# Extracting Log Patterns from System Logs in LARGE

Yining Zhao

Computer Network Information Center

Chinese Academy of Sciences

HPBDC,  27 May 2016

# Content

- ScGrid in CAS

- The System of LARGE

- Log Patterns – Why?

- 2 Algorithms and Optimizations of Extracting

  - Identical Word Rate

  - Tree-matching

  - Optimizations

- Comparisons on Performances
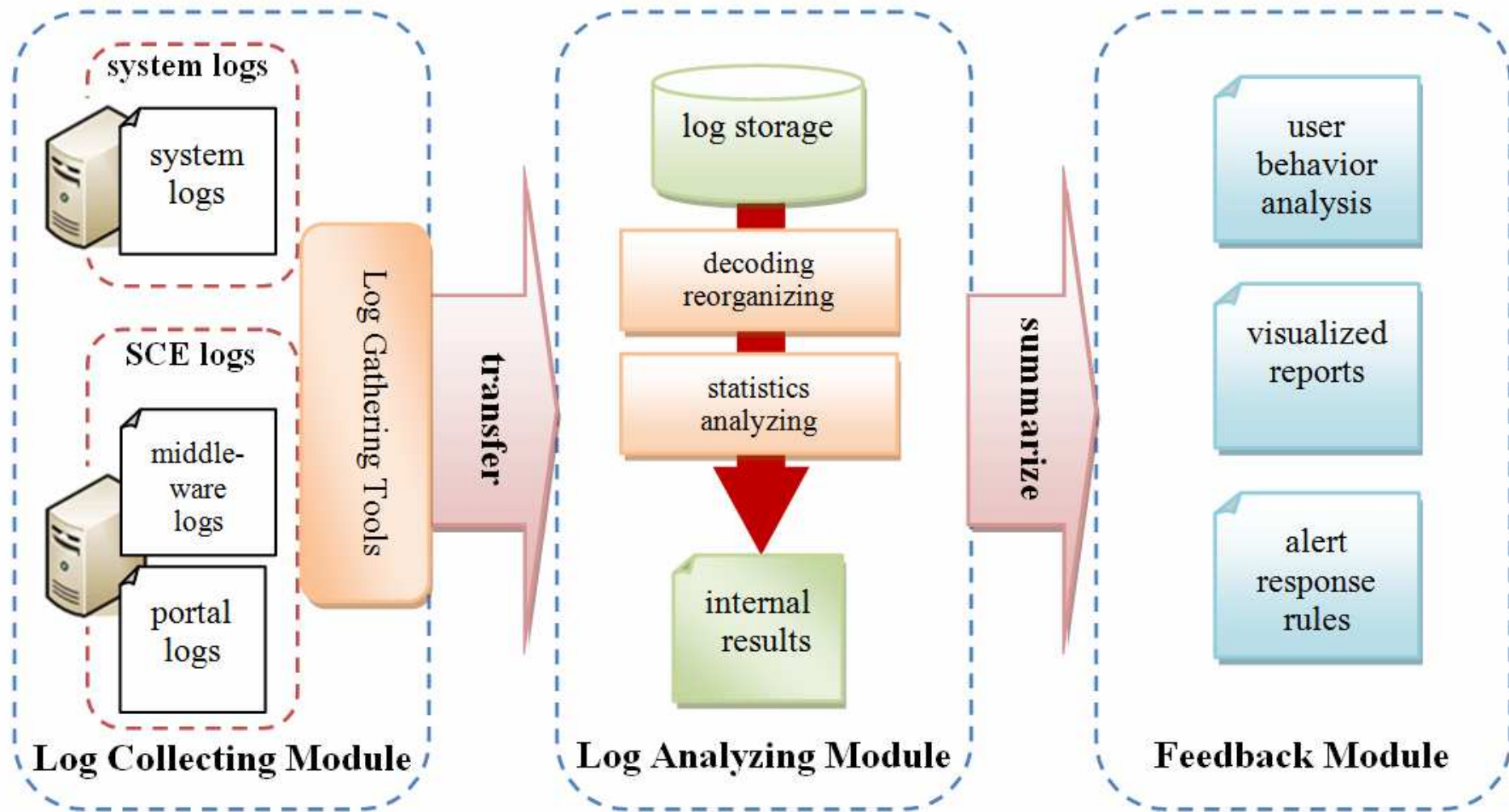
- Conclusion

# ScGrid in CAS

- Scientific Computing Grid Environment
- Integrated by many supercomputing centers in China
  - CNIC the head center
  - 8 national centers
  - 18 institute centers
- Using SCE middleware developed by SCCAS
- Provide computing resources
  to researches in various fields
  - Meteorology, Metal Forging,
    Fluid Mechanics, High Energy Physics,
    Computational Chemistry, Astrology…

