



Studies of Big Data metadata segmentation between relational and non-relational databases

**M.Golosova (NRC KI), G.Dimitrov (CERN), M.Grigorieva (NRC KI),
A.Klimentov (BNL), M.Potekhin (BNL), E.Ryabinkin (NRC KI)**

**National Research Centre «Kurchatov Institute»
CERN
Brookhaven National Laboratory**

CHEP2015, Okinawa, Japan, April 13-17 2015



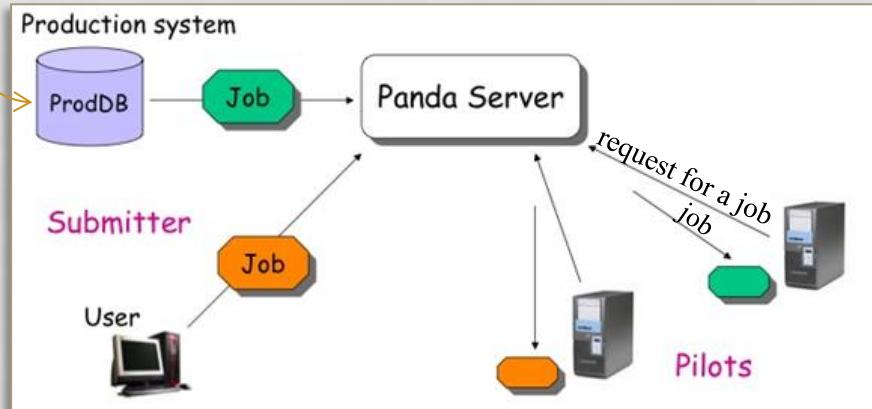
Outline

- **PanDA WMS**
 - Metadata
 - Monitoring
 - Data specification
- **NoSQL as archive**
 - Precalculation for speedup
- **Study**
 - Test set-up
 - Scaling
 - Precalculation effect
 - NoSQL archive performance



PanDA WMS

Metadata storage

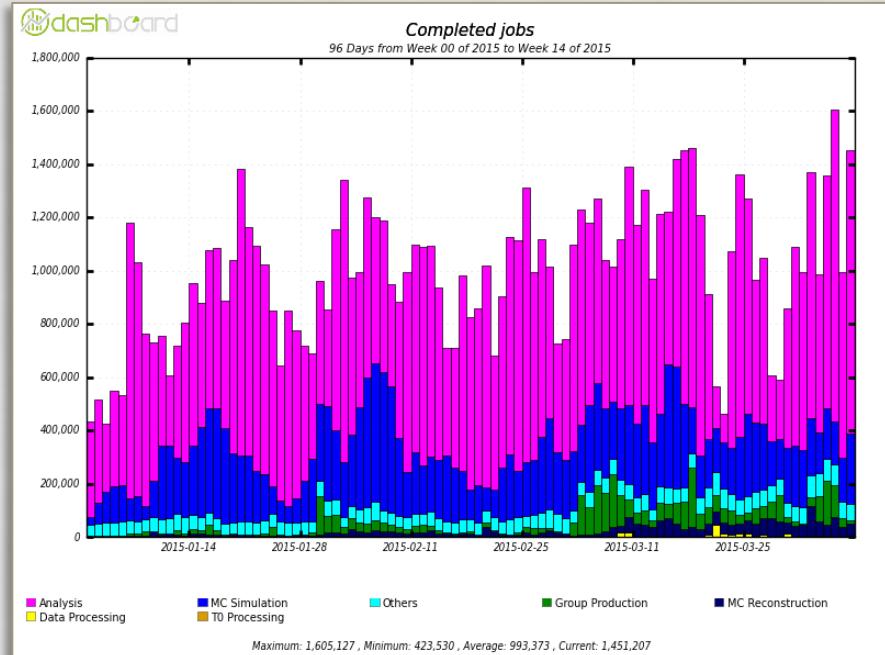


courtesy picture from T.Maeno

Metadata analysis and monitoring: PanDA Monitor

The screenshot shows the PanDA Monitor interface. At the top, there's a navigation bar with links like ATLAS PanDA, Dash, Tasks, Jobs, Errors, Users, Sites, Incidents, Search, Admin, Prodsys, Services, VO, Help, and a status bar indicating 'Built 16:44, cache 3'. Below the navigation is a table titled 'Job details for PanDA job 2444524017'. The table has columns for Pandaid, Owner, Request Task ID, Status, Created, Time to start d:h:m:s, Duration d:h:m:s, Modified, Cloud Site, and Priority. One row is shown: 2444524017, Andrey Minaenko, 10739, defined, 2015-04-08 16:43, 0:00:44, 04-08 16:44, DE ANALY_FREIBURG, 834. Below the table is a message about job name and datasets. The bottom half of the screen shows a 'Site error summary' table with columns for ANALY_CONNECT_SHORT, count, and error details. The table lists various errors such as 'exe:1137 Put error: Error in copying the file from job workdir to localSE' and 'jobdispatcher:100 lost heartbeat : 2015-04-08 13:35:18'.

