

# Integrating DIRAC workflows in Supercomputers

Status and next steps

Alexandre Boyer - Universite Clermont Auvergne, CERN  
alexandre.franck.boyer@cern.ch



# Introduction

## **The DIRAC WMS implements the Pilot-Job paradigm**

- Able to federate a large variety of heterogeneous computing resources
- Mainly Grid Sites, Clouds, but also opportunistic resources. What about Supercomputers?

## **Supercomputers represent an important computing power**

- Different from the Grid Sites: integrating VO-specific workflows on such machines through DIRAC requires work
- Each machine is unique and the landscape quickly evolves



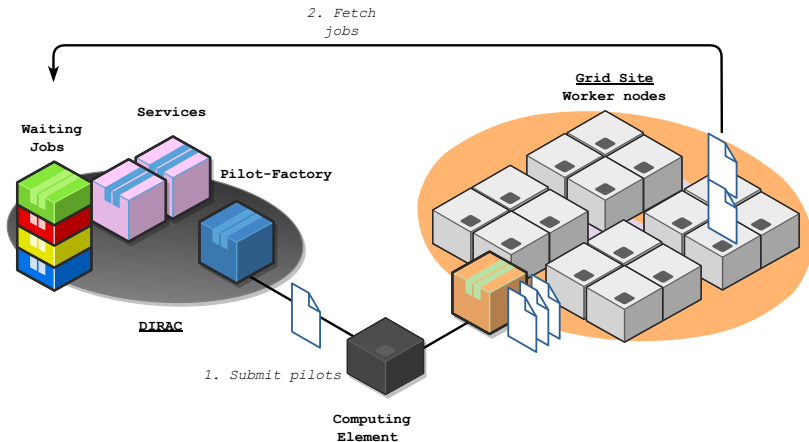
# Table of Contents

- **DIRAC WMS on Grid Sites**
- **DIRAC WMS and Supercomputers**
- **Tackling the distributed computing challenges**



# ●●●● DIRAC WMS on Grid Sites

## ●● WMS Workflow



# ●●●● DIRAC WMS on Grid Sites

## ●● "Typical" job requirements

### 4 SP Jobs running on a Grid Site

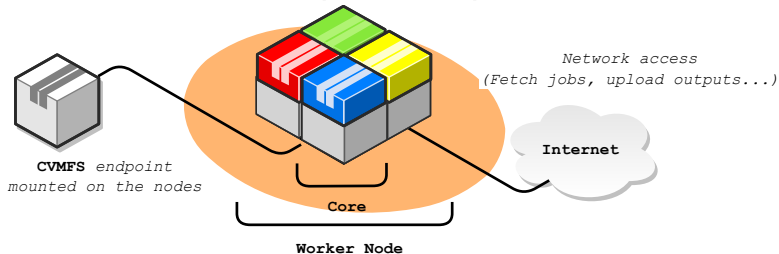
Single-Core  
allocation

x86  
architecture

SLC6/CC7  
compatible

>2Gb RAM  
per core

LRMS accessible  
from outside



# ●●●● DIRAC WMS and Supercomputers

## ●●●● Presentation

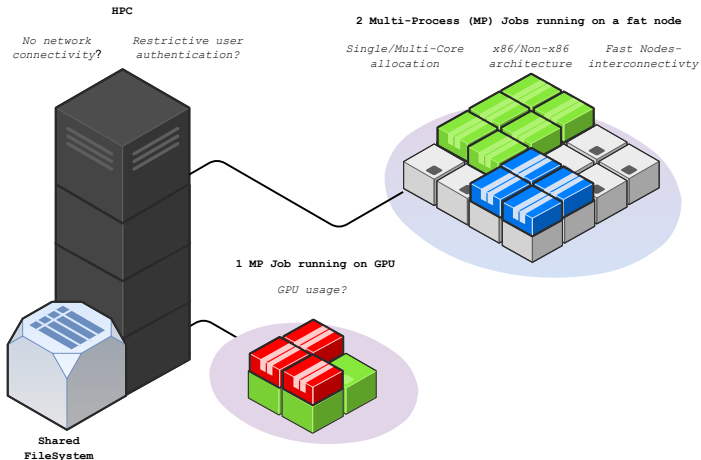
**Definition** : A mainframe computer that is among the largest, fastest, or most powerful of those available at a given time.

