

@GridPP

@twhyntie

GridPP Infrastructure and Approaches

T. Whyntie*

** Queen Mary University of London*

SKA/GridPP F2F, University of Manchester

Wednesday, 2nd November 2016



Outline of the talk

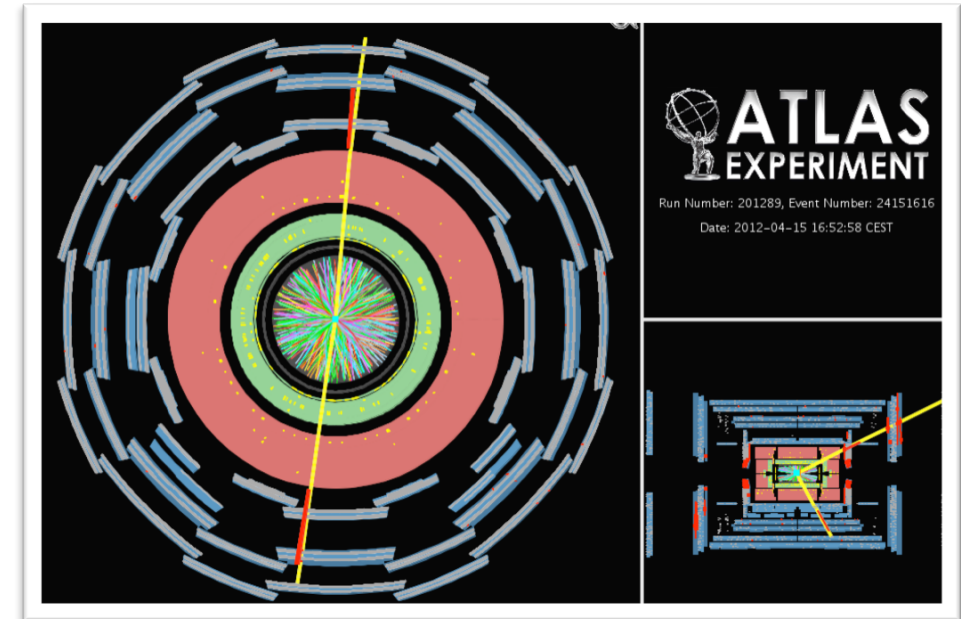
- GridPP: an introduction:
 - *GridPP and the Worldwide LHC Computing Grid (WLCG); But what is a Grid? When is a Grid useful? Can I use the Grid?*
- Engaging with the Grid:
 - *Infrastructure for non-LHC VOs; documentation; advances approaches for large VOs.*
- Selected case studies:
 - *GEANT4 simulation campaigns; GalDyn (Galaxy Dynamics); the Large Synoptic Survey Telescope.*
- Summary and conclusions.

GridPP: an introduction

GridPP and the Worldwide LHC Computing Grid (WLCG); But what is a Grid? When is a Grid useful? Can I use the Grid?

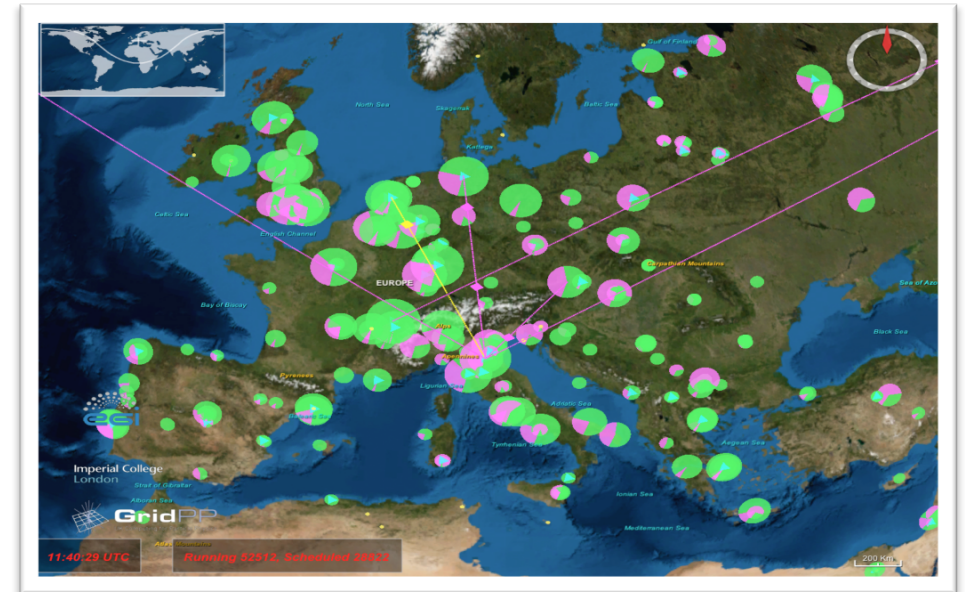
GridPP and the WLCG

- The Worldwide LHC Computing Grid (WLCG) was developed to meet the challenges presented by data from CERN's Large Hadron Collider (LHC):
 - *40 million particle collisions per second;*
 - *150 million channels in ATLAS/CMS detectors;*
 - *At least 15 PB of data per year;*
 - *Expect a few per million of e.g. Higgs events.*
- GridPP (the UK Grid for Particle Physics) represents the UK's contribution:
 - *A collaboration of 20 institutions, 100+ people;*
 - *101k logical CPU cores, 37 PB storage;*
 - *Accounts for ~11% of WLCG resources.*

