

# Nordic Data Lakes Success Story

NeIC NT1 Manager  
Mattias Wadenstein  
<[maswan@ndgf.org](mailto:maswan@ndgf.org)>

2023-05-11  
CHEP, Norfolk, USA



# Overview

- Concept
- Components
- Collaboration
- Current status
- Challenges
- Conclusion





# Concept

- Starting with a distributed Nordic tier-1 site
  - Supporting ALICE and ATLAS
- One storage element spanning many sites in different countries
- Many independent computing resources connected to the same storage
- Could we build a wider data lake on this with further collaboration?



# Concept

- Nordic
  - Involving the Nordic countries
  - For us: Denmark, Finland, Norway, Sweden
  - Note: This work has not received funding by the European Union or the ESCAPE project
- Data Lakes
- Success
- Story





# Concept

- Nordic
- Data Lakes
  - A data lake is a repository of data
  - That can be transformed [into science]
  - Covers a geographical area
- Success
- Story





# Concept

- Nordic
- Data Lakes
- Success
  - For a research infrastructure improvement project:
  - Provides greater value to researchers at the same cost
  - At a lower cost for the same value
  - Or both
  - Can only really be evaluated after it is used in production
- Story



# Concept

- Nordic
- Data Lakes
- Success
- Story
  - A narrative, an account of events
  - Here told in slides and spoken word



# Component: Distributed dCache

- dCache architecture

- Lots of microservices that can communicate over WAN
- Local “movers” at data storage pools where data is transferred
- Common namespace and authorization components

