# It Can Understand the Logs, Literally

**Aidi Pi**, Wei Chen, Will Zeller and Xiaobo Zhou

UCCS
University of Colorado
Colorado Springs

IPDPSW'19 @ Rio de Janeiro

# Outline

- Introduction to distributed system logs

- Challenges

- NLog: A NLP based log analysis approach

- Evaluation

- Conclusion

# Logging in general

- Logging is a general approach to record events in a system

- System logs are critical for understanding and troubleshooting targeted systems

# Challenges in log analysis

- Large number of log files

- Rich information in log messages

  - Identifiers, entities, events, etc.

- Effectiveness in information extraction

  - A single log message contains multiple fields

  - Multiple log messages can contain information about the same object

# A motivation example

**Task 39 force spilling in-memory map to
disk and it will release 159.6 MB memory**

- Existing approaches only extract identifiers and numeric values

- NLP approaches can extract events from logs

# Logs in natural languages

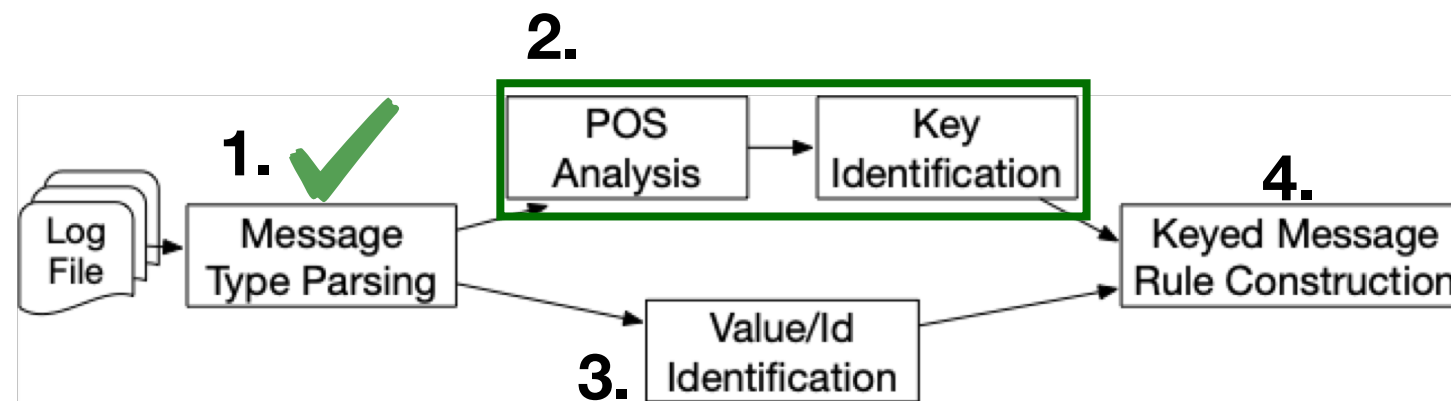| Frameworks | NL logs | Total logs | % of NL logs |
|:---:|:---:|:---:|:---:|
| Yarn | 84652 | 88628 | 99.5% |
| Spark | 106686 | 106686 | 100% |
| MapReduce | 85752 | 92648 | 92.6% |
| **Average** | **-** | **-** | **97.4%** |

- Our observation finds that most logs of data analytics frameworks are written in a *natural language*

# NLog

- NLog: a Natural Language Processing (NLP) based approach

- It can identify objects and events even without identifiers in logs

- Targeted systems: distributed data analytics frameworks

# NLog overview



1. Message type parsing: a solved problem by Spell*

2. Identification of key objects

3. Finding identifiers and numeric values

4. Storing parsing results in keyed messages**

* M. Du and F. Li,"Spell: Streaming parsing of system event logs" in *proc of ICDM'17*.

**A. Pi, W. Chen, X. Zhou, and M. Ji, "Profiling distributed systems in lightweight virtualized environments with logs and resource metrics" in *proc of HPDC'18*

