# Automatically pinpointing original logging functions from log messages for network troubleshooting
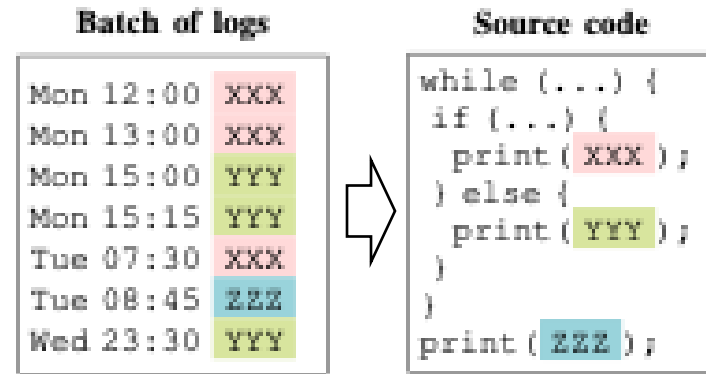
Gaspard Damoiseau-Malraux[1], Satoru Kobayashi[2], Kensuke Fukuda[3]

1: Sorbonne Université, 2: Okayama University, 3: NII/Sokendai

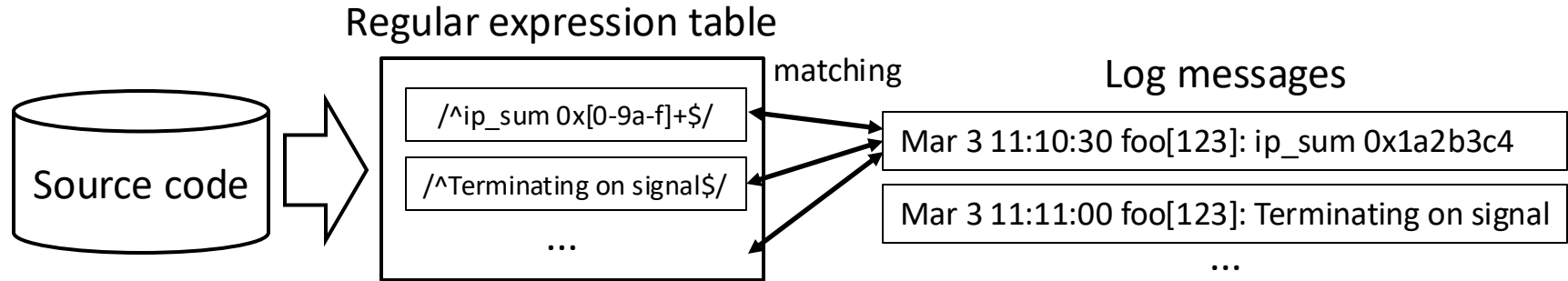COMPSAC NCIW 2025, July 11, 2025

1

# Introduction

- Network log messages
  - Huge amount
  - Often lack details to explain failures
  - ➤ Further analysis is needed
- Source code analysis to find log origins
  - Identifying <u>original logging functions</u>
  - Effective to understand system behaviors
    - e.g., What happens in the system when the log message appear?

**Batch of logs**

```
Mon 12:00  XXX
Mon 13:00  XXX
Mon 15:00  YYY
Mon 15:15  YYY
Tue 07:30  XXX
Tue 08:45  ZZZ
Wed 23:30  YYY
```