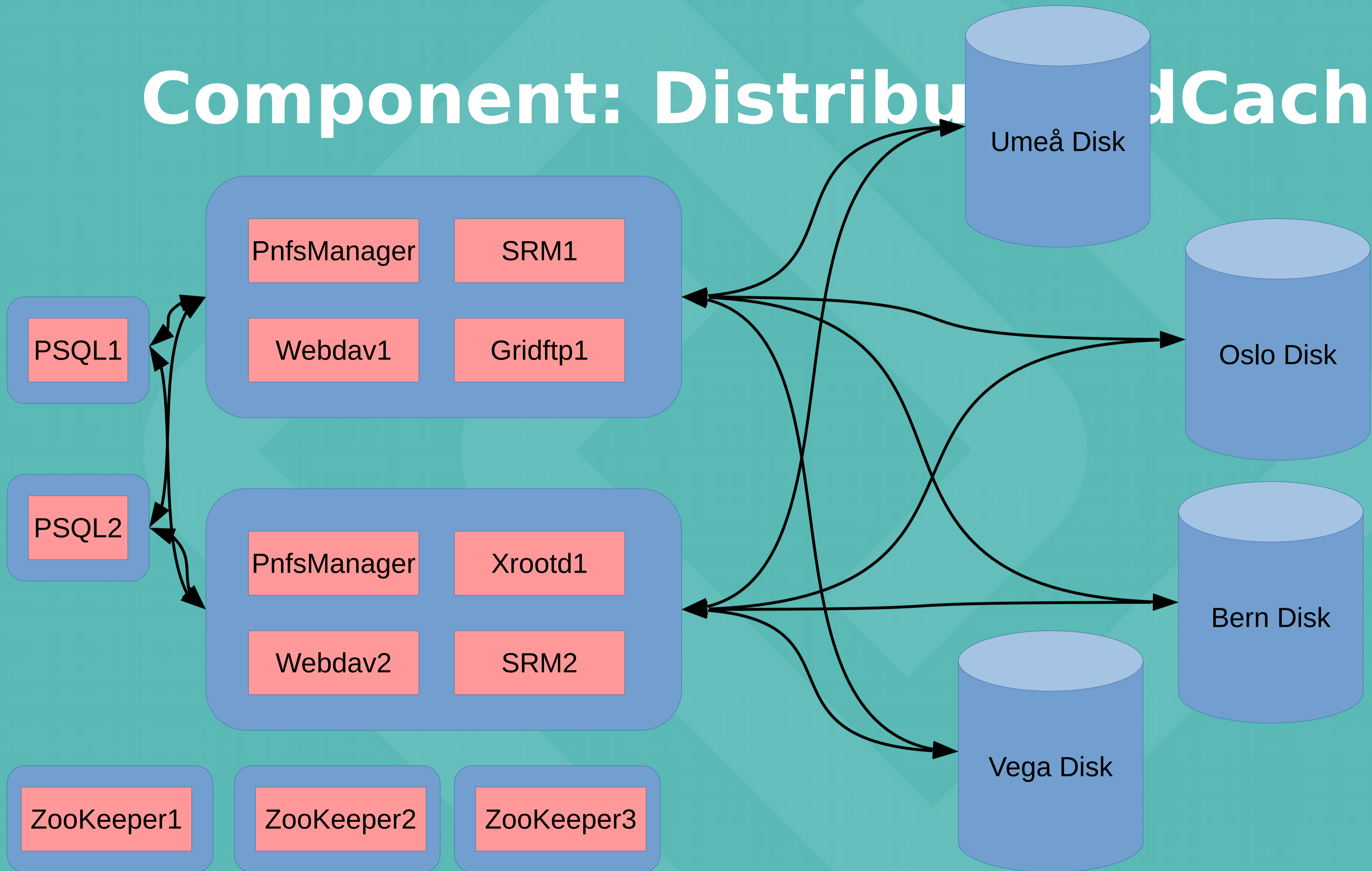
- Nordic innovation

- Multiple tape backends (hsminstances)
- Better redirect support to pool movers for serveral protocols
- High Availability improvements for core component upgrades without user impact





# Component: Distributed Cache



# Component: Distributed dCache

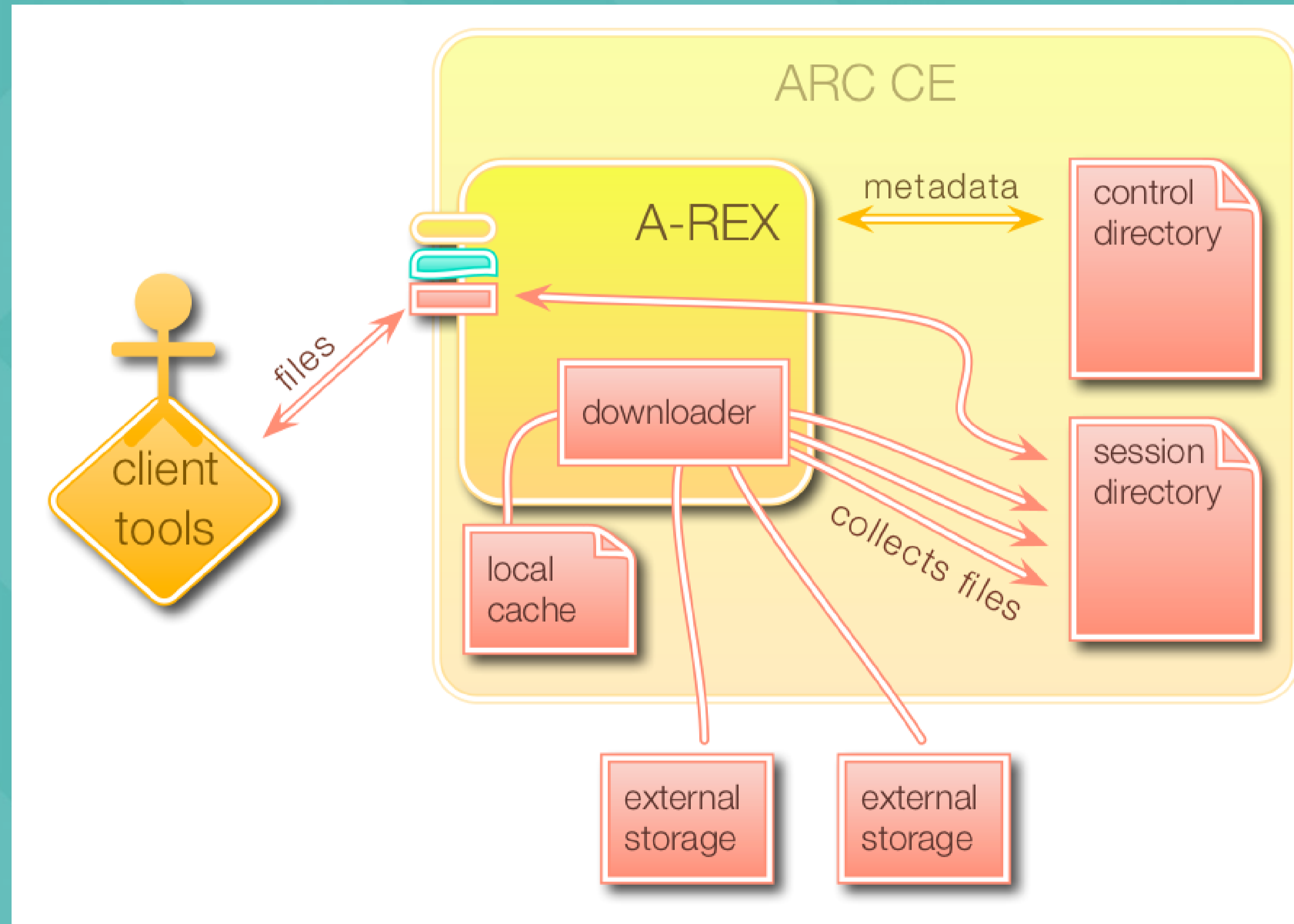
- Large set of guidelines for buying and running pools for our site admins:
  - [https://wiki.neic.no/wiki/DCache\\_Pool\\_Hardware](https://wiki.neic.no/wiki/DCache_Pool_Hardware)
    - RAM, CPU, MB/s/TB, etc. Supposed to be easy to turn into procurement.
  - [https://wiki.neic.no/wiki/DCache\\_Pool\\_installation](https://wiki.neic.no/wiki/DCache_Pool_installation)
    - Checklist for everything needed to be a pool in operations
  - [https://wiki.neic.no/wiki/Operations\\_Tuning\\_Linux](https://wiki.neic.no/wiki/Operations_Tuning_Linux)
    - Large TCP windows + BBR, Linux VM settings, some HW raid quirks
- Local site admins runs hardware, storage, and OS
- Central ops run dCache in unprivileged user
  - Ansible for handling pool tasks with good scalability