- Twice a year, top500.org releases the list of the most powerful SC of the world.
- #1 Fugaku is composed of ARM processors and contains  $\sim 7$  million cores
- #2 and #3 leverage IBM Power processors and Nvidia GPUs, and contain  $\sim 1.5$ -2 million cores
- In comparison, WLCG provides  $\sim 1$  million cores (many additional parameters have to be taken into account for a fair comparison though)



# ●●●● DIRAC WMS and Supercomputers

## ●●●● Features of Supercomputers



# ●●●● DIRAC WMS and Supercomputers

## ●●● Challenges

### Software architecture (VO)

- SC are many-core architecture
- They can include non x86 CPUs (ARM, AMD, Power), GPUs...
- They might contain less than 2Gb/core

⇒ SC are all made differently: hard to build a unique solution for all of them

### Distributed computing (DIRAC)

- SC policies may differ from those of HEP Grid Sites.
- They might lack of CVMFS, outbound connectivity, external access to the LRMS...





# ●●●● Tackling the distributed computing challenges

## ●●●● Overview

- 1 main variable directly affects the chosen solution (push, pull):
    - + Do WNs have an external connectivity? yes (or only via the edge node), no
  - Other variables generate variations that can be added up to the proposed solution:
    - + Is CVMFS mounted on the WNs? yes, no
    - + Is LRMS accessible from outside? yes, no
    - + What type of allocations can we make? single-core, multi-core, multi-node
- ⇒ We will go through different cases: from the easiest to the hardest one

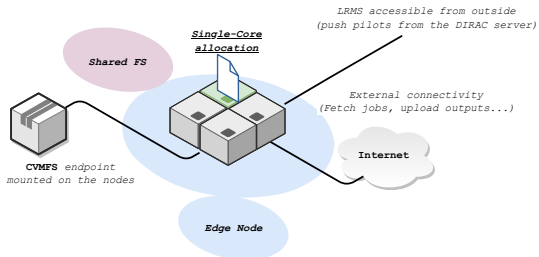


# ●●●● Tackling the distributed computing challenges

## ●●●● Pull model: single-core allocation

### Similar to a Grid Site

- Uncommon for a SC.
- Often need to collaborate with the system administrators.



# ●●●● Tackling the distributed computing challenges

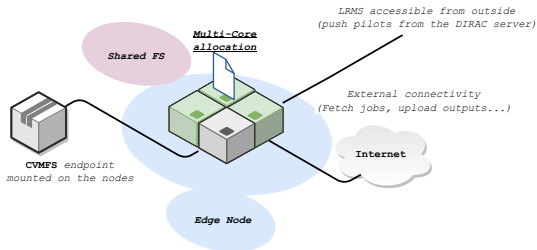
## ●●● Pull model: Multi-core allocation

Integrated since v7r0

SC often require their users to allocate many cores or even nodes to run a program (queue configuration).

### Fat node partition [3]

- One pilot per fat node: execute several SP/MP jobs per allocation.
- In the Queue conf, add:  
`LocalCETType=Pool`  
and  
`NumberOfProcessors=N`.



# ●●●● Tackling the distributed computing challenges

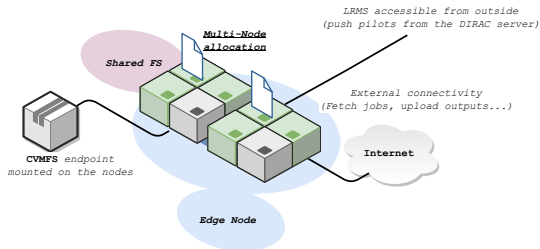
## ●●● Pull model: Multi-node allocation

Almost integrated in v7r1

Allows to get a large number of resources with a small number of allocations.