**Source code**

```
while (...) {
 if (...) {
  print( XXX );
 } else {
  print( YYY );
 }
}
print( ZZZ );
```
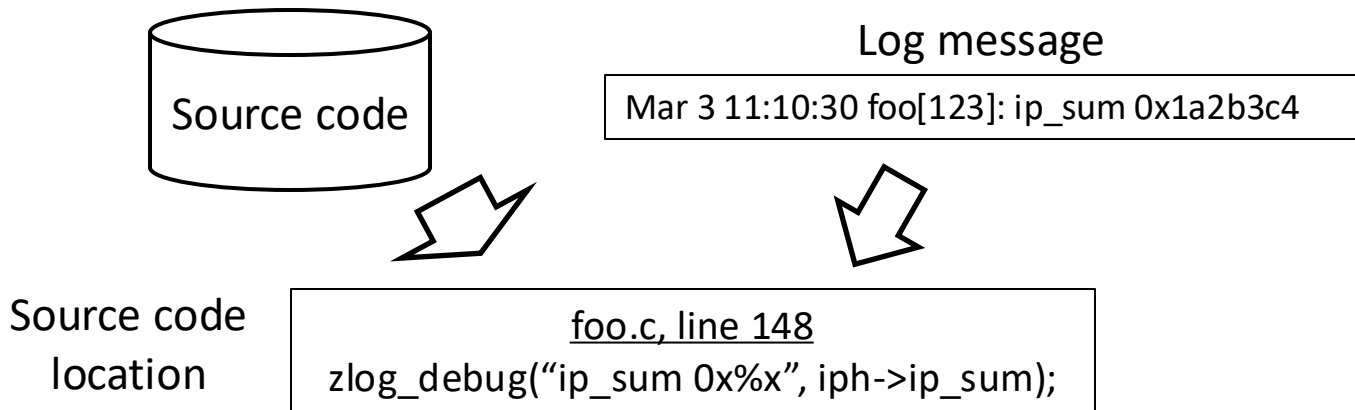
# Existing approach and challenge

- Match input log messages with <u>regular expression (regex) patterns</u> extracted from logging functions in source code [3]

- Challenge: <u>Large processing time</u> due to matching all regex patterns for each log input

Regular expression table

matching

Log messages

| Source code | Regular expression table | Log messages |

/^ip_sum 0x[0-9a-f]+$/

/^Terminating on signal$/

...

Mar 3 11:10:30 foo[123]: ip_sum 0x1a2b3c4

Mar 3 11:11:00 foo[123]: Terminating on signal

...

[3] W. Xu, et al. "Detecting large-scale system problems by mining console logs," Proceedings of SOSP '09, pp.117–132, 2009.
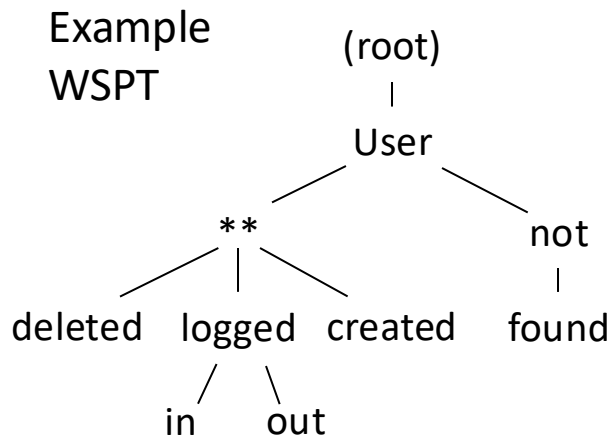
# Research goal

- Design and develop a system to automatically identify original logging functions from source code
  - Fast log matching with new hybrid approach
  - Implement SCOLM (Source Code Origins of Log Messages) for C programs

Log message

Source code

Mar 3 11:10:30 foo[123]: ip_sum 0x1a2b3c4

Source code location

foo.c, line 148
zlog_debug("ip_sum 0x%x", iph->ip_sum);
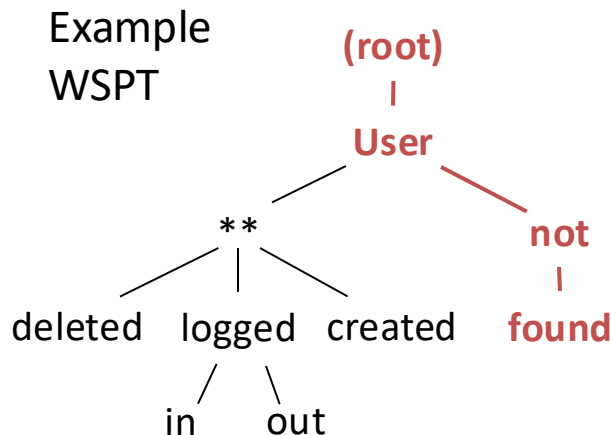
4

