

# *Lecture 2 - Domain Drivers*

Professor Richard O. Sinnott

Director, eResearch  
University of Melbourne

9<sup>th</sup> March 2017

[rsinnott@unimelb.edu.au](mailto: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
  - That 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...
    - Bioinformatics 生物信息学
    - Electronics
    - Astrophysics 天体物理学
    - High Energy Physics
    - Arts and Humanities
    - Biomedical and bioinformatics domain
- BREAK
  - Social Sciences
  - Urban research domain
  - Clinical Domain

## Completed

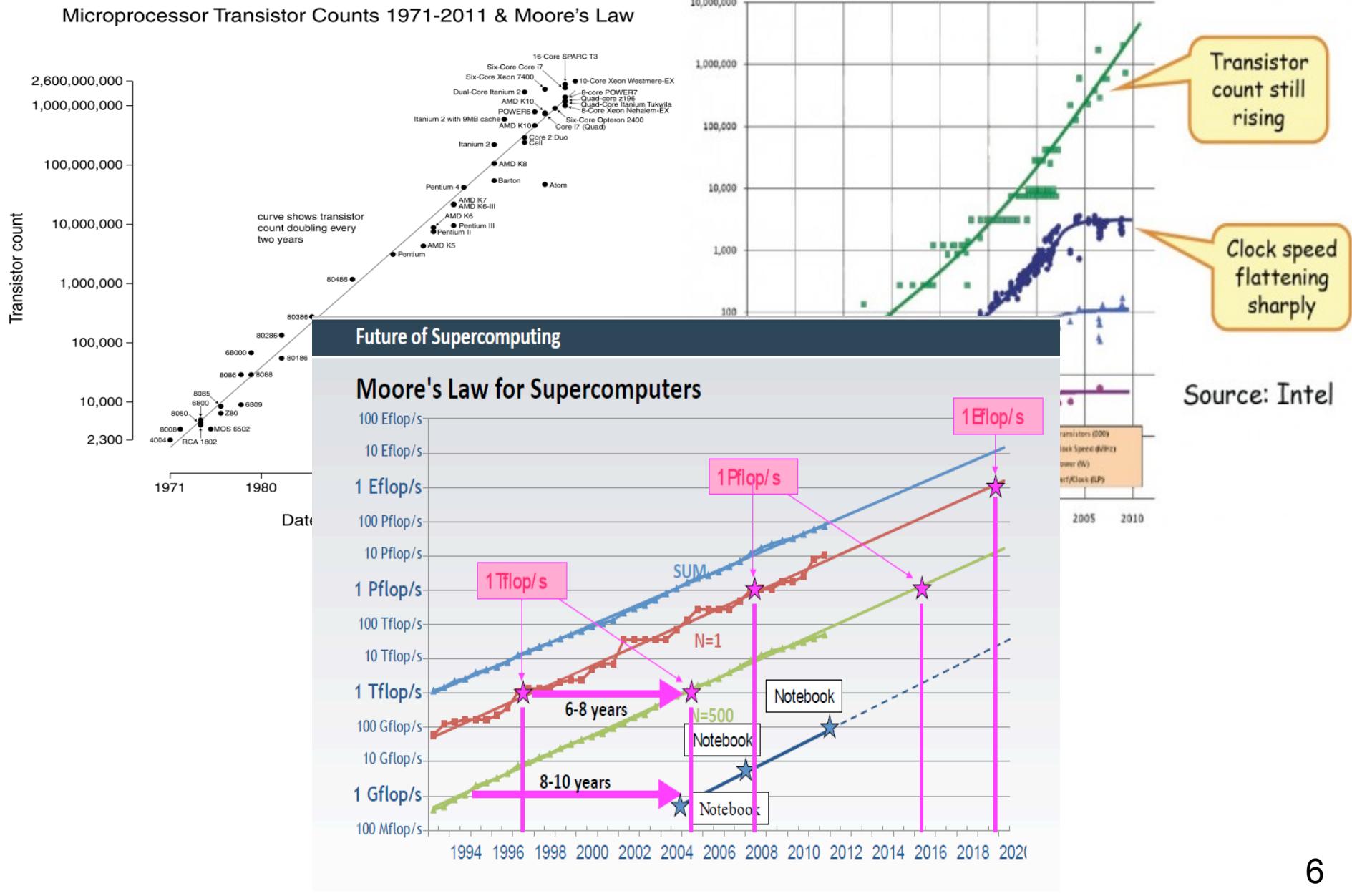
- National e-Science Centre (I, II, III)
- Dynamic Virtual Organisations for e-Science Education
- Biomedical Research Informatics Delivered by Grid Enabled Services
- GridNet, GridNet 2
- Grid Enabled Microarray Expression Profile Search
- Glasgow early adoption of Shibboleth
- Joint Data Standards 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

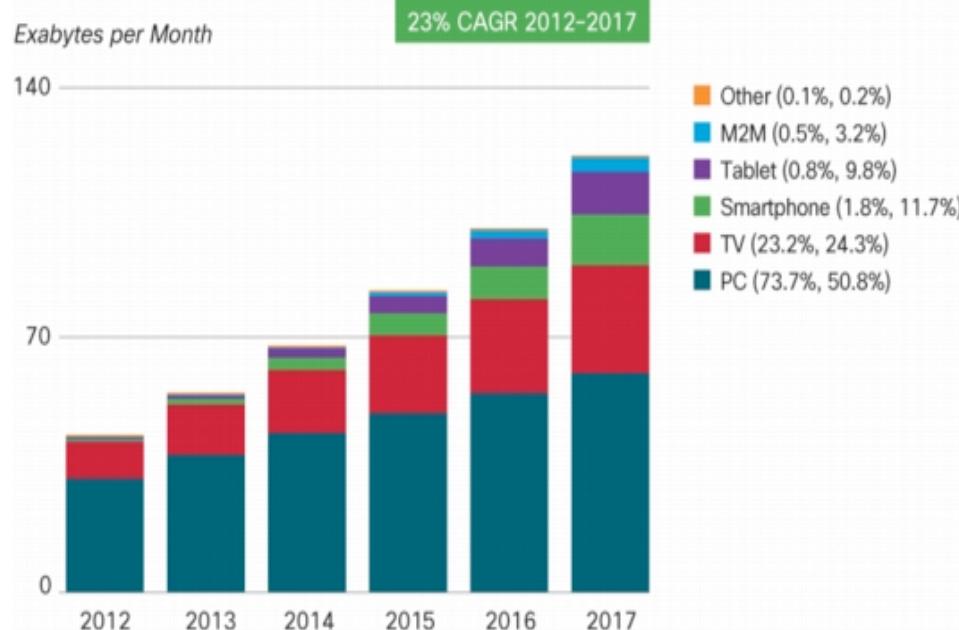
### On-Going

- 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 Phaeochromocytoma (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)
- 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 – Phase II (ADDN2)
- Helicopter advanced training system, Australian Department of Defence
- 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
- VicSpin Victoria-wide Flu Surveillance System
- ElectraNetLIDAR // VectorNZ Lidar

# Compute Scaling



# Network Scaling



Source: Cisco VNI, 2013

The percentages within parenthesis next to the legend denote the relative traffic shares in 2012 and 2017.

**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

Source: Cisco VNI, 2013

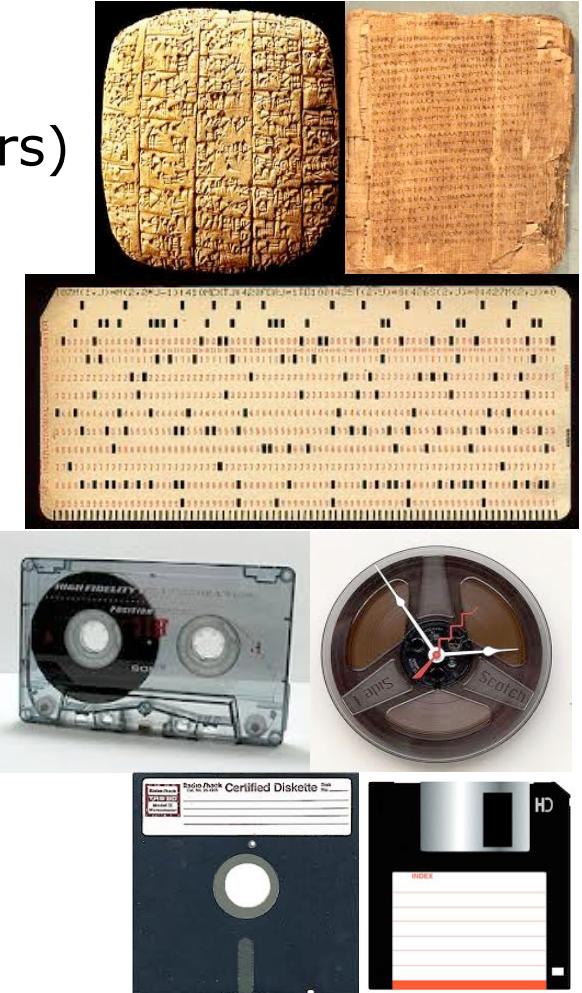
**Table 6. Table A-1 Global IP Traffic, 2012-2017**

	2012	2013	2014	2015	2016	2017	CAGR 2012-2017
<b>IP Traffic, 2011-2016</b>							
By Type (PB per Month)							
Fixed Internet	31,339	39,295	47,987	57,609	68,878	81,818	21%
Managed IP	11,346	14,679	18,107	21,523	24,740	27,668	20%
Mobile data	885	1,578	2,798	4,704	7,437	11,157	66%
<b>By Segment (PB per Month)</b>							
Consumer	35,047	45,023	56,070	68,418	82,683	98,919	23%
Business	8,522	10,530	12,822	15,417	18,372	21,724	21%
<b>By Geography (PB per Month)</b>							
Asia Pacific	13,906	18,121	22,953	28,667	35,417	43,445	26%
North America	14,439	18,788	23,520	28,667	34,457	40,672	23%
Western Europe	7,722	9,072	10,568	12,241	14,323	16,802	17%
Central and Eastern Europe	3,405	4,202	5,167	6,274	7,517	8,844	21%
Latin America	3,397	4,321	5,201	5,975	6,682	7,415	17%
Middle East and Africa	701	1,049	1,483	2,013	2,659	3,465	38%
<b>Total (PB per Month)</b>							
Total IP traffic	43,570	55,553	68,892	83,835	101,055	120,643	23%

Source: Cisco VNI, 2013

# Data Past

- From tablets, to papyrus, to books
  - (quite adequate for several thousand years)
- Enter silicon transistors <sup>纸莎草 硅晶体管</sup> <sup>大约</sup> 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

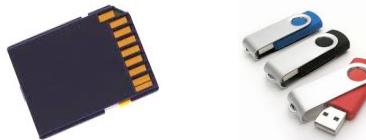


- Google



- Clouds

- ...