Scale of the challenge: up to 1.7M completed jobs per day



Related CHEP talks: T.Maeno,
[The Future of PanDA in ATLAS](#)
[Distributed Computing](#) (CHEP ID 144)



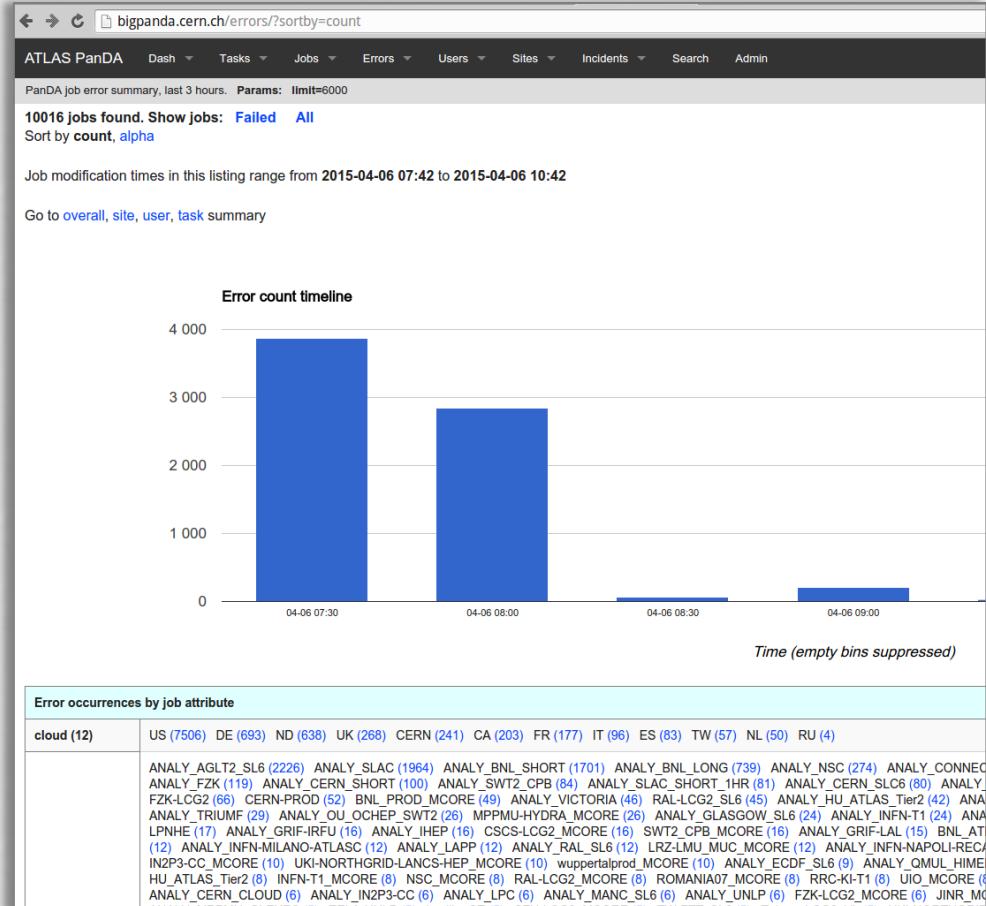
PanDA monitoring

- **PanDA Monitor**

Web-based analytical interface providing information about jobs and tasks within the system

- **Page «Errors»**

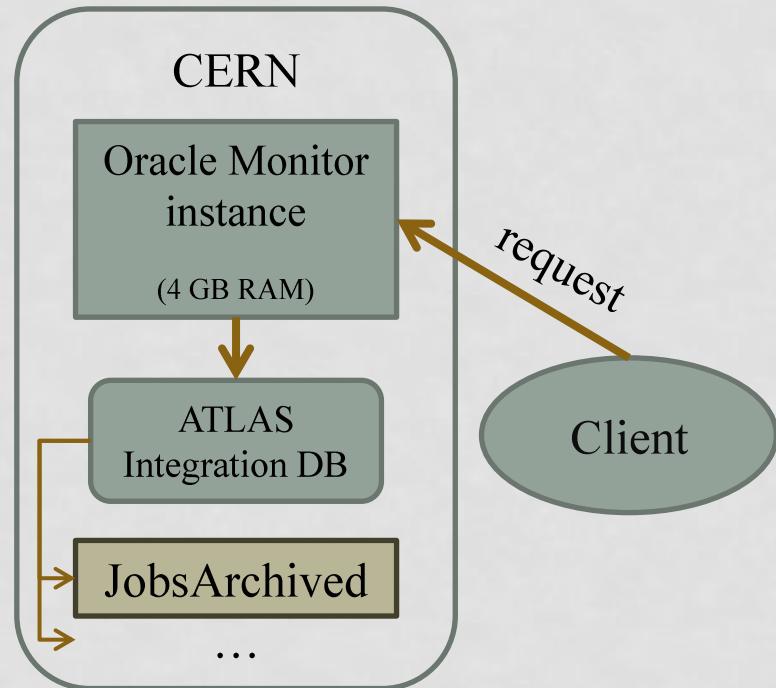
Page of PanDA Monitor providing information about job errors that have occurred in the system



