

Domain Drivers – Lecture 2

Professor Richard O. Sinnott
Director, Melbourne eResearch Group
University of Melbourne
rsinnott@unimelb.edu.au

Objectives

- To give the “big picture” of why we need Cluster and Cloud Computing
 - This lecture is not focused on technologies, but on giving examples of how challenges are shaping the technological landscape
 - ...and how on-going/completed projects have met/are meeting those challenges
 - Many perspectives
 - Big Data – and the hype!
 - Big Compute
 - Big Distribution
 - Big Collaboration
 - Big Security
 - ...

Noting...

- Often similar challenges facing many research domains
- Tools, technologies and methodologies have been/can/are evolving to tackle these challenges
 - There is a huge amount of work still to be done
 - Don't believe the hype!!!
 - The pace of research evolution FAR outweighs the pace of IT know-how to deal with the challenges
 - Domain knowledge!!!

Focal Point

- Examples from different research domains
 - {Computational/Data/Distributed/Collaboration/Security} bound...
 - High Energy Physics
 - Astrophysics
 - Electronics
 - eHealth & Bioinformatics
- BREAK
 - Clinical/biomedical research
 - Social Sciences/Urban research domain
 - (prelude to workshop and assignment 2)

Completed

- National e-Science Centre (I, II, III)
- Dynamic Virtual Organisations for e-Science Education
- Biomedical Research Informatics Delivered by Grid Enabled Services
GridNet, GridNet 2
- Grid Enabled Microarray Expression Profile Search
- Glasgow early adoption of Shibboleth
- Joint Data Standards Survey
- ESP-Grid
- HPC Compute cluster award // Sun industrial sponsorship
- OGC Collision
- OMII-Security Portlets // OMII-RAVE
- Integrating VOMS and PERMIS for Superior Grid Authorization
- NCeSS
- CESSDA PPP
- Pharming of Therapeutic RNA
- Grid Enabled Occupational Data Environment
- Towards an e-Infrastructure for e-Science Digital Repositories
- Grid enabled Biochemical Pathway Simulator
- Virtual Organisations for Trials and Epidemiological Studies
- A European e-Infrastructure for e-Science Repositories
- Modelling, Inference and Analysis for Biological Systems up to the Cellular Level
- Drug Discovery Portal
- Parliamentary Discourse
- Scots Words and Placenames
- Qvolution stress management survey system
- Advanced Grid Authorisation through Semantic Technologies ShinTau
- AlstromUK VRE
- Grid-enabled Virtual Safe Settings
- Clinical Streaming Transcription Software
- Enhancing Repositories for Language and Literature Researchers (ENROLLER)
- Proxy Credential Auditing Infrastructure for the NGS
- Scottish Bioinformatics Research Network (SBRN)
- Generation Scotland Scottish Family Health Study
- Breast Cancer Tissue Biobank
- Data Management through e-Social Science (DAMES)
- Meeting the Design Challenges of nanoCMOS Electronics (nanoCMOS)
- EU FW7 AvertIT
- EU FW7 EuroDSD
- NeSC Research Platform (NRP)
- NeSC Information Network (NIN)
- ESF Network for Study of Adrenal Tumors
- Scottish Health Informatics Platform for Research (SHIP)
- National E-Infrastructure for Social Simulation (NeISS)
- EU R4SME Diagnosis of Parkinsons Disease (DiPAR)
- Automating River Pollution Detection (CAPIM)
- Endocrine genomics Virtual Laboratory (endoVL)
- DSDNetwork Australasia

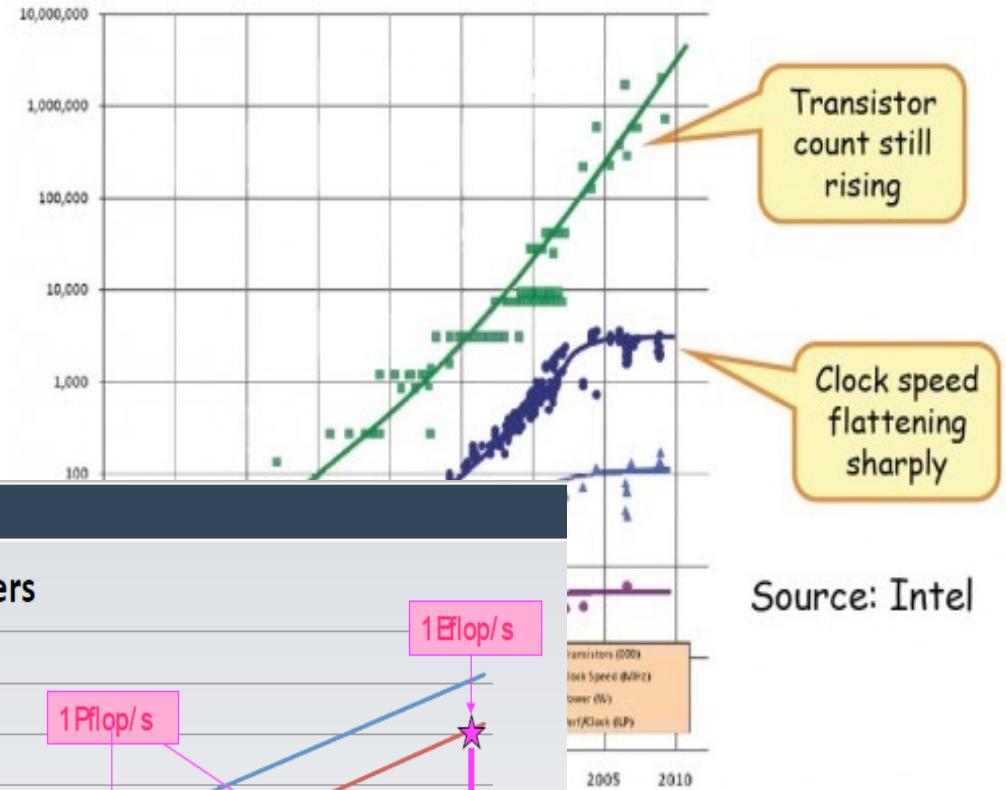
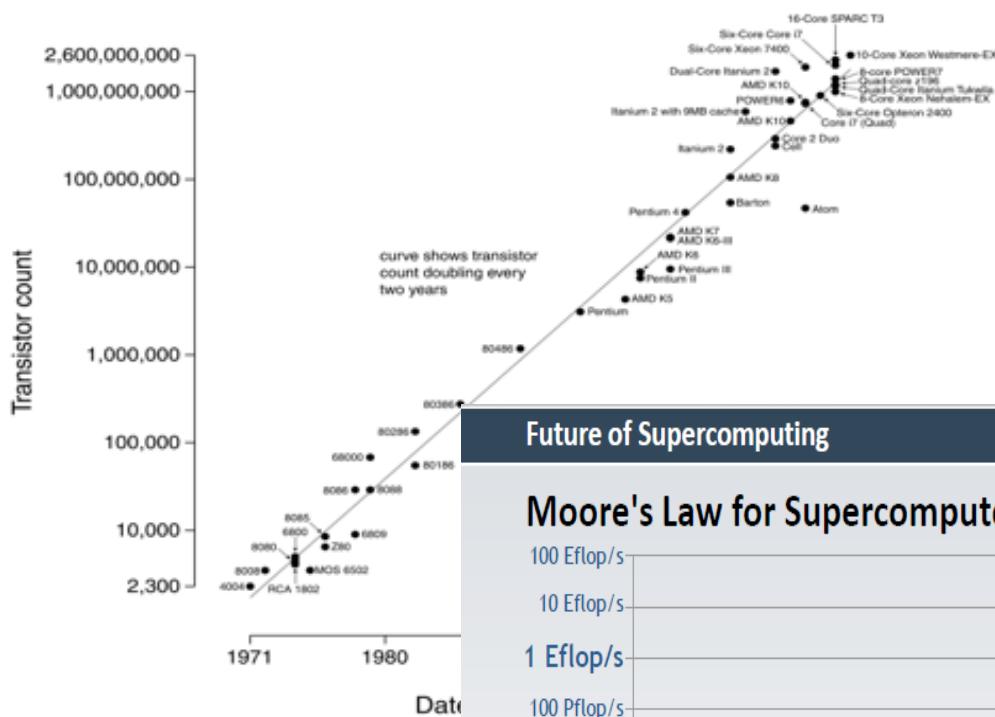
Project Portfolio

Subset of On-Going

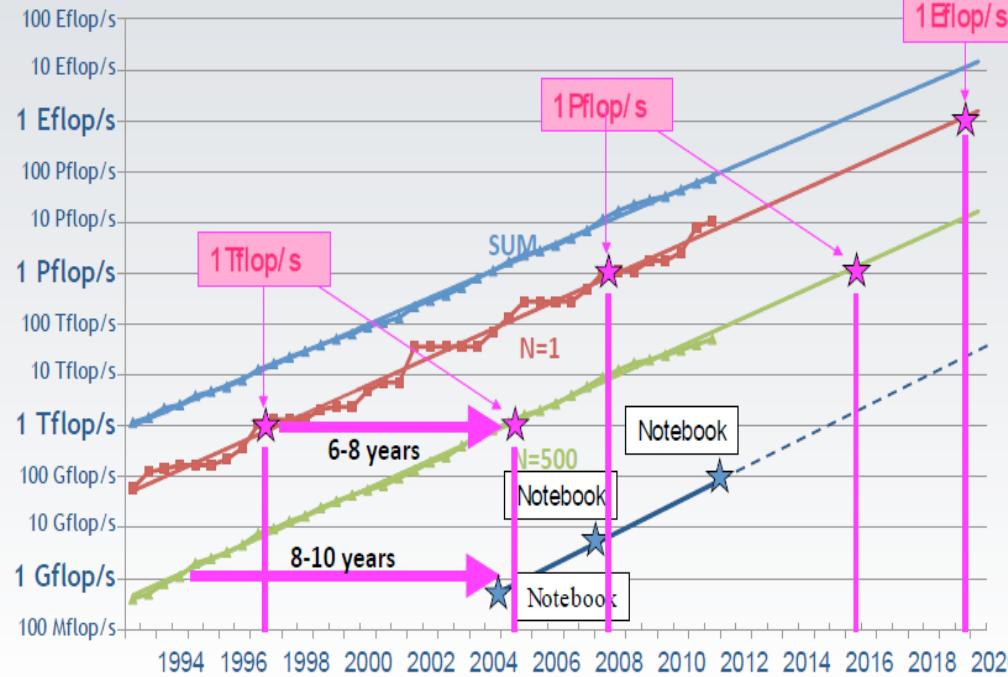
- EU European Platform for Study of Wolfram, Alstrom, Bardet Biedl (EuroWABB)
- Multicenter prospective study of biochemical profiles of monoamine-producing tumors (PMT Study)
- European Society of Hypertension Study on Pheo/PGL
- International DSD
- EU FW7 European Network for Study of Adrenal Tumors Cancer Research Platform (ENSAT-CANCER)
- VicHealth Health Indicators and Spatial Objective Data
- National Spinal Injury Research Platform
- Australian Urban Research Infrastructure Network (AURIN)/Spatial Urban Data Observatory (SUDO)
- Epilepsy e-Learning portal
- Type-1 Diabetes study of environmental factors on onset of T1D
- Australian Diabetes Data Network (ADDN)
- International Niemann-Pick A, B and C Registry
- Carlton Connect Data Journalism in the Big Data Era
- FAMIAN - Combined 18F-fluorodeoxyglucose positron emission tomography and 123I-Iodometomidate Imaging for Adrenal Neoplasia
- Melbourne Genomics Health Alliance (variant DB)
- NeCTAR Cloud Encryption/Decryption and Secure Deletion
- CRE for Protection of Pancreatic Beta Cells
- Airbox (Atmospheric Physics and Climate Research)
- NESP Clean Air and Urban Environments
- Application of omics-based strategies for improved diagnosis and treatment of endocrine hypertension
- Youth alcohol consumption database and mobile app
- LIDAR Data Analytics Research Environment
- Type-1 Diabetes Clinical Research Network
- American Asian Australian Adrenal Alliance
- International League Against Epilepsy
- Platform for Research Software Solutions (PRESS)
- Mobile applications for the Environmental Determinants of Islet Autoimmunity
- Secure Data Solutions for the Biomedical Communities of the Cloud
- Metabolomics Sample Management and Processing Platform
- Linked Data PolicyHub Stage II: Urban & Regional Planning & Communications
- Australian Genomics Health Alliance
- Melbourne Genomics Health Alliance
- Australian Diabetes Data Network
- Helicopter advanced training system, Australian Department of Defence
- Hort-eye Cloud analytics
- Public Records Office Victoria Data Management Solutions
- Complex System Modelling Platform and GPU utilisation
- Public Records Office Victoria Data Management Solutions Follow-Up Grant
- VicHealth 2016 Indicators API
- Helicopter advanced training system Phase II, Australian Department of Defence
- Twitter data analytics for business
- Mobile Applications for Patients with Neuroendocrine Tumours
- Systems Genomics Support Platform
- SWARM: Smartly-aggregated Wiki-style ARgument Marshalling (SWARM)
- ORCA Cognitive Assessment Platform
- 88days Backpacker app
- VicSpin Victoria-wide Flu Surveillance System
- ElectraNetLIDAR/VectorNZ Lidar
- Growing Landscape Carbon
- Replicats
- Bushfire data management platform

Compute Scaling

Microprocessor Transistor Counts 1971-2011 & Moore's Law



Future of Supercomputing
Moore's Law for Supercomputers



Network Scaling

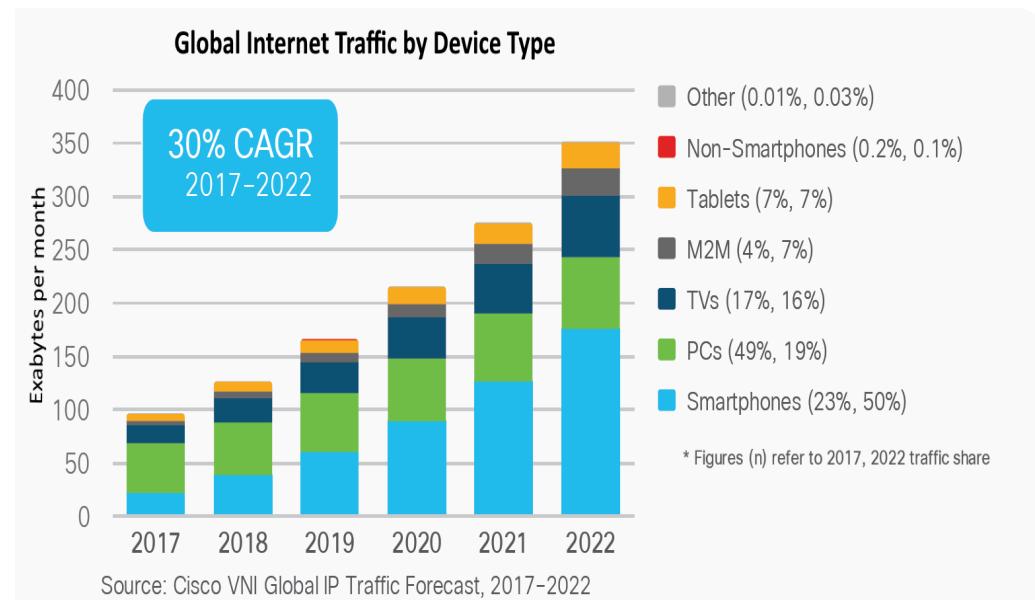
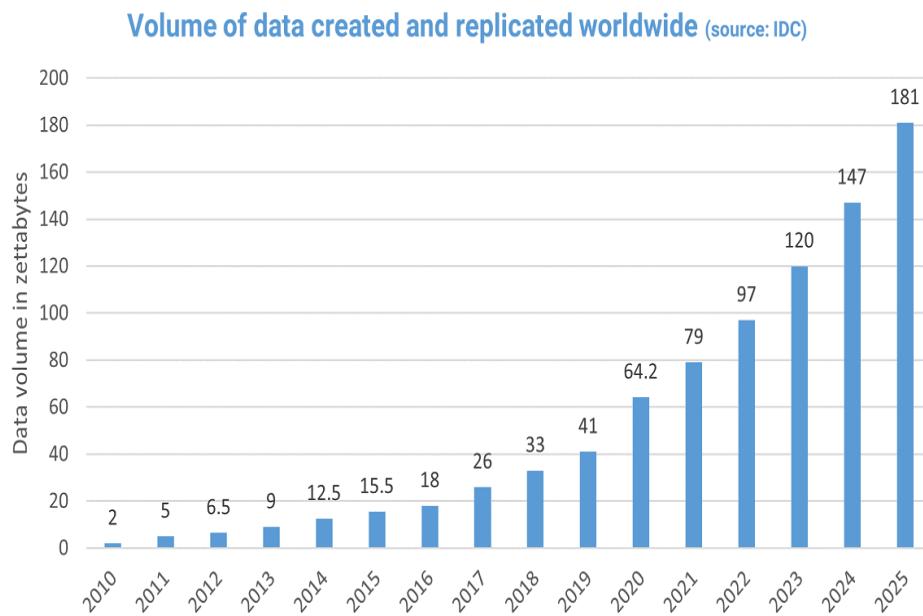


Table 1. The VNI Forecast Within Historical Context

Year	Global Internet Traffic
1992	100 Gigabytes per Day
1997	100 Gigabytes per Hour
2002	100 Gigabytes per Second
2007	2,000 Gigabytes per Second
2012	12,000 Gigabytes per Second
2017	35,000 Gigabytes per Second

Data Past

- From tablets, to papyrus, to books
 - (quite adequate for several thousand years)
 - Enter silicon transistors circa 1960
 - punch cards,
 - punched streamer tape,
 - magnetic tape,
 - floppies,
 - ...
- ~RIP!



