# imense

## the GRID as a catalyst for new business

## Dr David Sinclair

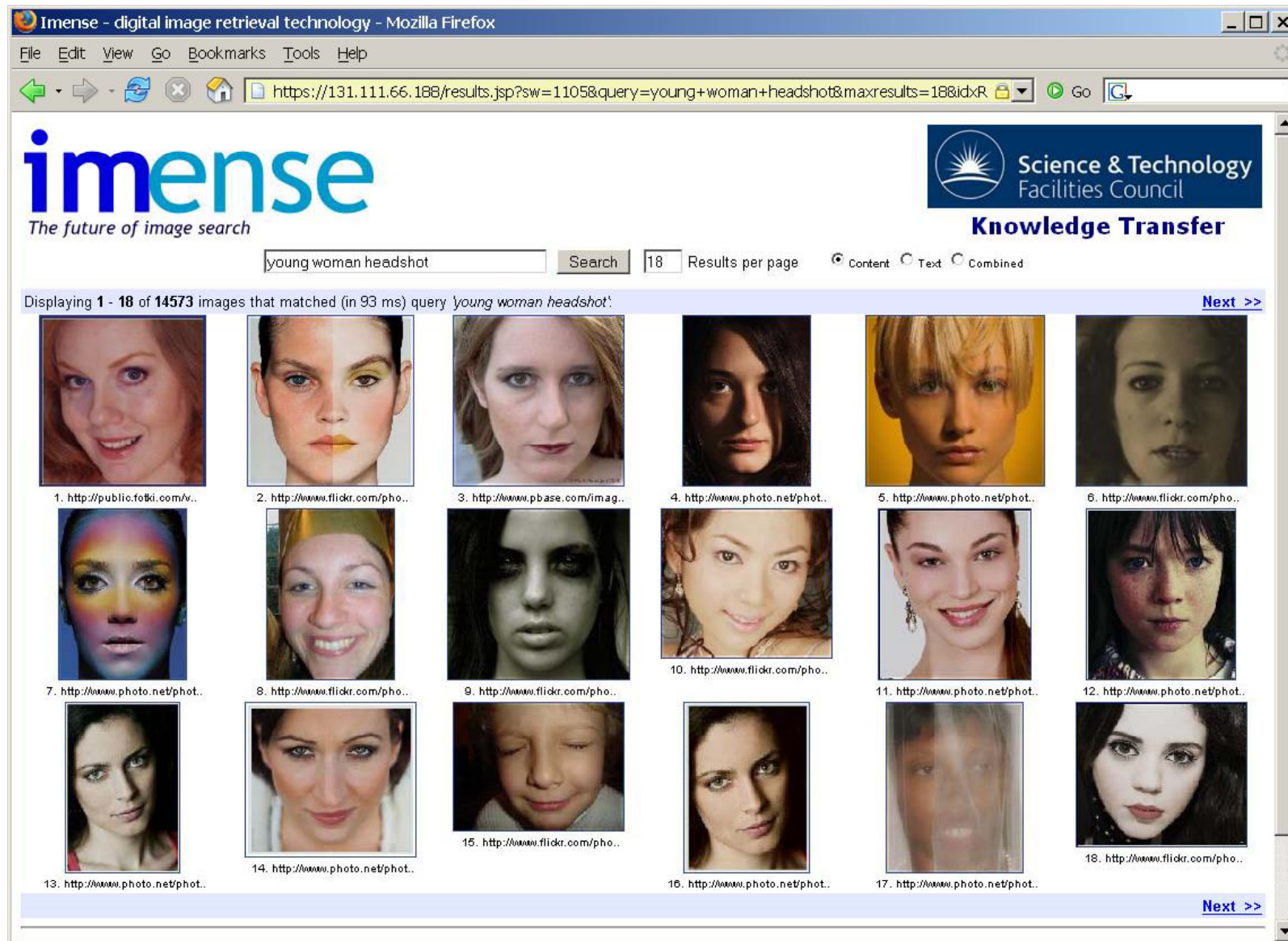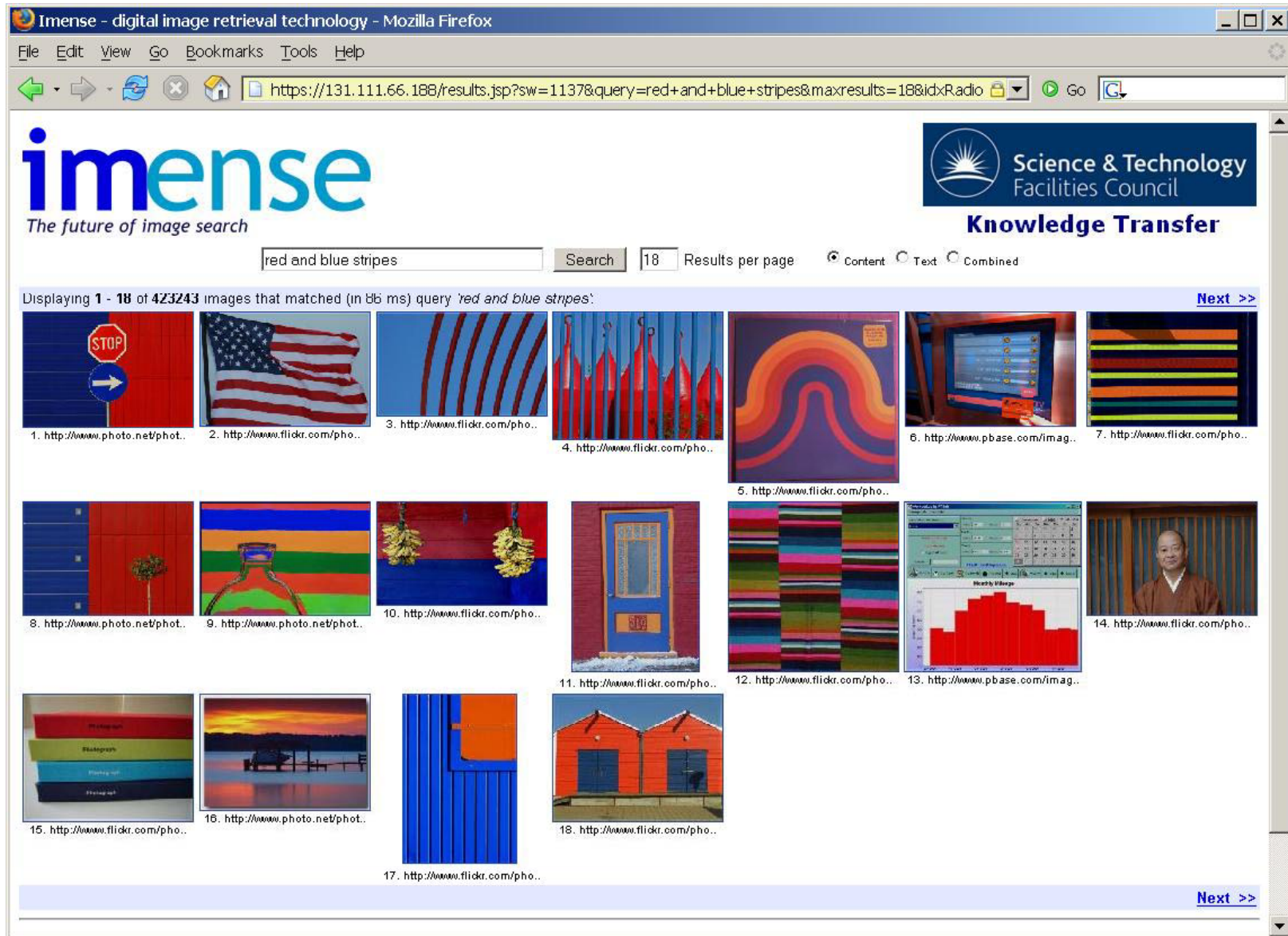david.sinclair@imense.com