### Sub-Pilots (specific to SLURM currently)

- One sub-pilot per fat node allocated: pilots sharing a same id, status and output.
- In the Queue conf, add:  
ParallelLibrary=PL  
and  
NumberOfNodes=N<-M>.



# ●●●● Tackling the distributed computing challenges

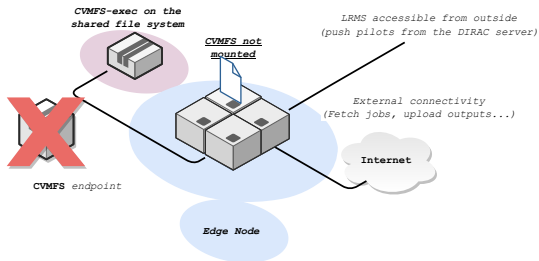
## ●●● Pull model: CVMFS not mounted on WN

Not Integrated VO action

SC, by default, do not provide CVMFS on the WNs.

### CVMFS-exec on the shared FS [2]

- Mount CVMFS as an unprivileged user.
- Purely a site/admin/VO action actually: might need to add a parameter in DIRAC to ease the process.



# ●●●● Tackling the distributed computing challenges

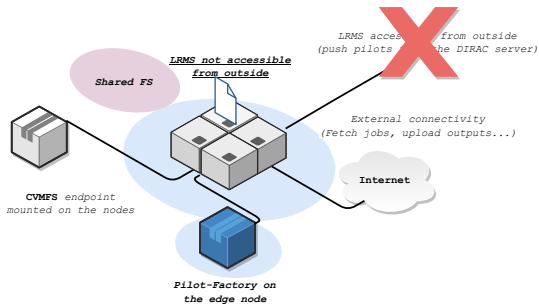
## ●●● Pull model: No remote access to LRMS

Integrated since v7r0 VO action

Some SC can only be accessed via a VPN (No CE, no direct SSH).

### Site Director on the edge node

- Directly submit pilots from the edge node.
- Would need to be allowed to execute agents on the edge node.
- Would need to be updated manually.



# ●●●● Tackling the distributed computing challenges

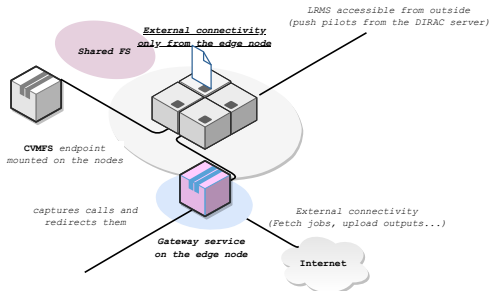
## ●●●● Pull model: Ext. connectivity only from the edge node

Not Integrated VO action

Some SC only provide external connectivity from the edge node. Pilots cannot directly interact with DIRAC services in this case.

### Gateway

- Would be installed on the edge node (if possible)
- Would capture the Pilot and Job calls and would redirect them



# ●●●● Tackling the distributed computing challenges

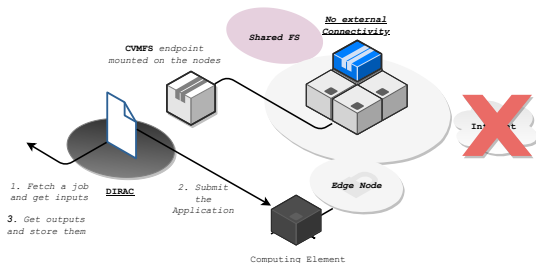
## ●●●● Push model: No Ext. connectivity

In Progress: v7r2?

Some SC do not provide any external connectivity at all, neither on the WNs or the edge node.