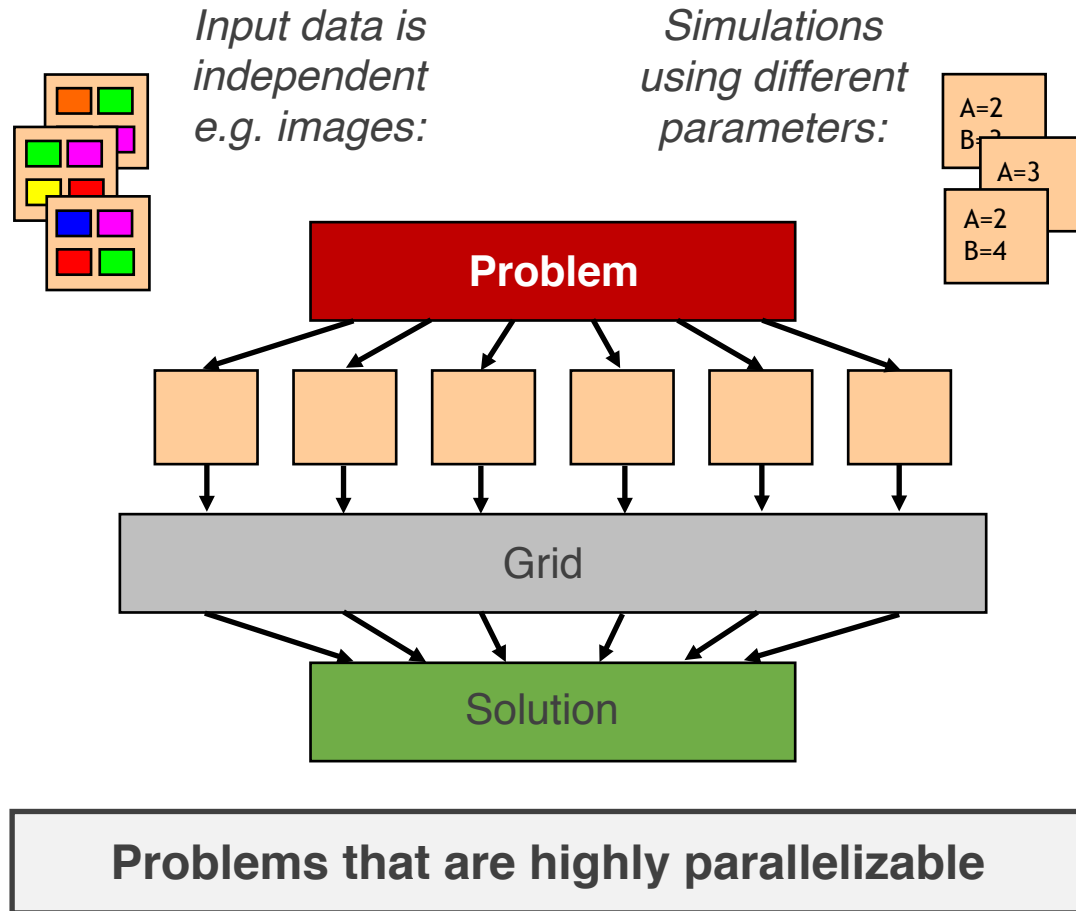


But what is a Grid?

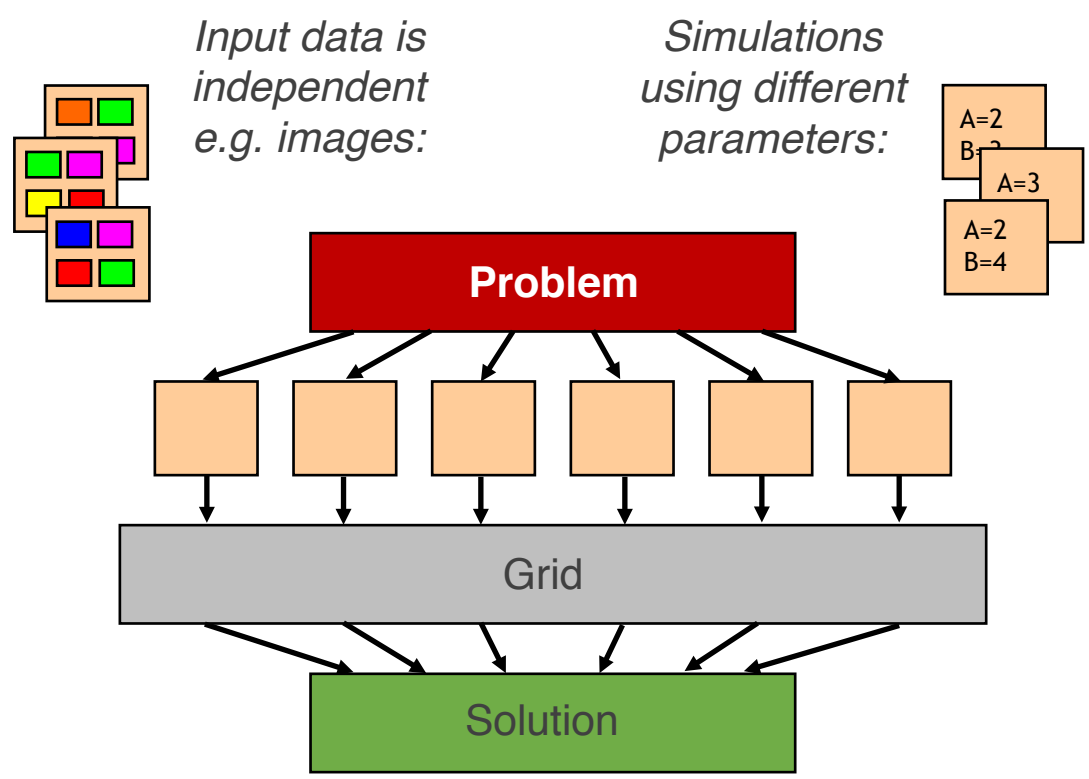
- ‘**Grid**’ computing – key concepts:
 - *All processing power and storage always on, always available, whenever it is needed;*
 - *The end user doesn’t know or care about where or how (c.f. electricity grid) thanks to middleware technologies.*
- The **WLCG** itself is distributed computing/High Throughput Computing (HTC) on a huge scale:
 - *As of August 2013, 152 sites in 36 countries, 365k logical CPUs, 210 PB storage;*
 - *By number of cores (not a fair measure), it would rank 3rd in the current top-10 super computers worldwide;*
 - *As acknowledged by CERN DG, crucial in the discovery of the Higgs boson in July 2012.*



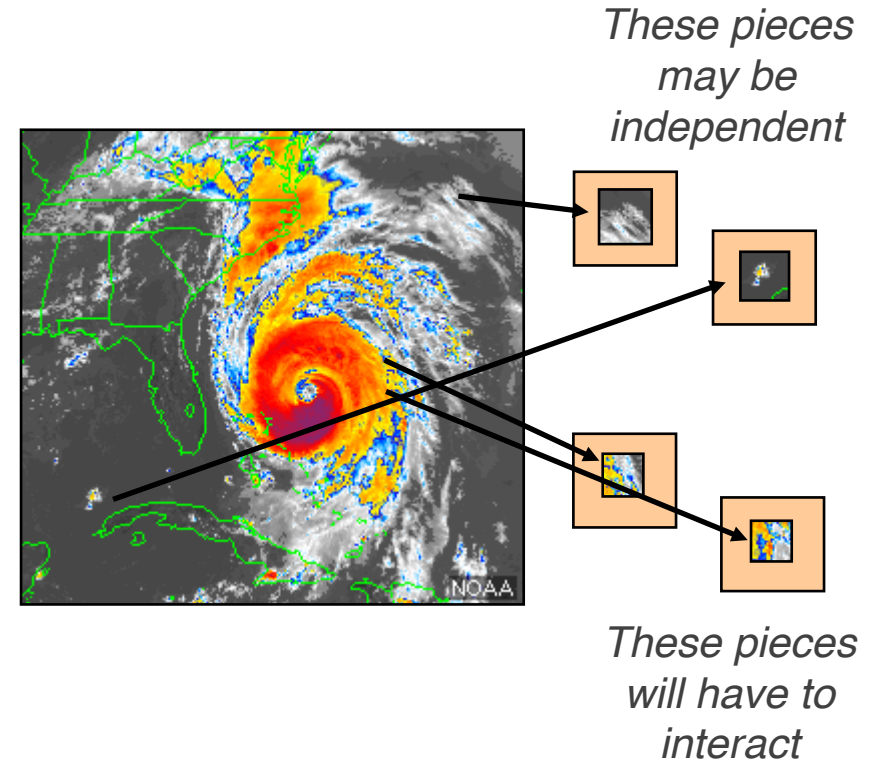
When is a Grid useful?



