

# LogParse: Making Log Parsing Adaptive through Word Classification

**Weibin Meng**, Ying Liu, Federico Zaiter, Shenglin Zhang, Yihao Chen, Yuzhe Zhang  
Yichen Zhu, En Wang, Ruizhi Zhang, Shimin Tao, Dian Yang, Rong Zhou, Dan Pei



清华大学  
Tsinghua University



南开大学  
Nankai University



UNIVERSITY OF  
TORONTO



吉林大学



HUAWEI

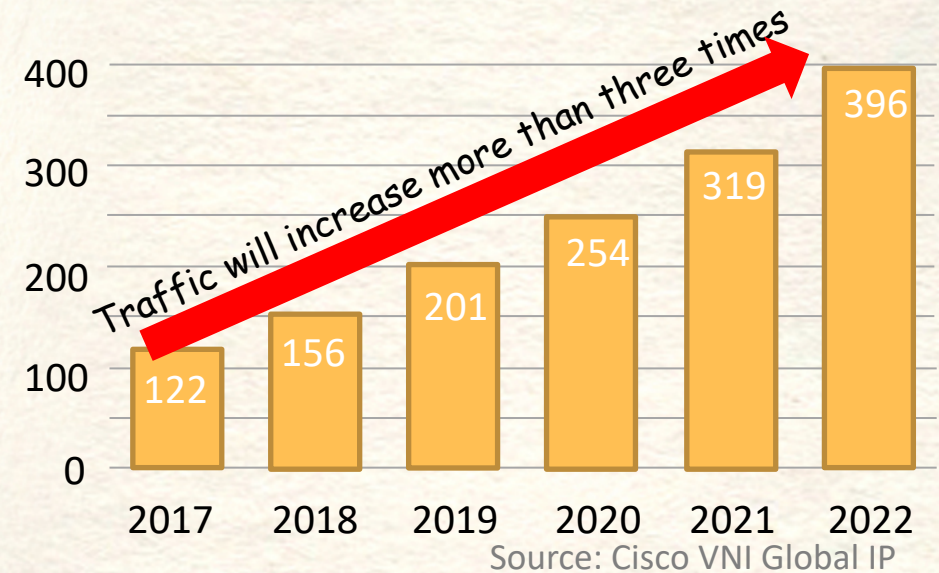






# Background

- Internet provide various types of services
- The traffic is growing rapidly.









# Background

- Stability of services are becoming more and more important.

Monitor services to keep stability

COMPANY	ESTIMATED ANNUAL	REVENUE LOSS PER	REVENUE LOSS PER
 HomeDepot.com			
 BestBuy.com	\$6,126,000,000.00	\$698,832.00	\$11,647.20
 Costco.com	\$6,108,500,000.00	\$696,852.00	\$11,614.20

### Delta Says Computer Breakdown

August, in part because of the outage and subsequent recovery efforts, the carrier said in a statement Friday. The breakdown reduced unit revenue, as the measure is also known, by two percentage points, Delta said.

The country's second-largest airline earlier forecast that third-quarter unit revenue would fall 4 percent to 6 percent.

A [power-control module](#) at Delta's Atlanta computer center failed and caught fire Aug. 8, shutting down electricity to the system. About 300 of the airline's 7,000 servers weren't wired to backup power, the company had said.

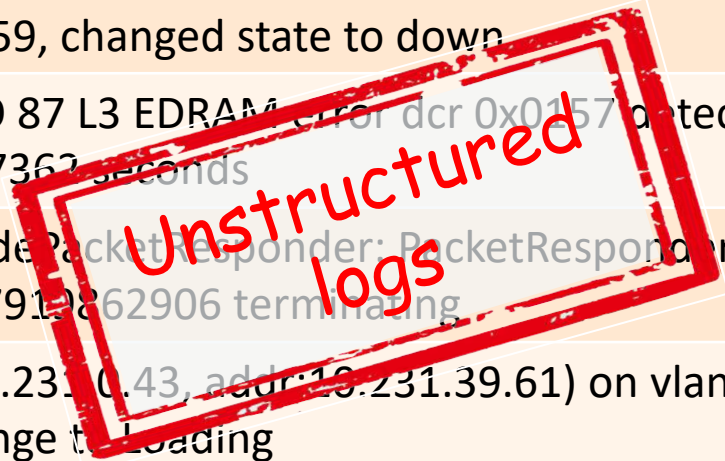




# Service logs

- Logs are the most valuable data for service management
  - Logs record a vast range of events (7\*24) of services
  - Every service generates logs

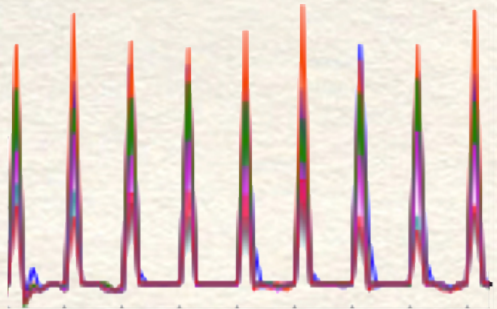
Types	Timestamps	Detailed messages
Switch	Jul 10 19:03:03	Interface te-1/1/59, changed state to down
Supercomputer	Jun 4 6:45:50	RAS KERNEL INFO 87 L3 EDRAM error dcr 0x0157 detected and corrected over 27362 seconds
HDFS	Jun 8 13:42:26	INFO dfs.DataNodePacketResponder: PacketResponder 1 for block blk_-1608999687912862906 terminating
Router	Jul 11 11:05:07	Neighbour(rid:10.231.0.43, addr:10.231.39.61) on vlan23, changed state from Exchange to Loading



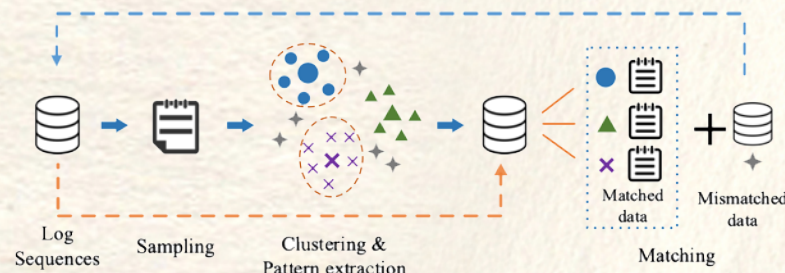


# Log parsing

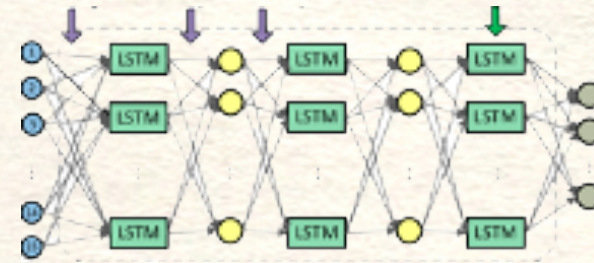
- Log analysis → Log-based service management
- Log analysis contains two steps<sup>[1]</sup>:
  - Log Parsing and Log Mining
- Log parsing effects the performance of log analysis



