# Analyzing log analysis:
# an empirical study of user log mining

Sara Alspaugh

@salspaugh

AMPLab, EECS Department, UC Berkeley
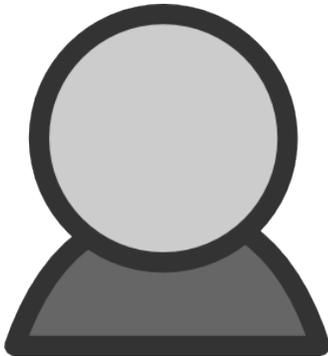
*in collaboration with:*
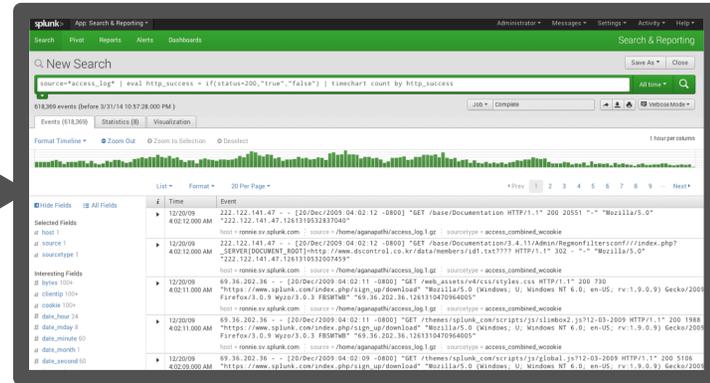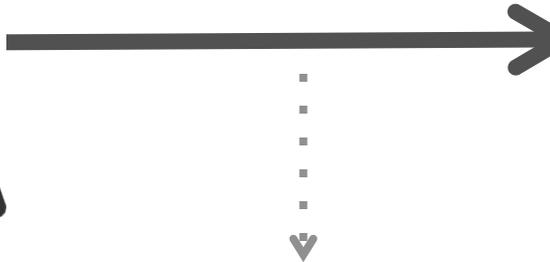
Archana Ganapathi (Splunk),

Beidi Chen, Jessica Lin, Marti Hearst, Randy Katz (UC Berkeley)

**USENIX LISA 2014**
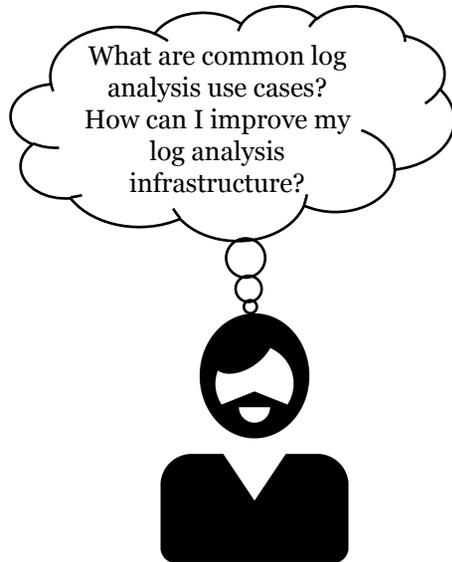
# Data collected

**User**

**Log analysis tool**

**Queries**

```
search index=os eventtype=linux-password-change-failed
search index=os eventtype="Failed_SU" index="os" sourcetype="interfaces" host=* | multikv
fields name, inetAddr, RXbytes, TXbytes | streamstats current=f last(TXbytes) as lastTX,
last(RXbytes) as lastRX by Name  | eval time=_time | strcat Name "-" inetAddr "@" host
Interface_Host | eval RX_Thruput_KB = (lastRX-RXbytes)/1024 | eval TX_Thruput_KB =
(lastTX-TXbytes)/1024 | timechart eval(sum(TX_Thruput_KB)/dc(time)) by Interface_Host
search index=os sourcetype=openPorts | MULTIKV | STATS count BY Port | SORT count
search index=os source=ps | multikv | timechart avg(VSZ_KB) by USER useother=F limit=10
"" | strcat source "@" host changelist | timechart count by changelist
search sourcetype=syslog error OR failed OR severe NOT assignment starthoursago=1 |
fields +_raw
search index=os source=ps | multikv | timechart avg(RSZ_KB) by COMMAND
search index=os source=iostat | multikv | timechart avg(rReq_PS) avg(wReq_PS)
search index=os source=lsof | multikv | timechart count(USER) by USER
search index=os source=vmstat | multikv | timechart avg(memTotalMB) by host
```
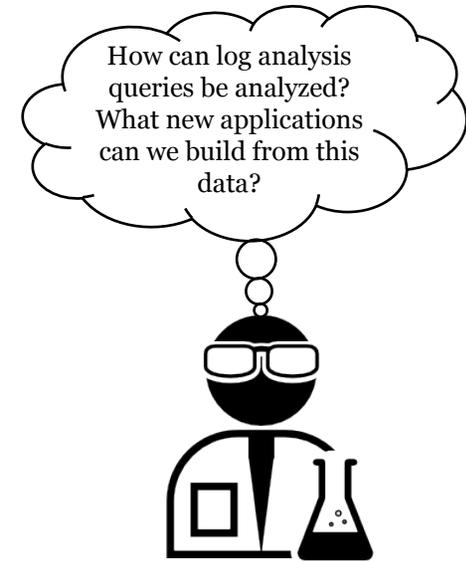
# Motivation

What are common log analysis use cases? How can I improve my log analysis infrastructure?

