# NDGF Site report: Experiences with ARC/dCache

*Josva Kleist*
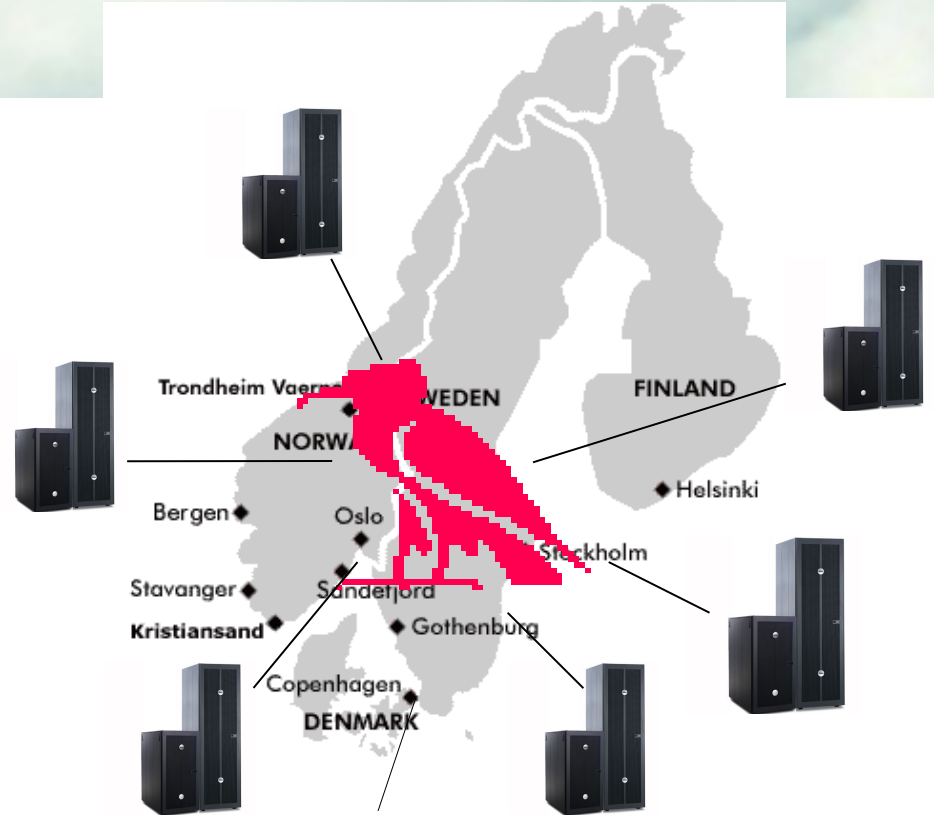
*Technical Director, NDGF*

*Amsterdam*, June 17, 2010

# NDGF is different

- A distributed virtual infrastructure (NDGF does not own any resources).

- "The Tier-1 is only about storage".

- Storage and compute resources are placed at national scientific compute centers.

- Computational resources are shared with other research communities.

- All computations carried out at "Tier-2" sites.