### PushJobAgent

- Works like a pilot outside of the SC
- Fetches jobs, deals with inputs and outputs, submits the application part to a SC
- Require a direct access to the LRMS





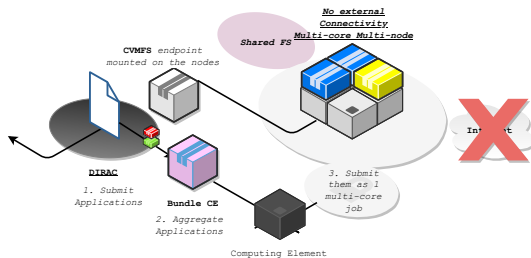
# ●●●● Tackling the distributed computing challenges

## ●●●● Push model: No Ext. connectivity, Multi-core/node

In Progress: v7r2?

### BundleCE

- Would aggregate multiple applications into one allocation



# ●●●● Tackling the distributed computing challenges

## ●●●● Push model: No Ext. connectivity, No CVMFS

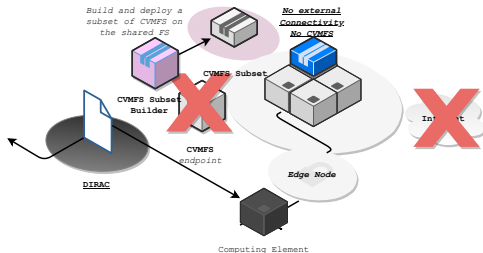
In Progress: v7r2?

VO action

As it was already said, SC do not provide CVMFS by default.  
CVMFS-exec cannot be used in this context.

### Subset-CVMFS-Builder

- Run & extract CVMFS dependencies of given jobs
- Use CVMFS-Shrinkwrapper [1] to make a subset of CVMFS
- Test it & deploy it on the SC shared FS



# Conclusion

## Status

- Support many SC with external connectivity (multi-core allocations)
- Tools to exploit SC with no external connectivity are in progress

## Next Steps

- Provide the push model solution and its variations
- Work on DB12 (CPU Power computation): support for multi-core allocations
- Provide a complete documentation about SC integration
- Provide side projects to minimize VO actions (Subset-CVMFS-Builder)



# Thanks

**Any questions? Comments?**





CVMFS. *cvmfs-shrinkwrap utility*.

<https://cvmfs.readthedocs.io/en/stable/cpt-shrinkwrap.html#cpt-shrinkwrap>. Online; accessed 4 May 2021. 2021.



CVMFS. *cvmfsexec*.

<https://github.com/cvmfs/cvmfsexec>. Online; accessed 4 May 2021. 2021.



Federico Stagni, Andrea Valassi, and Vladimir Romanovskiy. “Integrating LHCb workflows on HPC resources: status and strategies”. In: *arXiv:2006.13603 [hep-ex, physics:physics]* (June 2020). arXiv: 2006.13603. URL: <http://arxiv.org/abs/2006.13603>.



# Backup Slides

**Real use-cases we have in LHCb**



## ●●●● LHCb-supercomputers collaboration

### ●●●● Available HPCs

- Piz Daint in CSCS (Suisse)
- Marconi-A2 in CINECA (Italy) – not used anymore
- SDumont in LNCC (Brazil)
- MareNostrum in BSC (Spain)



## ●●●● LHCb-supercomputers collaboration

### ●●●● Piz Daint, CSCS

- Ranked 12th in Top500 (Nov. 2020)
  - 387,872 cores (Nov. 2020)
  - + Collaboration with the local System Administrators allows a traditional Grid Site usage
- ⇒ No change required





