



SentiLog: Anomaly Detecting on Parallel File Systems via Logbased Sentiment Analysis <u>Di Zhang¹, Dong Dai¹, Runzhou Han², Mai Zheng²</u>

¹University of North Carolina at Charlotte



²Iowa State University







Importance of Anomaly Detection

Fugaku

Summit



00000100:00080000:0.0:1583538533.632216:0:3384:0:(import.c:681:ptlrpc_connect_import()) fff8f3dc0323000 lustre-MDT0000_UUID: changing import state from DISCONN to CONNECTING

00000100:00080000:0.0:1583538533.632225:0:3384:0:(importe:524:import electrom()) lustre-MDT0000-lwp-OST0002: connect to NID 10.24.86.160@tep tase actempt os971.8959

00000100:00080000:0.0:1583538533.632228:0:338.0:(import.c:568:import_select_connection()) lustre-MDT0000-lwp-OST0002: tried all connections, increasing latency to 50s

00000100:00050000:0.0:158355 558.63232:0:3384:0:(pinger.c:217:ptlrpc_pinger_process_import()) lustre-MDT0 00-lwp-OST0002_UUI ->lustre-MDT0000_UUID: level CONNECTING/4 force 0 force_next 0 deactive 0 pingable 1 suppress 0

00000100:00130000:0.0:1583538538.632342:0:3384:0:(pinger.c:230:ptlrpc_pinger_process_import()) lustre-MDT0000 lwp-OST0002 (TD->lustre-MDT0000_UUID: not pinging (in recovery or recovery disabled: CONNECTING.

00000100:00080000:0.0:1583538583. lustre-MDT0000-lwp-OST0002: tried a. 2:3384:0:(import.c:568:import_select_connection()) ctions, increasing latency to 60s



Logs of PFSes

LASSERT(imp_conn->oic_last_attempt); CDEBUG(D_HA, "%s: tried all connections, increasing latency 'to %ds\n", imp->imp_obd->obd_name, at_get(at));

Log Content

00000100:00080000:0.0:15835

ustre-MDT0000_UUID: changing import state from DISCON

Log Level

Log Index

0.0:1583538533.632225:0:3384:0:(import.c:524:import_select_connection())

ononcemport.c:681:ptlrpc_

lustre-MDT0000-lwp-OST0002: connect to NID 10.24.86.168@tcp last attempt 6557138959

00000100:00080000:0.0:1583538533 632228:0:3384:0:(import.c:568:import_select_connection())

lastre-MDT0000-lwp-OST0002: tried all connections, increasing latency to 50s CDEBUG

00000100:00080000:0.0:1583538558.632332:0:3384:0:(pinger.c:217:ptlrpc_pinger_process_import()) lustre-MDT0000-lwp-OST0002_UUID->lustre-MDT0000_UUID: level CONNECTING/4 force 0 force_next 0 deactive 0 pingable 1 suppress 0

00000100:00080000:0.0:1583538558.632342:0:3384:0:(pinger.c:230:ptlrpc_pinger_process_import())
lustre-MDT0000-lwp-0ST0002_UUID->lustre-MDT0000_UUID: not pinging (in recovery or recovery
disabled: CONNECTING)

00000100:00080000:0.0:1583538583.632228:0:3384:0:(import.c:568:import_select_connection()) lustre-MDT0000-lwp-0ST0002: tried all connections, increasing latency to 60s

Existing Work: Three Different Ways



[1] Deeplog: Anomaly detection and diagnosis from system logs through deep learning.

[2] Mining Invariants from Console Logs for System Problem Detection.