# The System of LARGE

- Log Analyzing fRamework in Grid Environments
- Processing logs produced by the environment
  - gathering logs
  - processing, doing statistics and analyzing
  - producing feedbacks
- Two major types of logs
  - system logs – by log service in operating systems
  - SCE logs – by SCE middleware and job scheduling processes
- Helping the environment run correctly and steadily
  - generate alerts for particular patterns of logs
  - provide data for system analysis and maintenance
    - user behavior report
    - system errors and faults

# The System of LARGE

# Log Patterns – Why?

- We want to be alerted for logs in certain patterns, but…
  - too many logs for human to read
  - need to summarize patterns before defining alert rules
- Set of log patterns in our context:
  - patterns are different from each other
  - covering all logs in original set
  - significantly less than original
- The process of using log patterns
  - filter and remove frequent normal logs
  - use algorithms of extracting log patterns to get the set of patterns
  - manually check the set and pick out abnormal patterns
  - define rules to generate alerts for these patterns

# Algorithm of IWR

- Algorithm of identical word rate – a straight forward way
  - identical words
    - 2 words that are identical
    - and in the same position in 2 original logs
  - identical word rate: (number of identical words) / (total words)
  - predefined threshold $t$
- If IWR is greater than $t$, the two logs are in one pattern
- Logs with different length has IWR of ZERO!

```
It is a good day
It isn't a bad day
```

$t = 0.66$
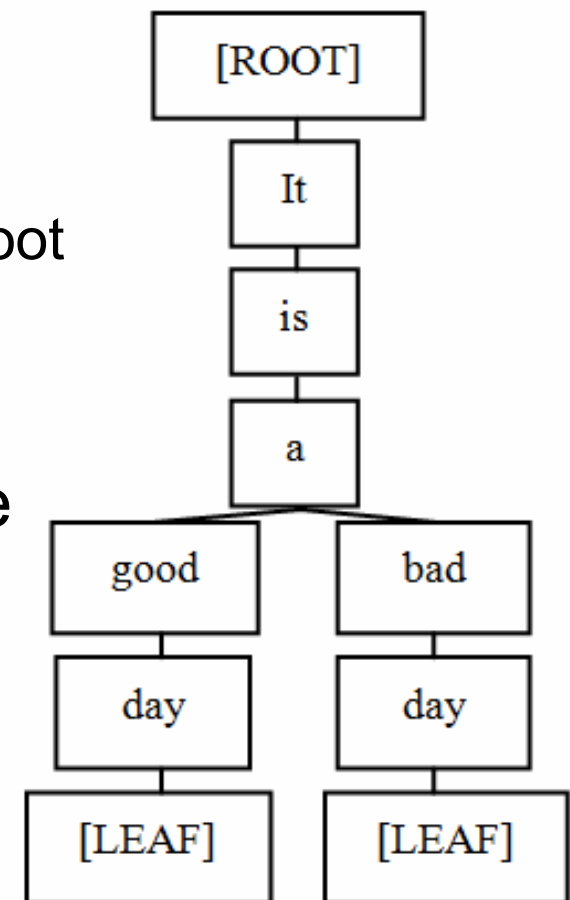IWR = 60%
NOT in the same pattern!

# Algorithm of IWR

- Process of algorithm of IWR
  - set threshold *t* and initial empty pattern set *P*
  - for each new incoming logs, compute IWR with each pattern in *P*
  - if pattern matched, skip to next; if none matched, add to *P*
- *P* will be affected by order of input
  - 3 logs, *L1:{a, b, c} L2:{a, b, d} L3:{d, b, c}, t* = 0.6
  - in order of *l1, l2, L3*, only 1 pattern left
  - in order of *l2, l1, l3*, 2 patterns left
- Not ideal in complexity
  - O(n²)

