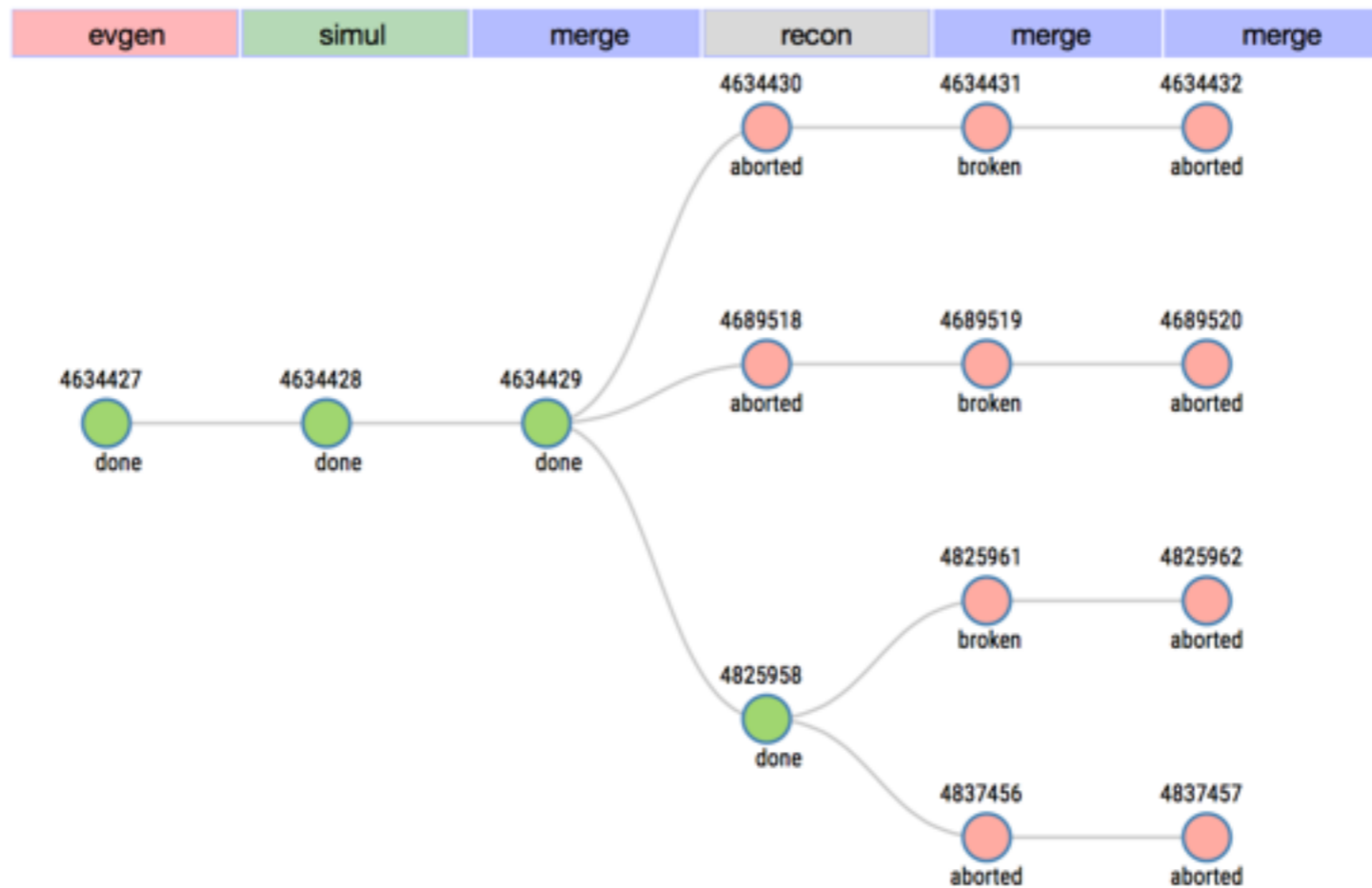# BigPanDA monitor

Sergey Padolski
BNL

# Status

- Recently: More data, more plots, more pages, more fixes

- Feedback: monitor is great, but performance is a problem. + More links to offline analytics is needed.