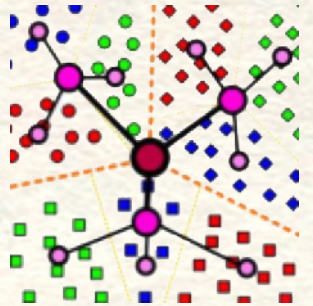
Performance monitoring  
[INFOCOM'19]



Problem Identifying  
[FSE'18]



Anomaly detection  
[IJCAI'19]



Failure prediction  
[SIGMETRICS'18]

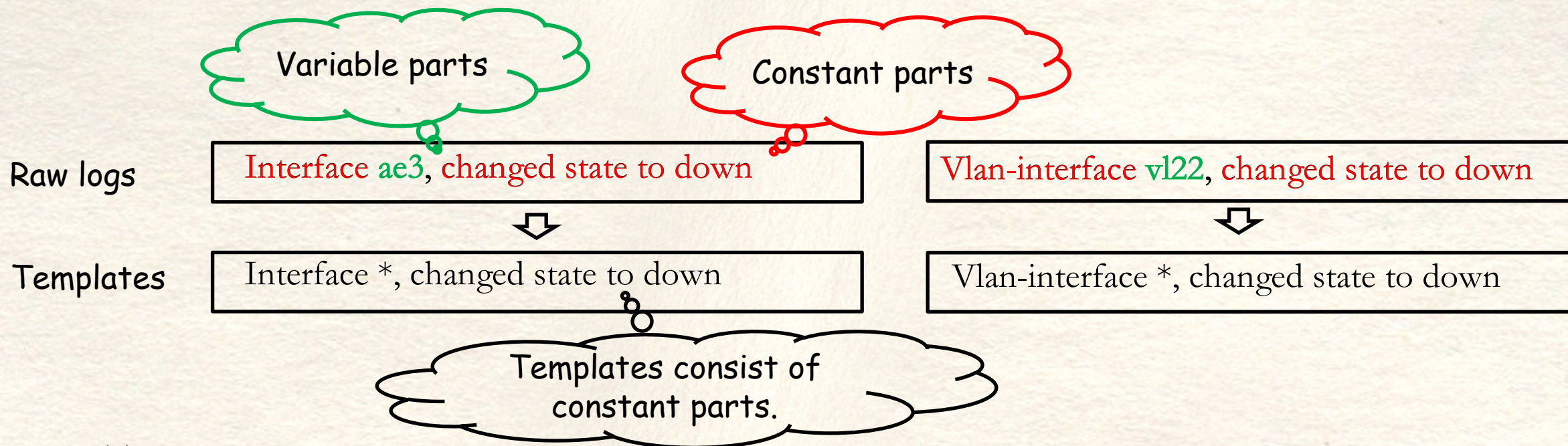
[1] Pinjia He, Jieming Zhu, et al. An Evaluation Study on Log Parsing and Its Use in Log Mining. DSN'16





# Log parsing

- An unstructured log is “printf”ed by services
- The goal of log parsing is to distinguish between constant part and variable part.





# Adaptiveness

---

- Adaptiveness is important for log parsing

- Goal: match any types of logs

- Intra-service adaptiveness ❌

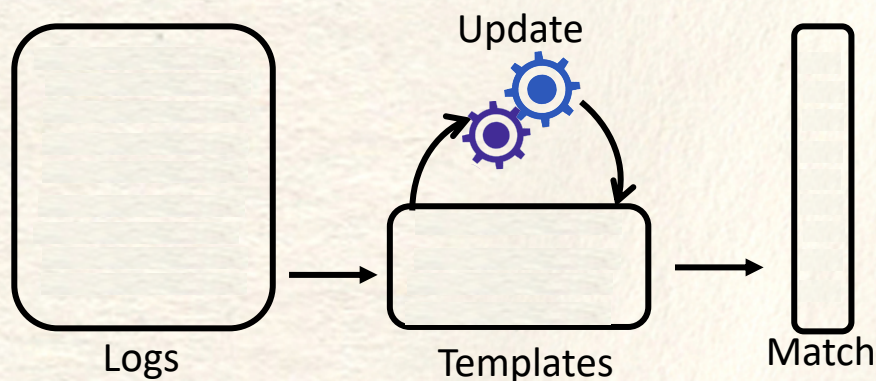
- Cross-service adaptiveness ❌

Traditional log parsing approaches or don't support intra-service adaptiveness, or do not support cross-service adaptiveness, or both.



# Intra-service Adaptiveness

- Intra-service adaptiveness
  - Software/firmware upgrades can generate new types of logs
  - New logs cannot match any existing templates



## Historical logs:

L<sub>1</sub>. Interface ae3, changed state to down  
L<sub>2</sub>. Vlan-interface vl22, changed state to down  
L<sub>3</sub>. Interface ae3, changed state to up  
L<sub>4</sub>. Interface ae1, changed state to down

## Real-time logs:

L<sub>5</sub>. Interface ae1, changed state to up  
L<sub>6</sub>. Vlan-interface vl22, changed state to up



## Template extraction:

T<sub>1</sub>. Interface \*, changed state to down  
T<sub>2</sub>. Vlan-interface \*, changed state to down  
T<sub>3</sub>. Interface \*, changed state to up

## Template update:

T<sub>4</sub>. ? ? ?

## Template match:

L<sub>1</sub>->T<sub>1</sub>,ae3   L<sub>2</sub>->T<sub>2</sub>,vl22   L<sub>3</sub>->T<sub>3</sub>,ae3  
L<sub>4</sub>->T<sub>1</sub>,ae1   L<sub>5</sub>->T<sub>3</sub>,ae1   L<sub>6</sub>-> ? ? ?



# Intra-service Adaptiveness

## Traditional log parsing methods:

- Drain (ICWS'17), FT-tree (IWQoS'18) which claimed to support template update
- LogSig (CIKM'11), Spell (ICDM'16), IPLoM (KDD'09) don't support template update

### Historical logs:

L<sub>1</sub>. Interface ae3, changed state to down  
L<sub>2</sub>. Vlan-interface vl22, changed state to down  
L<sub>3</sub>. Interface ae3, changed state to up  
L<sub>4</sub>. Interface ae1, changed state to down

### Real-time logs:

L<sub>5</sub>. Interface ae1, changed state to up  
L<sub>6</sub>. Vlan-interface vl22, changed state to up

### Template extraction:

T<sub>1</sub>. Interface \*, change  
T<sub>2</sub>. Vlan-interface \*, change  
T<sub>3</sub>. Interface \*, change

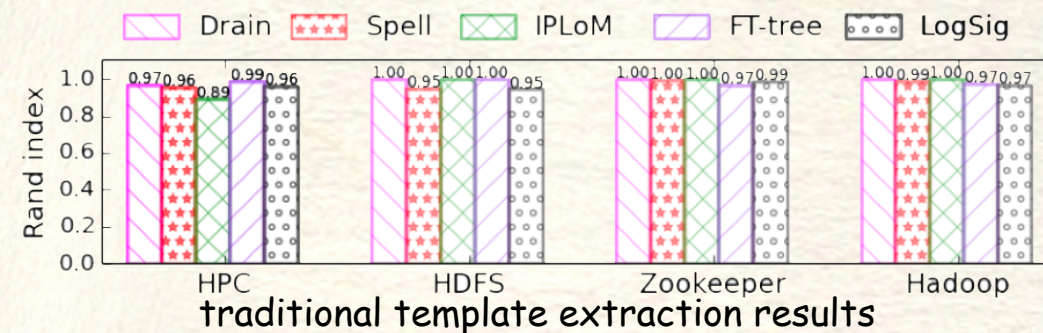
### Template update:

T<sub>4</sub>. ? ? ?