泛滥

# Data Deluge

- The combination of mobile devices and sensors approximated



at

(March 2007)  
digital Universe

- The total amount of data created in 2012 is approximately 48 exabytes. This is 48 times more than it was in 2003.

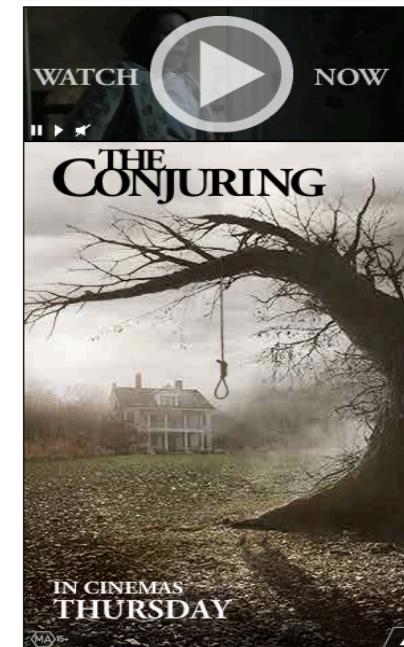
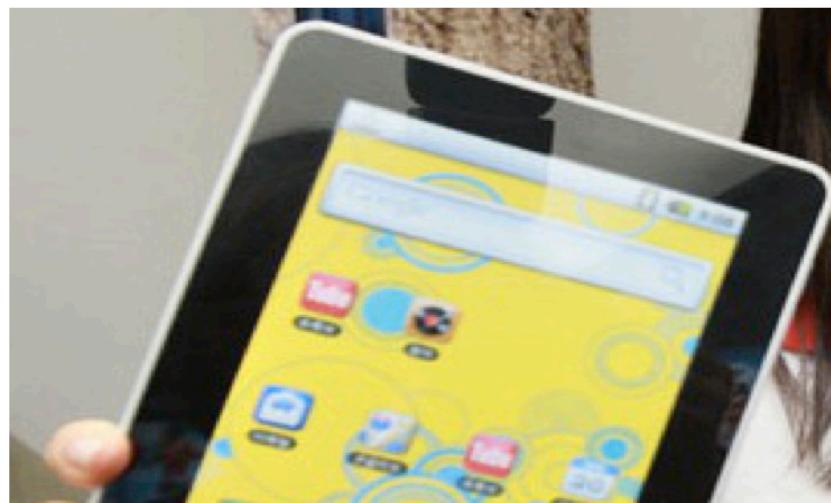


HAVE A TIP, PITCH OR  
GUEST COLUMN? TELL US.

2012.

ta Corporation

- In 2008, A



January 2010).  
n Consumers"

- By 2015, the world will have 48 exabytes of data. This is 48 times larger than it was in 2003.

uTube,  
s larger

h\_zettabyte/

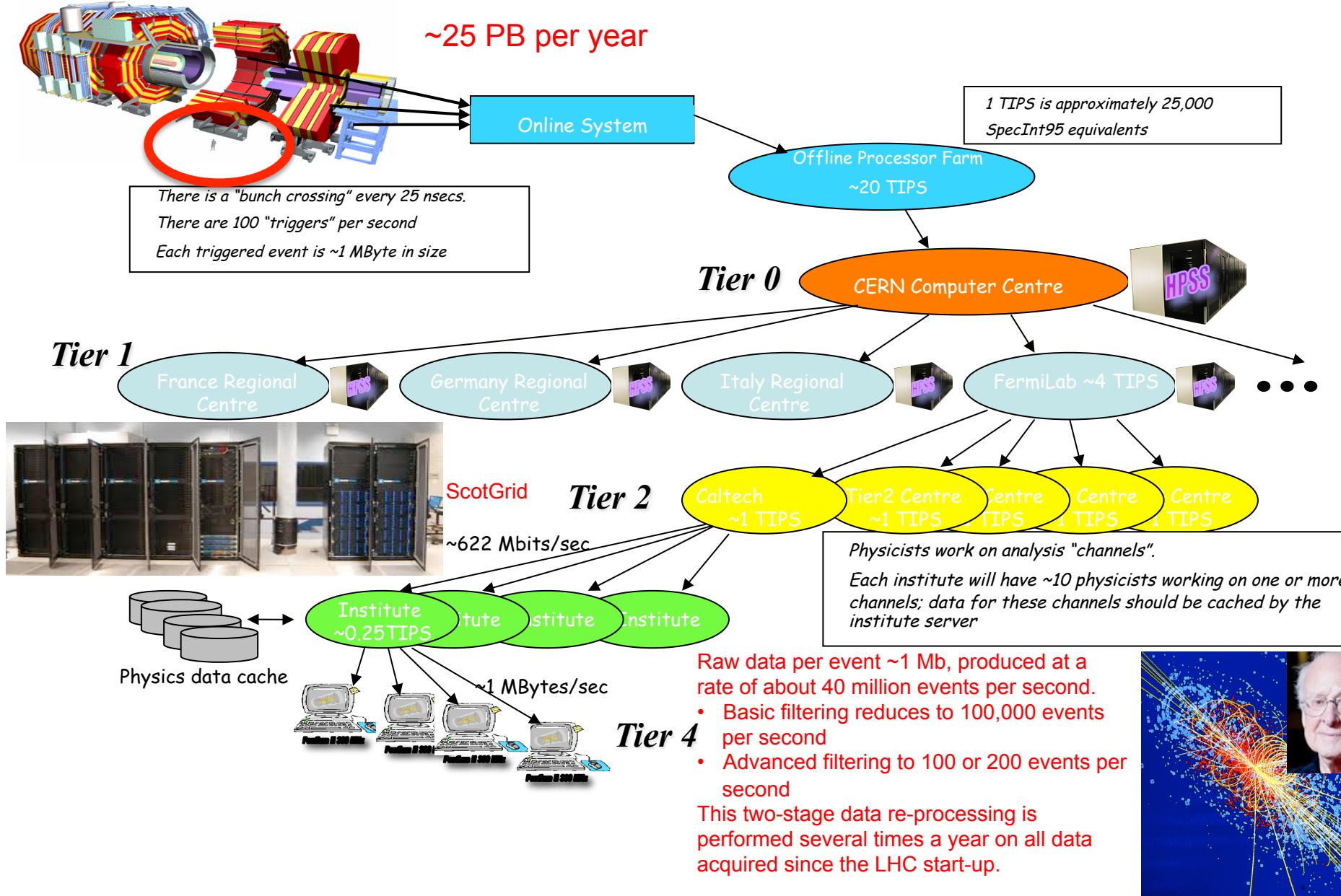
Naturally, all of this information helps Google. But he cautioned that just because companies like his can do all sorts of things with this information, the more pressing question now is if they *should*. Schmidt noted that while technology is neutral, he doesn't believe people are ready for what's coming.

*"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.

