

Adaptive Pilot Framework for Distributed Workload Execution in the SPD Online Filter System

Romanychev Leonid
JINR MLIT
romanychev@jinr.ru

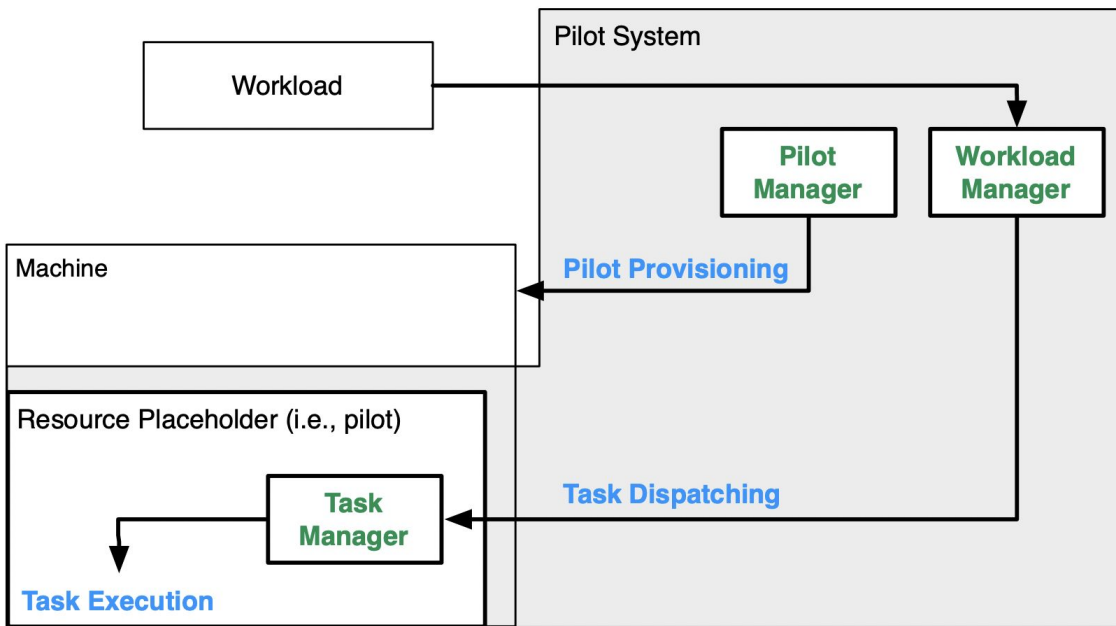
October 27-31, 2025

Introduction: Role of Pilot Framework

- Provide a flexible mechanism for execution of computational tasks.
- Widely used in high-throughput computing (HTC) systems for scientific data processing (ex. LHC computing).
- Issue: lack of unified abstraction and best practices leads to a variety of implementations.

Core Components of Pilot Software

- **Pilot Manager:** Launches pilots (resource placeholder: on computing resource; Interfaces with SLURM, HTCondor, etc).
- **Workload Manager:** Organizes task queue (dependencies, priorities, resource readiness).
- **Task Manager:** Executes tasks on pilot-reserved resources; Manages task lifecycle (launch, restart, monitor, error handling).



Functionality & Architecture

Functional Stages:

- Provisioning (Acquiring & deploying resources)
- Dispatching (Assigning tasks to pilots)
- Execution (Running tasks on resources)

Architectural Features:

- Multi-level scheduling
- Communication Models (Master-worker, Broker-oriented)
- Flexibility (integrates with various DCRs: clusters, grids, clouds)

