

vl-e



virtual laboratory for e-science

e-Science and e-Bioscience

The VL-e approach

L.O. (Bob) Hertzberger

Director of Virtual Laboratory for e-Science (VL-e)

Adjunct director Netherlands Bio Informatics Centre (NBIC)

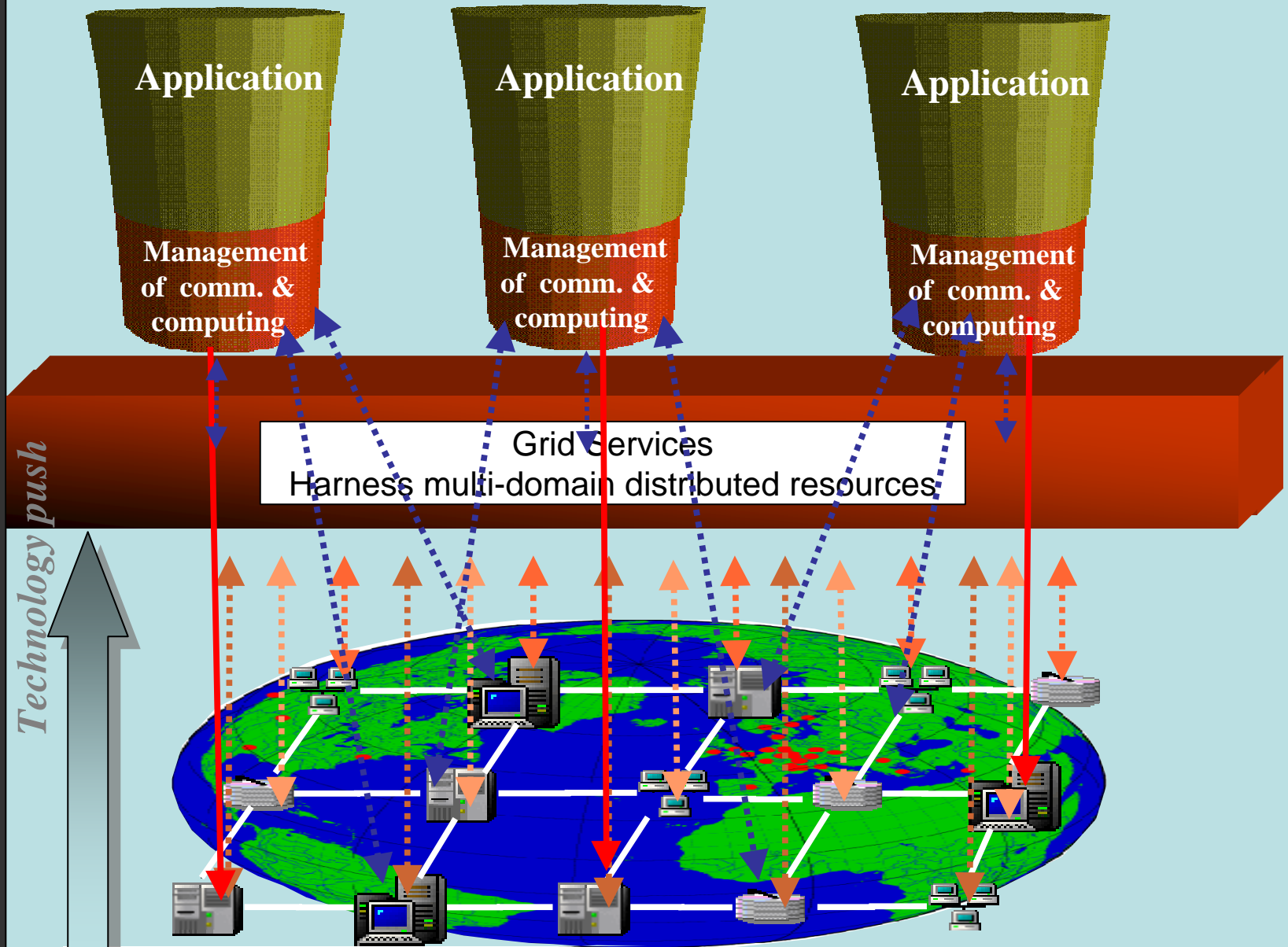
One of the directors of BIG GRID

bob@science.uva.nl

Content

- Some issues on Grid important for e-Science
- Developments towards e-Science
- The VL-e approach
- Some examples
- Conclusions

What is Grid



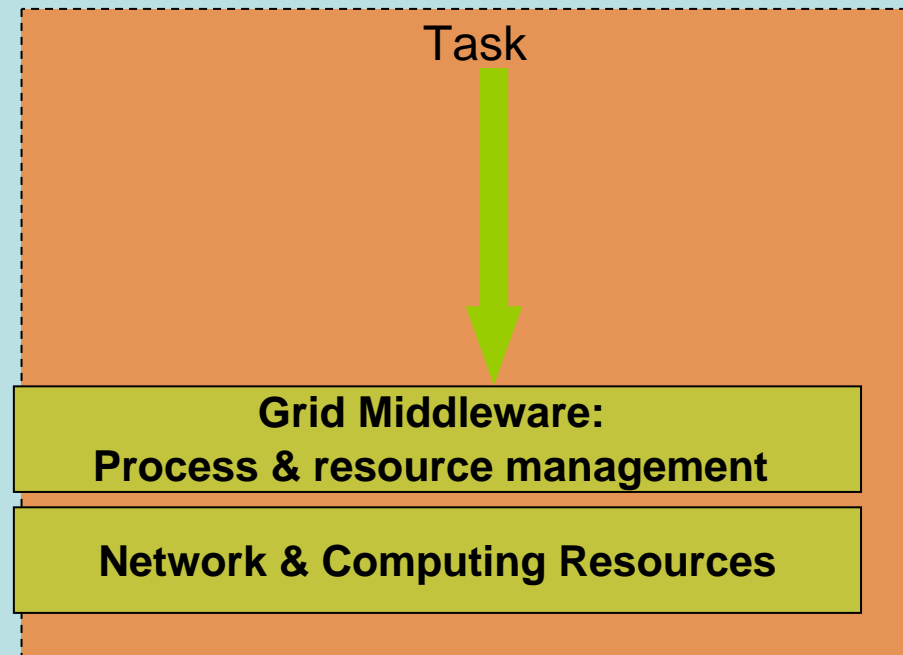
Grid before WSRF/OGSA

The word '*grid*' has been used in many ways

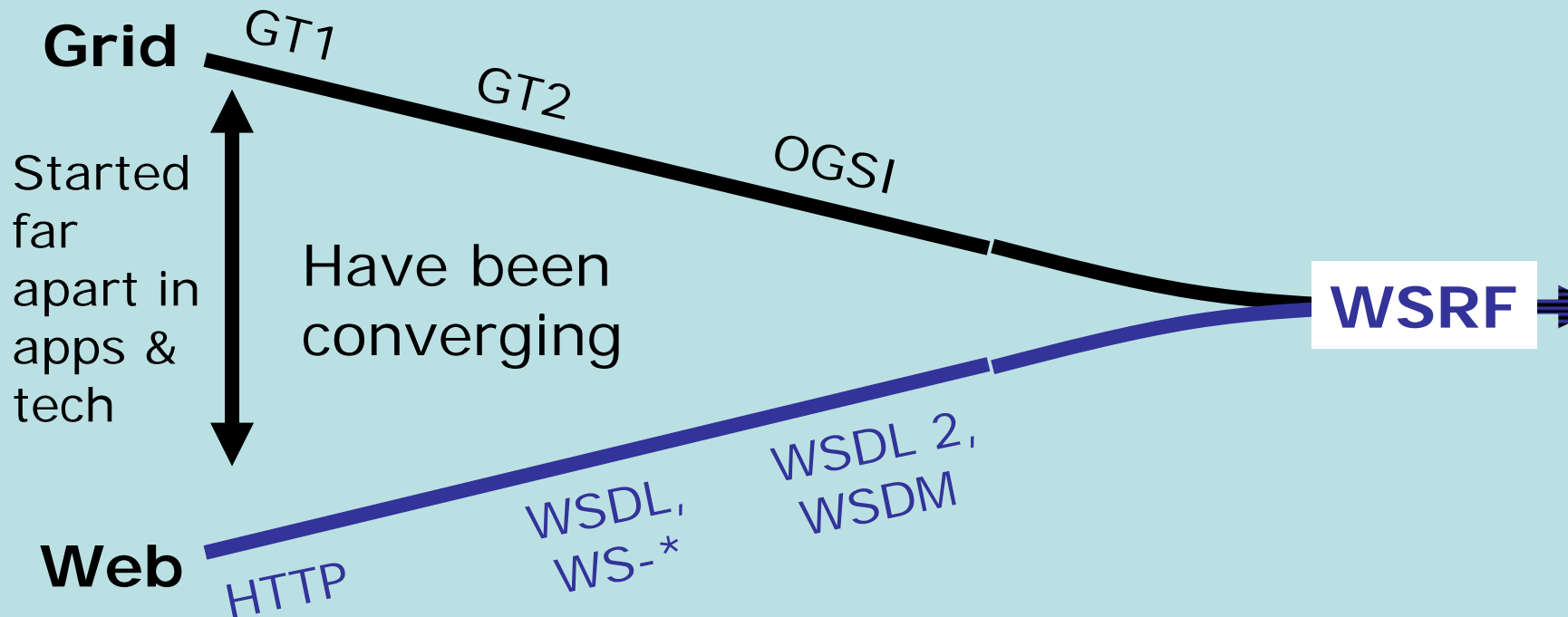
- ✓ cluster computing
- ✓ cycle scavenging
- ✓ **cross-domain resource, data and information sharing**

A definition for what we mean with grid

- Coordinates resources not subject to centralised control
- Using standard, open and generic protocols & interfaces
- Provides non-trivial qualities of collective service
- Virtualization of resources via among others Virtual Organizations



Grid and Web Services Convergence



Definition of **Web Service Resource Framework(WSRF)** makes explicit distinction between **"service"** and stateful entities acting upon service i.e. the **resources**

Means that Grid and Web communities can move forward on a common base!!!

Grid after WSRF & OGSA

- Important aspects:

- ✓ Uniform syntax & use of WSDL or other

- ✓ **Problems**

- ✓ Level of abstraction of service

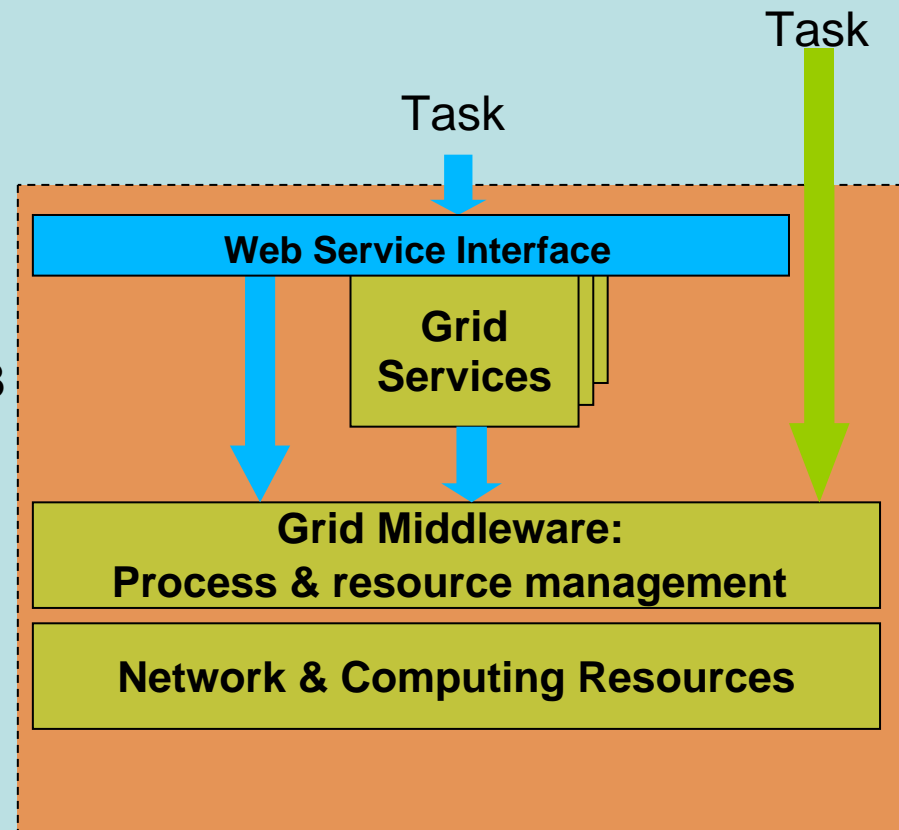
- ✓ How does service react

- ✓ Quick standardization

- First implementation in 2008 might be late

- **Advantage:**

- ✓ OGSA offers a coherent set of services



What is next ???

Knowledge and the knowledge producing & consuming protocols & patterns are already in Grid Middleware and Grid Applications

Embedded in middleware code, in schemas, in catalogues, in applications and in practice
They are often called metadata.

Carol Goble.

Issue is to make this knowledge explicit

Semantic WEB/GRID

Levels of Grid abstraction

Knowledge (Semantic) Web/Grid

Information Web/Grid

Data Grid

Computational Grid

Background information experimental sciences

- There is a tendency to look ever deeper in:
 - ✓ Matter e.g. Physics
 - ✓ Universe e.g. Astronomy
 - ✓ Life e.g. Life sciences
- Instrumental consequences are increase in detector:
 - ✓ Resolution & sensitivity
 - ✓ Automation & robotization
- Therefore experiments become increasingly more complex

Impact in the life sciences

- Impact of high throughput methods e.g. **Omics** experimentation
 - ✓ genome ==> genomics
- **Instrumentation being used in omics experimentation:**
 - ✓ Transcriptomics via among others; micro-arrays
 - ✓ Proteomics via among others; Mass Spectroscopy (MS)
 - ✓ Metabolomics via among others; MS & Nuclear Magnetic Resonance (NMR)

Background information experimental sciences

- There is a tendency to look ever deeper in:
 - ✓ Matter e.g. Physics
 - ✓ Universe e.g. Astronomy
 - ✓ Life e.g. Life sciences
- Instrumental consequences are increase in detector:
 - ✓ Resolution & sensitivity
 - ✓ Automation & robotization
- Therefore experiments become increasingly more complex
- Results in an increase in **the amount and the complexity of data**

Data explosion

- In 2005 more data is being produced than during existence of human
- In 2005 increase from 3 Mexabytes towards more than 40 Mexabytes

Source: Disaster Recovery Journal
Autum 2004

- Impact on science
- How do we deal with
- How to extract information



Data explosion results in the Life sciences in a Paradigm shift

- Past experiments where **hypothesis driven**
 - ✓ Evaluate hypothesis
 - ✓ Complement existing knowledge
- Present experiments are **data driven**
 - ✓ Discover knowledge from large amounts of data
 - ✓ More **integration** necessary

Background information experimental sciences

- There is a tendency to look ever deeper in:
- Instrumental consequences are increase in detector:
- Therefore experiments become increasingly more complex
- Results in an increase in the amount and complexity of data
- Something has to be done to harness this development
 - ✓ Virtualization of experimental resources enabling sharing & leading to e-Science

