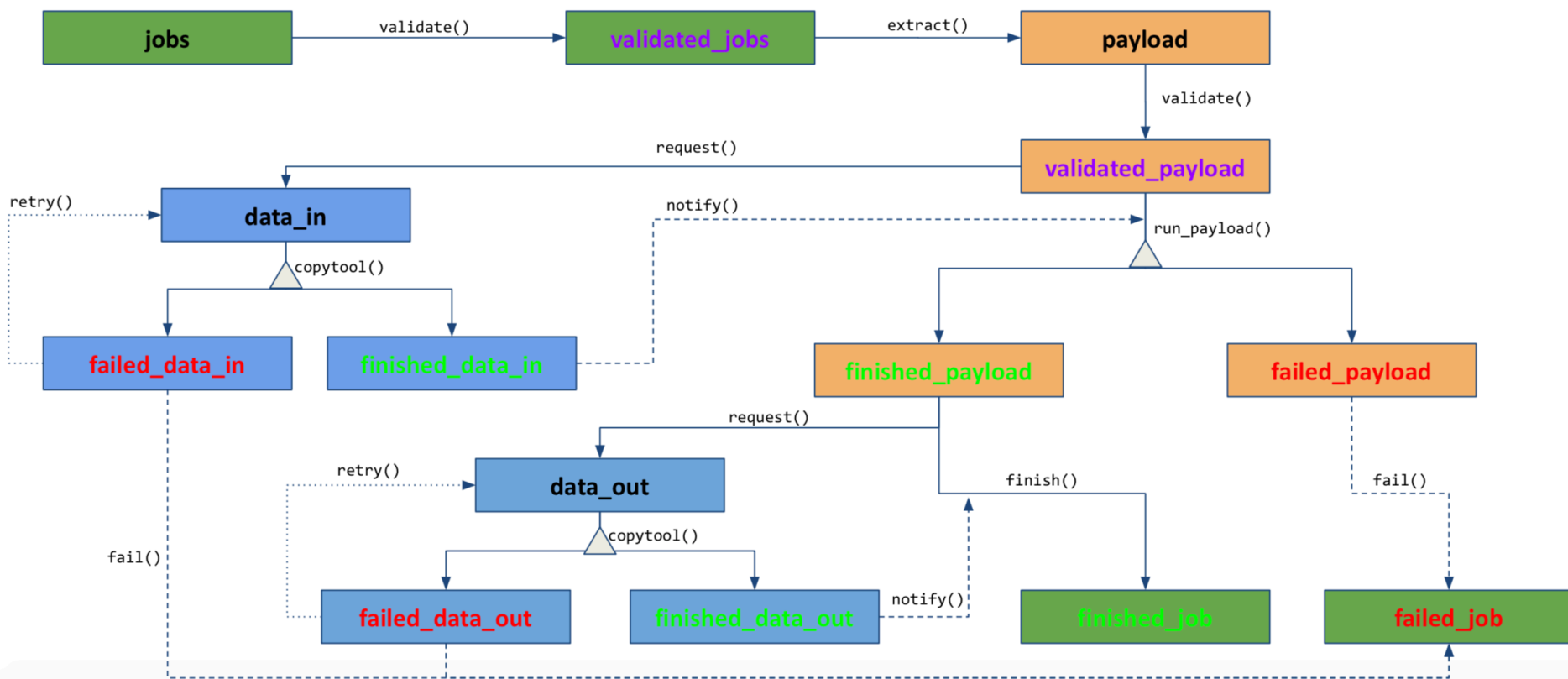Paul Nilsson

# PILOT 2

# Pilot 2 Project Overview

- Decision to replace the aging PanDA Pilot was taken in 2016
  - Old Pilot has served us well but after over a decade of development it has become difficult to keep building on
  - Pilot 2 is a complete rewrite
- Development of Pilot 2 is already in its second year and is going forward
  - Development is often done in sprints
  - Note: Pilot 1 is still often highest priority
- Component Model chosen last year
  - Project is in implementation stage since early this year
- First components delivered in March 2017
  - Harvester is using Pilot 2 Data API with stage-in/out (rucio)
  - Other APIs to follow soon (later slide)

**BROOKHAVEN**
NATIONAL LABORATORY

# Internal Flow of the Jobs Objects

- Job objects are kept in a job queue and are handled by the different pilot components