*log analysis practitioners*

How can we be more relevant to our customers? How can we improve customer experience?

*data analysis and management companies*

How can log analysis queries be analyzed? What new applications can we build from this data?

*computer science researchers*

What do people use a highly popular commercial purpose-built log analysis tool to do?

# Splunk screenshot

# Example Splunk query

```
search "error"
| stats count by status
| lookup statuscodes status OUTPUT statusdesc
```

| 0.0 | - | **error** | 404 |
|-----|---|-----------|-----|
| 0.5 | - | OK | 200 |
| 0.7 | - | **error** | 500 |
| 1.5 | - | OK | 200 |

→ `search "error"` →

*stage 1*

| 0.0 | - | **error** | 404 |
|-----|---|-----------|-----|
| 0.7 | - | **error** | 500 |

↓

`stats count`
`by status`

*stage 2*

↓

| count | status | statusdesc |
|-------|--------|------------|
| 1 | 404 | Not Found |
| 1 | 500 | Internal Server Error |

← `lookup statuscodes`
`status OUTPUT statusdesc` ←

*stage 3*

| count | status |
|-------|--------|
| 1 | 404 |
| 1 | 500 |

# Practitioner viewpoint

- What are the primitives of log analysis?
  - commands or transformations
- What are the main tasks of log analysis?
  - additional detail
- Why do users analyze logs?
  - context, roles, goals, use cases

# Practitioner viewpoint

- ## What are the primitives of log analysis?
  - commands or transformations
- ## What are the main tasks of log analysis?
  - additional detail
- ## Why do users analyze logs?
  - context, roles, goals, use cases

# Splunk commands in order of frequency of appearance

| command | count |
|---|---:|
| search | 232373 |
| eval | 178080 |
| stats | 75927 |
| table | 44967 |
| fields | 37803 |
| rename | 35919 |
| where | 32402 |
| inputlookup | 30490 |
| sort | 30442 |
| lookup | 28620 |
| outputlookup | 27042 |
| dedup | 22731 |

*... snip ...*

| | |
|---|---:|
| localop | 27 |
| reverse | 15 |
| abstract | 10 |
| map | 7 |
| anomalies | 3 |
| extract | 2 |
| outlier | 2 |
| datamodel | 2 |
| format | 1 |
| outputtext | 1 |
| dbinspect | 1 |

Approach: impose hierarchical organization into tasks, sub-tasks, lower-level activities

addinfo
appendcols
bin
bucket
eval
extract
iplocation
kv
outputtext
rangemap
rex
spath
strcat
xmlkv