<http://bigpanda.cern.ch/errors/>



Job errors monitoring

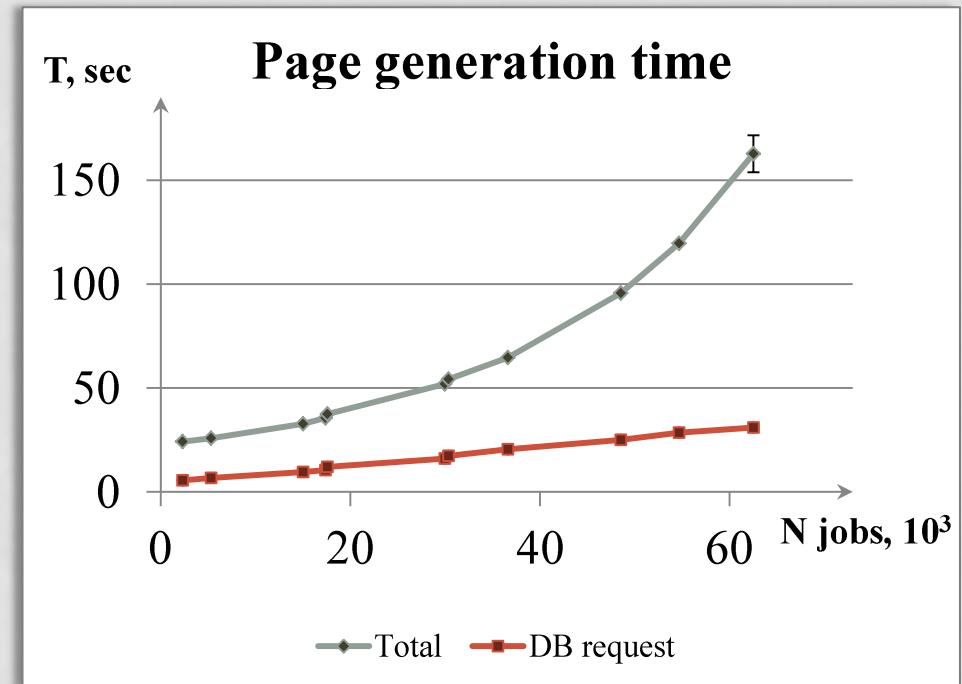


1) Get list of jobs
from DB (according
to request
parameters)

⇒ 2) Summaries by
parameters (overall,
computing site, user,
task)