# Algorithm of Tree-matching

- Different in storing and matching structure with IWR
- Words stored orderly in a tree
  - branches for different words
- When performing the algorithm
  - compare each word in the incoming log from root
  - if successfully matched to leaf, check next
  - if unmatched found, create a new branch
- To get pattern set, use depth-first traverse
- Better complexity, but worse result
  - O(n $log$(n))
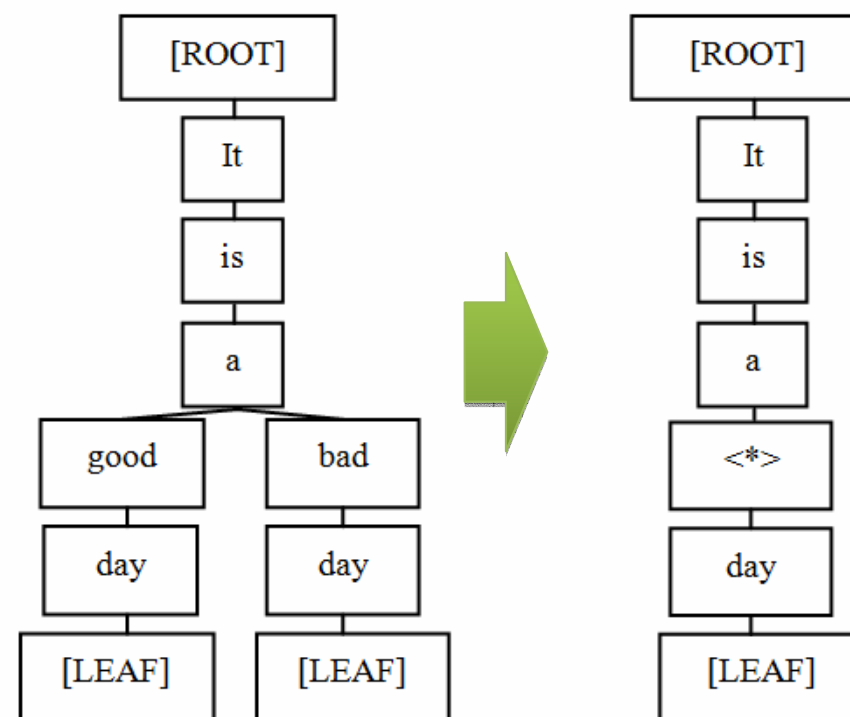  - unacceptable number of patterns
  - need optimizations

# Algorithm Optimizations

- Tree pruning
  - if two nodes has same subtrees, merge to a key node (<*>)
  - after previous step, merge all key nodes to keep uniqueness
  - key node can be matched to any word, but only the last option