When is a Grid useful?



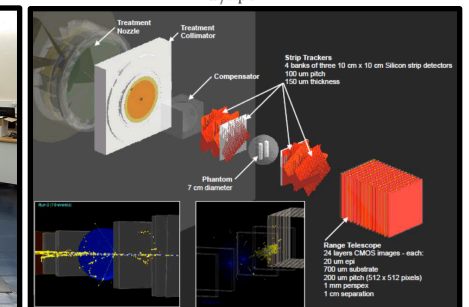
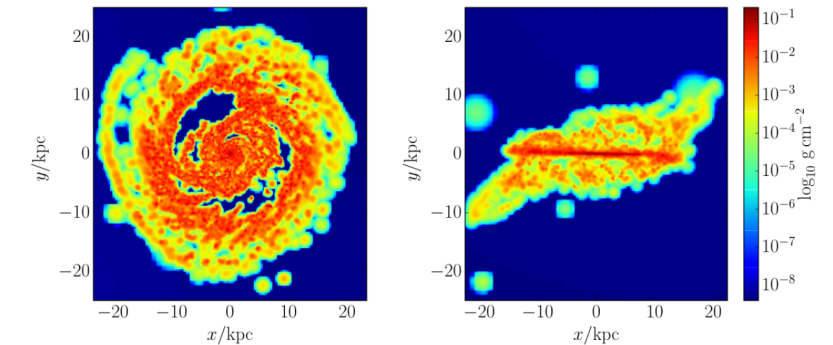
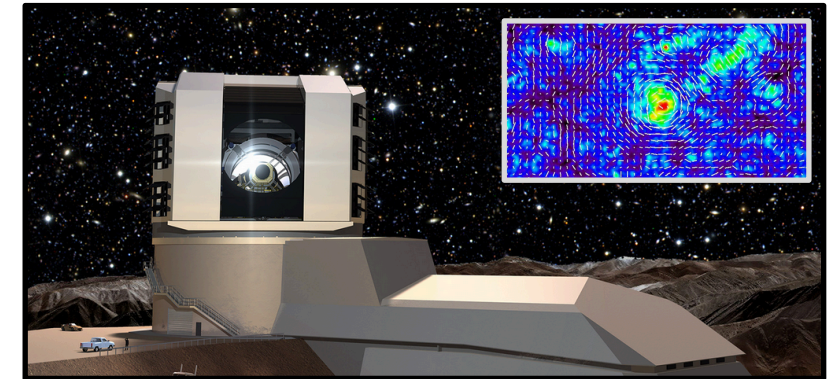
Problems that are highly parallelizable



Not so good for closely coupled problems

Can I use the Grid?

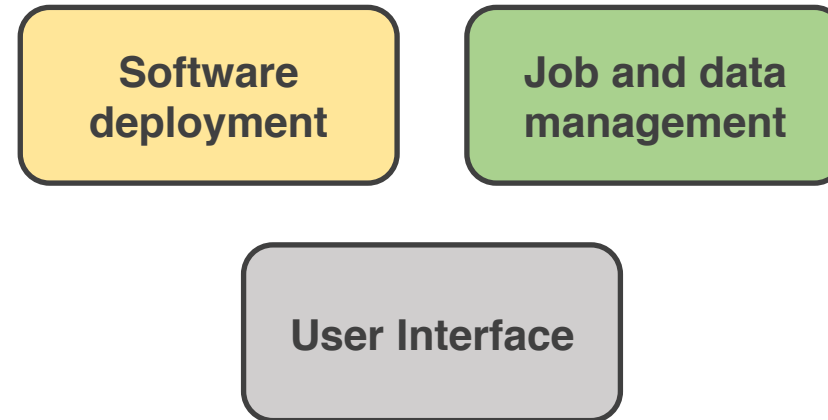
- Yes! GridPP offers up to 10% of its resources to non-LHC Virtual Organisations (VOs):
 - VO – community of users grouped by e.g. experiment;
 - All users need a Grid certificate (X509);
 - Problem: smaller groups tend to lack the resources needed to develop the infrastructure required to engage with the Grid, e.g. Uis, middleware framework.
- GridPP's New User Engagement Programme:
 - Provide infrastructure, tools and documentation for engaging with GridPP resources;
 - Present a standardised approach for small VO's.
- Key components:
 - A multi-VO job and data management system (DIRAC), use of a software distribution system (CVMFS), and standard User Interface (Ganga, CernVM) for small/non-LHC user communities;
- Results – many non-LHC/non-HEP user communities now using GridPP resources.

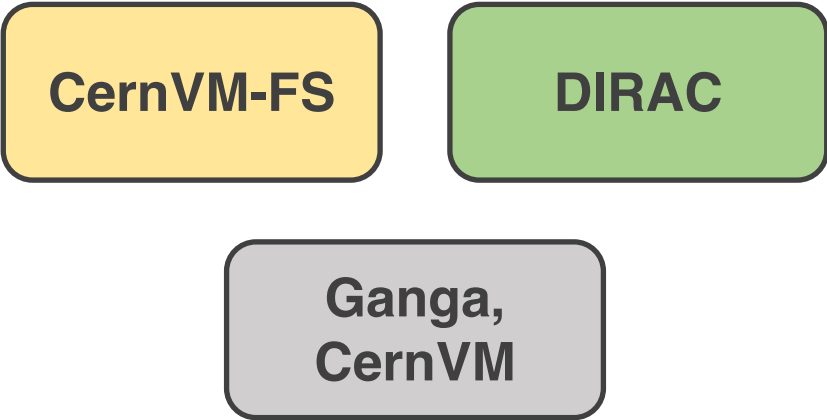


Engaging with the Grid

Infrastructure for non-LHC VOs; documentation; advanced approaches for larger VOs.