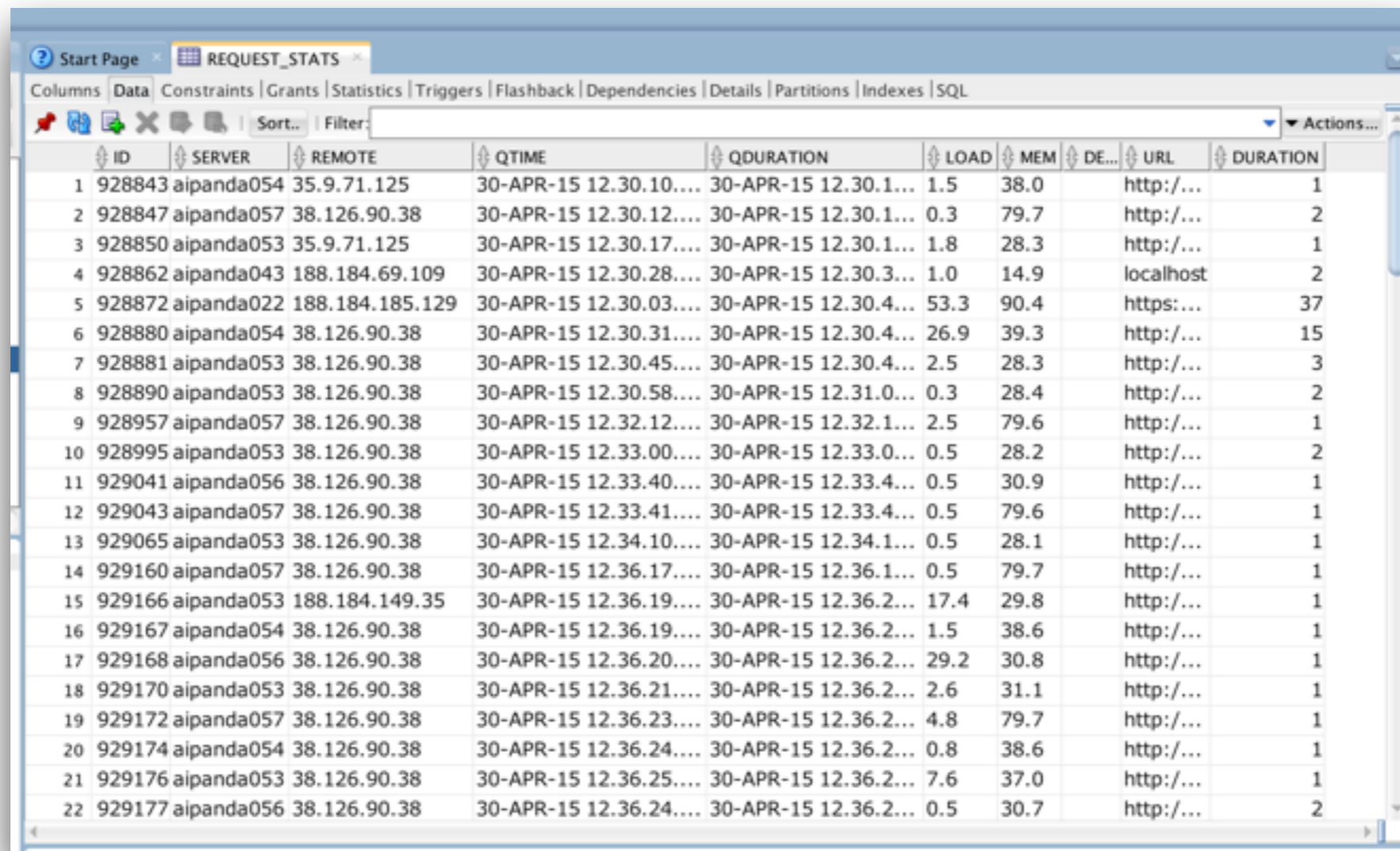- Near future: performance reworking, numbers revision.

# Status



**Chain of Tasks**

New feature from Kurchatov Institute

# Task List

| ✓ | Title | Short description | Who provided input | Assigned priority | Time estimation | Progress Milestones | ID | Comments |
|---|-------|-------------------|--------------------|--------------------|-----------------|---------------------|-----|----------|
| x | ~~ATLASPANDA-255~~ | ~~ATLASPANDA-255~~ | ~~David Cameron~~ | ~~Major~~ | | ~~14/04/2016 registered~~ | ~~104~~ | |
| | Add task age histogram and events values to taskList page | Add task age histogram and events values to taskList page | Andrej Filipcic/Ale/Tatiana | Normal | | 15/04/2016 registered | 105 | Tatiana |
| | Add nolimit button | Add nolimit button | Ale | Normal | | 15/04/2016 registered | 106 | |
| x | ~~plots on the ES task info page~~ | ~~plots on the ES task info page~~ | ~~Rod~~ | ~~Normal~~ | | ~~15/04/2016 registered~~ | ~~107~~ | ~~Tatiana~~ |
| | maxpss per core and HS06s plots | add maxpss per core and HS06s plots to taskInfo pages (both ordinary and es) | Rod | Normal | | 15/04/2016 registered | 108 | Tatiana |
| | ATLASPANDA-256 | ATLASPANDA-256 | Alessandra Forti | Minor | | 15/04/2016 registered | 109 | |
| x | ~~ATLASPANDA-257~~ | ~~ATLASPANDA-257~~ | ~~Enrico Tassi~~ | ~~Blocker~~ | | ~~16/04/2016 registered~~ | ~~110~~ | |
| | | | | | | | 111 | |
| | | | | | | | 112 | |
| | | | | | | | 113 | |
| | | | | | | | 114 | |
| | | | | | | | 115 | |

https://docs.google.com/spreadsheets/d/1cSz2j4iI-zIapLc8Ap9uOd4riODo0hVTHCJ9DYhWsXM/edit#gid=0