- we can do this because…

```
Login failed from user alice: password failed
Login failed from user bob: password failed
Login failed from user chris: password failed
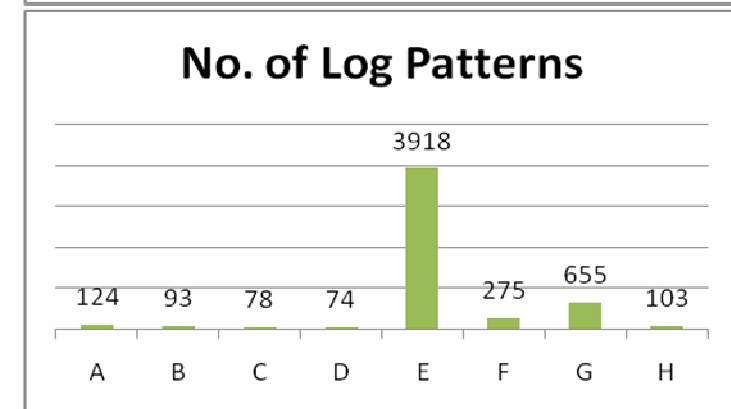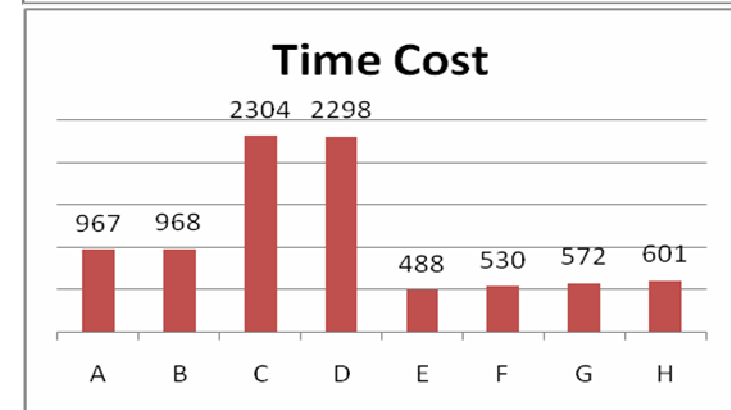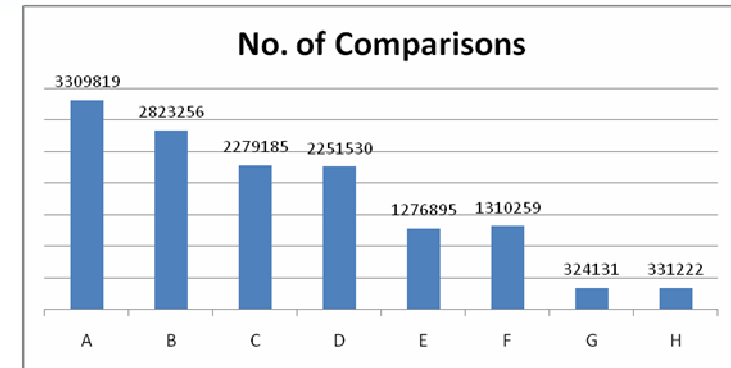```

key information position

# Algorithm Optimizations

- Word converting function
  - a preprocessing optimization
  - converting commonly appeared expressions to predefined strings

- We can do this because…
  - At this stage, IPs, usernames, etc. are not essential
  - could be a distraction for extracting and enlarge pattern set

- Helpful for tree-matching, but not IWR
  - IWR has higher tolerance for differences of words
  - Tree-matching is more sensitive

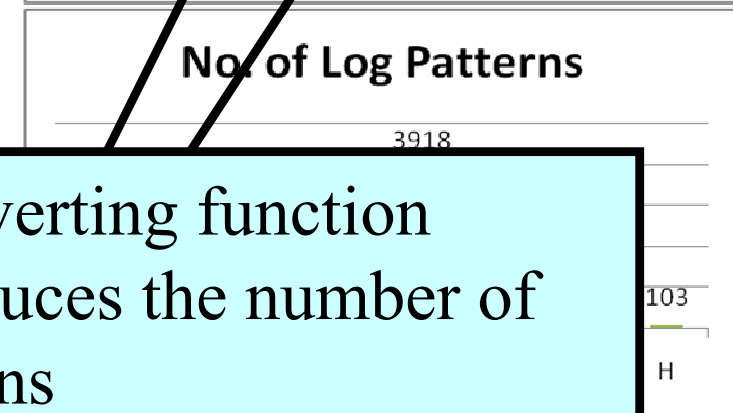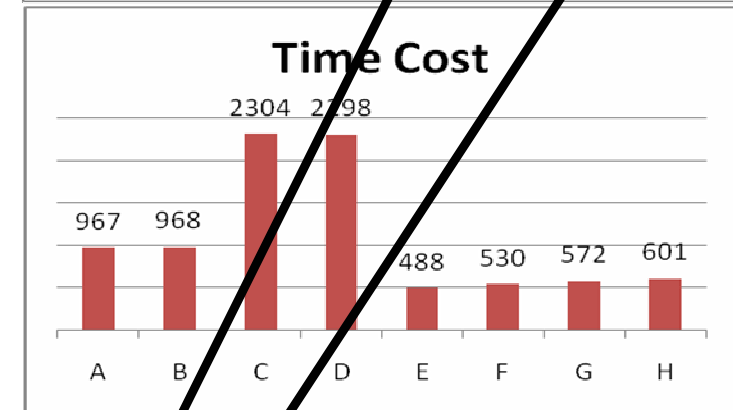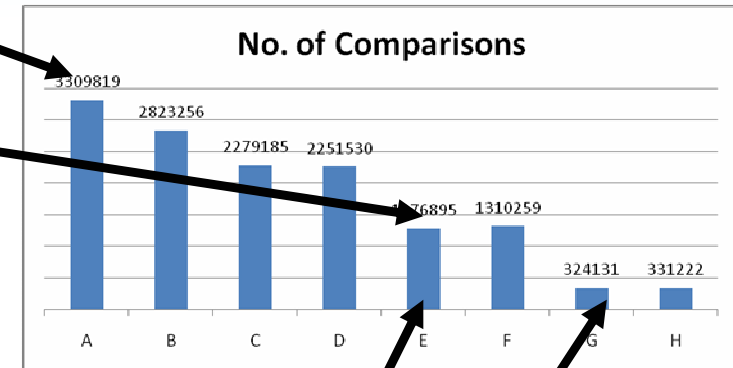| Word Format in Regular Expression | Converted Result |
|---|---|
| [0-9]+.[0-9]+ .[0-9]+ .[0-9]+ | <IP> |
| [0-9]+ | <NUMBER> |
| name=[A-Za-z][A-Za-z0-9_]+ | name=<USER> |
| UID=[0-9]+ | UID=<NUMBER> |
| GID=[0-9]+ | GID=<NUMBER> |

# Comparisons on Performances

original: 49079 input logs

A. IWR, t = 0.66

B. IWR, t = 0.6

C. IWR + word converting, t = 0.66

D. IWR + word converting, t = 0.6

E. Tree-matching

F. Tree-matching + pruning

G. Tree-maching + word converting

H. Tree-matching + both

# Comparisons on Performances

**Tree-matching has significantly reduced the number of comparisons**

o

A. IWR, t = 0.66
B. IWR, t = 0.6