The what of e-Science

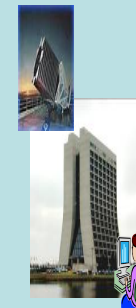
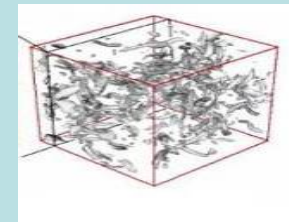
- e-Science is **a new science paradigm** complementing theoretical and experimental science
 - ✓ More than only coping with data explosion
 - ✓ A multi-disciplinary activity combining human expertise & knowledge

e-Science: A New Science Paradigm

- Thousand years ago:
Experimental Science
 - description of natural phenomena
- Last few hundred years:
Theoretical Science
 - Newton's Laws, Maxwell's Equations ...
- Last few decades:
Computational Science
 - simulation of complex phenomena
- Today:
e-Science or Data-centric Science
 - unify theory, experiment, and simulation
 - using data exploration and data mining
 - Data captured by instruments
 - Data generated by simulations
 - Data generated by sensor networks
 - ✓ Scientist analyzes databases/files



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$



The what of e-Science

- e-Science is **a new science paradigm** complementing theoretical and experimental science
 - ✓ More than only coping with data explosion
 - ✓ A multi-disciplinary activity combining human expertise & knowledge
- e-Science should
 - ✓ Apply and integrate Web/Grid methods where and whenever possible
 - ✓ Will also be a driver for new Web/Grid methods (semantics)
 - ✓ Apply an optimal ICT infrastructure
- e-Science demands a different approach to science
 - ✓ At the end-user level: application or domain **system level science**
 - ✓ At the e-Science technology level: **integrative science**

e-Science Objectives

- It should enhance the scientific process by:
- Stimulating collaboration by sharing data & information
 - ✓ Result is re-use of data & information

The data sharing potential for Cognition

- **Collaborative** scientific research
 - ✓ Information sharing
 - ✓ Metadata modeling
- Allows for **experiment validation**
 - ✓ Independent confirmation of results
- Statistical methodologies
 - ✓ Access to large collections of data and metadata
- **Training**
 - ✓ Train the next generation using peer reviewed publications *and* the associated data

Sharing neuroimaging studies of human cognition

John Darrell Van Horn, Scott T Grafton, Daniel Rockmore & Michael S Gazzaniga

After more than a decade of collecting large neuroimaging datasets, neuroscientists are now working to archive these studies in publicly accessible databases. In particular, the fMRI Data Center (fMRIDC), a high-performance computing center managed by computer and brain scientists, seeks to catalogue and openly disseminate the data from published fMRI studies to the community. This repository enables experimental validation and allows researchers to combine and examine patterns of brain activity beyond that of any single study. As with some biological databases, early scientific, technical and sociological concerns hindered initial acceptance of the fMRIDC. However, with the continued growth of this and other neuroscience archives, researchers are recognizing the potential of such resources for identifying new knowledge about cognitive and neural activity. Thus, the field of neuroimaging is following the lead of biology and chemistry, mining its accumulating body of knowledge and moving toward a 'discovery science' of brain function.

The observation that changes in regional cerebral blood flow accompany neural activity during cognition¹⁻³ has been a boon to the cognitive and brain sciences, most notably through the use of brain mapping technologies such as functional magnetic resonance imaging (fMRI). Current research efforts for imaging the brain 'in action' are underway to rigorously explore the structure and function of cognitive brain processes, thereby characterizing the mental properties that make us uniquely human⁴. The fMRI studies range from the examination of familiar cognitive processes such as human memory and language processing to novel studies of racial threat⁵ and the neurofunctional components of humor⁶.

This increasing dependence on brain mapping for exploring cognition has led to an unprecedented data explosion that is pressing neuroscientists to manage and analyze data on scales never before imagined. Complete fMRI study data sets now routinely reach several gigabytes in size, with the amount of brain image data collected in some articles^{7,8} beginning to rival the current size of many biological and physical science databases^{9,10}. What is more, the size of fMRI studies has grown over time, and what is now considered a large

fMRI study will seem relatively small within only a few years, as new technological developments occur in scanner physics, engineering and protocol design.

Unfortunately, despite this progress, much of these fMRI data are not readily available to anyone beyond the original research team that collected them. There are several reasons behind the fact that other investigators do not typically get to work with the actual data that went into the heavily processed images appearing in a published article: (i) limitations of publication space on the complete representation of fMRI methods and findings, (ii) the proprietary feelings of investigators against letting others view their data, (iii) the immensity of data set size and (iv) the convention of only reporting tabular representations of activity in individual image voxels. However, given recent success stories from genomics¹¹ and proteomics¹² for organizing, archiving and mining large amounts of data from their communities, it may come as no surprise that cognitive neuroscientists are now looking to unfettered data sharing and study archiving to better understand these rich collections of dynamic brain data.

Data sharing sociology in neuroimaging

In 2000, with precisely such a goal, we founded the fMRIDC (www.fmriddc.org) at Dartmouth College. We sought to facilitate progress in understanding cognitive processes through the collection, archiving and open distribution of neuroimaging data sets in the peer-reviewed literature¹³. We reasoned that there could be several positive outcomes to making the complete study data sets available to others. First, the study findings could be independently confirmed, helping to strengthen the findings drawn by the original authors. Second, new statistical methodologies could be applied to the data, providing novel insights into cognitive processes. Different studies could be compared, possibly identifying unanticipated functional homologies between seemingly different cognitive tasks. Moreover, these studies could be used to train the next generation of neuroscientists by using fMRI data that had already undergone interpretation by those who collected it and had published it in leading journals. We decided to focus on fMRI data from published articles and not to be concerned with unpublished data. This allowed us to focus the enormous chore of collecting and managing the data, as well as to construct an archive that was representative of the field's body of work.

We approached the editors of several leading journals and were pleased by their initial support. To form the first corpus of data sets and accompanying study material, a special issue of the *Journal of Cognitive Neuroscience* (JOCN) was published containing a collection of articles from leading laboratories (Vol. 12, Suppl. 2, 2000). The authors of these articles generously provided the raw, processed and

structural images and study meta-data

All authors are at Dartmouth College, Hanover, New Hampshire 03755, USA. John Darrell Van Horn and Scott T Grafton are at the Center for Cognitive Neuroscience, the Dartmouth Brain Imaging Center and the fMRI Data Center, Daniel Rockmore is in the Department of Mathematics and the fMRI Data Center, and Michael S. Gazzaniga is at the Center for Cognitive Neuroscience and the fMRI Data Center.
e-mail: John.D.Van.Horn@dartmouth.edu



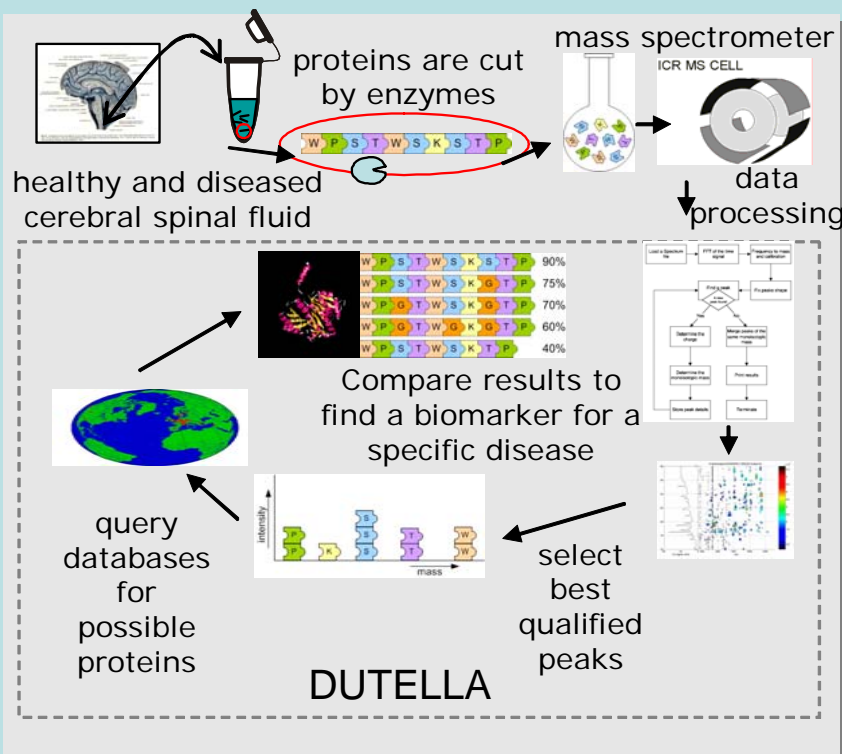
e-Science Objectives

- It should enhance the scientific process by:
- Stimulating collaboration by sharing data & information
 - ✓ Improve re-use of data & information
- **Combing data and information from different modalities**
 - ✓ **Sensor data & information fusion**

In Biomarker Research Multiple Resources have to be integrated



- Samples
 - ✓ Patients
 - ✓ Hospitals
 - ✓ Diseases
- Analytical instrumentation
- Researchers and expertise
- Processing tools
- Information
- Data storage facilities
- Computational resources



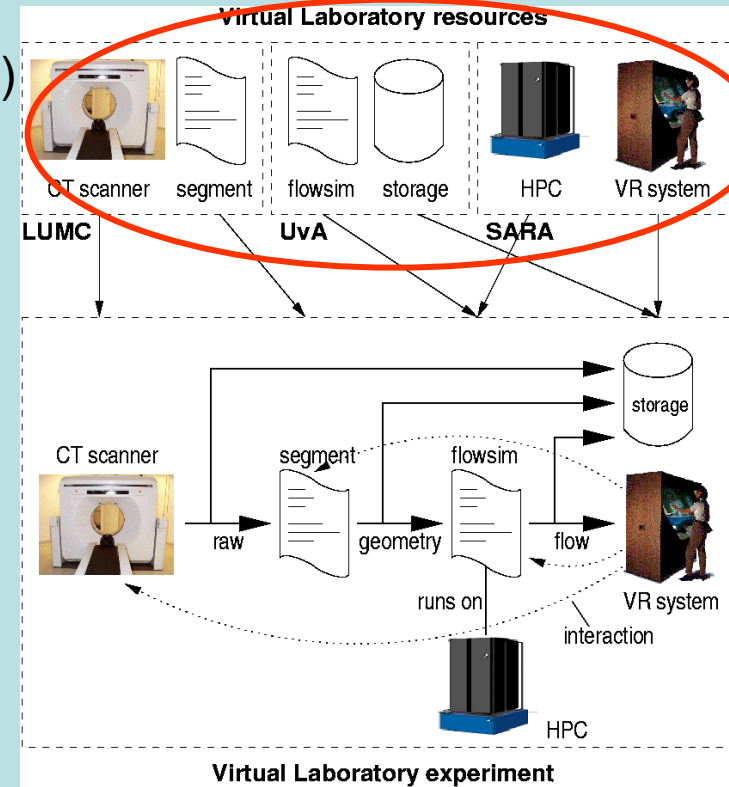
e-Science Objectives

- It should enhance the scientific process by:
- Stimulating collaboration by sharing data & information
 - ✓ Improve re-use of data & information
- Combing data and information from different modalities
 - ✓ Sensor data & information fusion
- Realize the combination of real life & (model based) simulation experiments

Simulated Vascular Reconstruction in a Virtual Operating Theatre

- *An example of new models and interaction*
 - patient specific vascular geometry (from CTA)
 - segmentation
 - blood flow simulation (**Lattice Boltzmann**)
- Pre-operative planning (**interaction**)
- Suitable for parallelization through functional decomposition

Grid resources



Patient's vascular geometry (CTA)

Simulated "Fem-Fem" bypass

e-Science Objectives

- It should enhance the scientific process by:
- Stimulating collaboration by sharing data & information
 - ✓ Improve re-use of data & information
 - ✓ Combing data and information from different modalities
 - Sensor data & information fusion
- Realize the combination of real life & (model based) simulation experiments
- Allow for modeling of dynamic systems