augment

dedup
head
regex
search
tail
uniq
where

filter

addcoltotals
counttable
eventcount
geostats
stats
timechart
top

aggregate

# Top log analysis transformations
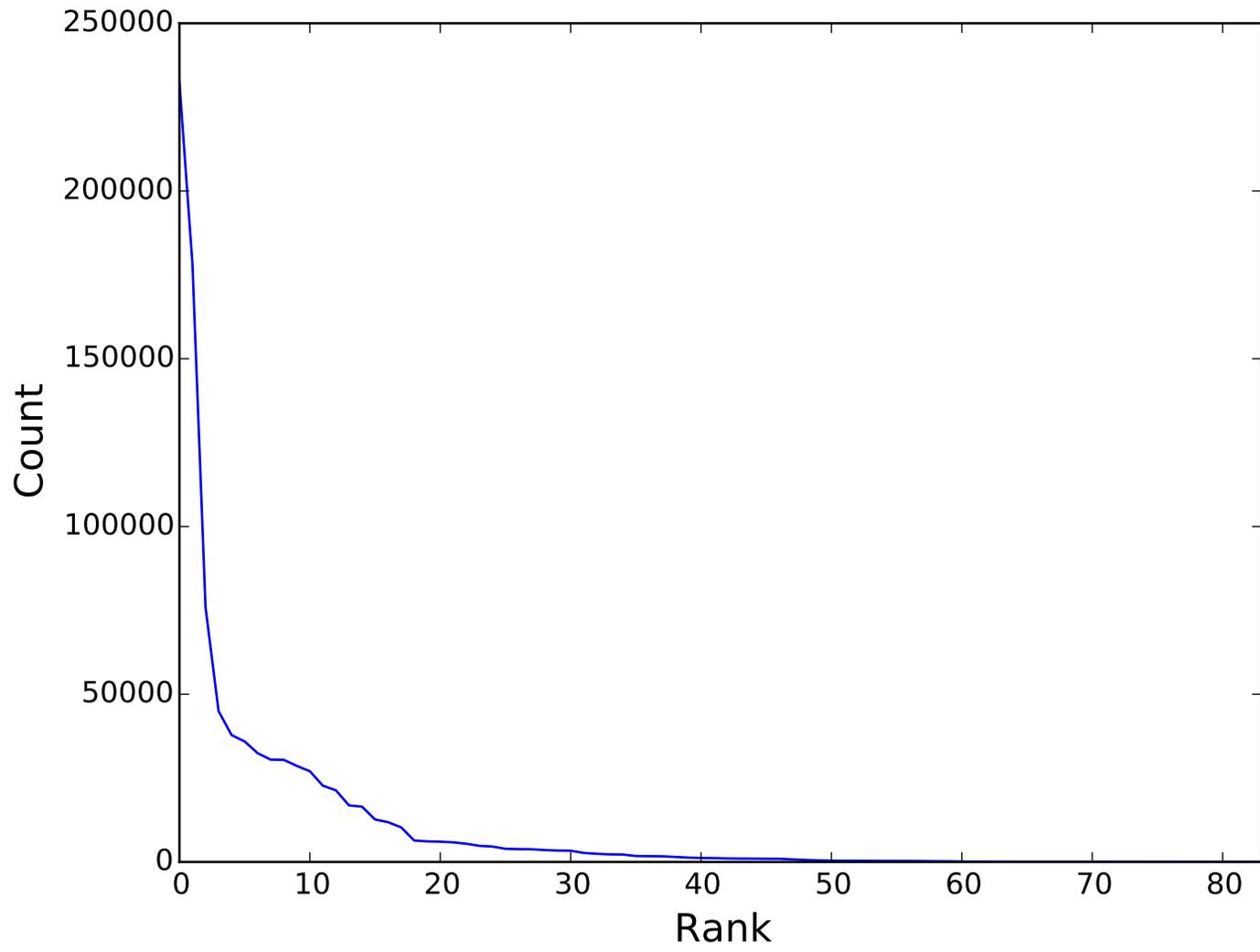
# Practitioner viewpoint

- What are the primitives of log analysis?
  - commands or transformations
- **What are the main tasks of log analysis?**
  - **additional detail**
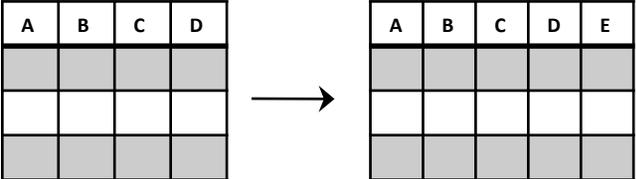- Why do users analyze logs?
  - context, roles, goals, use cases

# Transformation details

0. Split queries into stages.
1. Featurize each stage in given category (i.e., filter, augment, aggregate).
2. Perform PCA for dimensionality reduction.
3. Perform t-SNE to visualize.
4. Label clusters.

Types of **Filter** transformations

| | Percent |
|---|---|
| Filters by long logical condition | `search (x>10 and x< 20) or ...` |
| Filters by specifying fields | `search user=salspaugh` |
| Filters by string contains | `search "error"` |
| Selects | `search index=os` |
| Deduplicates | |
| Filters by time range | `search daysago=7` |
| Uses macro | |
| Filters by result of function | |
| Uses subsearch | |
| Filters by index | |
| Filters by regex | |

Types of **Augment** transformations

*Percent*

| | 0 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |

- String manipulation — `eval user=lower(user)`
- Conditional statement — `eval weekend=if(day="Sat" or day="Sun",1,0)`
- Arithmetic — `eval percent=cnt/total*100`
- Date or time calculation — `eval hour=strftime(time, "H%")`
- Multi-value operation — `eval first=mvindex(lst, 0)`
- Assigns to group
- Assigns simple value — `eval lastseen=_time`
- Uses subsearch