# Component: ARC with datastaging

- ARC-CE can do data staging
  - Prepares all input files needed by the job before submission to batch system
  - Saves all requested outputs to remote storage afterwards
  - Cache for reuse of input files between jobs



# Component: ARC with datastaging

- ARC in data caching mode
  - Each job description has a list of input and output files (rucio://...)
  - The CE stages all these files to local cache and links them in the session directory
  - The job is submitted to batch system and runs on local files only
  - Afterwards the listed output files are uploaded to SEs
  - Transfers over https, so same path as data movement
- Caches are normal shared filesystems
  - NFS, CephFS, GPFS, Lustre, etc
  - Size reasonable for SSD for ATLAS: 20TB + 5TB/1kcore



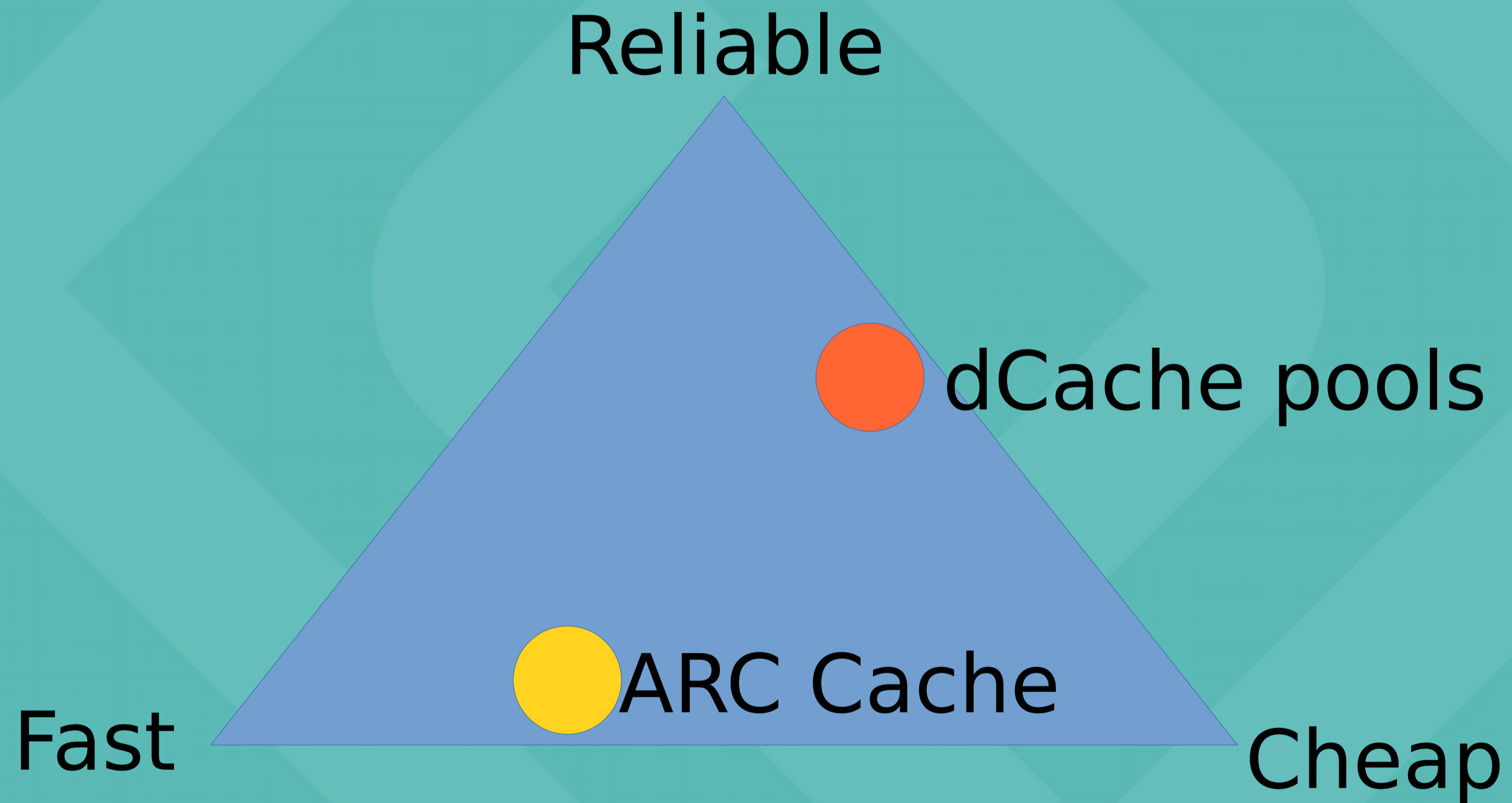


# Component: ARC with datastaging

- Overall efficiency
  - Data access is on low-latency local filesystems
  - Download before submission to batch system → better CPU efficiency
  - E.g. 47% → 90% CPU efficiency [M Pedersen, CHEP 2019]
- Non-local storage
  - Like NDGF-T1 with distributed storage
  - Or a “compute only” site
- Limited external connectivity
  - Like HPC sites where external connectivity might be blocked or only available through a slow NAT