Bird behaviour in relation to weather and landscape

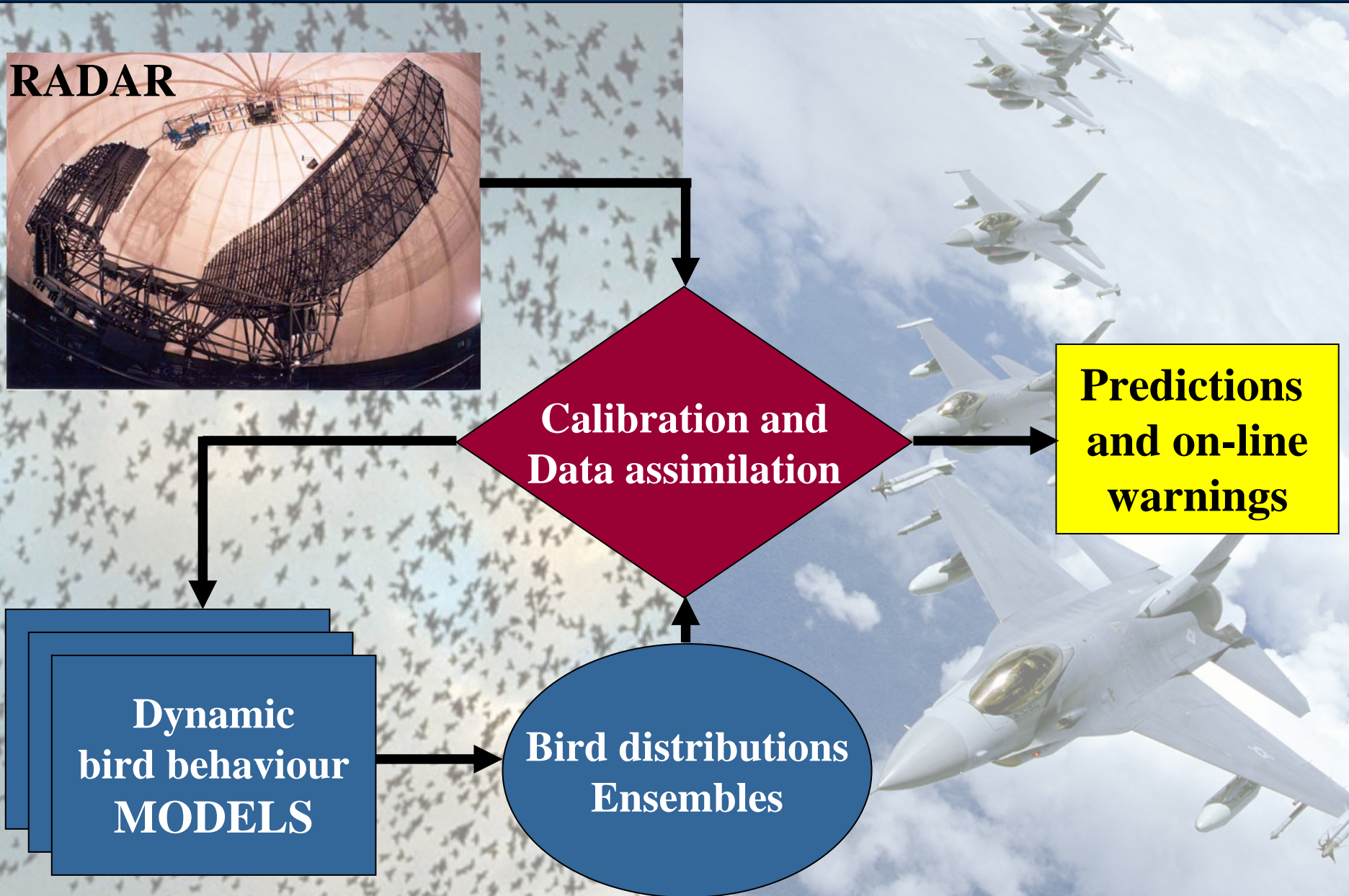


**Calibration and
Data assimilation**

**Predictions
and on-line
warnings**

**Dynamic
bird behaviour
MODELS**

**Bird distributions
Ensembles**

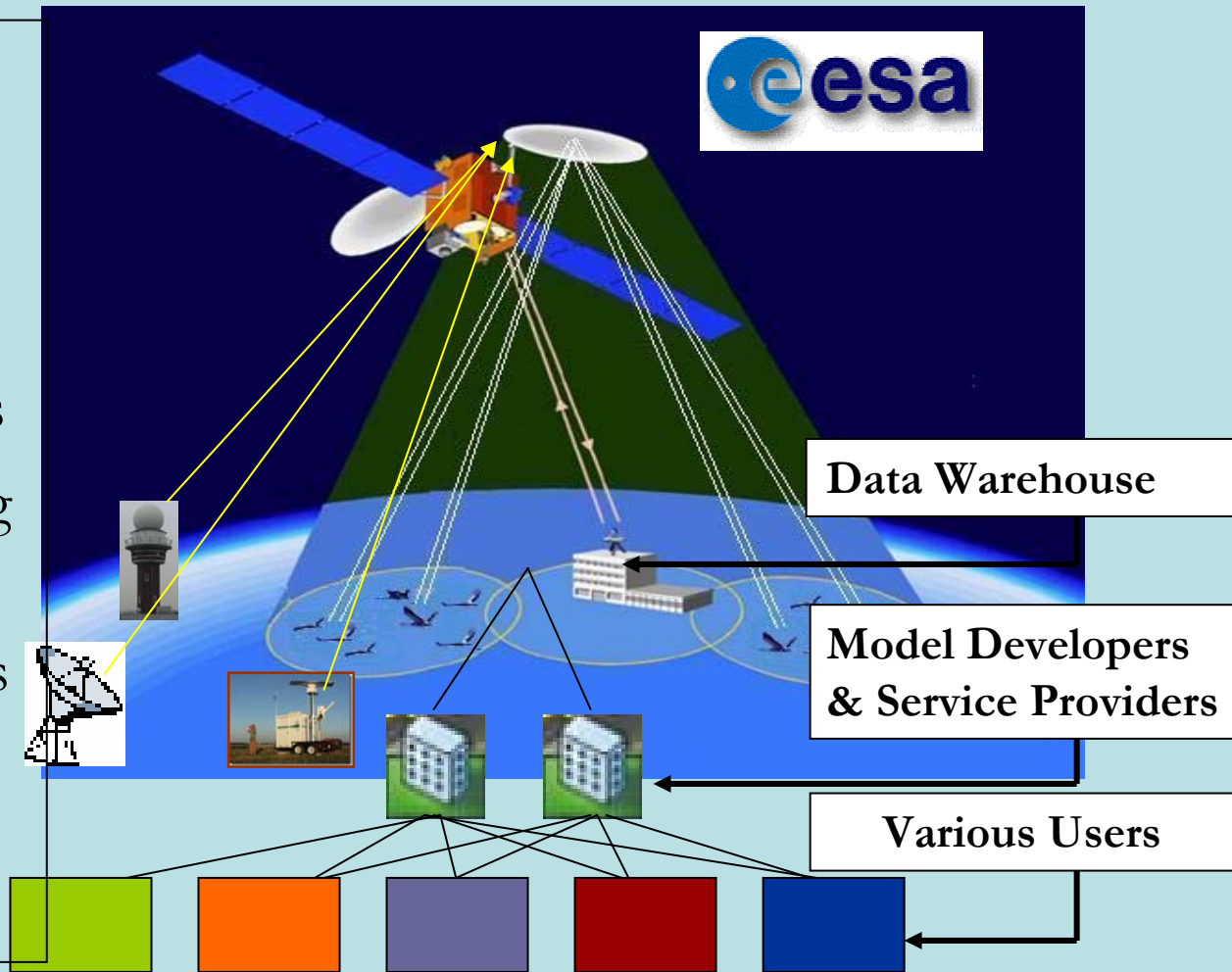


System of Systems: Avian Allert (2006-2013)

Monitoring and Modelling Bird Behaviour and Migration

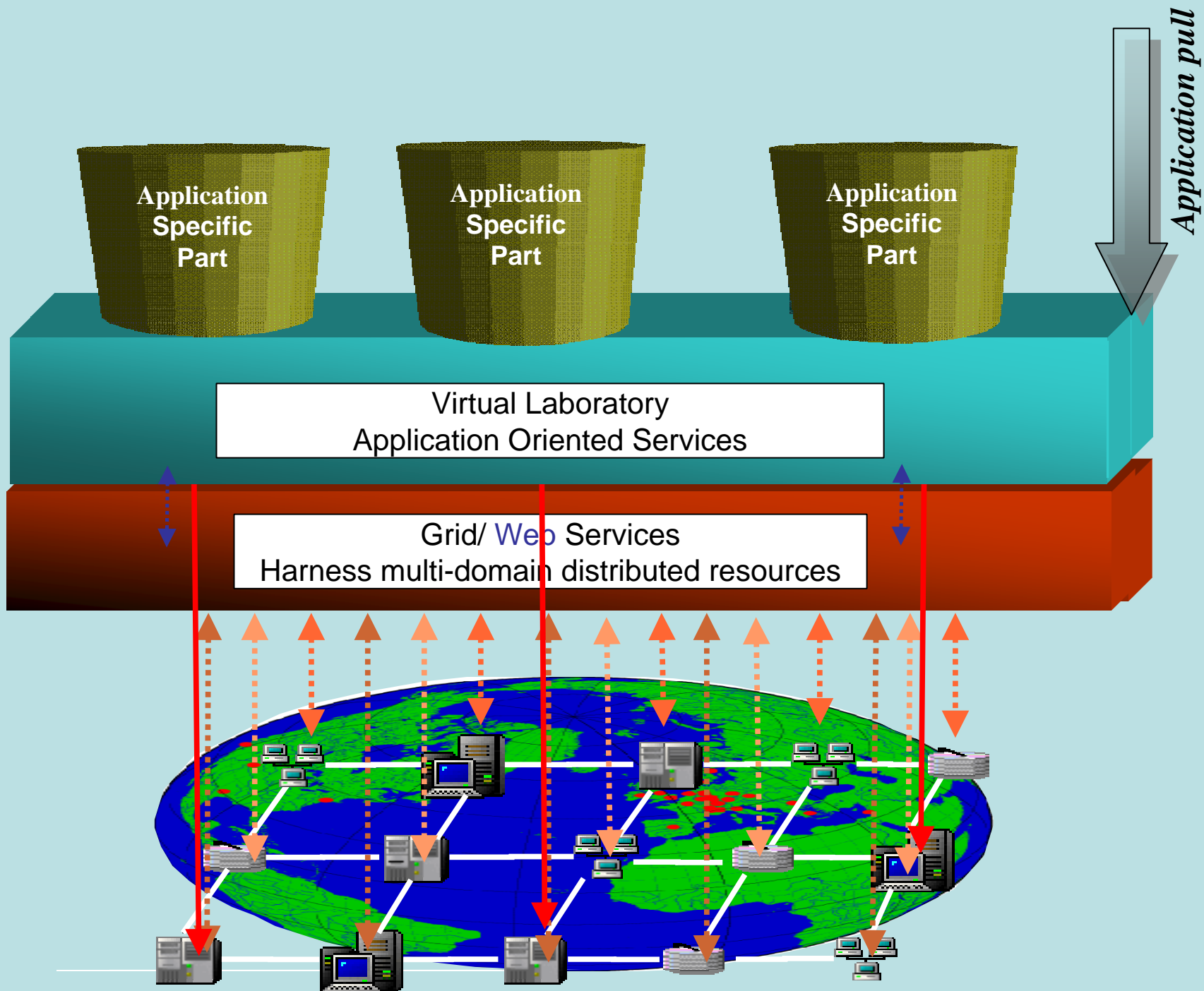
Systems

1. Military Surveillance Radars
2. Weather Radars
3. Remote Sensing
4. GPS on Individual Birds
5. Virtual Lab. for Model Development



Virtual Lab for e-Science research Philosophy

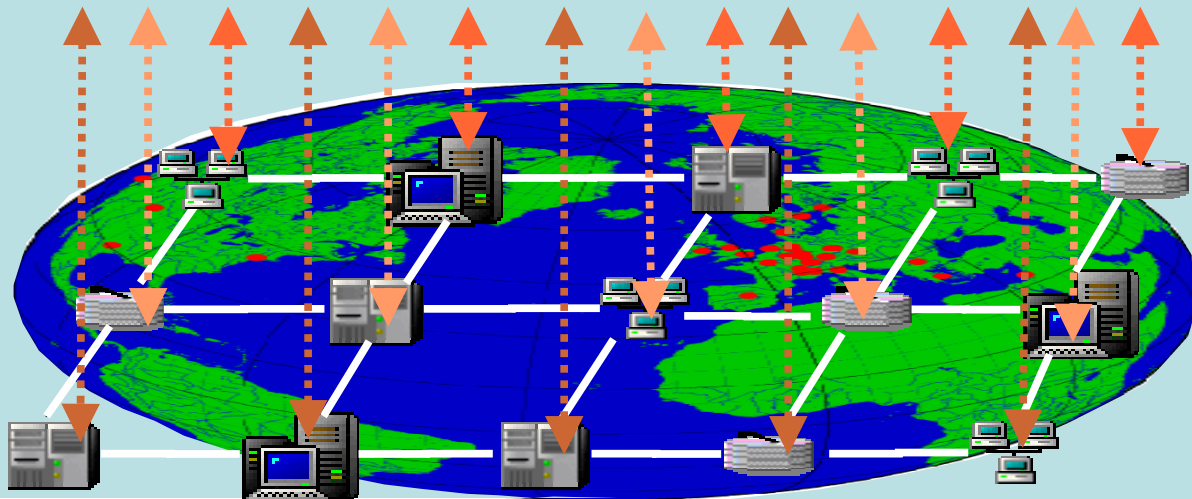
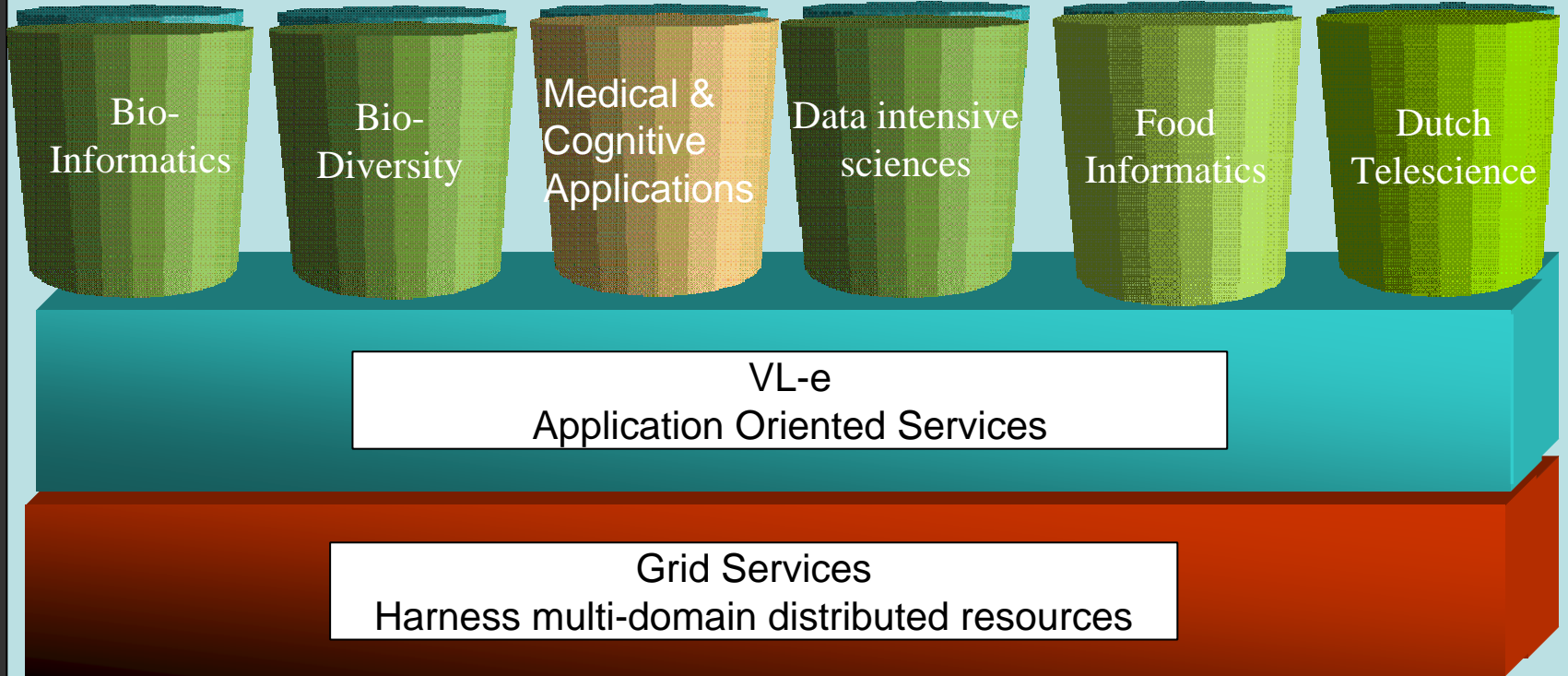
- Generic application support
 - ✓ Application cases are drivers for computer & computational science and engineering research
 - ✓ Problem solving partly generic and partly specific
 - ✓ Re-use of components via generic solutions whenever possible



Potential for generic e- Science services

- Virtual Reality Visualization & user interfaces
- Modeling & Simulation
 - ✓ Interactive Problem Solving
- Data & information management
 - ✓ Data modeling
 - ✓ dynamic work flow management
- Content (knowledge) management
 - ✓ Semantic aspects
 - ✓ Meta data modeling
 - Ontologies
- Wrapper technology
- Design for Experimentation