Infrastructure for non-LHC VOs





Distributed Infrastructure with Remote Agent Control
A software framework for distributed computing with grid resources.

See <http://diracgrid.org/>

CernVM-FS

DIRAC

**Ganga,
CernVM**

CernVM File System
A network file system for delivering experiment software in a scalable, fast, reliable way via http.

See the website [here](#).

Distributed Infrastructure with Remote Agent Control
A software framework for distributed computing with grid resources.

See <http://diracgrid.org/>

CernVM-FS

DIRAC

**Ganga,
CernVM**

CernVM File System
A network file system for delivering experiment software in a scalable, fast, reliable way via http.

See the website [here](#).

Distributed Infrastructure with Remote Agent Control
A software framework for distributed computing with grid resources.

See <http://diracgrid.org/>

CernVM-FS

DIRAC

**Ganga,
CernVM**

Supported by Imperial College London, Uni. Birmingham:
<http://ganga.readthedocs.io>

<http://cernvm.cern.ch/>

Ganga and the GridPP CernVM
Ganga is a Python-based interface for distributed computing. The CernVM is a baseline Virtual Software Appliance developed for the participants of CERN LHC experiments with built-in CVMFS access.

RAL hosts a CernVM-FS Stratum 0 for non-LHC VO software repositories.

See the website [here](#).

The GridPP DIRAC instance is hosted and supported by Imperial College London.

See <https://dirac.gridpp.ac.uk/>

CernVM-FS

DIRAC

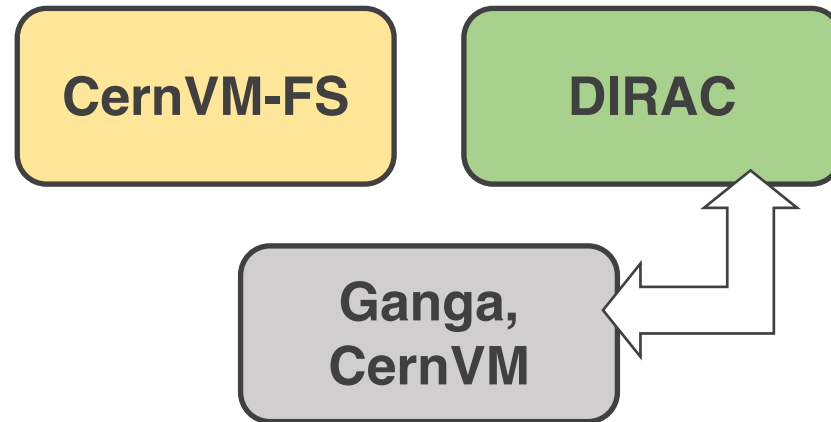
**Ganga,
CernVM**

*Supported by Imperial College London, Uni. Birmingham:
<http://ganga.readthedocs.io>*

<http://cernvm.cern.ch/>

Ganga and the GridPP CernVM

Ganga is available to anyone with CVMFS access via a grid-enabled cluster or a GridPP CernVM. GridPP offers a contextualised CernVM suitable for new users if cluster access cannot be arranged.



Ganga is a Python-based interface for distributed computing with local, batch, or grid running.

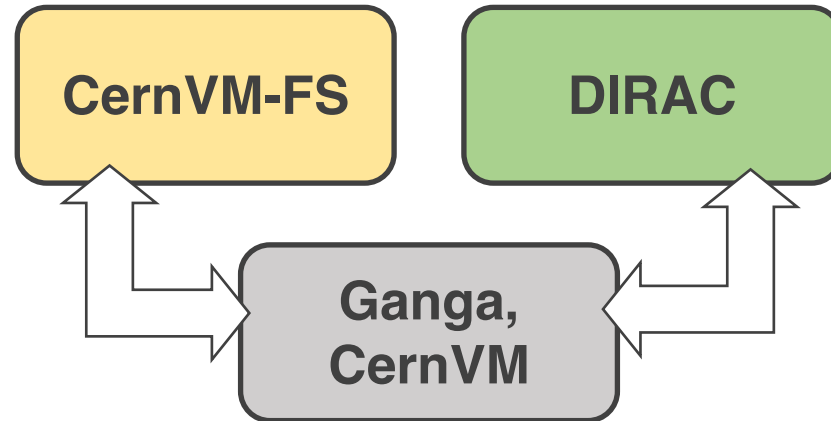
Ganga can be configured to use the GridPP DIRAC system as a backend.

Switching from local to batch to grid is trivial thanks to CVMFS, making local testing before Grid running easy.

Users can build their software (executables and libraries) on their local SL6 cluster or GridPP CernVM ready for deployment to the grid.

Users can upload their software to their own CernVM-FS repository from using the `gsi` tools in the repository `/cvmfs/grid.cern.ch`*

Custom CernVM-FS repository software can be accessed from CVMFS-enabled cluster or a contextualised GridPP CernVM.



Ganga is a Python-based interface for distributed computing with local, batch, or grid running.

Ganga can be configured to use the GridPP DIRAC system as a backend.

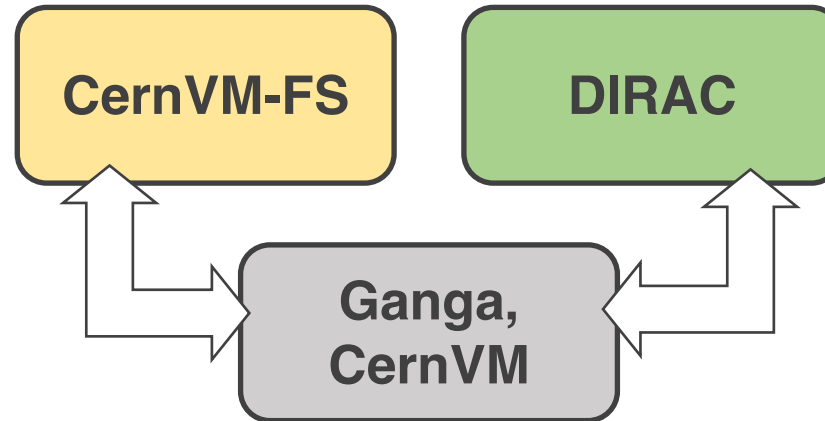
Switching from local to batch to grid is trivial thanks to CVMFS, making local testing before Grid running easy.

Users can build their software (executables and libraries) on their local SL6 cluster or GridPP CernVM ready for deployment to the grid.