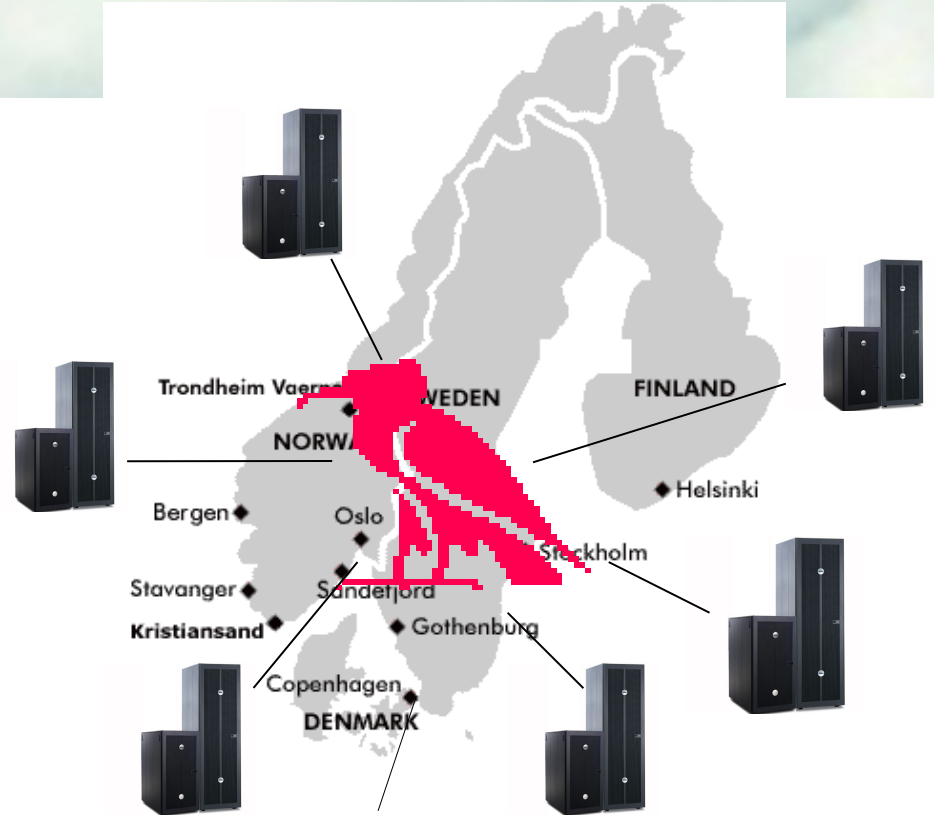- Relatively fat network pipes to most sites.