### Template match:

L<sub>1</sub>->T<sub>1</sub>, ae3    L<sub>2</sub>->T<sub>2</sub>, vl22    L<sub>3</sub>->T<sub>3</sub>, ae3  
L<sub>4</sub>->T<sub>1</sub>, ae1    L<sub>5</sub>->T<sub>3</sub>, ae1    L<sub>6</sub>-> ? ? ?

When face new types of logs,  
all traditional methods  
achieve low accuracies



traditional template extraction results  
when all the logs are used for training



traditional template extraction results  
when only 10% of logs are used for training



# Cross-service Adaptiveness

## ■ Observation:

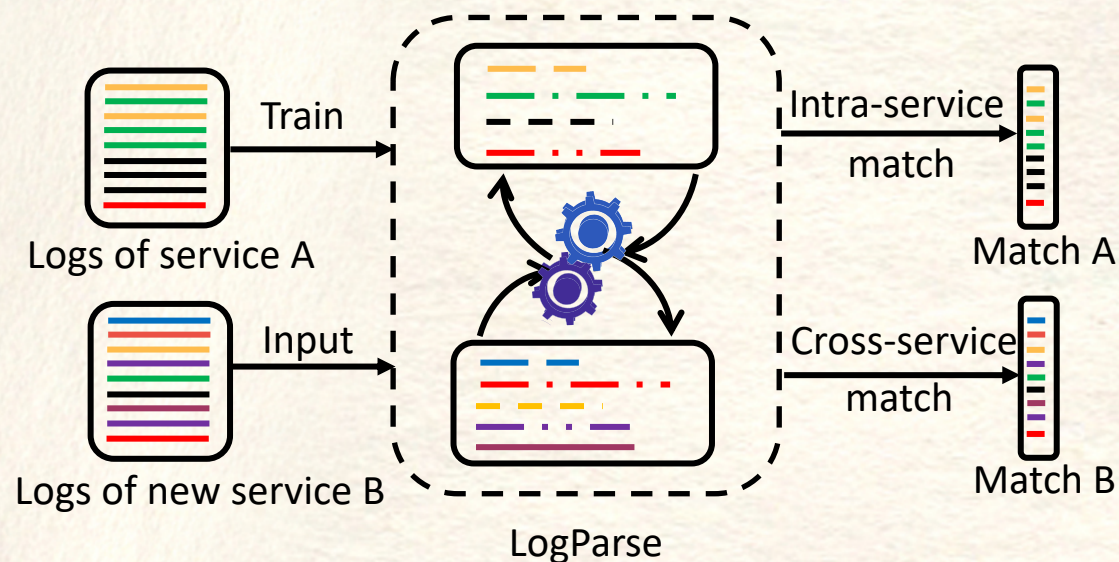
- No enough historical logs when a brand new service goes online

## ■ Aim:

- **A model** trained by service A is also suitable for service B



- Cross-service adaptive is for models rather than template sets.
- Templates are generated by trained model.







# Log compression

- Log parsing **without adaptiveness limit many log analysis** applications

- Requires a corresponding template for any given log

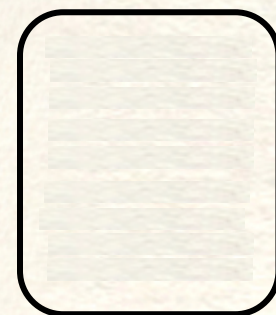
- Log Compression (An application)

- Short-term storage (real-time queries)

- Template index + Variables

- Long-term storage

- LogParse + Traditional method (e.g., zip)



Historical logs:

L<sub>1</sub>. Interface ae3, changed state to down  
L<sub>2</sub>. Vlan-interface vl22, changed state to down  
L<sub>3</sub>. Interface ae3, changed state to up  
L<sub>4</sub>. Interface ae1, changed state to down

Real-time logs:

L<sub>5</sub>. Interface ae1, changed state to up  
L<sub>6</sub>. Vlan-interface vl22, changed state to up

Short term storage by LogParse:

L<sub>1</sub>->T<sub>1</sub>,ae3  
L<sub>2</sub>->T<sub>2</sub>,vl22  
L<sub>3</sub>->T<sub>3</sub>,ae3  
L<sub>4</sub>->T<sub>1</sub>,ae1  
L<sub>5</sub>->T<sub>3</sub>,ae1  
L<sub>6</sub>->T<sub>4</sub>,vl22

Long-term compression







# Idea

## ■ Observation:

- Operators usually distinguish variables based on features of words

Mixed characters  
and numbers are  
usually variables

### Historical logs:

- L<sub>1</sub>. Interface ae3, changed state to down
- L<sub>2</sub>. Vlan-interface vl22, changed state to down
- L<sub>3</sub>. Interface ae3, changed state to up
- L<sub>4</sub>. Interface ae1, changed state to down

### Real-time logs:

- I<sub>5</sub>. Interface ae1, changed state to up
- L<sub>6</sub>. Vlan-interface vl22, changed state to up