# User Behavior Analysis



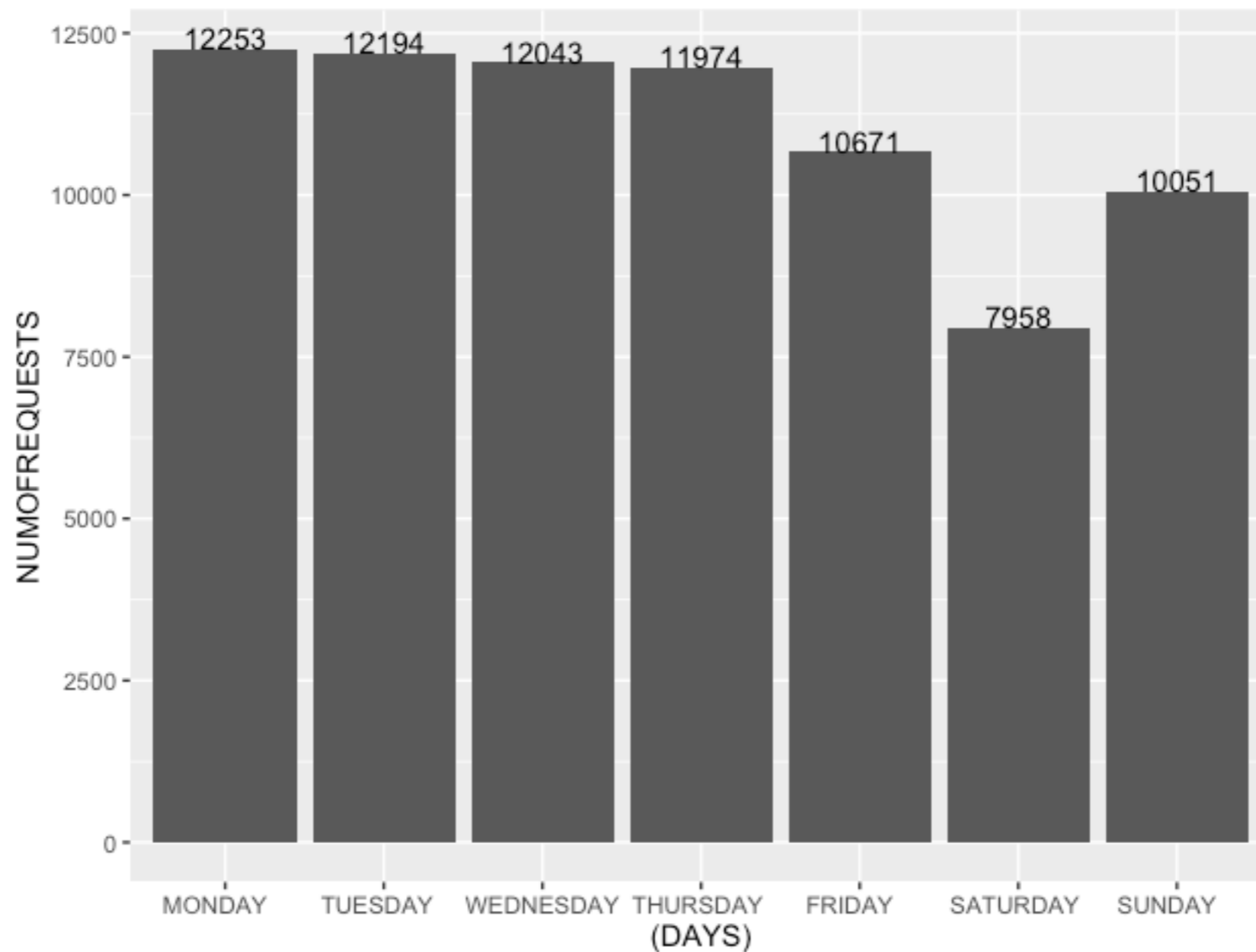New analytic views were recently developed
(SESSIONANALYSIS1, SESSIONANALYSIS2 )