## ●●●● LHCb-supercomputers collaboration

## ●●●● Marconi-A2, CINECA

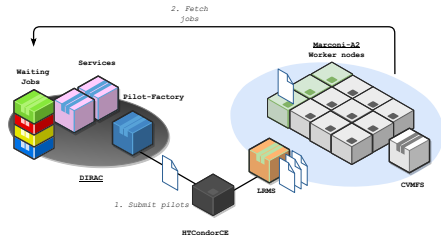
- Ranked 19th in Top500 (Nov. 2019)
- + External connectivity from the WNs
- + CVMS mounted on the WNs
- + Accessible via a CE
- 348,000 cores (Nov. 2019)
- Multi-core allocations: 272 logical cores per node (Intel KNL)
- Low memory/core



# ●●●● LHCb-supercomputers collaboration

## ●●●● Marconi-A2, CINECA: Development

- External Connectivity: Use the pull model
- Multi-core allocations: Use the *fat node partitioning* variation



## ●●●● LHCb-supercomputers collaboration

### ●●●● Marconi-A2, CINECA: Status

- ⇒ More details about LHCb work on CINECA: [3]
- ⇒ Marconi-A2 has been replaced by Marconi-100: V100 GPUs and Power9 CPUs cluster

#### Done

- Exploited 68/272 cores per node: not enough memory for more jobs

#### To be done

- Nothing to do, Marconi-A2 disappeared
- LHCb software not ready for GPUs...



## ●●●● LHCb-supercomputers collaboration

### ●●●● SDumont, LNCC

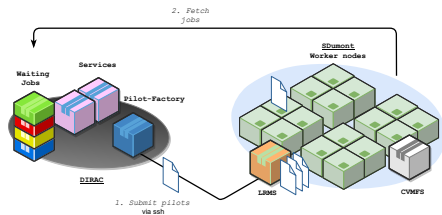
- Ranked 277th in Top500 (Nov. 2020)
- + External connectivity from the WNs
- + CVMFS mounted on the WNs
- + Accessible via SSH (special access)
- 33,856 cores (Nov. 2020)
- Multi-core allocations: 24 or 48 logical cores per node
- Multi-node allocations: 21 nodes per allocation required by some queues



# ●●●● LHCb-supercomputers collaboration

## ●●●● SDumont, LNCC: Development

- External Connectivity: Use the pull model
- Multi-core allocations: Use the *fat node partitioning* variation
- Multi-node allocations: Use the *sub-pilots* variation



●●●● LHCb-supercomputers collaboration

●●●● SDumont, LNCC: Status

### Done

- Exploit 24/24 and 48/48 cores per node

### To be done

- Multi-node allocation: should have results soon
- Dirac Benchmark: not adapted to multi-core allocations, 20% of the jobs run out of time



## ●●●●● LHCb-supercomputers collaboration

### ●●●●● Mare Nostrum, BSC

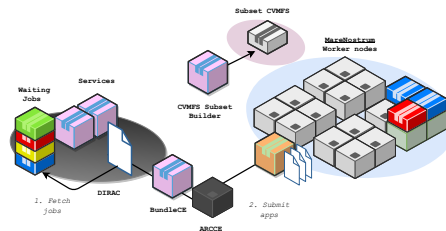
- Ranked 42<sup>nd</sup> in Top500 (Nov. 2020)
- + Accessible via a CE (ARC) and also SSH
- + Single-core allocation possible but not recommended
- 153,216 cores (Nov. 2020)
- No network connectivity
- CVMFS not mounted on the WNs



# ●●●●● LHCb-supercomputers collaboration

## ●●●●● Mare Nostrum, BSC: Development

- No external connectivity: Use the push model
- No CVMFS mounted on the WNs: Use the *Subset-CVMFS-Builder* variation
- To get multi-core allocations: Use the *BundleCE* variation





## ●●●●● LHCb-supercomputers collaboration

## ●●●●● Mare Nostrum, BSC: Status

### Done

- Prototype to run simple submissions (Hello World)
- First version of the Subset-CVMFS-Builder

### To be done

- CE configuration to run jobs within Singularity
- BundleCE to aggregate multiple jobs in an allocation