⇒ 3) Get additional
information
from DB

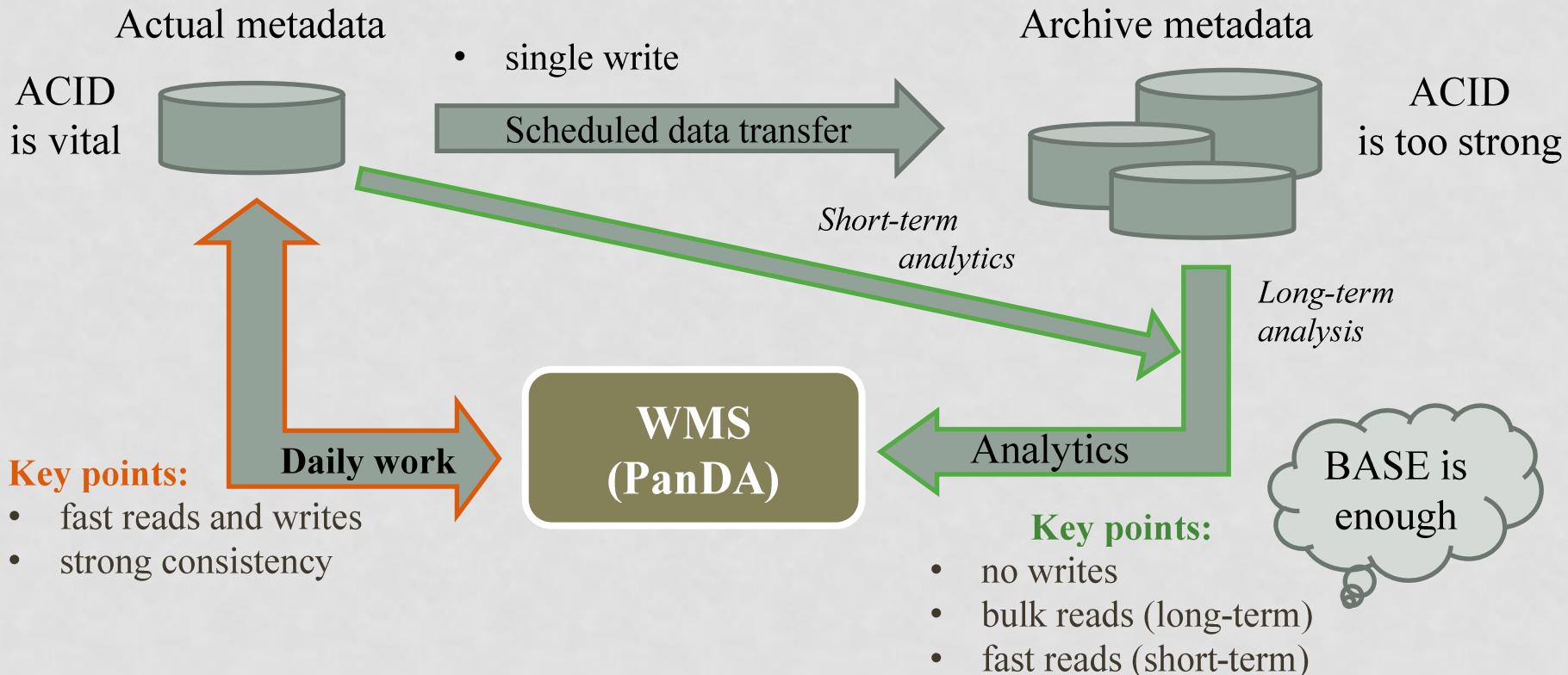
⇒ 4) Generate output





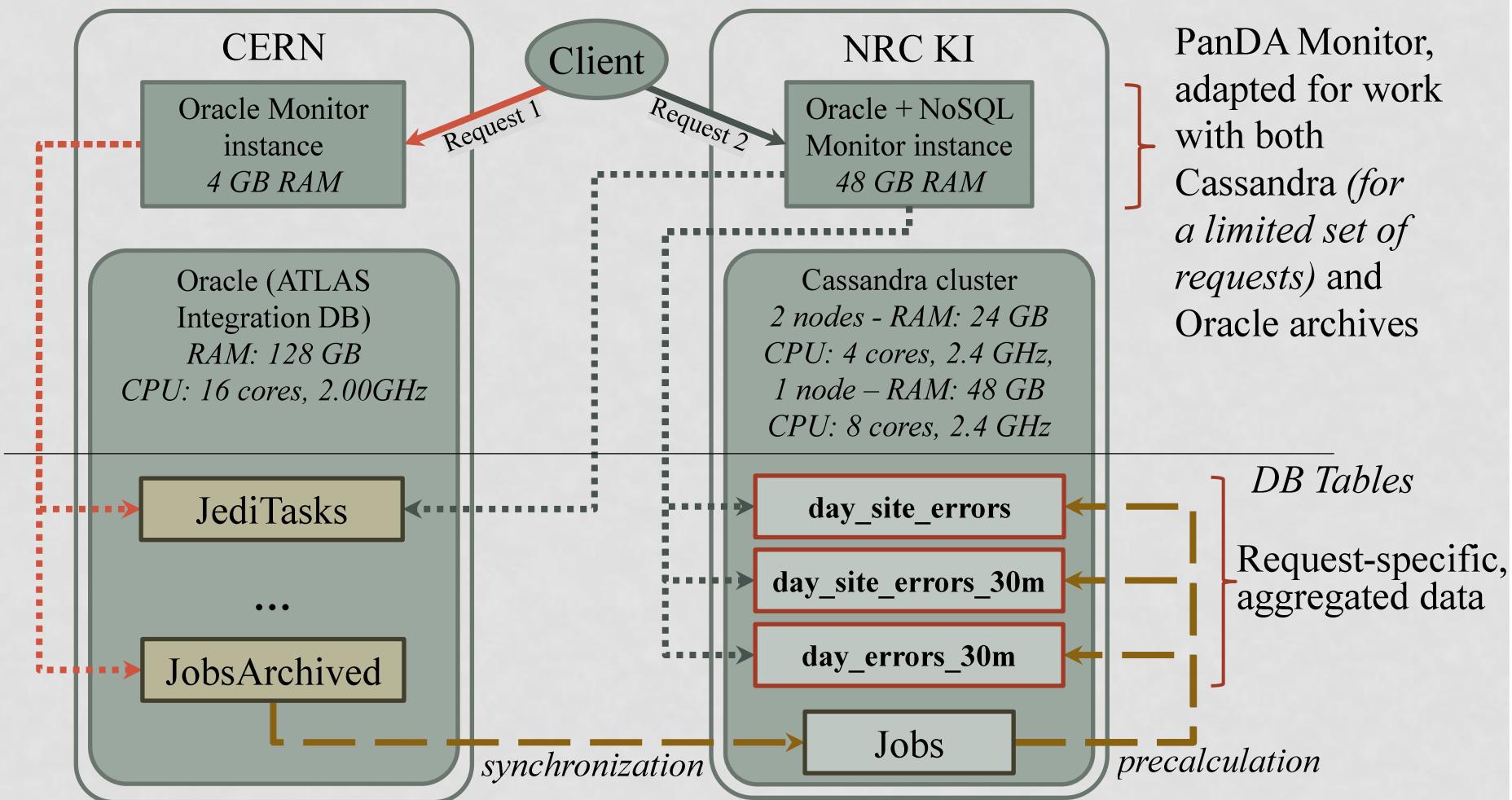
Specificity of the metadata storage

- Number of requested records (ATLAS PanDA Archive: ~900M of jobs)
- RDBMS as archive backend (ACID standard)