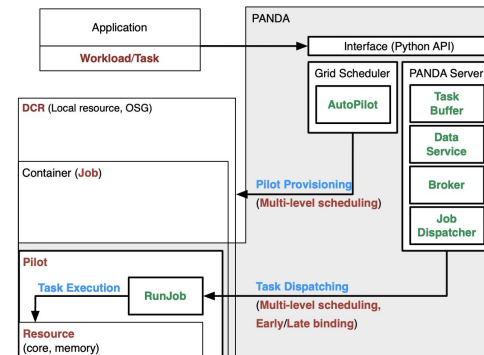
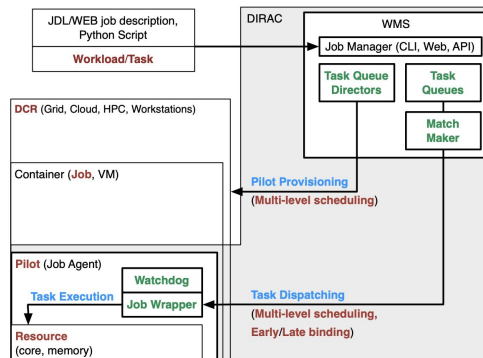
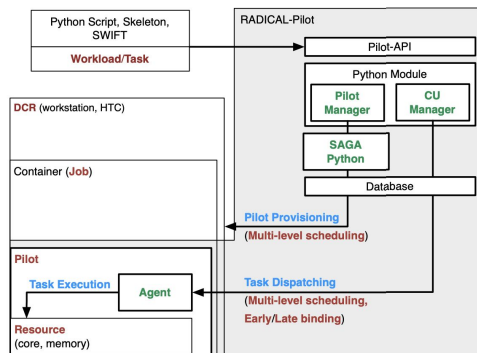
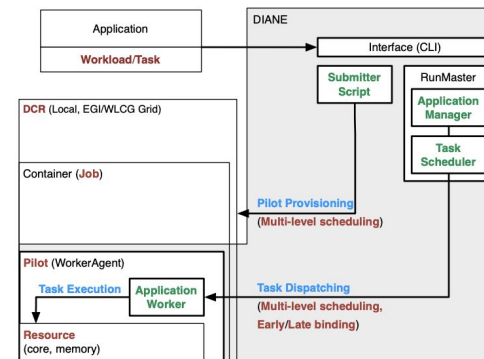
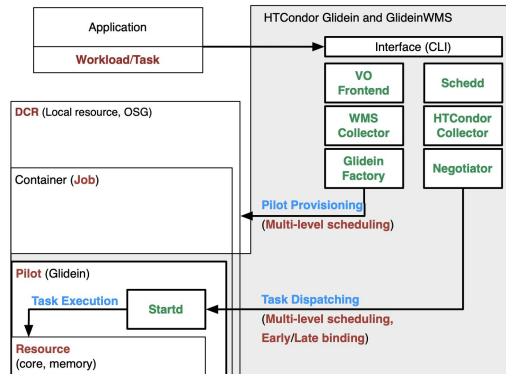
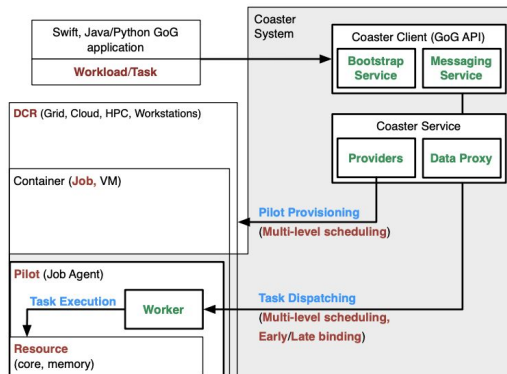
Late Binding Mechanism

Definition:

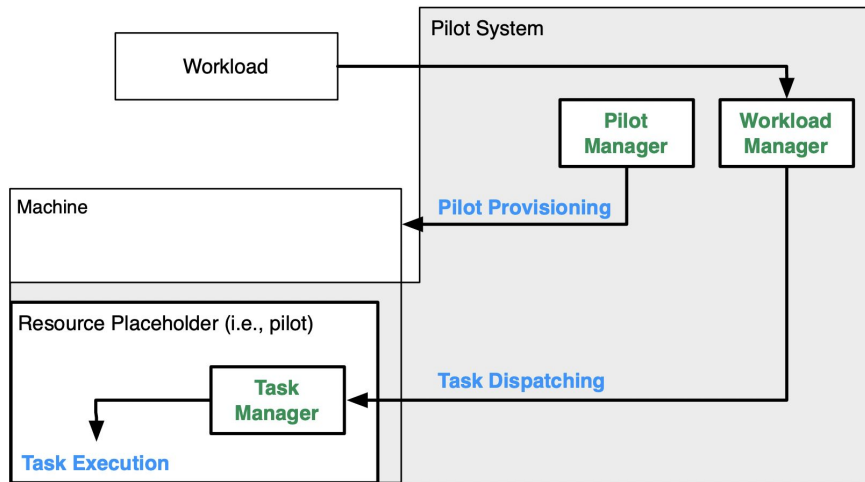
Late binding is the process of assigning tasks to active pilots at the moment of availability, unlike early binding, where tasks are tied to inactive pilots.

