# LHCb Production System

Federico Stagni, on behalf of the LHCbDirac team

# Overview

- ▸ What's this

- ▸ DIRAC workflows

- ▸ LHCbDirac Bookkeeping

- ▸ Application Steps and Production Requests

- ▸ Transformation System

  - ▸ LHCb extensions

- ▸ Production system in action
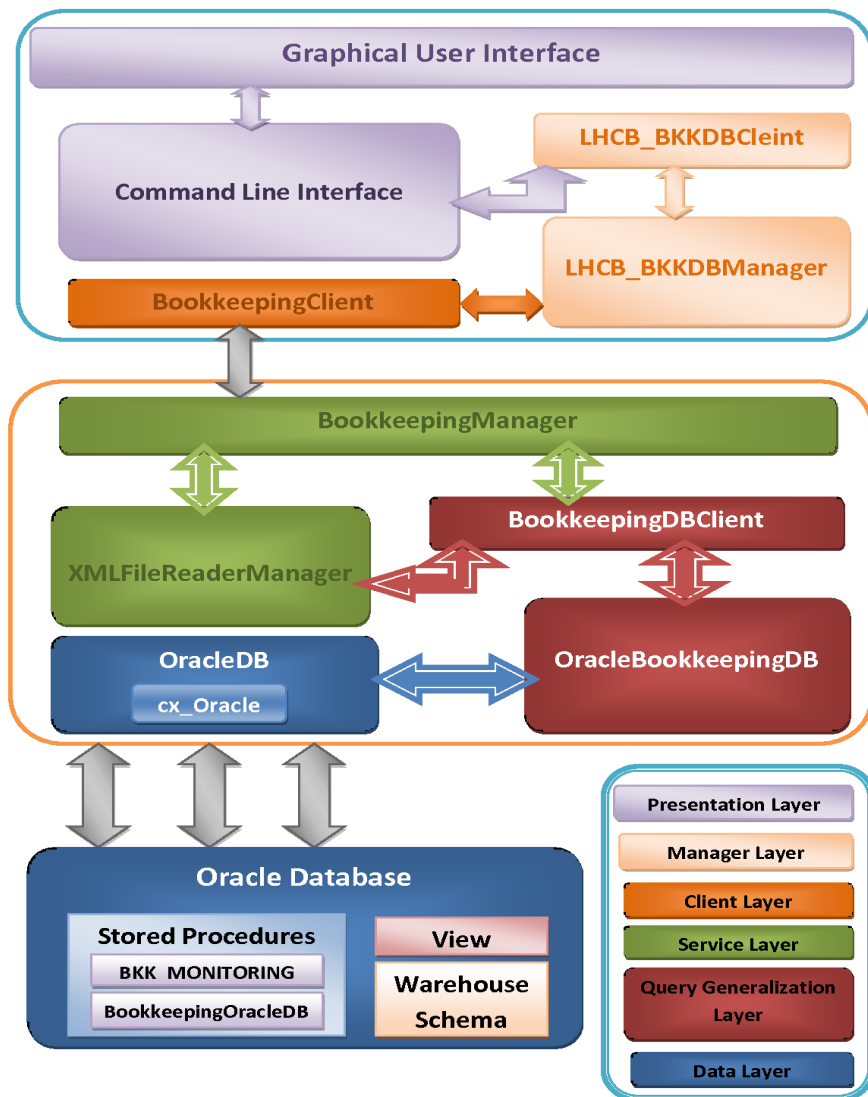
- ▸ Issues, Limitations, ToDo

# What's this

- The LHCb "Production system" is a large part of LHCbDirac

- Handles different types of productions:
  - MonteCarlo Simulations (MC)
  - Reconstruction (and Reprocessing)
  - Stripping (Selection of Physics events)
  - Merging
  - Working Groups analysis

- Not only an extension of the Transformation System

# Bookkeeping

- Data provenance and dataset retrieval

- Not necessarily a tool for distributed computing

- Retrieving datasets for:

  - users (for analysis) and production system

    - Conditions (data taking, simulation)

    - Processing (applications, detector condition parameters)

    - Event type

    - File type

- Integrated in LHCbDIRAC as a "Catalog" (LHCbDIRAC.Resources.Catalog)

- Oracle backend (DIRAC.Core.Utilities.OracleDB)

# BKK: design



- Layered design

- Focus on scalability

- One independent GUI, and a page integrated in the web portal

# BKK: GUI

# DIRAC Workflows

<span style="color:orange">DIRAC user meeting 2011</span> (click!)