C. IWR + word converting, t = 0.66
D. IWR + word converting, t = 0.6
E. Tree-matching
F. Tree-matching + pruning
G. Tree-maching + word converting
H. Tree-matching + both

**No. of Comparisons**

3309819
2823256
2279185  2251530
76895  1310259
324131  331222

A  B  C  D  E  F  G  H

**Time Cost**

2304  2298
967  968
488  530  572  601

A  B  C  D  E  F  G  H

**No. of Log Patterns**

3918
103
H

**Word converting function further reduces the number of comparisons**
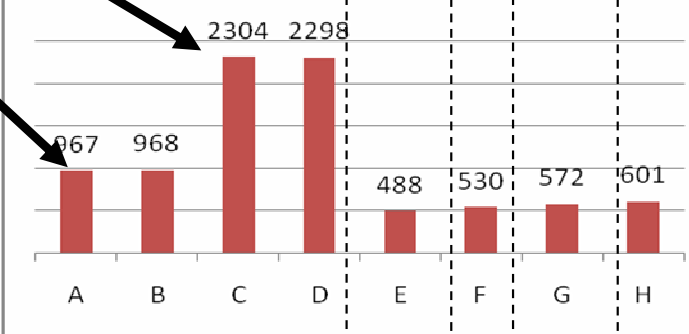
Word converting function is quite time costly

orig...

A. IWR, t = 0.66

B. IWR, t = 0.6

C. IWR + word converting, t = 0.66

D. IWR + word converting, t = 0.6

E. Tree-matching

F. Tree-matching + pruning
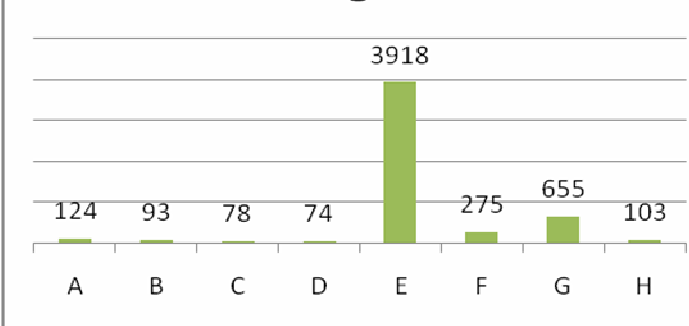
G. Tree-maching + word converting

H. Tree-matching + both

**No. of Comparisons**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 3309819 | 2823256 | 2279185 | 2251530 | 1276895 | 1310255 | 324131 | 331222 |

**Time Cost**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 1067 | 968 | 2304 | 2298 | 488 | 530 | 572 | 601 |

**No. of Log Patterns**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 124 | 93 | 78 | 74 | 3918 | 275 | 655 | 103 |

**GRID** 中国国家网格 China National Grid · **SCGRID**

orig Extracting efficiency (time and comparison costs): tree-matching is better

A.

B. IWR, t = 0.6

C. IWR + word converting, t = 0.66
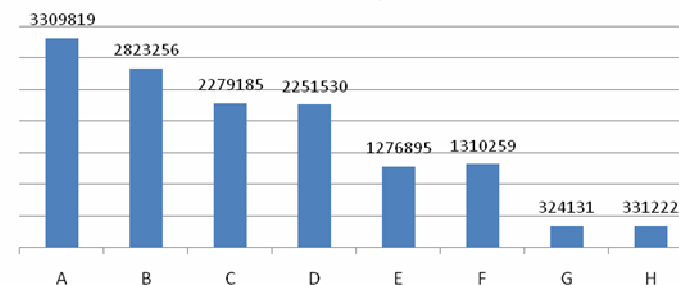
D. IWR + word converting, t = 0.6
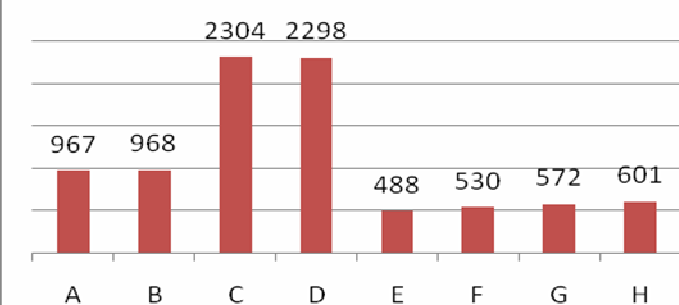
E. Tree-matching