Types of **Aggregate** transformations

*Percent*

| | 0 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |

- Aggregation — `stats count`
- Visualize aggregation over time — `timechart count`
- Visualize aggregation — `chart count`
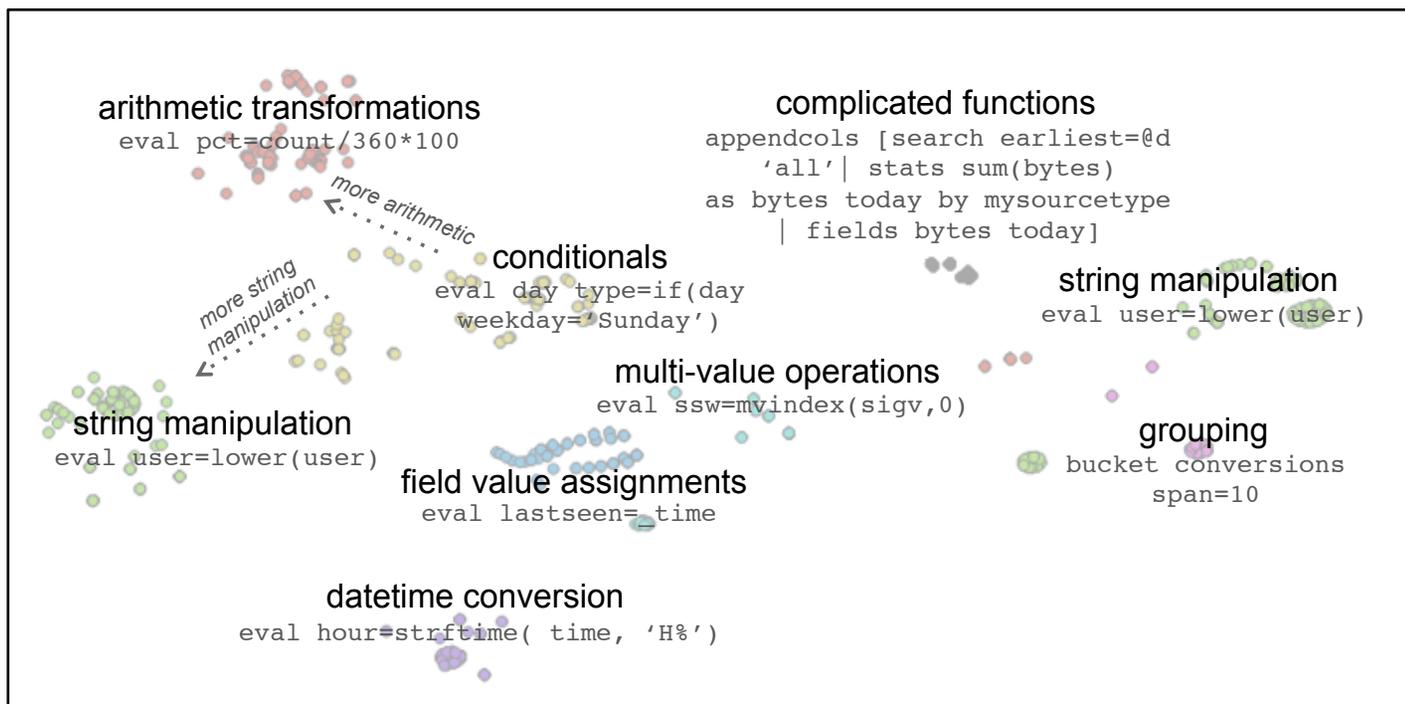- Aggregate, sort, limit — `top`
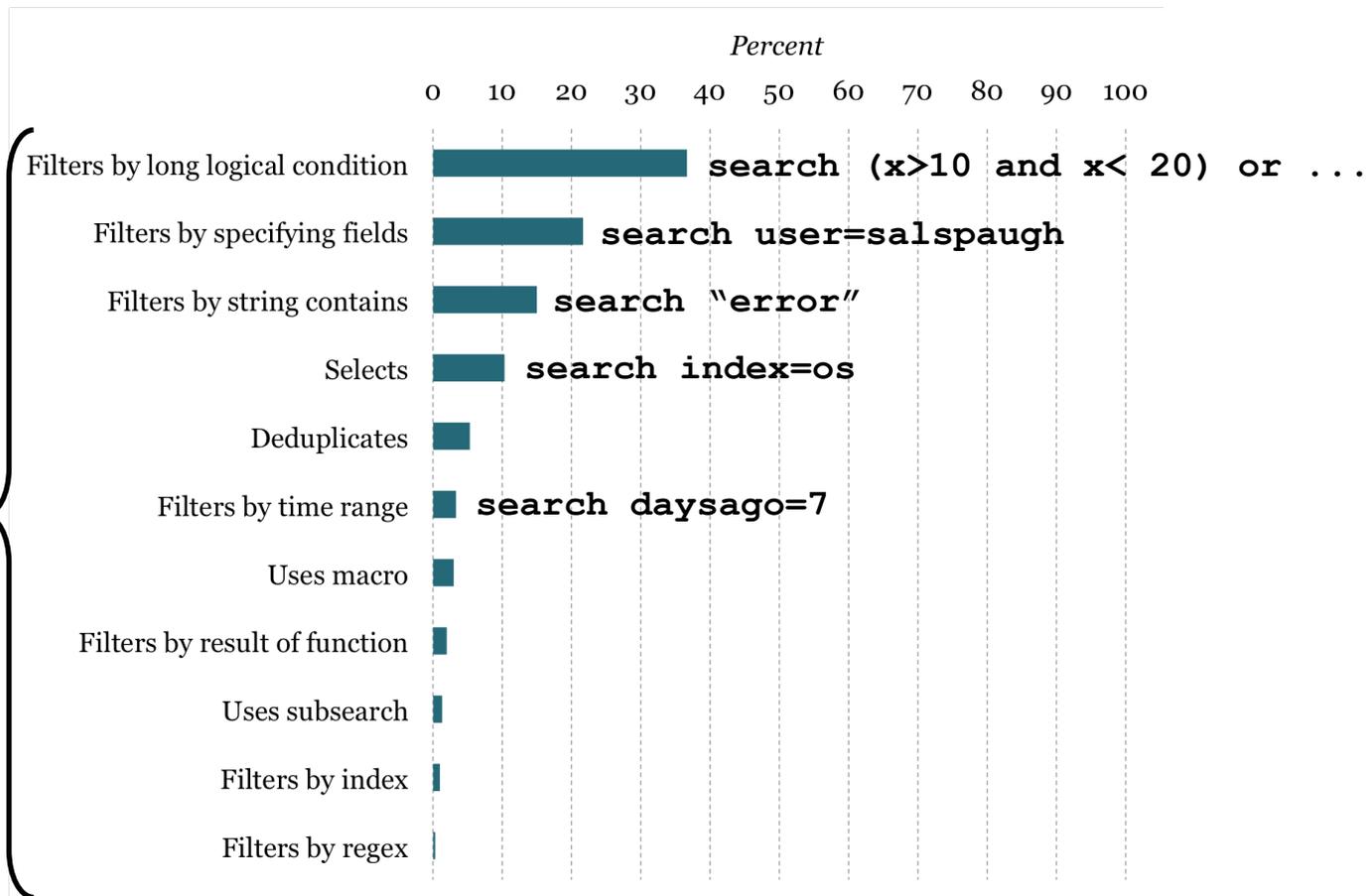- Group by time