Data Present

- Data Storage today

- CDs,
- DVDs,
- local (computer) hard disks,
- shared storage,
- tape storage.
- mobile storage,
- The Internet!

- Dropbox
- Google
- Clouds
- ...



Google



Data Scaling::The Deluge

The TechCrunch homepage features a prominent advertisement for the Dell Latitude 10 tablet. The ad includes the text "Built for business: the Dell Latitude™ 10 tablet." and "Let's Get Started >". Below the ad, there is a navigation bar with links to various tech topics like YAHOO, APPLE, FACEBOOK, TWITTER, GOOGLE, MICROSOFT, and NSA.

A banner for TechCrunch Disrupt SF 2013. It says "Limited Extra Early Bird Tickets Available Until July 15 Last Day!" and has a "GET TICKETS NOW" button.

Eric Schmidt: Every 2 Days We Create As Much Information As We Did Up To 2003

MG SIEGLER

Wednesday, August 4th, 2010

2 Comments



Naturally, all of this information helps Google. But he cautioned that just because companies like his can do all sorts of things with this information, the more pressing question now is if they *should*. Schmidt noted that while technology is neutral, he doesn't believe people are ready for what's coming.

"I spend most of my time assuming the world is not ready for the technology revolution that will be happening to them soon," Schmidt said.

The ScienceDaily homepage features a large logo with "Science" in blue and "Daily" in red. Below the logo, it says "Your source for the latest research news". There are tabs for News, Articles, Videos, Images, and Books. A sidebar mentions "Authorised by the Australian Government".

Science News

Big Data, for Better or Worse: 90% of World's Data Generated Over Last Two Years

May 22, 2013 — A full 90% of all the data in the world has been generated over the last two years. The internet companies are awash with data that can be grouped and utilised. Is this a good thing?

Share This:

Like 221

Tweet 209

Share 23

Share 59

An increasing amount of data is becoming available on the internet. Each and every one of us is constantly producing and releasing data about ourselves. We do this either by moving around passively -- our behaviour being registered by cameras or card usage -- or by logging onto our PCs and surfing the net.

The volumes of data make up what has been designated 'Big Data' -- where data about individuals, groups and periods of time are combined into bigger groups or longer periods of time.

Research advantages

Petter Bae Brandzaeg of SINTEF ICT points to the huge research centres now developed at internet companies such as Facebook and Google.

'The advantages they have is the enormous volume of data that other social researchers can only dream of,' he says. However, it has also changed the way SINTEF researchers work. Even those not working in the major internet companies can still access Big Data.

Brandzaeg has investigated a tool called Wisdom developed by the American-based company MicroStrategy, and has started applying it in the deITA-project which addresses young people's social activity on the internet.

'This gives me access to data about over 20 million people -- without making a single inquiry. I can analyse different preferences on Facebook and look at age and gender differences between various groups and nations across the world. So far I have compared gender differences in social activity on Facebook between people in Norway, Spain, England, USA, Russia, Egypt, India and China.'

Data protection is a problem we often associate with Big Data, but according to Brandzaeg, data from Wisdom is restricted to large groups and does not go down to 'individual level'. This makes it possible for him to compare large groups without any data protection problems.



Big Data makes it possible to achieve research results that cover a wide range of issues, and can tell us a great deal about developments in the world in many different areas. (Credit: © zothen / Fotolia)

Related Topics

Computers & Math

- ▶ Computers and Internet
- ▶ Information Technology
- ▶ Encryption

Articles

- ▶ List of data structures
- ▶ Computing
- ▶ Streaming media

Science & Society

- ▶ Security and Defense
- ▶ Privacy Issues
- ▶ Popular Culture

Articles

- ▶ Identity theft
- ▶ Data mining
- ▶ Wi-Fi

Corporate Compliance

www.gfi.com

Issues? Let GFI EventsManager help your compliance. Download trial!



Big Data Architecture

www.teradata.com/BigData

Extract insights from big data to create new revenue opportunities.

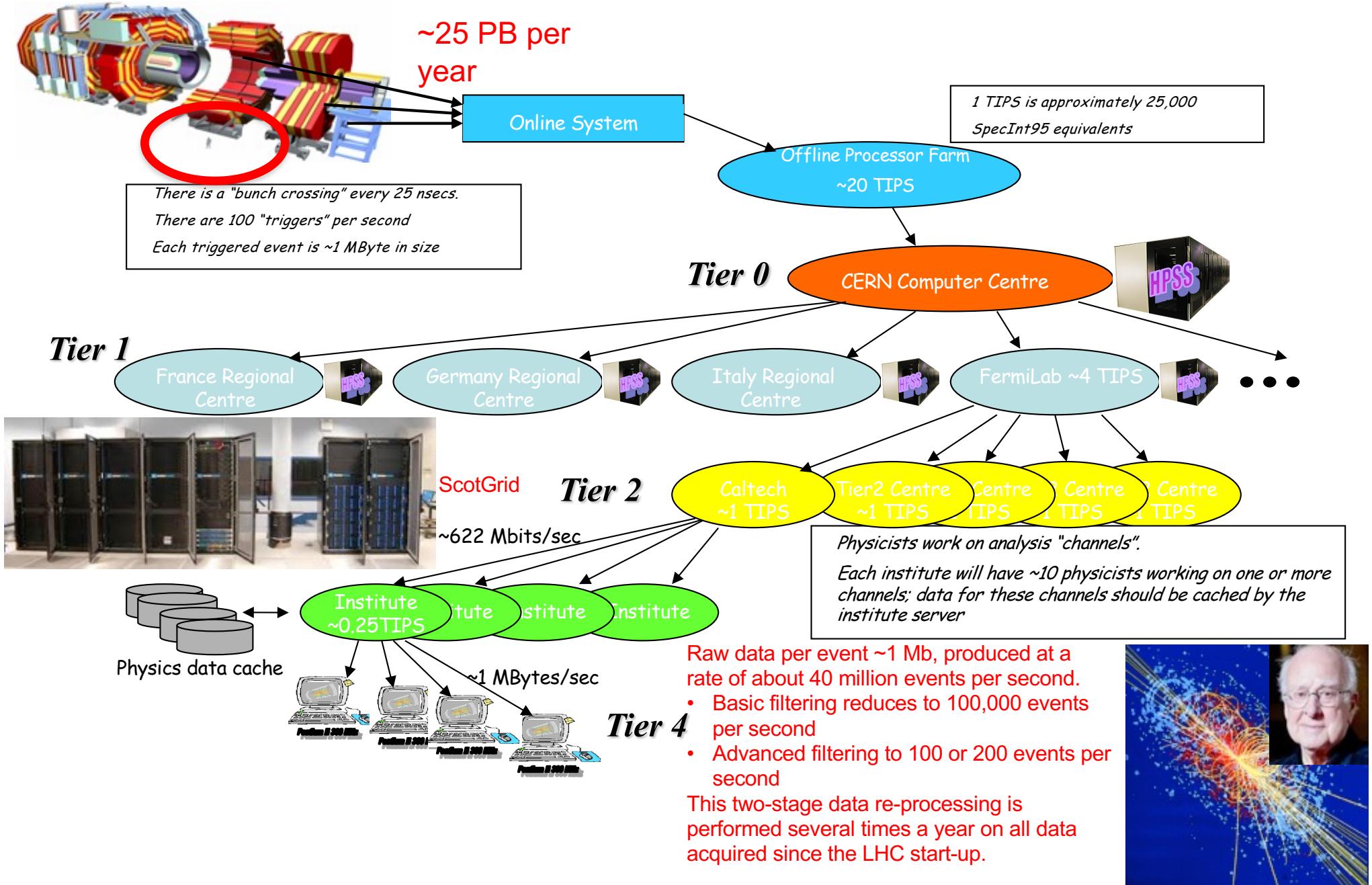


Data Intensive / Data driven Research

- Researchers need tools, methodologies
 - To search for/discover data
 - To use/analyse data
 - To share data
 - To store data
 - To track data
 - To destroy data
 - To move data around
 - To check authenticity of data
 - To visualise data
 - To overcome issues of data heterogeneity
 - ...

... and this should be tailored to the researchers needs!!!

Compute Infrastructure for High Energy Physics



Mapping the Skies



"Chipsets needed to process data access and SKA applications will need to be capable of 20-25 exaflops of processing power", according to IBM Research's Ton Engbersen, DOME scientist and project leader. "Take the current global daily Internet traffic, double it, and you are in the range of the data set that the SKA will collect each day." This would equate to around 40,000Pb every 24 hours.

Meeting the Design Challenges of Nano-CMOS Electronics

e-Science Pilot Project (EPSRC)