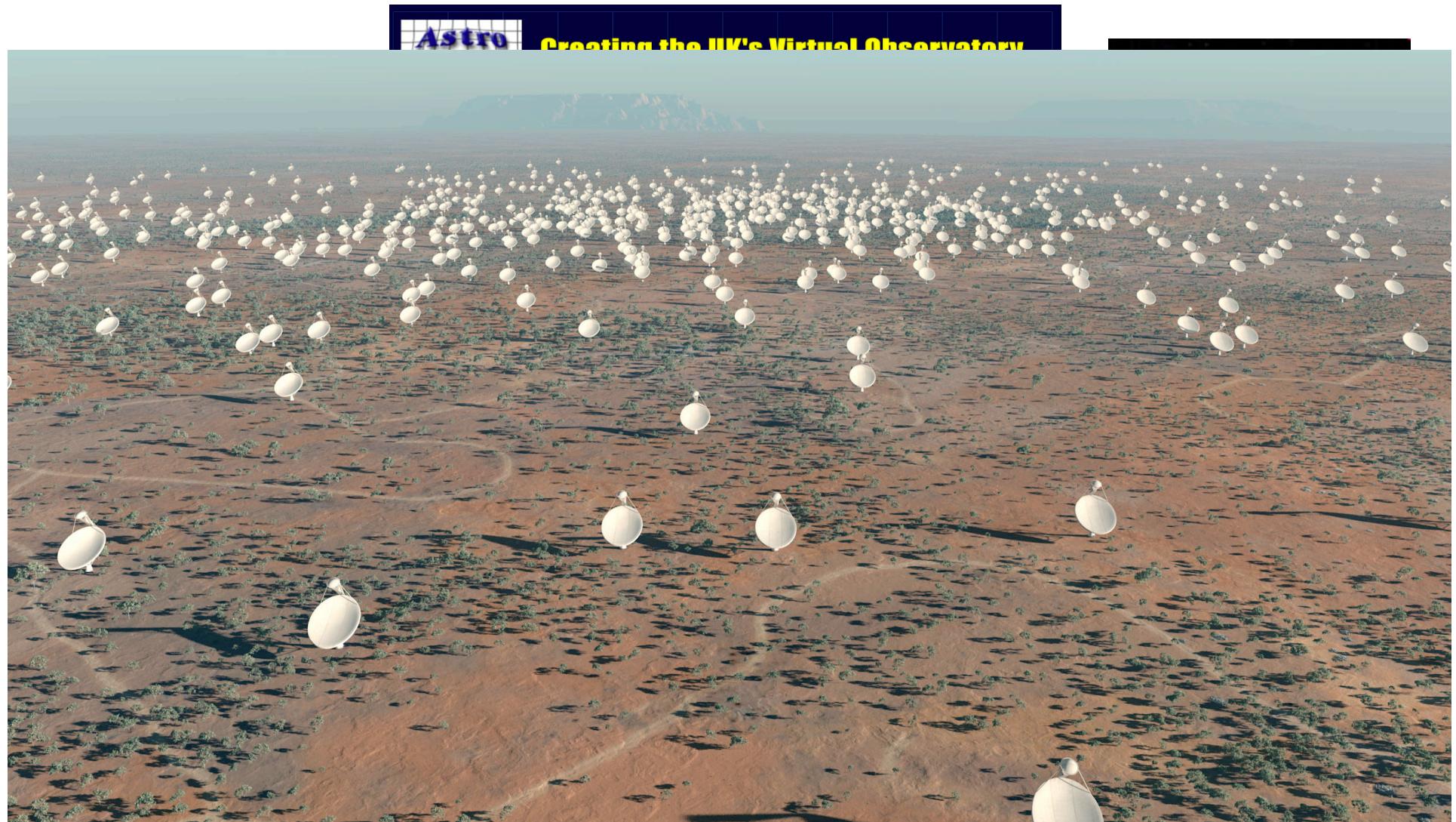
# Data Intensive / Data driven Research

- 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

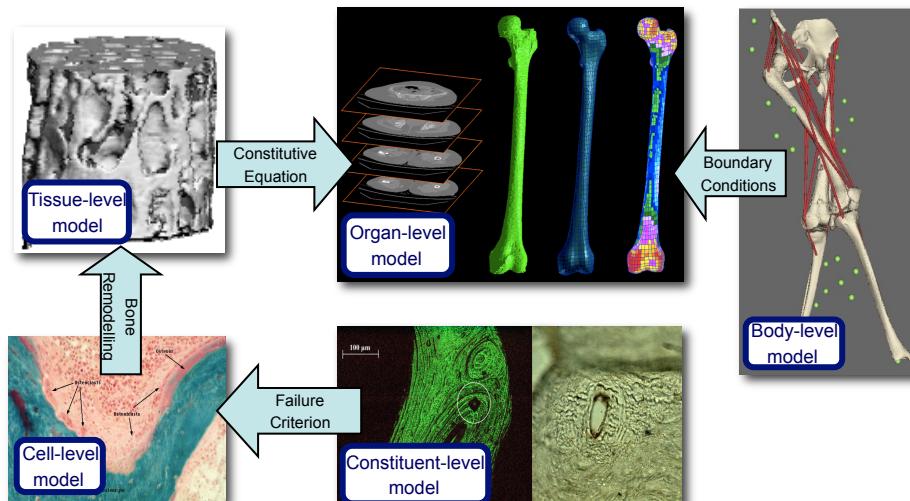
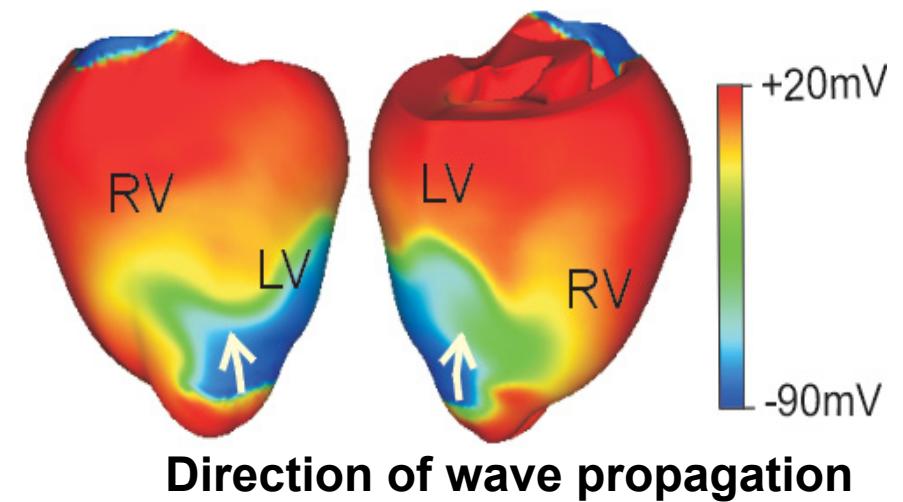
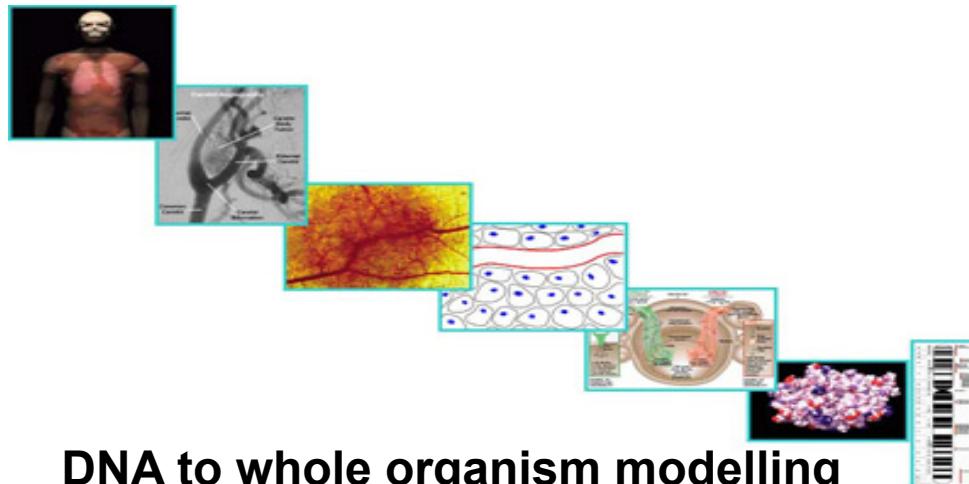


# Mapping the Skies

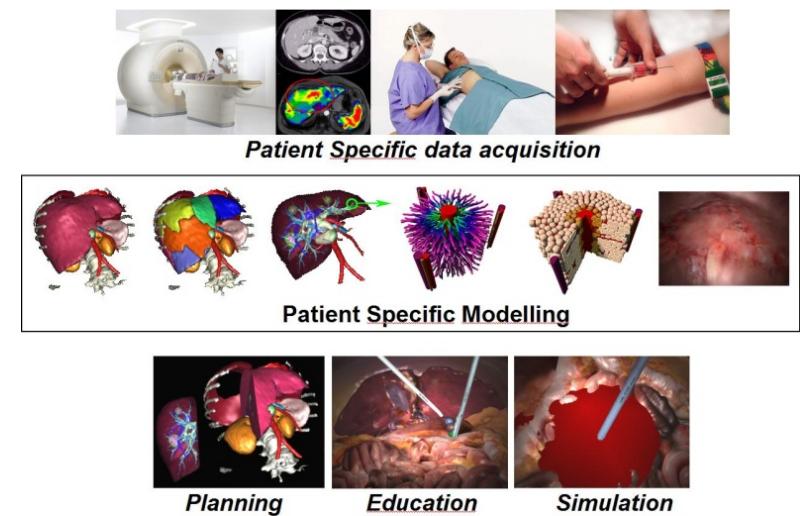


"Chipsets needed to process data access and SKA applications will need to be capable of 20-25 exaflops of processing power", according to IBM Research's Ton Engbersen, DOME scientist and project leader. "Take the current global daily Internet traffic, double it, and you are in the range of the data set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

# Macro-micro Simulations

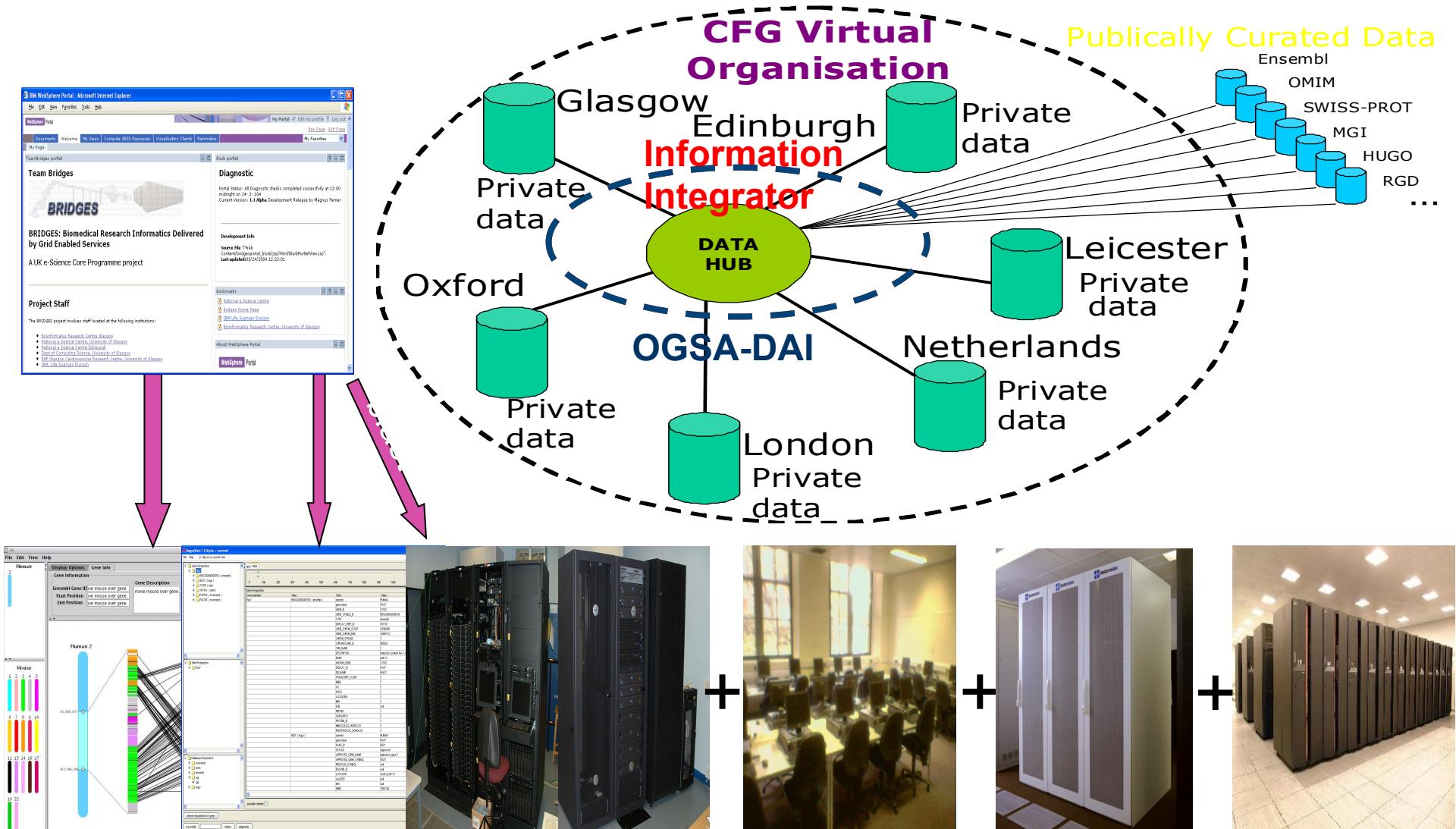


Develop the 1st multiscale, patient-specific model of the musculoskeletal system and use model for diagnosis, prognosis and treatment of osteoporotic fractures