# Practitioner viewpoint

- What are the primitives of log analysis?
  - commands or transformations
- What are the main tasks of log analysis?
  - additional detail
- **Why do users analyze logs?**
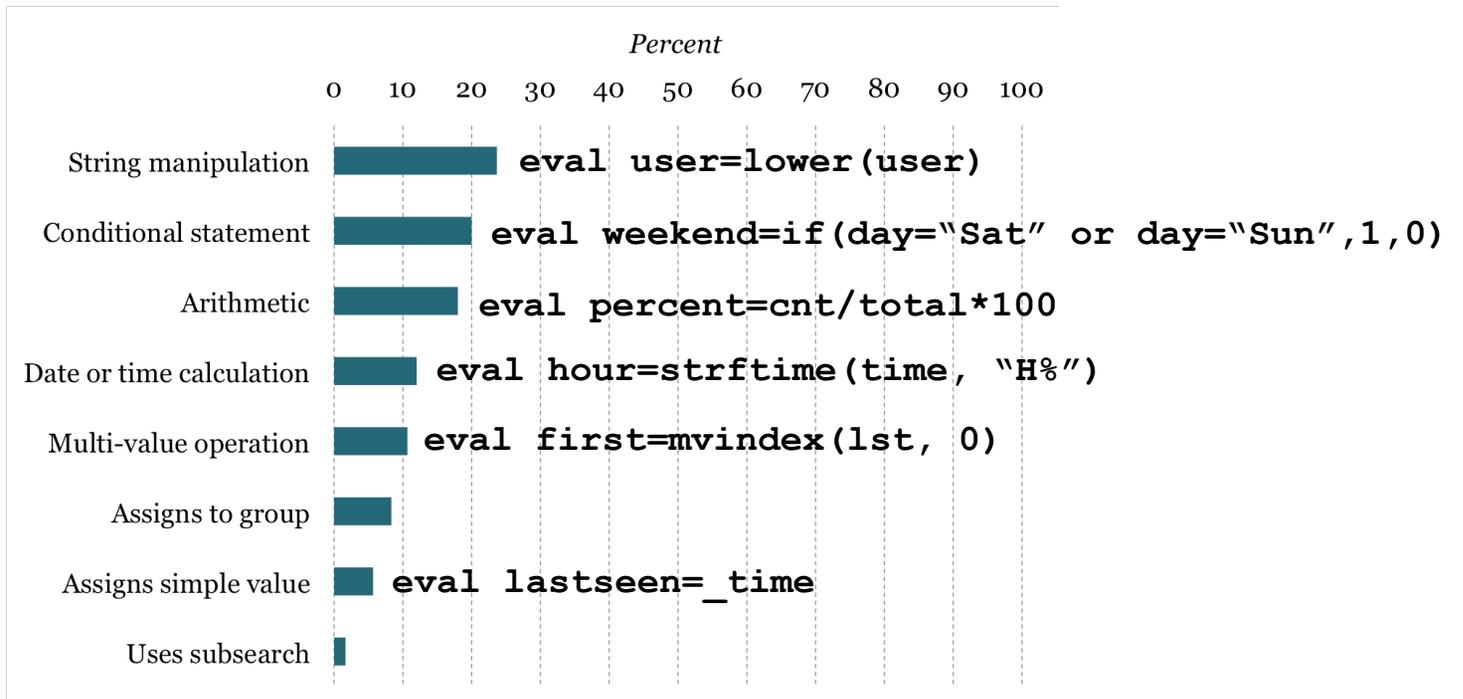  - **context, roles, goals, use cases**