Benefits:

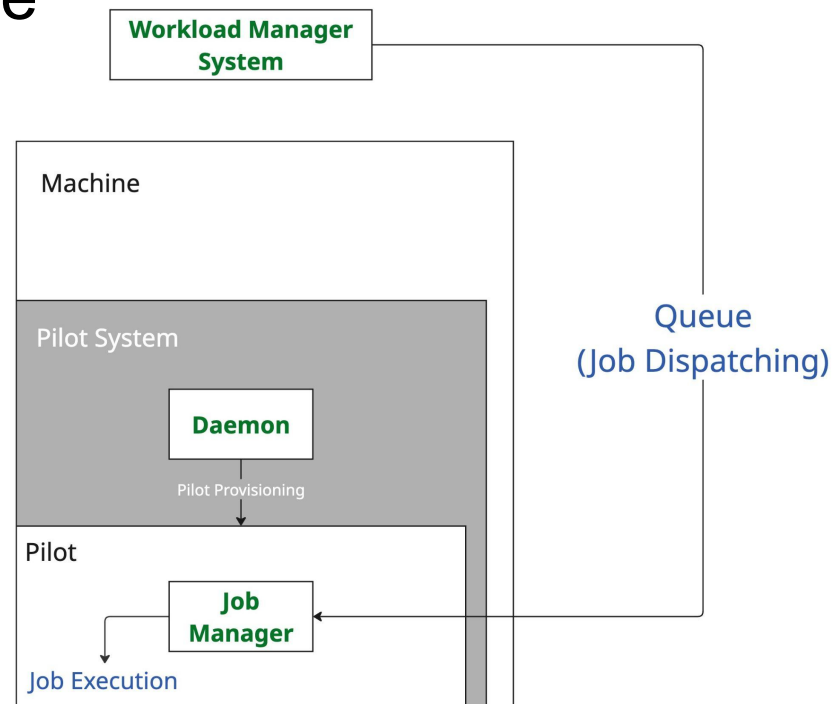
- Dynamic task allocation improves resource utilization efficiency.
- Reduces queue wait times, critical for high-performance systems.
- Enables high throughput (e.g., up to 1 million tasks per day for ATLAS).