R. Sinnott (e-Science Director)

Resources

£3.7M EPSRC; £4.4M FEC
£5.8M incl. industrial contributions

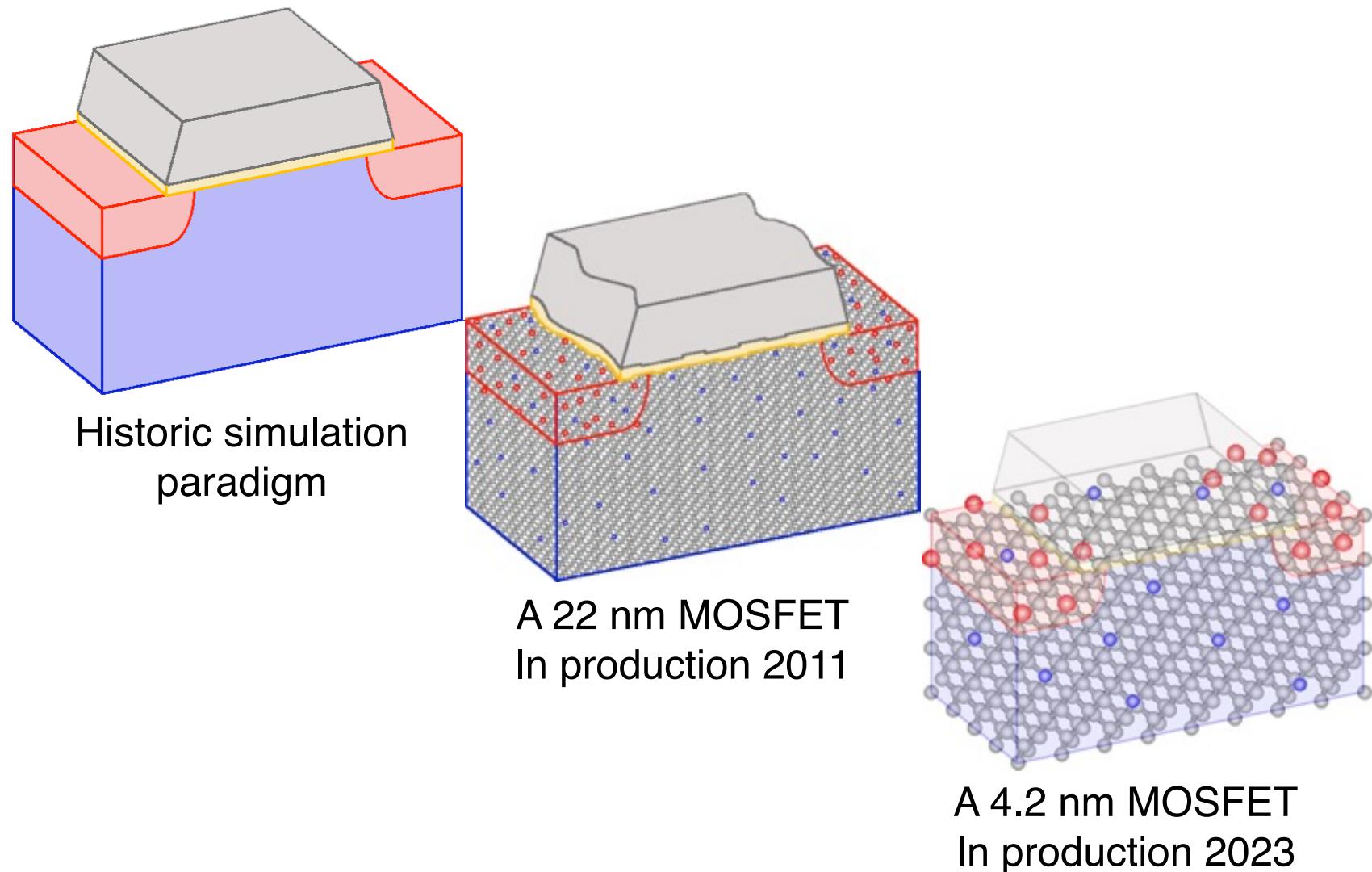
11 PDRAs (7 + 4)
7 PhD

Industrial partners

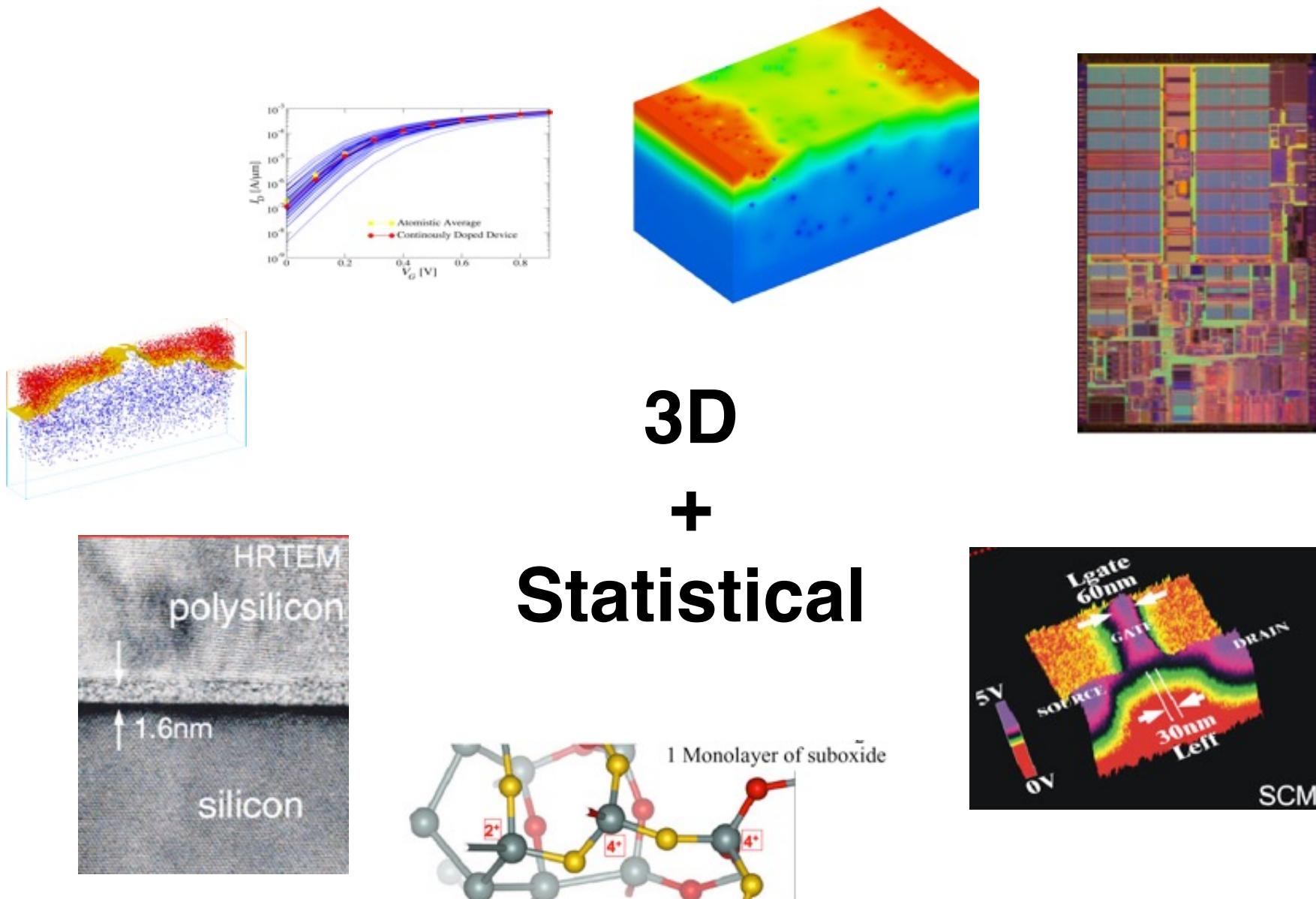


University partners

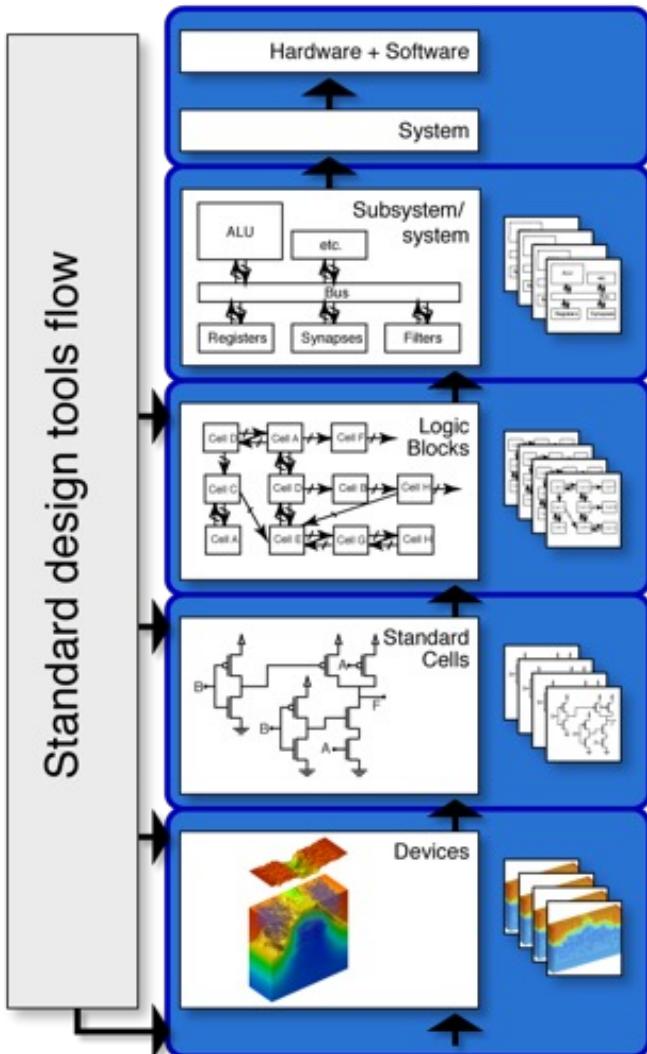
Meeting the Design Challenges of Nano-CMOS Electronics



Challenges of NanoCMOS Design



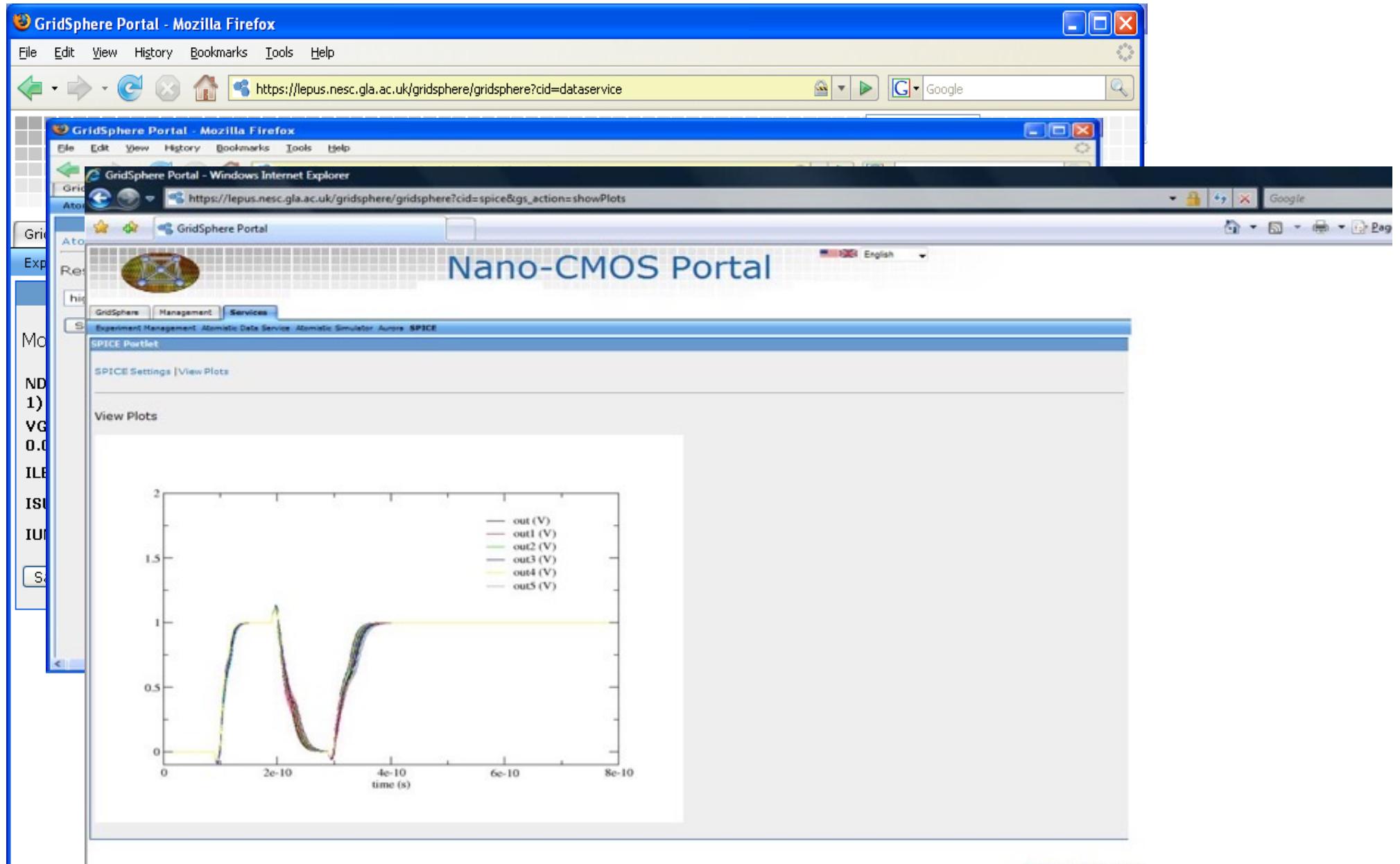
Challenges of Hierarchical statistical system simulations



- ❑ Very large device and circuit simulations
3D devices
 10^5 circuit components
- ❑ Large statistical samples
1000 - 100000 3D simulations - 4D
1000 - 100000 circuit simulations
- ❑ Complex flow and storage of data
Many files per simulation
Metadata capture and data provenance
- ❑ Collaboration between 5 partners
Multidisciplinary background
Complex data exchange
- ❑ Stringent security requirements
Commercial IP
Expensive software licenses

e-Xperiences

- We started with a secure portal and a wiki!!!



But ended up with...

... a command-line based solution

This community

Security solution

Secure, distribut

Meta-data captu

Job submission (massive (at the t

- *ScotGrid, NCSA, ...*
- *Millions of Computation*

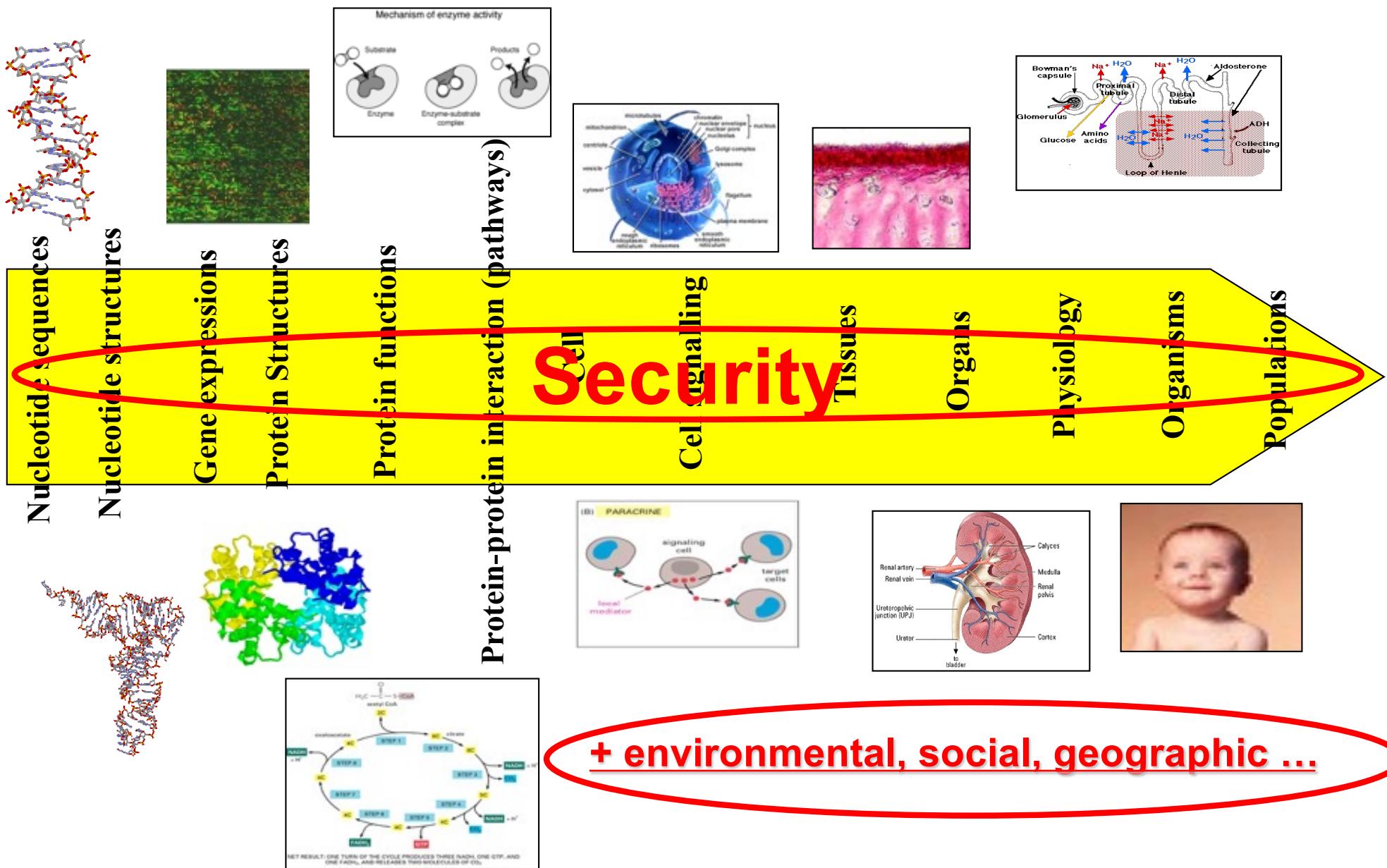
The -g flag!!!

The screenshot shows a Mozilla Firefox browser window with a title bar "DataService.swf (application/x-shockwave-flash Object) - Mozilla Firefox". The menu bar includes File, Edit, View, History, Bookmarks, Tools, and Help. On the left, a sidebar lists "File Record Details", "Core Data", "Name", "Value", "recordtype", "FileRecord", "location", and "afs/nesci.scc.ac.uk/nanocmostest/job/items/95ed550d-3bfc-42c8-84c6-6e3edc91fe03". The main content area displays a large 3D watermark reading "Data vs Metadata". Below the watermark, there is a code editor window titled "metadata/802" showing XML code:

```
<geronimo>
<metadata>
<created-time value="14:48.31" />
<created-date value="05.11.2008" />
<execution-host value="node021.cvos.cluster" />
<subjob-partition value="3" />
<job-uri value="https://nanodata.vidar.ngs.manchester.ac.uk/nanocmostest/job/items/95ed550d-3bfc-42c8-84c6-6e3edc91fe03/" />
<job-id value="95ed550d-3bfc-42c8-84c6-6e3edc91fe03" />
<experiment-uri value="https://nanodata.vidar.ngs.manchester.ac.uk/nanocmostest/experiment/items/00000000-0000-0000-0000-000000000000" />
<experiment-id value="00000000-0000-0000-0000-000000000000" />
<task-uri value="https://nanodata.vidar.ngs.manchester.ac.uk/nanocmostest/task/items/16d68fffc-1027-45fa-a7c9-8627d3695a75/" />
```

At the bottom, there are buttons for "Done" and the URL "nanodata.vidar.ngs.manchester.ac.uk" with a lock icon.

The e-Health Future...



Life Sciences

- Extensive Research Community
 - Parkville Precinct for example
- Many people care about them
 - Health, Food, Environment – truly interdisciplinary!
- Interacts with virtually every discipline
 - Physics, Chemistry, Maths/Stats, Nano-engineering, ...
- Thousands of databases relevant to bioinformatics (and growing!)
 - Heterogeneity, Interdependence, Complexity, Change, ...
- Some of the Big Questions/Challenges
 - How does a cell work?
 - How does a brain work?
 - How does an organism develop?
 - Why do people who eat less tend to live longer?
 - ...

More (and more and more) genomes...



Distributed, completely heterogeneous data

The image is a composite of three distinct sections. The top section shows a computer monitor displaying a terminal window with a massive amount of JSON data, illustrating the scale and complexity of the distributed data. The bottom section features a woman holding a newborn baby, with a small device attached to the baby's chest, likely a glucose monitor, symbolizing the physical reality of the data being analyzed. The middle section is a screenshot of the ENDIA study website. It features a logo of a stylized flower or leaf inside a circle. The page has a light blue header with the text "environmental determinants of islet autoimmunity". Below the header are two buttons: "Participant Login" and "Staff Login". The main content area includes navigation links for "Home", "About Us", "What's Involved", "Regional Program", "Health Professionals", "Researcher Resources", "News", and "Contact". A prominent blue text overlay on the left side of the website asks, "Why are more children getting Type 1 Diabetes? Please note that recruitment has closed". At the bottom left, there is a blue button labeled "Contact the Team". The overall theme is the connection between the raw data and its real-world application in medical research.