# User Behavior Analysis

# User Behavior Analysis

Session is user actions sequence with timeout less than 30 minutes

# User Behavior Analysis

| | |
|---|---|
| task | 17.97 |
| jobs | 17.32 |
| main | 11.46 |
| dash | 11.1 |
| tasks | 11 |
| user | 8.17 |
| errors | 5.3 |
| sites | 3.76 |
| users | 3.26 |
| site | 2.07 |

Access impact (%), only non cached pages

# User Behavior Analysis

# User Behavior Analysis

# User Behavior Analysis

# User Behavior Analysis

# User Behavior Analysis

# User Behavior Analysis

# System performance

# Bottleneck analysis

Jobs page

1. Retrieve jobs (4 active + possible archived)

2. If (dynamic conditions) retrieve JediJobRetryHistory, reconstruct chain

3. if user request conditions, go to JediDatasetContents

4. if (dynamic conditions) go to JediDatasets

5. go to Filestable4

6. if (dynamic conditions) go to FilestableArch

7. if (dynamic conditions) go to Filestable4

8. if (dynamic conditions) go to FilestableArch
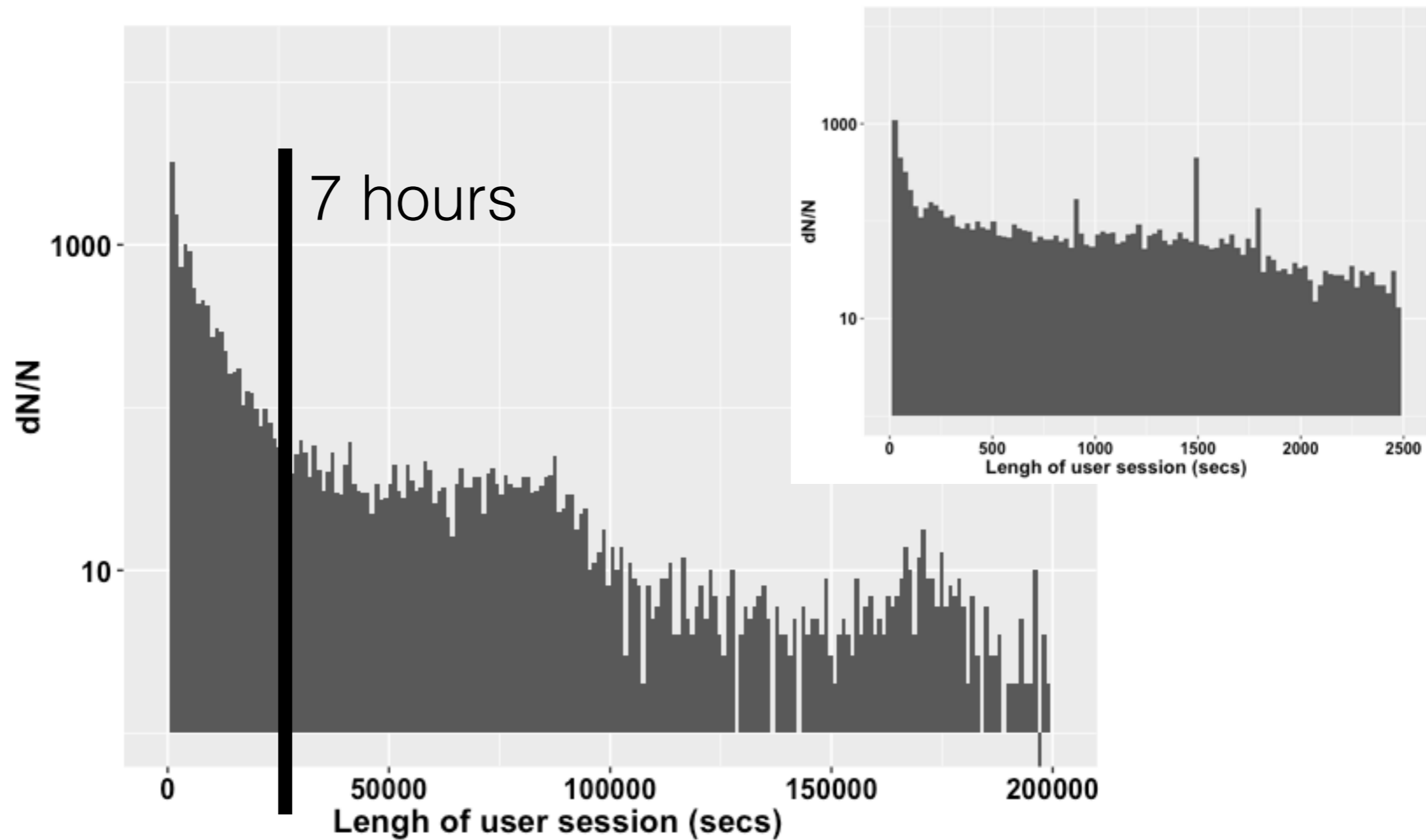
9. if (dynamic conditions) go to JEDI_EVENTS

