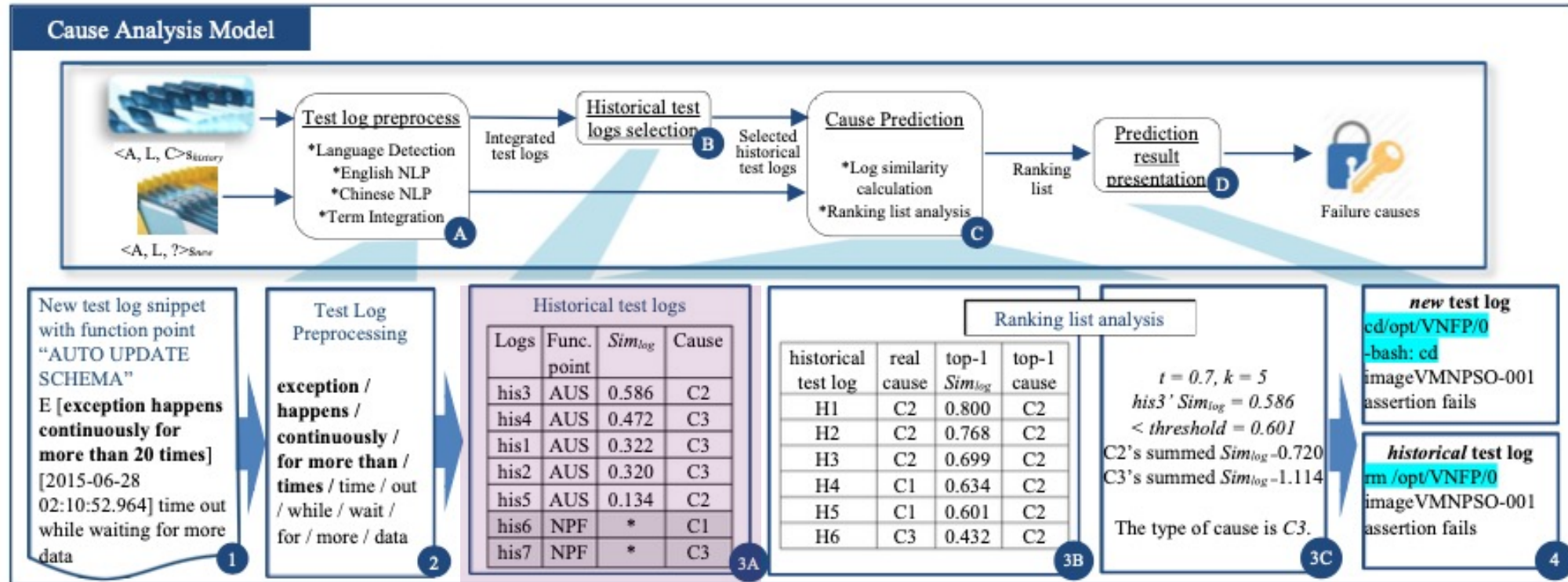
Bug Report



Task Formulation: Give $\langle A, L, ? \rangle$, predict correct C

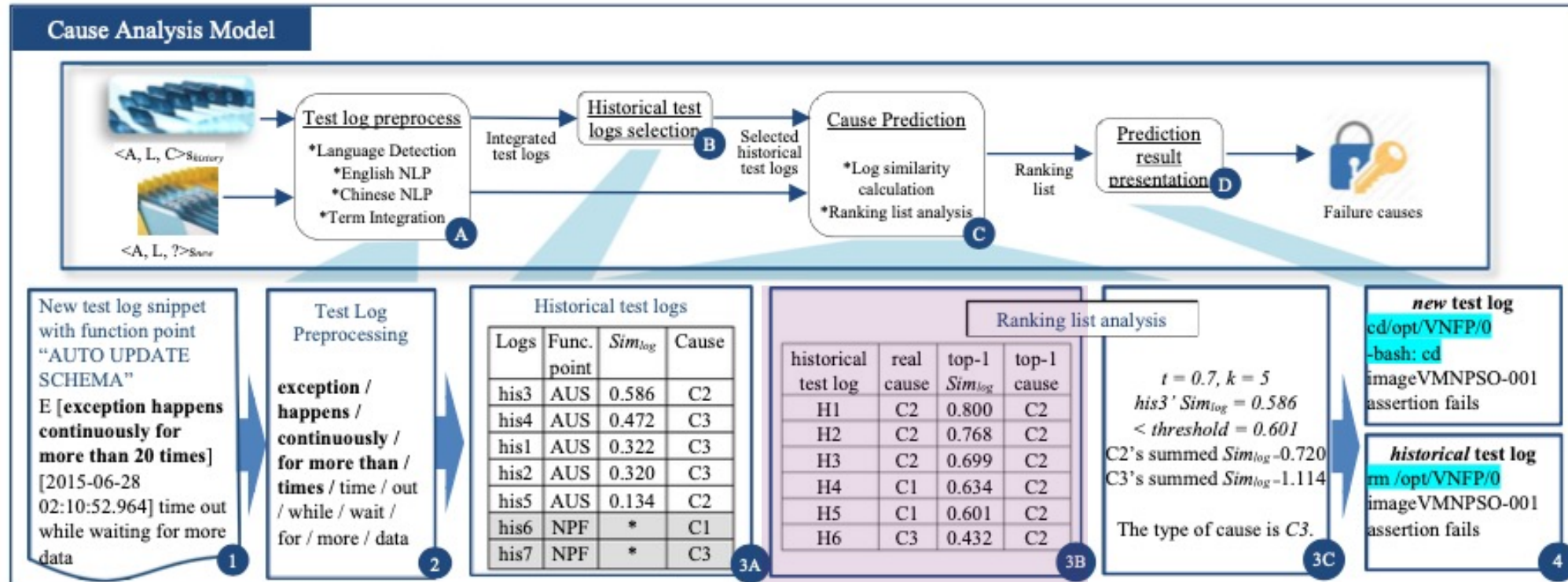
Evaluation: Test in two industrial testing projects at Huawei-Tech Inc.

Retrieval-based Methods