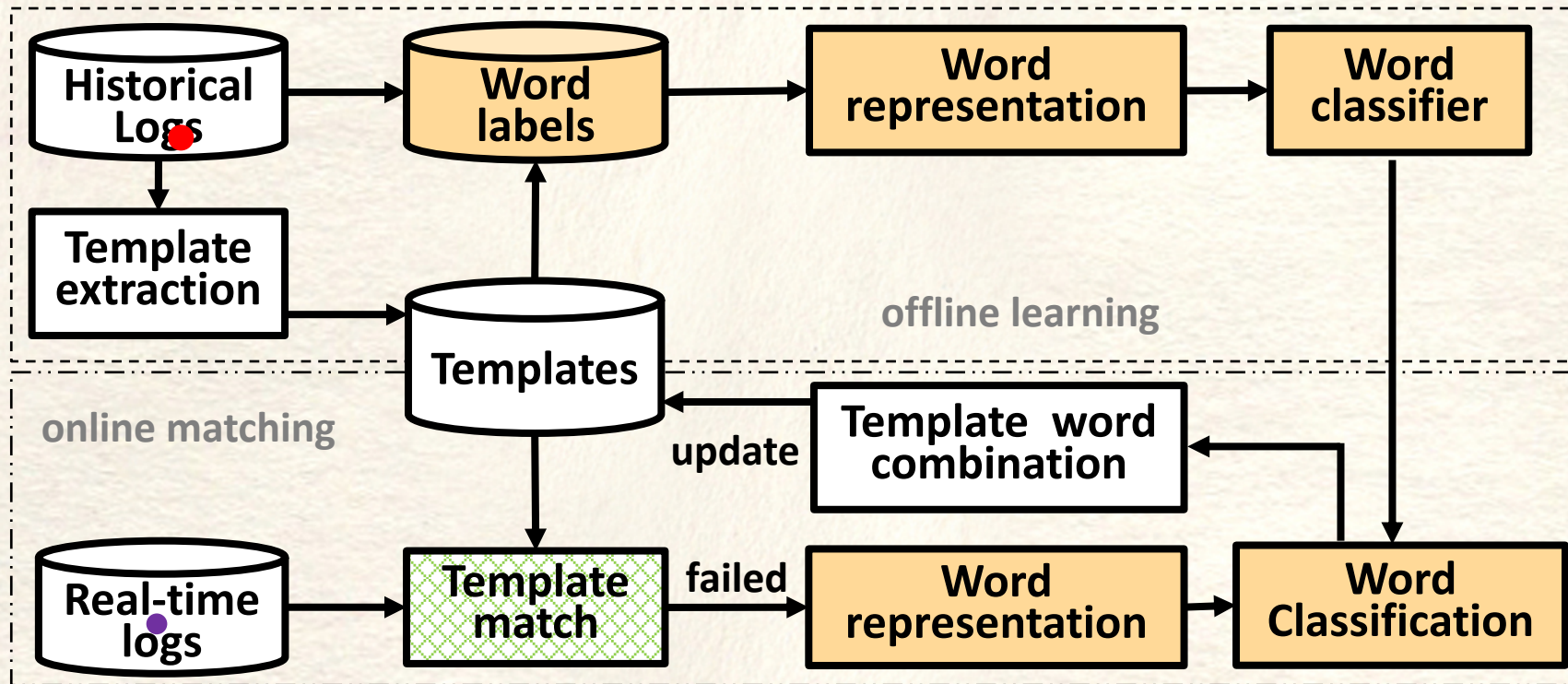
letters are  
usually template  
words

**A log parsing problem → A word classification problem**





# LogParse Workflow

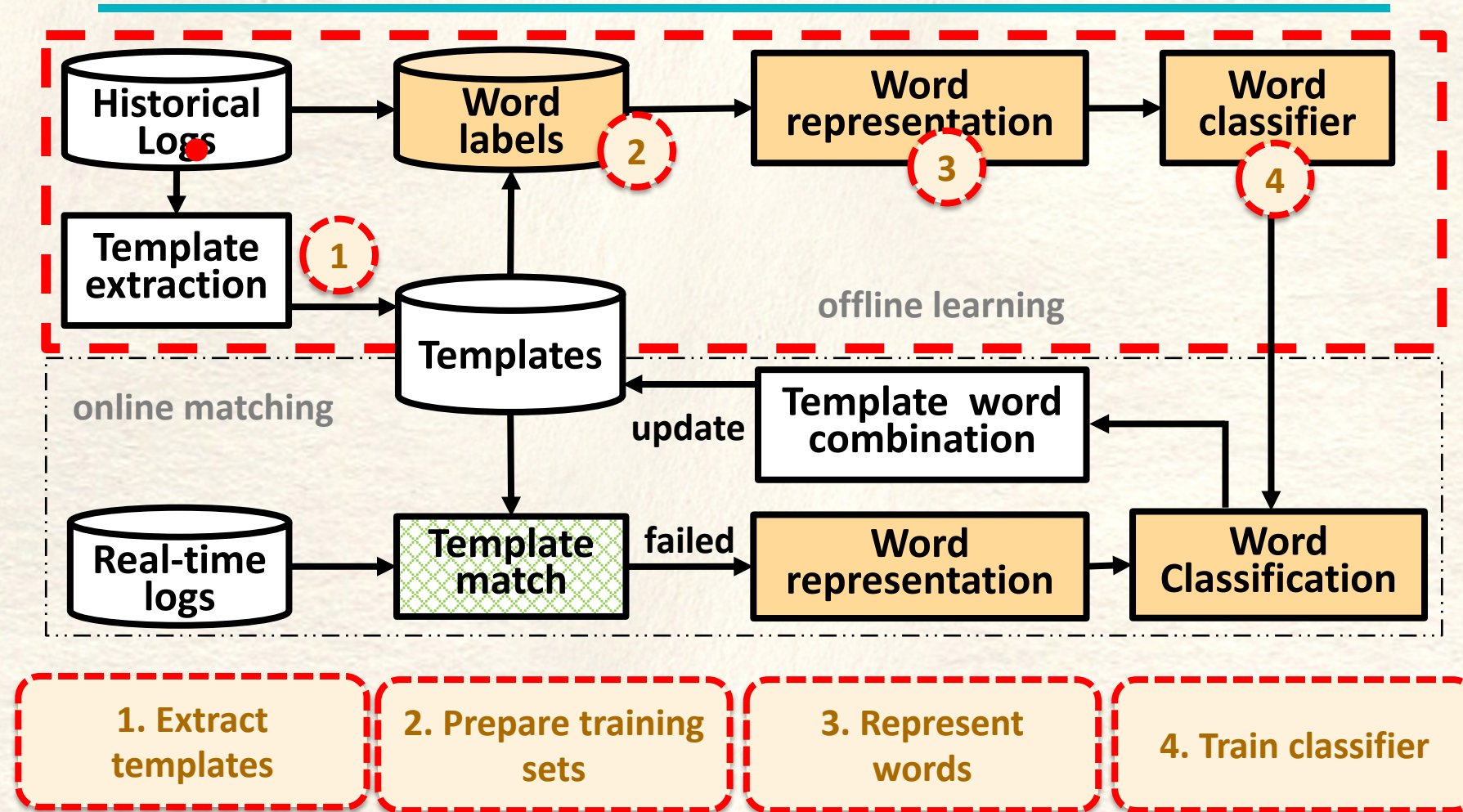


- Offline learning:
  - Prepare training word sets and train word classifier
- Online log parsing:
  - Match logs and update template sets

An adaptive framework for online log parsing



# Offline Learning







# Template extraction

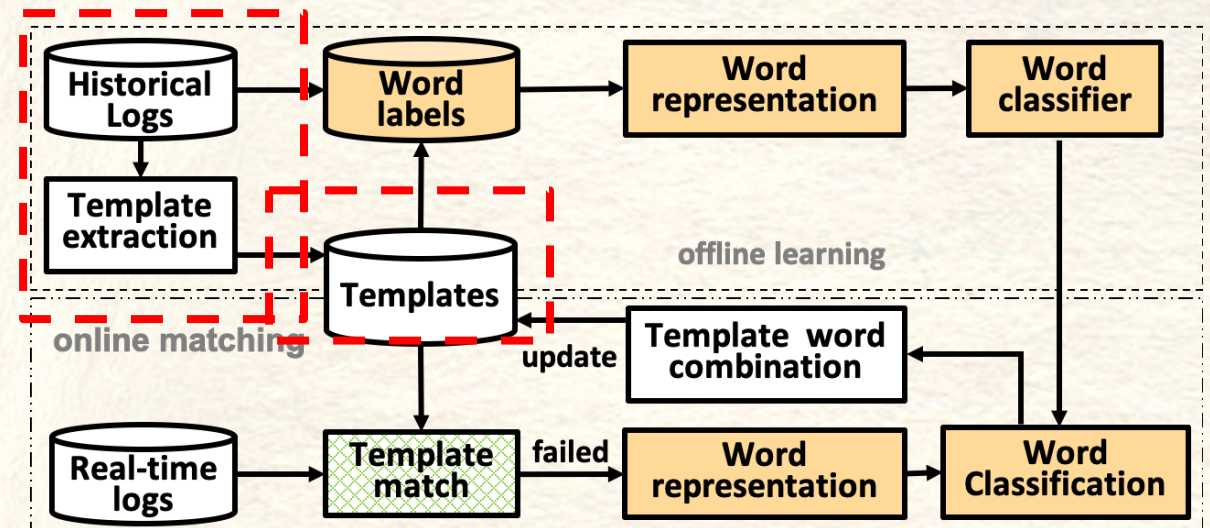
- Extract templates by traditional log parsing methods
  - Generate accurate templates (in offline stage)
  - Unsupervised methods
- Use the results as the initial template set