# Distributed dCache

- dCache Installation
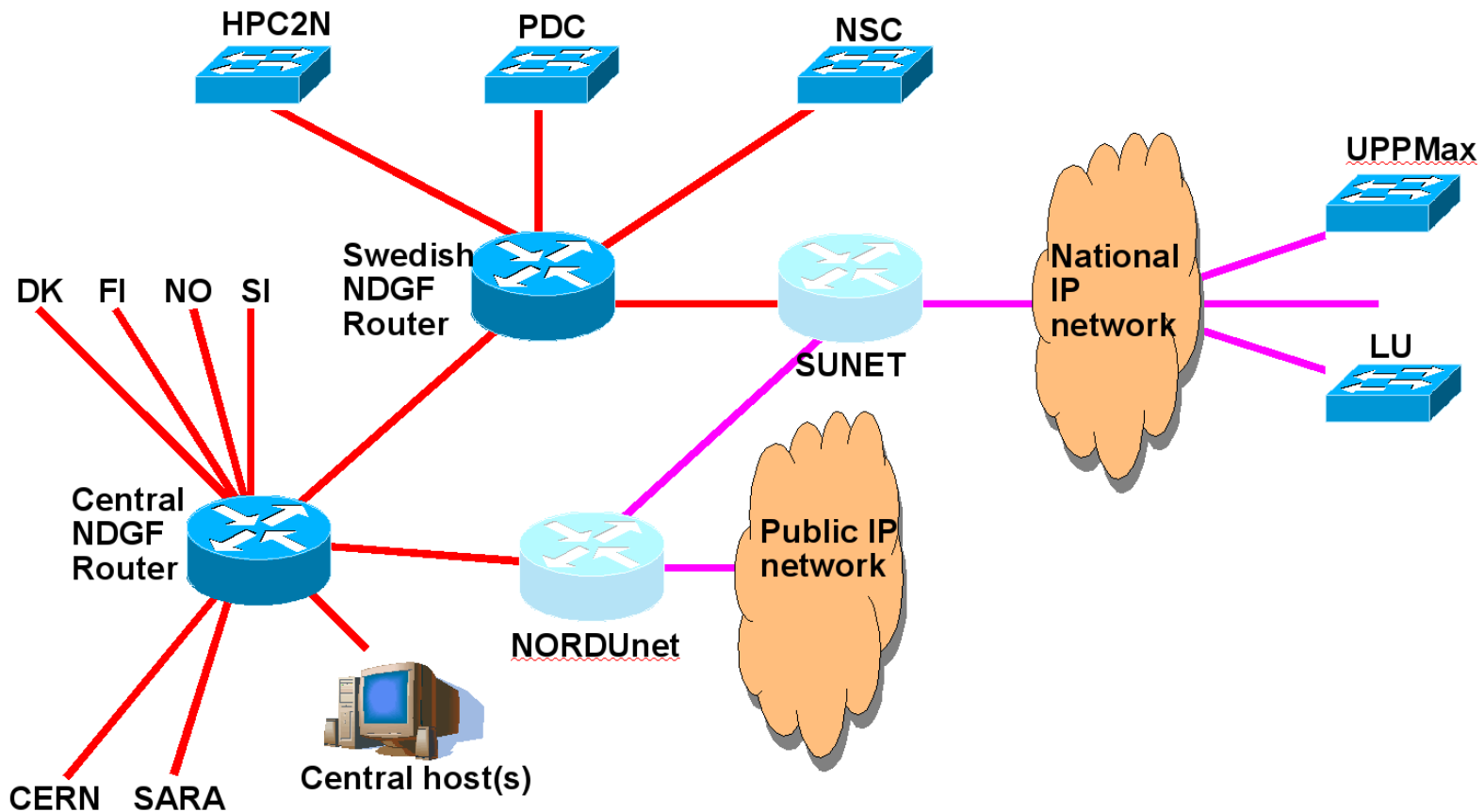- dCache head nodes at Nordic GEANT endpoint
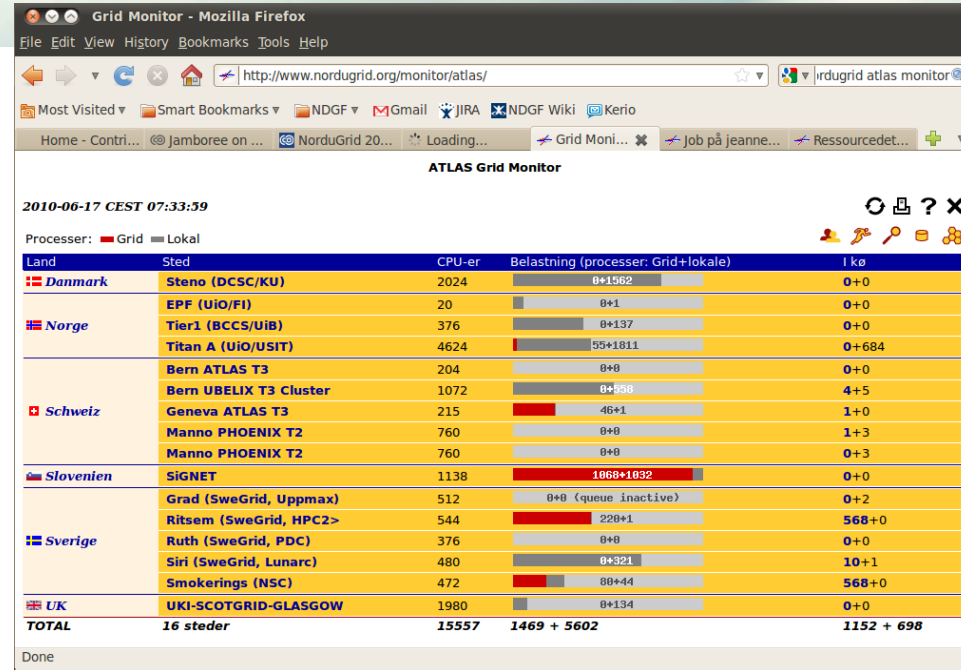- Pools at sites