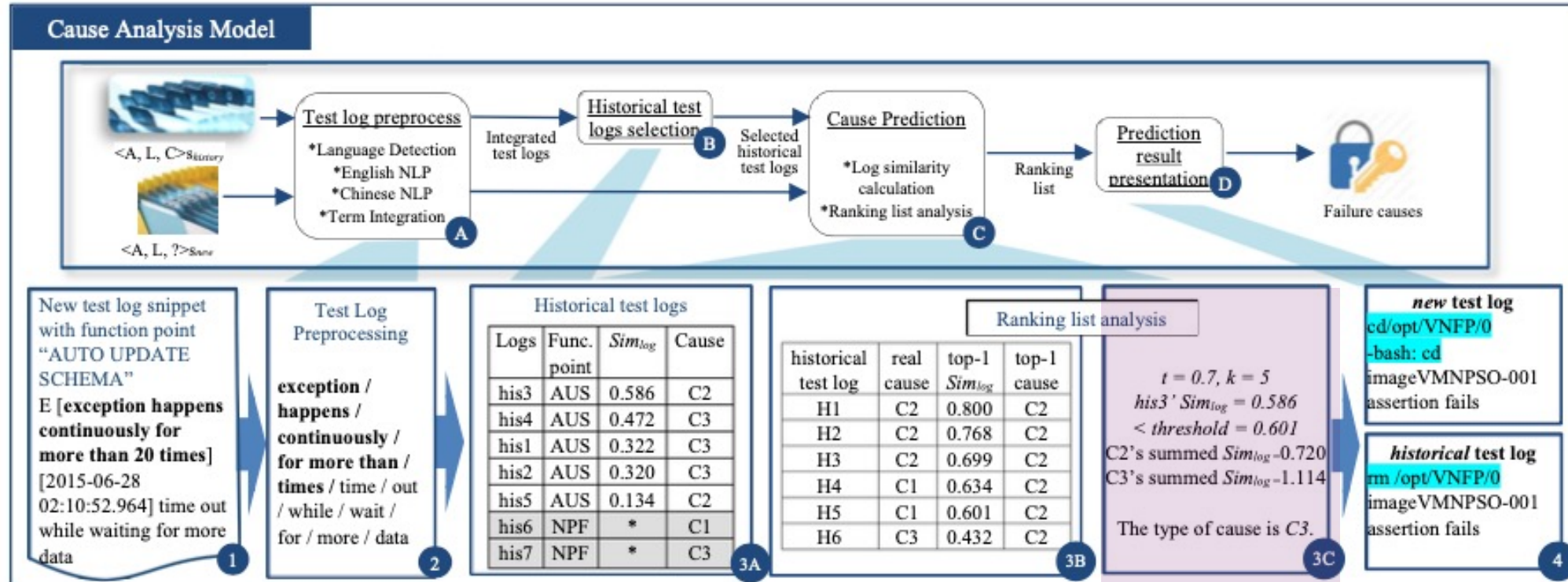
Selects historical test logs by examining the *function points of test scripts*, since the authors find that test scripts with the same function point usually target the same functionalities to check.