F. Tree-matching + pruning

G. Tree-maching + word converting
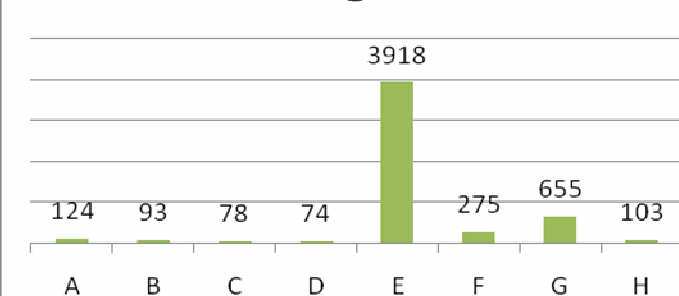
H.

Extracting effect (number of patterns): IWR is better

**No. of Comparisons**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 3309819 | 2823256 | 2279185 | 2251530 | 1276895 | 1310259 | 324131 | 331222 |

**Time Cost**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 967 | 968 | 2304 | 2298 | 488 | 530 | 572 | 601 |

**No. of Log Patterns**

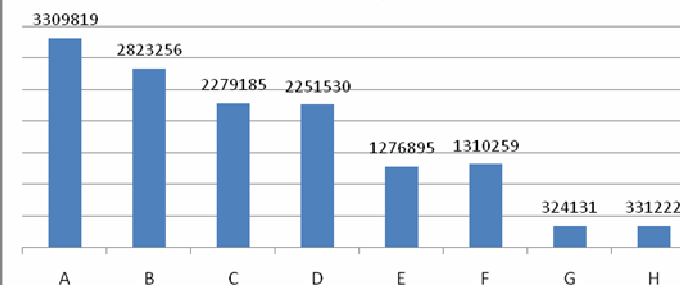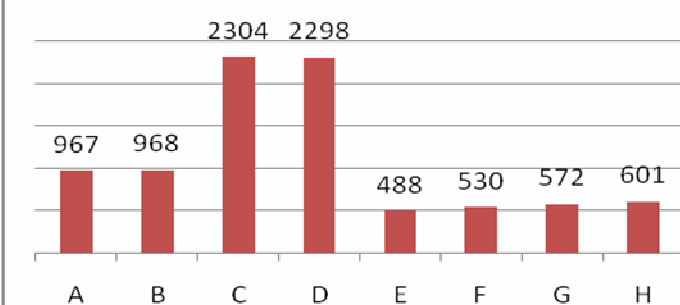| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 124 | 93 | 78 | 74 | 3918 | 275 | 655 | 103 |

# Comparisons on Performances

Lower threshold gives lesser patterns in IWR

A. IWR, t = 0.66
B. IWR, t = 0.6
C. IWR + word converting, t = 0.66
D. IWR + word converting, t = 0.6
E. Tree-matching
F. Tree-matching + pruning
G. Tree-maching + word converting
H. Tree-matching + both

**No. of Comparisons**
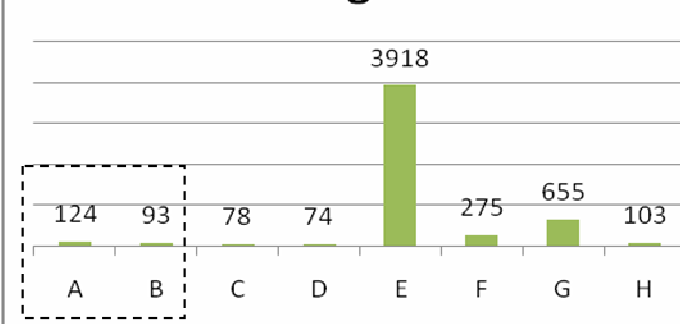
| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 3309819 | 2823256 | 2279185 | 2251530 | 1276895 | 1310259 | 324131 | 331222 |

**Time Cost**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 967 | 968 | 2304 | 2298 | 488 | 530 | 572 | 601 |

**No. of Log Patterns**

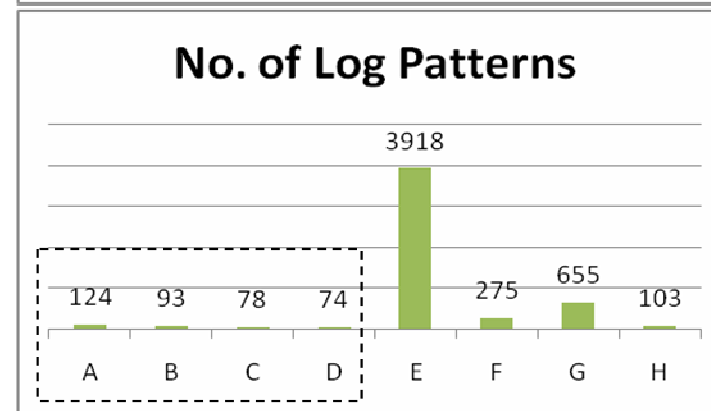| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 124 | 93 | 78 | 74 | 3918 | 275 | 655 | 103 |

Number of comparisons has nearly direct proportion to the number of extracted patterns

A. IWR, t = 0.66

B. IWR, t = 0.6