▸ Xml ↔python dict

▸ A workflow connects steps together

▸ dirac-jobexec aWorkflow.xml

Jdl:

▸ Executable = "$DIRACROOT/scripts/dirac-jobexec";

▸ Arguments = "jobDescription.xml -o LogLevel=verbose";

# Production Request System



- ▸ Application Managers defines application steps

- ▸ "What to run" to go from X to Y

- ▸ LHCb application

- ▸ A step "translates" in a workflow application step

# Production Request System /2



▸ Steps are combined in production requests (e.g. MC, or Reconstruction)

## Generate production script

### Please specify Production parameters

| Parameter ▲ | Value |
|---|---|
| GENERAL: Set True for EXPRESS (Run at C… | False |
| GENERAL: Set True for certification test | False |
| GENERAL: Set True for local test | False |
| GENERAL: Set True to create validation pro… | False |
| GENERAL: Use Oracle | True |
| GENERAL: Workflow string to append to pr… | 1 |
| GENERAL: Workflow system config e.g. x8… | ANY |
| PROD-RECO: DataReconstruction or DataRe… | DataReconstruction |
| PROD-RECO: Group size or number of files … | 1 |
| PROD-RECO: Max CPU time in secs | 1000000 |
| PROD-RECO: Number of Files | -1 |
| PROD-RECO: Output Data Storage Element | Tier1-RDST |
| PROD-RECO: ancestor production if any | 0 |
| PROD-RECO: dicrete list of run numbers (do… | |
| PROD-RECO: distribute output data True/Fal… | False |
| PROD-RECO: priority | 7 |
| PROD-RECO: production plugin name | AtomicRun |
| PROD-RECO: run end, to set the end of the … | 0 |
| PROD-RECO: run start, to set the start run | 0 |

« Previous | Next » | Generate | Preview | ScriptPreview | Cancel

▸ Production requests are submitted using production templates

▸ e.g.: priority, which plugin, where the outputs are stored, DIRAC CPU, etc.

▸ Each production is created using the Production API

# LHCbDirac TS

- Extension of the DIRAC TS, mostly for interacting with the BKK

  - DB:

    - Physics RUNs information

    - BKK queries (supersedes TransformationInputDataQuery)

  - Service and clients extended for the DB extension

- Agents

  - BookkeepingWatchAgent

    - Looks for BKK queries, and fills the TransformationFiles table

    - Threaded, uses pickle file for caching

  - DataRecoveryAgent

    - Resets input files in "Unused" status, in case the jobs failed

    - A counter is kept, with a maximum of re-trials

  - Extensions for cleaning, and closing productions

- Plugins (LHCbDIRAC.TransformationSystem.Agent.TransformationPlugins)

  - Many LHCb plugins coded

    - e.g. ByRun, AtomicRun, with flushing...

    - This is where you want to extend

- TaskManager

  - Extended to handle the inputs the LHCb way

# LHCb Transformation Monitor Web

- Expose functionalities to connect together TS, BKK and Production Request System

- Use LHCbJob.py (extension of DIRAC.Interfaces.API.Job.py) to create a DIRAC workflow, whose xml is uploaded to the Transformation DB

- python modules are run within the workflow, grouped within steps. Application steps and Finalization steps are present

Request

App Step 1
App Step 2
⋮
App Step n

Prod 1
Prod 2

Task 1 ⸺ Job 1
Task 2 ⸺ Job 2
Task 3 ⸺ Job 3
⋮
Task m ⸺ Job m

App Step 1
App Step 2
Finalization Step

Registered in the BKK

Production Templates

Transformation

DIRAC WF XML description

# Distributed Computing Coordinator

**PPG/OPG**

Physics and Operations Planning Group, follow activities from high level

**Application Managers**

Core software, application releases

**Production Managers**

Run, control and coordinate most of the production activities

**Productions Monitor**

**DQ**

Data Quality shifter

**Computing shifter**

Mainly look after productions activities

**Data Manager**

Data issues, distribution, removal

**GEOC**

Grid Expert On Call: referes to WLCG, run daily operations meeting, 1st contact for shifters