M. Lassnig

# Component Model Updates

- **Extended monitoring**
  - Pilot monitoring of internal threads
    - To be implemented (not urgent but could eventually be useful)
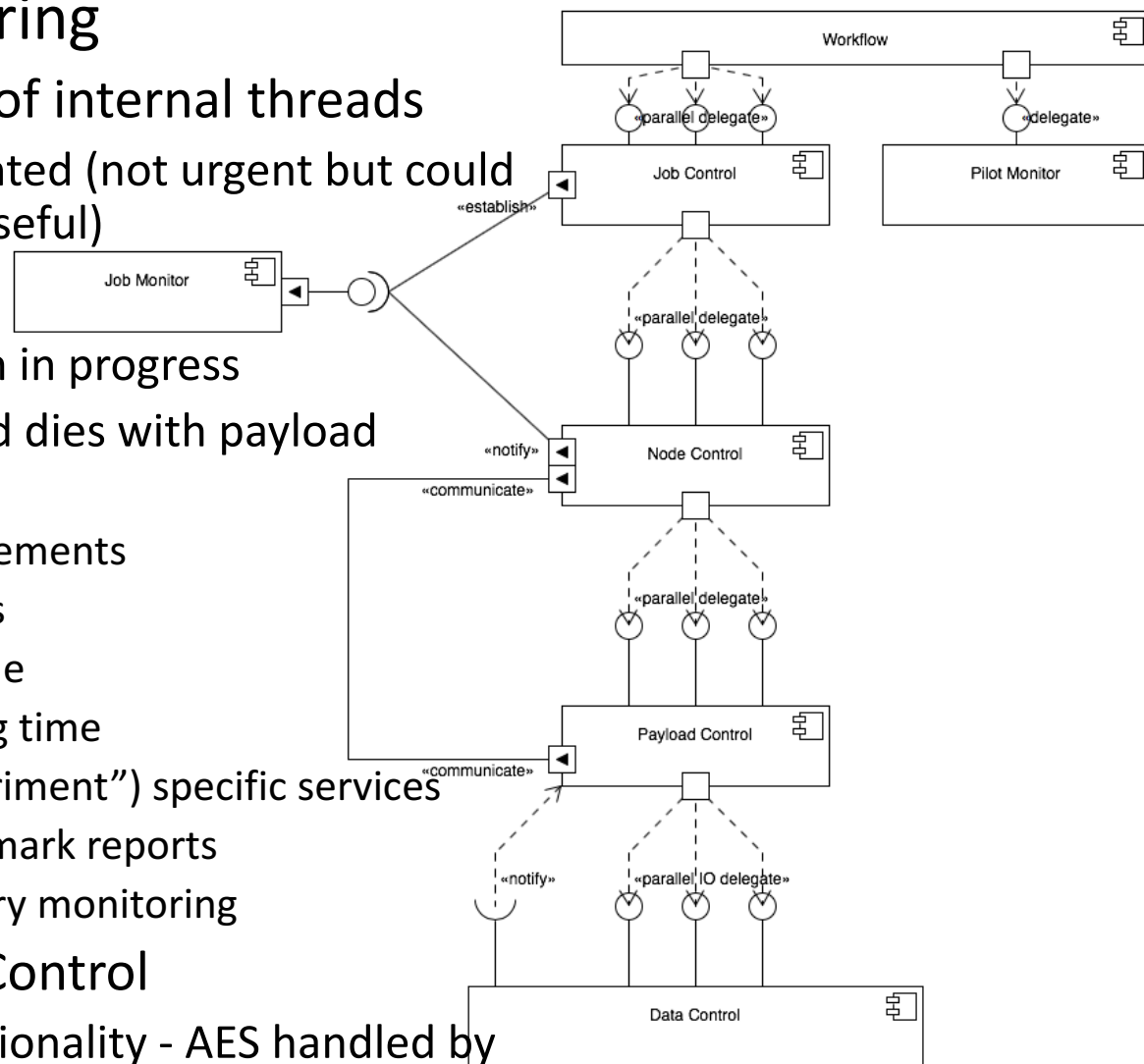  - Job monitor
    - Implementation in progress
    - Thread lives and dies with payload
      - Heartbeats
      - Size measurements
      - Looping jobs
      - Proxy lifetime
      - Pilot running time
      - User ("experiment") specific services
        - » Benchmark reports
        - » Memory monitoring
  - Removed Event Control
    - No loss of functionality - AES handled by dedicated workflow module

# Pilot Options (1/2)

- Implemented support in current development version (not final, might very well be changed)

| Option | Description |
|--------|-------------|
| -a | Pilot work directory |
| -d | Debug mode |
| -w | Force work flow; generic (implemented), event service, production or user jobs |
| -l | Maximum life time |
| -q | Queue name |
| -s | Site name (only needed by the dispatcher) |
| -j | Job label (user, managed, ptest) |
| --cacert | CA cert (not used on the grid) |
| --capath | CA cert path (not used on the grid) |

