

# Distributed Analysis with PANDA

Tadashi Maeno (BNL)

# New features of pathena

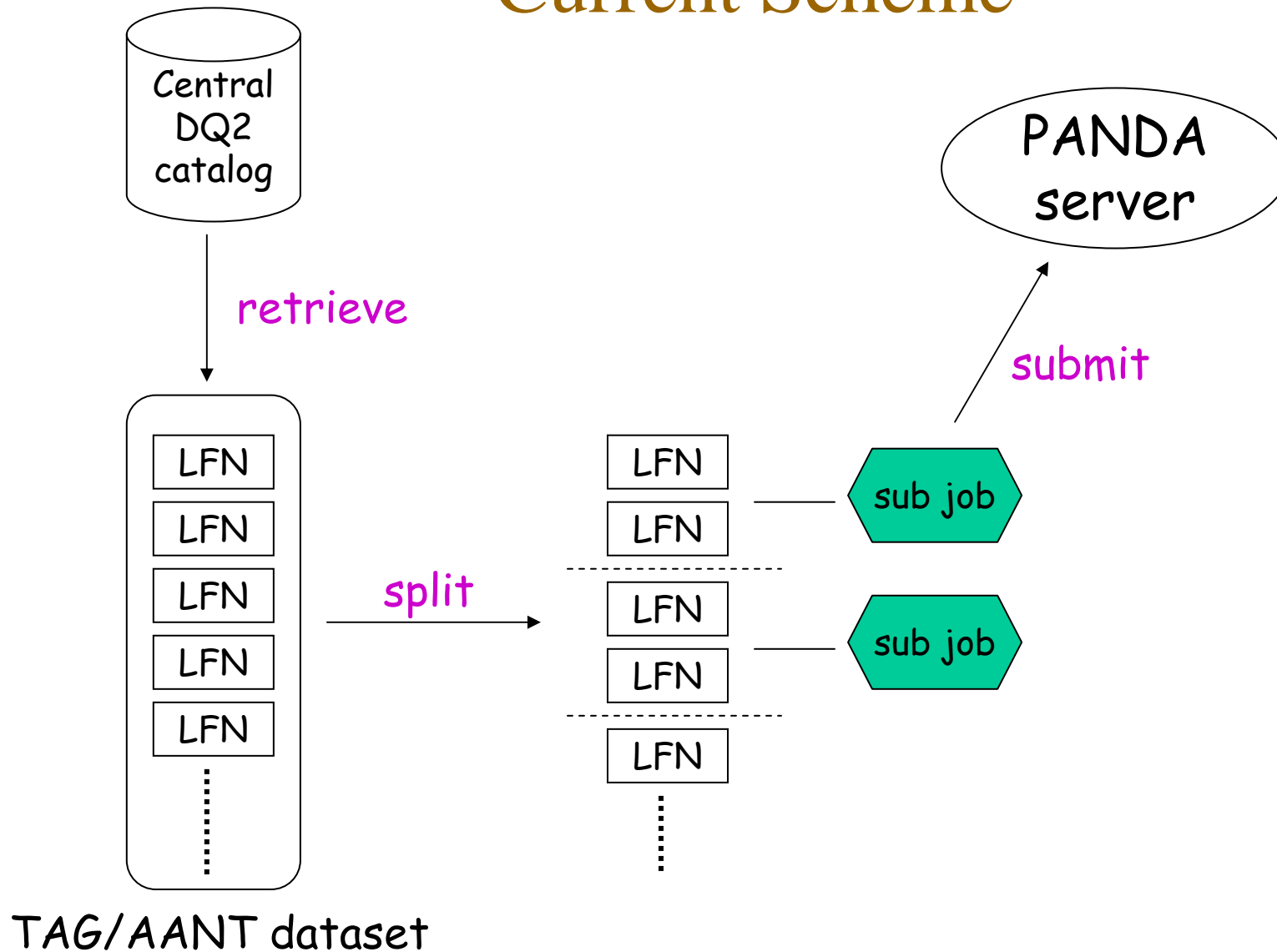
- Works with release kit in addition to AFS ver of Athena
- Supporting Group Area
- TAG(ROOT ver.)/AANT-based analysis
  - Back-navigation
- Merging AANTs
- Multiple input streams
  - Signal + Minimum-bias/Cavern
- Random seeds
- Repackaged for 12.0.3
  - From UserAnalysis to PhysicsAnalysis/DistributedAnalysis/PandaTools
  - Can be evolved independently of UserAnalysis

# Merging AANT

- addAANT has been included since 12.0.1/12.1.0  
<https://twiki.cern.ch/twiki/bin/view/Atlas/AthenaAwareNTuple#Merging>
- hadd in ROOT5 can merge AANTs but it breaks Athena-awareness
  - Athena will crash if the resultant file is used as input
- In 12.0.2, one has to copy AANTs to local PC first and run addAANT
- In 12.0.3 and onward, one can merge AANTs by submitting a job to PANDA