**Resources Manager**

Plans for resources

**DAST**

Distributed Analysis Support shifter

**Sites testing and monitoring**

**Dirac Experts**

Contact for all DIRAC-related issues

**DIRAC Config**

**Tier1 contacts**

A contact person for each Tier 1

**Sites Manager**

Mostly handles T2s

Taking the LHCb Reprocessing 2012 as example

# The Plan (source: Stefan Roiser)

- Reprocessing from mid Sep until Xmas break
  - In this time process all 2012 data up to Sep TS
  - Using a simplified workflow, only 1 output file of reco
- Operations on T1 sites + 20 % attached T2s (as in 2011)
- First pass processing at CERN + 2 T2s

# Let's think BIG (source: Stefan Roiser)

- Can we expand this model by using more Tier2s?

  – Preferably choose CVFMS sites

  – Because of easier software management (e.g. conditions deployment)

  – No preference on "site power"

  – Questions

  – Will the T1 storage sustain the load?

  – Network congestions?



2012

?

1/3 more Tier 2 sites (total 37 out of 86)

CVMFS site  Non-CVMFS site

# How?

- We "attached" T2s to T1s:
  - (output): /Resources/Sites/<yourSite>/AssociatedSEs
  - (input): /Operations/SiteLocalSEMapping

(some mix between Resources and Operations)

For scheduling we intervene on:

- /Operations/JobScheduling/MatchingDelay
- /Operations/JobScheduling/RunningDelay

# Current Status  (source: Stefan Roiser)

Running reconstruction jobs at T2s
during 2011 and 2012 reprocessing



2012

10k



2011

5k

All grid activities since start of the reprocessing campaign

- Running at "full steam"

- Enough room for user jobs

- Less Simulation jobs

- Peaks of 30k running jobs (usual 25k)

# Problems (source: Stefan Roiser)

- It's not all singing and dancing
  - Mostly scaling issues both at (LHCb)Dirac and sites
  - Because of high load, comparable to data post LS1
- Data moving is >just< following at some sites
  - Staging of tape data, moving between storages
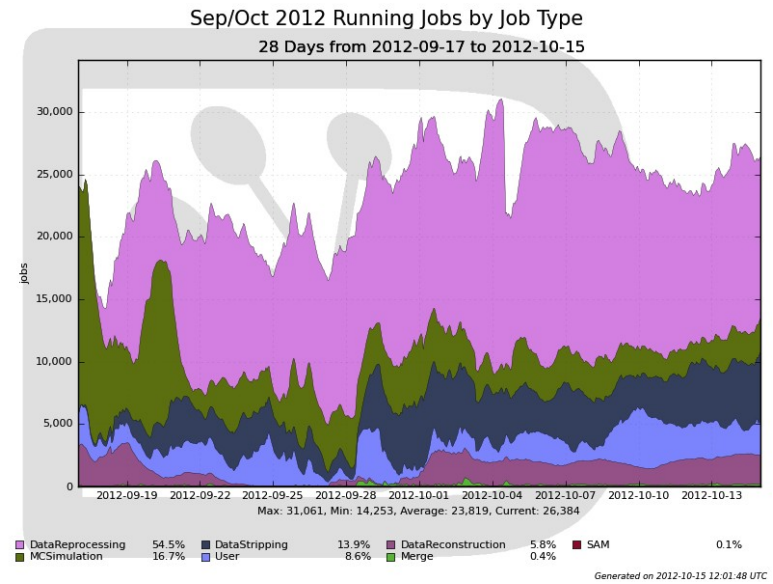- Several adjustments to Dirac to handle huge productions
  - Several module changes needed to speedup for handling large amount of jobs in a short period
  - Scaling also by multiplying data handling agents
- Production priorities didn't really seem to work!