SPD Online Filter Pilot Software



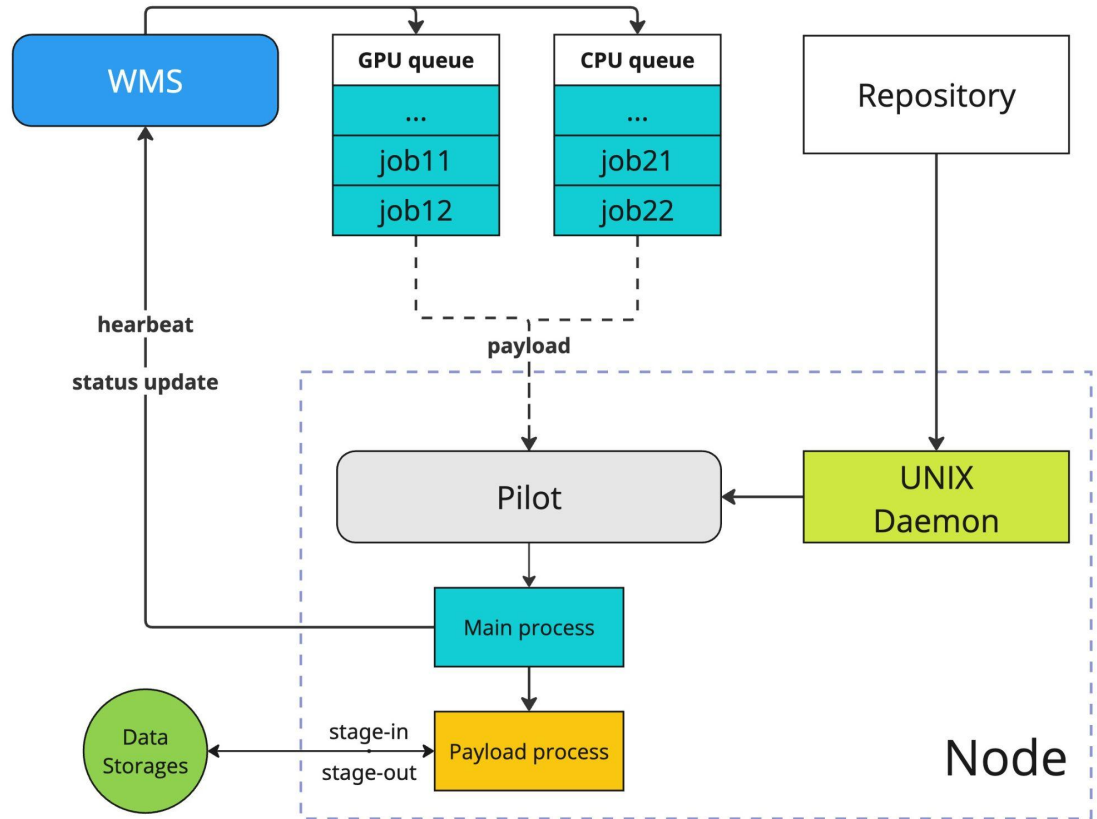
The general scheme



SPD Online Filter Pilot Scheme

SPD Online Filter Pilot Software

1. Two separate queues for CPU and GPU tasks.
2. The Pilot consists of two processes:
 - a. **Main Process:** communicates with the WMS (sends heartbeat and status updates).
 - b. **Payload Process:** executes the actual job payload.
3. Input/output data is transferred via Data Storages (NFS now).
4. WMS controls task distribution and monitoring.