# Pilot Options (2/2)

| Option | Description |
|---|---|
| -p | PanDA server port (default stored in config file) |
| --config | Config file path (default config stored in pilot source) |
| --countrygroup | Country group for getJob request |
| --workinggroup | Working group for getJob request |
| --allowothercountry | Is the resource allowed to be used outside the privileged group? |
| --allowsameuser | Multi-jobs will only come from the same taskID |
| --url | PanDA server URL (default stored in config file) |
| --pilotuser | Pilot user, e.g. name of experiment (used to select user specific code stored in special pilot directory) |

- More options will be implemented when they are needed

# Pilot 2 APIs

- Some Pilot functionality is exposed to external users by APIs; currently being planned for, or is already available
  - **Data API**                                                                    DELIVERED
    - Basic stage-in/out already used by Harvester
    - New request: asynchronous stage-in/out
  - **Communicator API**
    - Functions for communicating with PanDA server, Harvester, aCT, ..
    - API defined; contains functions for downloads/updates of jobs and event ranges
  - **Environment API**
    - Interface to the job execution environment on HPCs
  - **Services API**
    - Possible new API which could expose functionalities related to services run by the pilot (being discussed), see later slides
  - **Container API**
    - Possible new API, see next slide

# Container Support

- ## Work in progress
  - Usage of containers discussed in separate session (how exactly to implement; tests done with Pilot 1)

- ## Pilot 2 container module defined
  - May later be used via Container API in case of interest
    - Could be used e.g. by wrapper
    - Can be delivered soon if needed
  - Currently consists of single function for command execution
    - Function can be used to wrap/decorate a command with timer, time-out, container command, ..

# Utilities

- Pilot 1 has hundreds of major and minor functions
  - A large part of Pilot 2 development is to re-implement many of these (cut and paste in some cases)
- Pilot 2 has utilities organized in dedicated util/ folder
  - Current code base include functions in multiple modules
    - E.g. constants, disk, filehandling, https, information, ..
      - Preliminary information module presents interface to AGIS and schedconfig
        - » To be replaced by a full Information Service component (where AGIS/schedconfig are not hardcoded but accessed via user code)
        - » Development to start as soon as possible
    - Pilot 2 now supports standard configuration files
      - Config files are shipped with pilot source (default values), but can be preplaced either in /etc or in init directory

# Copy Tools and Data API

- Copy tools rucio and xrdcp copy are supported in Pilot 2 development version
  - Only rucio is available in current production version Data API
- We want to provide functionality for non-Rucio users/experiments
  - Current xrdcp copy tool implementation relies on rucio for PFNs
    - Can be made independent if pilot is given PFNs from somewhere else
  - Restructuring / refactoring of the code is underway to facilitate adding additional copy tools (< 2 weeks)
- About to start development for supporting additional copy tools
  - gfal-copy, lsm (mv/storm could be considered special cases of lsm)
- Support for asynchronous stage-in/out has been requested (currently only basic stage-in/out is implemented)
  - To be developed asap (to be decided)

# Event Service in Pilot 2

- A new algorithm / workflow for event service is being planned for Pilot 2 based on years of experience and evolution of AES

  - AES development for Pilot 1 continues but the rewrite of the algorithm is currently only planned for Pilot 2

- Implementation to begin shortly (within 2-3 weeks)



Wen Guan

# Benchmarking on HPCs and Tier-0

- Expose benchmarking functions (for execution, output dictionary manipulations, reporting, ..) in new Pilot 2 Services API
  - Note: pilot adds some info not present in the default benchmark dictionary, and removes some fields that are not wanted

- *Idea #1*:
  - Harvester can add benchmark step to worker / MPI rank #0 which will execute it before running the actual job

- *Idea #2:*
  - Introduce new Benchmark_tf.py transform that allows for having dedicated benchmark tasks running on the grid / HPCs whenever needed [ok with Graeme Stewart and Marcelo Vogel]
  - **Complications**: 1) want to run benchmark suite with same number of cores as payload, but if there's no Athena payload in these jobs then what (re-run the task) ;  2) we loose a valuable HPC queue slot for a mere benchmark job
  - **Advantage**: Allows for benchmarking on Tier-0 where pilot is not normally running

- When job has finished, Harvester will discover the output benchmark dictionary in the shared file system and again use Pilot 2 Services API functions to report results to intermediary Elastic Search service (which populates ES)
  - On Tier-0, corresponding wrappers take benchmark results from the transform and upload them post-facto using pilot functions [ok with Graeme Stewart and Marcelo Vogel]

