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



Weibin Meng



Internet provide various types of servicesThe traffic is growing rapidly.





Stability of services are becoming more and more important.

					Dolto Computor Drookdown	
		Moni	tor s	ervi	ces 1	to keep stability
	2	HomeDep ot.com		evenue	LOSS	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.
	BEST BUY <mark>-</mark>	BestBuy.com	\$6,126,000,000.00	\$698,832.00	\$11,647.20	Merril Lynch The country's second-largest airline earlier forecast that third-quarter unit revenue would fall 4 percent to 6 percent.
	Costco	Costco.com	\$6,108,500,000.00	\$696,852.00	\$11,614.20	A <u>power-control module</u> 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.

Weibin Meng



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 CHOP dcr 0x0157 detected and corrected over 27362 seconds
HDFS	Jun 8 13:42:26	INFO dfs.DataNodePacketheSonder: PacketResponder 1 for block blk1608999687913862906 terminating
Router	Jul 11 11:05:07	Neighbour(rid:10.231,0.43, addr:13.231.39.61) on vlan23, changed state from Exchange to Loading



- ■Log analysis \rightarrow Log-based service management
- Log analysis contains two steps^[1]:
 - Log Parsing and Log Mining
- Log parsing effects the performance of log analysis



[1] Pinjia He, Jieming Zhu, et al. An Evaluation Study on Log Parsing and Its Use in Log Mining. DSN'16 2021/1/27



An unstructured log is "printf"ed by services

The goal of log parsing is to distinguish between

constant part and variable part.





Adaptiveness is important for log parsing
 Goal: match any types of logs

Intra-service adaptiveness

Cross-service adaptiveness

Traditional log paring 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



8

>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





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.



2021/1/27

Log compression





Observation:

Operators usually distinguish variables based on features of words

Mixed characters and numbers are usually variables

5	Historical logs:		-
	L_1 . Interface <u>ae3</u> , changed state to down		
	L_2 . Vlan-interface <u>vl22</u> , changed state to dows	α	7
	L_3 . Interface <u>ae3</u> , changed state to up		6
)	L ₄ . Interface <u>ae1</u> , changed state to down		-
	Real-time logs:	2.145	
	L. Interfacoae1, changed state to up		
	L_6 . Vlan-interface vl22, changed state to up		

letters are usually template words

A log parsing problem \rightarrow A word classification problem

LogParse Workflow



•Offline learning:

Prepare training word sets and train word classifier

Online log parsing:

2021/1/27

Match logs and update template sets

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

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:

L₁. Interface ae3, changed state to down

 L_2 . Vlan-interface vlan22, changed state to down

L₃. Interface ae3, changed state to up.

L₄. Interface ae1, changed state to down

Templates:

T₁. Interface *, changed state to down

T₂. Vlan-interface *, changed state to down

T₃. 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

W	ec	can	
repre	se	nt a	ny
wor	rd e	ever	1
unsee	en	wor	ds

		Character count vector:		Word Word	
Templates words:		a - z, A - Z, 0 - 10, symbol	label	Historical Word representation classifier	
Interface changed state		$[1, 2, 4, x, 0, 0, \dots, x, x, x]$	template		
to Vlan-interface		$[x, x, x, x, x, 1, 2, \cdots, x, x, x]$	variable	Tamplata	
down up		$[x, x, x, x, x, 0, 0 \cdots, x, 1, x]$	template	extraction	1
Server and the server and	$ \rangle$	$[1, 2, 4, x, 0, 0, \dots, x, x, x]$	template	offline learning	
		$[x, x, x, x, x, 1, 2, \cdots, x, x, x]$	variable	Tomplate word	1
		$[x, x, x, x, x, 0, 0 \cdots, x, 1, x]$	template	online matching update combination	1
Variables words:		[x, ,x, x, ,x ,1 ,2 ···, x, x, x]	variable		
ae3, ae1, vlan22	-nret c	$[x, x, x, x, x, 0, 0 \cdots, x, 1, x]$	template	Template failed Word Word	1
			S.S.	Real-time match representation Classification	Ì

> Word classifier

Train <u>supervised</u> machine learning classifier

E.g., SVM, Random forest.

The whole framework of LogParse is unsupervised

The whole framework is still unsupervised

We used unsupervised methods to generated training set



> Online log parsing





Steps:

Classify each word by the trained word classifier.

Construct a new template by combining all template words



log parsing problem \rightarrow word classification problem



Weibin Meng

> Template matching

Build a prefix-tree for template matching

Each root-to-leaf path is a template



root



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.





data increases from 10% to 90%

> Evaluation on cross-service adaptive

	Training data	Testing data (service B)					
	(service A)	HPC	HDFS	Zookeeper	Hadoop		
Traine	d by HPC —	_	0.983 Match to	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	-		
2	Logs of service A	LogParse	Cross-service match Match B Weibin Meng	On avera ogParse ac a cross-set accuracy of	ige, hieves rvice 0.980		

25

> Evaluation on compression

Туре	Method	НРС	HDFS	Zookeeper	Hadoop	Average
	LogParse	13.0%	14.0%	19.6%	4.6%	12.8
Shore-term	bzip	9.6%	17.4%	9.7%	6.4%	10.8%
storage	7zip	9.7%	18.1%	9.4%	5.9%	10.8%
	zip	11.4%	20.9%	10.1%	7.2%	12.4%
Long-term	LogParse+bzip	1.4%	2.1%	2.4%	1.2%	1.8%
storage	LogParse+7zip	2.3%	2.6%	· .0/_		d'a
	LogParse+zip	2.2%	2.6%	LogPars	se is help	ful
			>	 to log c 	ompressi	on
					J	J



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