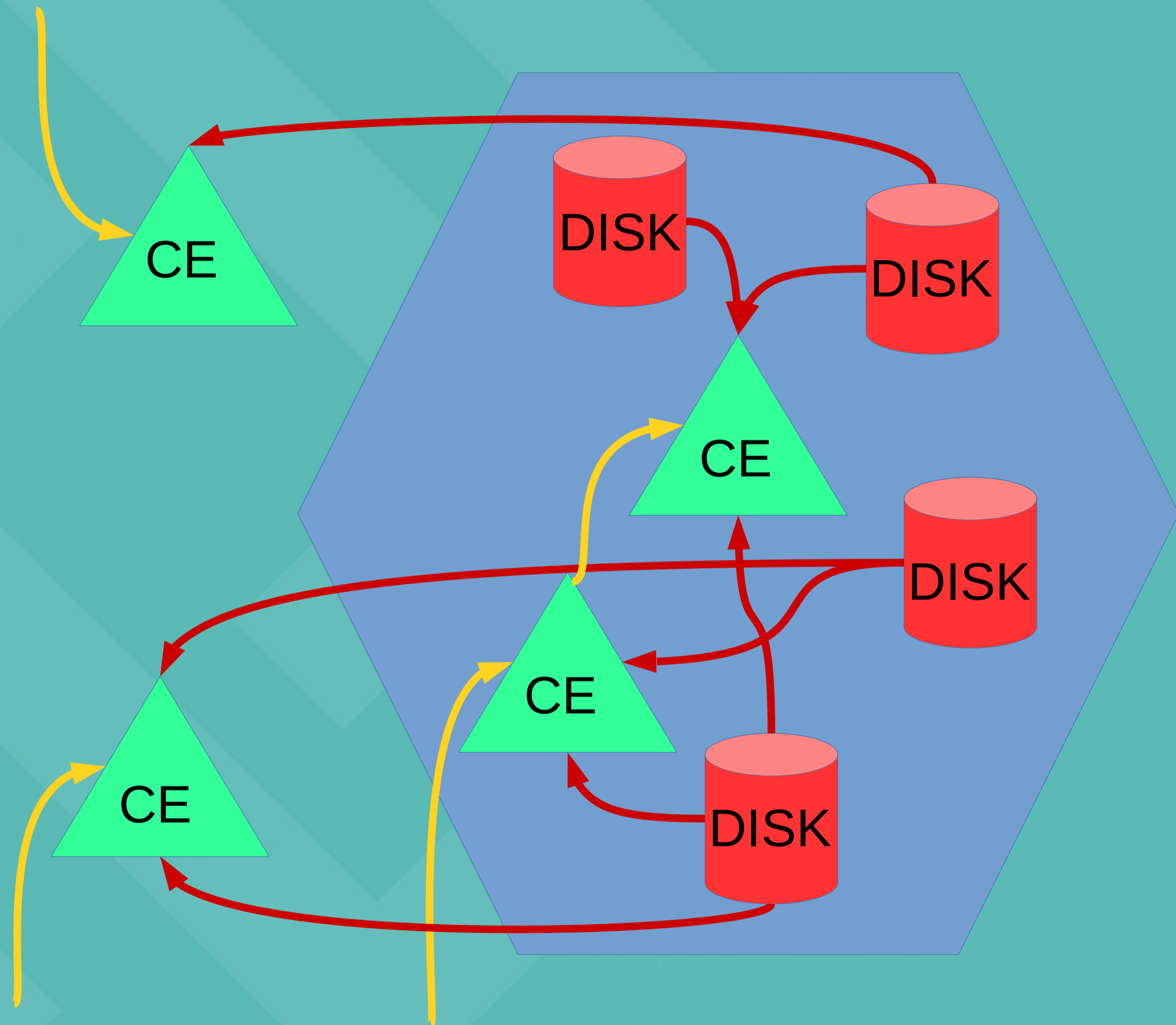
# Components





# A Hexagonal Data Lake

- Staging makes ARC location agnostic
- Setting to prefer “local” (T1) data
- No problem getting some data to/from other sites
- Fast internal network to keep CPUs full



# Collaboration

- **Motivation: Lower cost**
  - Managing storage elements including user support is non-trivial
  - The distributed nature has some overhead, but the reference comparison of 4-6 tier2 sites in the Nordic countries is on par
  - Adding more storage sites at very low marginal cost to NT1 saves on staffing, running pools (including procurement and commissioning) takes about 10% FTE.
- **Motivation: Better value**
  - Many small storage elements provide less value than a few large
  - Higher overall reliability, in particular for data taking (i.e. useful for job output destination)