Users can upload their software to their own CernVM-FS repository from using the `gsi` tools in the repository `/cvmfs/grid.cern.ch`*

Custom CernVM-FS repository software can be accessed from CVMFS-enabled cluster or a contextualised GridPP CernVM.

Combined with Ganga, this gives the user the ability to run local/batch jobs immediately (e.g. for testing) – without needing a grid certificate.



Ganga is a Python-based interface for distributed computing with local, batch, or grid running.

Ganga can be configured to use the GridPP DIRAC system as a backend.

Switching from local to batch to grid is trivial thanks to CVMFS, making local testing before Grid running easy.

DIRAC, Ganga and experiment software can be deployed via CernVM-FS for local running.

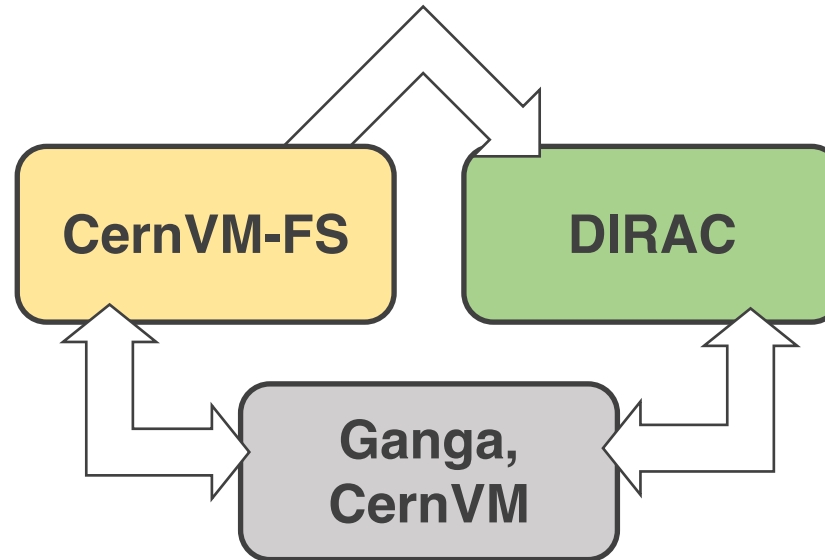
User software in (custom) CernVM-FS repositories can be used by Grid jobs managed by the GridPP DIRAC system (i.e. sites with CVMFS access).

Users can build their software (executables and libraries) on their local SL6 cluster or GridPP CernVM ready for deployment to the grid.

Users can upload their software to their own CernVM-FS repository from using the `gsi` tools in the repository `/cvmfs/grid.cern.ch`*

Custom CernVM-FS repository software can be accessed from CVMFS-enabled cluster or a contextualised GridPP CernVM.

Combined with Ganga, this gives the user the ability to run local/batch jobs immediately (e.g. for testing) – without needing a grid certificate.



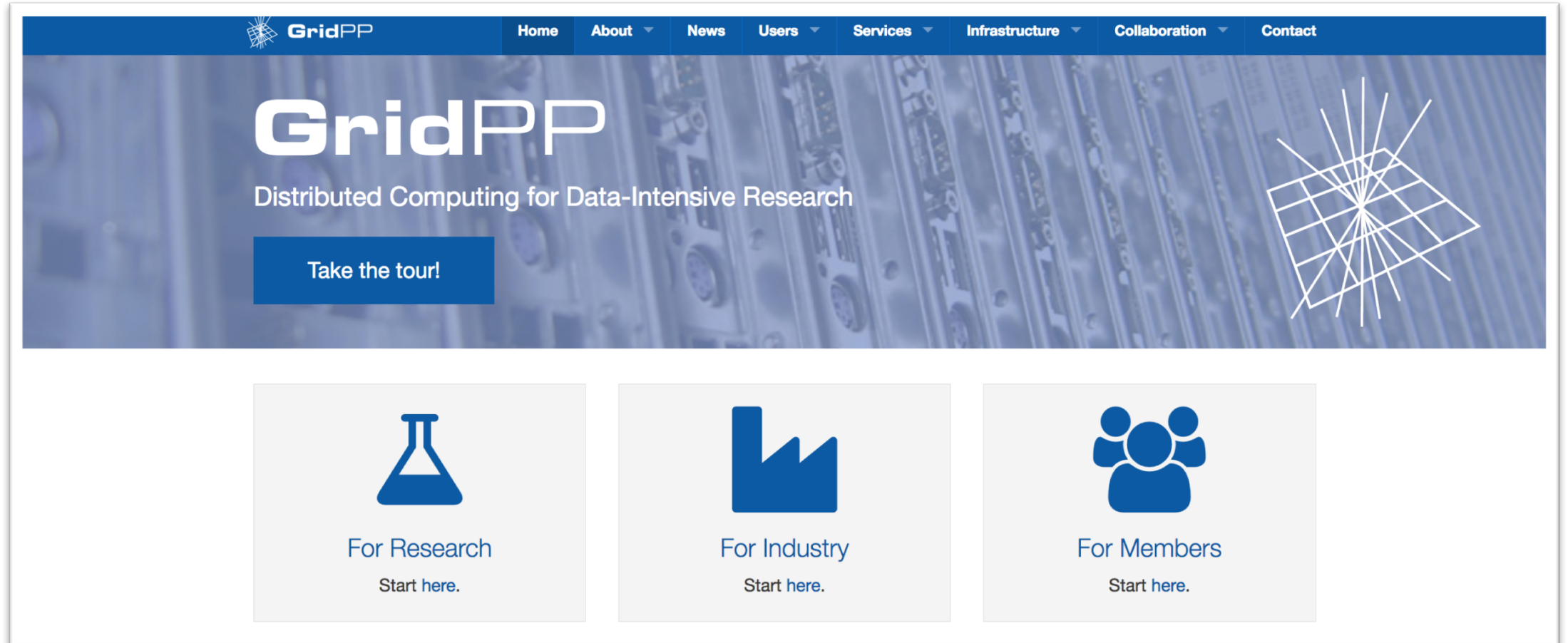
Ganga is a Python-based interface for distributed computing with local, batch, or grid running.

Ganga can be configured to use the GridPP DIRAC system as a backend.

Switching from local to batch to grid is trivial thanks to CVMFS, making local testing before Grid running easy.

DIRAC, Ganga and experiment software can be deployed via CernVM-FS for local running.