## Rawlogs:

L<sub>1</sub>. Interface ae3, changed state to down  
L<sub>2</sub>. Vlan-interface vlan22, changed state to down  
L<sub>3</sub>. Interface ae3, changed state to up.  
L<sub>4</sub>. Interface ae1, changed state to down

## Templates:

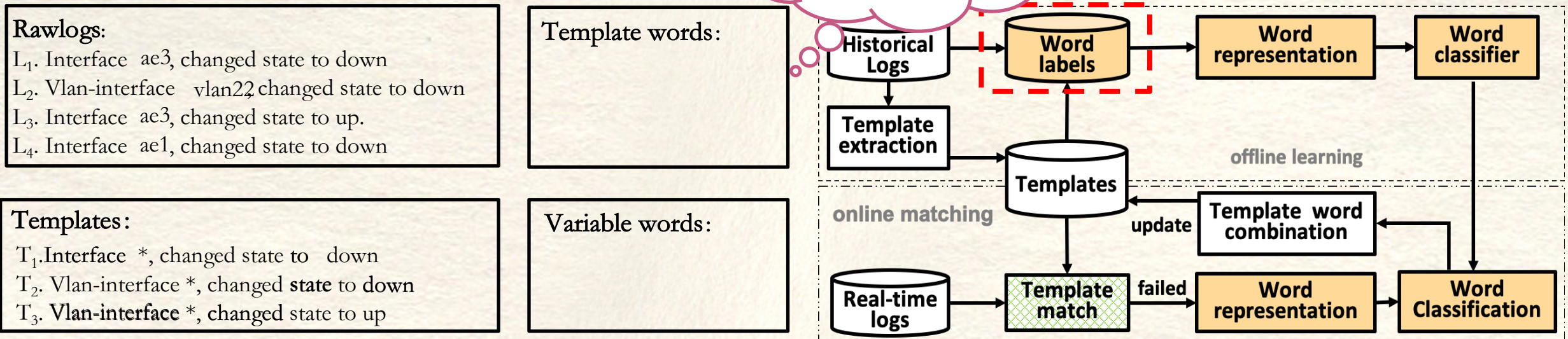
T<sub>1</sub>. Interface \*, changed state to down  
T<sub>2</sub>. Vlan-interface \*, changed state to down  
T<sub>3</sub>. Vlan-interface \*, changed state to up





# Prepare training sets

- Distinguish variable/template words
  - variable words: words in logs but not in templates
  - template words: words in templates

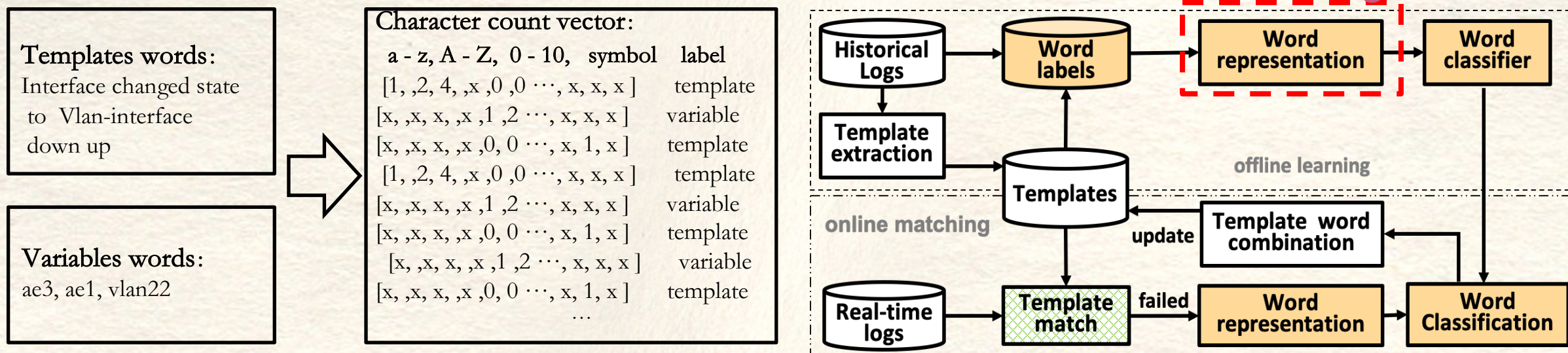




# Word representation

- Machine learning algorithms require **structured data**
- Present each word by using a character-level count vector
  - The set of characters is fixed -> fixed dimensionality
    - e.g., 128 characters in ASCII

We can represent any word even unseen words





# Word classifier

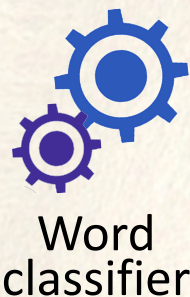
- Train supervised machine learning classifier
  - E.g., SVM, Random forest.
- The whole framework of LogParse is unsupervised
  - We used unsupervised methods to generated training set

The whole framework is still unsupervised

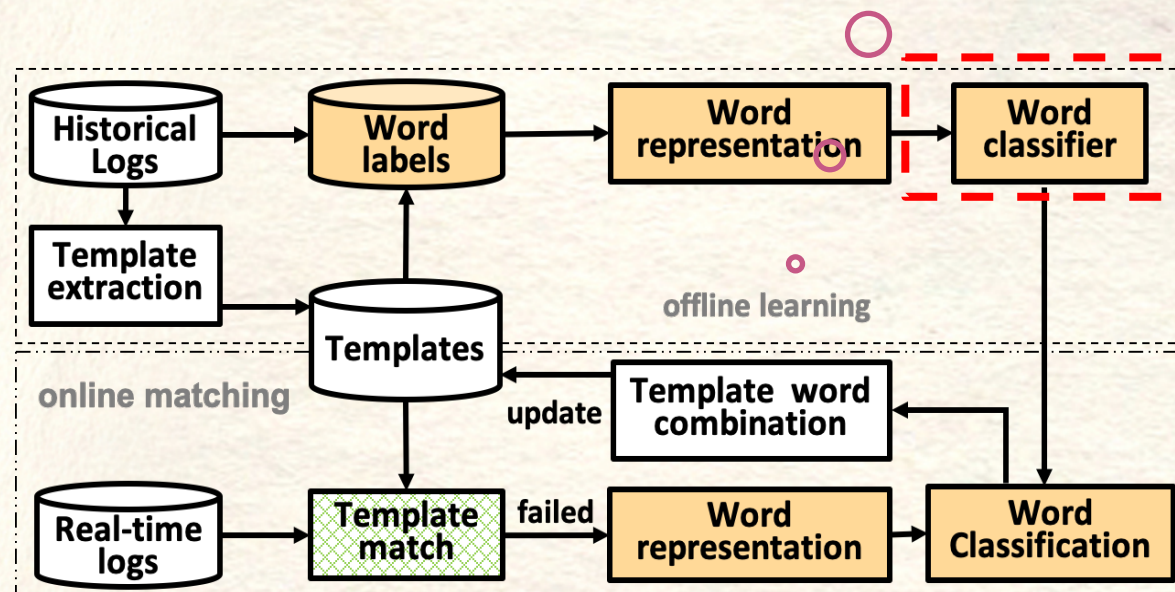
Character count vector:

a - z, A - Z, 0 - 10, symbol	label
[1, 2, 4, x, 0, 0 ..., x, x, x]	template
[x, x, x, x, 1, 2 ..., x, x, x]	variable
[x, x, x, x, 0, 0 ..., x, 1, x]	template
[1, 2, 4, x, 0, 0 ..., x, x, x]	template
[x, x, x, x, 1, 2 ..., x, x, x]	variable
[x, x, x, x, 0, 0 ..., x, 1, x]	template
[x, x, x, x, 1, 2 ..., x, x, x]	variable
[x, x, x, x, 0, 0 ..., x, 1, x]	template
...	

Machine learning



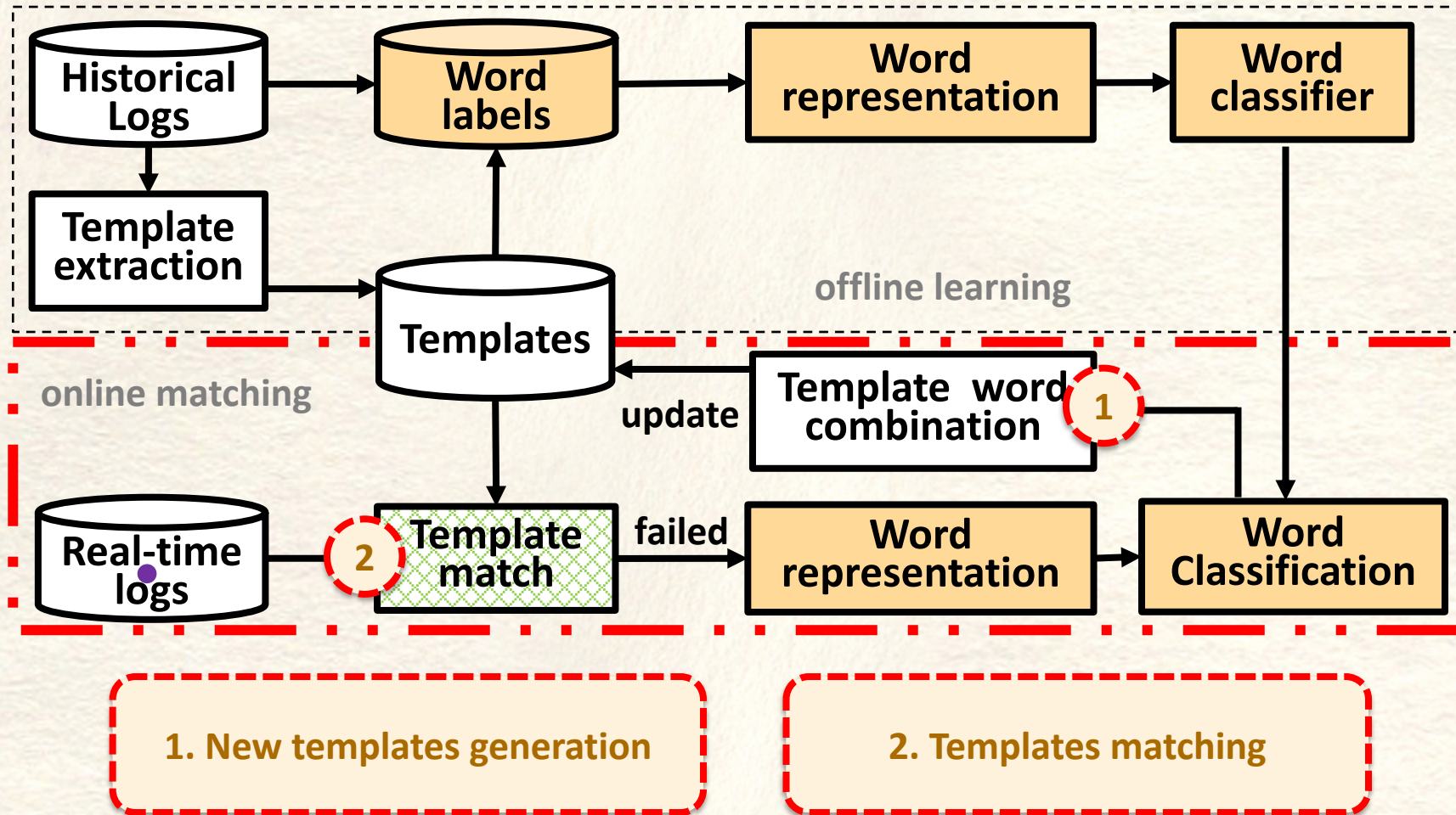
Word classifier







# Online log parsing





# New templates generation

## Steps:

- Classify each word by the trained word classifier.
- Construct a new template by combining all template words

### Historical logs:

L<sub>1</sub>. Interface ae3, changed state to down  
L<sub>2</sub>. Vlan-interface vl22, changed state to down  
L<sub>3</sub>. Interface ae3, changed state to up  
L<sub>4</sub>. Interface ae1, changed state to down

### Online logs (new type):

L<sub>5</sub>. Vlan-interface vl22, changed state to up

Template words

variable words

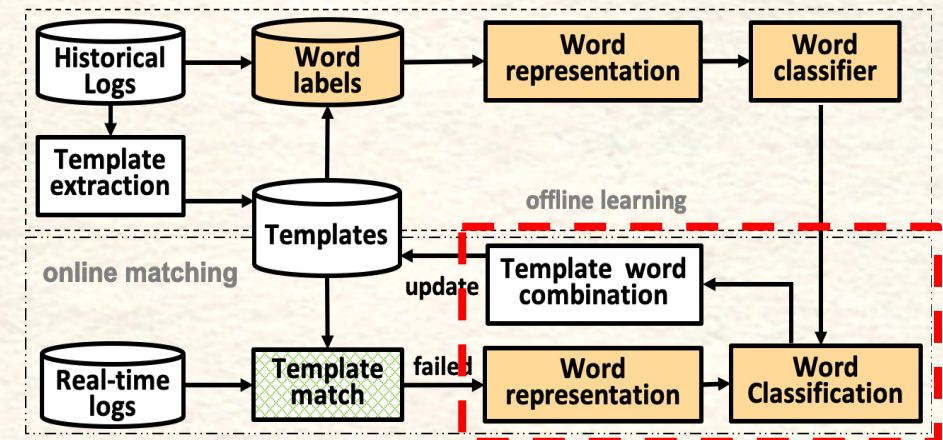
### Template set:

T<sub>1</sub>. Interface \*, changed state to down  
T<sub>2</sub>. Vlan-interface \*, changed state to down  
T<sub>3</sub>. Interface \*, changed state to up  
T<sub>4</sub>. \*,

### Template match:

L<sub>1</sub>->T<sub>1</sub>, ae3    L<sub>2</sub>->T<sub>2</sub>, vl22  
L<sub>3</sub>->T<sub>3</sub>, ae3    L<sub>4</sub>->T<sub>1</sub>, ae1

log parsing problem → word classification problem







# Template matching

- Build a prefix-tree for template matching
- Each root-to-leaf path is a template

Templates:

$T_1$ . Interface \*, changed state to down

$T_2$ . Vlan-interface \*, changed state to down

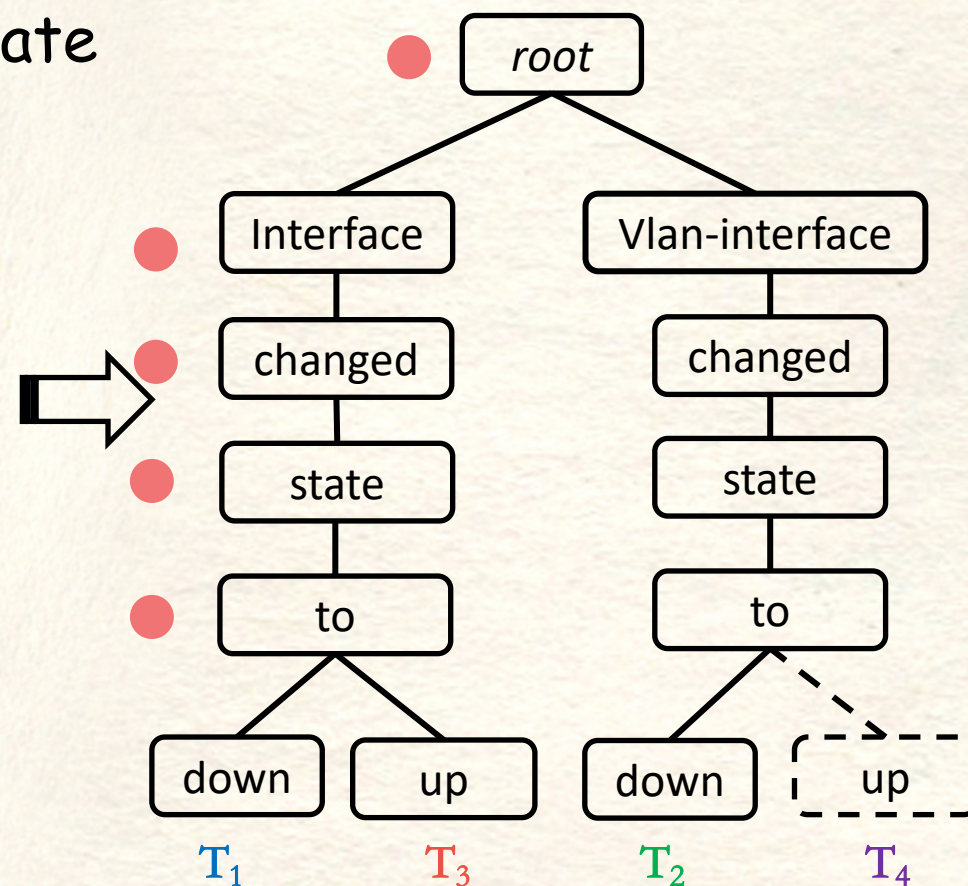
$T_3$ . Interface \*, changed state to up

Template update:

$T_4$ . Vlan-interface \*, changed state to up

Online log:

Interface ae3, changed state to down





# Datasets & Baselines

---

## ■ Datasets:

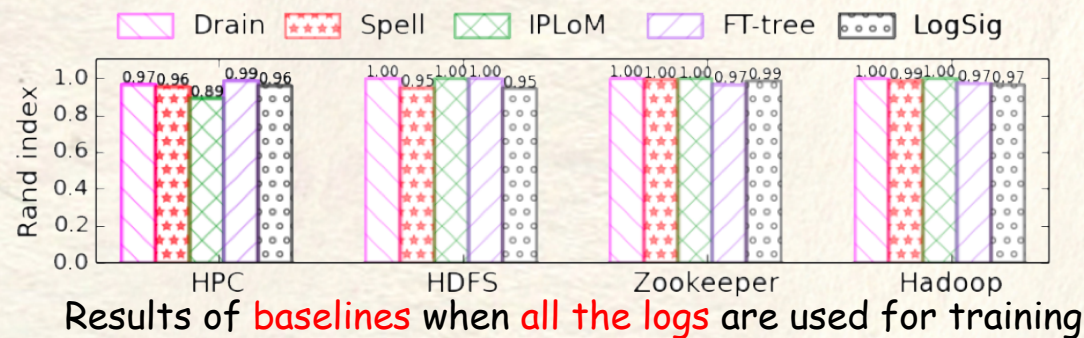
Datasets	Description	# of logs
HPC	High performance cluster	433,489
HDFS	Hadoop distributed file system	11,175,629
ZooKeeper	ZooKeeper service	74,380
Hadoop	Hadoop MapReduce job	394,308

## ■ Baselines:

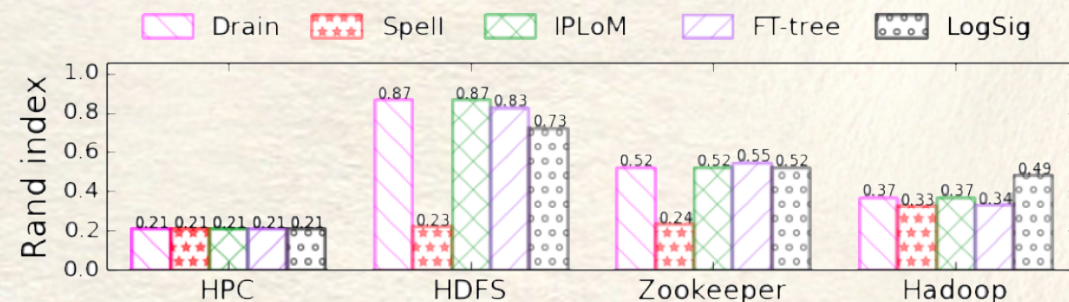
- Drain (ICWS'17), FT-tree (IWQoS'18) which claimed to support template update
- LogSig (CIKM'11), Spell (ICDM'16), IPLoM (KDD'09) don't support template update



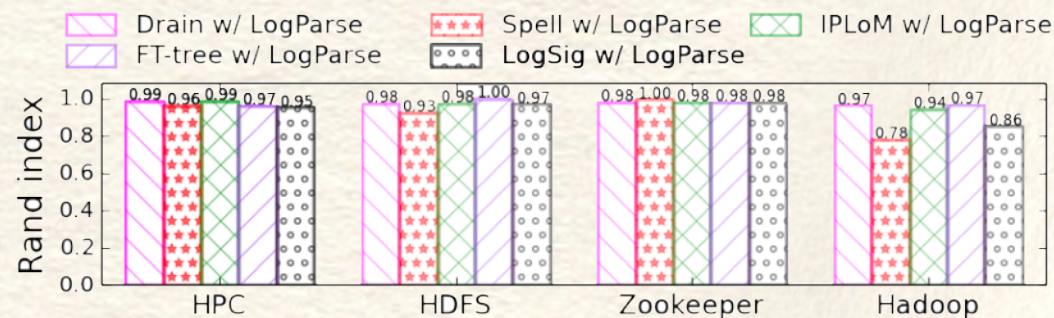
# Evaluation on Intra-service adaptiveness