Messy

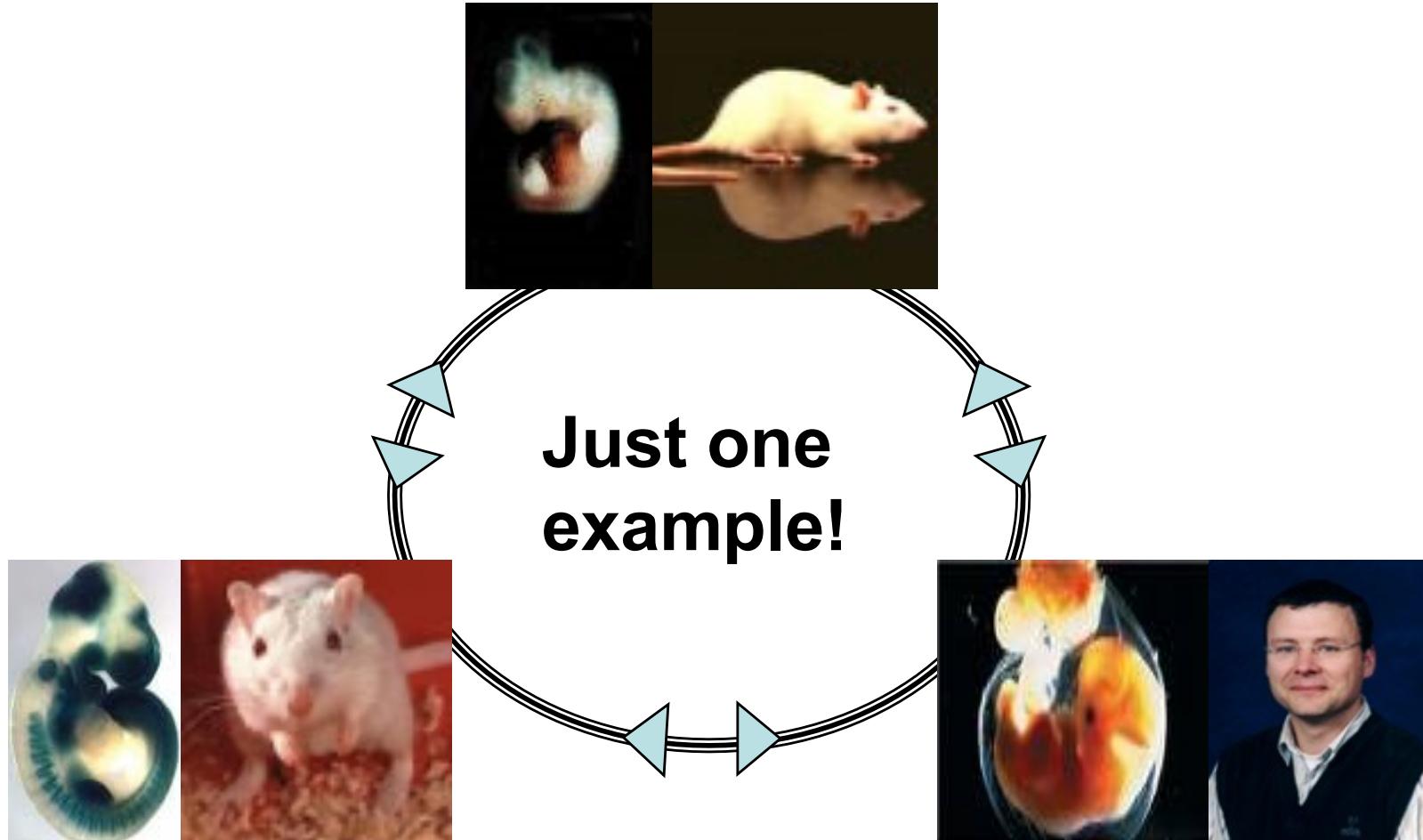
GPL96 - Notepad

File Edit Format View Help

%Annotation!Annotation_date = 09/20/2006 15:35:01!Annotation_platform = GPL96!Annotation_platform_title = Affymetrix GeneChip Human Genome U
n replication factor C, 40-kDa subunit (A1) mRNA, complete cds 1590810 M87338 RFC2 0117_at heat shock 70kDa protein 6 (HSP
5 01431_at cytochrome P450, family 2, subfamily E, polypeptide 1 Human cytochrome P450IE1 (ethanol-indu
an farnesyl-protein transferase beta-subunit mRNA, complete cds 292032 L00635 FNTB 0177_at phospholipase D1, phosphatidylc
86095 NM_000991 RPL28 0200004_at eukaryotic translation initiation factor 4 gamma, 2 Homo sapiens eu
0200010_at ribosomal protein L11 Homo sapiens ribosomal protein L11 (RPL11), mRNA 15431289 NM_0009
_002954 RPS27A 0200018_at ribosomal protein S13 Homo sapiens ribosomal protein S13 (RPS13), mRNA 1459191
55 (RPS5), mRNA 71164878 NM_001009 RPS5 0200025_s_at ribosomal protein L27 Homo sapiens ribosomal
(RPS11), mRNA 34335149 NM_001015 RPS11 0200032_s_at ribosomal protein L9 Homo sapiens ribosomal
200039_s_at proteasome (prosome, macropain) subunit, beta type, 2 Homo sapiens proteasome (prosome, macropain) subunit, beta type
member 1 Homo sapiens ATP-binding cassette, sub-family F (GCN20), member 1 (ABCF1), transcript variant
ing factor 2, 45kDa (ILF2), mRNA 24234746 NM_004515 ILF2 0200053_at CD40 ligand, membrane
SCC3L1 0200059_s_at ras homolog gene family, member A Homo sapiens Ras homolog gene family, member A, mRNA (c
2288 AF275719 HSP90AB1 0200065_s_at ADP-ribosyl cyclase 1, mRNA 3126877 AF061832
eterogeneous nuclear ribonucleoprotein M Homo sapiens Macrocistin, elavl-like, mRNA 3126877 AF061832
RNA (cDNA clone MGC:4498 IMAGE:2964510), complete cds 33873259 0200079_s_at lysyl-tRNA synt
014267 C1orf58 0200085_s_at Homo sapiens Macrocistin, elavl-like, mRNA 3126877 AF061832
rase, CAAX box, alpha (FNTA), transcript variant 1 Homo sapiens Macrocistin, elavl-like (FNTA), transcript variant 3, mRNA
omo sapiens ATPase, H⁺ transporting, mitochondrial membrane NM_003945 ATP6VOE
U 0200594_x_at Homo sapiens heat shock transcription factor A (HSF1), transcript variant
, mRNA 4507676 NM_003299 HSP90AA1 0200600_s_at heat shock transcription factor A (Grp94), member 1 Homo sapiens he
ein kinase, cAMP-dependent, regulat 0200601_s_at heat shock transcription factor A (Grp94), member 1, mRNA 47132579//4713
006265 RAD21 0200602_s_at WD repeat domain 1 (WDR1), transcript variant 1, m
rotein complex 2, beta 1 subunit (AP2B1), transcript variant 1 Homo sapiens WD repeat domain 1 (WDR1), transcript variant 1, m
n 3 (phosphorylase kinase, delta) Homo sapiens WD repeat domain 1 (WDR1), transcript variant 1, mRNA 58218967 NM_0051
mRNA 47419913 NM_004184 WARS 0200603_s_at SET translocation (myeloid leukemia-associated) Homo sa
Homo sapiens protein tyrosine phosphatase, receptor type, F (PTPRF), transcript variant 1, mRNA 109633040 NM_002840 PTPRF
yptophan 5-monoxygenase activation protein, zeta polypeptide Homo sapiens tyrosine 3-monoxygenase/tryptophan 5-monoxygenase activa
0646_s_at nucleobindin 1 Homo sapiens nucleobindin 1 (NUCB1), mRNA 39725676 NM_006184 NUCB1
d protein beta) Homo sapiens signal sequence receptor, beta (translocon-associated protein beta) (SSR2), mRNA 6552341 NM_003145
ide translocator), member 5 (SLC25A5), mRNA 4502098 NM_001152 SLC25A5 0200658_s_at prohibitin Homo sa
NM_006145 DNAJB1 0200665_s_at secreted protein, acidic, cysteine-rich (osteoneectin) Homo sapiens secreted p
g protein 1 (XBP1), mRNA 14110394 NM_005080 XBP1 0200671_s_at spectrin, beta, non-erythrocytic 1 Homo sapiens secreted p
7977 NM_003347 UBE2L3 0200677_s_at pituitary tumor-transforming 1 interacting protein Homo sapiens pi
nt 1, mRNA//Homo sapiens ubiquitin-conjugating enzyme E2L 3 (UBE2L3), transcript variant 2, mRNA 38157977//38157975 NM_003347//NM_
on factor 1 gamma Homo sapiens eukaryotic translation elongation factor 1 gamma (EEF1G), mRNA 83656774 NM_001404
D (Asp-Glu-Ala-Asp) box polypeptide 24 (DDX24), mRNA 14251213 NM_020414 DDX24 0200695_s_at protein phospa
mRNA 8051609 NM_006854 KDELR2 0200700_s_at KDEL (Lys-Asp-Glu-Leu) endoplasmic reticulum protein retention receptor
6127 NM_001959 EEF1B2 0200706_s_at lipopolysaccharide-induced TNF factor Homo sapiens lipopolysaccharide
bule-associated protein, RP/EB family, member 1 Homo sapiens microtubule-associated protein, RP/EB family, member 1 (MAPRE1), mRNA
phase kinase-associated protein 1A (p19A) (SKP1A), transcript variant 2, mRNA 25777710//25777712 NM_006930//NM_170679 SKP1A
t 2, mRNA 61676201//61676202 NM_005898//NM_203364 GPIAP1 0200723_s_at GPI-anchored membrane protein 1
apiens ARP2 actin-related protein 2 homolog (yeast) (ACTR2), transcript variant 1, mRNA//Homo sapiens ARP2 actin-related protein 2 homolog (ye
AX1 (hPTPCAA1) mRNA, complete cds 1777754 U48296 PTP4A1 0200734_s_at ADP-ribosylation factor 3 Homo sa
3 (S. cerevisiae) (SUMO3), mRNA 48928057 NM_006936 SUMO3 0200741_s_at ribosomal protein S27 (metallopanstimul
6 NM_002074 GNB1 0200747_s_at nuclear mitotic apparatus protein 1 Homo sapiens nuclear mitotic ap
e-serine-rich 2 Homo sapiens splicing factor, arginine-serine-rich 2 (SFRS2), mRNA 47271442 NM_003016 SFRS2
P-ribosylation-like factor 6 interacting protein 5 Homo sapiens ADP-ribosylation-like factor 6 interacting protein 5 (ARL6IP5), m
nce similarity 120A (FAM120A), mRNA 68299753 NM_014612 FAM120A 0200768_s_at methionine adenosyltransferase
rity 120A (FAM120A), mRNA 68299753 NM_014612 FAM120A 0200775_s_at heterogeneous nuclear ribonucleoprotein
491//NM_001008492//NM_004404//NM_006155 SEPT2 0200779_at activating transcription factor 4 (tax-responsive enhan
n-related protein 1 (alpha-2-macroglobulin receptor) Homo sapiens low density lipoprotein-related protein 1 (alpha-2-macroglobulin r
sapiens IQ motif containing GTPase activating protein 1 (IQGAP1), mRNA 57242794 NM_003870 IQGAP1 0200792_at

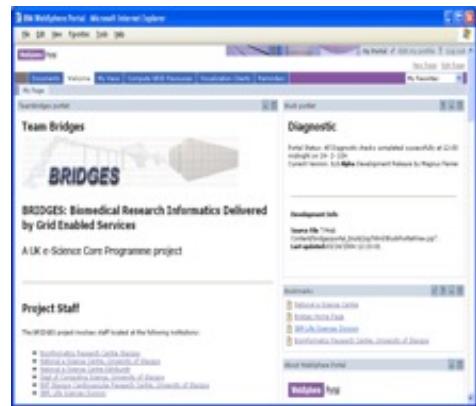
Next Generation Sequencers + 100 TB data