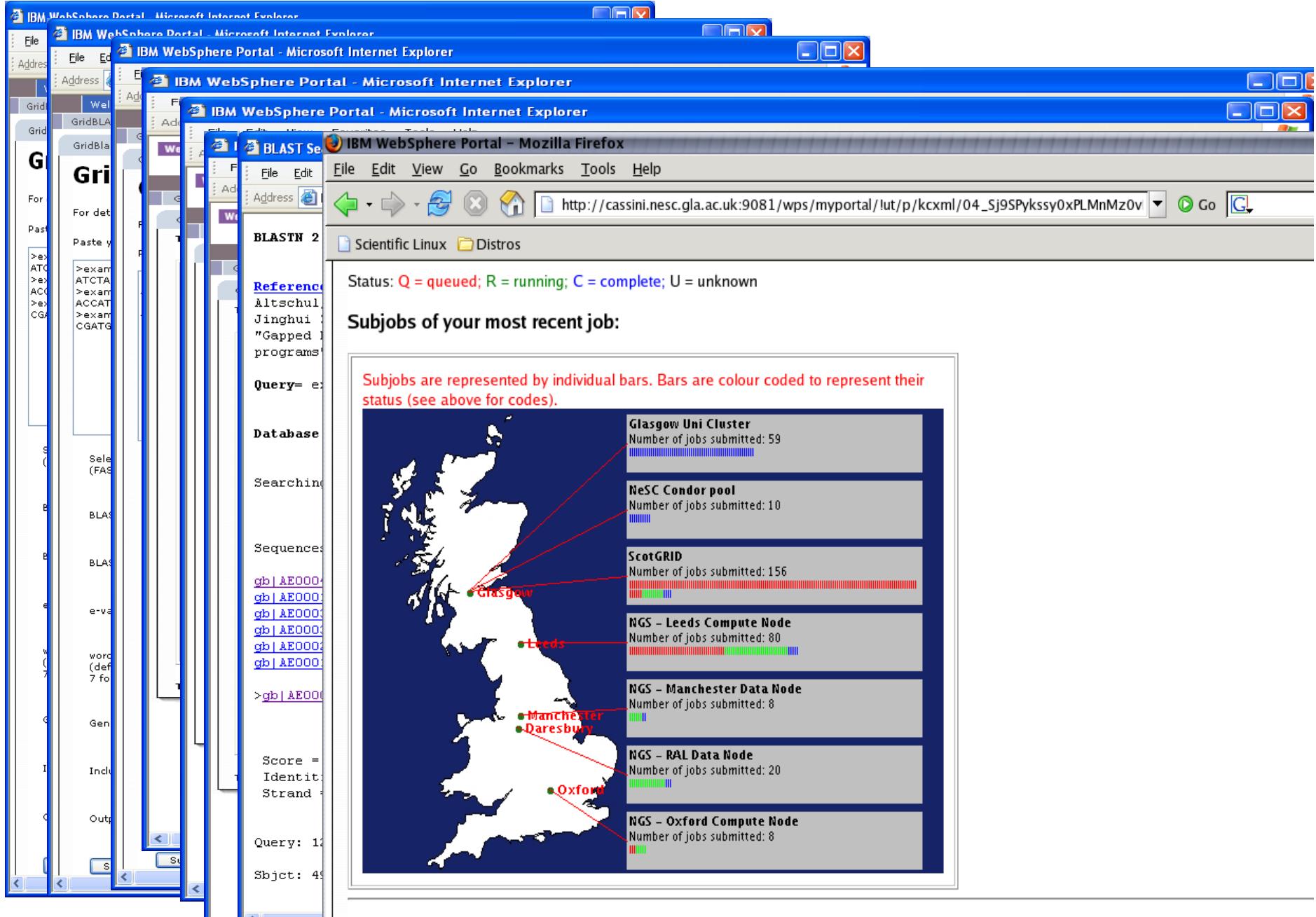


Patient Specific Simulation and PreOperative Realistic Training for liver surgery