C. IWR + word converting, t = 0.66

D. IWR + word converting, t = 0.6

E. Tree-matching

F. Tree-matching + pruning

G. Tree-maching + word converting
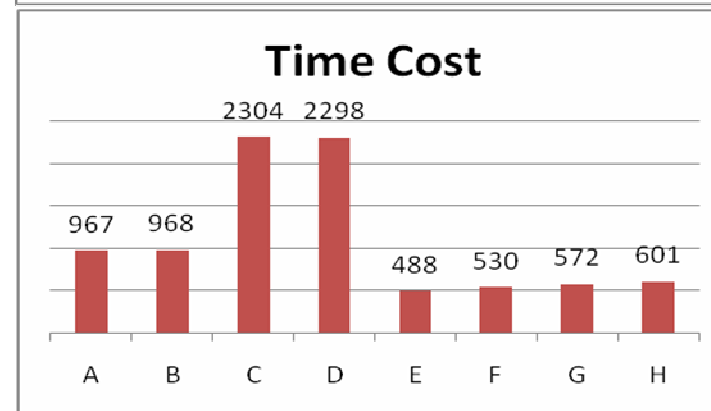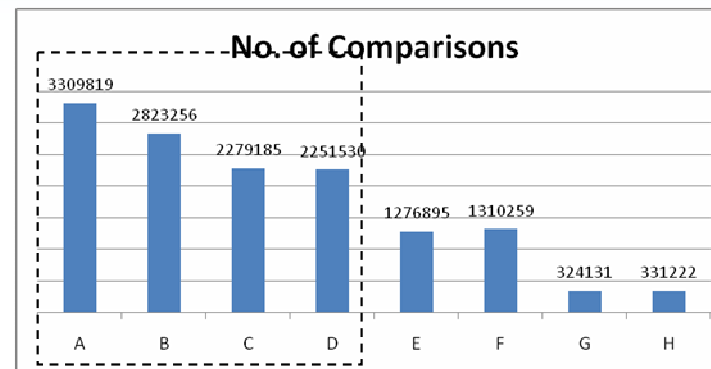
H. Tree-matching + both



No. of Comparisons

A 3309819, B 2823256, C 2279185, D 2251536, E 1276895, F 1310259, G 324131, H 331222

Time Cost

A 967, B 968, C 2304, D 2298, E 488, F 530, G 572, H 601

No. of Log Patterns

A 124, B 93, C 78, D 74, E 3918, F 275, G 655, H 103

Pure tree-matching:

good in extract efficiency

very bad in extracting effect

A.

B.  IWR, t = 0.6

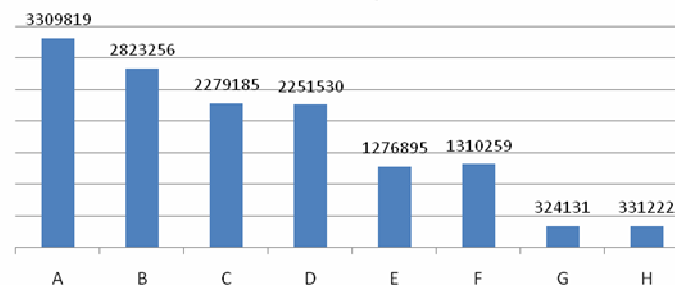C.  IWR + word converting, t = 0.66

D.  IWR + word converting, t = 0.6
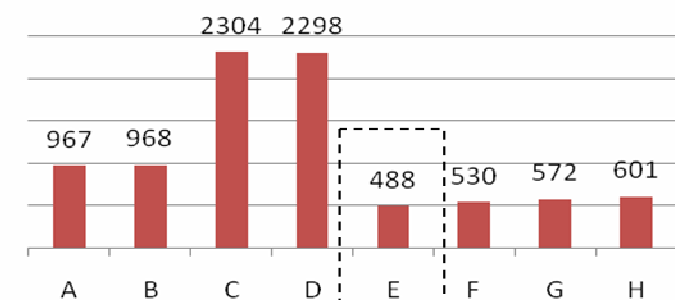
E.  Tree-matching

F.  Tree-matching + pruning

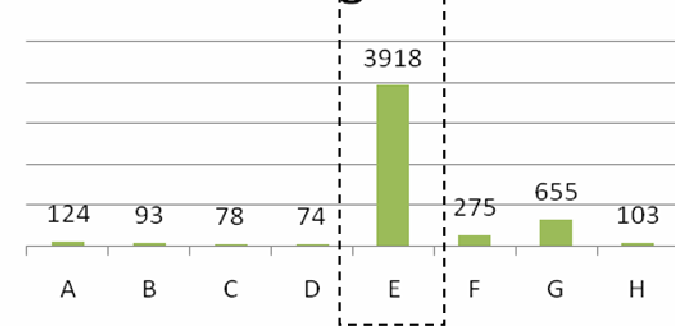G.  Tree-maching + word converting

H.  Tree-matching + both

**No. of Comparisons**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 3309819 | 2823256 | 2279185 | 2251530 | 1276895 | 1310259 | 324131 | 331222 |

**Time Cost**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 967 | 968 | 2304 | 2298 | 488 | 530 | 572 | 601 |

**No. of Log Patterns**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 124 | 93 | 78 | 74 | 3918 | 275 | 655 | 103 |

o Pruning has greatly improved the extracting effect

A. IWR, t = 0.66

B. IWR, t = 0.6

C. IWR + word converting, t = 0.66

D. IWR + word converting, t = 0.6

E. Tree-matching

F. Tree-matching + pruning