http://www.imense.com

# What does imense do?

- Sell image retrieval technology:
- Image content analysis
  - High reliability classifiers for visual content
    - Shapes, scenes, faces, colours etc.
- Indexing and query language provision
  - NLP parser for text query over image content
- Image search portal
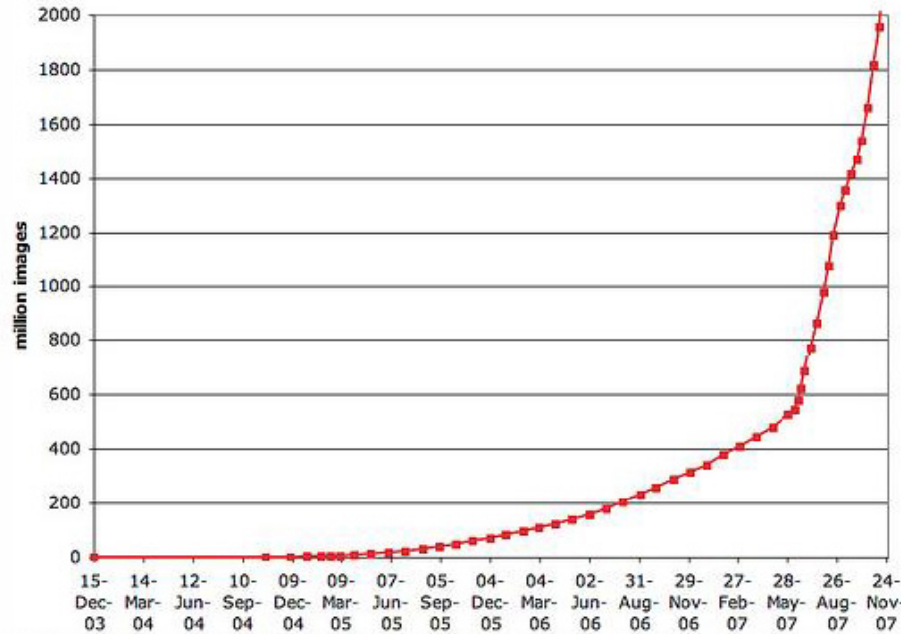  - Search photo-hosted and commercial images

# imense

# How did imense get here?

- Founded in 2004 by: Dr David Sinclair and Dr Chris Town on the back of a DTI grant for R&D.
  - Equipment: 1 laptop, 1 3.2GHz desktop PC, 300Gig of disk.
  - Processed and indexed ~50,000 images.
  - ~$5 \times 10^6$ images uploaded to flickr by early 2005.
    - 2.5 years processing with or available resources…
  - Google had ~200,000 servers in 2005.
- Mini PIPPS grant awarded carrying GRID access 2006.
  - imense indexed $3.5 \times 10^6$ images over the course of the grant
  - Indexing time roughly a month on a small fraction of the GRID.
  - $300 \times 10^6$ images uploaded to flickr by mid 2007.
- Angel investment of 535k generated in July 2007 as a consequence of demonstrating that our indexing technology scaled.
  - imense has now indexed $10 \times 10^6$ images.
  - $2 \times 10^9$ images now uploaded to flickr
  - $15+ \times 10^9$ images on the internet?
  - We still need the GRID for Internet scale image retrieval.
  - Google has 450,000+ servers (estimated 2008).
- Full PIPPS grant to look at improving access to the GRID

# Internet scale image search?



2,000,000,000 images
uploaded to flickr by Jan 2008

- There are more then $10 \times 10^9$ images hosted on the Internet!
  - 4439 cpu-years to index.
    - Or 16 days on a 100,000cpu cluster!
  - 300kBytes/ image $3 \times 10^{12}$ Bytes
    - $0.10/Gig $300,000 just for bandwidth to move the images about.
- We still need grid computing.