# Existing work on log matching

- ## WSPT (Word Segmented Prefix Tree) [4]
  - Prefix tree for words or wildcards
  - A wildcard node corresponds to only one word (i.e., no separators included)
    - WS-templates: satisfying this rule
  - Fast log matching
    - Reduce processing time for log classification by 1/10 compared to regex approach [4]

Example WSPT

```
                    (root)
                      |
                    User
                  /        \
               **           not
            /   |    \        |
      deleted logged created found
            /      \
          in       out
```

- User not found
- User Alice logged in
- User Bob created
- User Charlie deleted

[4] S. Kobayashi, et al. "amulog: A general log analysis framework for comparison and combination of diverse template generation methods," International Journal of Network Management, vol.32, no.4, p.e2195, 2022.
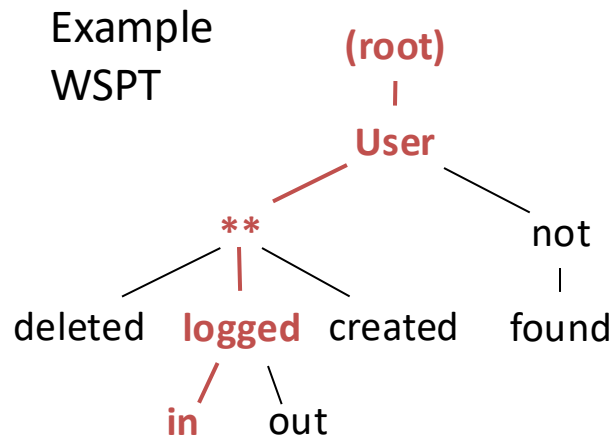
# Existing work on log matching

- ## WSPT (Word Segmented Prefix Tree) [4]

  - Prefix tree for words or wildcards

  - A wildcard node corresponds to only one word (i.e., no separators included)

    - WS-templates: satisfying this rule

  - Fast log matching

    - Reduce processing time for log classification by 1/10 compared to regex approach [4]

Example WSPT



- **User not found**
- User Alice logged in
- User Bob created
- User Charlie deleted

[4] S. Kobayashi, et al. "amulog: A general log analysis framework for comparison and combination of diverse template generation methods," International Journal of Network Management, vol.32, no.4, p.e2195, 2022.

# Existing work on log matching

- WSPT (Word Segmented Prefix Tree) [4]
  - Prefix tree for words or wildcards
  - A wildcard node corresponds to only one word (i.e., no separators included)
    - WS-templates: satisfying this rule
  - Fast log matching
    - Reduce processing time for log classification by 1/10 compared to regex approach [4]

Example WSPT



- User not found
- **User Alice logged in**
- User Bob created
- User Charlie deleted

[4] S. Kobayashi, et al. "amulog: A general log analysis framework for comparison and combination of diverse template generation methods," International Journal of Network Management, vol.32, no.4, p.e2195, 2022.

# Problem of WSPT

- Log formats from source code are <u>not available as WS-templates</u> for WSPT
  - One format specifier may embed multiple words
  - Not satisfying the assumption "<u>one wildcard corresponds to one word</u>"

Regex approach
Format specifiers can
<u>match multiple words</u>

Regular expression:

/^sshd: user ".*" login$/

✔ sshd: user "sat" login

✔ sshd: user "oka taro" login

WSPT approach
Wildcards can <u>match only one word</u>

Word-segmented template:

sshd: user <u>**</u> login

✔ sshd: user "<u>sat</u>" login

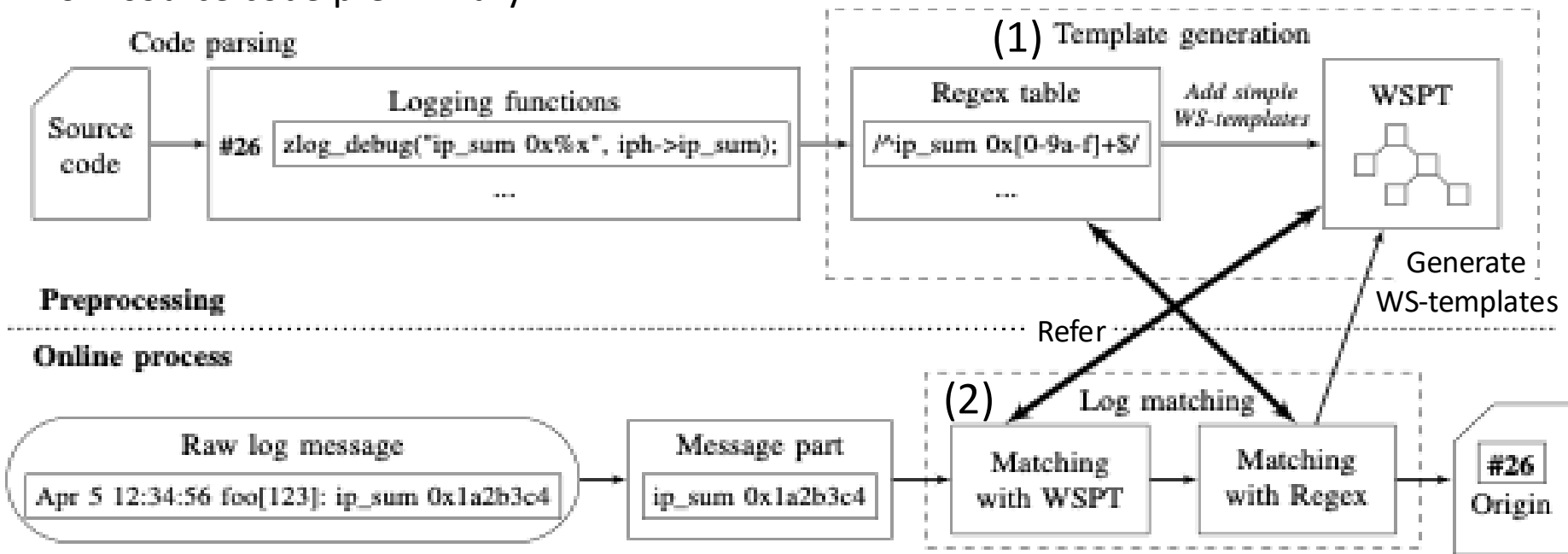✘ sshd: user "<u>oka taro</u>" login

# Proposed approach

- Combining regex tables and WSPT
  - Generate WS-template after matching regex patterns and log messages
  - ➢Number of words is determined for the messages
- <u>Fast processing of logs appearing after the second time</u>

| Method | Processing time | External templates | Overmatching |
|---|---|---|---|
| Regex table | Large | Available | Small |
| WSPT | Small | Not available | Medium |
| Proposed method | Small | Available | Small |

# Overview of proposed method

Generate regex patterns
from source code preliminary

Add simple WS-
templates to WSPT
(For fast processing)

Code parsing

(1) Template generation

| Source code | #26 | Logging functions | | Regex table | Add simple WS-templates | WSPT |
| | | zlog_debug("ip_sum 0x%x", iph->ip_sum); | | /^ip_sum 0x[0-9a-f]+$/ | | |
| | | ... | | ... | | |

Generate
WS-templates

**Preprocessing**

Refer

**Online process**

(2) Log matching

| Raw log message | | Message part | | Matching with WSPT | Matching with Regex | #26 Origin |
| Apr 5 12:34:56 foo[123]: ip_sum 0x1a2b3c4 | | ip_sum 0x1a2b3c4 | | | | |

Use message part (no header)
in the input log messages

Log matching combining
regex and WSPT

10

# (1) Template generation

1. Extract logging functions from source code
2. Replace format specifiers into regex patterns
3. Generate simple WS-template that has one wildcard for one format specifier