# Distributed dCache

- dCache Installation
- dCache head nodes at Nordic GEANT endpoint
- Pools at sites
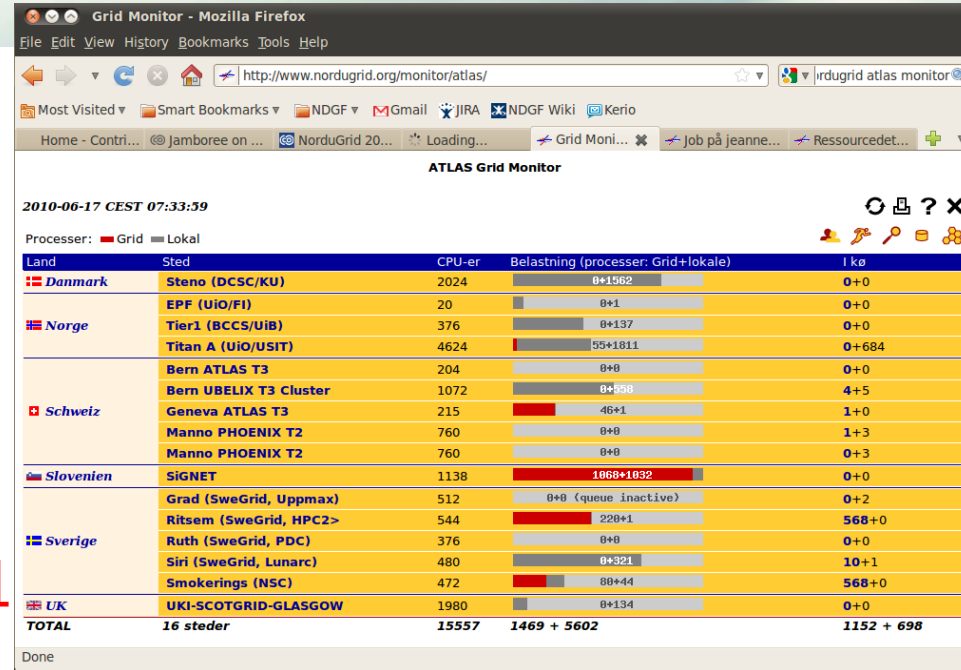- In this setting – Not much different from having storage at a single site.

# NDGF Compute model

- Data management is handled by the grid CE.

- Jobs access files on local file system (NFS/GPFS/LUSTRE).

- The CE has the ability to cache files.

# NDGF Compute model

- Data management is handled by the grid CE.

- Jobs access files on local file system (NFS/GPFS/LUSTRE).

- The CE has the ability to cache files.

- (Almost) all Tier-2 storage is used for cache (approx 1 PB).

- 100+ TB per site.

# NDGF Compute model

- Data management is handled by the grid CE.

- Jobs are not submitted to lrms before all input files are available.

- Without any knowledge the system will not perform...

# Where is a file cached?

- Without any knowledge the system will not perform…
- 20K analysis jobs consumes approx 400TB of data.
- But only 20-40TB of unique files.

- Without any knowledge the system will not perform...
- 20K analysis jobs consumes approx 400TB of data.
- But only 20-40TB of unique files.
- Without knowledge this took 4 days.

- Without any knowledge the system will not perform…
- 20K analysis jobs consumes approx 400TB of data.
- But only 20-40TB of unique files.
- Without knowledge this took 4 days.
- With the *cache index* only 10h.

- We do not need perfect knowledge.

- We do need fast lookup.

- We do need to be able to cope with failure (fast reconstruction of the index).



13

**NDGF**
NORDIC DATAGRID FACILITY

- We do not need perfect knowledge.

- We do need fast lookup.

- We do need to be able to cope with failure (fast reconstruction of the index).

- Bloom filters gives us this,
  - incremental update possible,
  - but with some probability of false positives.



14

# Data management at CE - Advantages

- Transfers can be scheduled and prioritized.
- Bandwidth usage can be throttled (e.g. user/VO limits).
- Cache management (LRA discard).
- Somewhere between *d* and *e* in Philippe's list of models.

- Preparation of cache possible.
- Cache2cache transfers can be implemented
  - authorization is the tricky bit.

# Data management at CE - Disadvantages

- Pilot jobs does not play well with CE handling data mgmt.

- ARC Ctrl Tower acts as pseudo pilot.

- In most cases only a few minutes delay between job submission before execution starts.

Thank you!

# Using the index

- Daemon rebuilds index at site and provides it on request to the index.

- Query is simple http request
  - \> curl -k  "https://cacheindex.ndgf.org:6443/data/index?url=http://www.nordugrid.org:80/data/echo.sh"

    {"http:\/\/www.nordugrid.org:80\/data\/echo.sh": ["benedict.grid.aau.dk"]}

- Removal of entries difficult → rebuild on a regular basis.

- Basically independent of ARC CE.

# Analysis jobs May

- titan            10167
- pikolit          157862
- grad             38508 (17k in one week)
- ritsem           35984
- pdc              6108
- siri             24085
- smokerings       3465
- swiss            11361

- **Total          287504**