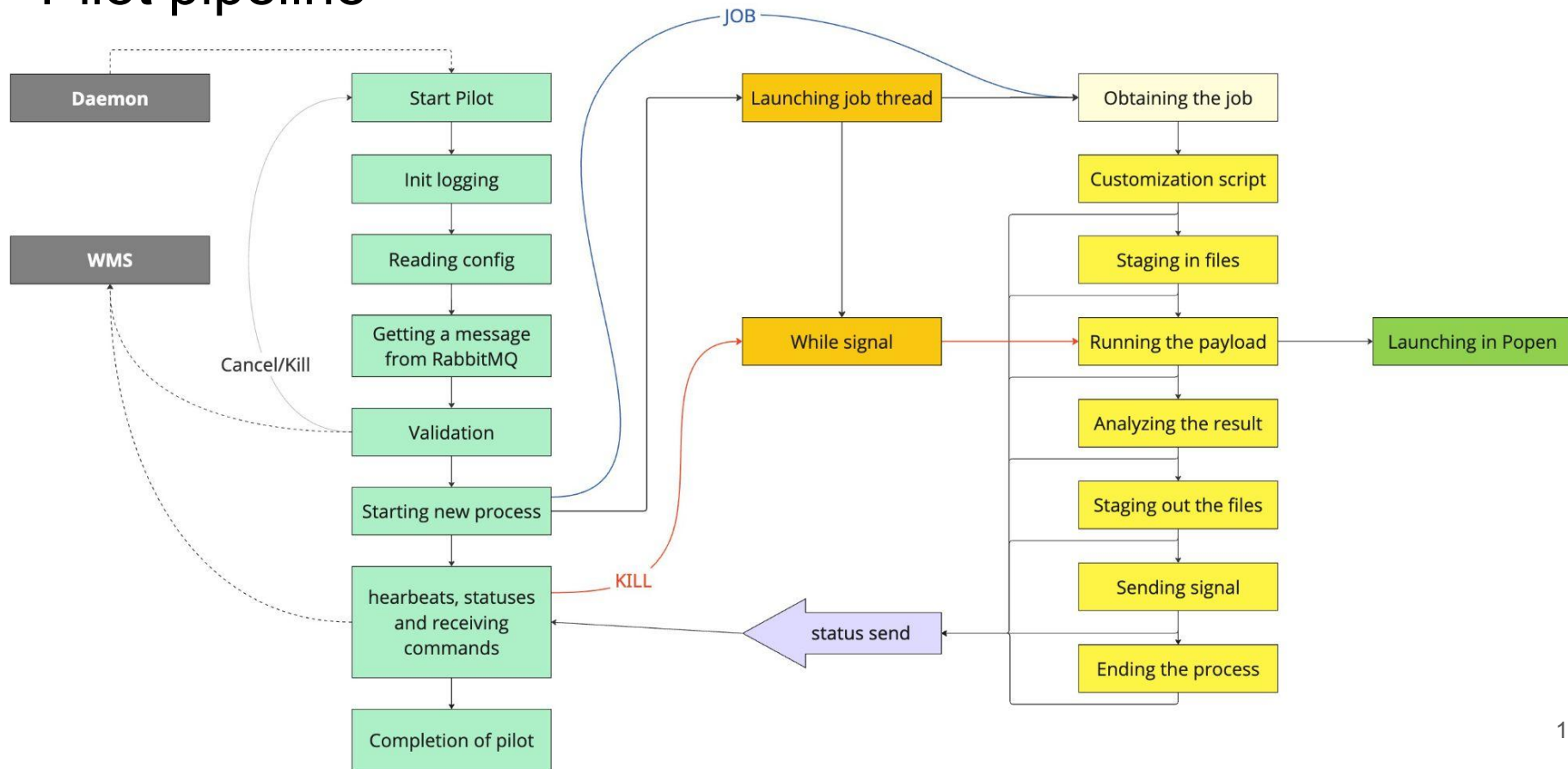
Job state model

- REGISTERED/READY
- EXTRACTED
- OBTAINED
- PRE-PROCESSING
- STAGE-IN
- **RUNNING**
- FINISHED
- POST-PROCESSING
- STAGE-OUT
- COMPLETED
- **FAILED**

set by the WMS

set by the Pilot

Pilot pipeline



Additional scripts

Name of script	Type	Description
update_pilot.py	Legacy	Downloads the latest pilot package from GitLab and updates it in NFS storage automatically.
registrator.py	Setup/Testing	Registers datasets and files in the API, calculates checksums, and manages builder directories.
daemons_runner.py	Testing	Manages multiple daemon processes: start, stop, and status. Handles config files and logging.
archiver.py	Setup	Archives a directory or copies a file to NFS storage. Supports both release and dev modes.
setup.sh	Setup	Installs dependencies, unpacks binaries/configs, and prepares the environment for running the daemon.

Configuration of the pilot and daemon

```
pilot > ≡ config_example.ini
1  [rabbit_settings]
2  RABBITMQ_USERNAME=user
3  RABBITMQ_PASSWORD=password
4  RABBITMQ_HOST=11.111.111.11
5  RABBITMQ_EXCHANGE=jobs
6  RABBITMQ_VIRTUAL_HOST=virtual_host
7  RABBITMQ_PORT=5672
8
9  [node_settings]
10 SERVER_ADRESS = http://11.111.111.11:8080
11 PROCESSOR_TYPE = cpu
12 PILOT_ID = 1
13 SCRIPT_PREFIX = /path_to_data/
14
15 [logging]
16 LOG_LEVEL = INFO
17 MAX_LOG_SIZE = 10485760  # 10MB
18 BACKUP_COUNT = 5
```

```
daemon > ≡ daemon_config.ini
1  [node_settings]
2  SCRIPT_DELAY = 10
```

“DAQ emulator”

1. Using SPD DAQ emulator, we’ve generated 50 files, each ~2Gb;
2. Input dataset has been registered with these files;
3. Task has been processed (or 50 jobs);
4. The payload for Pilot is simple: compute the MD5/BLAKE3 hash, as there is no actual computation involved at this stage.;
5. Generation of one files takes around ~7 min, using JINR Cloud VM: 12x 1-core Intel Xeon E5-2650
6. Registration of the entire dataset: ~10 sec

```
# Configuration file for SPD DAQ data
# 2023/03/01
```

```
#Data file name format: run-<run number>--<chunk
number>--<builder id>.spd
DataFileNameFormat = run-%06u-%05u-%02u.spd
```

```
#RND generator seed:
RandomSeed = 12345
```

```
#The size limit of the output data file in bytes:
DataFileSizeLimit = 2147483648
```

```
#debug mode for debugging front-end card. If it is 1
then generator will
#produce all data words (headers and trailers) even
if there are no hits,
#otherwise all empty data blocks are removing
DebugMode = 0
```

```
#Source ID(s) of the clock module(s) for
measurement start of frame time:
FrameClockID = 1000,1001
```

```
#Source ID(s) of the TDC module(s) for measurement
of the bunch crossing time:
BunchCrossingID = 1004
```

```
#Slice length in ns (must be less than smallest TDC
over-roll time (4.5 ms for RS)):
SliceLength = 10000
```

```
#Number of slices in a frame:
FrameLength = 100000
```

Conclusion

- Unified task execution interface
- Adaptability across platforms
- Reduced overhead (scheduling & execution)
- Scalability (to millions of concurrent tasks)

Thanks for your attention!