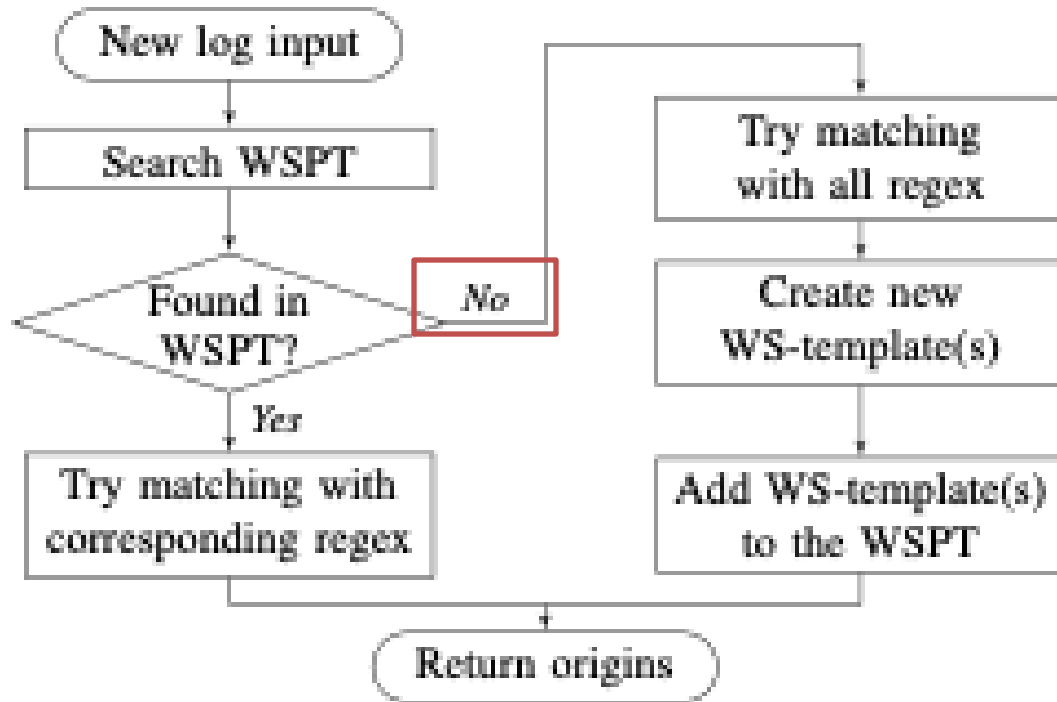logging function | zlog_debug("ip_sum 0x%x", iph->ip_sum)

regex pattern | /^ip_sum\ 0x[0-9a-f]+$
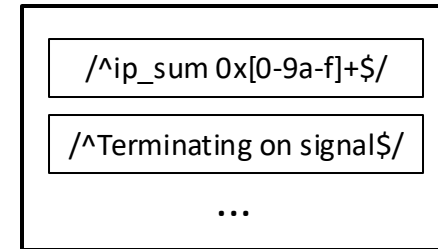
(simple) WS-template | ip_sum **

# (2) Log matching

New log input

Search WSPT

Found in WSPT?

No → Try matching with all regex → Create new WS-template(s) → Add WS-template(s) to the WSPT

Yes → Try matching with corresponding regex

Return origins

Input

ip_sum 0x1a2b3c4

ip_sum 0x2b3c4d5

…

Regex table

/^ip_sum 0x[0-9a-f]+$/

/^Terminating on signal$/

…

WSPT

no matching template

# (2) Log matching



New log input

Search WSPT

Found in WSPT?

No

Try matching with all regex

Create new WS-template(s)

Add WS-template(s) to the WSPT

Yes

Try matching with corresponding regex

Return origins

Input → ip_sum 0x1a2b3c4

ip_sum 0x2b3c4d5

...

Regex table

/^ip_sum 0x[0-9a-f]+$/    matching

/^Terminating on signal$/    matching

...    ...

matching all patterns

WSPT

(root)

ip_sum    ...

**

added    13

# (2) Log matching

New log input

Search WSPT

Found in WSPT?

*No*

*Yes*

Try matching with corresponding regex

Try matching with all regex

Create new WS-template(s)