Retrieval-based Methods



Predict the cause of a new test alarm by first ranking the selected historical test logs according to their similarities with the new test log.

Retrieval-based Methods



Normally predict the test alarm cause with the **top-1 cause** in the ranking list. If the similarity is low, use **KNN strategy**.

Retrieval-based Methods

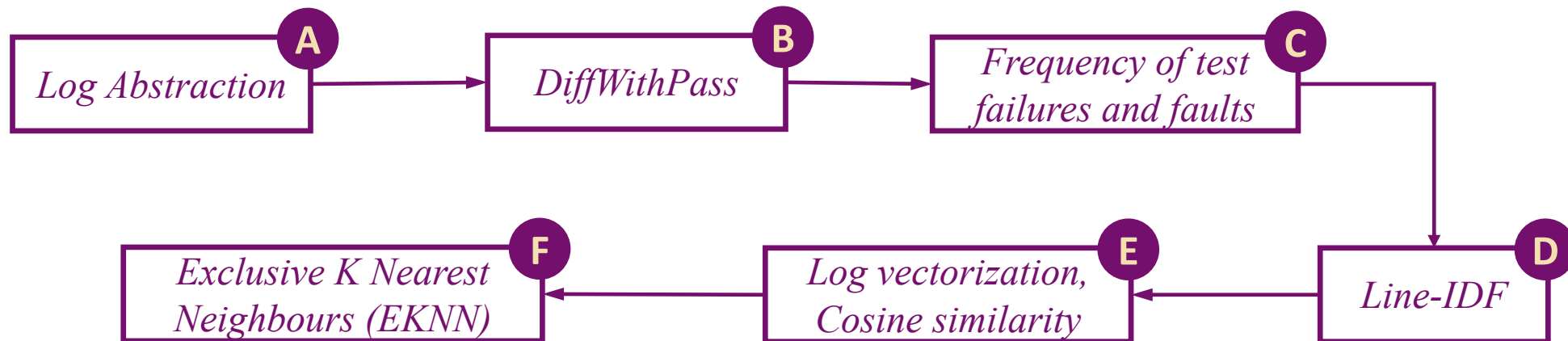
Previous limitation and new challenge:

- The complex test environment introduces many non-product test failures.
- The logs contain an overwhelming amount of detailed information.

#FaultsFound

#LogLinesFlagged

Approach:



Retrieval-based Methods

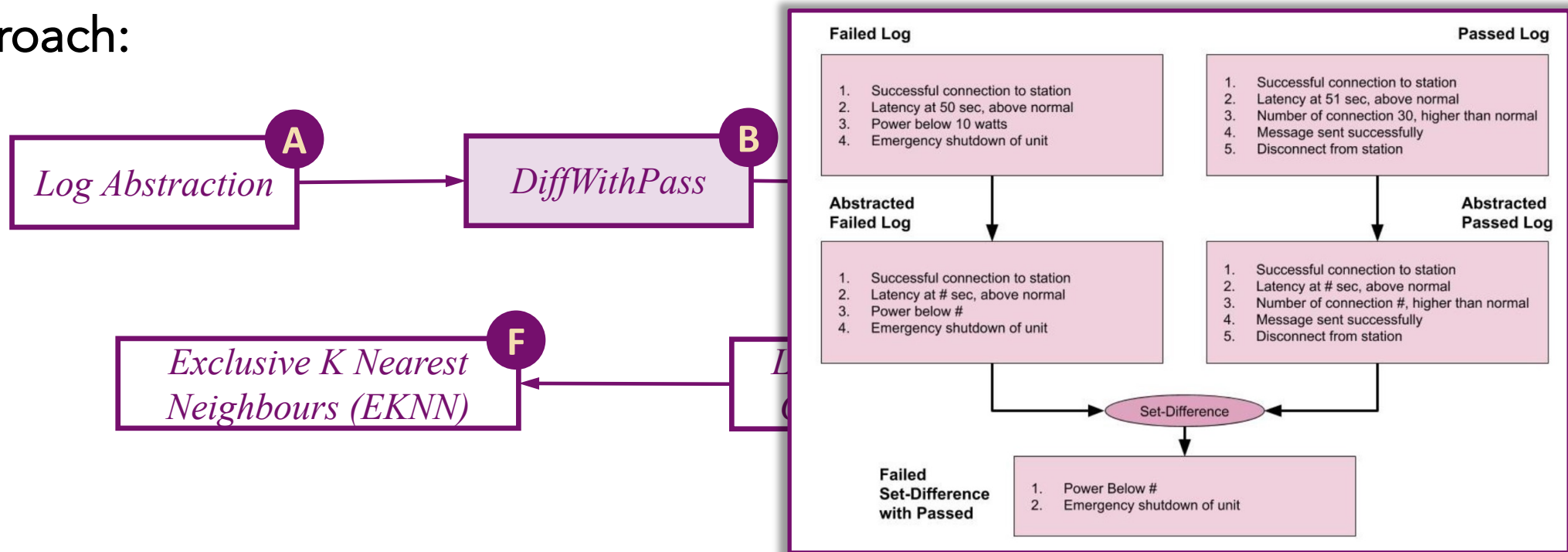
Previous limitation and new challenge:

- The complex test environment introduces many non-product test failures.
- The logs contain an overwhelming amount of detailed information.

#FaultsFound

#LogLinesFlagged

Approach:



Retrieval-based Methods

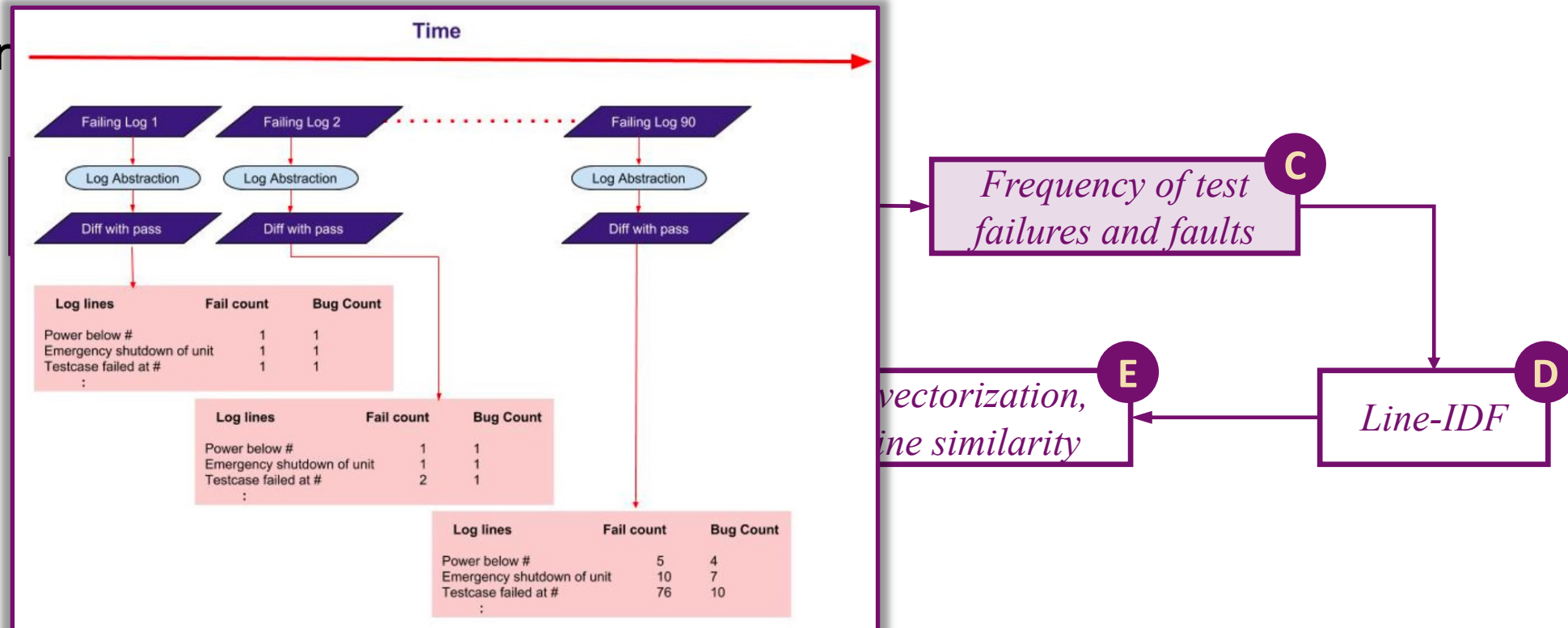
Previous limitation and new challenge:

- The complex test environment introduces many non-product test failures.
- The logs contain an overwhelming amount of detailed information.

#FaultsFound

#LogLinesFlagged

Appro



Retrieval-based Methods

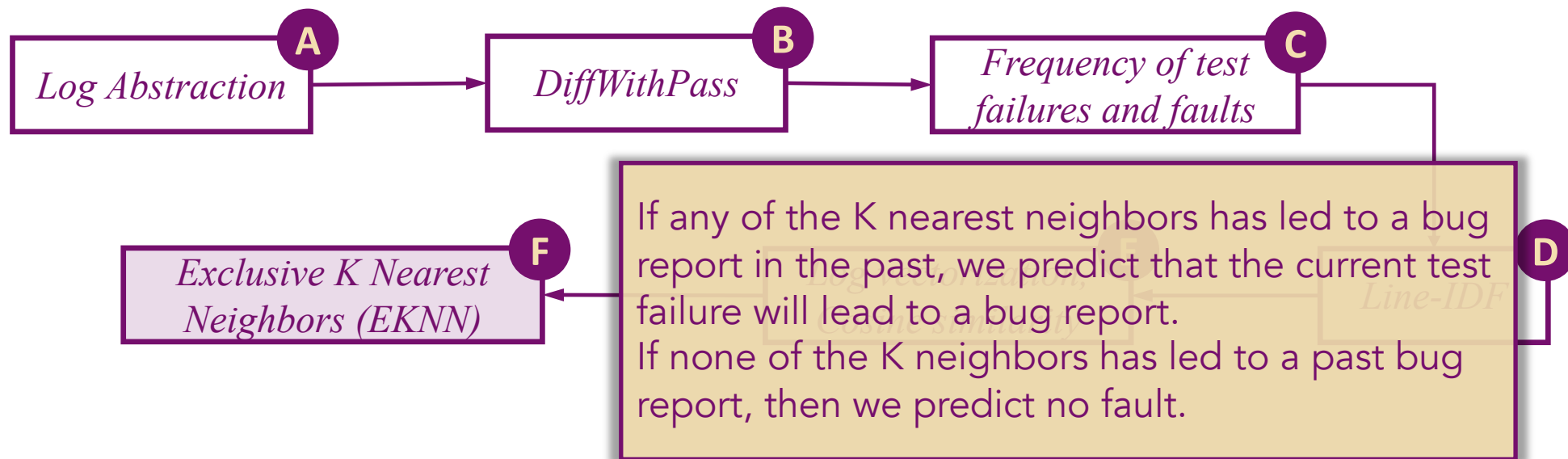
Previous limitation and new challenge:

- The complex test environment introduces many non-product test failures.
- The logs contain an overwhelming amount of detailed information.