# Sales engineer survey

- What are the roles of the primary Splunk users within the organization?
- Write-in answers:
  - manufacturing team
  - data analysts
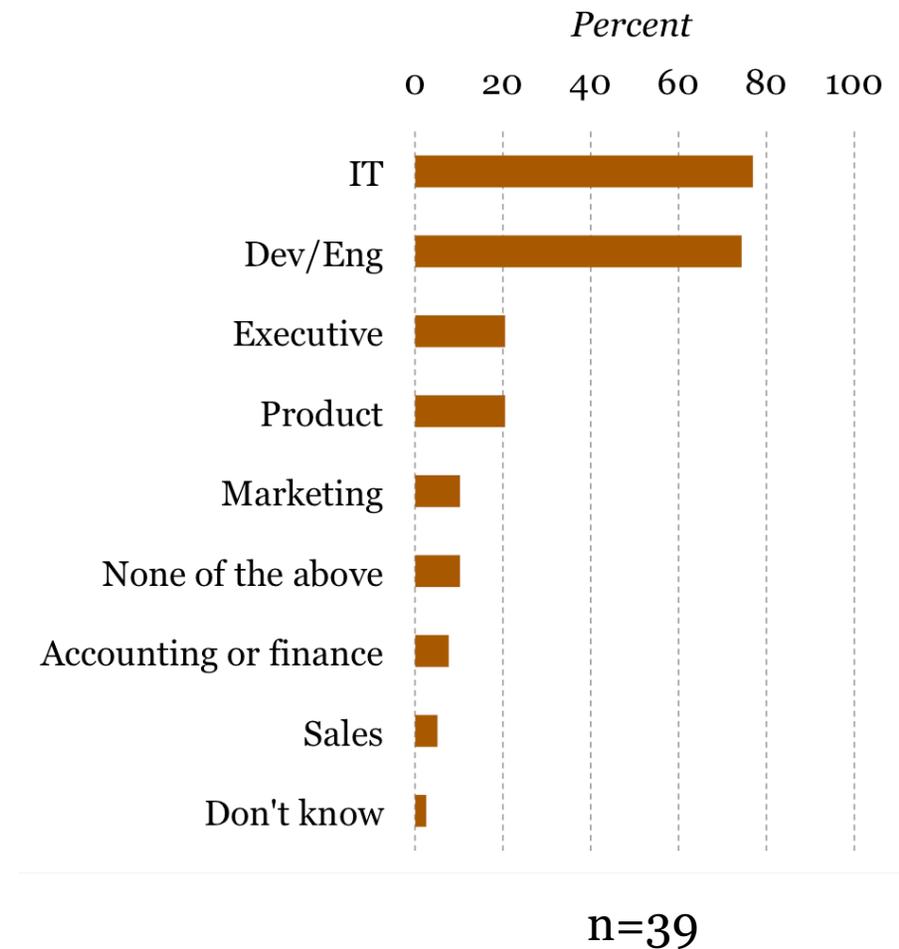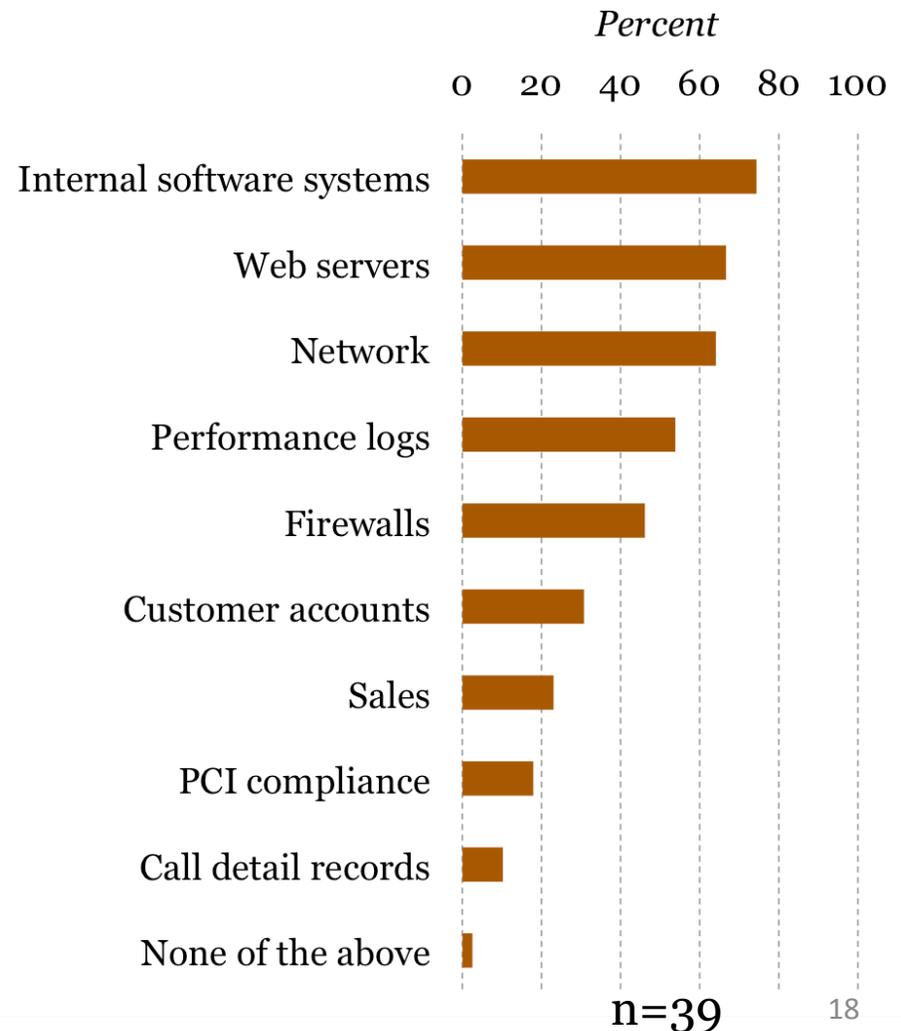  - compliance team
  - security team
  - email team