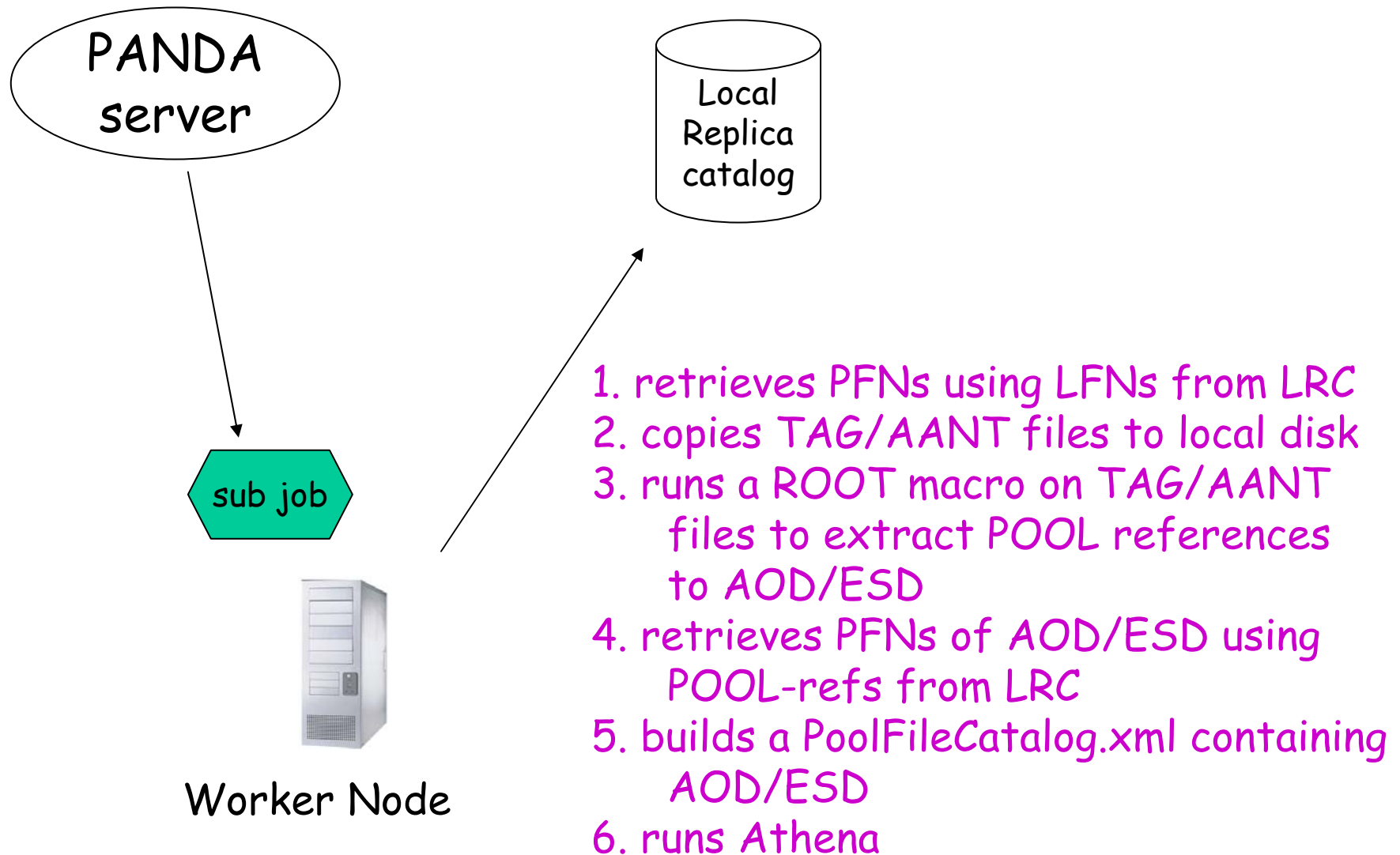
# TAG/AANT-based analysis

## – Current Scheme –



# TAG/AANT-based analysis

## – Current Scheme (cntd) –



# TAG/AANT-based analysis

– Coming Improvement for job-splitting –

- Currently one TAG file points to one AOD
- Job-splitting is based on the number of input files
  - One sub-job processes a certain number of files
  - It doesn't work if TAGs are merged to one file, for example
- Should be changed so that one sub-job processes a certain number of events
  - Sub-job 1 has something like  
EventSelector.Query = "EventNumber < 1000",  
Sub-job 2 has  
EventSelector.Query = "EventNumber >= 1000 AND  
EventNumber < 2000" ...

# TAG/AANT-based analysis

– with TAG database –

- Job-splitting would be based on the number of events
  - Each sub-job has a selection criteria
- Sub-jobs would local TAG DB
- Analysis TRF would retrieve a list of POOL-refs from TAG DB using the selection criteria to build a PoolFC.xml
- TNT : Tag Navigator Tool has been developed by Caitriana Nicholson at Glasgow  
<http://ppewww.ph.gla.ac.uk/~caitrian/tnt/>

# Personal history dataset

## ➤ Problem

Given that an user submitted a job on an official dataset (e.g., containing 1000files) and got a result

- The dataset is not closed in many cases, so new files (e.g., 100files) may be added after that. pathena\_utils allows the user to resubmit the same job on the dataset. However, the new job has to run on 1000+100files because it doesn't know how the previous job ran.
- Several files might be missing at the remote site where the user submitted the job. Most physicists may want to run their jobs even if datasets are incomplete. Analysis TRF is smart enough to skip files missing in the LRC. But currently there is no way for users to know how many files their jobs actually processed.



# Personal history dataset (cntd)

- Proposed solution
  1. Each new job creates an empty dataset
  2. Analysis TRF tells pilot what files it processed actually
  3. Pilot adds the files to the dataset
- When re-submit a job, pathena would compare input dataset and the history dataset and then instantiate sub-jobs for the increased files only
- This mechanism will be implemented by next month

## PANDA on LCG

- Server-pilot separation
- If pilots run in LCG, PANDA can use LCG resources
- Torre implemented a wrapper for pilot to run on CERN short queue and confirms it works properly
- Basic mechanism is already there. Will provide real service soon