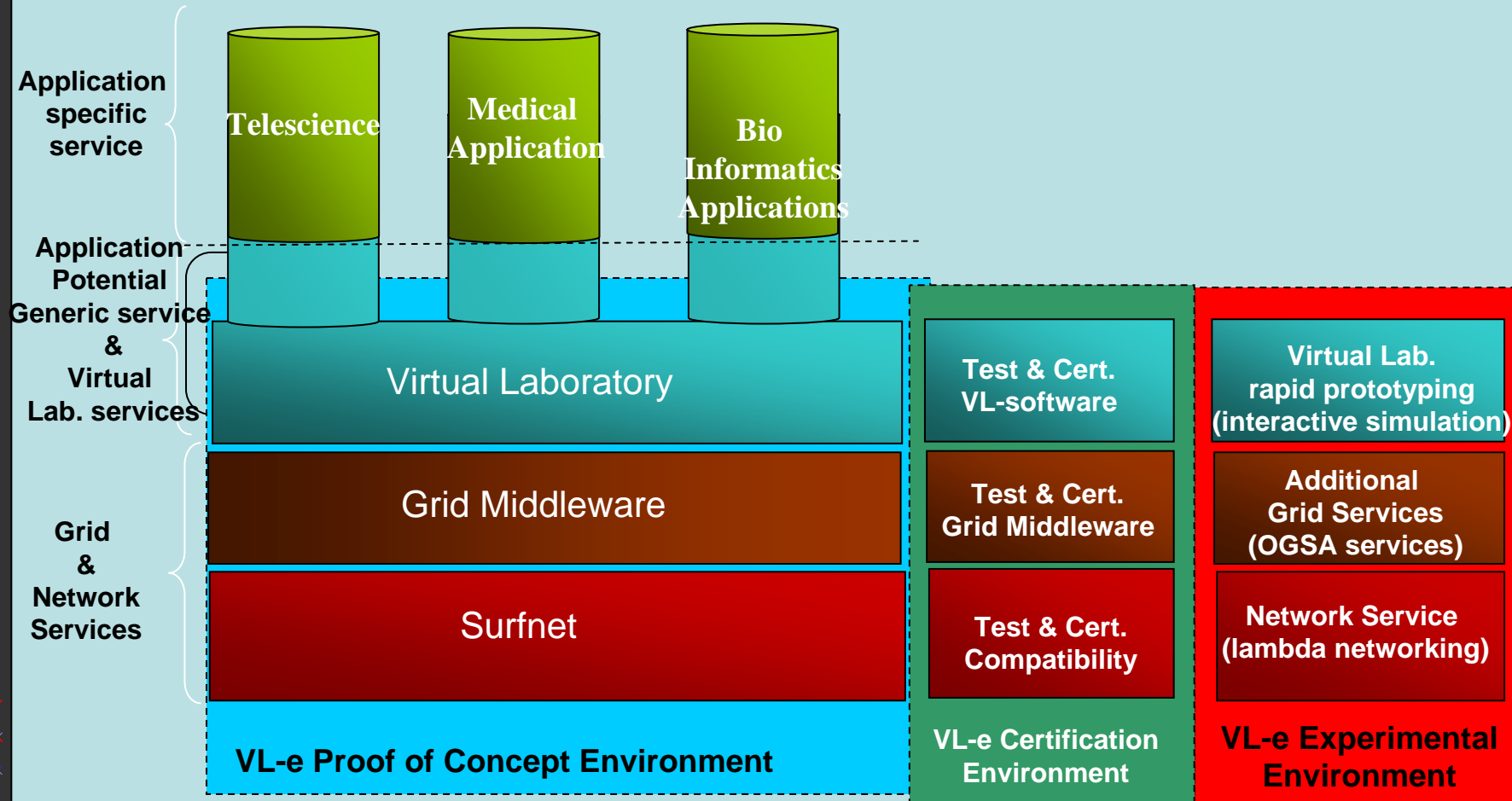
VL-e project

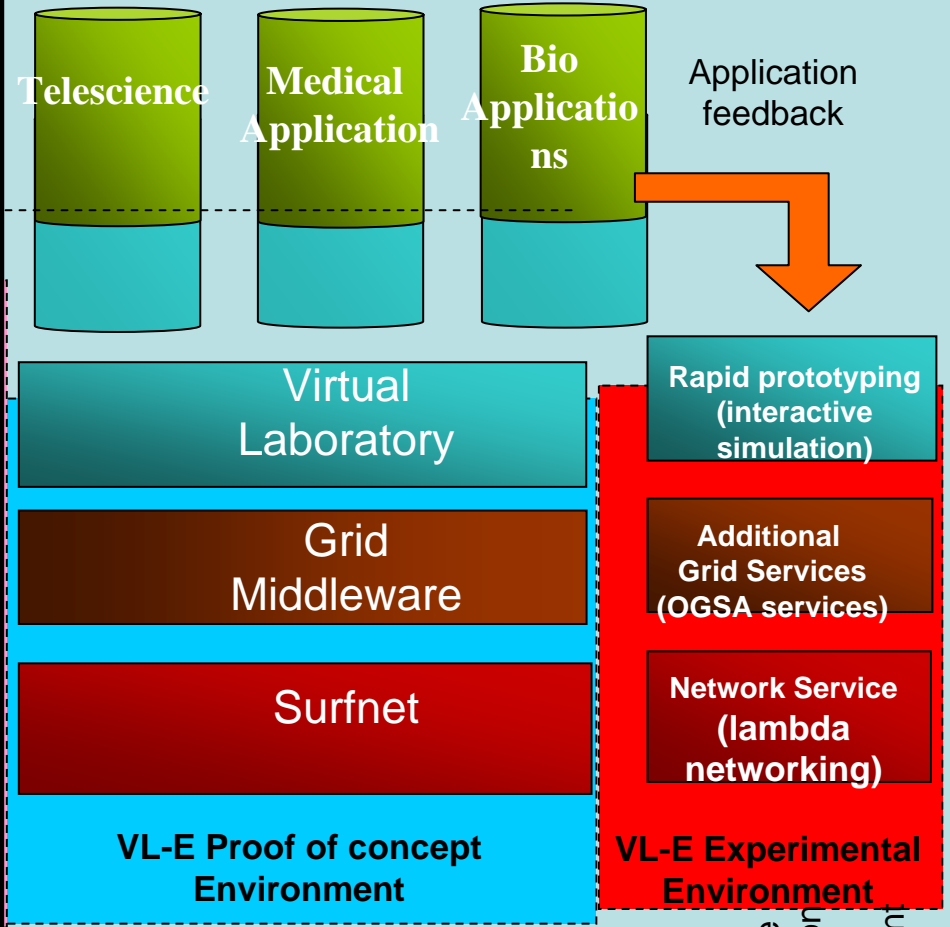
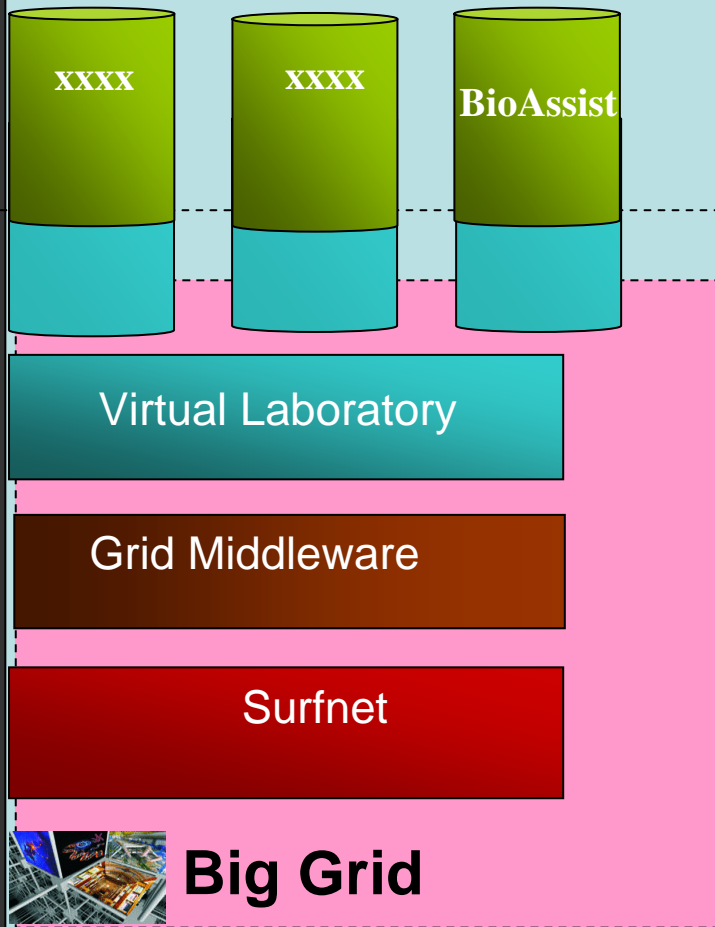


Virtual Lab for e-Science research Philosophy

- Multidisciplinary research and development of related ICT infrastructure
- Generic application support
 - ✓ Application cases are drivers for computer & computational science and engineering research
 - ✓ Problem solving partly generic and partly specific
 - ✓ Re-use of components via generic solutions whenever possible
- Rationalization of experimental process
 - ✓ Reproducible & comparable
- **Two research experimentation environments**
 - ✓ Proof of concept for application experimentation
 - ✓ Rapid prototyping for computer & computational science experimentation

The VL-e infrastructure





e-Science
Roll out
is next step

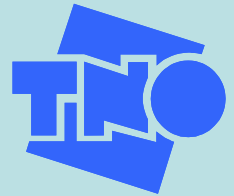
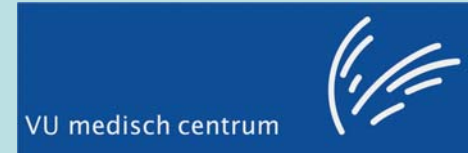
Stable
Application
& VL-e
component

Unstable
Application
& VL-e
component



Conclusion

- e-Science is not about porting current applications towards e-Science infrastructures
- To fully exploit the potential of e-Science and its ICT infrastructures one has to do **integrative** experiments
- Because of its potential to do **system level** end user science e-Science might well lead to a radical change in science methodology
 - ✓ The VL-e generic approach will help



AGROTECHNOLOGY & FOOD SCIENCES GROUP WAGENINGENUR





Netherlands
Bioinformatics
Centre

E-bioscience: a new way of life (science)

May 23, 2007

Antoine van Kampen
Scientific Director
Netherlands Bioinformatics Centre
Antoine.van.kampen@nbic.nl
www.NBIC.nl

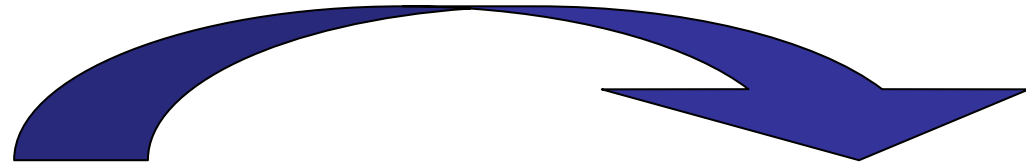


ICTDELTA
het ICT Innovatiecongres

What is Bioinformatics?

The development and application of informatics, mathematical, and statistical methods in life sciences.

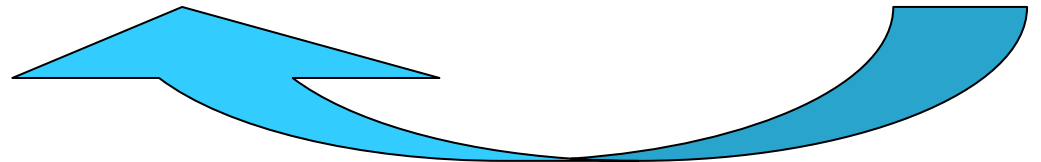
- *convert data to knowledge
- *generate new hypotheses



DATA

Knowledge

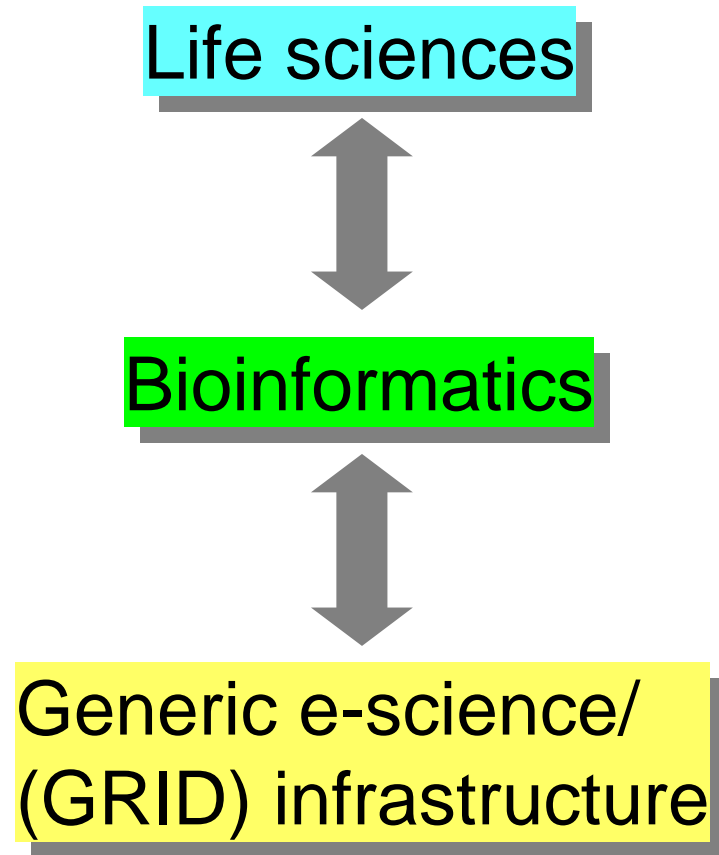
Enabling science for genomics



- *Design new experiments

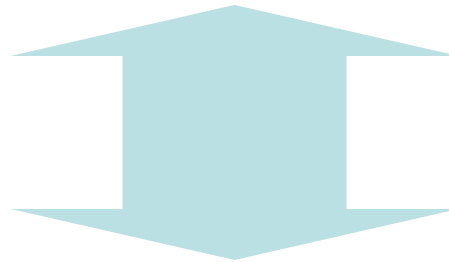
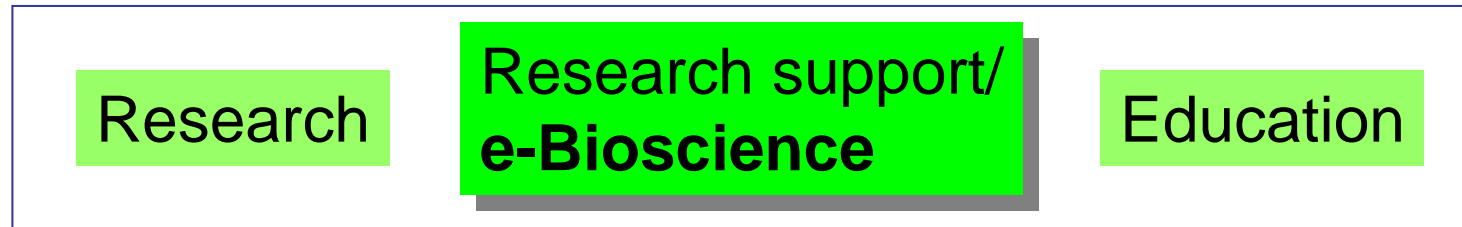
Bioinformatics as interface

- e-Bioscience
- How can we make generic e-science methodologies and (GRID) ICT infrastructure of benefit to life sciences?



Netherlands Bioinformatics Centre (NBIC)

Programmes



Life science researchers (end-users)

Deliver tools and databases to end-users

Research



Research support/
e-Bioscience



Life science researchers (end-users)

High-throughput experimental technologies in life sciences

Cell/Tissue

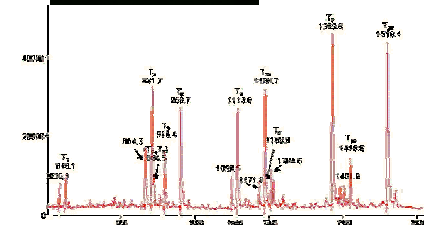
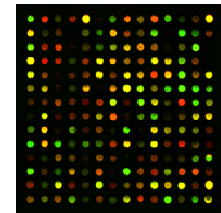


Determine complete DNA sequence of organism (e.g., mutations)

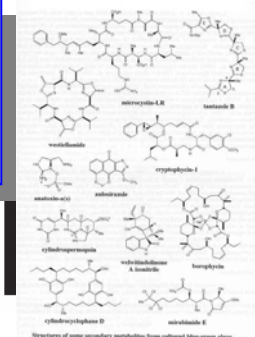
Measure expression of all genes.

Identify many proteins or their expression level.

Identify many metabolites or their concentration.