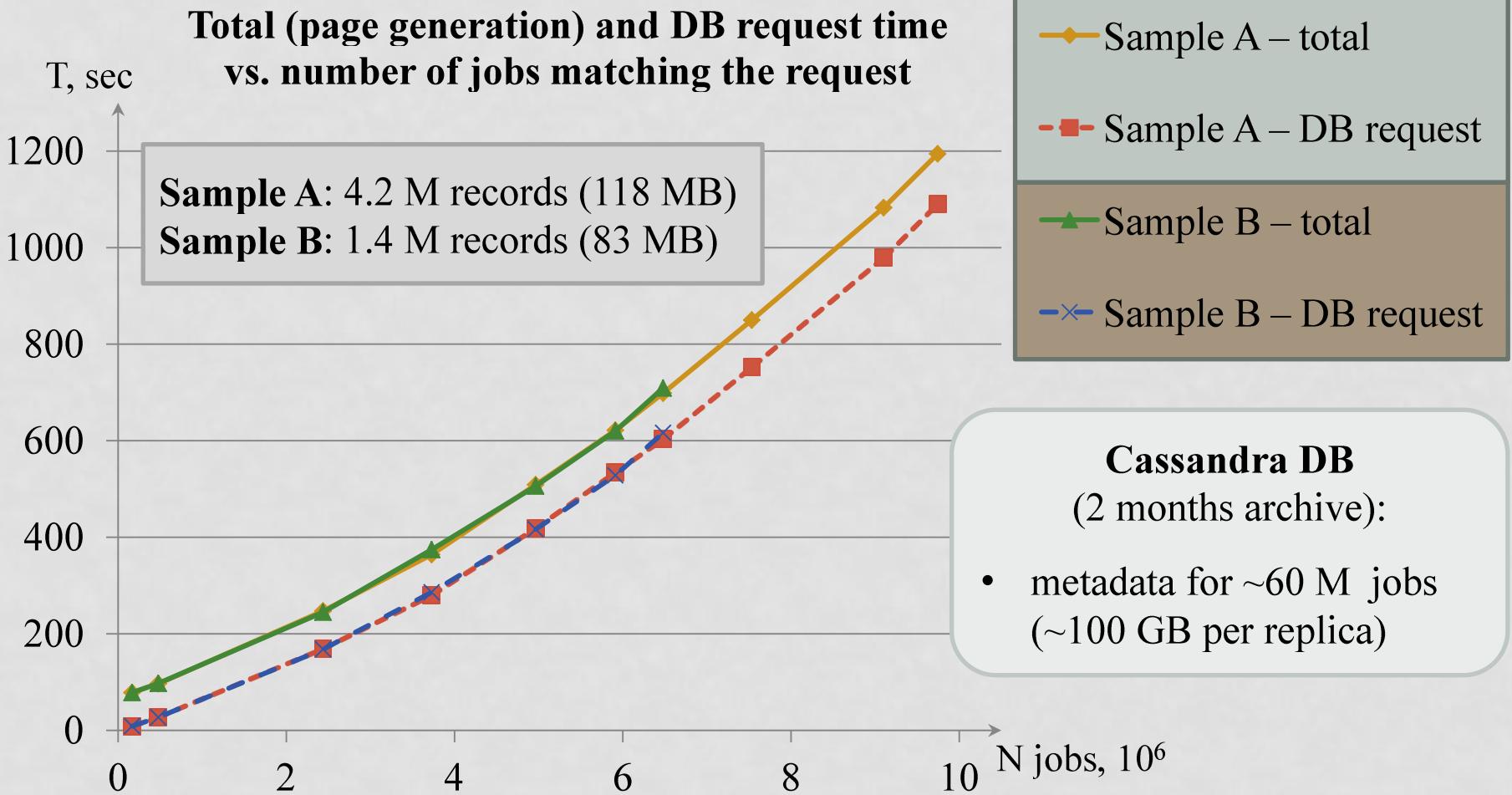


Test set-up



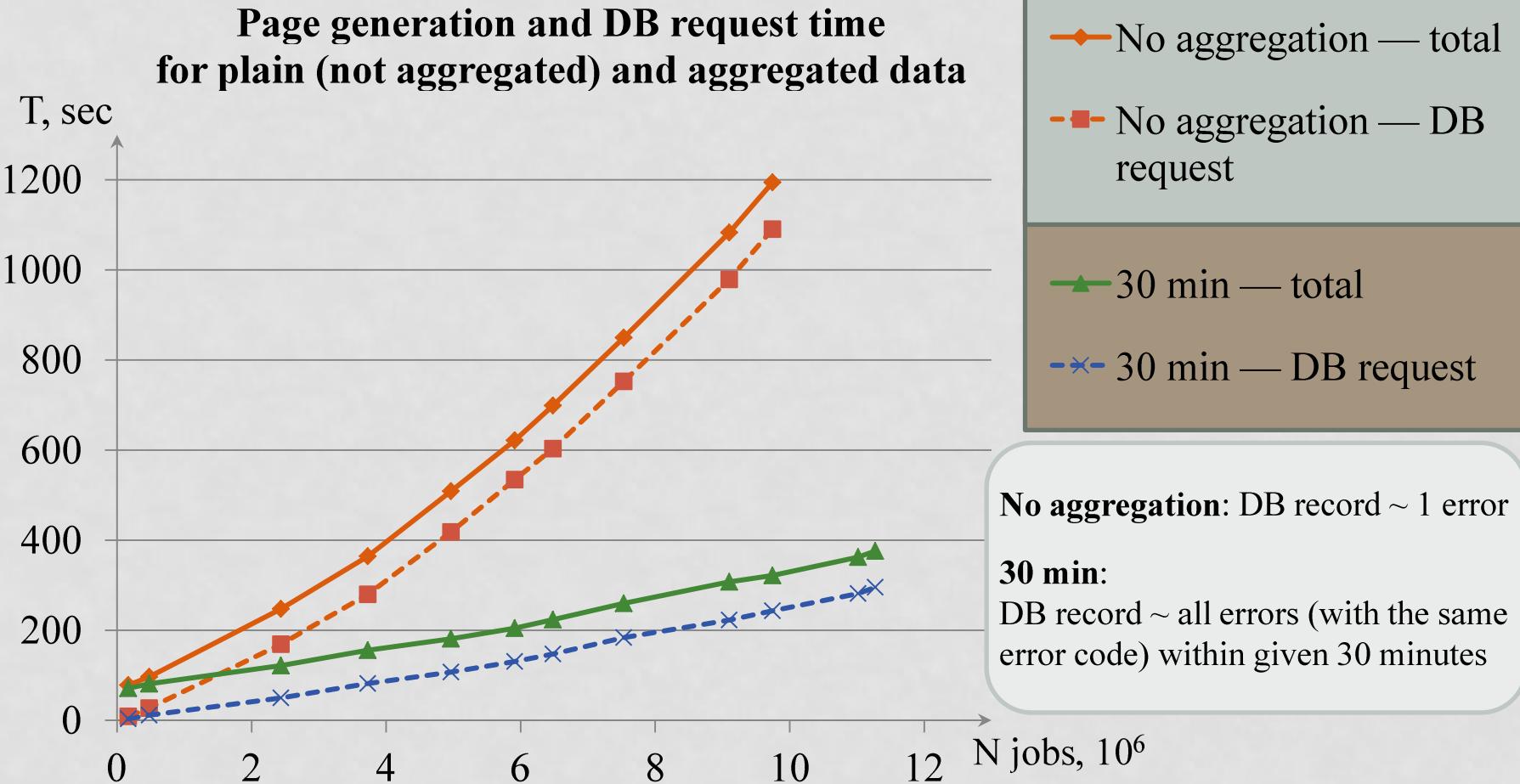