Translational Research

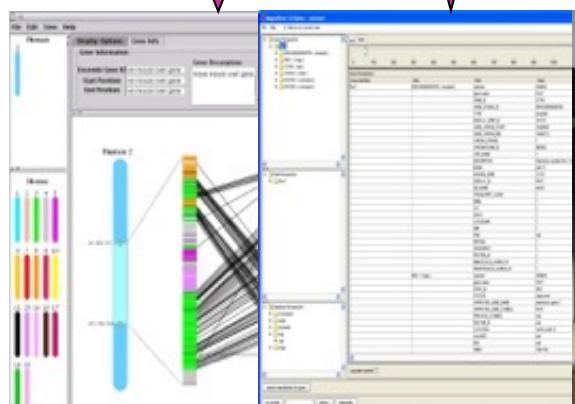


BRIDGES Project

VO Authorisation

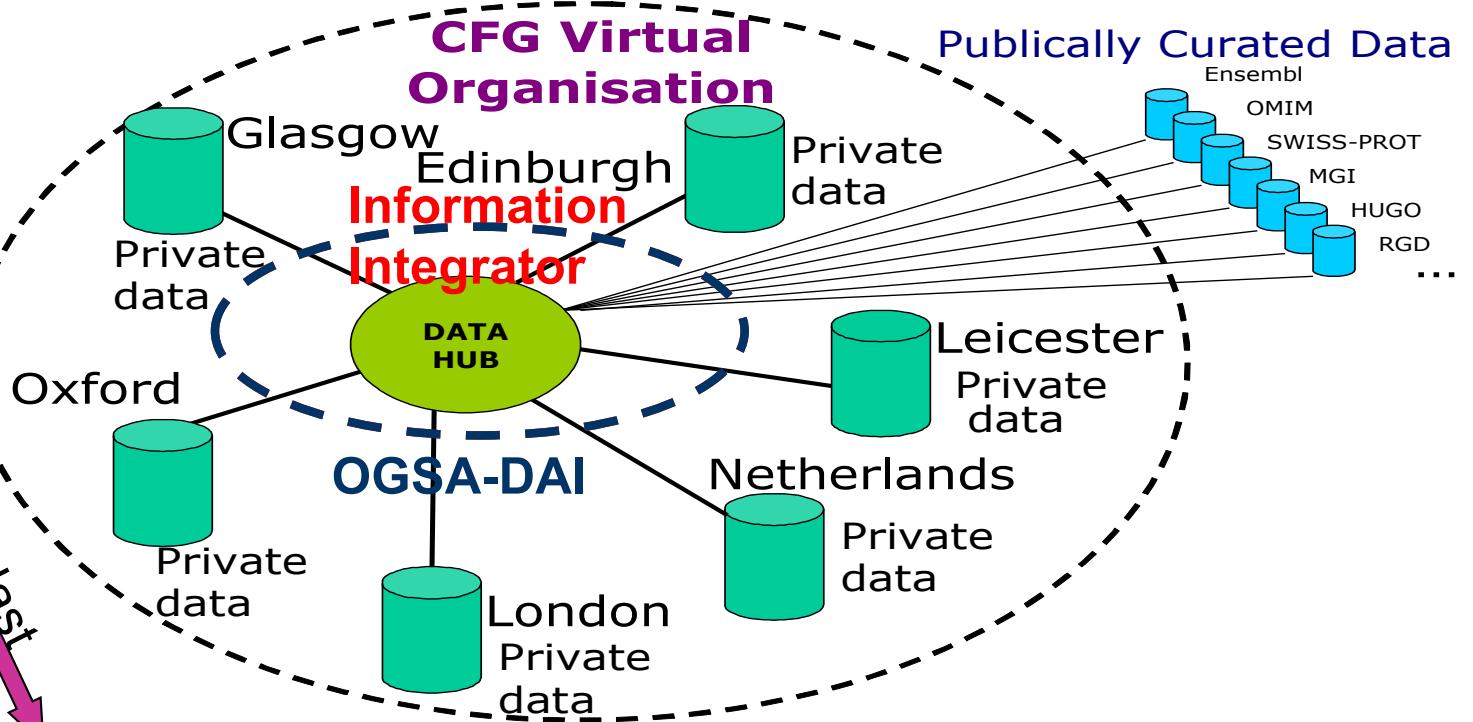


Synteny Service



Magna Vista Service

blast



+



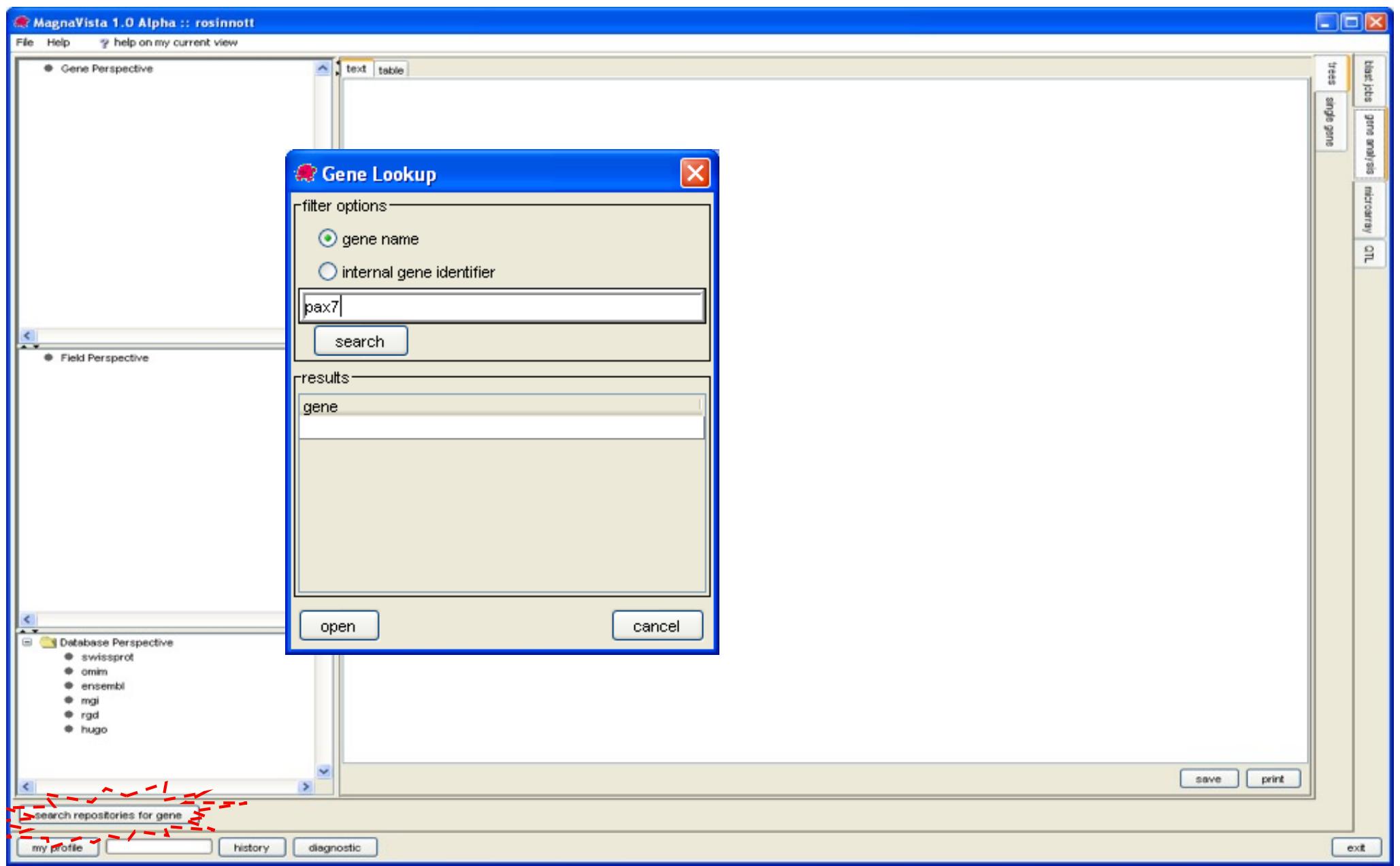
+



+



MagnaVista



MagnaVista

MagnaVista 1.0 Alpha :: rosinnott

File Help ? help on my current view

Gene Perspective
Pax7
ENSG00000009709 (ensembl)
8621 (hugo)
12185 (mgi)
167410 (omim)
P23759 (swissprot)
P47239 (swissprot)

Field Perspective
Pax7

Database Perspective
swissprot
omim
ensembl
mgi
rgd
hugo

Profile for user:rosinnott

Profile for user:rosinnott

Profile for user:rosinnott

My Profile Restore System Defaults

fields found for database: ensembl

Field	swissprot	omim	ensembl	mgi	rgd	hugo
species	<input type="checkbox"/>					
gene name	<input type="checkbox"/>					
GENE_ID	<input type="checkbox"/>					
GENE_STABLE_ID	<input type="checkbox"/>					
TYPE	<input type="checkbox"/>					
DISPLAY_XREF_ID	<input type="checkbox"/>					
GENE_CHROM_START	<input type="checkbox"/>					
GENE_CHROM_END	<input type="checkbox"/>					
CHROM_STRAND	<input type="checkbox"/>					
CHROMOSOME_ID	<input type="checkbox"/>					
CHR_NAME	<input type="checkbox"/>					
DESCRIPTION	<input type="checkbox"/>					
BAND	<input type="checkbox"/>					
KNOWN_GENE	<input type="checkbox"/>					
DISPLAY_ID	<input type="checkbox"/>					
DB_NAME	<input type="checkbox"/>					
TRANSCRIPT_COUNT	<input type="checkbox"/>					
EMBL	<input type="checkbox"/>					
GO	<input type="checkbox"/>					
HUGO	<input type="checkbox"/>					
LOCUSLINK	<input type="checkbox"/>					
MIM	<input type="checkbox"/>					
PDB	<input type="checkbox"/>					
REFSEQ	<input type="checkbox"/>					
SWISSPROT	<input type="checkbox"/>					
PROTEIN_ID	<input type="checkbox"/>					
MMUSCULUS_HOMOLOG	<input type="checkbox"/>					
RNORVEGICUS_HOMOLOG	<input type="checkbox"/>					

genes databases perspectives field cross references fields Probe Sets qtl preferences

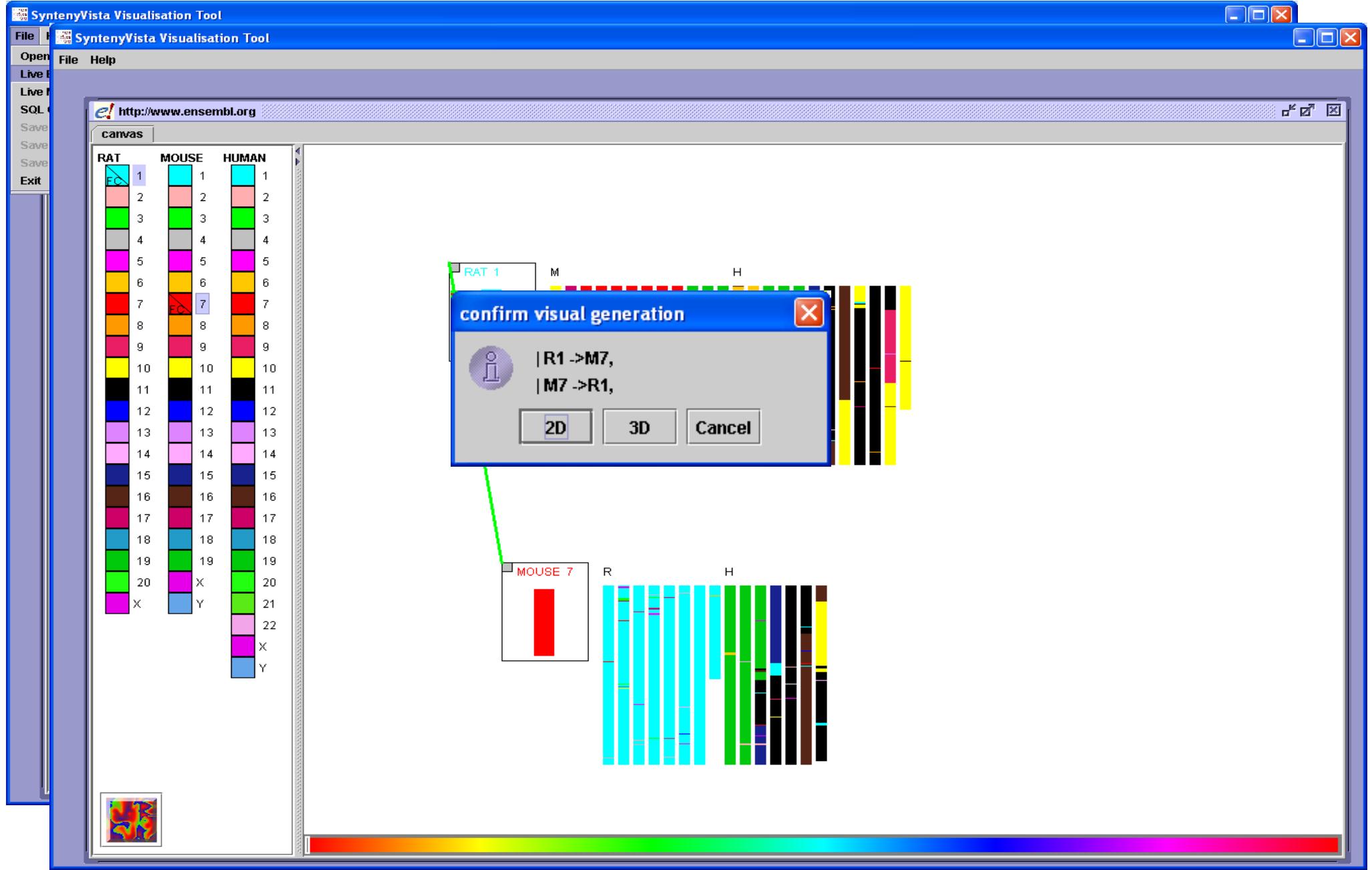
blast jobs gene analysis microarray QTL

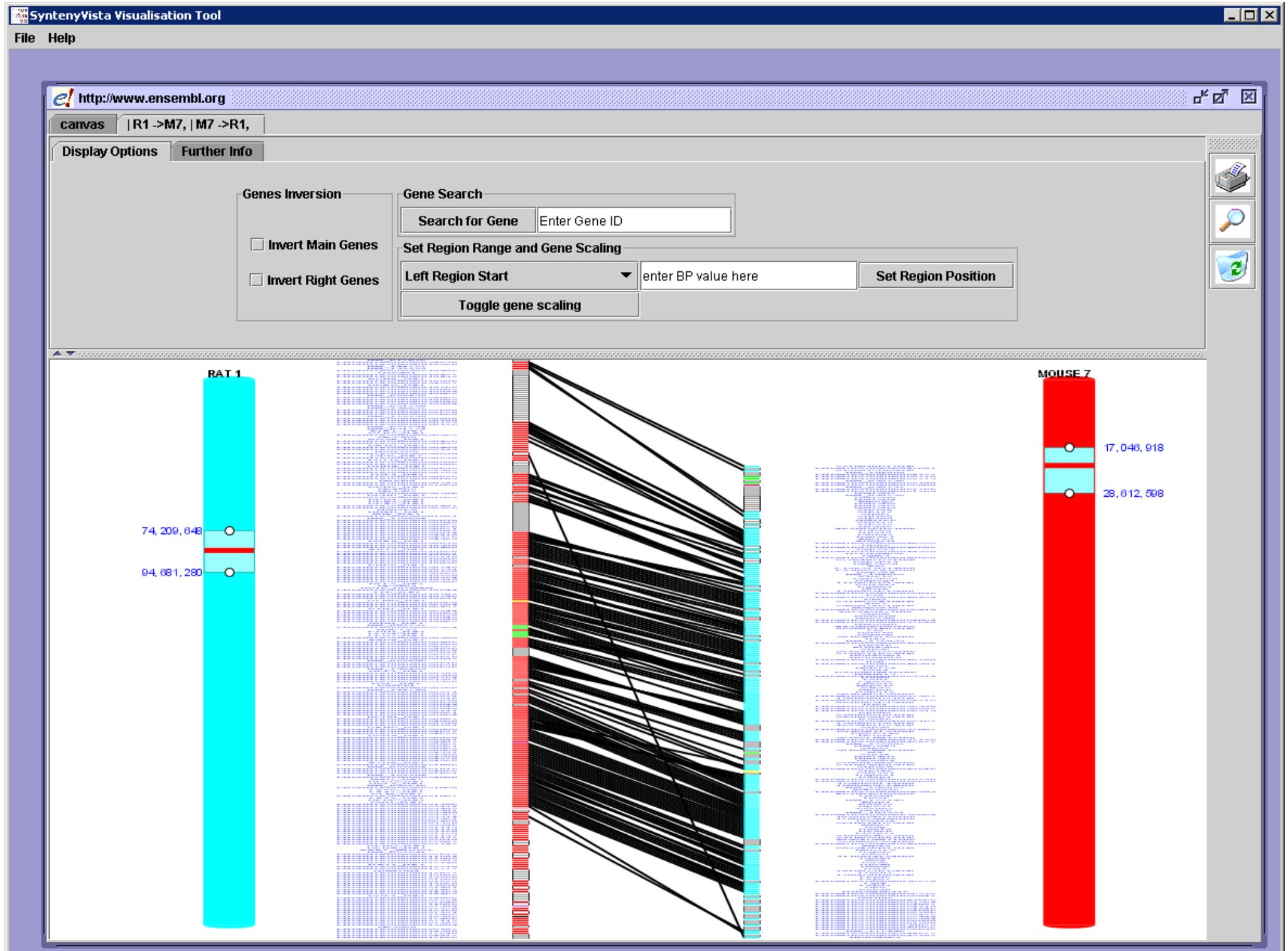
my profile history diagnostic

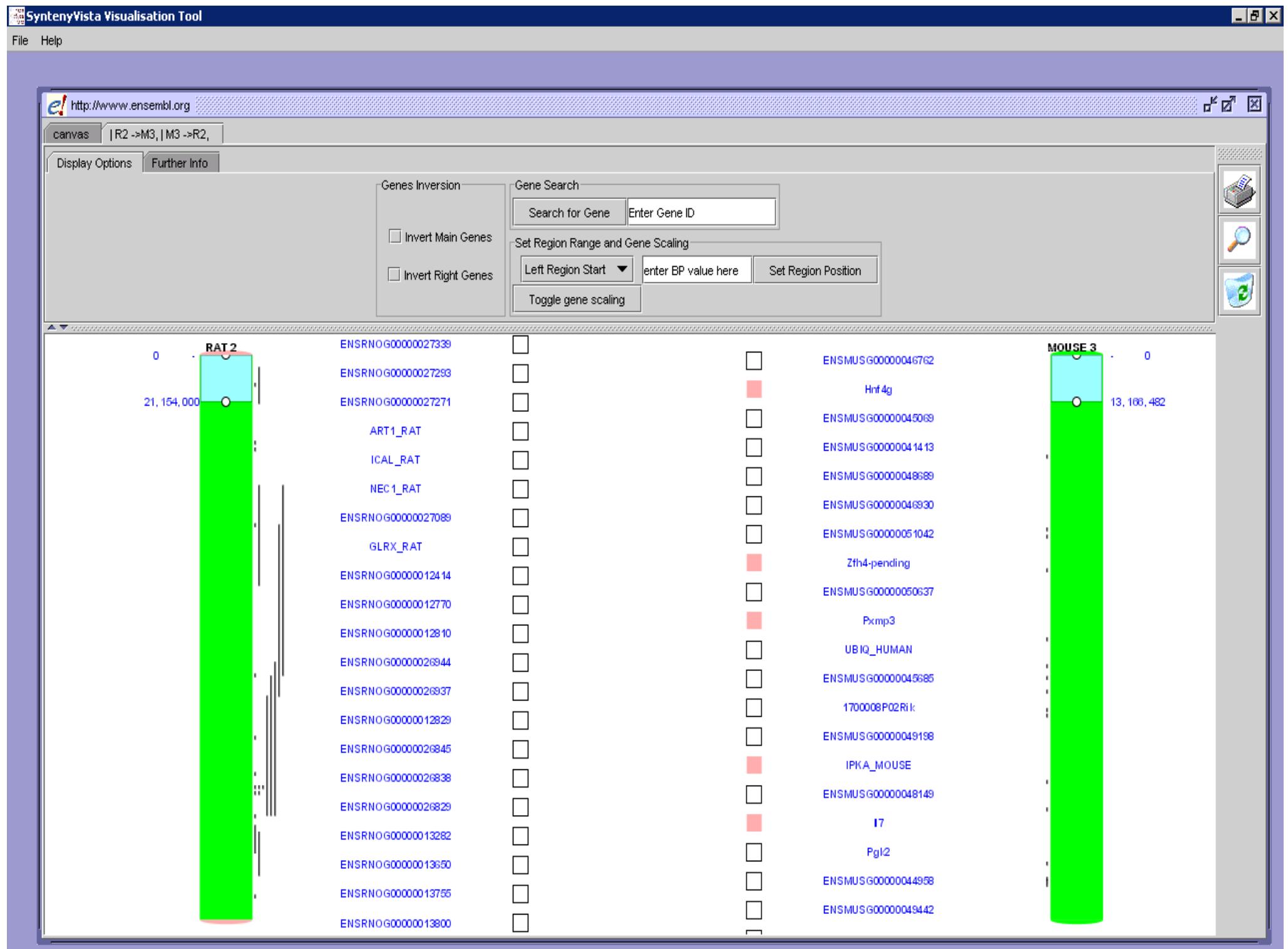
ok

exit

Importance of Data Visualisation







BREAK

Supporting Research in the Clinical/Biomedical Domain

- Different researchers/communities have different research computing needs
 - Based on my experiences
 - Data and data management
 - Inter-organisational collaborations
 - Security
 - Analytics, e.g. machine learning
 - Visualisation
 - *approximately in that order!*
 - Key to many research endeavours is Ethics
 - Governance, Trust, Policy, Procedures, GDPR, ...
 - All essential before any tech solutions applied!!!
 - Importance of relationships!!!
 - Trust, engagement, delivery, ...

