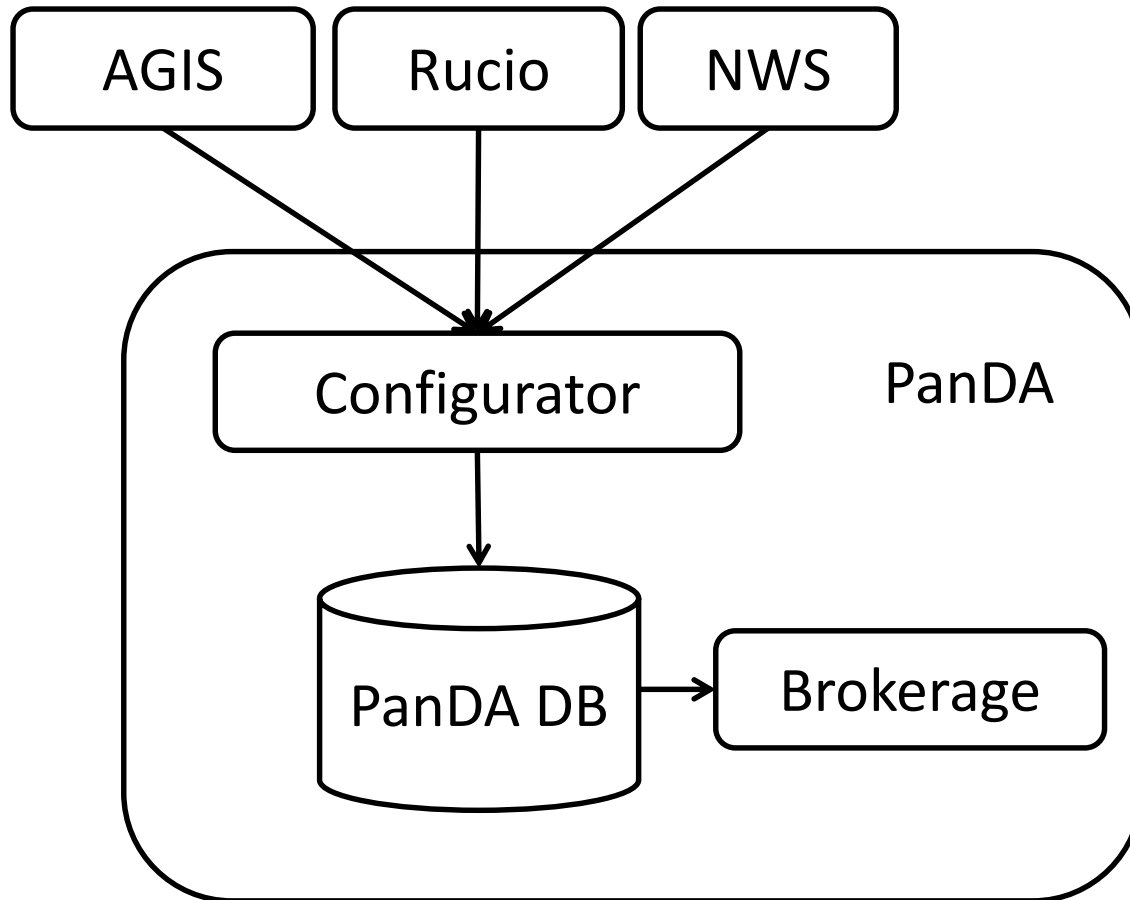# PanDA Configurator and Network Aware Brokerage

Fernando Barreiro Megino, Kaushik De, Tadashi Maeno

14 March 2015, US ATLAS Distributed Facilities Meeting, Clemson University

# Configurator: overview



- PanDA agent running every 30 minutes collecting information useful for brokerage, in particular regarding **WORLD** cloud migration

- Adding progressively new sources as we see the need

# WORLD cloud 101

WORLD is the evolution of MCP. Tasks are not confined to their cloud anymore.

**Nucleus**:
- PanDA task brokerage will assign tasks to Nuclei (=T1s and selected T2s).
- The output will be aggregated in the Nucleus.