# Example of BRIDGES Project



# Grid Blast Interface



# 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

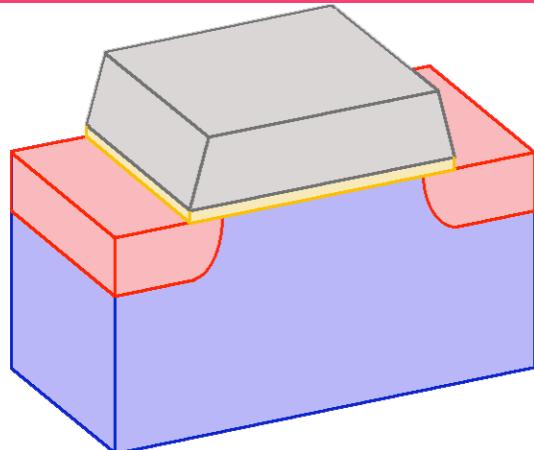
11 PDRAs (7 + 4)  
7 PhD

## Industrial partners

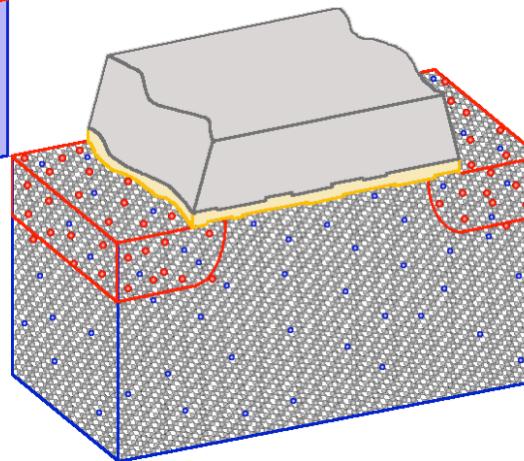


## University partners

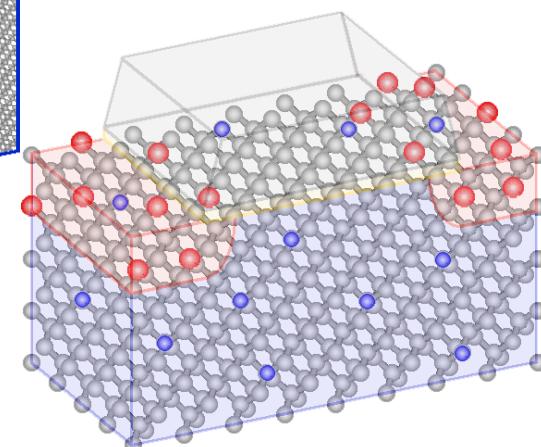
# Semiconductor device variability



Historic simulation paradigm

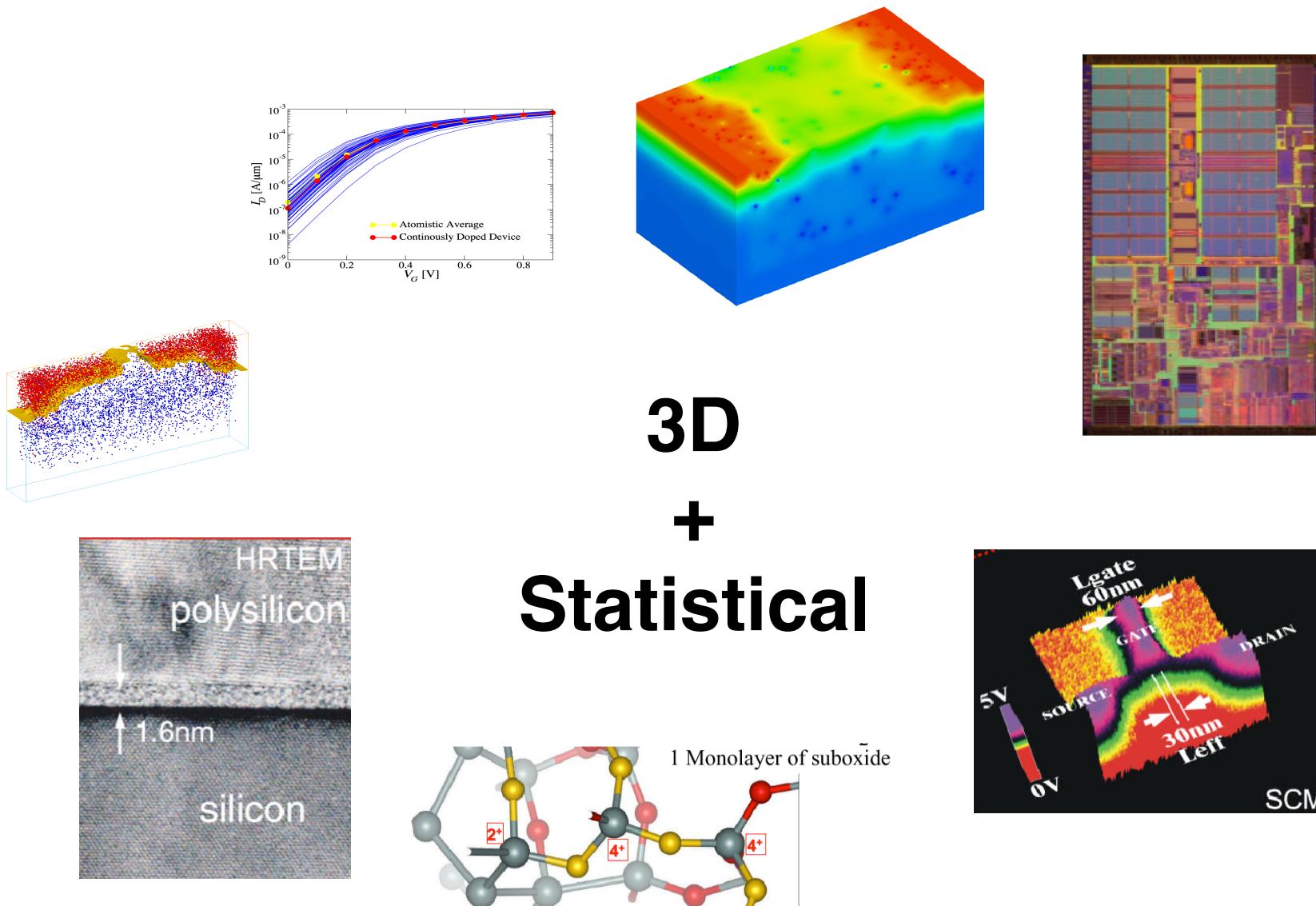


A 22 nm MOSFET  
In production 2011



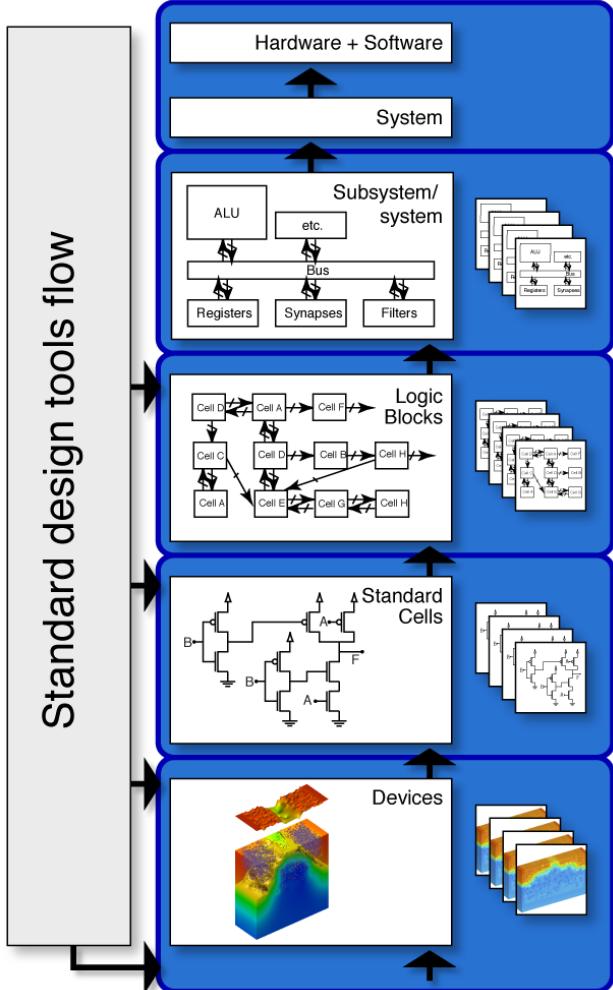
A 4.2 nm MOSFET  
In production 2023

# Challenges of NanoCMOS Design



# Challenges

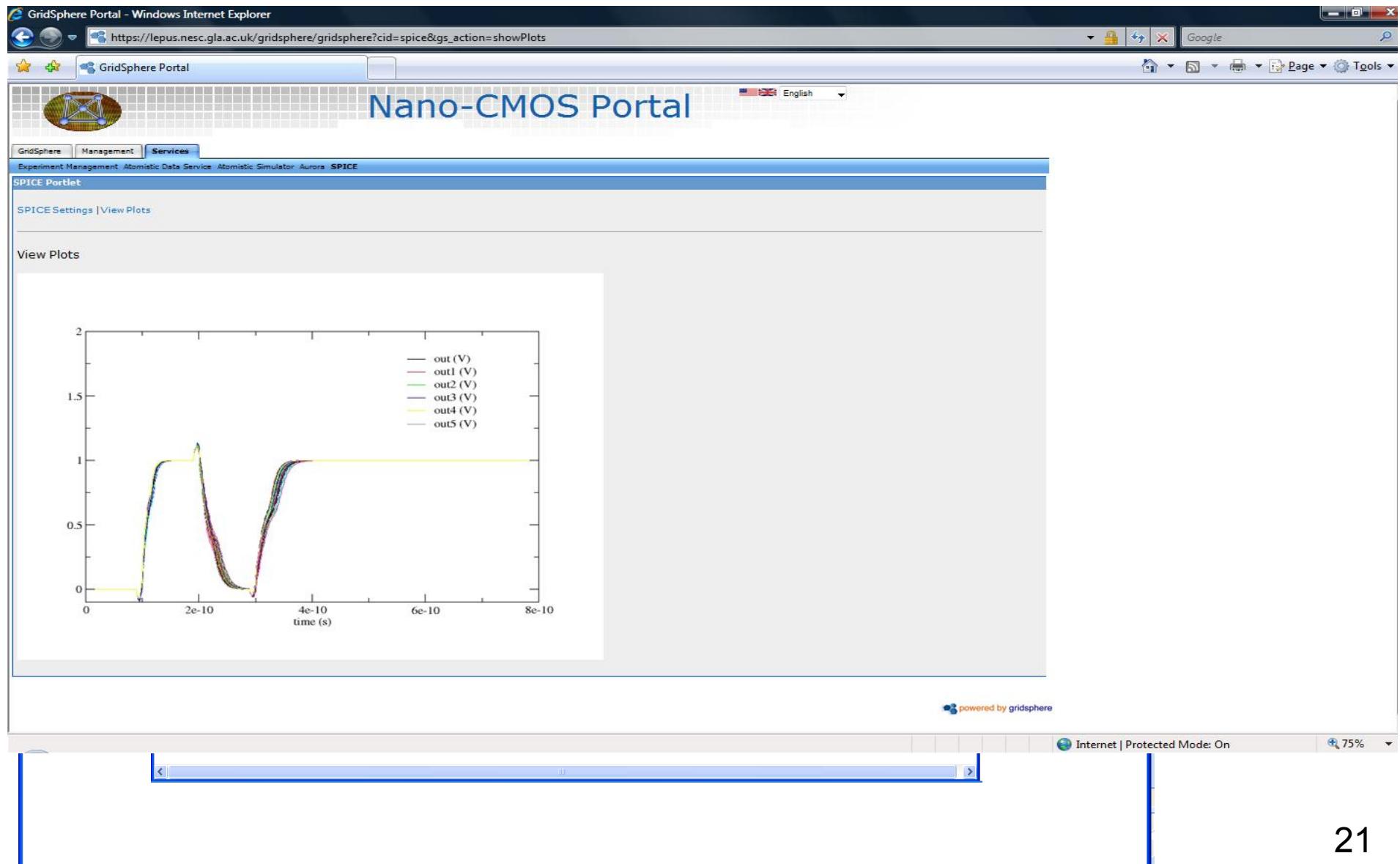
# Hierarchical statistical system simulations



- ❑ 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-Experiences

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



# But ended up with...

- ... a command-line based solution

This community

Security solution

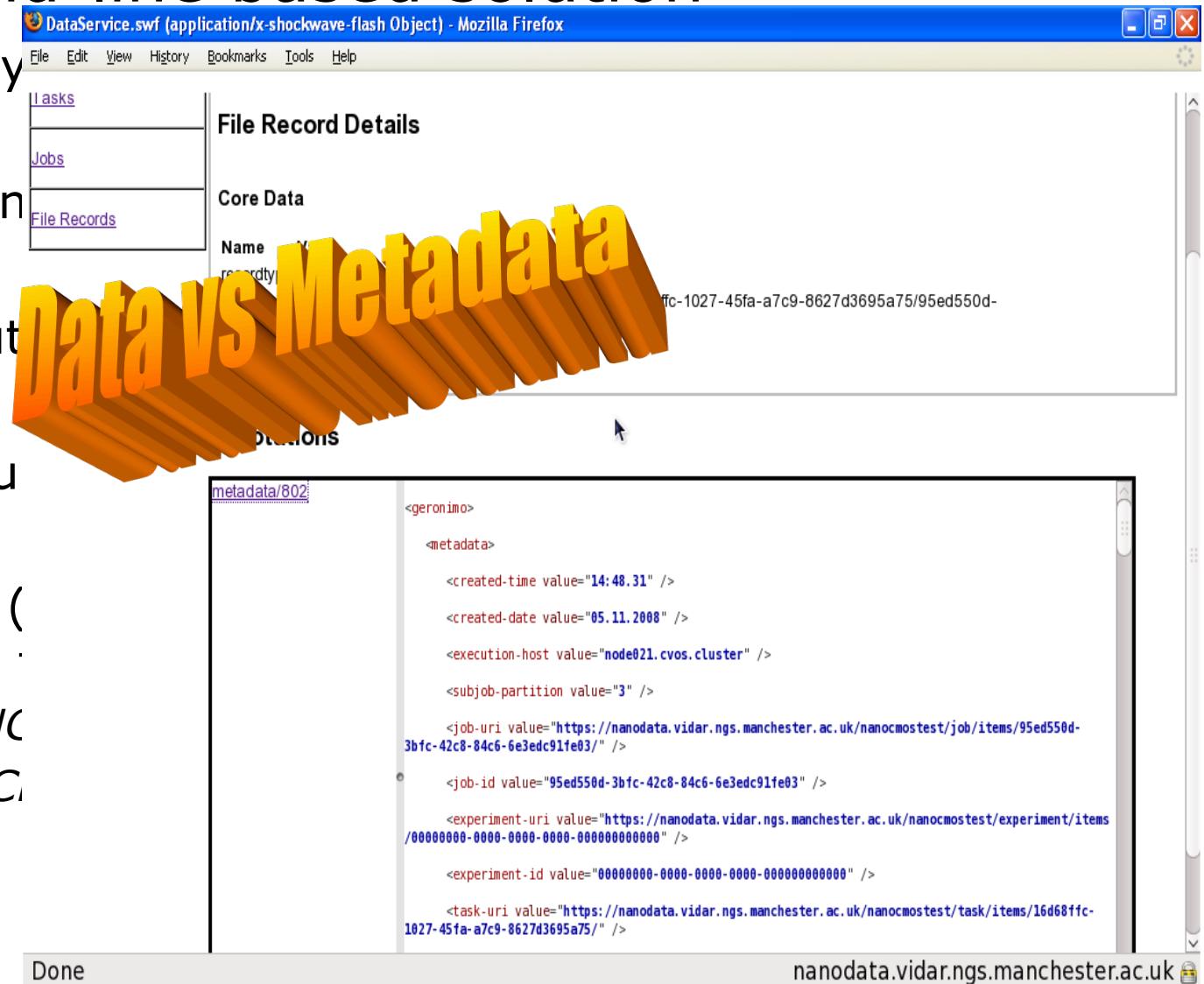
Secure, distribut

Meta-data captu

Job submission (massive (at the

- *ScotGrid, NC*
- *Millions of C*

*The -g flag!!!*



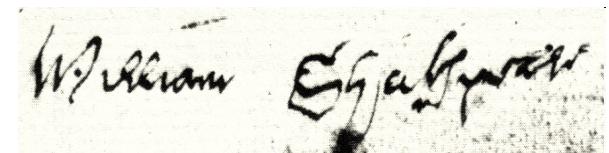
# Arts and Humanities

- Data resources <sup>大量</sup>galore
- Language and literature research
  - Linguists
  - Philologists
  - Historians
  - Societal observers
  - Culturists
  - ...
- Written resources
- Spoken resources
- Video resources

# Good old Bill



- To be or not to be...
    - Shakespeare, Shakespere, Shakespear, Shakspeare, Shackspeare, Shakspere, Shackespeare, Shackspeare, Shaxspere, Shaxpere, Shagspere, Shaksper, Shaxpeare, Shaxper, Shakespeare, Shakespe, Shakp
    - Inside London, Outside London,
    - Printed, cited, referenced
      - 8 signatures – all different!