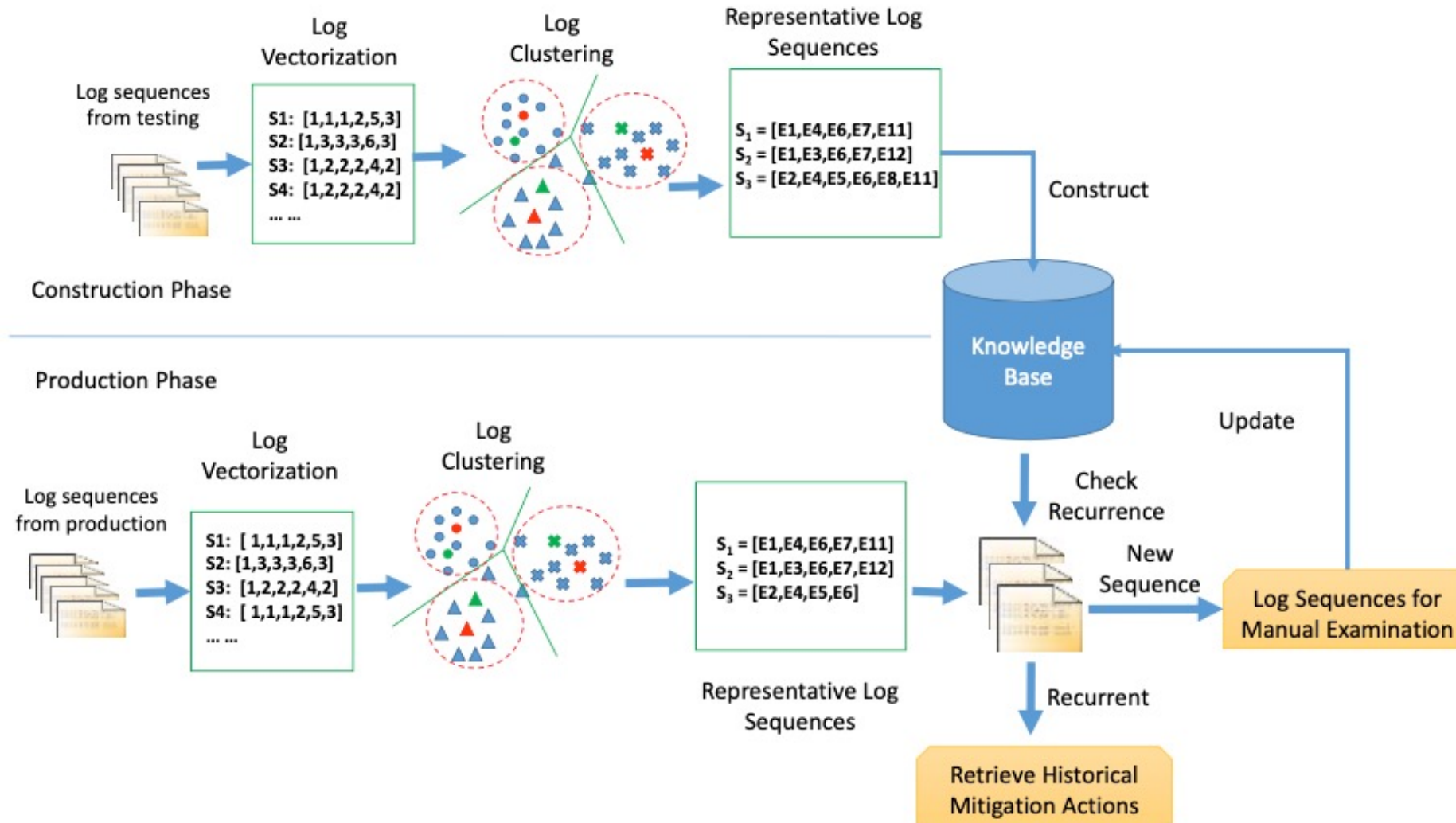
#FaultsFound

#LogLinesFlagged

Approach:



Retrieval-based Methods



Conclusion

#1: Execution Replay Methods

- Infer execution flow
- Trace back to software failure

#2: Model-based Methods

- Build reference model
- Check violate

#3: Statistics-based Methods

- Statistical analysis

#4: Retrieval-based Methods

- Information Retrieval

Thank you!



香港中文大學
The Chinese University of Hong Kong