2006 To F mass spectrum of the hydrolyzate of a high-molecular weight (105 000) glucan (D-71.0)

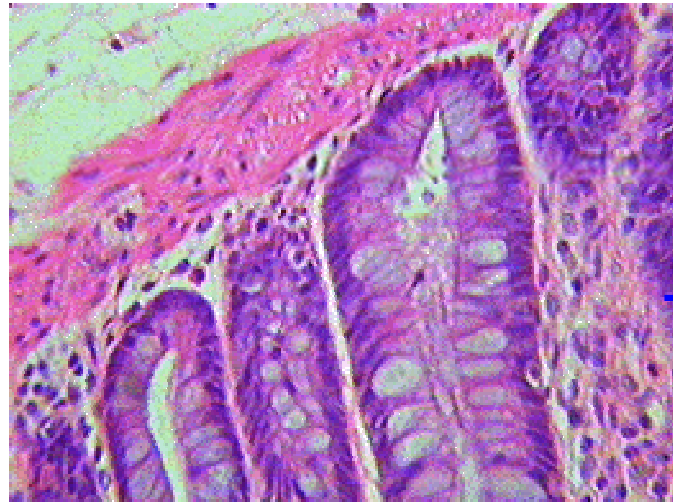


nl

lands
matics

Compare normal versus cancer

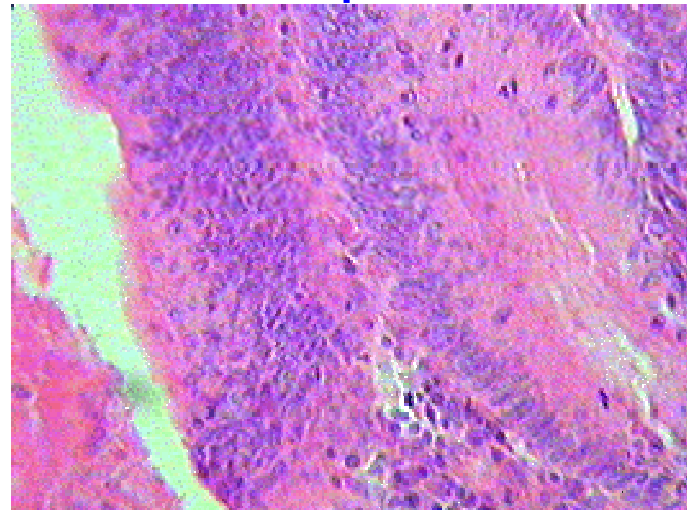
High magnification of a normal human colon cell



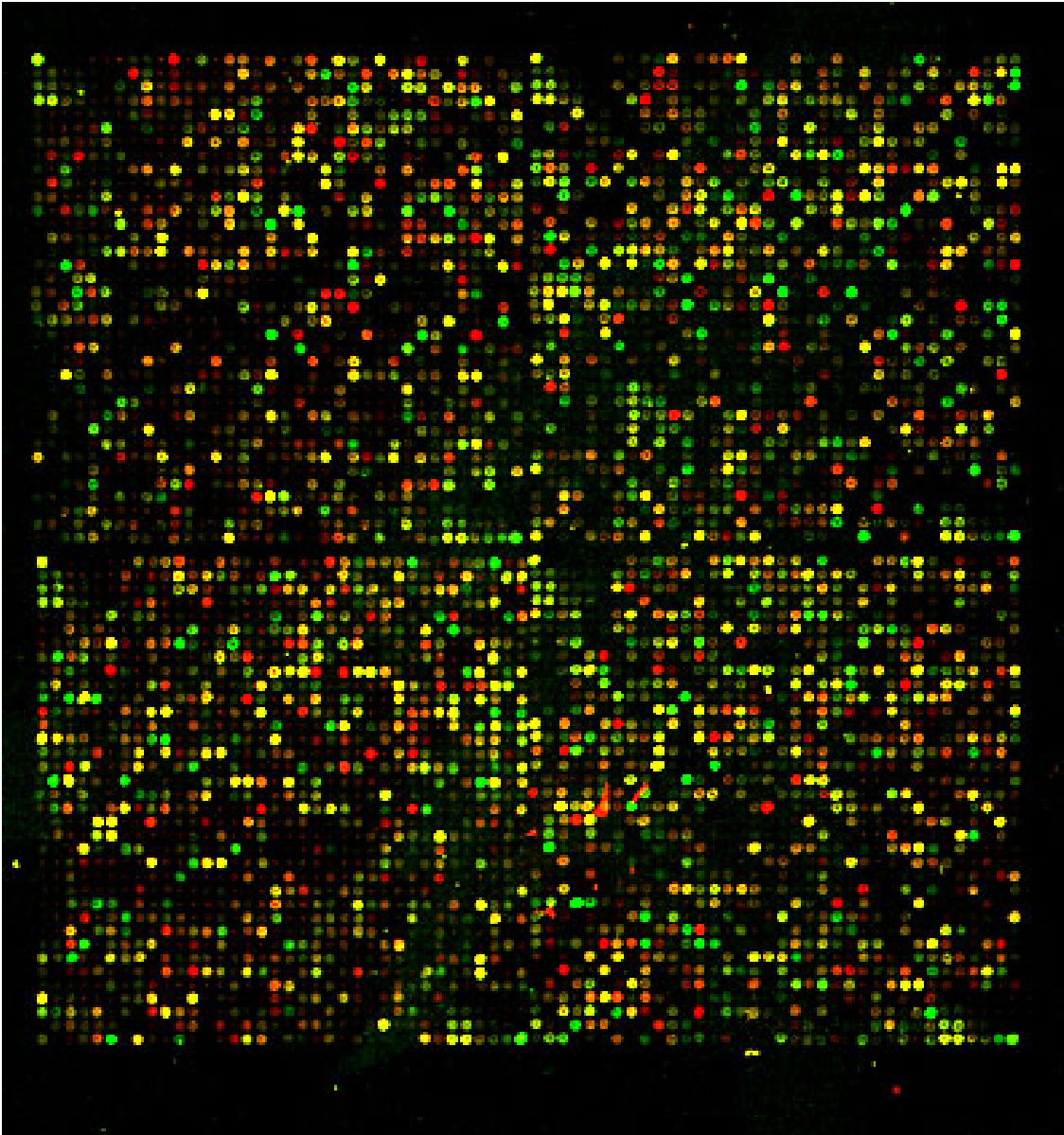
Data

Data

High magnification of a human colon cell with carcinoma



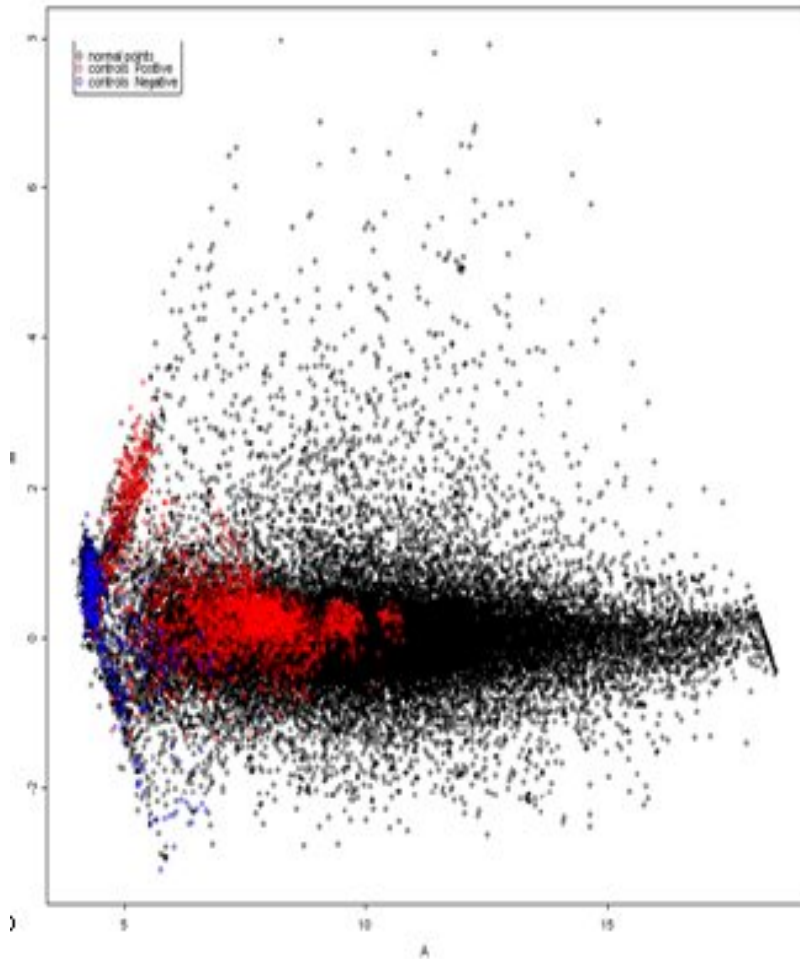
DNA microarrays



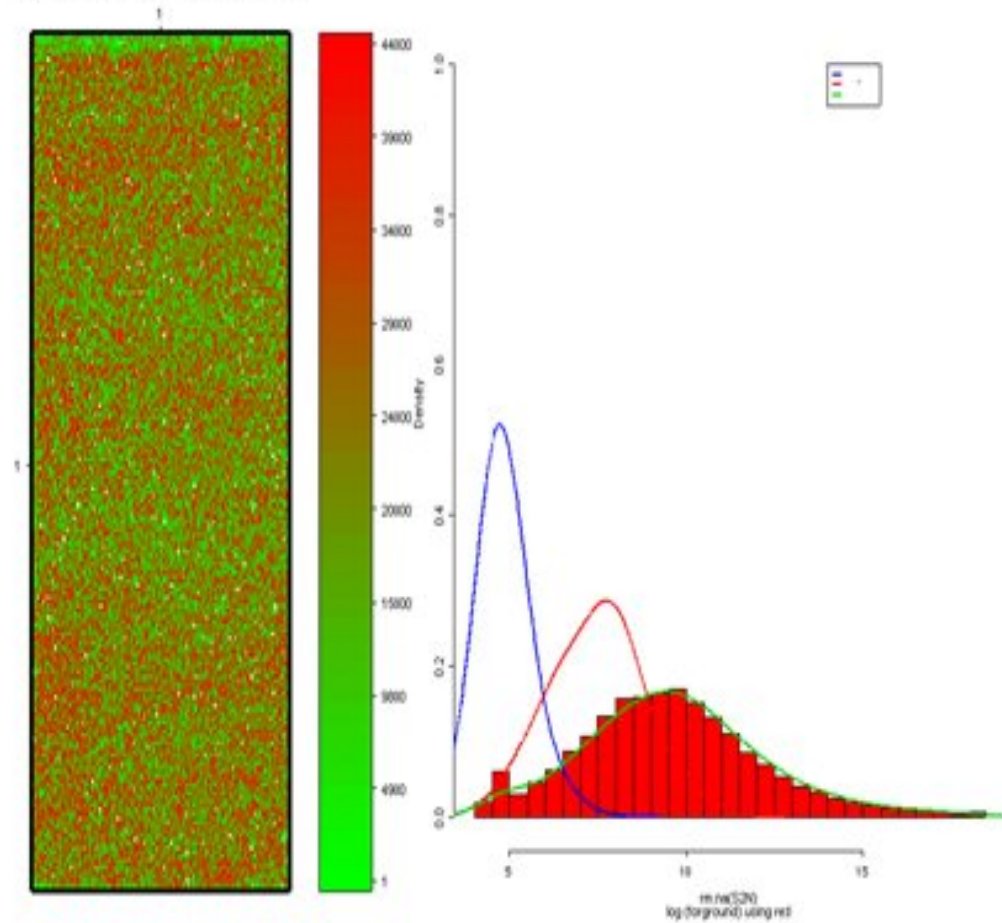
Single experiment:
30.000 – 40.000 genes

Requires dedicated
approaches for
analysis

Quality control

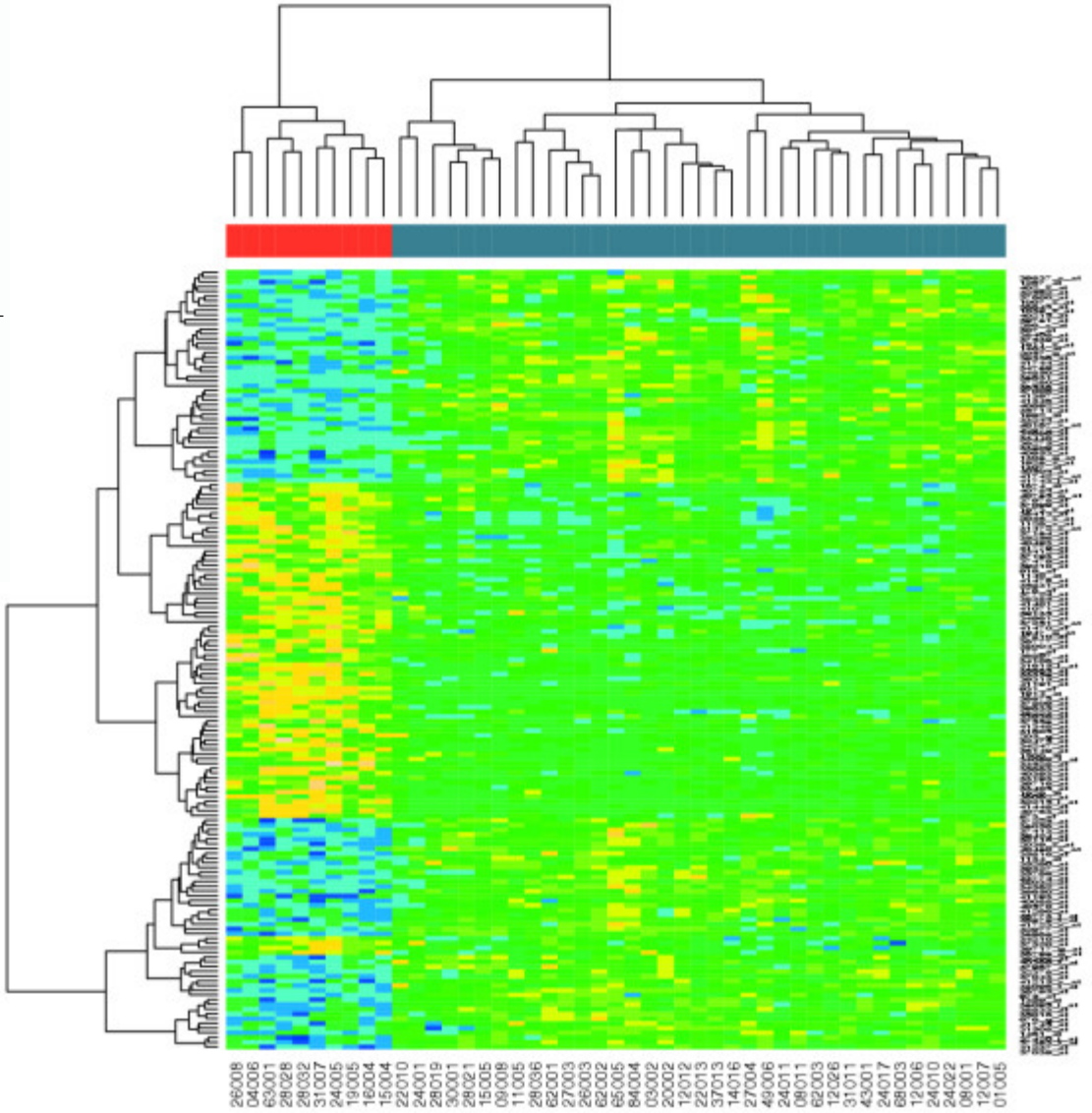
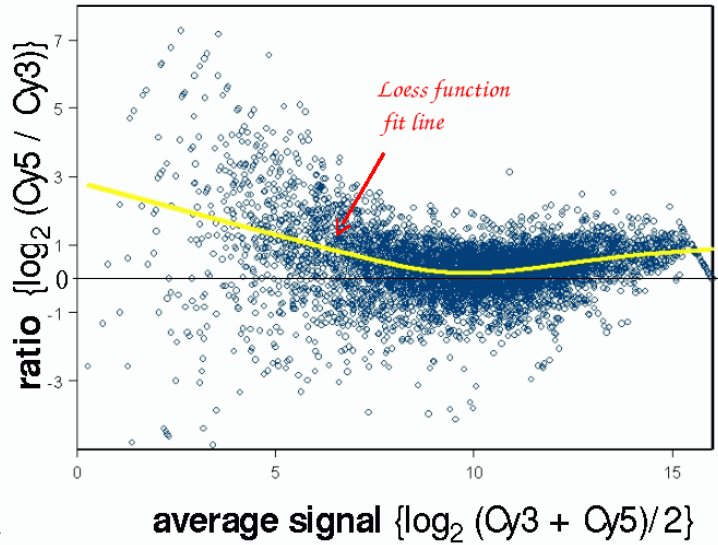