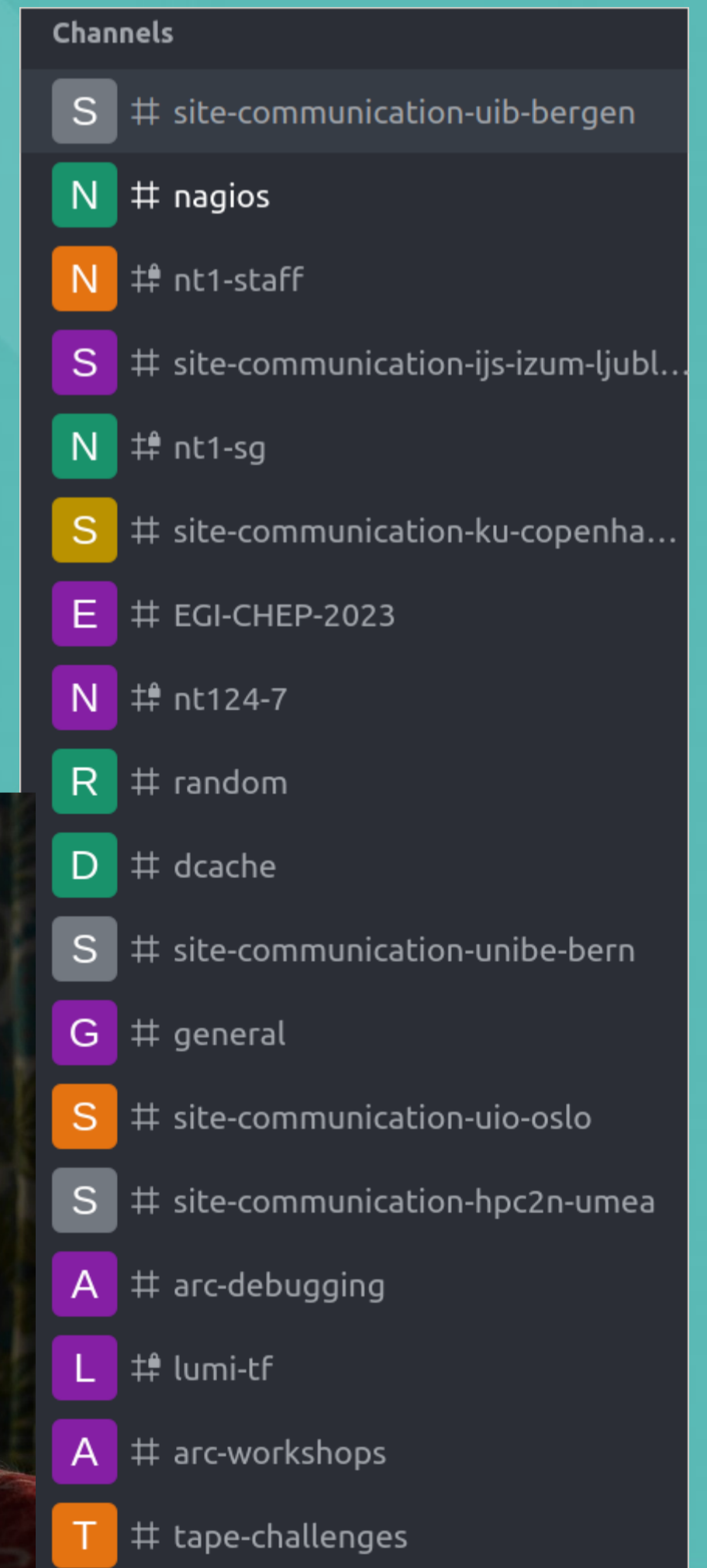
# Collaboration

- A successful data lake is a successful collaboration between:
  - Funding agencies – usually one in each participating country
  - Sysadmins – NeIC central team and site admins at each site
  - Physics projects and their PIs – one to two per country for us
  - Networking providers – NORDUNet, GEANT, CERN, plus all NRENs
  - Researchers – the entire purpose of research infrastructure
  - Experiment coordinators – ALICE and ATLAS currently
  - Scientific computing centers – Nine currently participating
  - Coordinating body – Nordic e-Infrastructure Collaboration, NeIC
  - etc
  - etc



# Collaboration

- Real-time communication in chat rooms for operational issues
- Regular meetings and other forums for coordinating with stakeholders
- Tickets, issues, applications, evaluations, ...
- Many emails