*Percent*

| Role | |
|---|---|
| IT | |
| Dev/Eng | |
| Executive | |
| Product | |
| Marketing | |
| None of the above | |
| Accounting or finance | |
| Sales | |
| Don't know | |

n=39

# Sales engineer survey

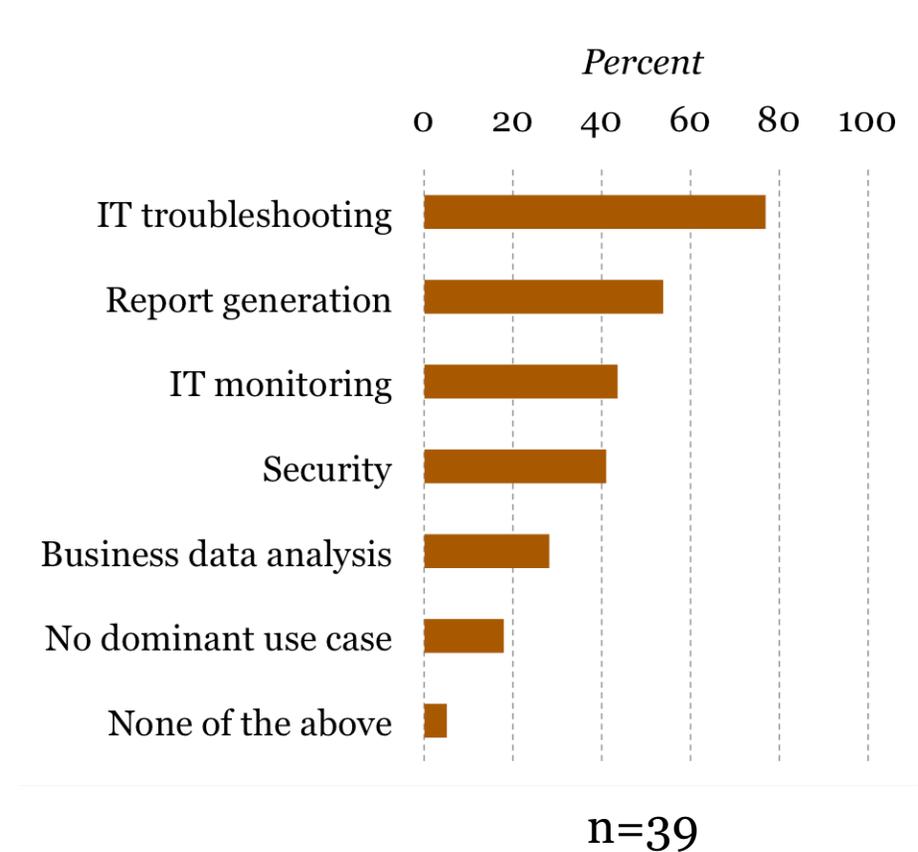- Roughly what types of data sources did each of these customers have?

- Write-in answers:
  - email
  - sensors
  - mobile apps
  - middleware
  - custom applications

*Percent*

| | 0 | 20 | 40 | 60 | 80 | 100 |
|---|---|---|---|---|---|---|

Internal software systems
Web servers
Network
Performance logs
Firewalls
Customer accounts
Sales
PCI compliance
Call detail records
None of the above

n=39

18

# Sales engineer survey

- What problems does the organization typically try to address with Splunk?

- Write-in answers:
  - app management
  - customer satisfaction
  - workflow automation
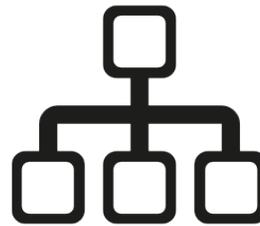  - email monitoring
  - manufacturing

*Percent*

| | 0 | 20 | 40 | 60 | 80 | 100 |

IT troubleshooting

Report generation

IT monitoring

Security

Business data analysis

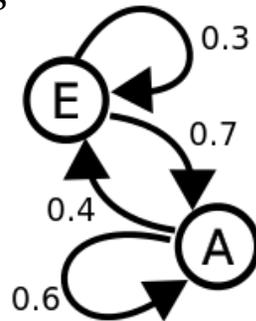No dominant use case

None of the above

n=39

# In the paper

How many different types of tasks are there? How is the frequency of each of these tasks distributed? What are the most and least common tasks?