**Satellites**:
- Run jobs and ship the output to the Nuclei.
- Satellites will be selected for each task, with a maximum of 10 satellites per task.
- Job brokerage will select satellites based on usual criteria (e.g. #jobs in different states, data availability, …)
- Job brokerage will not confine the task to a cloud, but will **increasingly** be based on the network connectivity and transfer queues between the sites.

# Configurator: topology data (reminder)



Base info / Space usage / Downtime info (ATLAS_PANDA.DDM_ENDPOINT)

Base info / Role (ATLAS_PANDA.SITE)

Base info / Locality (ATLAS_PANDA.PANDA_SITE)

# Configurator: static network data

- Configurator agent downloads and processes network information every 30 min from AGIS and NWS. Data is cached in a key-value table in PanDA DB
  - Table structure avoids adding/removing columns every time a new metric appears/disappears
- AGIS closeness matrix: static closeness between each source-destination pair:
  - Value between *1* (good) and *9* (bad) based on the hierarchic cloud model. Examples:
    - T1 → T1: *1*
    - T2 → T1 in same cloud: *2*
    - …
    - T2 → T2 in different clouds: *7*
    - T3 → T3 in different clouds: *9*
  - Special values:
    - -1 to blacklist a channel
    - 0 to define a combined site (in progress)

# Configurator: dynamic network data

- The Analytics platform contains a lot of raw network information (FTS, FAX, PerfSonar). We are working with the Rucio team to get [aggregates per source-destination pair](#) combined with Rucio queue data:

  - #files transferred in last 1 and 6 hours (by activity)
  - #files queued (by activity)

  Available in 1st NWS version

  - Throughput according to FTS, based on 1 week data
  - Throughput according to FAX
  - PerfSonar metrics (latency, packet loss, throughput)

  Available since 2nd NWS version (work in progress)

# WORLD cloud: satellite selection in JEDI production job brokerage

1. Filter out candidates with blacklisted AGIS closeness (closeness = -1)

2. Calculate network weight for remaining candidates, combining static and dynamic info

$$NW\_Weight_{static} = 1 + \frac{MAX\_CLOSENESS - closeness}{MAX\_CLOSENESS - MIN\_CLOSENESS}$$

$$NW\_Weight_{dynamic} = 1 + \frac{queued\_bestsource / done6h\_bestsource}{queued\_channel / done6h\_channel}$$