Add WS-template(s) to the WSPT
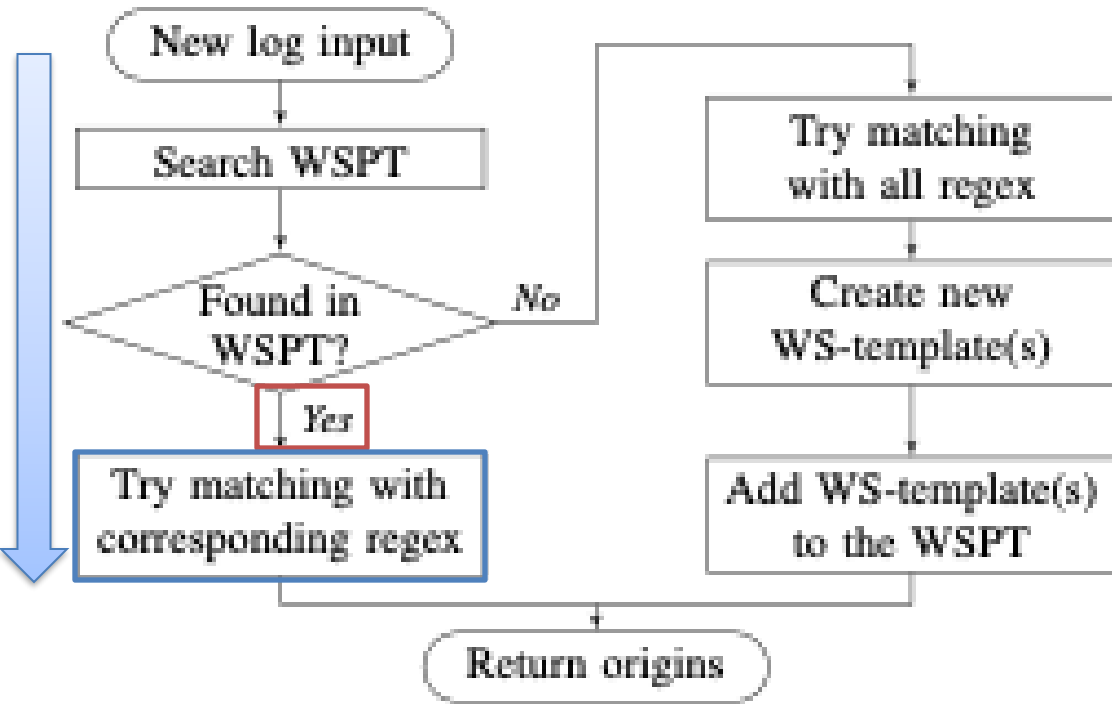
Return origins

Input

ip_sum 0x1a2b3c4

ip_sum 0x2b3c4d5

...

Regex table

/^ip_sum 0x[0-9a-f]+$/    matching

/^Terminating on signal$/

...