A word about data...

- Clinical/biomedical data is...
 - almost always different across sites,
 - inconsistent across sites,
 - captured at different times/for different reasons,
 - hosted on different systems,
 - uses different databases/data technologies,
 - uses different coding schemes (or none!),
 - has different end users in mind,
 - has different security/ethics/governance arrangements...

Such challenges must be overcome!

Rapid prototyping, community data agreements
=> gradual hardening of secure systems
for use in clinical/biomedical research

Biomedical Data Sharing Models

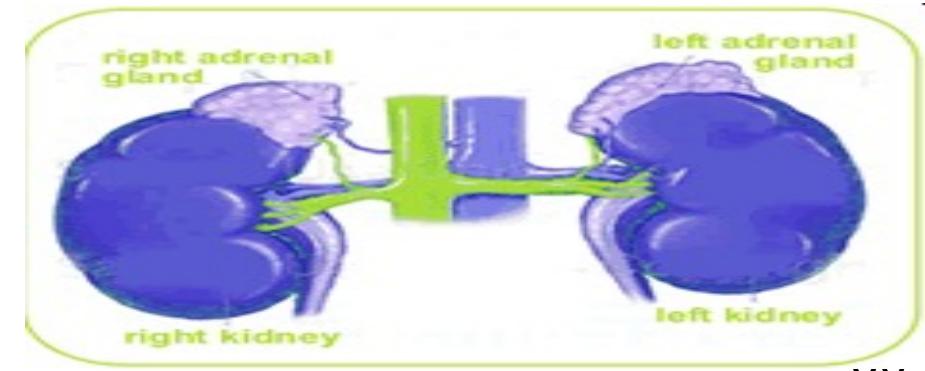
- Many models exist in the biomedical domain
 - There is no single solution that always applies
- Numerous models defined and supported
 - High-level architectural classification
 - Dependent, open-distributed, data sharing models;
 - direct distributed/federated querying (BRIDGES);
 - Independent centralised data sharing models;
 - Independent of existing hospital systems
 - Easiest to achieve (ENSAT-CANCER, ...)
 - Dependent push-based (out-going) data sharing models;
 - Re-use existing medical data
 - Achievable but sustainable? (ADDN, ...)
 - Dependent security-oriented pull-based data sharing models;
 - (asynchronous) data access
 - Dependent security-oriented push-based (incoming) clinical data sharing models
 - The major challenge (on many fronts, e.g. Optus/Medibank!!!)

Case Study Adrenal Cancer

(Independent centralised data sharing model)

- European Network for Study of Adrenal Tumours (ENSAT – www.ensat.org)

- 6m Euro; 5-year project;
 - Started 2011; officially ended 2016, but...
 - Funded one developer
 - Initial prototype in 6 weeks!
- Focus on 4 types of adrenal tumour
 - Adrenocortical carcinoma (ACC)
 - Aldosterone Producing Adenoma (APA)
 - Non-aldosterone cortical adrenal adenomas (NAPACA)
 - Pheochromocytomas and related paragangliomas (Pheo)
 - These are rare (ACC ~1-2 cases per million pop.)
- *Imaging...*
- *Biobanking...*
- *Clinical Trials...*



(Clinical) Data Aggregation

Aka how to improve data quality...

	ACC	Pheo	NAPACA	APA	Total
Records	4638	4986	7949	3983	22046
Biosamples	8562	6108	11009	13270	38949
Clinical Annotations	58760	27367	22169	13930	122226
Annotations Per Patient (Mean)	12.66	5.48	2.78	3.49	5.54
Biosamples Per Patient (Mean)	1.84	1.22	1.38	3.33	1.76
Active Centers	72	82	62	41	121

Associated Study and Registry Distribution

	ACC	Pheo	NAPACA	APA	Total	Principal Investigator	Study Protocols	Study sites/centers
EURINE-ACT	569	8	2130	237	2945	Wiebke Arlt	EURINE-ACT	
Ki-67	51	0	1	0	52	Martin Fassnacht		
Stage III/IV ACC	41	0	0	0	41	Eric Baudin		
PMT	0	259	162	2	423	Graeme Eisenhofer	PMT	PMT
TMA	8	188	22	0	218	Ronald Krijger		
Long-term PHPGL	0	282	0	0	282	Pierre-François Plouin		
AVIS-2	0	0	0	2	2			
PMT3	0	12	0	0	12	Graeme Eisenhofer		
ADIUVO	64	0	0	0	64	Massimo Terzolo		ADIUVO
ADIUVO Observational	111	0	0	0	111	Massimo Terzolo		
HairCo-2	9	0	1	0	10	Marcus Quinkler		
FIRST-MAPPP	0	53	0	0	53	Eric Baudin	FIRST-MAPPP	
German Cushing Registry	0	1	104	0	105	Martin Reincke		
German Conn Registry	0	1	0	1151	1154	Martin Reincke		German Conn Registry

Pheo PMT Study

	Full Study	Phase 1 (Screening)	Phase 2 (Clonidine)	Phase 3 (Pheo Characterization)	Phase 4 A (Excluded Follow-Up)	Phase 4 B (Pheo Follow-Up)
Patients	2238	454	17	6	1520	241

ENSAT-CANCER

ENSAT Registry

ENSAT Home | ACC Home | ACC Search | ACC Exported Data |

ACC Search Results

There are 60 records matching the following query:

Parameter	Condition	
year_of_birth	year_of_birth >= 1940 AND year_of_birth <= 1950	<input checked="" type="checkbox"/>
sex	M	<input checked="" type="checkbox"/>

[Repeat Search](#)

[Export these results](#)

[Export all your patient data](#)

[Run a new search](#)

ENSAT ID	Referral Doctor	Record Date	Date of First Registration	Sex	Year of Birth	Consent Level Obtained	
FRPA2-4	Jerome Bertherat (jerome.bertherat@cch.aphp.fr)	16 Apr 2007	16 Apr 2007	M	1944	National (France)	
FRPA2-51	Jerome Bertherat (jerome.bertherat@cch.aphp.fr)	09 Apr 2003	09 Apr 2003	M	1950	National (France)	
FRPA3-20	Eric Baudin (eric.baudin@igr.fr)	14 Dec 1989	14 Dec 1989	M	1946	National (France)	
FRPA3-24	Eric Baudin (eric.baudin@igr.fr)	03 May 2000	03 May 2000	M	1947	National (France)	
GYWU-671	Martin Fassnacht (fassnacht_m@medizin.uni-wuerzburg.de)	19 Apr 2011	26 Jan 2011	M	1940	National (Germany)	
GYWU-662	Martin Fassnacht (fassnacht_m@medizin.uni-wuerzburg.de)	19 Apr 2011	19 Jan 2011	M	1948	National (Germany)	
ITFL-3	Massimo Mannelli (m.mannelli@DFC.UNIFI.IT)	13 Dec 2010	14 Mar 2006	M	1947	Local	

Associated Record Search

[Surgery](#)

[Pathology](#)

[Biomaterial](#)

[Mitotane](#)

[Chemotherapy](#)

[Radiofrequency](#)

[Radiotherapy](#)

[Chemoembolisation](#)

[Follow-Up](#)

[ACC Home](#)

(Research) Impact

- Publications (100+)

- Nature Genetics (multiple)
- New England Journal of Medicine (multiple)
- Journal of Clinical Endocrinology and Metabolism
- Journal of Nuclear Medicine
- Journal of Endocrinology
- Journal of Human Molecular Genetics
- International Journal of Cancer
- Journal of Molecular and Cellular Endocrinology
- European Journal of Endocrinology
- Journal of Cancer Cell
- Journal of Clinical Cancer Research
- International Journal of Clinical Chemistry
- Journal of Hormone and Metabolic Research
- Journal of Modern Pathology
- Journal of Hypertension
- European Journal of Cancer
- Journal of Annals of Clinical Biochemistry
- Journal of Endocrine-related cancer
- Journal of the National Cancer Institute
- American Journal of Cancer Research
- ... Integrated genomic characterization of adrenocortical carcinoma

Guillaume Assié^{1-4,22}, Eric Letouzé^{5,22}, Martin Fassnacht⁶⁻⁸, Anne Jouinot¹⁻³, Windy Luscap¹⁻³, Olivia Barreau¹⁻⁴, Hanin Omeiri¹⁻³, Stéphanie Rodriguez¹⁻³, Karine Perlemoine¹⁻³, Fernande René-Corail¹⁻³, Nabila Elarouci¹⁵, Silviu Sbiera^{6,7}, Matthias Kroiss⁸, Bruno Allolio⁷, Jens Waldmann⁹, Marcus Quinkler¹⁰, Massimo Mannelli¹¹, Franco Mantero¹², Thomas Papathomas¹³, Ronald De Krijger¹³, Antoine Tabarin^{14,15}, Véronique Kerlan^{15,16}, Eric Baudin^{15,17}, Frédérique Tissier^{1-3,18}, Bertrand Dousset^{1-4,19}, Lionel Groussin¹⁻⁴, Laurence Amar²⁰, Eric Clauser²¹, Xavier Bertagna^{1-4,15}, Bruno Ragazzon¹⁻³, Felix Beuschlein⁶, Rossella Libé^{1-4,15}, Aurélien de Reyniès^{5,23} & Jérôme Bertherat^{1-4,15,23}

The NEW ENGLAND JOURNAL of MEDICINE

ORIGINAL ARTICLE

Constitutive Activation of PRKACA in Adrenal Cushing's Syndrome

Felix Beuschlein, M.D., Martin Fassnacht, M.D., Guillaume Assié, M.D., Ph.D., Davide Calebiro, M.D., Ph.D., Constantine A. Stratakis, M.D., D.Sc., Andrea Osswald, M.D., Cristina L. Ronchi, M.D., Ph.D., Thomas Wieland, M.Sc., Silviu Sbiera, Ph.D., Fabio R. Faucz, Ph.D., Katrin Schaak, Ph.D., Anett Schmittfull, M.S., Thomas Schwarzmayr, M.Sc., Olivia Barreau, M.D., Ph.D., Delphine Vezzosi, M.D., Ph.D., Marthe Rizk-Rabin, Ph.D., Ulrike Zabel, Ph.D., Eva Szarek, Ph.D., Paraskevi Salpea, Ph.D., Antonella Forlino, Ph.D., Annalisa Vetro, Ph.D., Orsetta Zuffardi, Ph.D., Caroline Kisker, Ph.D., Susanne Diener, M.Sc., Thomas Meitinger, M.D., Martin J. Lohse, M.D., Martin Reincke, M.D., Jérôme Bertherat, M.D., Ph.D., Tim M. Strom, M.D., and Bruno Allolio, M.D.

Somatic mutations in *ATP1A1* and *ATP2B3* lead to aldosterone-producing adenomas and secondary hypertension

Felix Beuschlein, Sheeraz Boulkroun, Andrea Osswald, Thomas Wieland, Hang N Nielsen, Urs D Lichtenauer, David Penton, Vivien R Schack, Laurence Amar, Evelyn Fischer, Anett Walther, Philipp Tauber, Thomas Schwarzmayr, Susanne Diener, Elisabeth Graf, Bruno Allolio, Benoit Samson-Couterie, Arndt Benecke, Marcus Quinkler, Francesco Fallo, Pierre-François Plouin, Franco Mantero, Thomas Meitinger, Paolo Mulatero, Xavier Jeunemaitre et al.

nature
genetics

Biomedical Data Sharing Models

- Many models exist in the biomedical domain
 - There is no single solution that always applies
- Numerous models defined and supported
 - High-level architectural classification
 - Dependent, open-distributed, data sharing models;
 - direct distributed/federated querying (BRIDGES);
 - Independent centralised data sharing models;
 - Independent of existing hospital systems
 - Easiest to achieve (ENSAT-CANCER, ...)
 - Dependent push-based (out-going) data sharing models;
 - Re-use existing medical data
 - Achievable but sustainable? (ADDN, ...)
 - Dependent security-oriented pull-based data sharing models;
 - (asynchronous) data access
 - Dependent security-oriented push-based (incoming) clinical data sharing models
 - The major challenge (on many fronts!)

Case Study: Australasian Diabetes Data Network (ADDN)

(Dependent push-based (out-going) data sharing model)

- National, longitudinal database for adult/child patients with T1D



- Used to
 - Increase patient centred research
 - Support Australasia-wide collaboration
 - Benchmarking across diabetes centres to drive quality improvement
 - Inform policy development for diabetes management

- Direct data entry not feasible/viable

- Currently
 - 12 adult sites, 23 pediatric sites
 - many other sites under discussion
 - Extremely heterogeneous systems



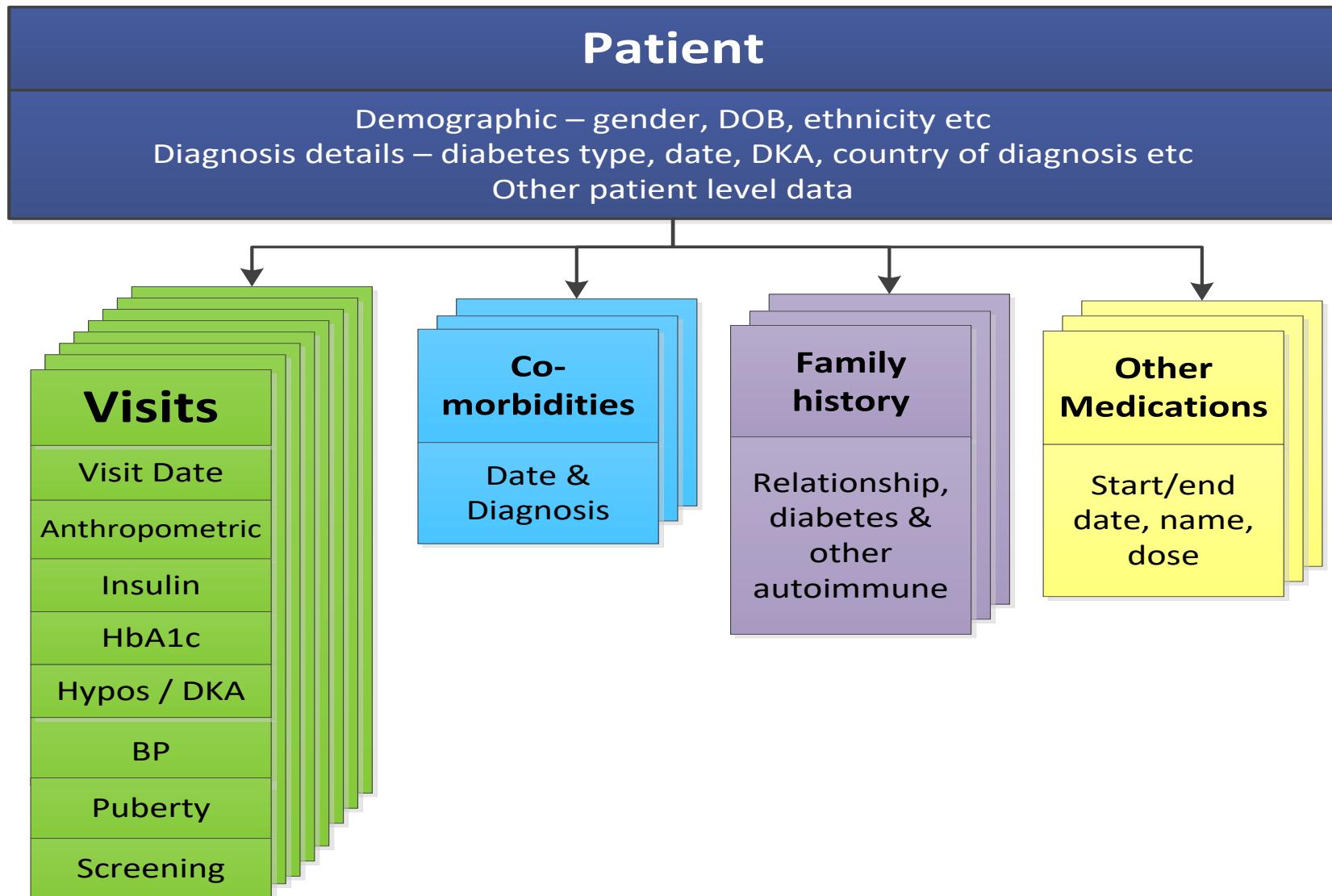
ADDN: Australasian Diabetes Data Network

There are more than 120,000 people living with Type 1 Diabetes (T1D) in Australia, with over half of them children or young people. The clinical management of T1D involves the replacement of insulin by injection or infusion and self-monitoring of blood glucose to balance food intake and exercise. The overall aim of treatment is to normalise blood glucose levels to prevent the complications that develop when blood glucose levels are high whilst reducing the major side effect of insulin therapy, low blood glucose or 'hypoglycaemia'.