# Memory Monitoring on HPCs (Yoda)

- ***Idea***:
  - Expose memory monitoring functions (for execution, parsing, ..) in new Pilot 2 Services API
    - More useful for Yoda than standard HPC job since Yoda should be able to monitor each rank, while this is not true for standard HPC job

- Yoda adds memory monitor step to worker / MPI rank #0 which will execute it in parallel with the actual job
  - Memory info available during running, can be reported to Harvester and PanDA server

- When job has finished, Harvester will discover the final memory monitoring dictionary in the shared file system and again use Pilot 2 Services API functions (to add it to jobReport.json / upload it to intermediary ES service)

# Pilot Wrappers

- Known changes for production wrappers
  - Will need to migrate to use new pilot options
    - Details to be discussed/worked out
    - Wider grid testing towards the end of the year (2017)
    - HC tests in (early?) 2018
  - Be able to run pilot within container
    - Being discussed
    - Possibly use Pilot 2 Container API (Singularity and Docker?)
- Development version of Pilot 2 uses special dev wrapper
  - Practical to have our own wrapper
    - E.g. Pilot 2 uses different/new pilot options, has ES functionality, and the wrapper is minimal
    - Currently located in personal GitHub repo(s)
  - Currently used for tests on EU and US resources
  - Not intended for production use or as replacement of other wrappers

**BROOKHAVEN**
NATIONAL LABORATORY

# MiniPilot

- Special pilot version for minimalistic workflows
  - Paper: http://ceur-ws.org/Vol-1787/197-201-paper-33.pdf
  - Intended for testing during early Pilot 2 development
  - Note 1: we do not intend to make major developments on this code - others may modify it as they like [see below..]
- Recently migrated into final Pilot 2 GitHub repository
  - Currently using its own Pilot modules, i.e. no current usages of Pilot 2 code
- Special GitHub branch for version used by Harvester
  - minipilot-harvester [to be created shortly]
  - Developed by D. Benjamin, T. Childers

# GitHub Technology

- Pilot 2 GitHub is using Travis CI for automatic code verification/validation and unit tests
  - GitHub pull request into Pilot 2 repo triggers external service (runs pep8, flake8 and unit tests)



D. Drizhuk

# Pilot 2 – Sphinx Documentation

- Semi-automatic code documentation using Sphinx
  - Module to be documented must be accompanied by related sphinx file
  - Pull request followed by [currently] local sphinx script execution which builds the documentation
  - Output need to be moved to www server

- Service can be added to GitHub (hosted)
  - To be investigated further



D. Drizhuk

# Pilot 2 – ACAT 2017 Poster

- "**Global heterogeneous resource harvesting: the next-generation PanDA pilot for ATLAS**"

- Presented at ACAT 2017, Seattle US

- https://indico.cern.ch/event/567550/contributions/2627120/attachments/1516202/2366285/ACAT_Poster_copy.pdf

- Proceedings paper to be written shortly

# Current/Near Term Development

- Full chain of generic workflow running
  - Currently ironing out bugs
  - To be completed asap to facilitate for the more complicated AES workflow
- Event Service
  - In design stage to be implemented shortly
- APIs
  - Either already in design or implementation stage
- Job Client manager
  - Returns the proper job client object depending on the user
    - Needed by APIs
  - Code written, to be migrated into Pilot 2 within weeks
- Data Control extensions
  - xrdcp support recently added, but data module now needs some refactoring before additional copy tools will be added
  - Asynchronous stage-in/out requested (no estimate for delivery yet)
- Function development
  - Ongoing (bulk part of Pilot 2 development)

# Time to Completion

- At least one more year of development until Pilot 1 can be fully retired

  - Pilot 1 feature requests are often a top priority and slows down other developments

- Focus now is on supporting existing workflows with more features (e.g. monitoring) and implement new ones (read: event service)

  - In particular, API development is a priority since it provides Pilot 2 functionality before the 'final' product is ready

  - Pilot 2 can gradually take responsibilities away from Pilot 1, e.g. event service

# Pilot OTP – Jun-Dec (Projected)

\* Developers involved with both Pilot 1 and Pilot 2

| Name | Time | Task |
| --- | --- | --- |
| Paul Nilsson (*) | 90% | Project leader; core software and coordination |
| Danila Oleynik | 60% | Pilot HPC, Pilot 2 payload specifics |
| Wen Guan (*) | 30% | OS/Pilot general development |
| Daniil Drizhuk | 25% | General pilot development |
| Alexey Anisenkov | 15% | AGIS/Pilot interactions |
| Kaushik De | 10% | Coordination |

(*) Event service contributions not included here but counted separately
NB: Contributions from Mario Lassnig / Tobias Wegner counted as DDM