# Collaboration

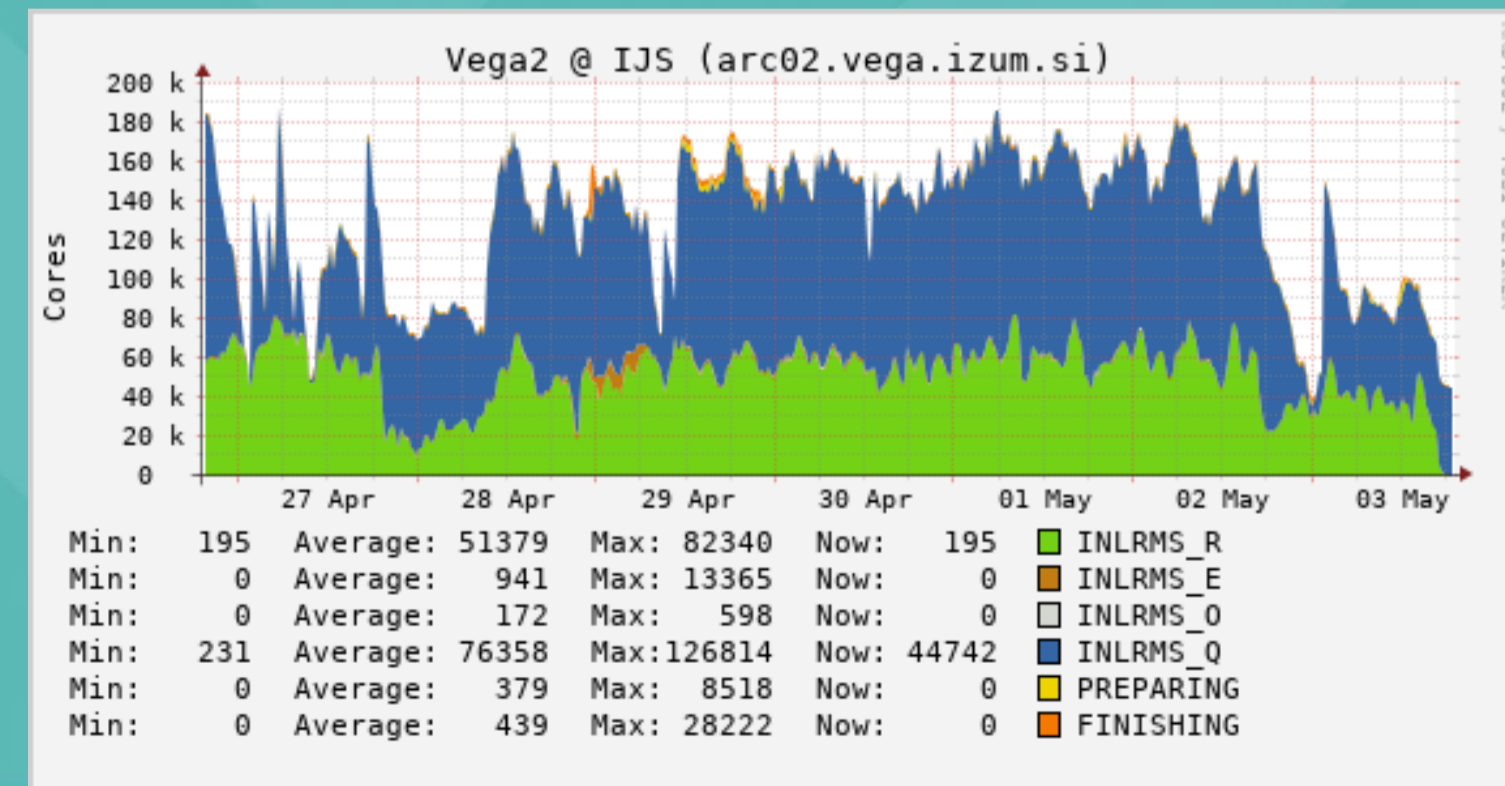
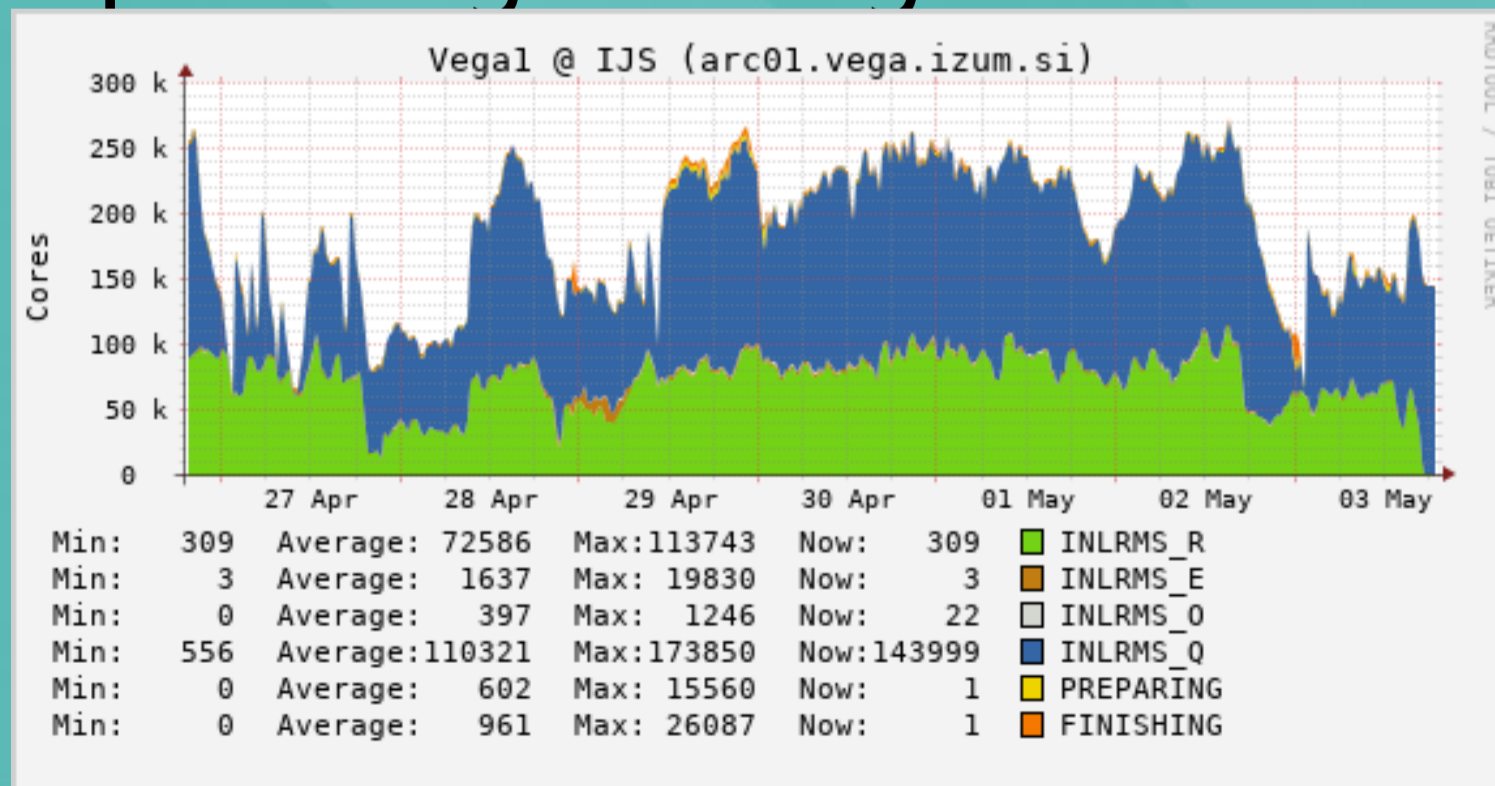
- Most recent onboarding: University of Bern
  - ATLAS Tier-2
  - 1.8 PB
  - Was running DPM, this a the DPM migration path
  - For process details, see HEPiX Spring 2023 presentation:  
<https://indico.cern.ch/event/1222948/contributions/5320953/>
- Other tier-2s integrated:
  - Slovenia (IJS and Vega)
  - Sweden (pledges both part of T1 and T2)