*tasks*

Within the high-level transformation tasks, what sub-tasks are performed? What is the frequency of each of these subtasks?

*sub-tasks*

How are sequences of tasks statistically distributed? What type of tasks usually come first? What comes last? What tasks typically follow a given other task? How many tasks are performed in an average query? What are common subsequences of tasks?
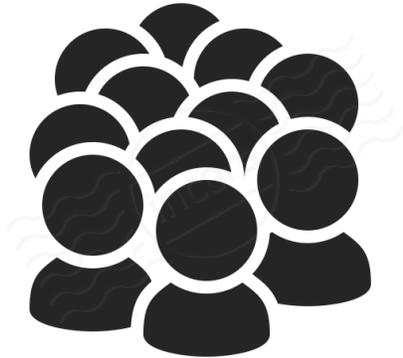
0.3
E
0.7
0.4
A
0.6

*pipelines*

What are the primary roles of Splunk users? What problems do they use Splunk for? What other software do they use along with Splunk? What is their level of Splunk expertise? How technical are Splunk users?
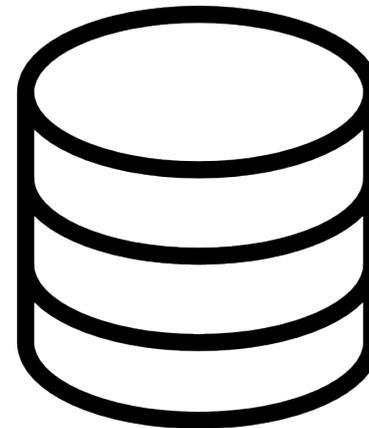
SURVEY

*ecosystem*

# Try on your own Splunk data



*users*

How many users are there? How are user arrivals distributed? What is the user interarrival rate? What are the basic properties of user sessions? What is the average session length? How many queries are there per session?
How do tasks vary by user?

How do transformation frequencies vary with respect to data source? Are some tasks more common in certain contexts than others? How do pipelines vary with respect to data source? How similar is usage with respect to data source type?



*data*

# Analysis challenges

- Messy, complex, with dirty provenance
- Discussion on improving logging: `www.eecs.berkeley.edu/~alspaugh/papers/alspaugh-idea2014-final.pdf`
- No query input/output, data, metadata, system details
- Very large command space with skewed distribution
- Representation mismatch between query language and analysis questions due to functionality overloading
- Caution: research-quality code!
  - query parser: `salspaugh.github.io/splparser`
  - query utilities (including command taxonomy): `salspaugh.github.io/queryutils`
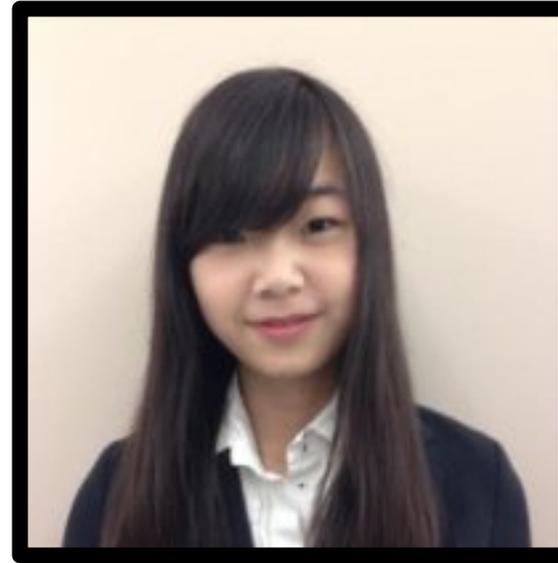  - paper code: `github.com/salspaugh/lupe`
- **`alspaugh@eecs.berkeley.edu`**

# Analysis summary

- 87% of log analysis: data cleaning and troubleshooting
  [40%] **filter**: field-value match, string match, selection
  [ 15%] **add a column**: that is function of other columns (string manipulation, categorization, simple arithmetic, datetime conversion, etc.)
  [ 15%] **aggregate**: count, average, max, min
  [  9%] **rename** columns
  [  8%] **project** columns (filter by column)
- Use cases: troubleshooting, security, report generation, monitoring, business intelligence
- Log analysis is not just for IT: management, product, marketing, sales looking at logs

# Berkeley CS: graduating May 2015!



`jessica.lin@berkeley.edu`



`bettychen824@berkeley.edu`

## LISA Lab Office Hours: 4:00-4:45 Thursday

Questions, comments, code requests:

# alspaugh@eecs.berkeley.edu