Technological advances have enabled the establishment of clinical databases which capture health information about people with diabetes as they visit their diabetes clinic. The Australasian Diabetes Data Network (ADDN) has



ADDN Dataset



ADDN Data Quality

- Centres submit data twice yearly

The screenshot shows a web-based data cleaning tool for ADDN. The top navigation bar includes links for ADDN, DATA CLEANING (which is active), Dashboard, Import Files, Create Loads, Review Loads, Manage Errors, Error Codes, Download Data, a user profile icon, and Logout.

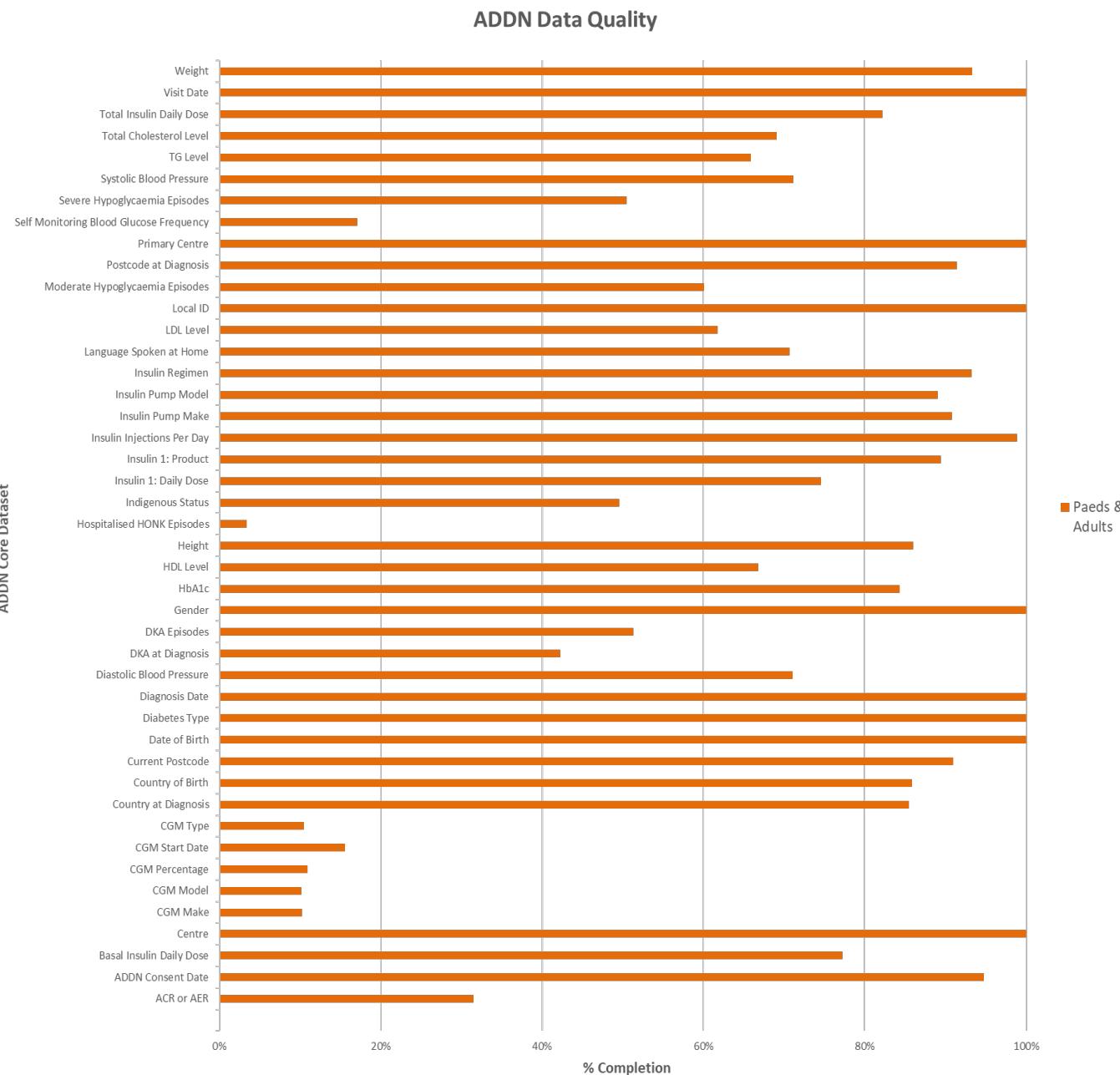
The main content area is titled "Loads for [REDACTED]" and displays a table of data quality metrics. The table has columns for Creation Date, Critical, Significant, Warning, Unprocessed, Status, and Actions.

Creation Date	Critical		Significant		Warning			Unprocessed	Status	Actions
	Tot.	Corr.	Tot.	Corr.	Tot.	Corr.	Ver.			
09/08/2018	1	1	0	0	1,516 -18	33	1,465 -10	18 -42	Under Review	
01/08/2018	1 +1	0	0	0	1,534 +160	0	1,475 +108	60 +53	Superseded	
14/03/2018	0	0	0	0	1,374 -8	0	1,367 +3	7 +1	Review Completed	
02/03/2018	0	0	0	0	1,382 +7	12	1,364 +20	6 -25	Superseded	
23/02/2018	0	0	0	0	1,375 +94	0	1,344 +69	31 +25	Superseded	
15/02/2018	0	0	0	0	1,281	0	1,275	6	Review Completed	

- *Clean* their data using on-line data cleaning tool
- Final dataset created – used for benchmarking reports, analysis and research data requests

Data Completeness

ADDN Core Dataset

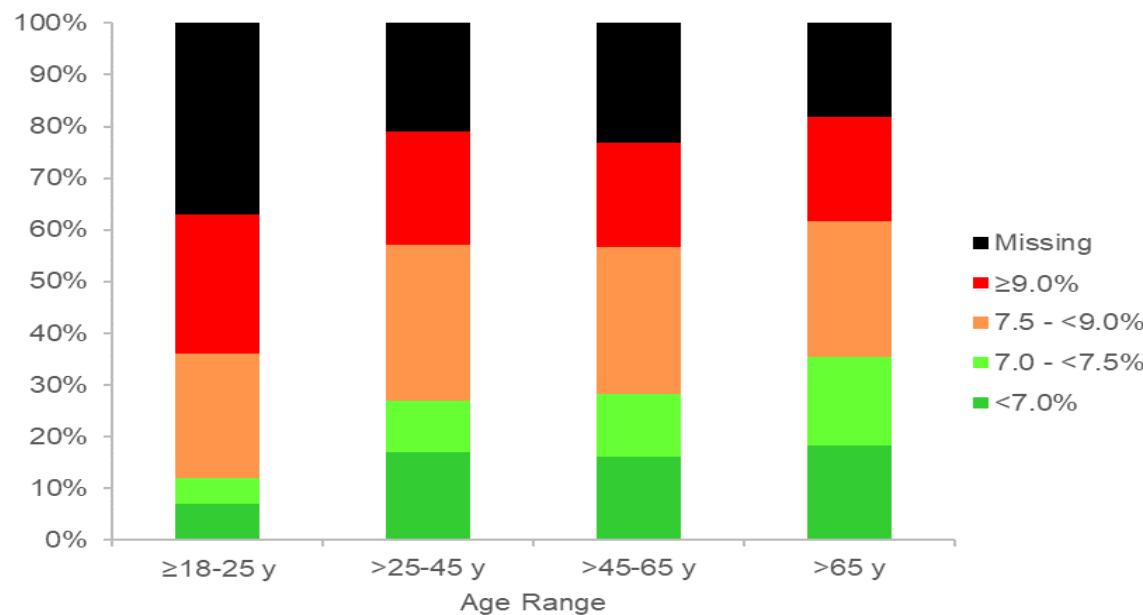


ADDN dataset =
131 data items

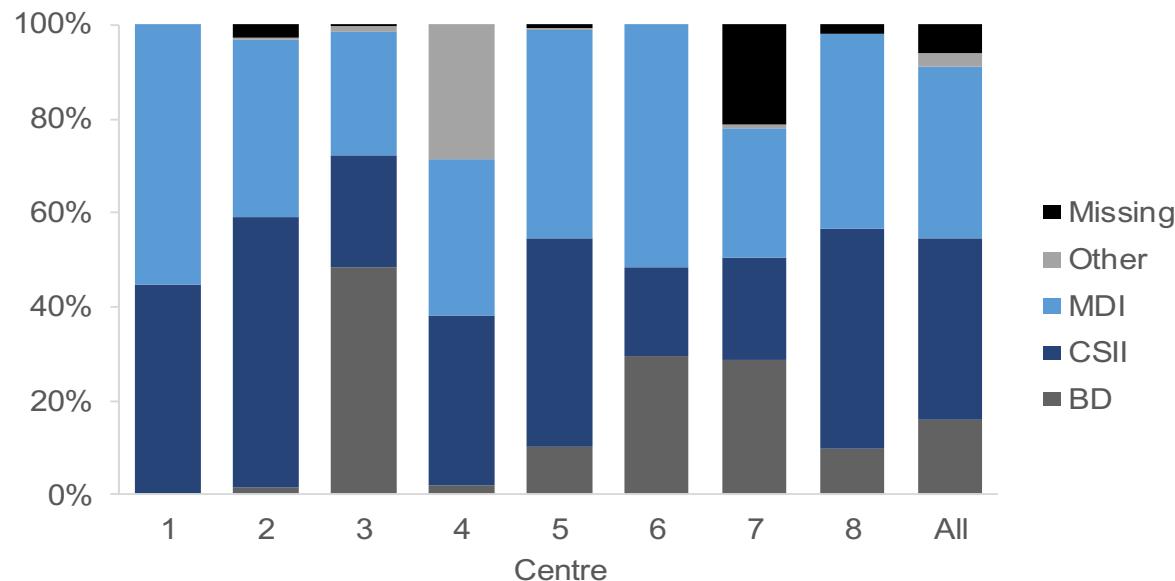
Core dataset =
44 data items

ADDN Example Analytics

Adult HbA1c
(average glycaemia over the preceding 6–8 weeks)



Paed Insulin Regimen
(Multiple Daily Injections, Continuous Subcutaneous Insulin Infusion, Bi-daily)



- Australian Urban Research Infrastructure Network (AURIN) federally funded project
 - ~\$70m funding
 - www.aurin.org.au
 - University of Melbourne are lead agent
 - 12+ years of my life!
- Establishing an e-Infrastructure for Urban and Built Environment Researchers
 - Distributed, (completely!) heterogeneous datasets
 - Data interrogation services
 - Security (unit level data, health data, com
 - Online analysis tools
 - Collaboration!!!