# Current Status

- Four Nordic countries plus Slovenia (IJS & Vega) and Switzerland (Bern) connected to one dCache
  - 8 PB ALICE disk
  - **23 PB ATLAS disk**
  - 19 PB ALICE&ATLAS tape
- Serving 50k-200k cores compute, T1+T2
  - depending on Vega fill situation

BNL-OSG2_DATADISK	29.92 PB
NDGF-T1_DATADISK	23.34 PB
RAL-LCG2-ECHO_DATADISK	20.58 PB
IN2P3-CC_DATADISK	18.83 PB





# Current Status

- **ALICE: Normal Nordic Tier-1, not widely distributed**
  - No local caching → worse CPU efficiency (avg a few percent)
  - Would get a bit worse if we had ALICE disk in southern Europe (RTT)
    - Could possibly be mitigated with Xcache
- **ATLAS: Large data lake for disk**
  - Much larger tier-1 disk area than our funded ambition of 6% of ATLAS tier-1 resources (currently second largest ATLASDATADISK area)
    - Tape is normal Nordic pledge of ~6% of tier-1 requests
  - Reliability usually on par with normal Tier-1s
    - A subset of transfer errors is shown harder to track down due to the distributed nature
    - On the other hand, a compute room power outage won't affect data taking
  - ATLAS finds more value in larger and more reliable storage elements



# Challenges

- **Reduced visibility for contributors**
  - Is a share of tier-1 storage as visible as a dedicated tier-2?
  - SRR feature should be able to handle this for WLCG accounting
  - Challenging to implement: little documentation and complex interactions between different systems (SRR, WSSA, CRIC, ..)
  - This is the first production deployment
- **Lowest performance needs to be good**
  - Slowest pool/site per TB determines average throughput
  - Running out of site bandwidth or buying a batch of slow servers





# Challenges

- Central operations needs long-term funding and continuity
  - Our funding agencies like to have competitive calls ever 4/5 years, NeIC has 6 of them (4 relevant for tier-1 central operations)
- Engaging new sites possible but usually non-trivial
  - Agreements and trust needed on several levels
  - Technical compatibility with local site admins
- ATLAS only LHC experiment using local ARC cache
  - Fixing payloads to read local filesystem files probably less complex than some of the heavy lifting to run jobs on inconvenient HPCs
  - Other caching solutions might be viable



# Conclusions

- We consider this a success: more value at lower cost
- We could integrate ~10 more storage sites into a single distributed storage for WLCG
  - Possibly more, but somewhere we start needing more staff than needed just to deal with the distributed Nordic sites → who pays?
- Supporting new experiments possible
  - Likely higher load for central team
  - Increased Nordic funding probably requires Nordic demand
- ARC with chaching for good compute efficiency
  - Even with storage far away
- Continious improvement for a smoother future







Questions?