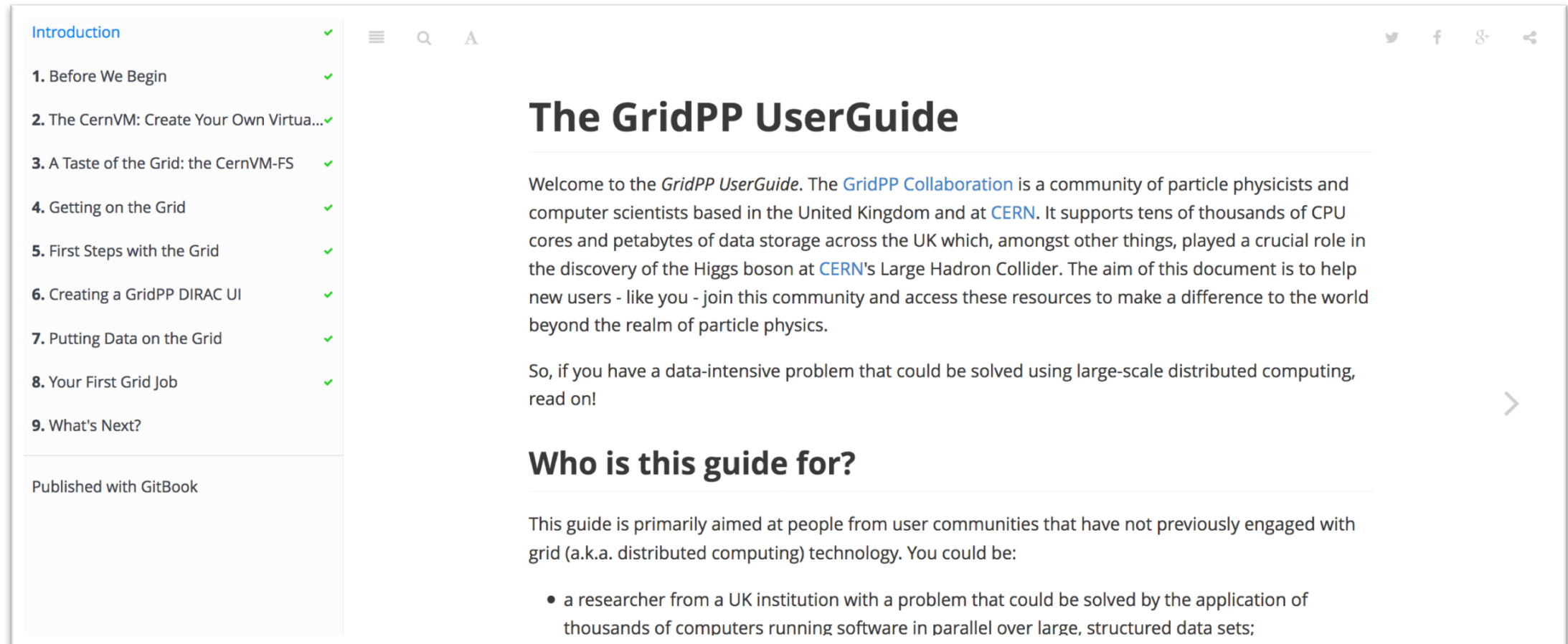
Documentation



The screenshot shows the GridPP website homepage. At the top is a dark blue navigation bar with the GridPP logo on the left and menu items: Home, About, News, Users, Services, Infrastructure, Collaboration, and Contact. Below the navigation bar is a large hero section with a background image of server racks. The hero section contains the GridPP logo, the tagline "Distributed Computing for Data-Intensive Research", and a blue button that says "Take the tour!". To the right of the hero section is a white grid icon. Below the hero section are three light gray boxes, each with a blue icon and text: "For Research Start here." (with a flask icon), "For Industry Start here." (with a factory icon), and "For Members Start here." (with a group of people icon).

<http://www.gridpp.ac.uk>

Documentation



The screenshot shows a web page for the GridPP UserGuide. On the left is a navigation sidebar with a table of contents. The main content area has a title, a welcome message, and a section titled 'Who is this guide for?' with a list of target users. The page includes a search bar, social media icons, and a 'Published with GitBook' notice.

Introduction ✓

1. Before We Begin ✓
2. The CernVM: Create Your Own Virtua... ✓
3. A Taste of the Grid: the CernVM-FS ✓
4. Getting on the Grid ✓
5. First Steps with the Grid ✓
6. Creating a GridPP DIRAC UI ✓
7. Putting Data on the Grid ✓
8. Your First Grid Job ✓
9. What's Next? ✓

Published with GitBook

The GridPP UserGuide

Welcome to the *GridPP UserGuide*. The [GridPP Collaboration](#) is a community of particle physicists and computer scientists based in the United Kingdom and at [CERN](#). It supports tens of thousands of CPU cores and petabytes of data storage across the UK which, amongst other things, played a crucial role in the discovery of the Higgs boson at [CERN's](#) Large Hadron Collider. The aim of this document is to help new users - like you - join this community and access these resources to make a difference to the world beyond the realm of particle physics.

So, if you have a data-intensive problem that could be solved using large-scale distributed computing, read on!

Who is this guide for?

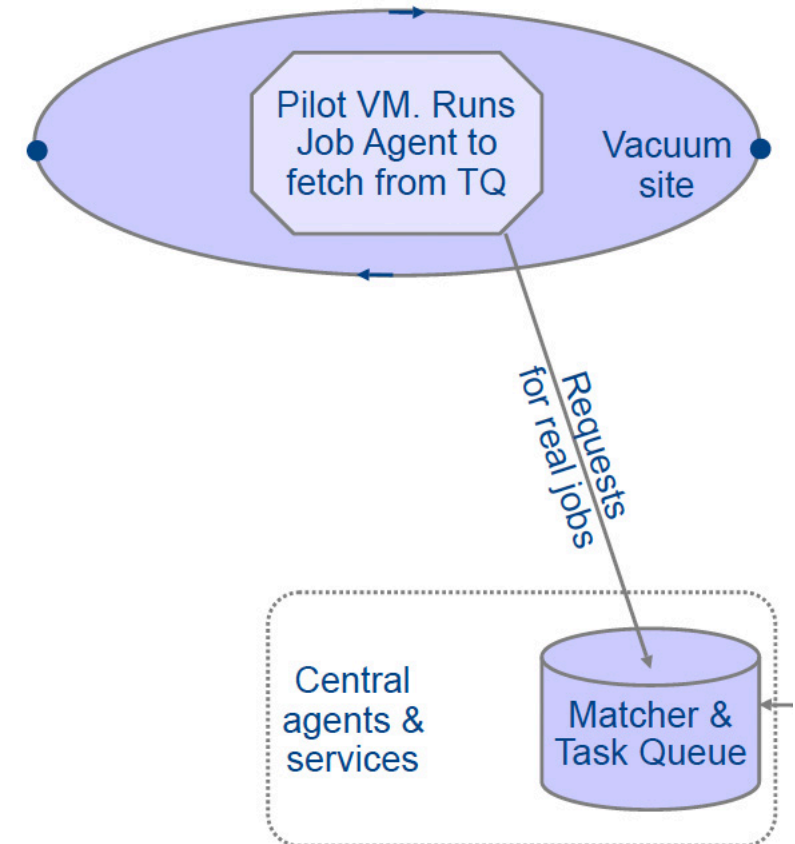
This guide is primarily aimed at people from user communities that have not previously engaged with grid (a.k.a. distributed computing) technology. You could be:

- a researcher from a UK institution with a problem that could be solved by the application of thousands of computers running software in parallel over large, structured data sets;

<http://www.gridpp.ac.uk/userguide>

Advanced approaches for larger VOs

- GridPP DIRAC is optimised for multiple, smaller VOs with HTC in mind;
- Larger VOs (i.e. those with resources for development work) could take advantage of additional DIRAC functionality:
 - *e.g. LHCb DIRAC.*
- Some possibilities:
 - *Rather than jobs, create custom Virtual Machines/clusters on the fly – the VAC model;*
 - *Integration with HPC systems;*
 - *Dedicated storage areas and functionality.*
- See Andrew McNab's talk...

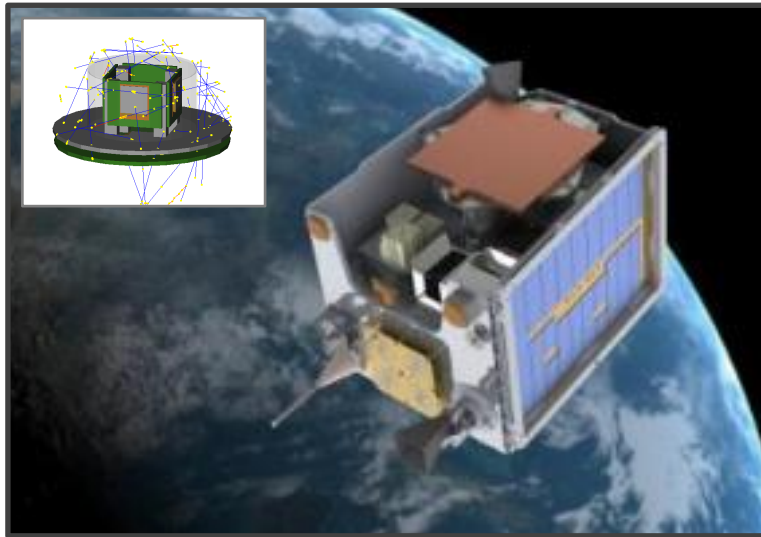


Selected case studies

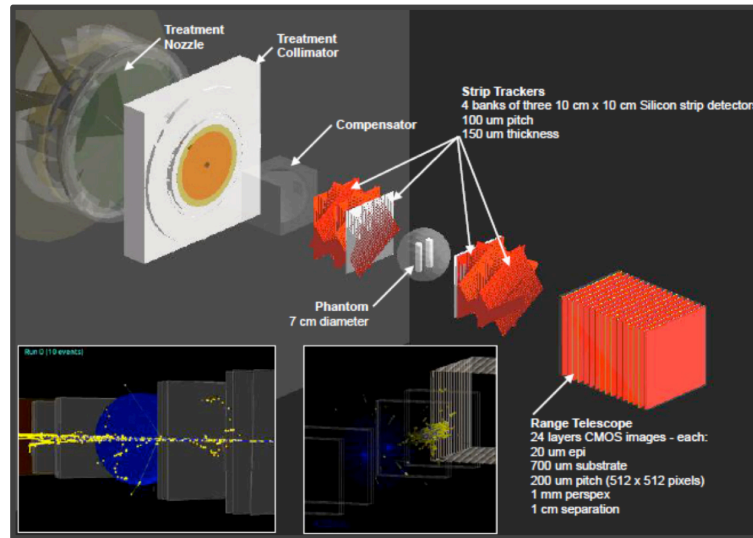
GEANT4 simulations; GalDyn (Galaxy Dynamics); LSST.

GEANT4 simulation campaigns

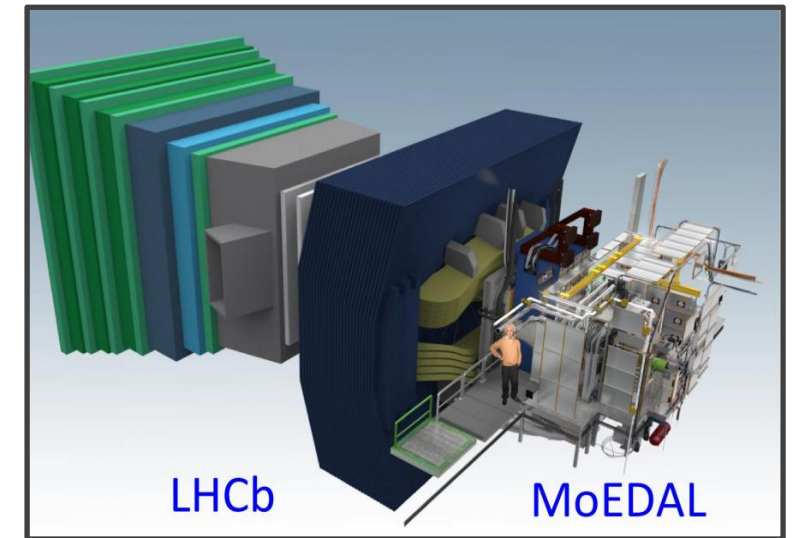
The Grid is ideal for running large-scale simulation campaigns over a large parameter space, e.g. GEANT4 particle transport simulations requiring millions of individual (independent) events. Some examples:



LUCID: estimation of data rates in Low Earth Orbit (LEO).



PRaVDA: optimisation of new proton therapy systems.