Spatial plot using Ranked M value values

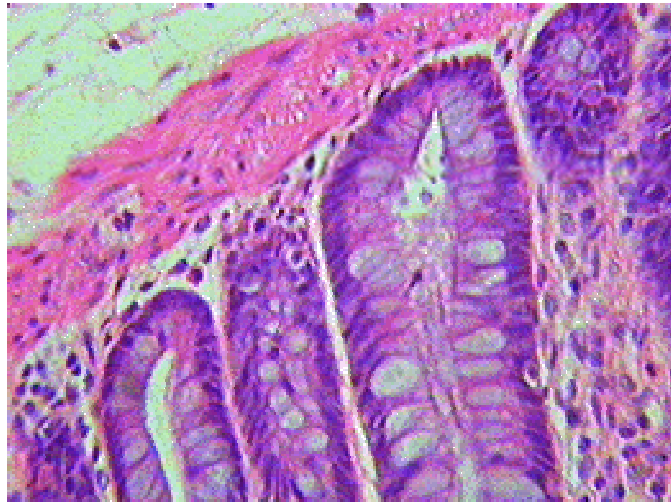


Normalization and statistical analysis

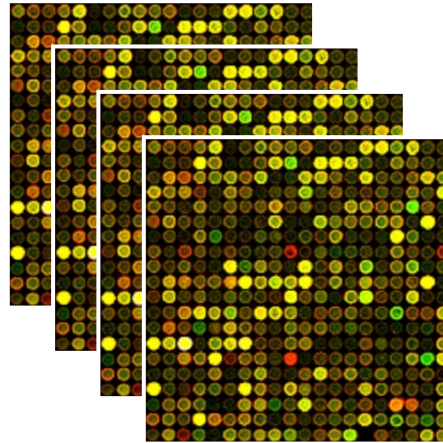
Loess Function



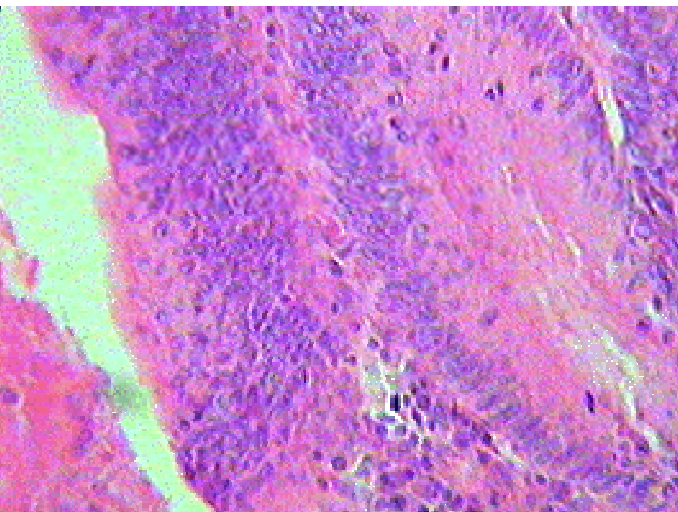
Molecular diagnostics



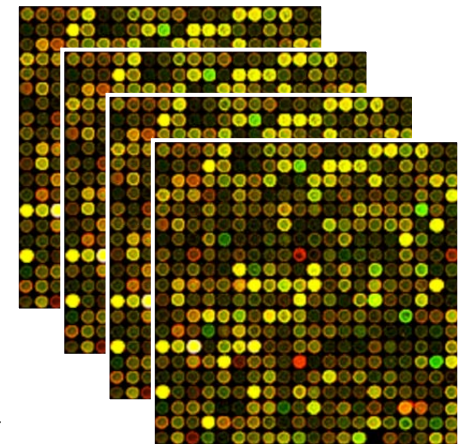
High magnification of a normal human colon cell



Which genes discriminate between normal & patient

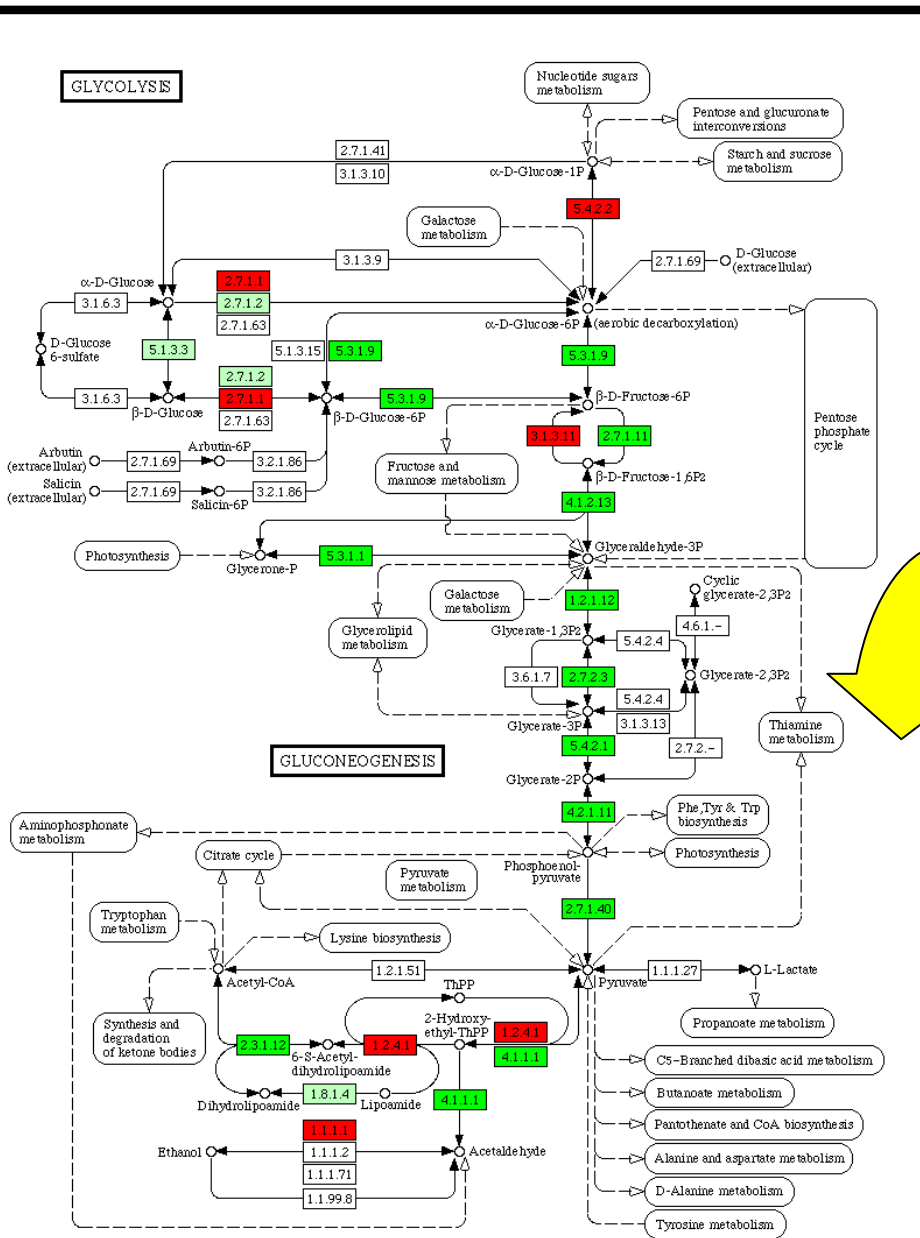


High magnification of a human colon cell with carcinoma

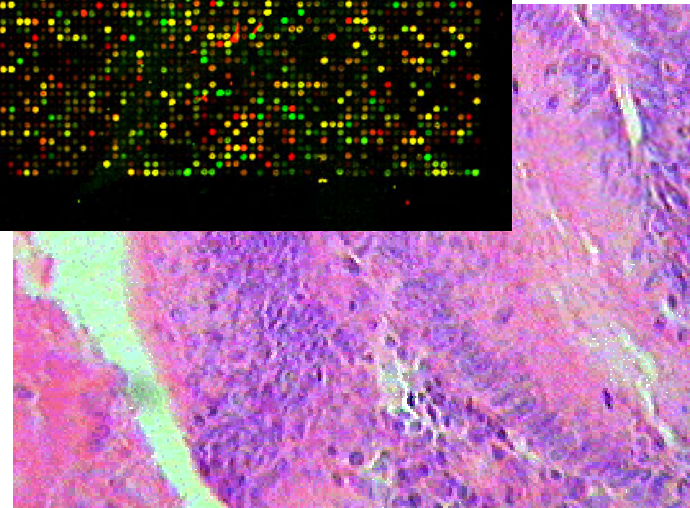
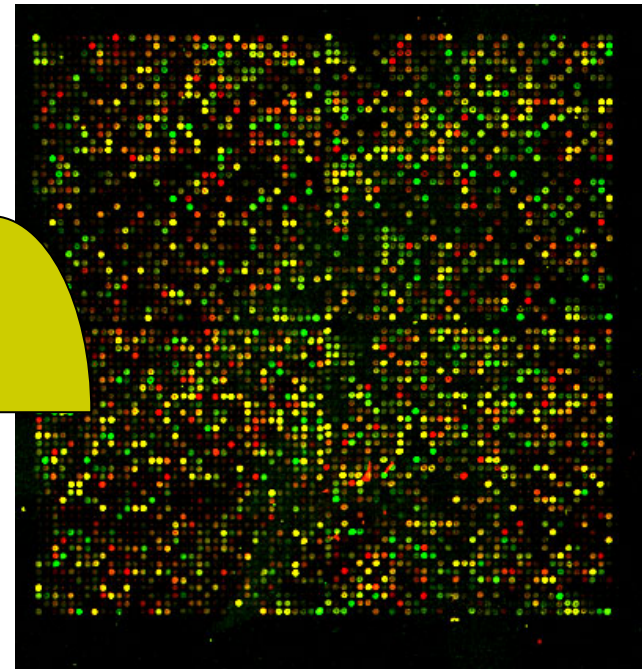


This requires statistical analysis of the data. Complex!