~ 7 different paths of DB access branched logic

**~13 individual DB queries, ~30 MB retrieved from DB,**

**DB data transfer ~ 1M/sec, 98% of time posix.waitpid**

# Requirements for the next iteration

Significantly reduce frequency of data access operations and amount of transferred bytes

Most of processing logic should be implemented nearby the data

# Tools we (can) try

- Google BigQuery

- Indexed data storage on top of Apache Lucene (Elastic Search, Solr)

- Spark

# Google BigQuery

- Pros:

  - Developed by Google (expecting reliability, high performance)

  - Highly scalable

- Cons:

  - UDF is a fresh technology (less than year - small knowledge base)

  - Critical part of monitoring is going to proprietary technology, outsourcing to external company, on shared resource.

  - Every user query costs money (rough estimation: few cents for each processing a large table). The more complex algorithm of the data processing - more it costs. New metric - cost of request. New error - billingTierLimitExceeded.

  - Network connectivity between CERN and BigQuery warehouse become critical part of the ability of monitoring to deliver fresh data to the end user.

  - Deploying monitor to another experiments involves accepting proprietary technologies, payments to google for processing, storage, import, etc, additional technology to support

# Apache Lucene solutions (ES, Solr)

A new prototype was developed in 2015:

- Cheap and quick approached - architecture has been kept

- A class to translate Django Query to Solr language developed

- Performance didn't increase significantly and R&D stopped

# Apache Spark

- Pros:

    - Highly scalable (Tencent example with 8000 nodes and 150PB +1 PB a day)

    - Easily deploying on virtual or physical clusters

    - Advanced underlying technologies (full Java support)

    - Native integration with R (one of the most popular analytic and machine learning tool)

    - Open Source

- Cons:

    - Have to maintain cluster or buy service

# Spark prototype

Goal: Check is technology stack applicable to monitoring, assess speed

Scope:

- Program skeleton to run embedded Spark server and «user queries»

- Establish mutual data stream

- Perform preprocessing

- Perform fast selection

- Perform complex data aggregation

# Spark prototype

- Stand alone Java application

- Few simultaneous version of data

- As new version of data is imported, it become actual and programs address queries to this dataset

- Only skimmed data is transferring (thanks to ORA_ROWSCN)

- Preprocessing is applied

- Speed of data aggregation is 1000 rows / core / sec

# Possible scenarios

- Big query:
  - Initial technology assessments, debugging, data streaming, delivery to web server (1-3 weeks)
  - Implementing monitor real case (2-4 weeks)
  - Assessment speed, cost, reliability, working on development plan (1 week)

- Spark
  - Setup the development spark cluster (~1w).
  - Setup the uninterrupted data stream (Oracle-Spark) (~1-2 w).
  - Implementing backend for Errors Page (~1-1.5 months).
  - Implementing backend for DashBoards (~2-2.5 months).
  - Implementing backend for Jobs (~3 months).
  - Implementing backend for Tasks page(~3.5 months).
  - Other pages which cooked with huge delays (~4.5 months).