# Most “Excellent” Dude

- seldlic syndorlic til wlitig special breme gradely noble  
thriven and thro burly singular dainty gentle proper  
beforepassing goodly daintiful thriven thriving  
priseprise vounde virtuous curious principal fine gay  
rare singlar egregious gallant eximious jolly jelly  
braw brave stamming surprising excelling phoenix  
royal of worth rare admirable sublimated valiant  
excellent twanging topgallant lovely prestantious  
splendid sterling spanking hogenhogan pure licking  
rattling tearing soaring famous yrare pure and daisy  
immense capital elegant trimming gallows budgeree  
crack dandy smicksmack
- 02.01.15.07.08.06 adj, Historical Thesaurus of English

# I ain't “Through” yet...

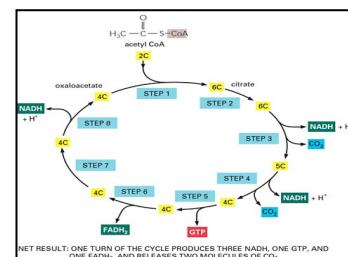
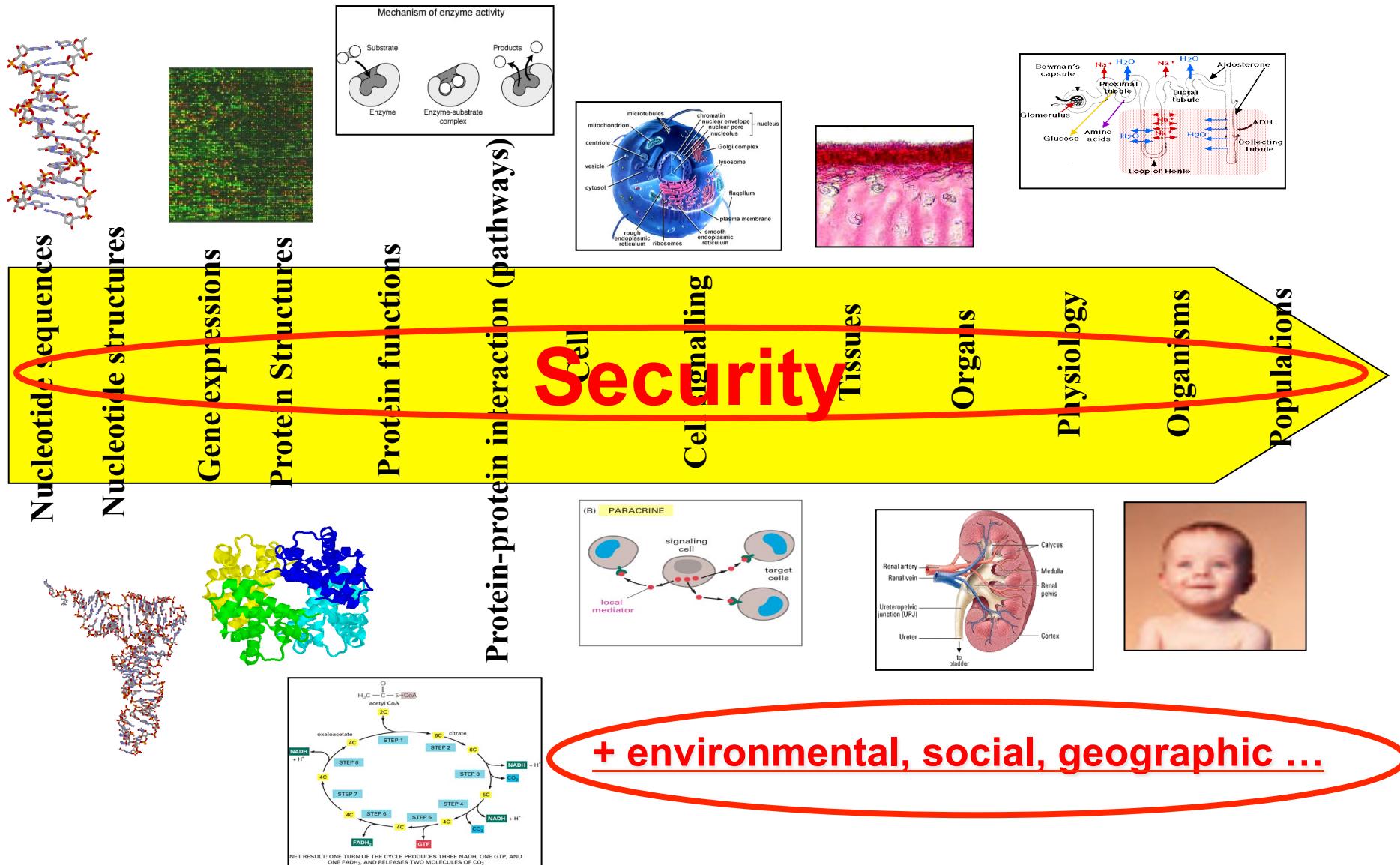
*Thro(u)ch, Through(e, Throw(e, Thorow(e,  
prep. Also: throuche, throwch(e, throwgh(e,  
thro(u)cht, throwcht, throught, threwch,  
thrwch, throuthe, thruch(t, thrugh(t, thruth(t,  
throche, throgh(e, throiche, throicht,  
throcge, thro(u)g, throu, throue, threu,  
threw, thru, thrw, thro, troch(t, trouch, trew,  
troithe, trhow, thorowch, thorycht, thorrow,  
thoro, thurrow*

Thro(u)ch (preposition), DOST

# ENROLLER

The screenshot shows a Mozilla Firefox browser window displaying the ENROLLER search interface. The address bar shows the URL <https://enroller.nesc.gla.ac.uk/web/guest/portal>. The page title is "Search - Liferay". A colorful abstract logo is on the left, followed by the word "ENROLLER". The top navigation bar includes links for "Search", "Advanced Grid Search", "My ENROLLER", and "About Us". A welcome message "Welcome richard sinnot!" is visible on the right. The main search area has a text input field with placeholder text "Enter a search term, or multiple search terms separated by commas:" and a "Search" button. Below the input field is a section titled "Choose the datasets you would like to search across:" with several checkbox options. Some checkboxes are checked (e.g., "The Historical Thesaurus of English", "The Scottish Corpus of Text and Speech", "The Dictionary of the Older Scottish Tongue") while others are unchecked (e.g., "The Oxford English Dictionary", "The Newcastle Electronic Corpus of Tyneside English", "The Corpus of Modern Scottish Writing", "The Middle English Grammar Corpus", "The Thesaurus of Old English").

# The e-Health Future...

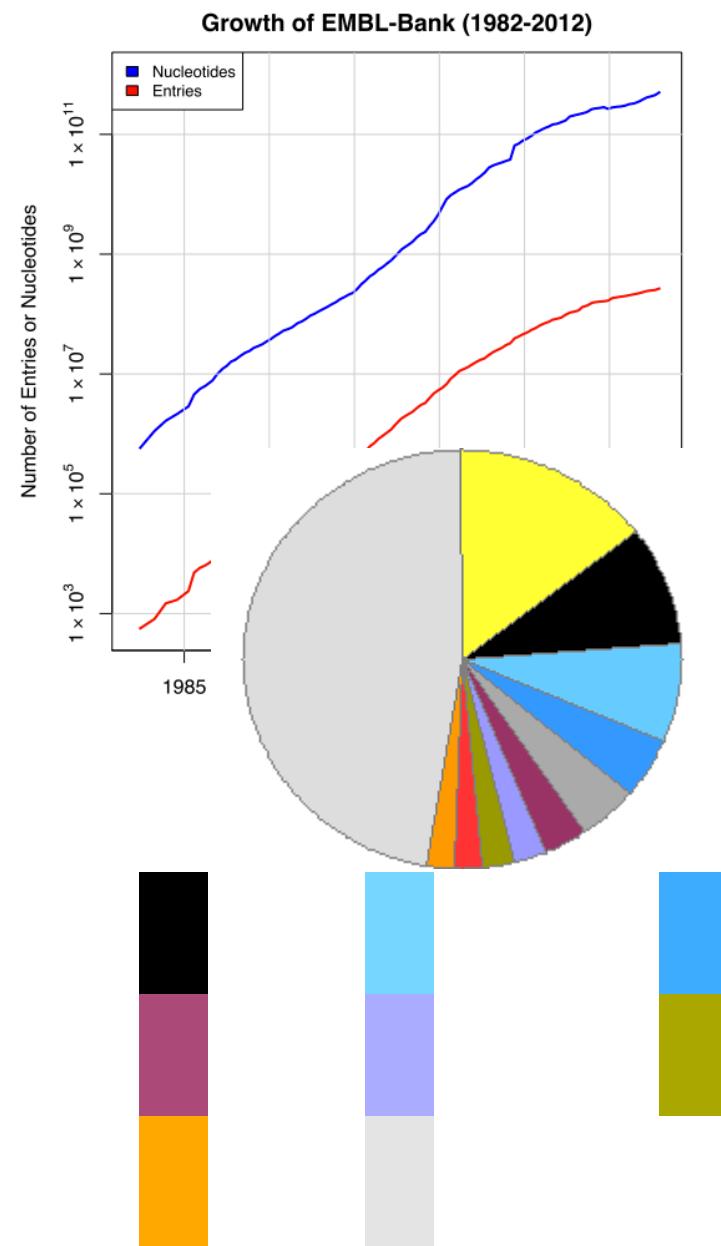
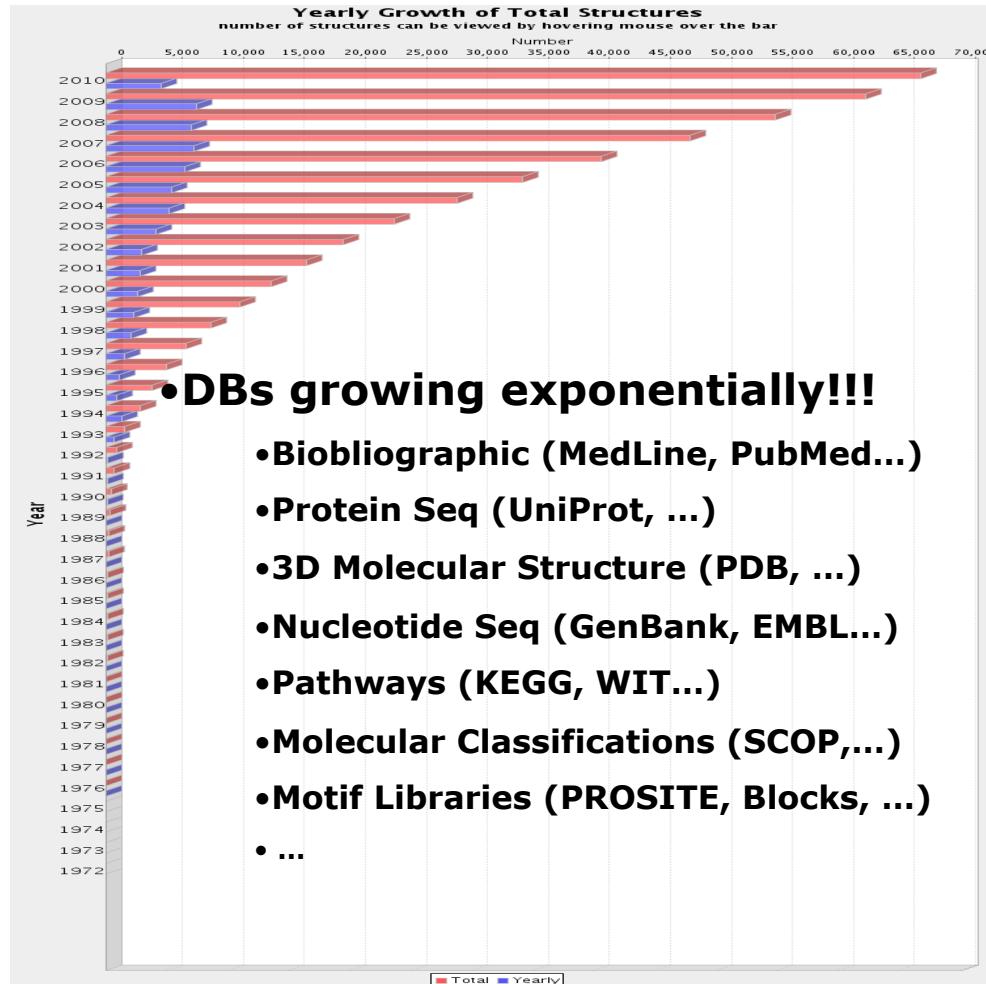