Cassandra scaling test



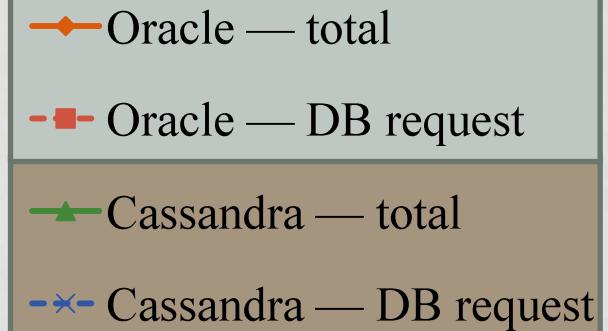
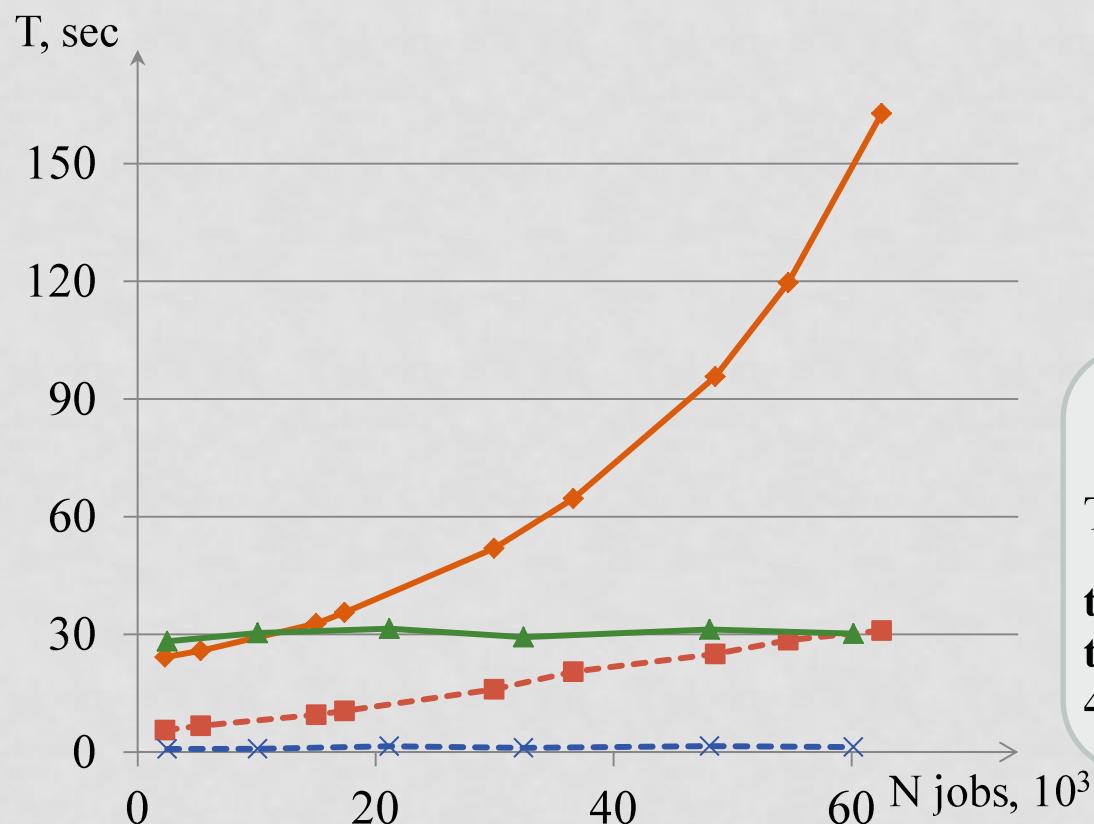