  The best-case estimate: 3 months, the most likely estimate: 4.5 months, the worst-case estimate: 6 months, weighted average (E = (best + 4mostlikely + worst) / 6) = 4.5
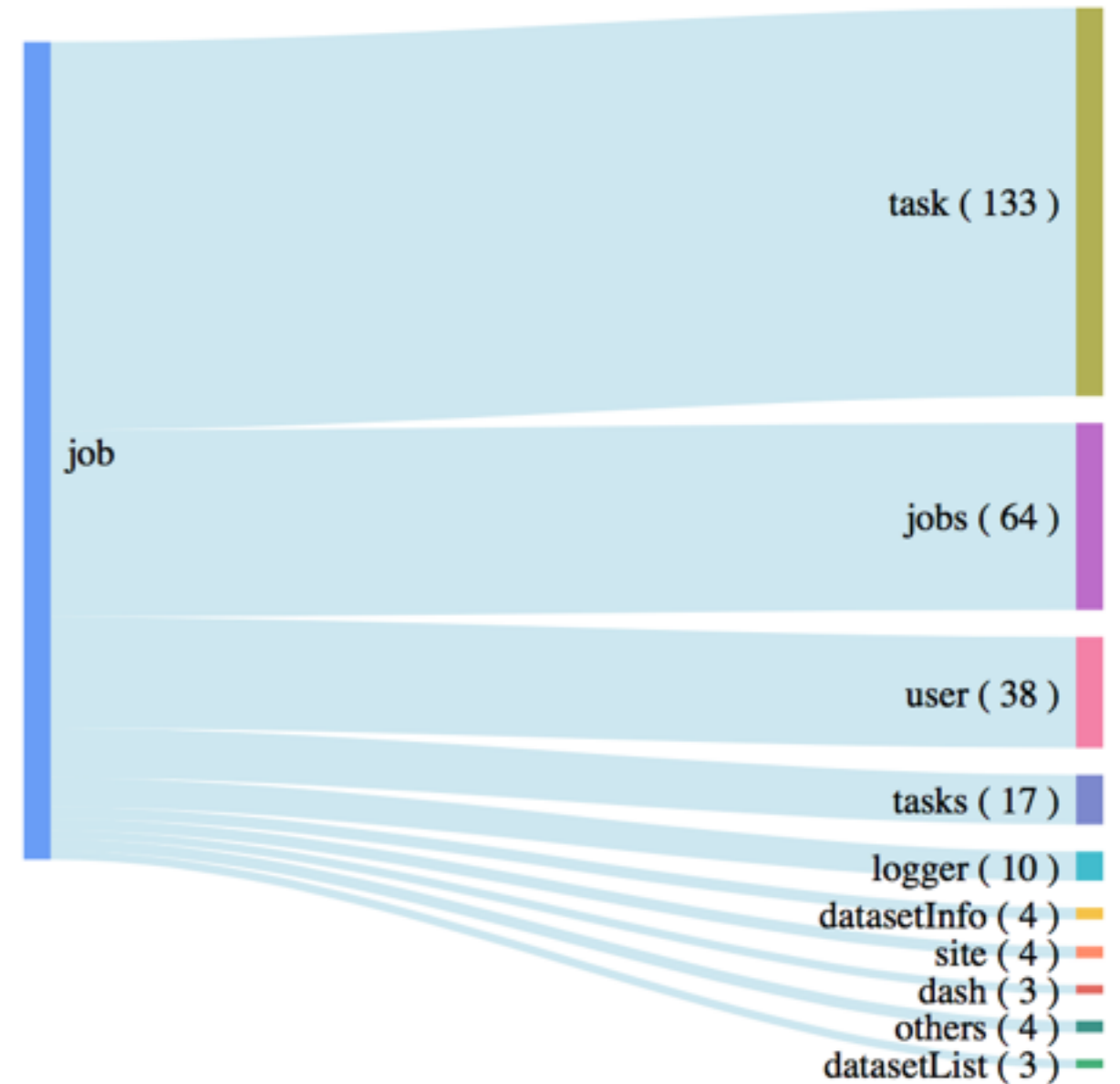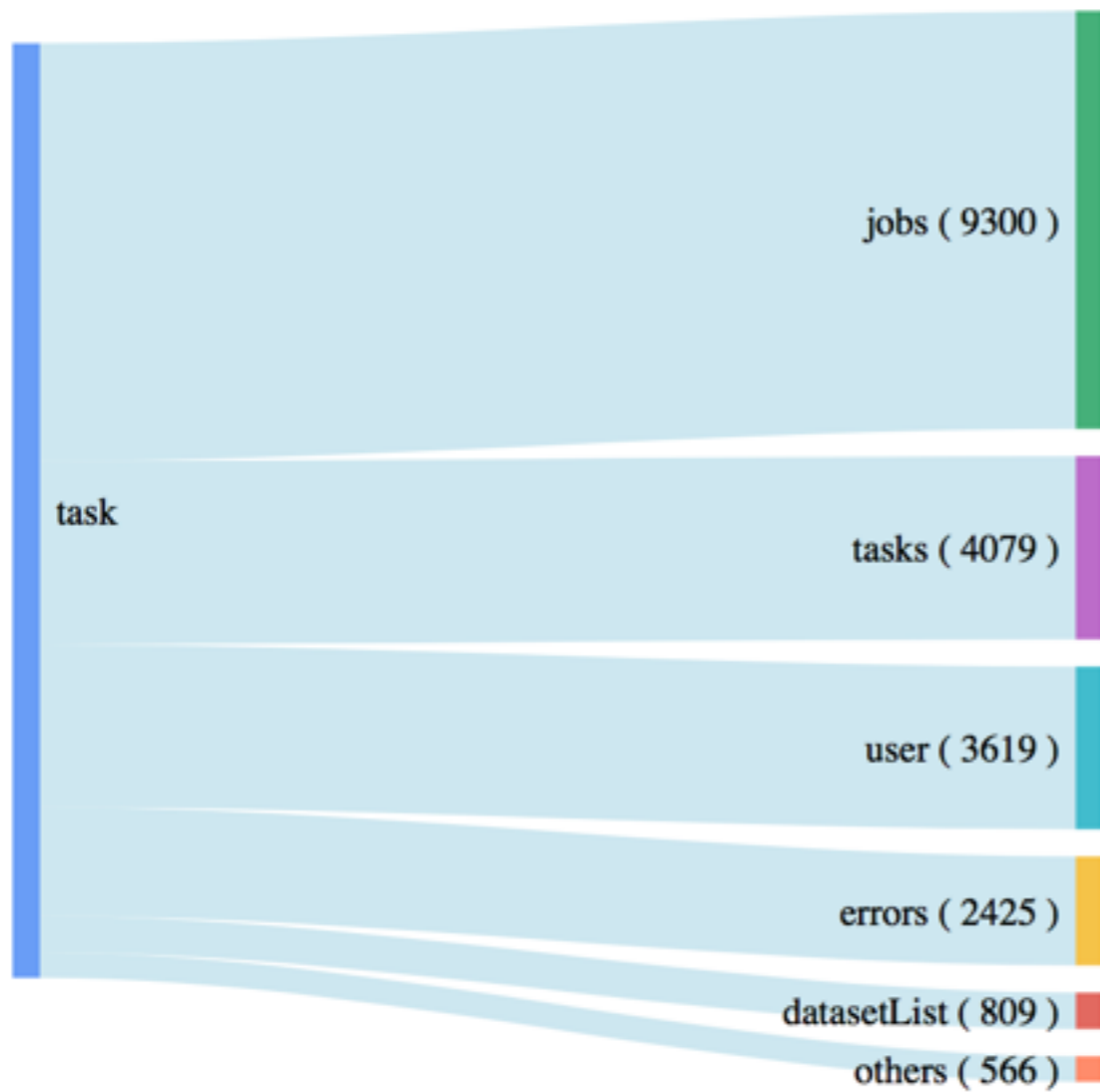
- Elastic search