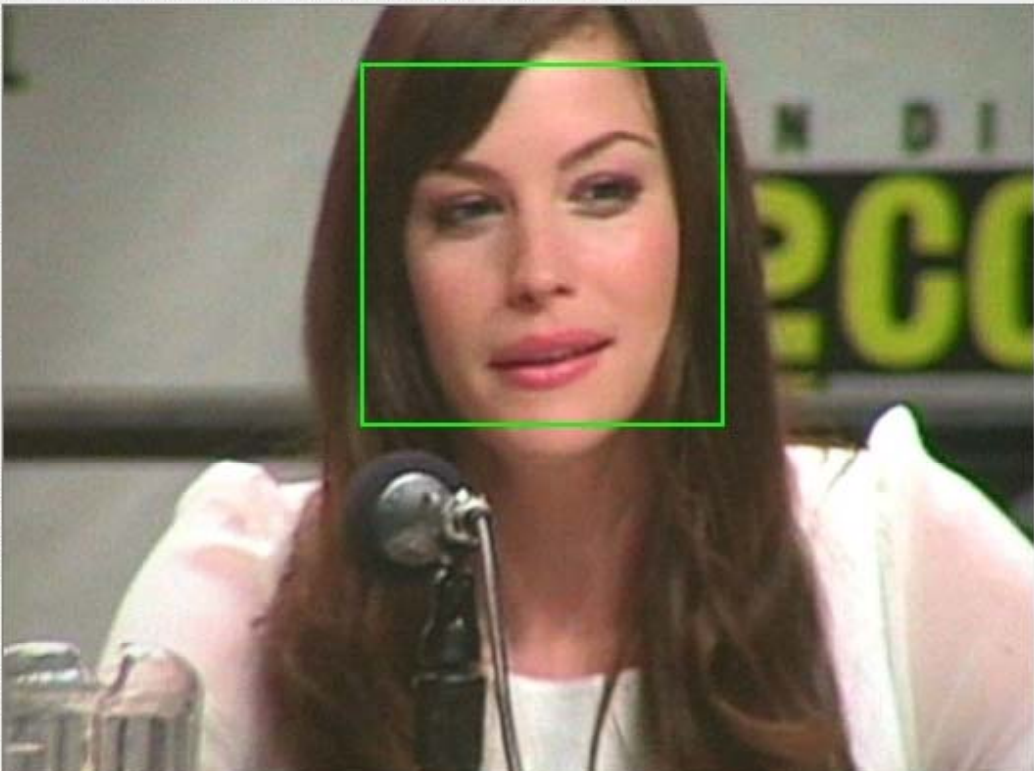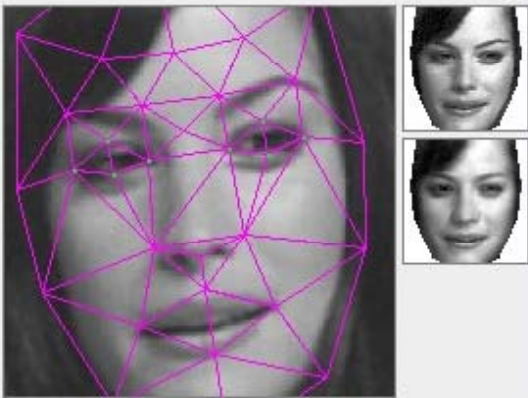Segmentation                    1 second per image.

# imense                                        problem

- Index all of the images on PhotoHosting websites (not even the whole web)
  - Rough estimate 10,000,000,000 images: (flickr.com, photobucket.com, webshots.com etc.)
  - 14sec per image…..4439 cpu_years.
    - 16days on a 100,000cpu cluster!
  - 300kBytes/ image 3x10^12 Bytes
    - $0.05/Gig $150,000 just for bandwidth to move the images about.

# Interacting with GridPP

- I have unashamedly stolen the next few slides for Karl Harrison who has contributed greatly to our work with GridPP.

# Turning a good idea into a working system

Four basic steps to enabling searches based on image content

**Image location**
- Images may be in an archive stored on disk, or may be distributed between web sites

**Image retrieval**
- Retrieve images from storage location to processing node

**Image analysis**
- Perform feature extraction

**Indexing**
- Collate and store analysis results

- Bulk of processing requirement is in analysis step: typically a few seconds per image

- Proof of principle based on several thousands of events is straightforward using minimal resources

- Building up index for many millions of images is more challenging

- Images are analysed independently of one another, so massive parallelisation is possible