G. Tree-maching + word converting

H. Tree-matching + both

**No. of Comparisons**

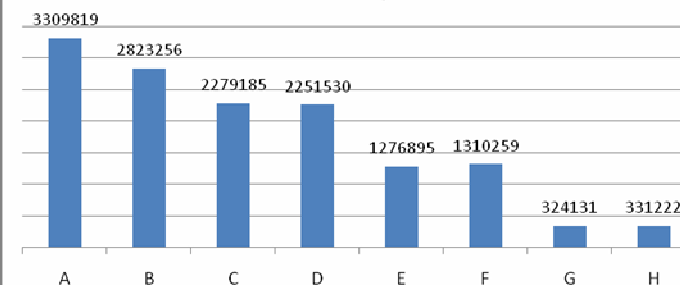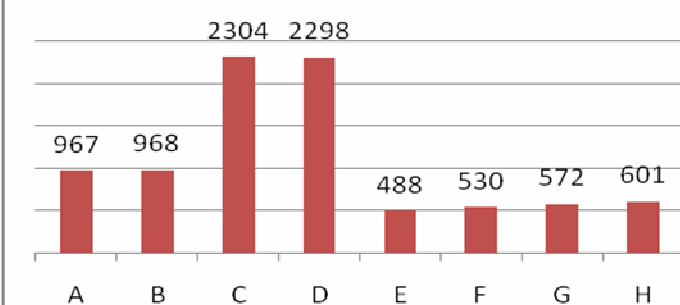A: 3309819, B: 2823256, C: 2279185, D: 2251530, E: 1276895, F: 1310259, G: 324131, H: 331222

**Time Cost**

A: 967, B: 968, C: 2304, D: 2298, E: 488, F: 530, G: 572, H: 601

**No. of Log Patterns**

A: 124, B: 93, C: 78, D: 74, E: 3918, F: 275, G: 655, H: 103

# Comparisons on Performances

original: 49079 input logs

A. I...

B. I...

C. I...

D. IWR + word converting, t = 0.6

E. Tree-matching

F. Tree-matching + pruning
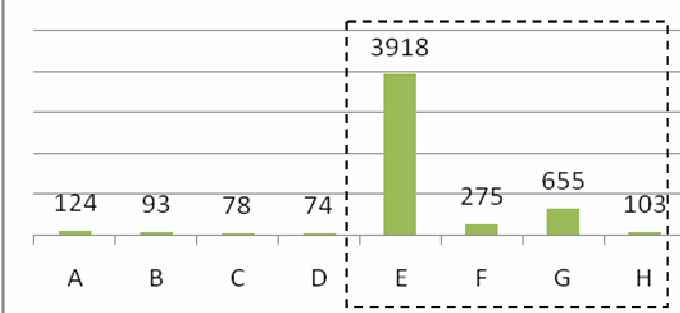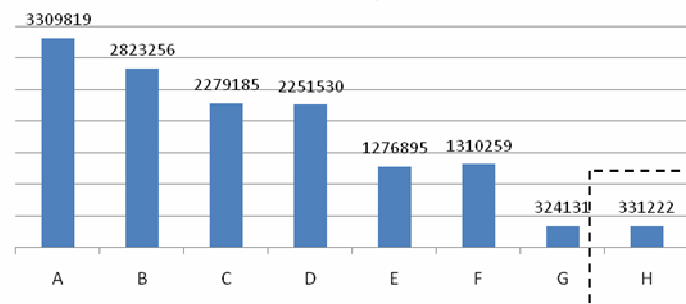
G. Tree-maching + word converting

H. Tree-matching + both

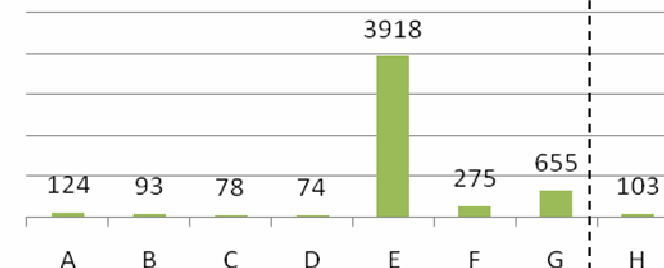Tree-matching with both optimizations: good in extracting efficiency and effect, satisfying algorithm!

### No. of Comparisons

3309819 — A
2823256 — B
2279185 — C
2251530 — D
1276895 — E
1310259 — F
324131 — G
331222 — H

### Time Cost

967 — A
968 — B
2304 — C
2298 — D
488 — E
530 — F
572 — G
601 — H

### No. of Log Patterns

124 — A
93 — B
78 — C
74 — D
3918 — E
275 — F
655 — G
103 — H

# Conclusion

- LARGE is a log analyzing system
- When monitoring system logs, we need to extract log patterns
- Two algorithms: IWR and tree-matching, plus optimizations
  - tree-matching with two optimizations looks good

- Future work
  - what if more than one key position in tree-matching algorithm?
  - we may use parallel computing to accelerate log processing

# Thank you!