  Development schema transformation, preprocessing, start importing data, working case prototype (1-2 months)
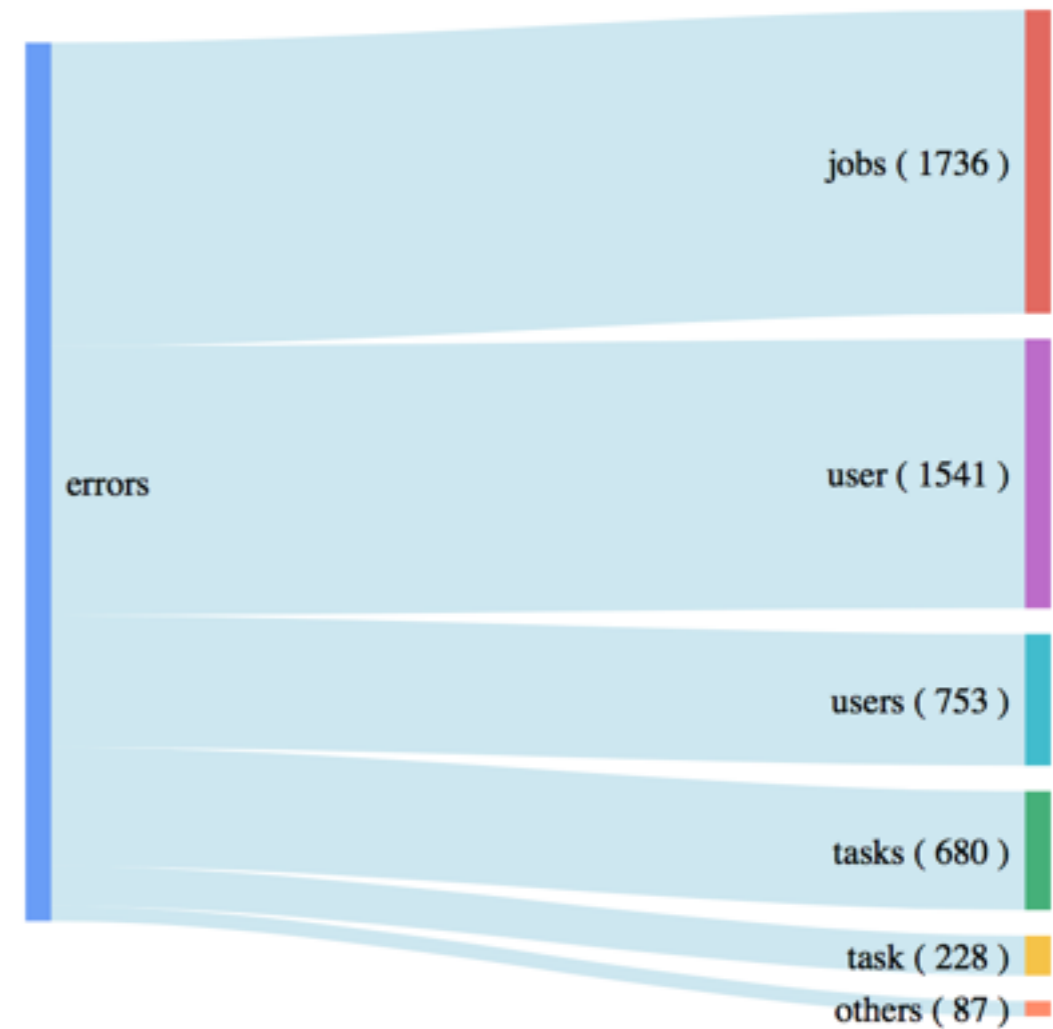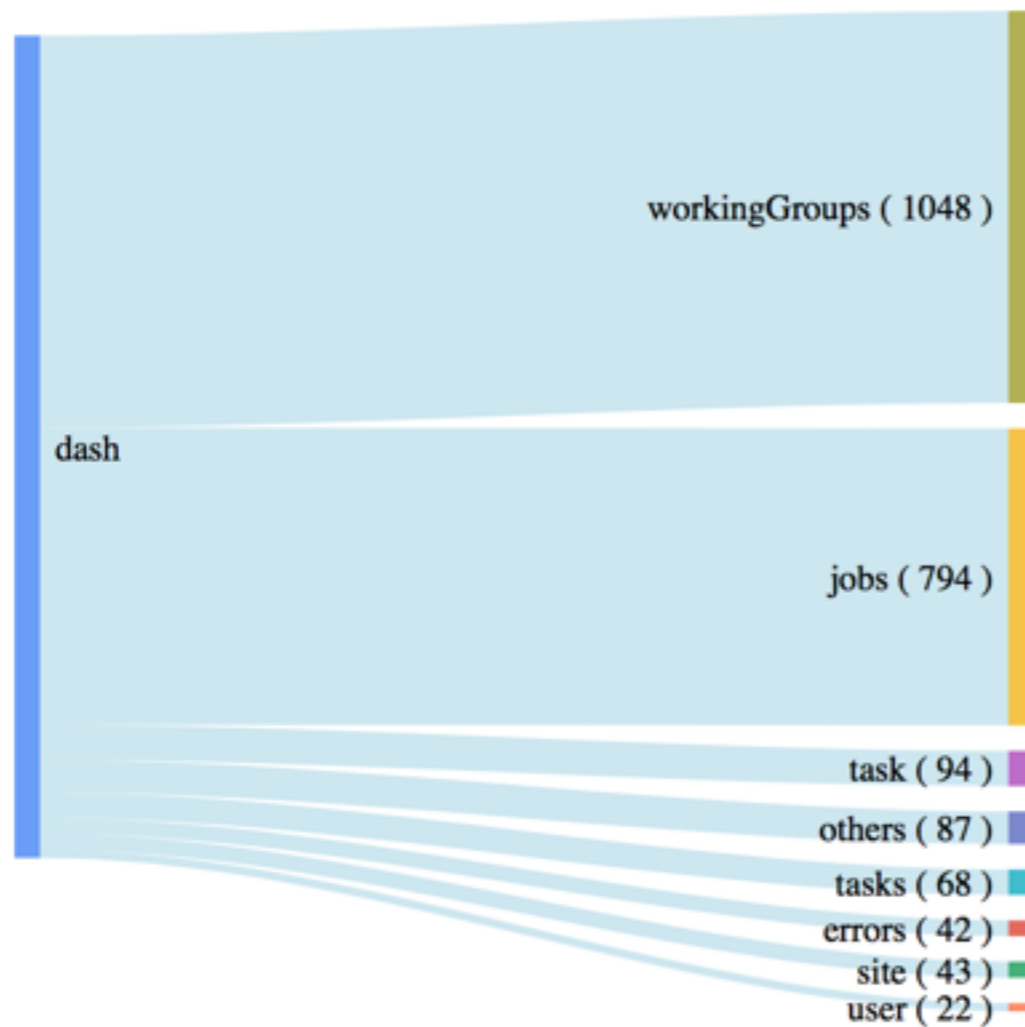
# Backup

# User Behavior Analysis

# User Behavior Analysis

# BigQuery cost estimation

Lets say we update a biqquery warehouse once a 10 minutes. That correspond to 220M of data transfer only for the jobs table. $0.01 (per 200M) x 6 (times per hour) x 24 (hours per day) x 365 = 525.6$. Just for refreshing jobs table

In the most simple case when we retrive jobs from jobsarchine table processing only MODIFICATIONTIME field:
10GB*0.001*5 = 5 cents per one query to jobsarchive.