# Step 1: message type parsing

- Message type: the static string sequence of in a corresponding log printing statement

```
fetcher 4 about to shuffle          fetcher * about to shuffle
output of map attempt_1      →      output of map *
decomp: 1965 len 1969 to MEMORY     decomp: * len * to MEMORY
```
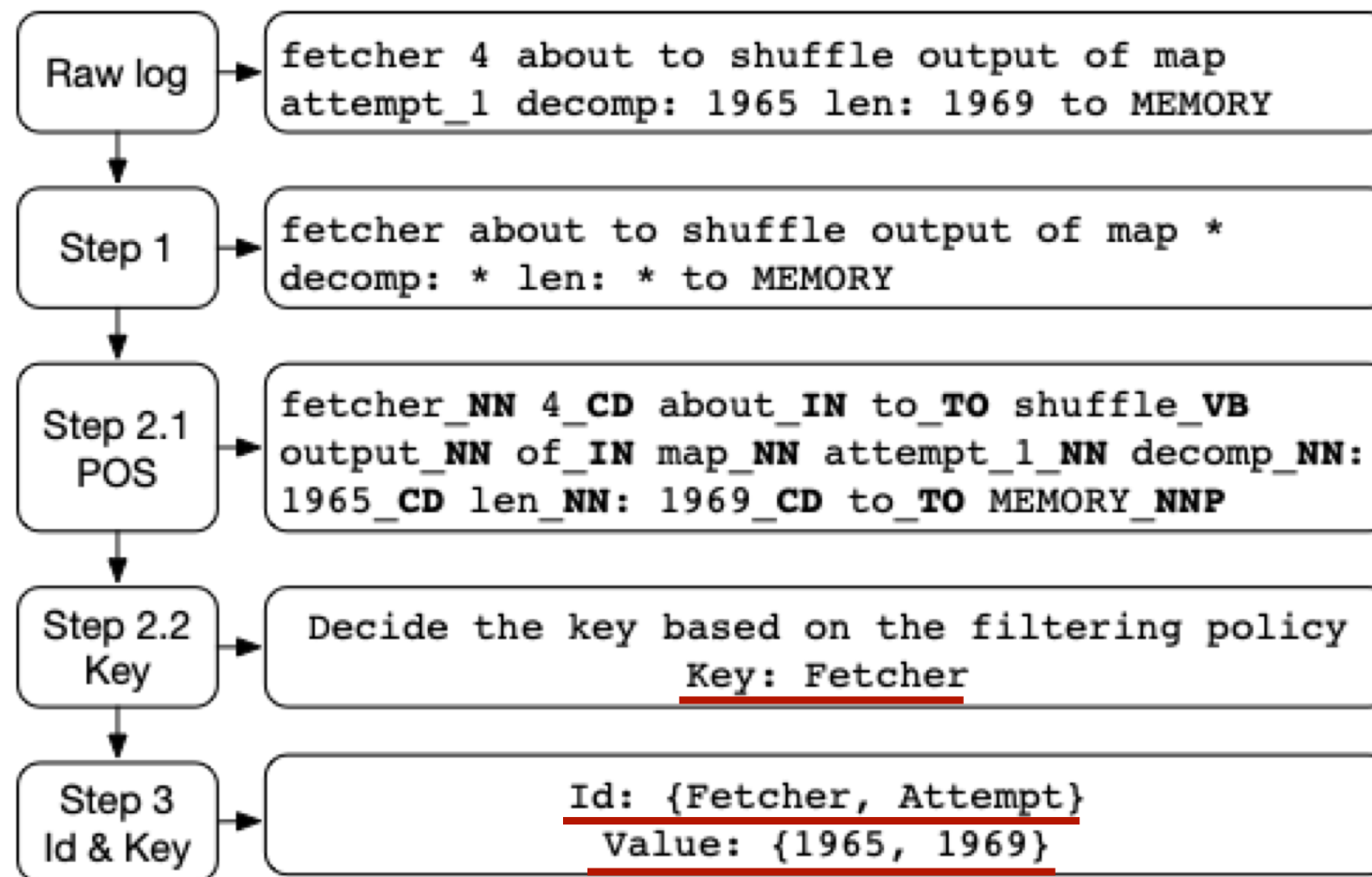
# Step 2: objects & event extraction by NLP

- Part-of-speech analysis: tag each word in a log message with its part-of-speech

- Find all the noun words

- Filter noun words with a top $\alpha$ frequency

  - Key object words have higher frequencies

- Assign key objects as keys of a log message

# Step 3: identifiers & values

- **Identifiers:** Numeric following a noun word

- **Values:** All other numeric value

  - Numeric values followed by units e.g. `kb` or `ms`

# An example : put it all together

| | |
|---|---|
| **Raw log** | `fetcher 4 about to shuffle output of map attempt_1 decomp: 1965 len: 1969 to MEMORY` |
| **Step 1** | `fetcher about to shuffle output of map * decomp: * len: * to MEMORY` |
| **Step 2.1 POS** | `fetcher_NN 4_CD about_IN to_TO shuffle_VB output_NN of_IN map_NN attempt_1_NN decomp_NN: 1965_CD len_NN: 1969_CD to_TO MEMORY_NNP` |
| **Step 2.2 Key** | `Decide the key based on the filtering policy` `Key: Fetcher` |
| **Step 3 Id & Key** | `Id: {Fetcher, Attempt}` `Value: {1965, 1969}` |