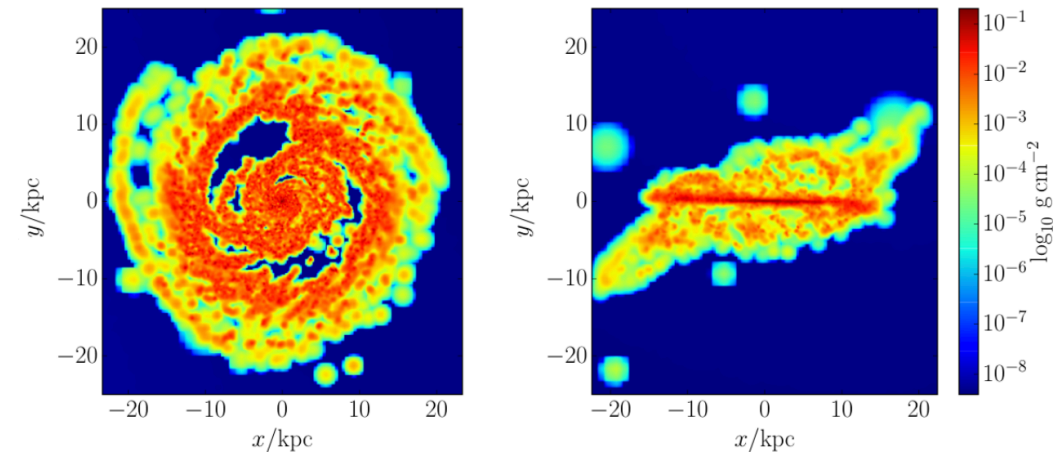


MoEDAL: detector acceptances for magnetic monopole searches.

In each case, GridPP DIRAC was used to manage thousands of jobs and data from millions of events across multiple Grid sites, and users reported reductions in running times of **months/weeks** to a **few days**.

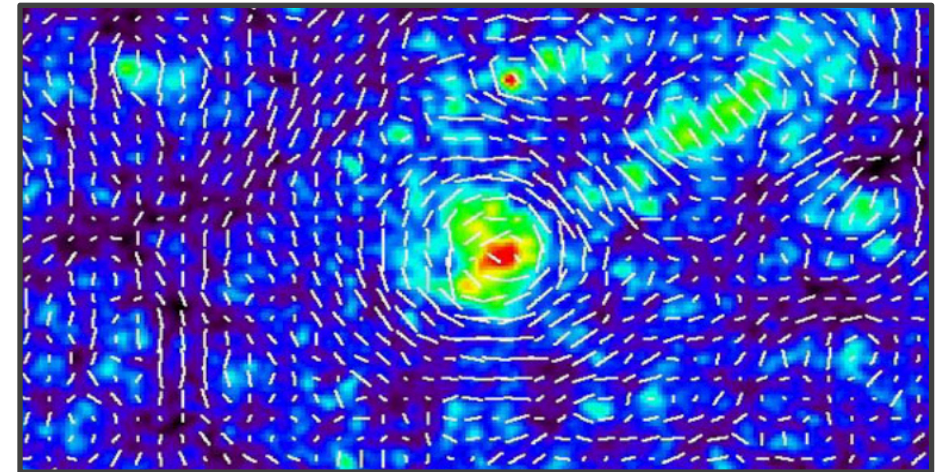
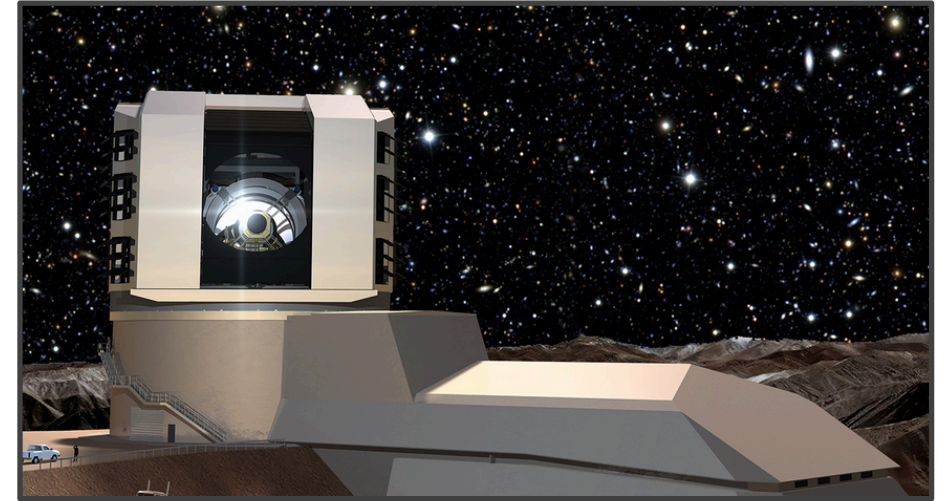
GalDyn – Galaxy Dynamics (UCLan)

- Galactic Dynamics group at UCLan have been simulating orbits of galactic matter:
 - *Focus on particles within galaxies (stars, dark matter, etc.) with thousands of parameters to vary – so highly parallelisable.*
- Use of GridPP infrastructure:
 - *Supported by Lancaster and Liverpool Tier-2s;*
 - *Software and workflows tested with a CernVM;*
 - *Jobs run via GridPP DIRAC on NorthGrid VO.*
- Impact:
 - *Thousands more parameters – new studies;*
 - *‘We can run our code on the Grid and are looking to move into production for our final results. It’s nice to know the facility is there and quick to get setup.’*



The Large Synoptic Survey Telescope

- Manchester LSST researchers have been using the Grid to study **cosmic shear**:
 - *Pilot exercise: use **im3shape** software to analyse $O(100k)$ Dark Energy Survey images.*
- Use of GridPP infrastructure:
 - *Software tested on local (Manc.) Grid cluster;*
 - *Analysis jobs managed with Ganga, locally to begin with before switching to GridPP DIRAC;*
 - *Data managed with the DIRAC File Catalog.*
- Impact:
 - *Results 'significantly faster', no longer relying on highly-contested HPC resources;*
 - *GridPP/LSST collaboration continuing.*



Summary and conclusions

- The Worldwide LHC Computing Grid (WLCG) offers substantial computing resources for HTC using the ‘Grid’ concept:
 - *Computing and data always there, where and how we do not care;*
 - *Crucial to the success of the Large Hadron Collider’s physics programme.*
- GridPP represents UK’s contribution to the WLCG (~11%):
 - *20 institutes, 100+ people, 101k logical CPU cores, 37 PB storage;*
 - *Commitment to make 10% of resources available to non-LHC users.*
- We can offer infrastructure, tools and documentation for new users:
 - *GridPP DIRAC for job/data management, Ganga for UI, CVMFS for software;*
 - *Potential to collaborate to develop solutions for larger communities (later talk);*
 - *Case studies: <https://www.gridpp.ac.uk/users/case-studies/>*

Huge thanks to the GridPP Collaboration, particularly Imperial College London GridPP DIRAC team ([GitHub](#)) and Imperial College London/Birmingham Ganga team ([GitHub](#)).

@GridPP

@twhyntie

Thanks for listening! Any questions?

T. Whyntie*, †

** Queen Mary University of London*

SKA/GridPP F2F, University of Manchester

Wednesday, 2nd November 2016