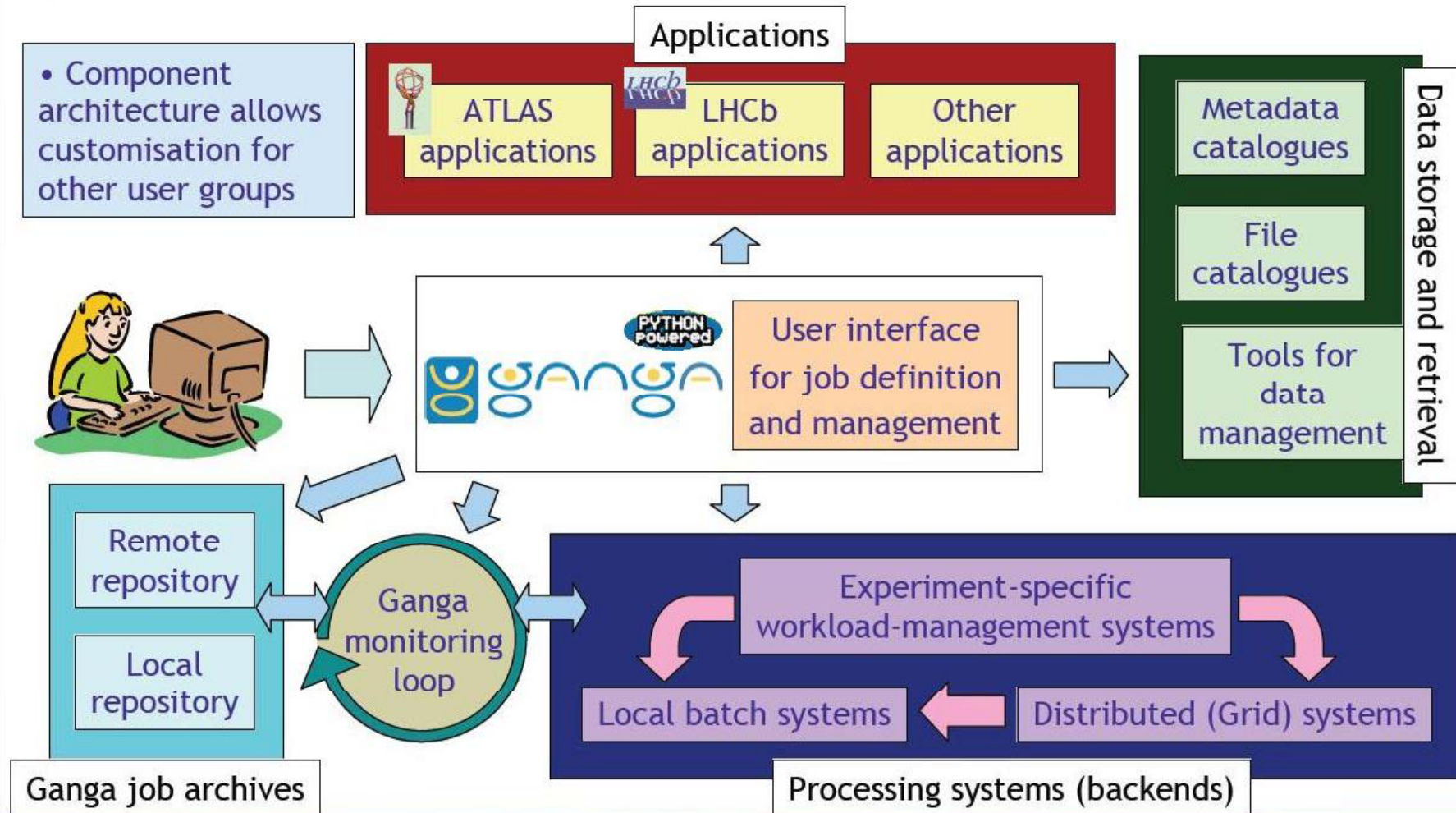  ▸ This is the type of problem where Grid solutions works well

EGEE 08

# Getting image processing onto the particle-physics Grid

- STFC knowledge-transfer projects set up to investigate Grid solutions for large-scale image processing
  - ▸ November 2006 - June 2007: mini-PIPSS award ⇒ feasibility study
  - ▸ October 2007 - April 2009: PIPSS award ⇒ optimised system
    - ▸ Collaboration between Imense Ltd, University of Cambridge High-Energy Physics Group and Cambridge e-Science Centre
    - ▸ Continued involvement from former Cambridge researchers now based at Birmingham
- New Virtual Organisation (camont) set up, and enabled at seven GridPP sites
  - ▸ Access to more sites possible if needed
  - ▸ Help with teething problems from GridPP experts and site managers
- Grid effectively providing computing on demand
  - ▸ Highest number of parallel jobs so far is about 150
  - ▸ Often useful at present to be able to run a few tens of parallel jobs
  - ▸ Aim to ramp up to larger samples later in the year