- The parsing results are in key-value format

- Users use queries on the results for troubleshooting purposes

# Evaluation setup

- Setup

  - Evaluation is conducted on a 25-node cluster

  - Four Xeon E5-2640 v3 CPU and 128GB memory per node

  - Cluster is connected by 10-Gbps Ethernet

  - Yarn-3.0.0-alpha, Spark-2.1.0

- Log files

  - Randomly choose 20 MB of of 2GB files

# Accuracy of object identification

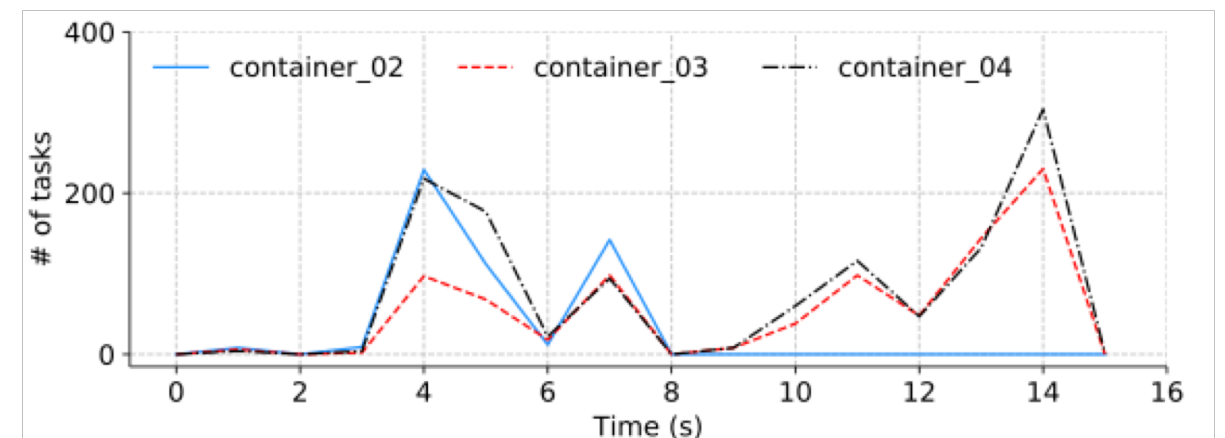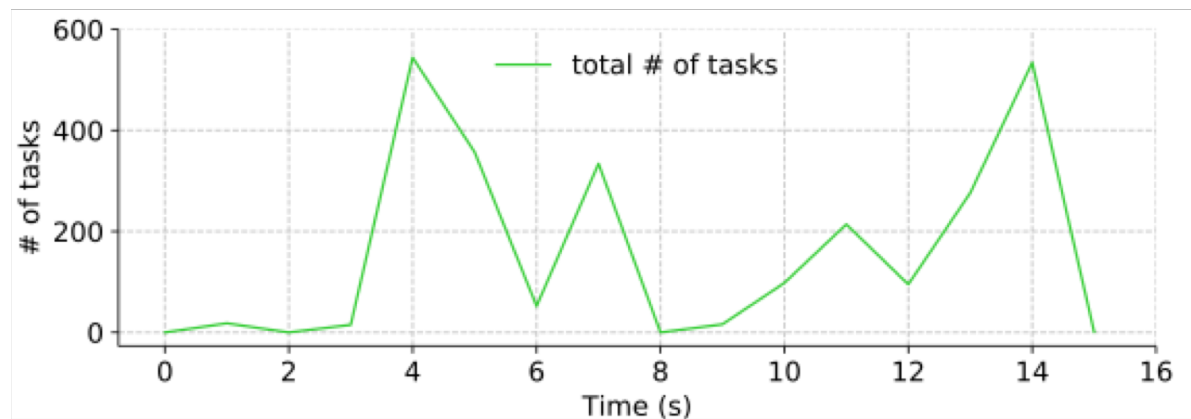| Frameworks | Total | Correct | Accuracy |
| --- | --- | --- | --- |
| Yarn | 115 | 99 | 85.3% |
| Spark | 34 | 32 | 94.1% |
| MapReduce | 92 | 86 | 93.5% |

- Inaccurate message types

  - All of its keys have too general meanings e.g. `service`

  - None of the keys includes the key objects

# A case study

**Spark TPC-H job**

**Inspect the number of tasks during job execution**

```
                    Related message types:
Q1: key: task       1. Got assigned task *
    aggregator: count  2. Running task * in stage * (TID *)
                    3. Task * force spilling in-memory map to
Q2: key: task       disk and it will release * memory
    aggregator: count  4. Finished task * in stage * (TID *)
    groupBy: container 5. Executor killed task * in stage * (TID *)
```





**Number of concurrently running tasks vary during job lifetime**

**Containers receive uneven number of tasks**

**The uneven task number distribution is caused by bug in Spark**

# Conclusion

- NLog, a NLP-based approach to identify key objects, identifiers and values in logs

- It is accurate in key object extraction

- It is helpful in understanding and troubleshooting targeted systems

# IntelLog

- IntelLog: a comprehensive NLP-based log analysis approach

- Objectives:

  - Information extraction

  - Automatic workflow reconstruction

  - Automatic problem detection

- IntelLog will be published in HPDC'19, Phoenix, AZ, USA

# Thank you!
## Q & A