$\Longrightarrow$

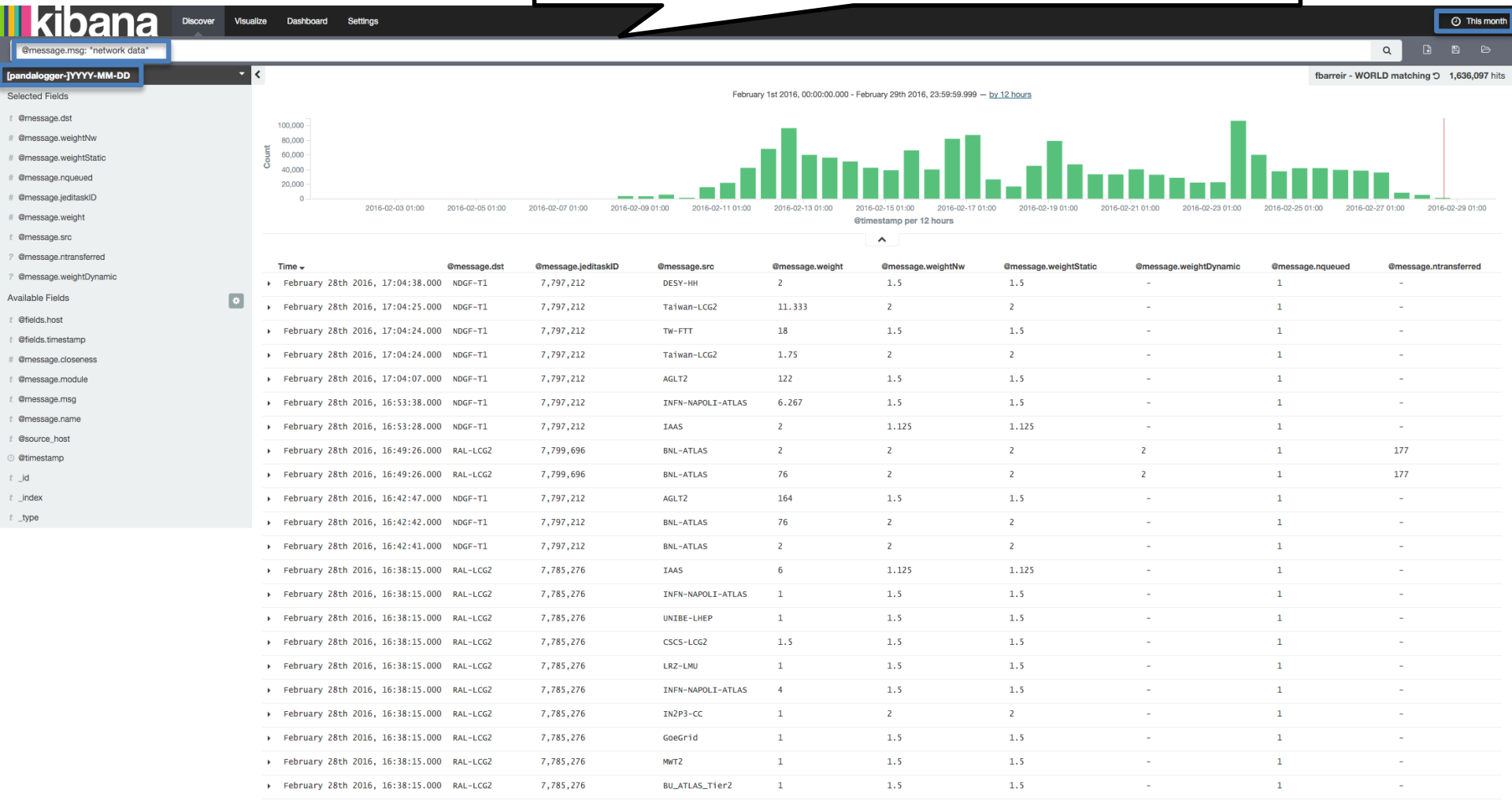$$NW\_Weight = 0.7 \times NW\_Weight_{static} + 0.3 \times NW\_Weight_{dynamic}$$

(exceptions apply)

3. Multiply traditional weight (based on data availability, #jobs in different stages, etc.) by network weight

Currently we are running in passive mode and sending the network brokerage decisions to Analytics platform, so we can tune the network model and algorithm