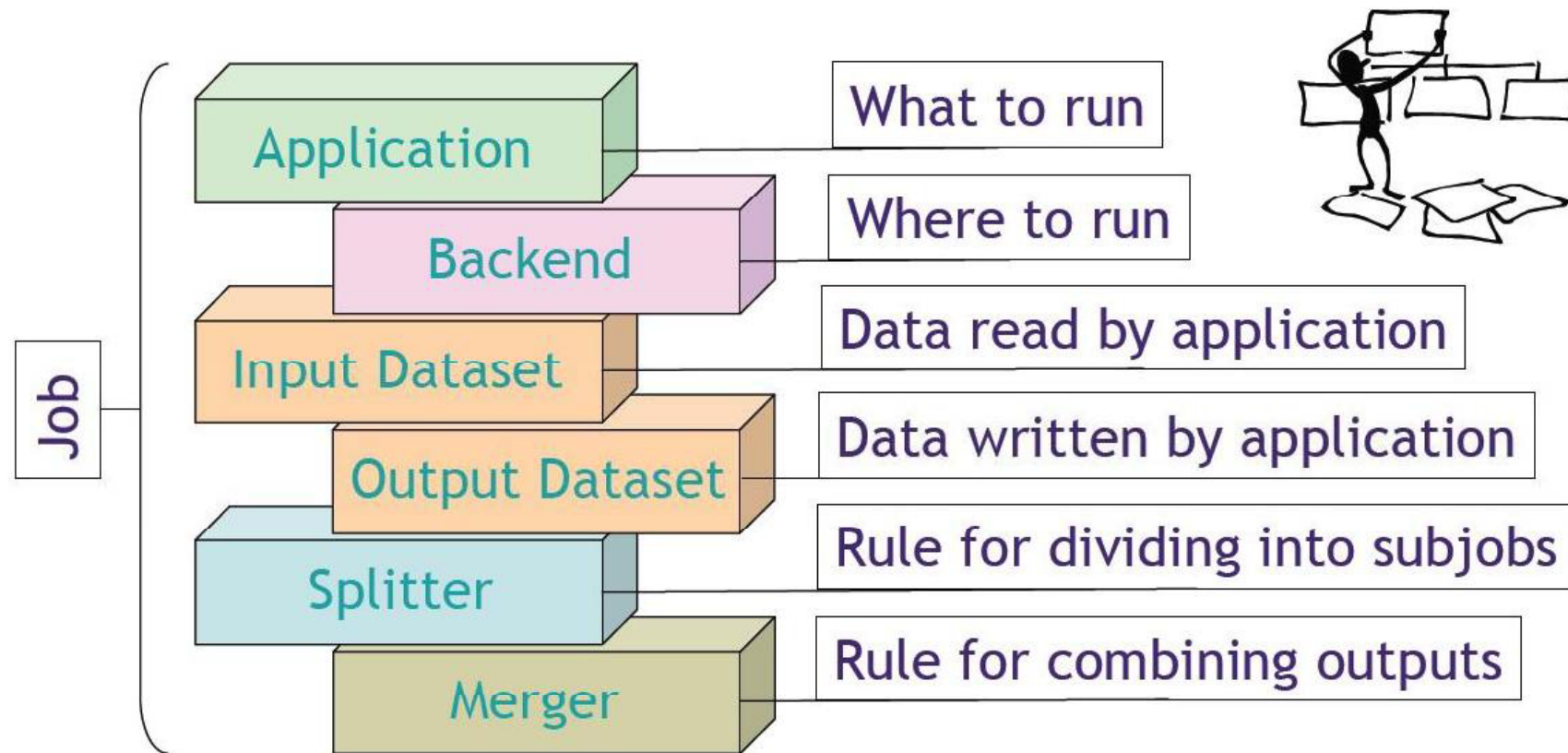EGEE 08

# Job-management sytem and Grid user interface

- Use Ganga system, developed to support particle-physics experiments (ATLAS and LHCb)