# Life Sciences

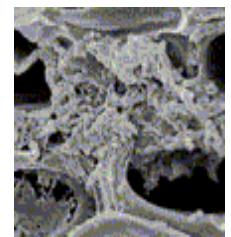
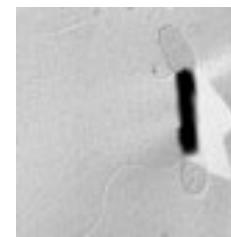
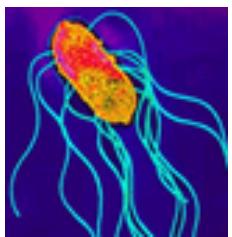
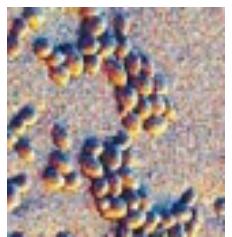
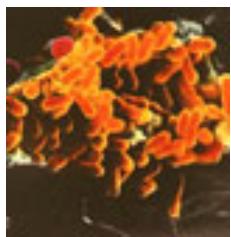
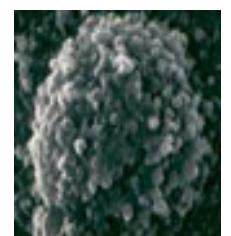
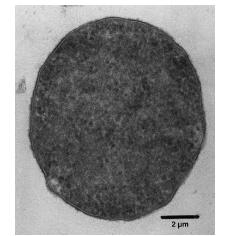
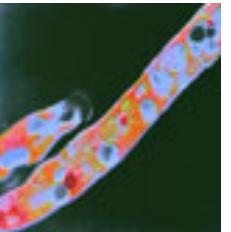
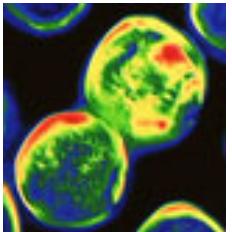
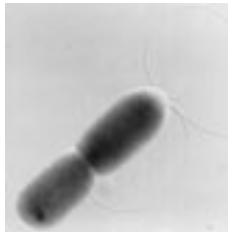
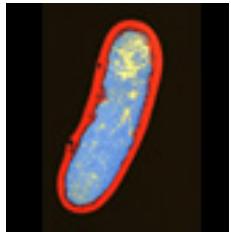
- Extensive Research Community
  - Parkville Precinct for example
- Many people care about them
  - Health, Food, Environment – truly interdisciplinary!
- Interacts with virtually every discipline
  - Physics, Chemistry, Maths/Stats, Nano-engineering, ...
- Thousands of databases relevant to bioinformatics (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?
  - ...