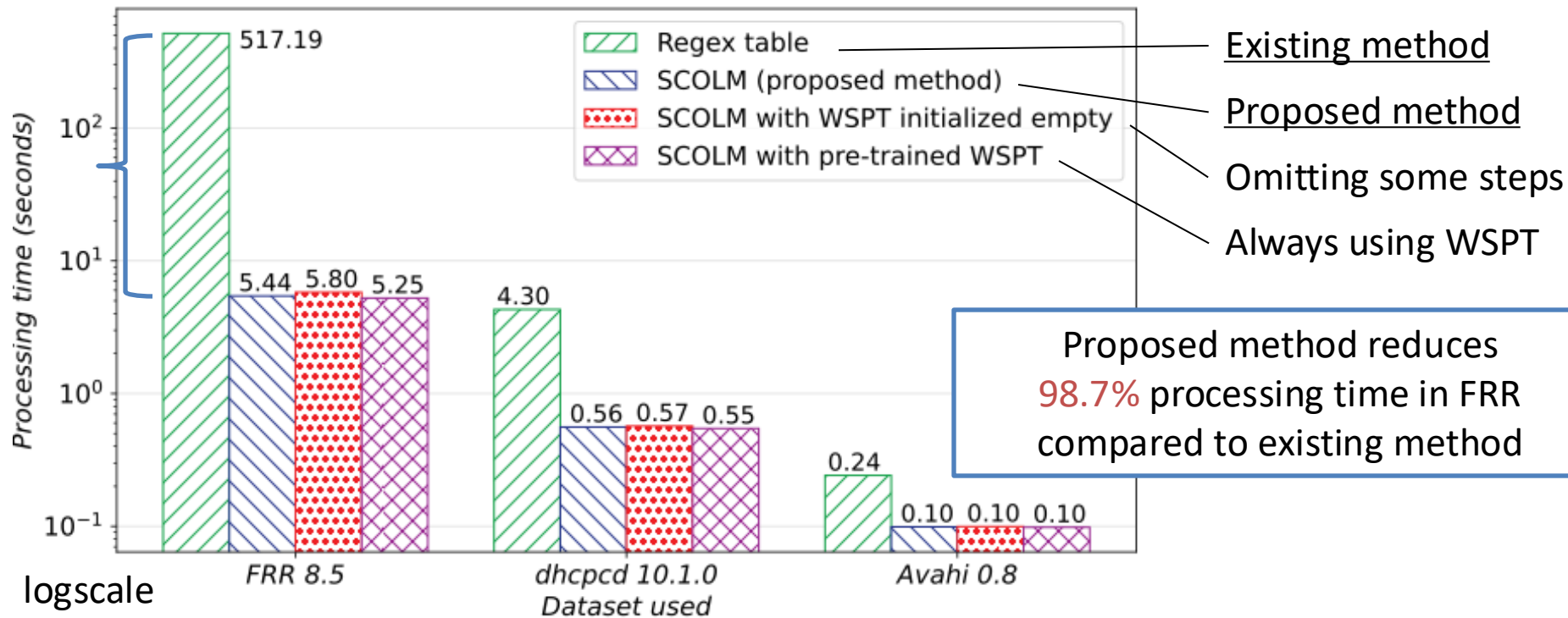matching corresponding patterns

WSPT

(root)

ip_sum          ...

**

found                                    14

# Evaluation setup

- Implementation
  - SCOLM: Implemented in Python 3.10, available in GitHub

- Datasets

http://github.com/3atlab/scolm

| Software | # of regex tpl. | # of WS tpl. | # of logs available | How to collect logs |
|---|---|---|---|---|
| FRRouting (FRR) 8.5 | 6,481 | 6,346 | 89,856 | Network emulation [23] |
| dhcpcd 10.1.0 | 586 | 580 | 5,509 | Home server |
| Avahi 0.8 | 367 | 363 | 1,971 | Lab server |

[23] C. Regal-Mezin, et al. "netroub: Towards an emulation platform for network trouble scenarios", in CoNEXT-SW, pp.17-18, 2023.

# Comparison on processing time



Regex table — Existing method
SCOLM (proposed method) — Proposed method
SCOLM with WSPT initialized empty — Omitting some steps
SCOLM with pre-trained WSPT — Always using WSPT

Proposed method reduces 98.7% processing time in FRR compared to existing method

logscale

Measurement in AMD Ryzen 7 7700, Ubuntu Server 22.04 (x86_64), 64GB memory

# Processing time for each log message



Converge to pure WSPT performance with sufficient input data

Requires large processing time for input log that does not match WSPT (for full regex match)

log-log-scale

# Number of obtained candidates

| # of candidates | Existing method Regex table only | SCOLM (proposed method) |
|---|---|---|
| | FRR 8.5 | |
| 0 candidates | 0 (00.00%) | 0 (00.00%) |
| 1 candidate | 80904 (90.04%) | 88310 (98.28%) |
| 2 candidate | 6972 (07.76%) | 1546 (01.72%) |
| 3 candidates | 1980 (02.20%) | 0 (00.00%) |
| Total match rate | 89856 (100.0%) | 89856 (100.0%) |
| | dhcpcd 10.1.0 | |
| 0 candidates | 498 (09.04%) | 498 (09.04%) |
| 1 candidate | 4882 (88.62%) | 5011 (90.96%) |
| 2 candidate | 129 (02.34%) | 0 (00.00%) |
| Total match rate | 5011 (90.96%) | 5011 (90.96%) |
| | Avahi 0.8 | |
| 1 candidate | 1971 (100.0%) | 1971 (100.0%) |
| Total match rate | 1971 (100.0%) | 1971 (100.0%) |

WSPT removes candidates with different number of words in variable part

SCOLM reduces False Positives than Regex approach

```
demo@computer $ python
Python 3.10.12 (main, Nov  6 2024, 20:22:13) [GCC 11.4.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> |
```

# Demonstration

# Summary

- Design and develop a system to automatically identify original logging functions for helping network troubleshooting
  - Fast log matching combining regex tables and WSPT
- Evaluate processing time and number of obtained candidates using SCOLM with three datasets
  - Reduce 98.7% processing time than regex table approach
  - Reduce false positives of obtained results