- Component architecture allows customisation for other user groups

**Applications**

ATLAS applications

LHCb applications

Other applications

Metadata catalogues

File catalogues

Tools for data management

Data storage and retrieval

**PYTHON powered** ganga User interface for job definition and management

Remote repository

Local repository

Ganga job archives

Ganga monitoring loop

Experiment-specific workload-management systems

Local batch systems

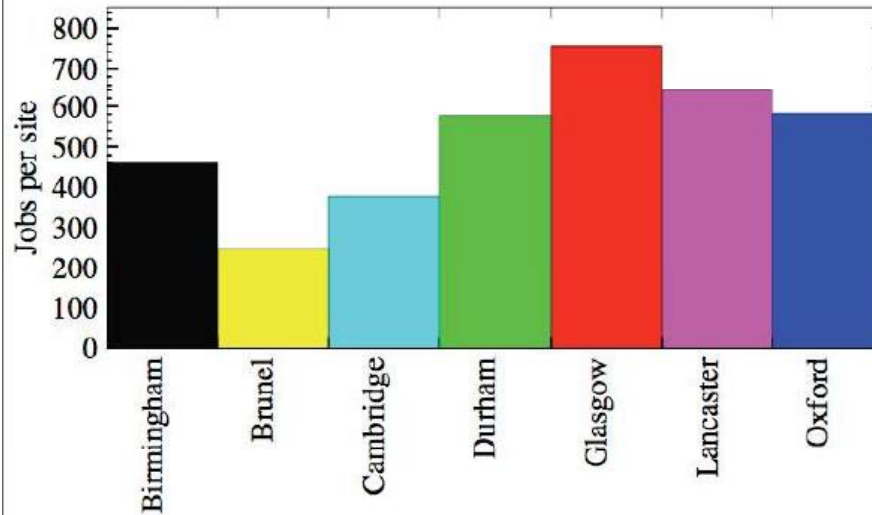Distributed (Grid) systems

Processing systems (backends)

EGEE 08

# Ganga job abstraction

A job in Ganga is constructed from a set of building blocks, not all needed for every job



| Block | Description |
|-------|-------------|
| Application | What to run |
| Backend | Where to run |
| Input Dataset | Data read by application |
| Output Dataset | Data written by application |
| Splitter | Rule for dividing into subjobs |
| Merger | Rule for combining outputs |

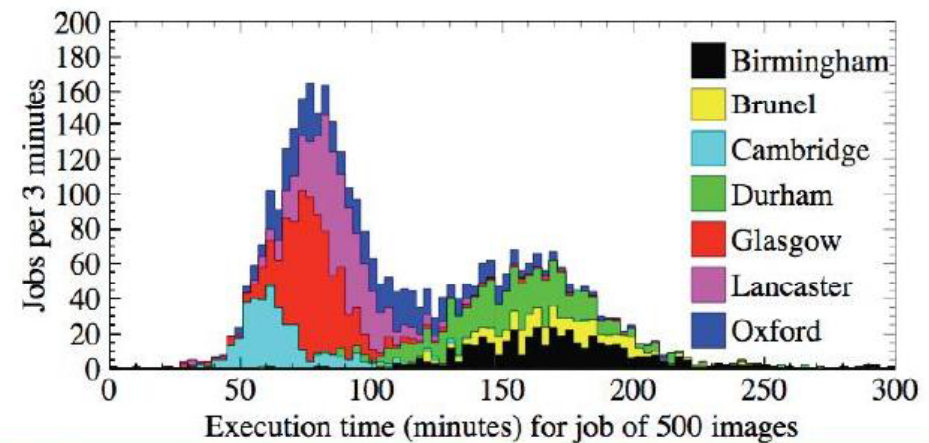UNIVERSITY OF BIRMINGHAM

EGEE 08

# Job destinations and execution times

Results for 3638 jobs submitted over four-week period, July-August 2008



Destination chosen by Resource Broker of Workload Management System, based on minimum estimated waiting time