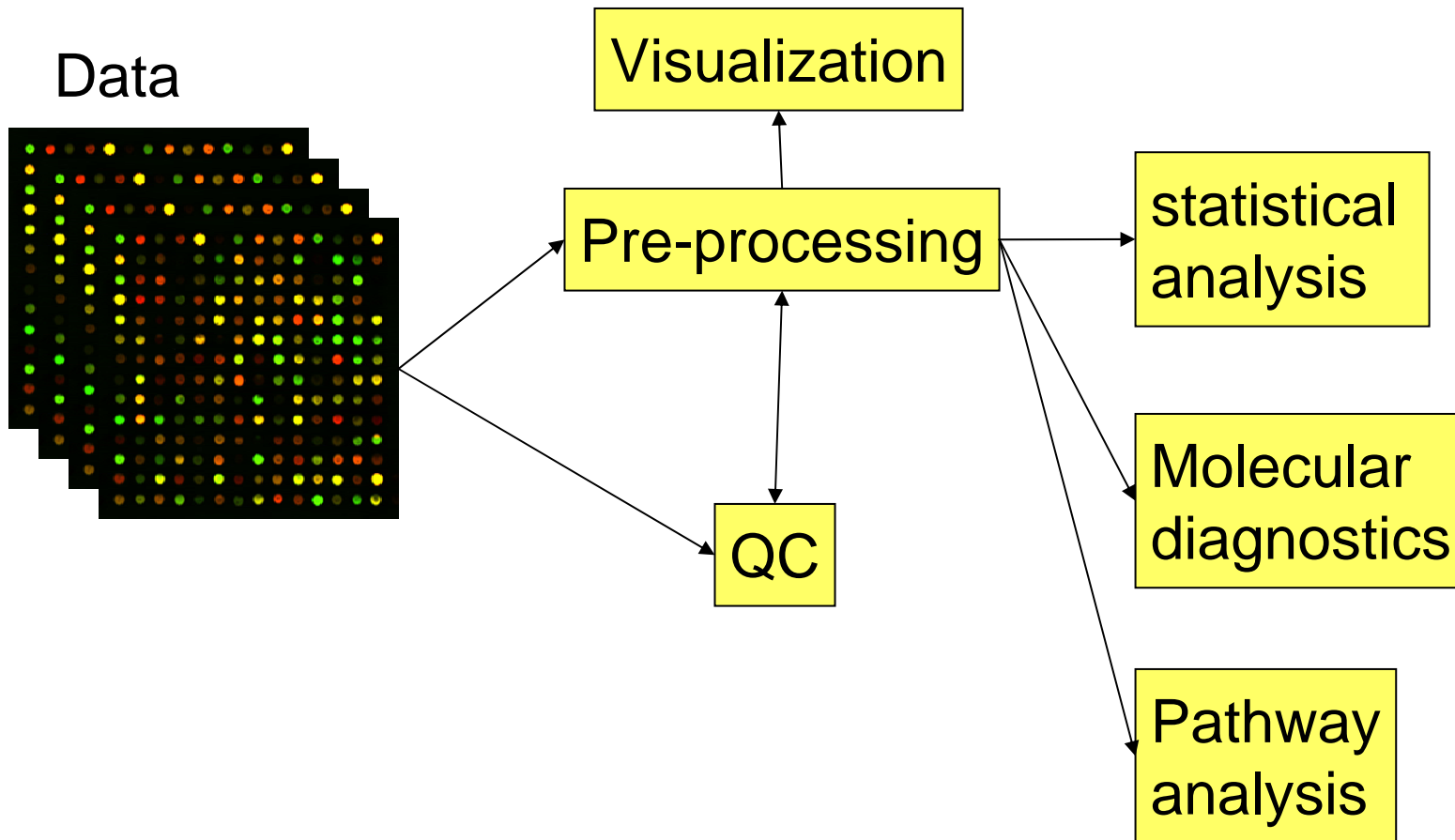
Understanding molecular processes

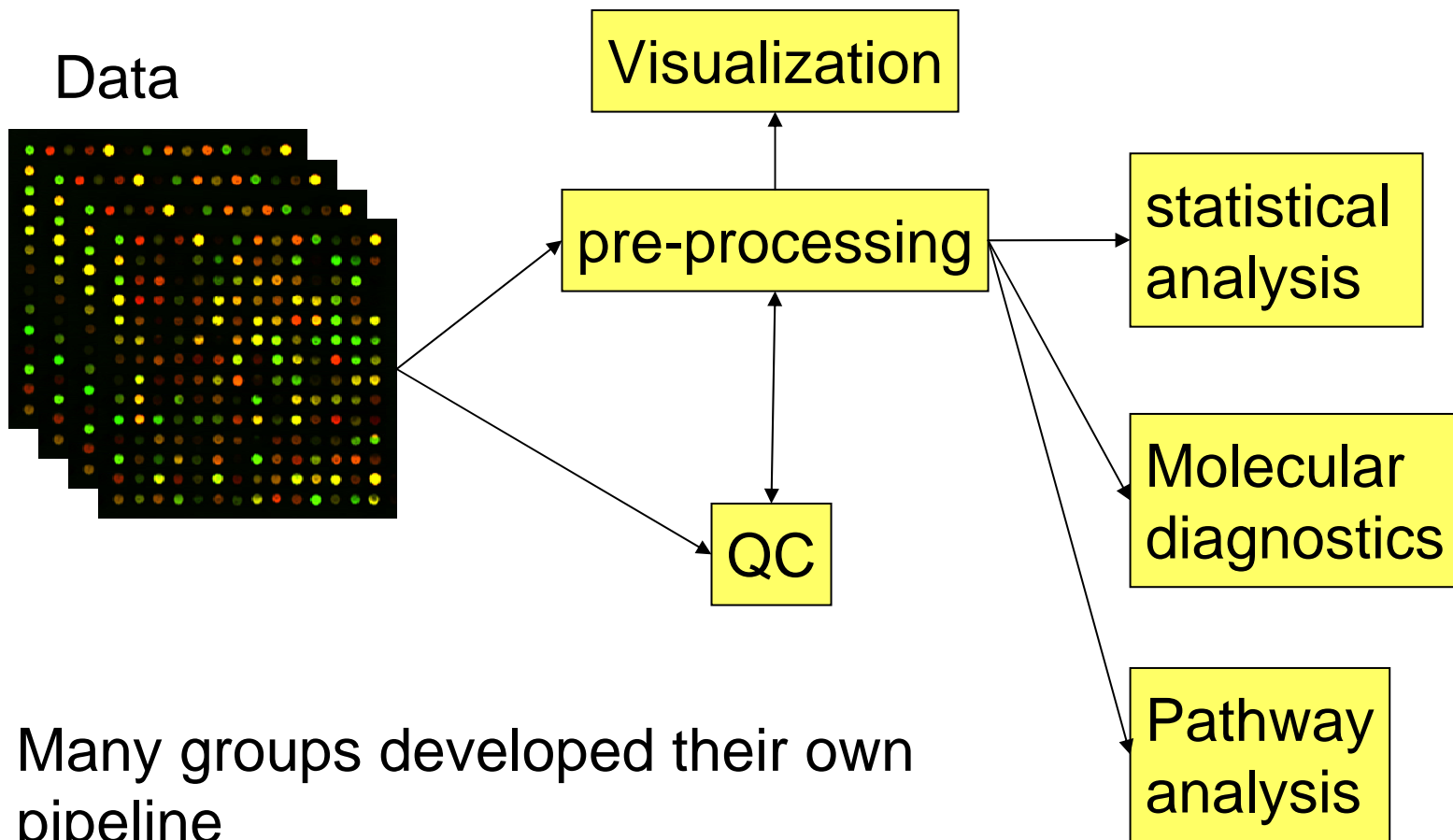


Integration with pathway databases (eg KEGG)



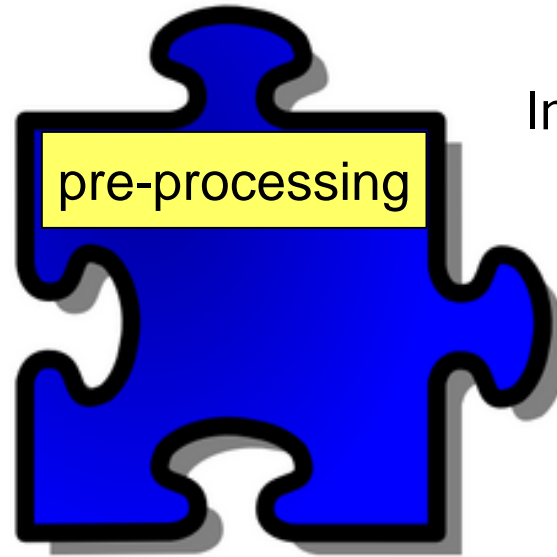
Microarray in-silico experimentation pipeline



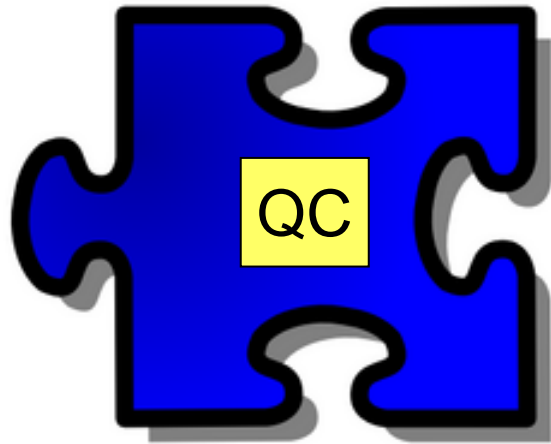


- Many groups developed their own pipeline
- Large effort
- Development of modules may require specific expertise
- Difficult to use (state-of-the-art) methods of other groups

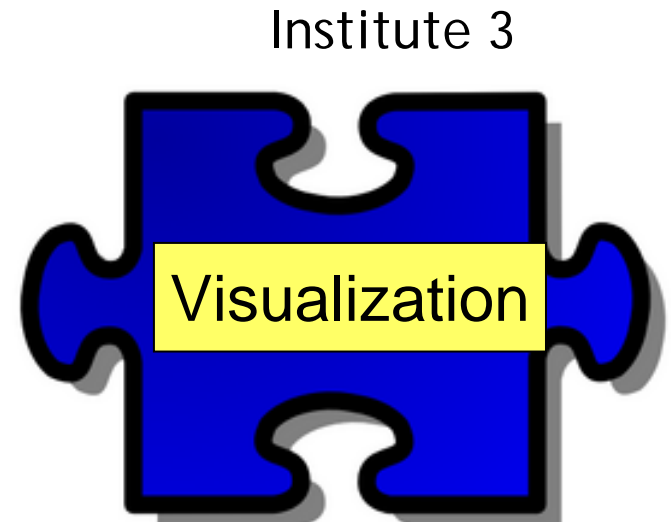
...but how to share tools, data, expertise?
...how to jointly solve problems?



Institute 2



Institute 1



Institute 3

....e-Bioscience



Collaborate to develop experimentation pipeline

Service oriented architecture

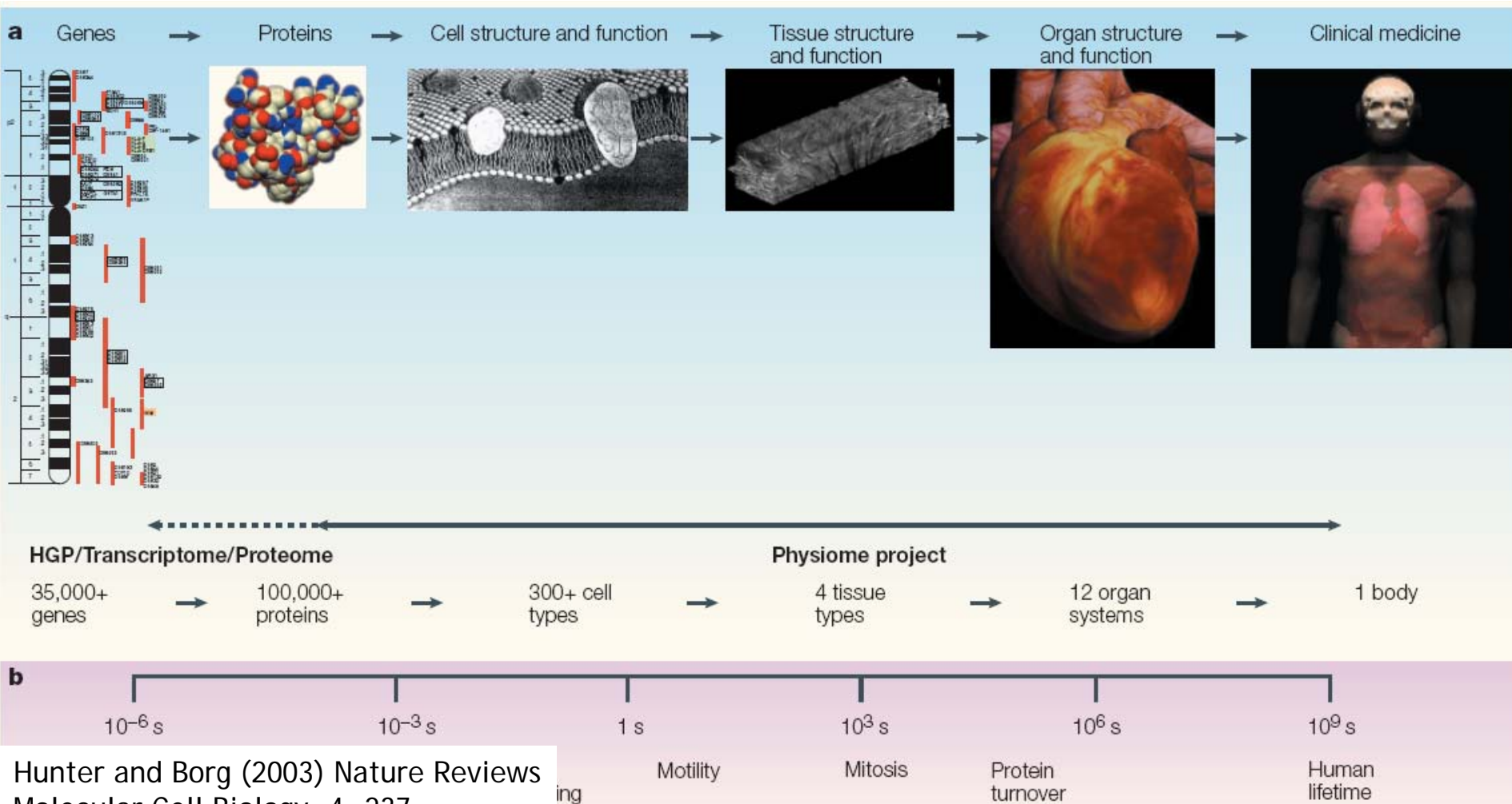
Share environment with *de facto* standards; use common approaches

Generic e-Science infrastructure (VLe)

Life sciences GRID (NCF pilot, BIG GRID)

Basic infrastructure (SURFnet, Gigaport)

Moreover, science is becoming increasingly complex and multi-disciplinary



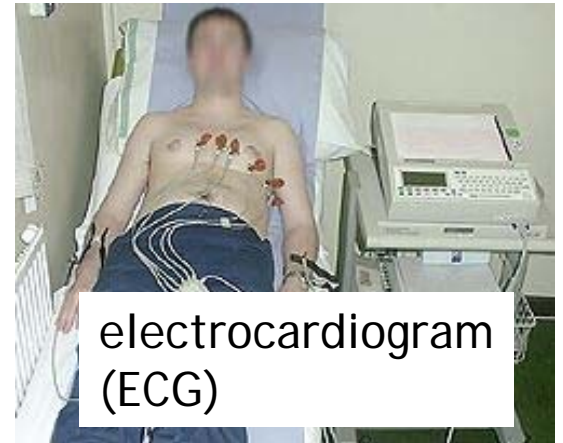
Clinical chemistry



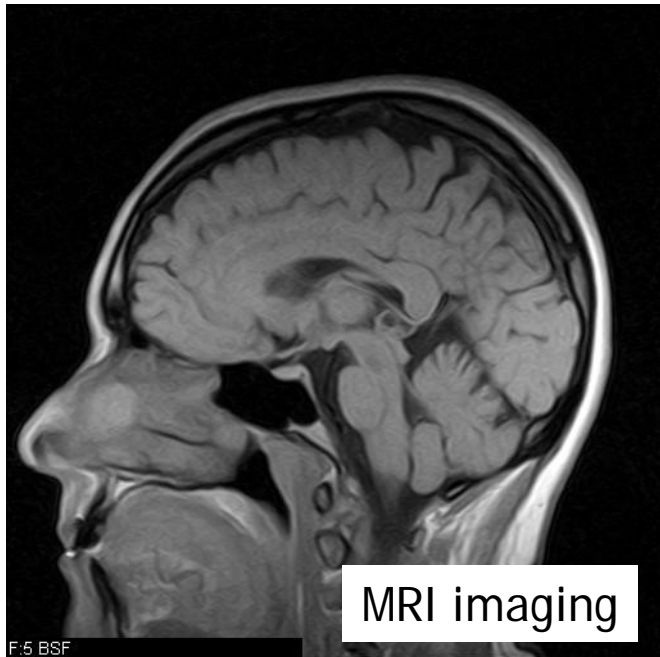
Life style



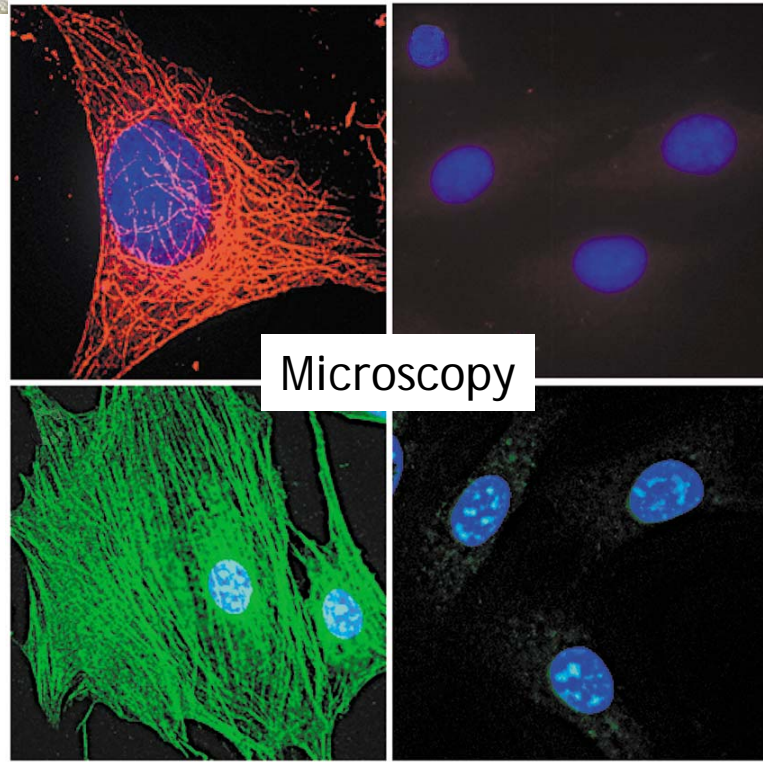
electrocardiogram (ECG)