# Handling the scalability issues
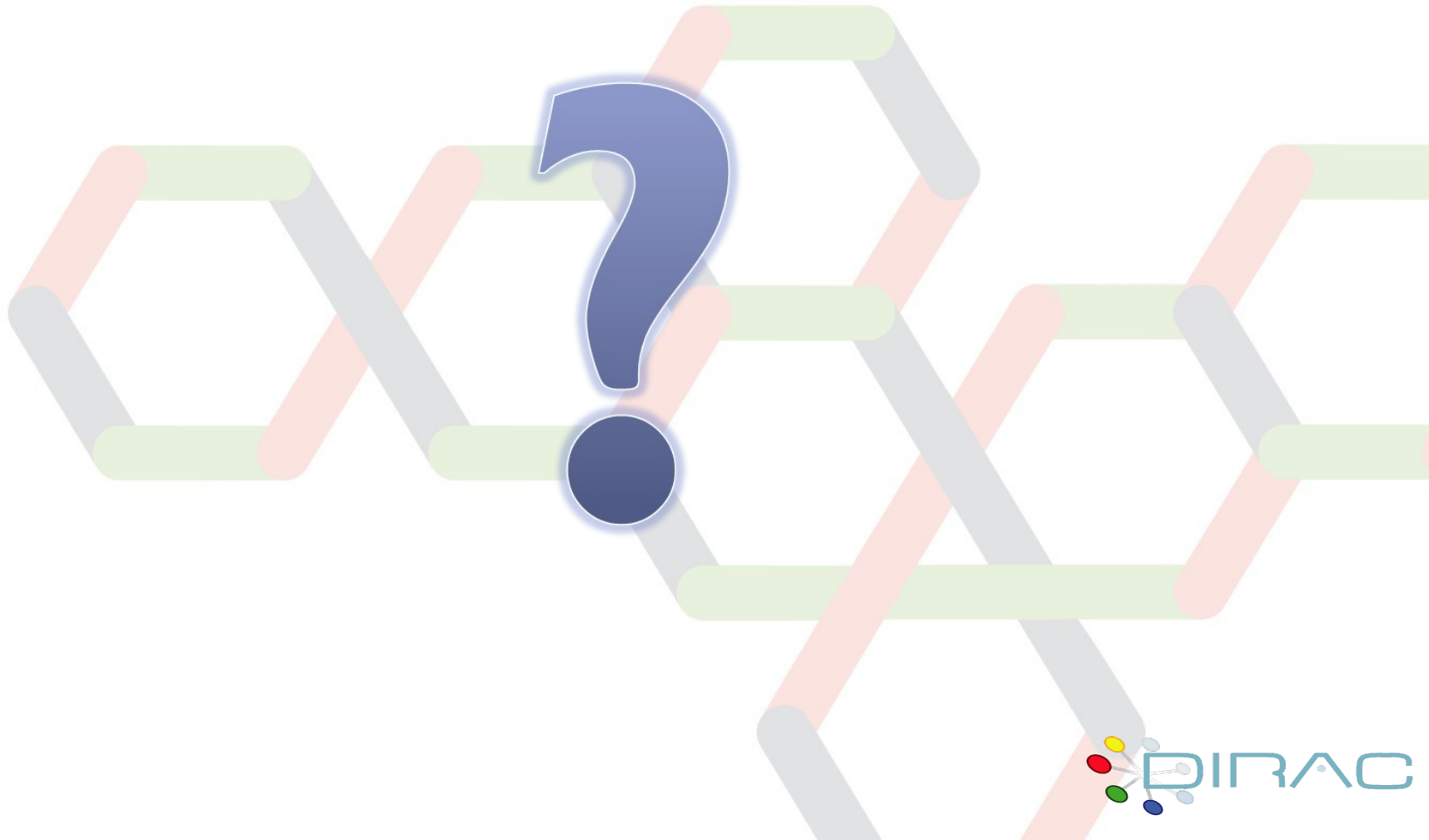
▷ e.g. productions with ~400k input files

▷ 2 TransformationManager and 2 BookkeepingManager (services) running in parallel on different machines

▷ The TransformationClient splits large queries in peaces.

▷ Large use of caching in TransformationAgent and BookkeepingWatchAgent

- Core, Framework

  - Use MySQL Transactions? SQLAlchemy?

- Request System

  - Complete review started... maybe using the executors framework?

  - Few steps made for the monitoring

- Scripts

  - Still feeling like if there is a little mess...  review?

- Testing and Certification

  - Made some steps in LHCb, adoption in DIRAC (jenkins)?

  - Documentation, documentation, documentation...

# Questions

# Setups

- DIRAC setups:

  - Development, Certification, Production

- LHCbDirac Production setup:

  - 1 machine for agents

  - 1 machine for security services

  - 2 machines for services

  - 1 for DIRAC SEs (logs, and Sandboxes)

  - 2 web servers (1 at CERN, 1 in Barcelona)

  - 6 MySQL DB machines (3 for accounting)

  - An Oracle DB maintained by CERN IT

  - 1 machine for special needs