All baselines perform good in offline stage



All baselines perform bad for online matching and update

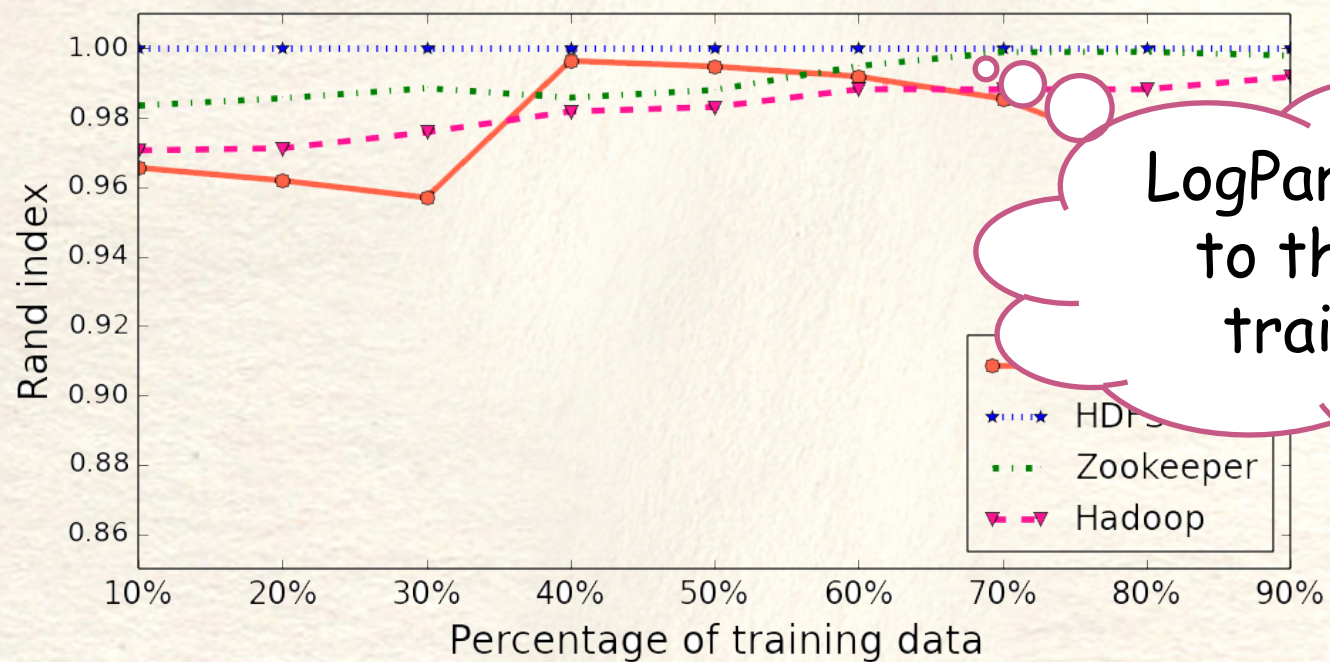


Accuracy of LogParse is even higher than baselines trained by all logs.





# Evaluation on stability



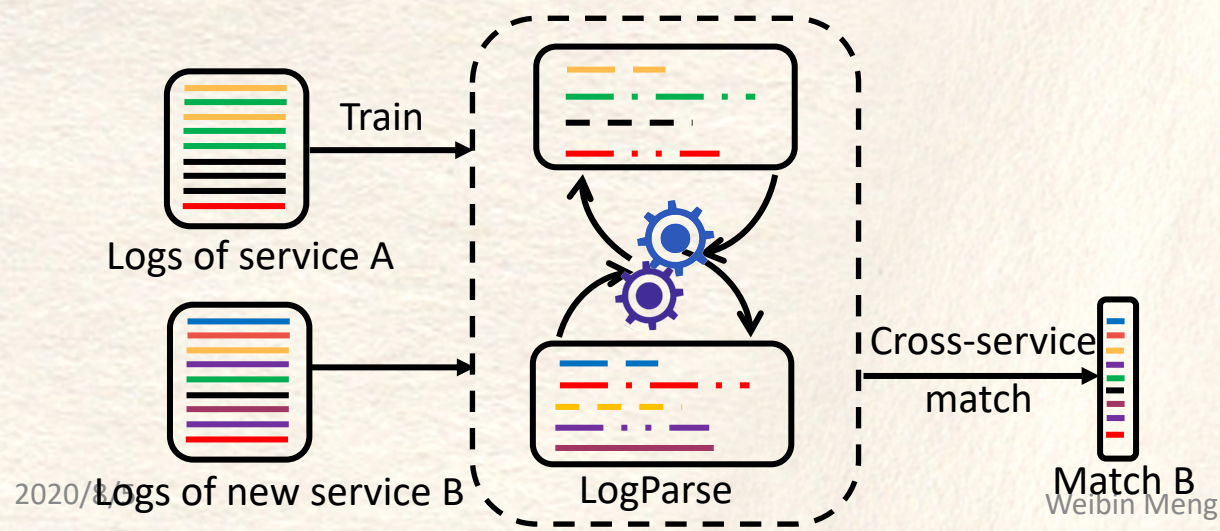
LogParse is stable  
to the scale of  
training data

template accuracy as the percentage of training  
data increases from 10% to 90%



## ➤ Evaluation on cross-service adaptive

Training data (service A)	Testing data (service B)			
	HPC	HDFS	Zookeeper	Hadoop
Trained by HPC	-	0.983	0.999	0.923
HDFS	0.982	-	0.993	0.974
ZooKeeper	0.993	1.0	-	0.937
Hadoop	0.983	0.999	0.999	-



On average,  
LogParse achieves  
a cross-service  
accuracy of 0.980



## Evaluation on compression

Type	Method	HPC	HDFS	Zookeeper	Hadoop	Average
Shore-term storage	LogParse	13.0%	14.0%	19.6%	4.6%	12.8
	bzip	9.6%	17.4%	9.7%	6.4%	10.8%
	7zip	9.7%	18.1%	9.4%	5.9%	10.8%
	zip	11.4%	20.9%	10.1%	7.2%	12.4%
Long-term storage	LogParse+bzip	1.4%	2.1%	2.4%	1.2%	1.8%
	LogParse+7zip	2.3%	2.6%	2.1%	1.5%	2.1%
	LogParse+zip	2.2%	2.6%	2.1%	1.5%	2.1%

LogParse is helpful  
to log compression





# Conclusion

---

LogParse, an adaptive log parsing method

- Intra-service
- Cross-service

Log compression, an application of LogParse

- Assign template for any given log

An open-source toolkit



# THANKS

## Q&A

mwb16@mails.tsinghua.edu.cn

Toolkit: <https://github.com/WeibinMeng/LogParse>