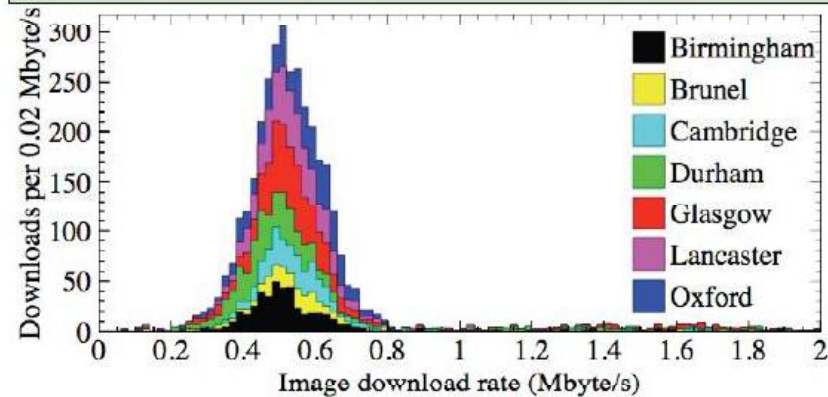Significant differences in execution times reflect inhomogeneity of site resources
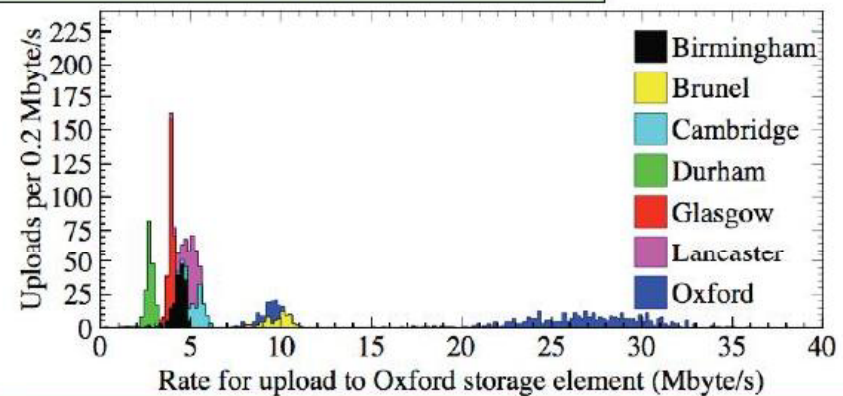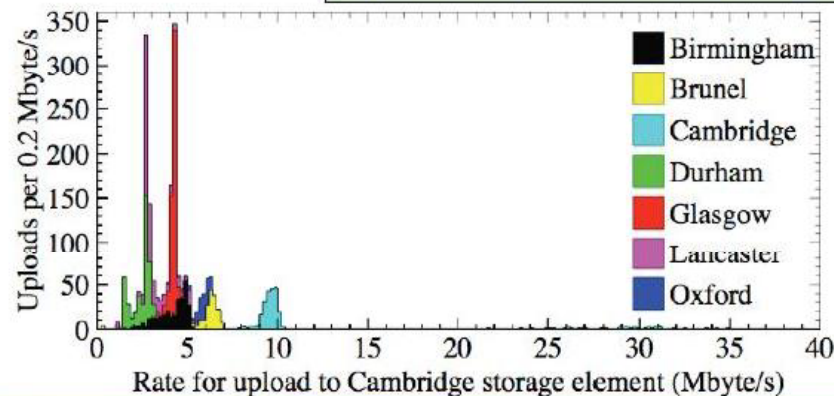
# Data-transfer rates

Image downloads from hosting site, using wget



Upload of results to Grid storage elements, using globus-url-copy

# Grid overheads

-Useful time is when job is downloading and processing images
- Grid overheads come from: startup time, system time for logging job completion, result upload and retrieval
- For jobs of 500 images, average start-to-finish time is 164 minutes, with 39 minutes spent on Grid overheads: 73% useful time
$\Rightarrow$ Timing distributions non-Gaussian, with long tails
$\Rightarrow$ Need to increase processing load to increase fraction of useful time

EGEE 08

# Experience of image processing on the Grid

- Grid has been successfully used to process several million images
  - ▸ Have processed both images from a disk archive and images retrieved directly to Grid nodes from image-hosting web sites
  - ▸ This has contributed to launch of beta version of new image search engine: http://imense.com/
- Have automated system, based on Ganga, for job submission and output retrieval
  - ▸ Makes keeping track of thousands of jobs and millions of images completely painless
- Job failure rates have been at 2% level, with two main causes
  - ▸ Proxy credential of submitting user expires before job starts
  - ▸ Network failures, preventing upload of results to storage element
- Positive experience with using the Grid for image retrieval and processing has prompted interest in using the Grid also for image location
  - ▸ Grid-enabled web crawler now at testing phase

EGEE 08

# imense                                              thanks

- STFC for providing
  - Technical vote of confidence in our project
  - Facilities to demonstrate our technology
  - Expertise in GRID processing
  - money
- Personal thanks for invaluable support with our project to:
  - Prof Michael, A. Parker of Cambridge University, High Energy Physics Group.
  - Dr Alexander Efimov of Qi3
  - Dr Karl Harrison regrettably now of Birmingham University.