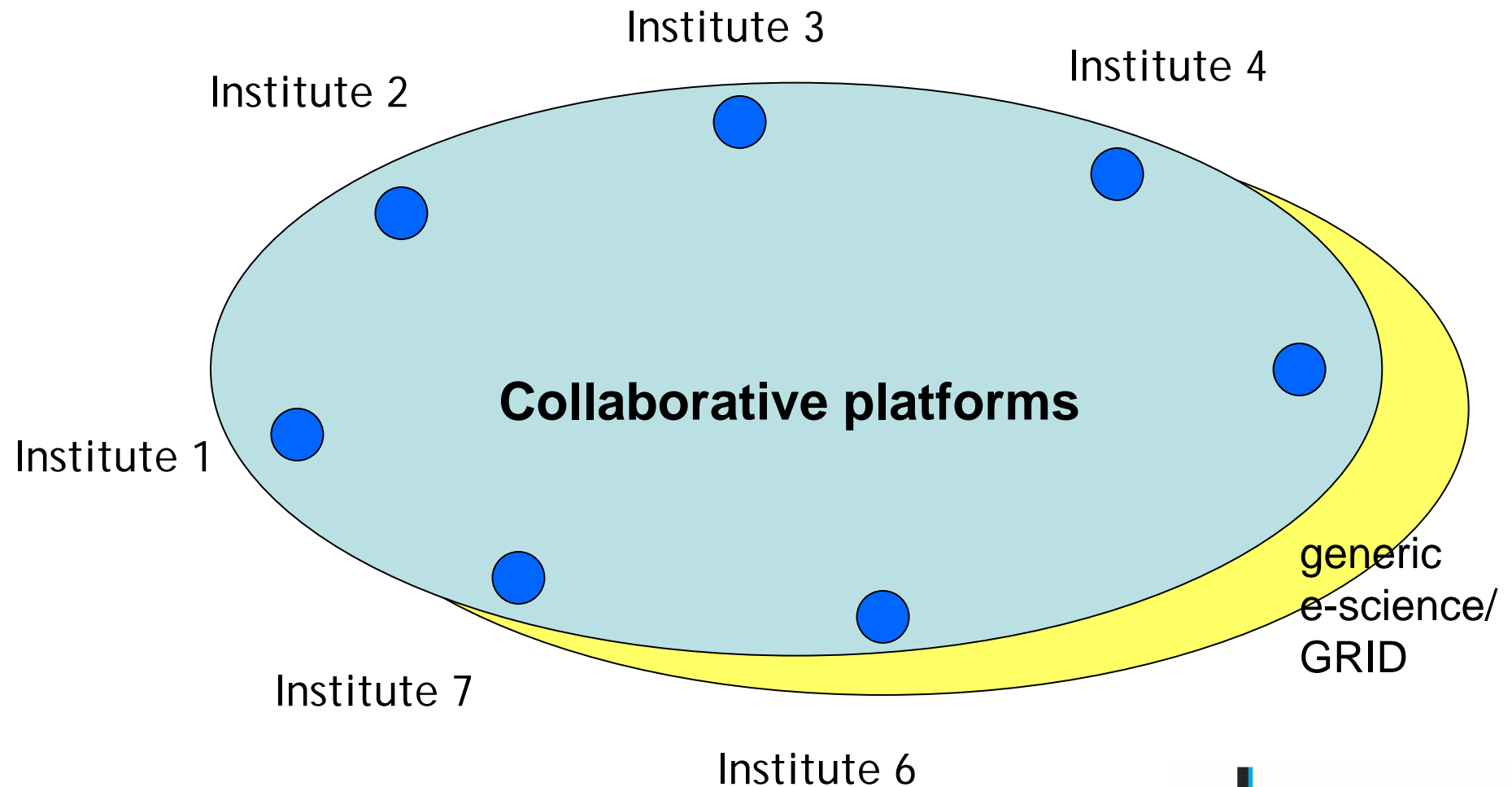
MRI imaging



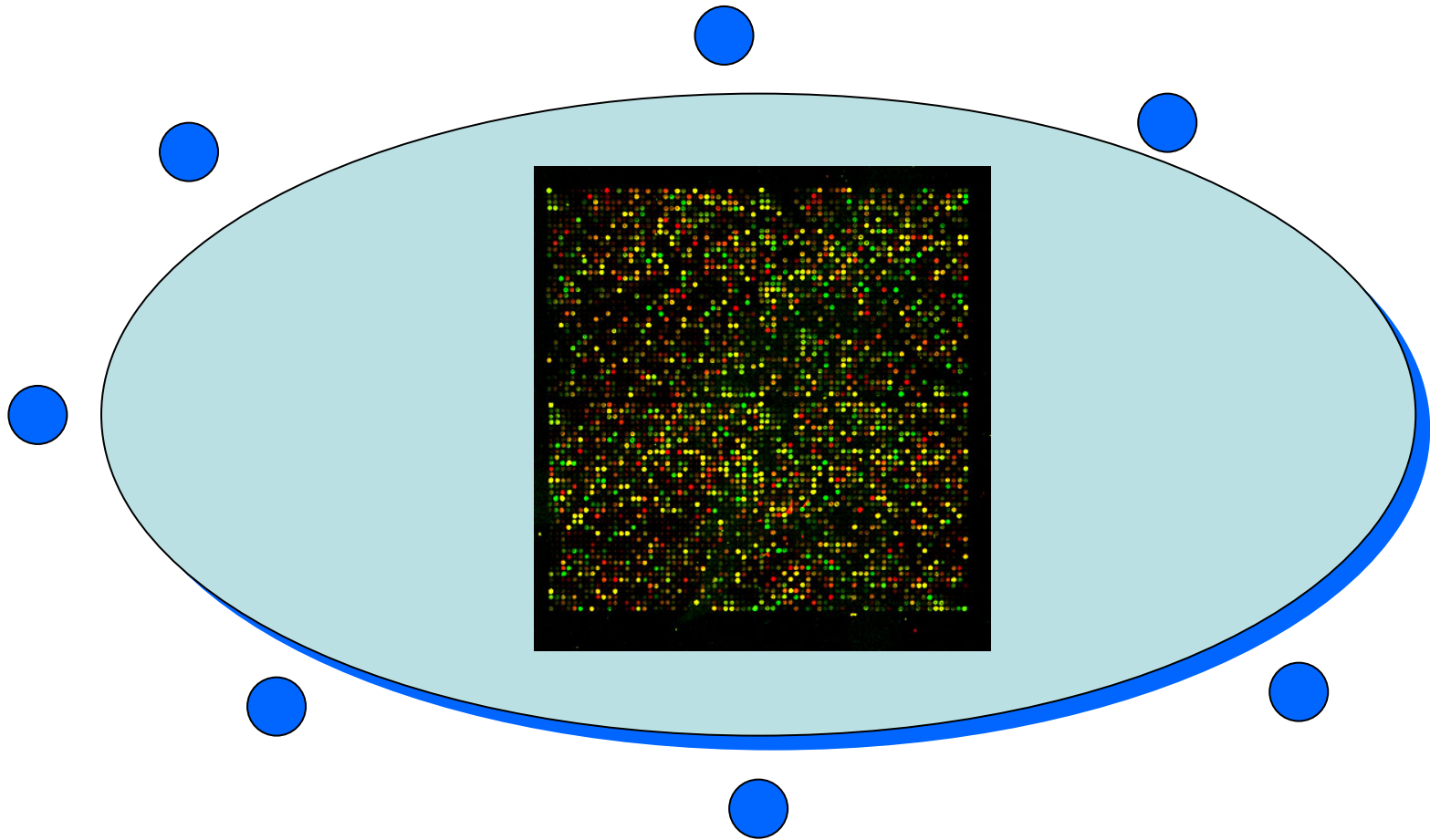
Microscopy



Collaborative platform to address research questions

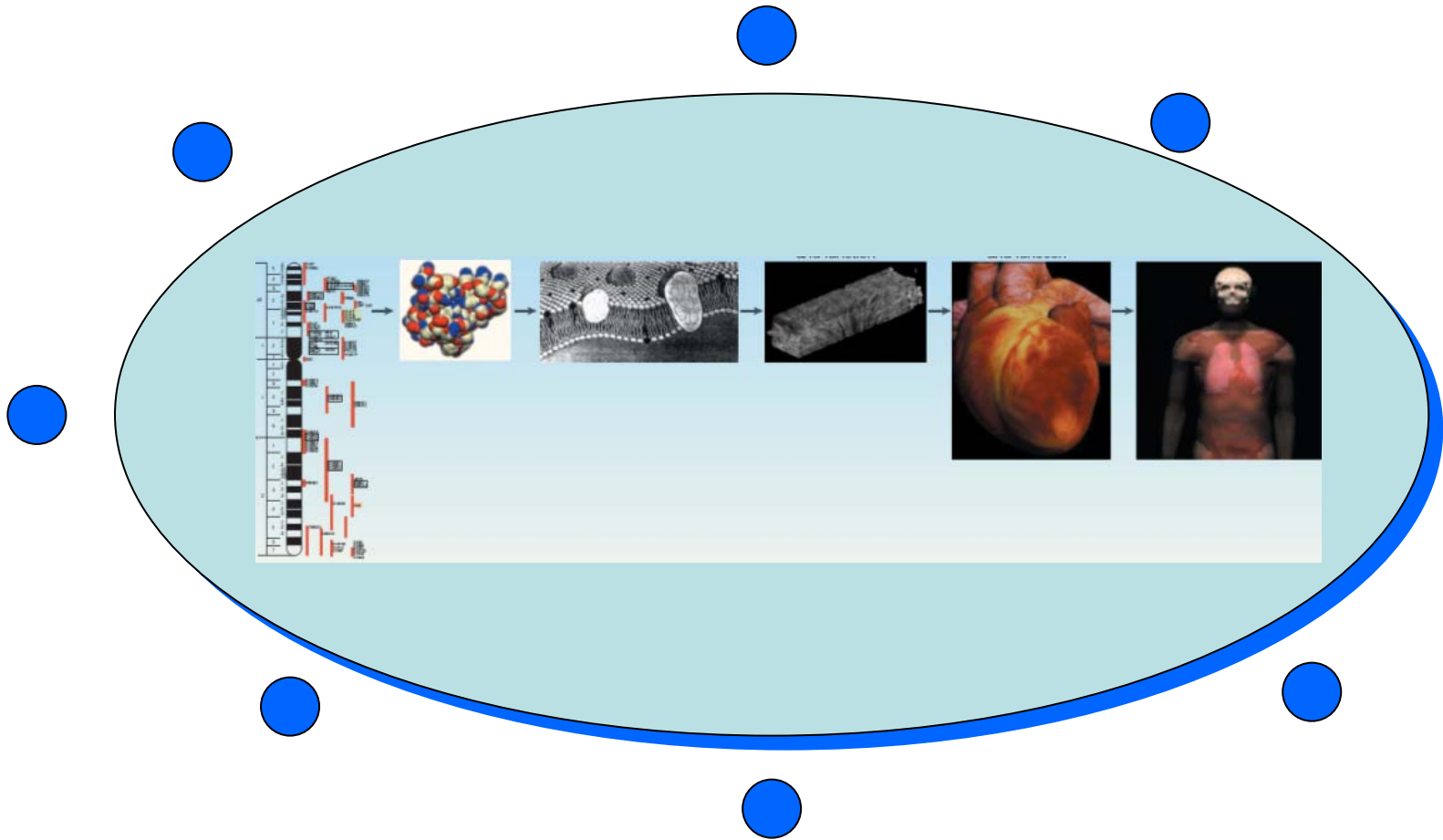


Collaborative platform for microarray research



Collaborative platform for systems biology

Truly multi-disciplinary!

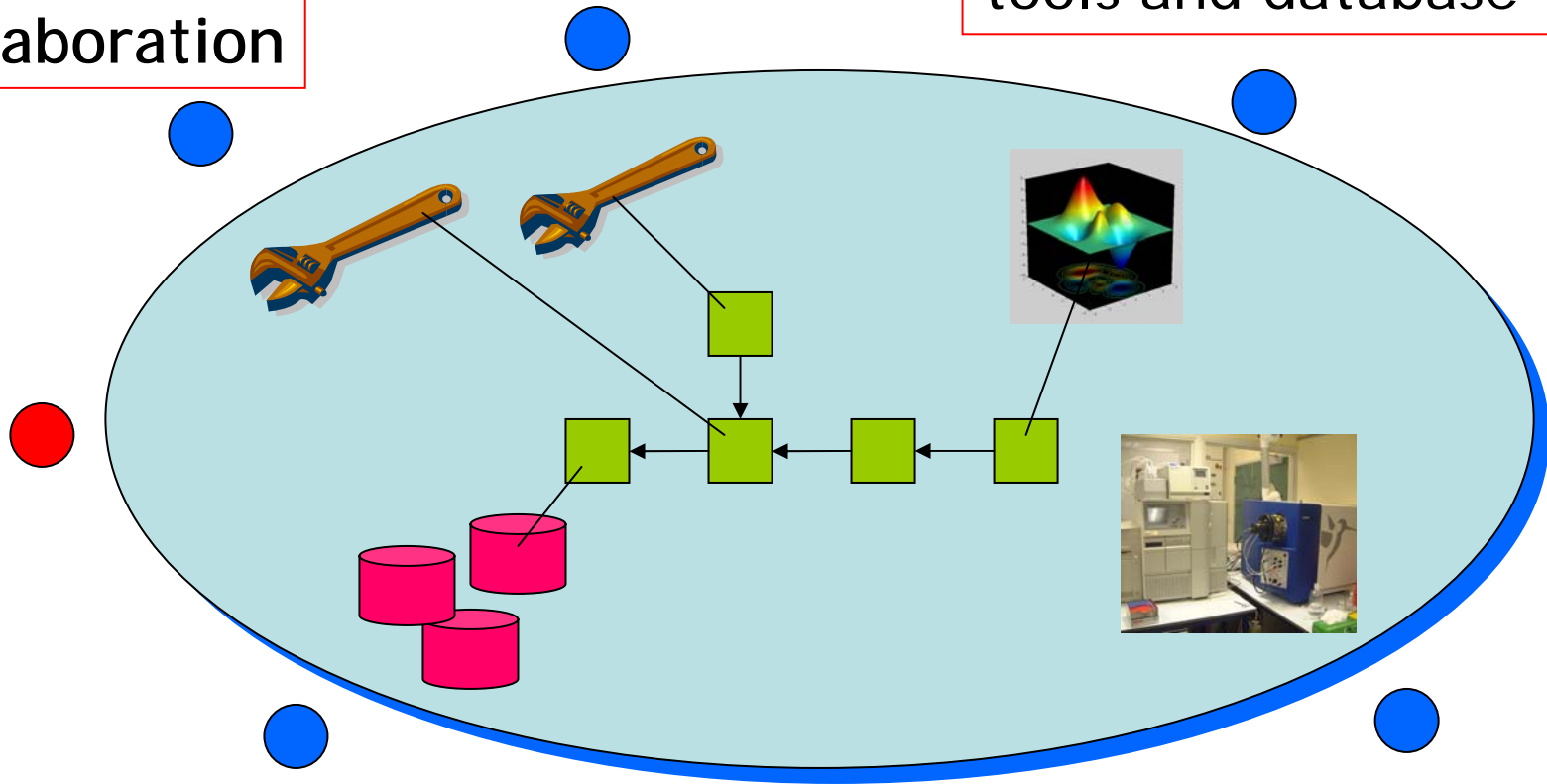


e-Bioscience

Exchange of data & tools & expertise

Define standardized workflows that connects tools and database

Collaboration



Accelerate research,
avoid redundancy,
reduce costs

Users

Solve the big
scientific questions

The e-Bioscience challenge

E-(bio)science/GRID are not production systems, instead

- developments/research on e-(bio)science and GRID is ongoing.
- Experience from current and future cases will mature this approach
- Collaborative platforms require sufficient time to be designed and implemented
- Requires specific expertise
- Investments (hardware, software, personnel)
- Willingness to collaborate

Life sciences



Bioinformatics



Generic e-science/
(GRID) infrastructure

Acknowledgements

University of Amsterdam

- Prof. dr. L.O. Hertzberger
- Dr. T. Breit
- Dr. M. Bouwhuis

CERN

- Dr. C. Jones

Virtual Laboratory for e-Science (VLe; www.vl-e.nl)

BIG GRID (www.nikhef.nl/grid/BIG)

NBIC (www.nbic.nl)