Aggregation effect test





NoSQL Archive performance



Values for Cassandra

$$T = 4 * (t - t_{add}) + t_{add}, \text{ where:}$$

t – total page generation time

t_{add} – not a summary-specific DB requests

4 – number of specific summaries



Summary

- It is hardly possible to perform long-term metadata analysis without any precalculation.
- Replacing RDBMS with NoSQL in archive part of metadata storage can be used to improve availability of historical data.
- Prototype of NoSQL (Cassandra) archive was created and tested on a 2-month slice of metadata from ATLAS PanDA Archive.
- First results look promising. NoSQL archive will be extended to confirm this belief in new tests.
- Adaptation of PanDA Monitor for work with NoSQL archive will be continued.



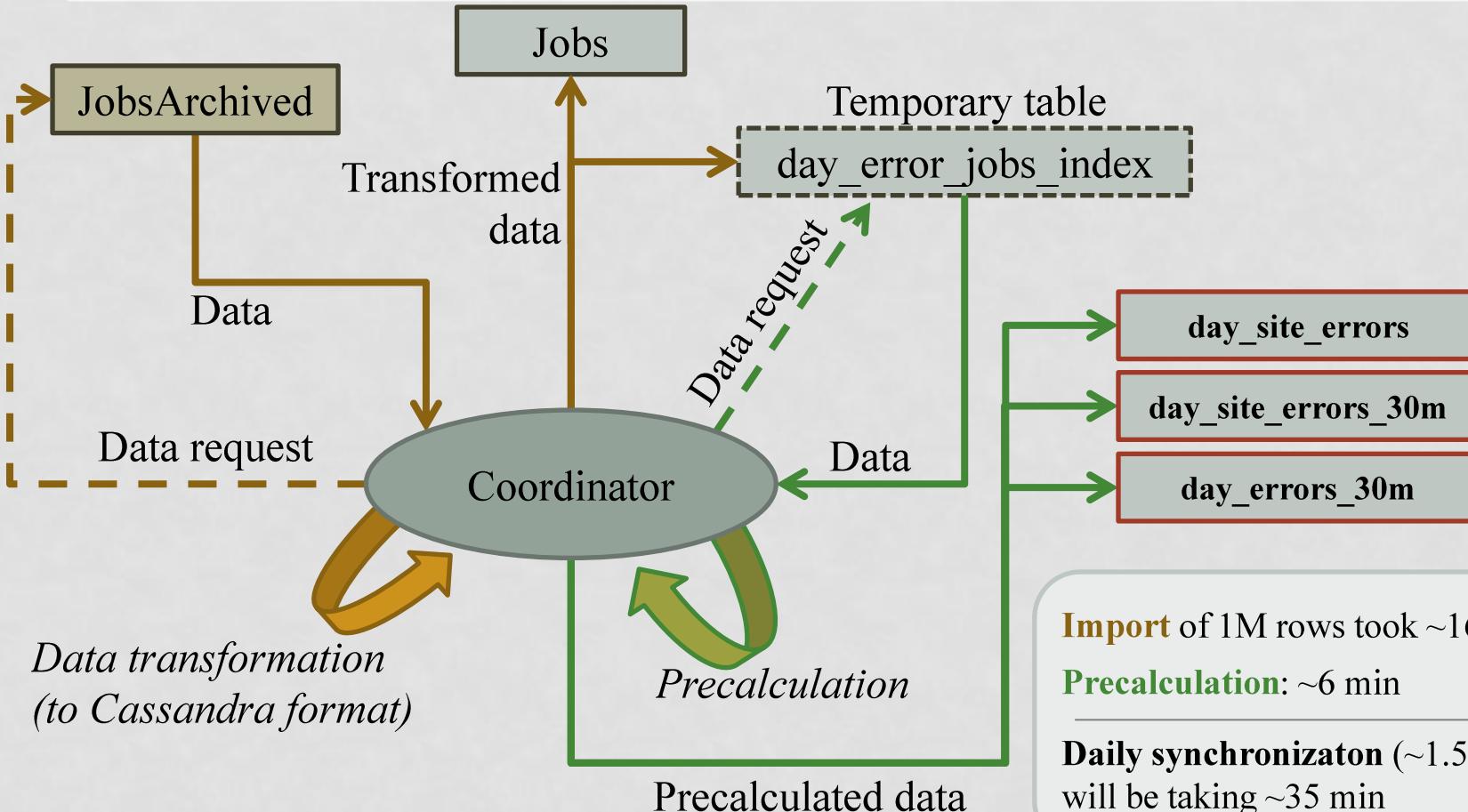
Acknowledgment

- This talk drew on presentations, discussions, comments, input from many our CERN, ATLAS, NRC-KI Colleagues, thanks to all people working in PanDA project