# Monitoring NW brokerage: ES messages

You can explore the data yourself: https://aianalytics01.cern.ch/
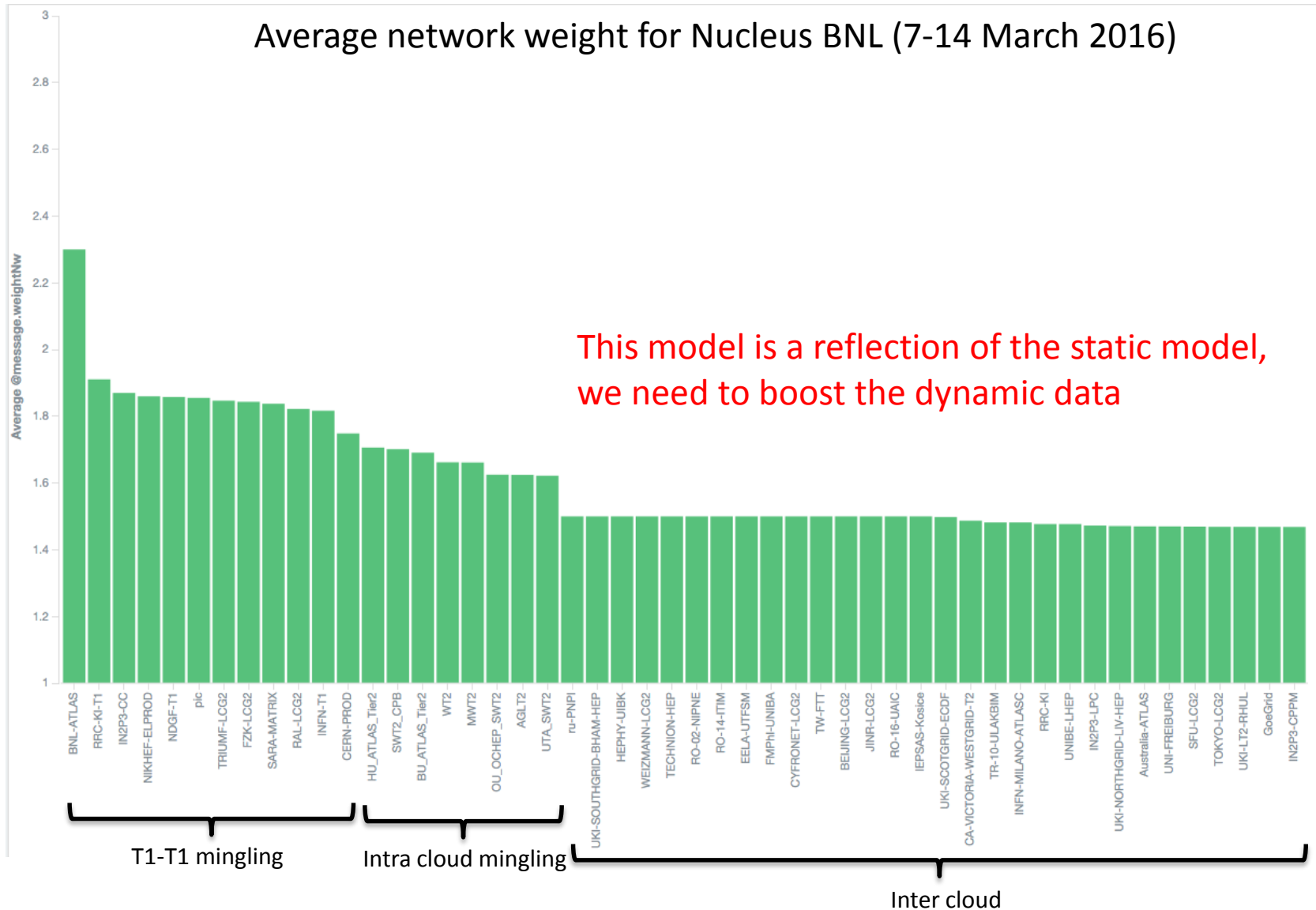


There is no network information for most links over the last 6h, so static closeness will prevail. We should also include metrics that cover longer time periods

# Example: Favorite satellites for nucleus BNL



Average network weight for Nucleus BNL (7-14 March 2016)

This model is a reflection of the static model, we need to boost the dynamic data

# Observations (1)

- This work is fairly new – we started after the Sitges TIM in December
  - We have improved considerably the data transmission model and are starting to use network data for a very broad case in PanDA
  - **We have powerful building blocks, but we need to get them operational and tune the algorithms**
  - We need a good data analysis of all the available information and recommendations (WIP by Mario et al.)

- AGIS closeness is too much of a reflection of the MONARC model
  - IMHO the data should reflect the reality, not a theoretical, obsolete model
    - Simple, semi-static classification?

- Aggregation of data from NWS is work in progress

# Observations (2)

- Verify, activate, improve algorithm for nuclei-satellite matching
  - Start using the second version of aggregated data, containing more info than nqueued and ntransferred:
    - FTS Mbps over last week to have dynamic data over longer period
    - PerfSonar data
  - Boost the dynamic data
  - Analyze if the network weight should be stronger
- Extend to other Rucio and FAX use cases

# Other possible network brokerage use cases

- Network weight for input file transfers (AKA Rodney Walker's case):
  - Input data is in site A and B, but sites are busy
  - Sites C and D are free, but don't have input data
  - → Consider network for the brokerage
- FAX network weight for Event Service jobs
- **BUT:** Both cases require that PanDA and the respective DM system (Rucio, FAX) follow a similar source selection logic
  - Otherwise PanDA might be taking useless decisions
- Review Overflow jobs?