



# Pilot Applications for Distributed Job Execution in the SPD Online Filter System

Romanychev Leonid JINR MLIT romanychev@jinr.ru

11th International Conference "Distributed Computing and Grid Technologies in Science and Education" (GRID'2025) July 7-11, 2025



### Introduction: Role of Pilot Applications

- 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

**Workload Manager** 

## 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







#### Additional scripts

Name of script	Туре	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.

11

#### Configuration of the pilot and daemon

```
pilot \geq \equiv \text{config}_{example.ini}
       [rabbit_settings]
       RABBITMQ_USERNAME=user
       RABBITMQ_PASSWORD=password
       RABBITMQ_H0ST=11.111.111.11
       RABBITMQ_EXCHANGE=jobs
       RABBITMQ_VIRTUAL_HOST=virtual_host
       RABBITMQ PORT=5672
       [node_settings]
       SERVER_ADRESS = http://11.111.111.111.8080
       PROCESSOR_TYPE = cpu
 11
       PILOT_ID = 1
 12
 13
       SCRIPT_PREFIX = /path_to_data/
 14
 15
       [logging]
       LOG\_LEVEL = INFO
 16
 17
       MAX_LOG_SIZE = 10485760 # 10MB
       BACKUP COUNT = 5
```

daemo	n > ≡ daemon_config.ini
1	[node_settings]
2	SCRIPT_DELAY = 10

## "DAQ emulator"

- 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.;
- 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 debuging 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 modulue(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!