- This work was funded in part by the Russian Ministry of Science and Education under Contract N14.Z50.31.0024



Back up slides

Synchronization & Precalculation



Import of 1M rows took ~16 min

Precalculation: ~6 min

Daily synchronization (~1.5M jobs)
will be taking ~35 min



Cassandra Column Families

Primary Key

- Partition Key
- Clustering Key

Columns

- Common
- Counter

day_site_errors_cnt_30m

date	computingSite	base_mtime	errcode	errdiag	err_count	job_count
------	---------------	------------	---------	---------	-----------	-----------

day_site_errors

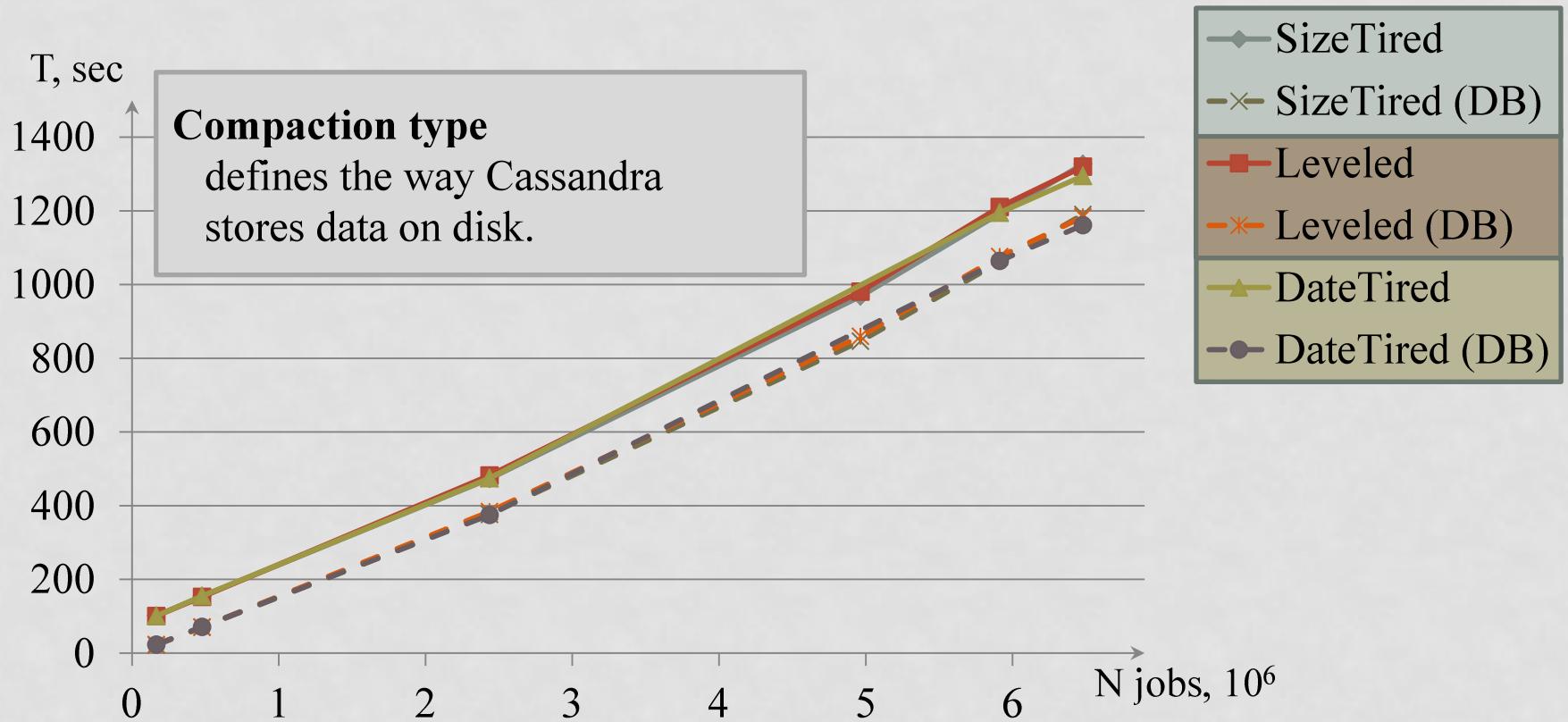
date	computingSite	errcode	pandaid	errdiag
------	---------------	---------	---------	---------

day_errors_30m

date	base_mtime	count
------	------------	-------



Compaction types





Aggregation effect test (DB records)