[3] LogAnomaly: Unsupervised Detection of Sequential and Quantitative Anomalies in Unstructured Logs



Challenge 1: Difficult to build appropriate sessions

Lustre

00000100:00080000:0.0:1607448618.327577:0:2290:0:(recover.c:58:ptlrpc_initiate_recovery ()) lustre-<u>OST0000_</u>UUID: starting recovery

00000100:00080000:0.0:1607448618.327580:0:2290:0:(import.c:681:ptlrpc_connect_import()) ffffa139cab87800 lustre-<u>OST0000_</u>UUID: changing import state from DISCONN to CONNECTING

00000100:00080000:0.0:1607448618.327589:0:2290:0:(import.c:524:import_select_connection ()) lustre-<u>OST0000</u>-osc-<u>MDT0000</u>: connect to NID <u>10.0.0.8</u>@tcp last attempt 4296114409

00000100:00080000:0.0:1607448618.327593:0:2290:0:(import.c:568:import_select_connection ()) lustre-<u>OST0000</u>-osc-<u>MDT0000</u>: tried all connections, increasing latency to 11s

HDFS

081109 203518 143 INFO dfs.DataNode\$DataXceiver: Receiving block <u>blk_</u>-<u>1608999687919862906</u> src: /10.250.19.102:54106 dest: /10.250.19.102:50010 081109 203518 35 INFO dfs.FSNamesystem: BLOCK* NameSystem.allocateBlock:

/mnt/hadoop/mapred/system/job_200811092030_0001/job.jar. <u>blk_-1608999687919862906</u>

081109 203519 143 INFO dfs.DataNode\$DataXceiver: Receiving block <u>blk_-</u> 1608999687919862906 src: /10.250.10.6:40524 dest: /10.250.10.6:50010

081109 203519 145 INFO dfs.DataNode\$PacketResponder: PacketResponder 1 for block <u>blk_-</u> <u>1608999687919862906</u> terminating

Challenge 1: Difficult to build appropriate sessions



Challenge 1: Difficult to build appropriate sessions



- Not able to build sessions based on their identifiers.
- Can only build sessions based on timestamps:
 - But, how to choose suitable time windows?
 - A small window may not include the relevant indices.
 - A large window have too many indices, which makes it difficult to discover the dependencies or invariants.

Log Index

Challenge 2: Log level may be inaccurate

Log_CRITICAL:Dec14 23:06:03 Main [App] » BeeGFS Helper Daemon Version: 7.2 Log_WARNING:Dec15 16:12:27 Main [App] » LocalNode: beegfs-mgmtd osboxes [ID:1] Log_WARNING:Dec15 15:58:37 Worker1 [Node registration] » New node: beegfsclient 435-5FD9237D-osboxes [ID: 2]; Source: 10.0.0.121:59206

- The three log entries above are simply reporting variable values, but they are labeled as 'Warning' or 'Critical' instead of normal by the developers.
- Previous study^[1] actually suggested the such variable printing logs were reported as normal level in 95% of the time in multiple opensource software.

[1] DeepLV: Suggesting Log Levels Using Ordinal Based Neural Networks.

Challenge 3: Labeled data is difficult to obtain



 Which lines are associated with anomalies and which are not?

How does SentiLog solve these challenges?





SentiLog Overview



Sentiment Analysis

"...before Portals cleanup: kmem %lld..."

"...Invoked LNET debug log upcall %s ..."

"cfs fail timeout ..."



Dec14 22:54:24 Main [App] » Unable to create subdir: buddymir/inodes/C/60 Dec16 15:39:34 Main [MgmtdTargetStateStore.cpp:446] » Could not read states. node-Type: beegfs-meta; Error: Path does not exist



Why multiple systems?

Biased Log Level
Correct Log Level







Comparing with Existing Solutions



- Lack of sequence info in PFSes logs makes DeepLog not suitable.
- DeepLog has too many false positives.

Comparing with Direct-Lookup



 Direct-Lookup: simply look up its corresponding logging statement in the source code and use its logging level to decide whether it is anomaly or not

Generality Evaluation



Multi source codes vs. Single source code

• SentiLog-Self: trained SentiLog using only the target PFS



Conclusion and Future Work

- Conclusion:
 - We propose to use sentimental analysis to understand log contents and detect the anomaly and show its effectiveness
 - We propose to train sentimental model using source code from multiple systems to solve the issue of lack of training data and to avoid bias of each system.
- Future Work:
 - Explore the possibility to consider more features besides the log statement description.
 - Conduct more experiments to validate and quantify the generic sentiment across different software.





Q&A Thank you!