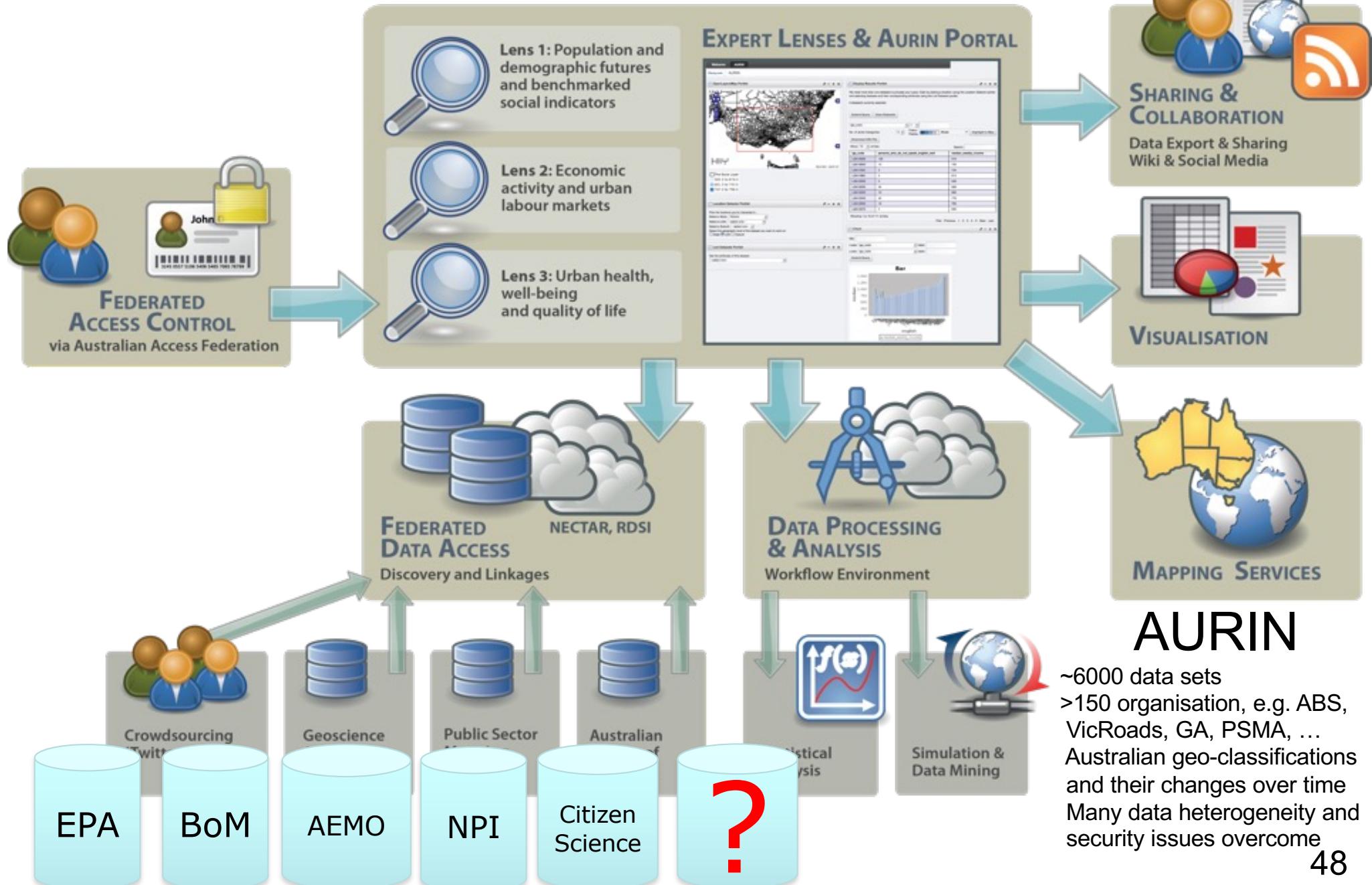


The Urban Context

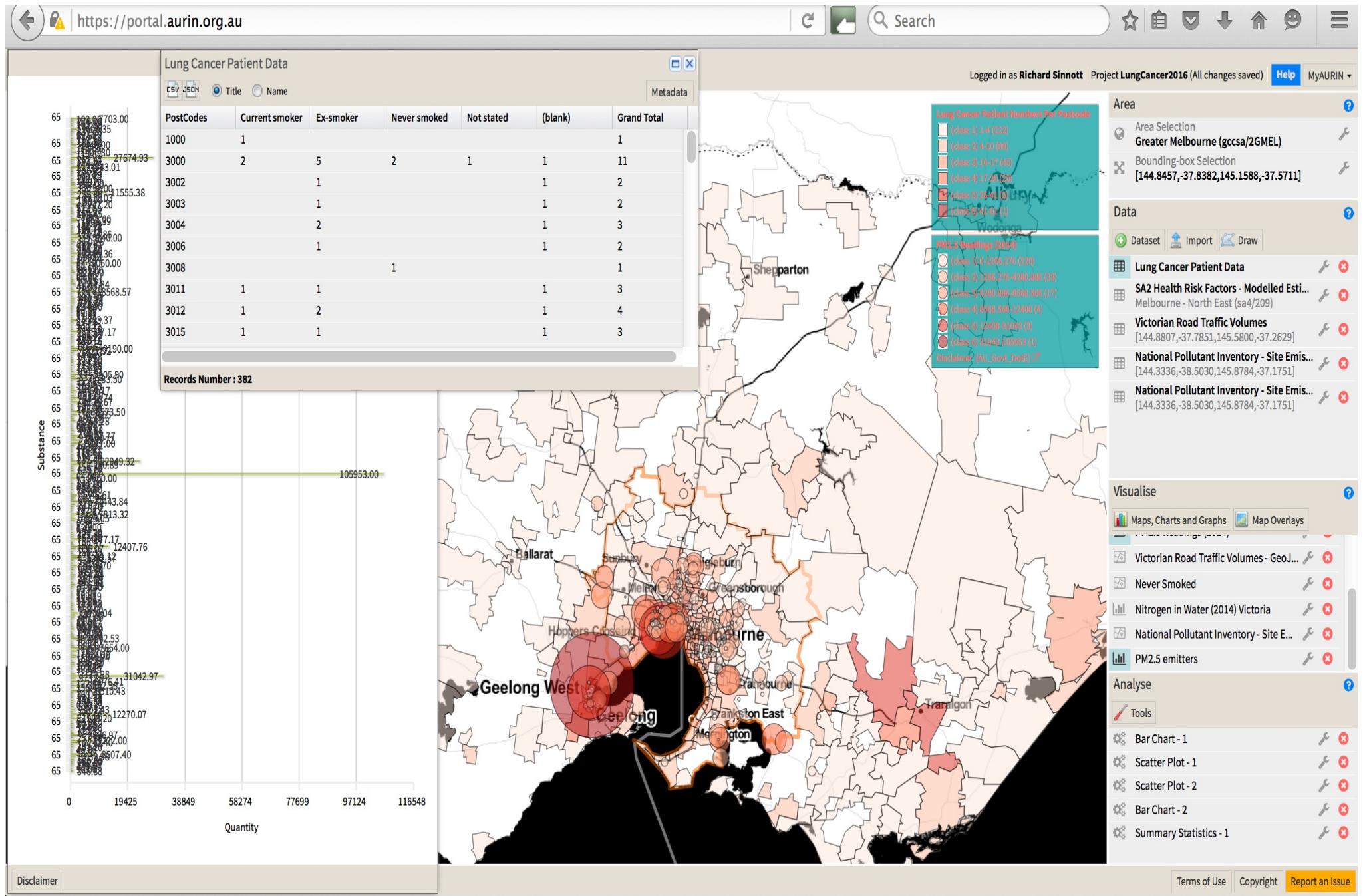
- Urban and built environment is extremely broad
 - health,
 - transport,
 - future population,
 - liveability,
 - crime,
 - housing,
 - design,
 - ...
- Much research depends on access to data
 - There is LOTS (and LOTS) of data often
 - Completely heterogeneous, e.g. geospatial
 - ...
 - Data is more often than not silo'd
- Requires tools to find, interrogate, analyze and visualize data and enforce good research methodologies
 - Consolidate tools and best practice/community know-how!
 - Allow researchers to share results, interact and collaborate
 - No single expert!
- Allow data providers to keep control of their data and its use
 - Authentication and authorisation (and auditing/accounting)



Platform Simplified

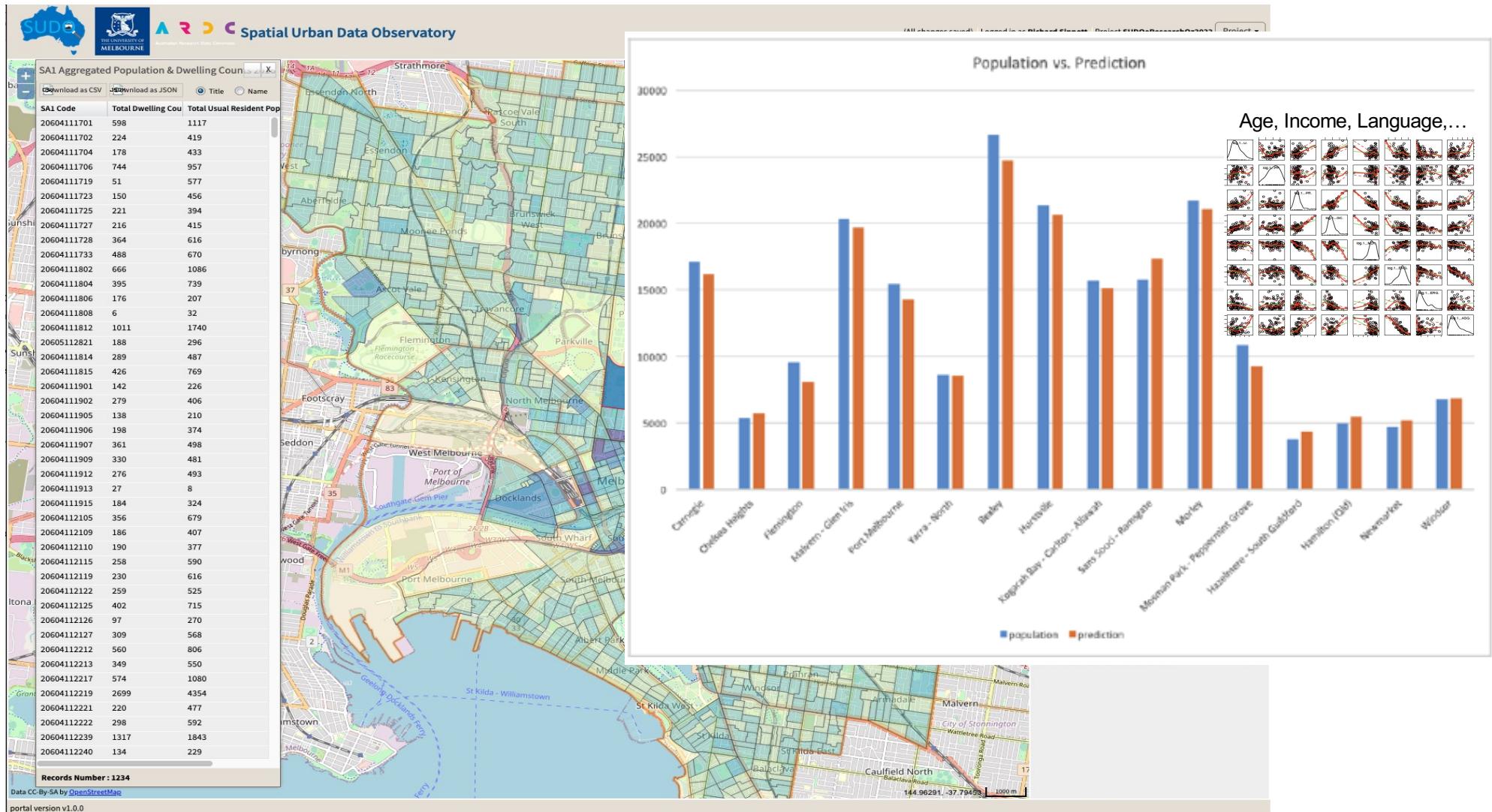


AURIN Example (in one slide!)



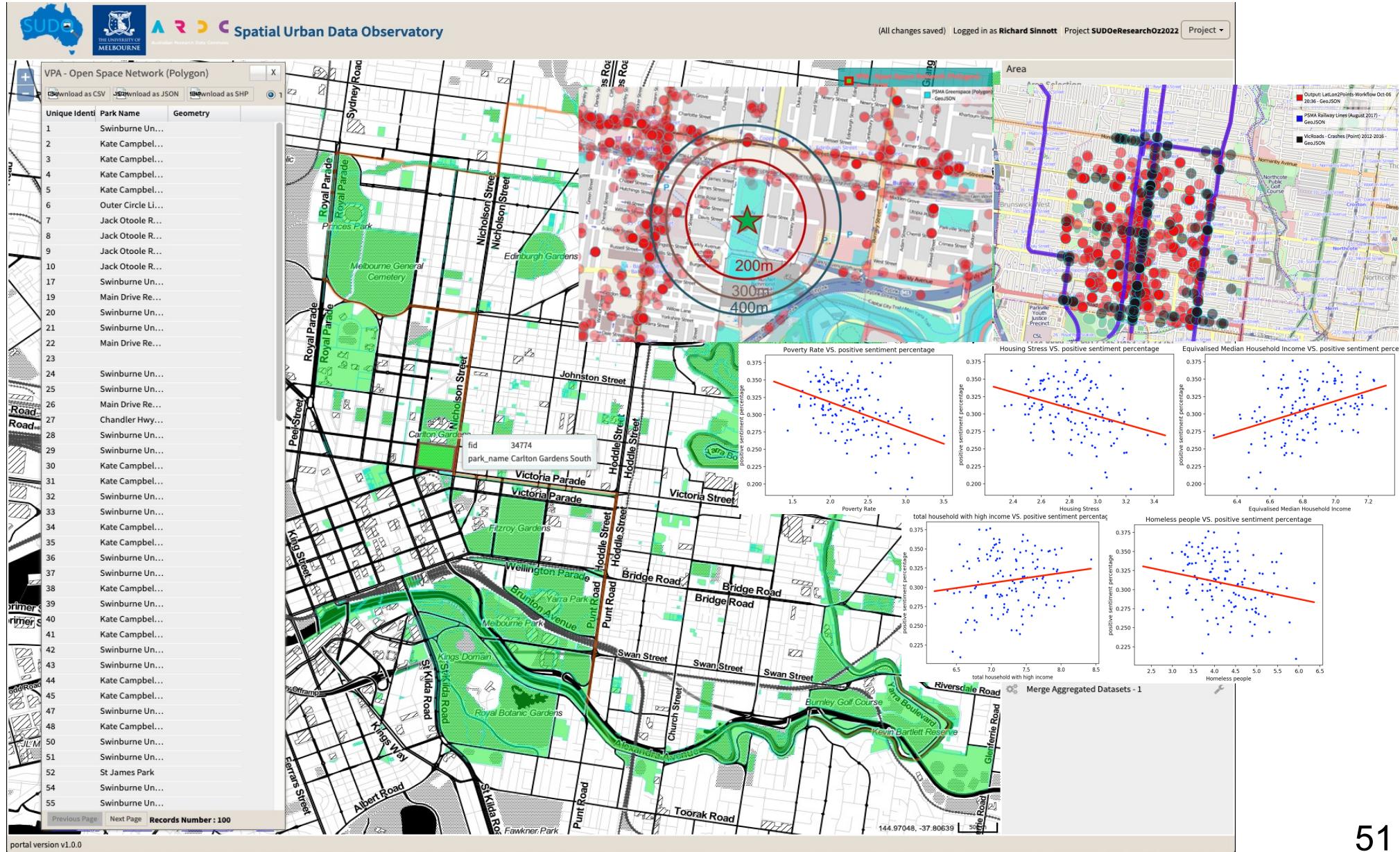
Population Estimation

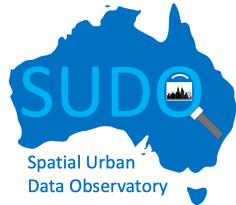
(student exemplars for ass2 inspiration)



Green Space and Happiness

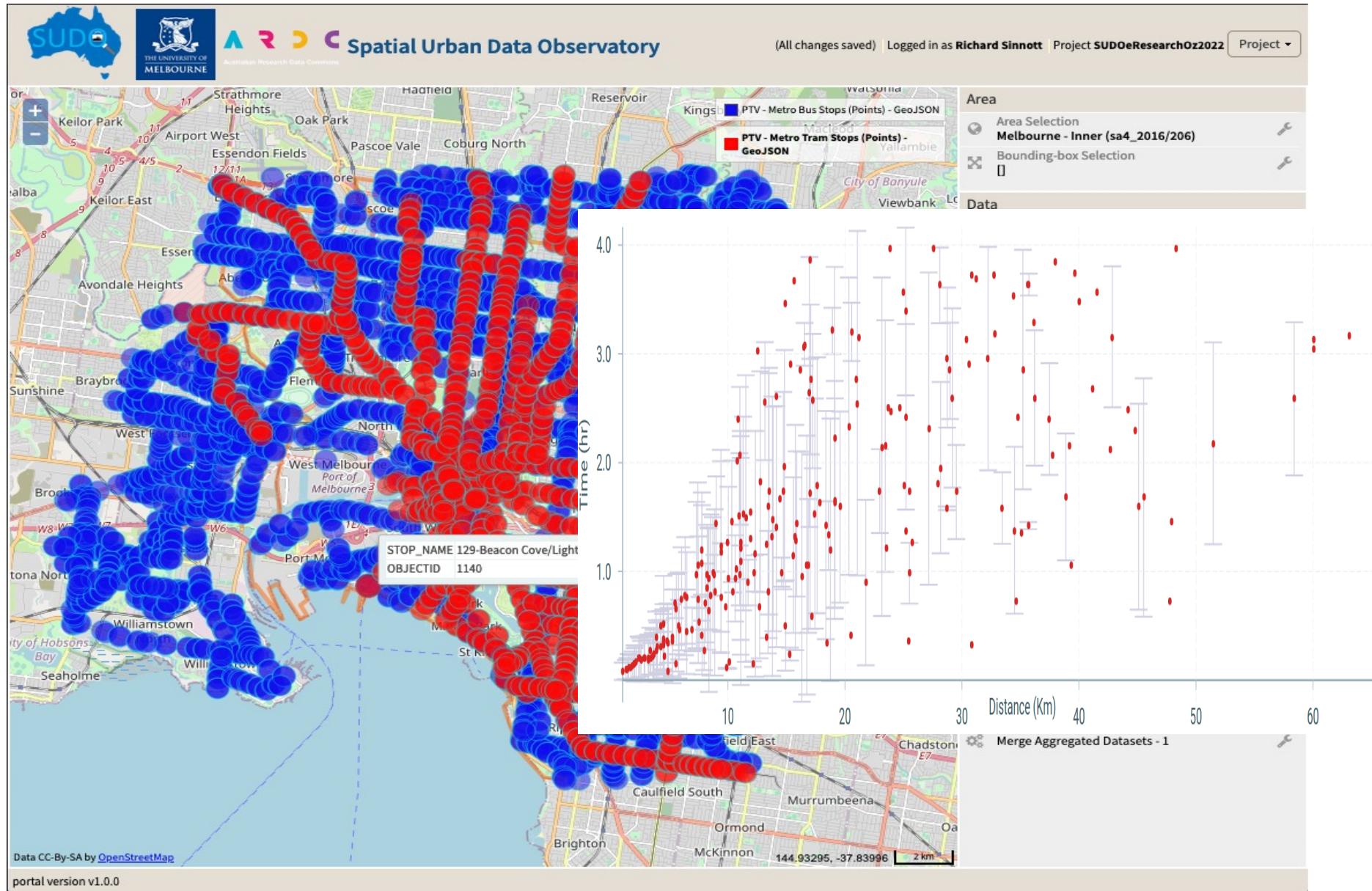
(student exemplars for ass2 inspiration)





Commuting Around Melbourne

(student exemplars for ass2 inspiration)



Demonstration of SUDO in Workshop

(Essential for Assignment II)

Homework

(<https://sudo.eresearch.unimelb.edu.au>)

Find out how many public toilets there are/were in Carlton
(not assessed, but useful for assignment 2!)

1. Log in to SUDO (select UniMelb)
2. Navigate to Australia::Victoria::SA4::SA3
3. Browse for Public Toilets (e.g., UQ_ERG)
 4. Select any/all attributes
5. Select some *other* random data set at SA2 level, e.g., from the ABS
 6. Spatialise the *other* data set

(Tools::Spatial Tools::Spatialise Aggregated Data Set and select the *other* data set) –
adds the polygon to the data

7. Count toilets

(Tools::Spatial Tools::Count Points in Polygons using the output Spatialised Data Set
and the data points from the Toilet Data Set)

Questions ... ?