# Database Growth

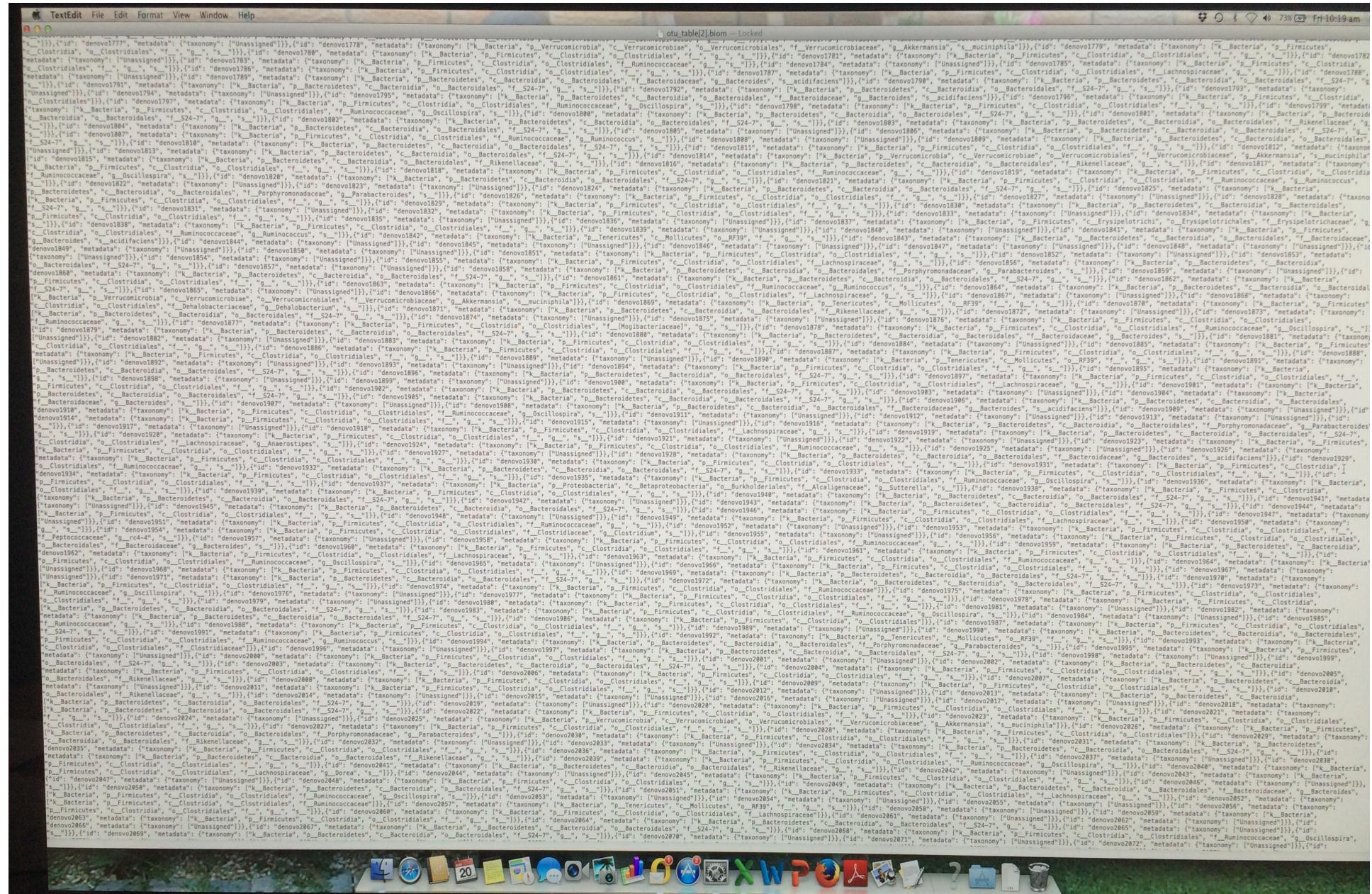
## PDB Content Growth



# More (and more and more) genomes...



# Distributed, completely heterogeneous data



# Messy

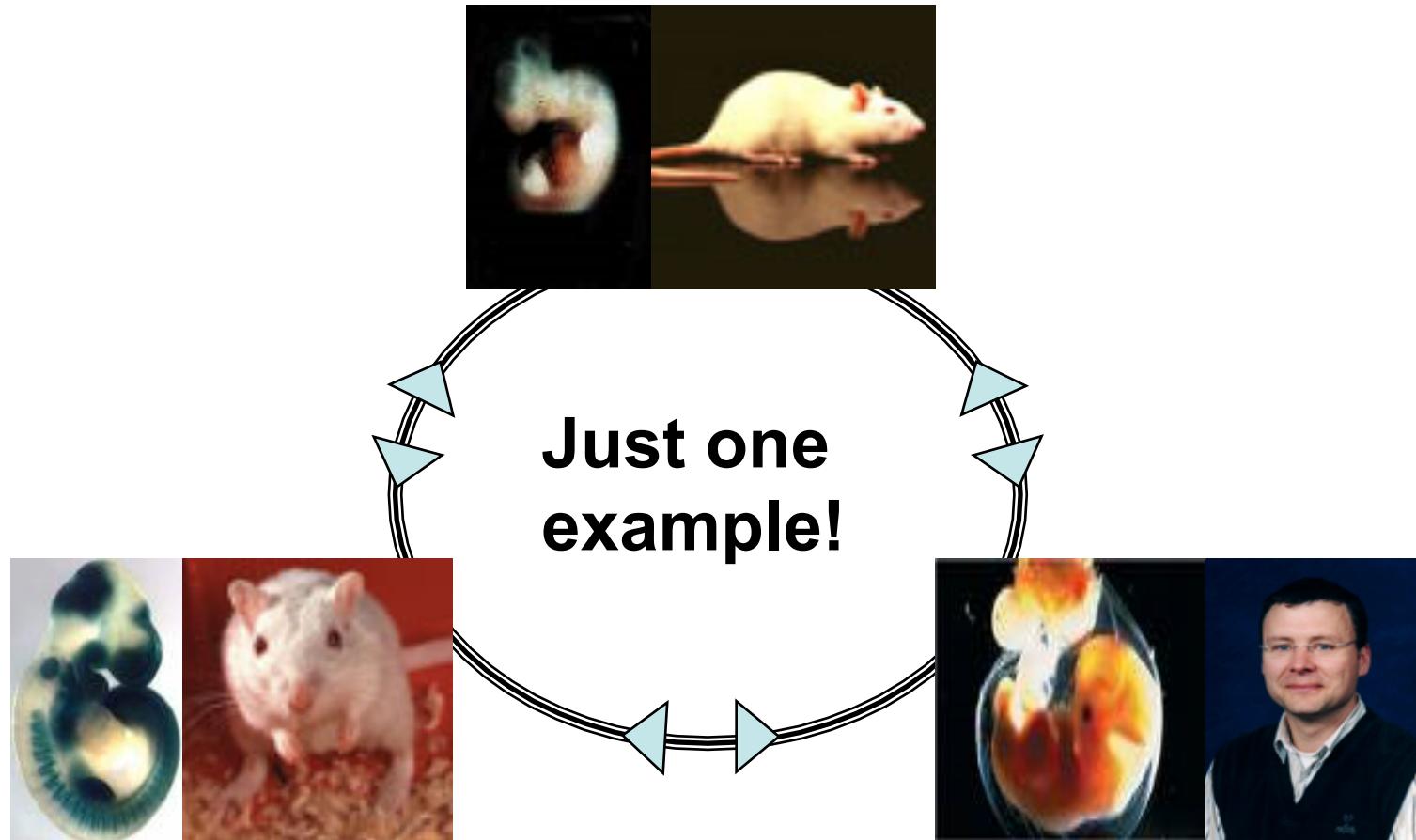
GPL96 - Notepad

File Edit Format View Help

Annotation!Annotation\_date = 09/20/2006 15:35:01!Annotation\_platform = GPL960!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 P450IIE1 (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  
002954 RPS27A 0200010\_at ribosomal protein L11 Homo sapiens ribosomal protein L11 (RPL11), mRNA 15431289 NM\_0009  
55 (RPS5), mRNA 71164878 NM\_001009 RPS5 0200018\_at ribosomal protein S13 Homo sapiens ribosomal protein S13 (RPS13), mRNA 1459191  
(RPS11), mRNA 34335149 NM\_001015 RPS11 0200025\_s\_at ribosomal protein L27 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 NM\_004515 0200032\_s\_at ribosomal protein L31 Homo sapiens ribosomal  
ing factor 2, 45kDa (ILF2), mRNA 24234746 NM\_004515 0200039\_s\_at ras homolog gene family member 1 Homo sapiens ras homolog gene family member 1, mRNA (c  
SCC3L1 0200059\_s\_at ras homolog gene family member 1 NM\_004515 0200040\_s\_at elongin B1 Homo sapiens elongin B1, mRNA (c  
2288 AF275719 HSP90AB1 0200059\_s\_at ras homolog gene family member A, mRNA (c  
eterogeneous nuclear ribonucleoprotein, alpha (FNTA), mRNA 34335149 NM\_004515 0200040\_s\_at elongin B1, mRNA (c  
RNA (cDNA clone MGC:4498 IM4498) NM\_004515 0200079\_s\_at lysyl-tRNA synt  
014267 C11orf58 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
rase, CAAx box, alpha (FNTA), mRNA 34335149 NM\_004515 0200079\_s\_at lysyl-tRNA synt  
omo sapiens ATPase, H+ trans 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
U 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
, mRNA 4507676 NM\_003299 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
ein kinase, cAMP-dependent, regul NM\_003299 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
006265 RAD21 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
rotein complex 2, beta 1 subunit (AP2B1), mRNA 4507676 NM\_003299 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
n 3 (phosphorylase kinase, delta) NM\_004184 WARS 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
mRNA 47419913 NM\_004184 WARS 0200040\_s\_at elongin B1, mRNA (c NM\_004515 0200079\_s\_at lysyl-tRNA synt  
Homo sapiens protein tyrosine phosphatase, receptor type, F (PTPRF), transcript variant 1, mRNA 109633040 NM\_002840 PTPRF  
yptophen 5-monooxygenase activation protein, zeta polypeptide 0646\_s\_at nucleobindin 1 Homo sapiens tyrosine 3-monooxygenase/tryptophan 5-monooxygenase 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 02000658\_s\_at prohibitin Homo sa  
NM\_006145 DNAJB1 02000665\_s\_at secreted protein, acidic, cysteine-rich (osteonectin) 02000658\_s\_at prohibitin Homo sa  
g protein 1 (XBP1), mRNA 14110394 NM\_005080 XBP1 0200671\_s\_at spectrin, beta, non-erythrocytic 1  
7977 NM\_003347 UBE2L3 0200677\_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\_at protein phospa  
mRNA 8051609 NM\_006854 KDELR2 0200700\_s\_at KDELR (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 (hPTCAAX1) 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 P-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

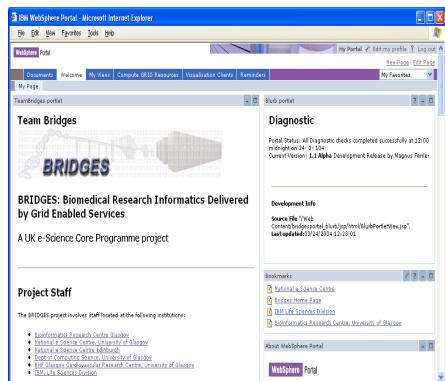
**Next Generation Sequencers ~1TB data**

# Translational Research

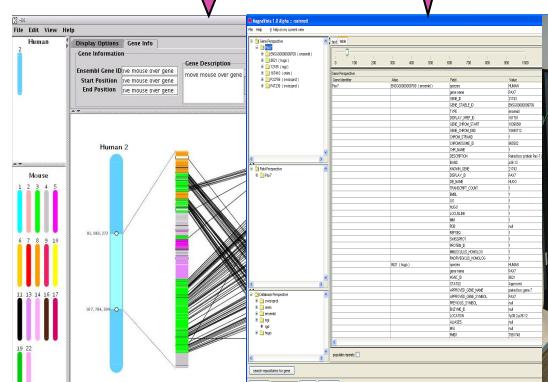


# BRIDGES Project (again)

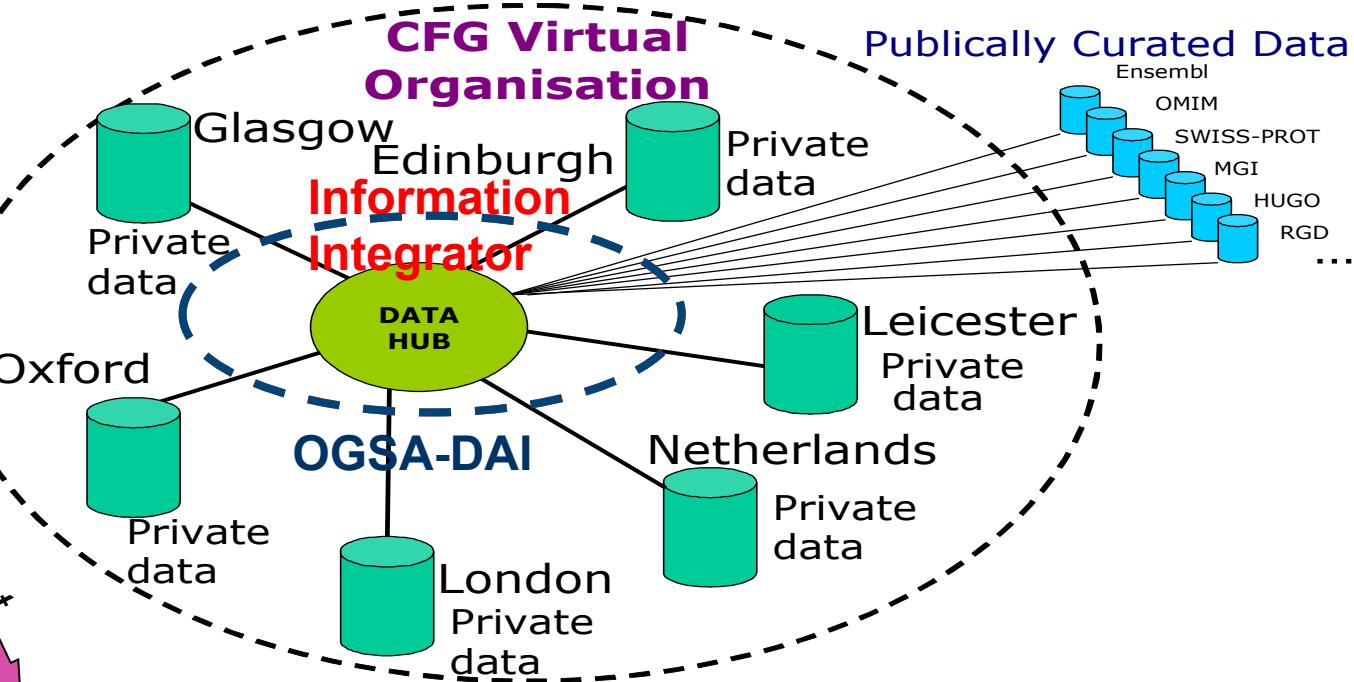
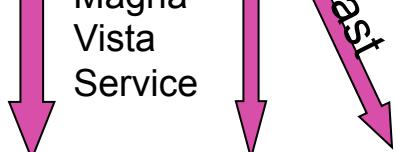
## VO Authorisation



## Synteny Service



## Magna Vista Service



+



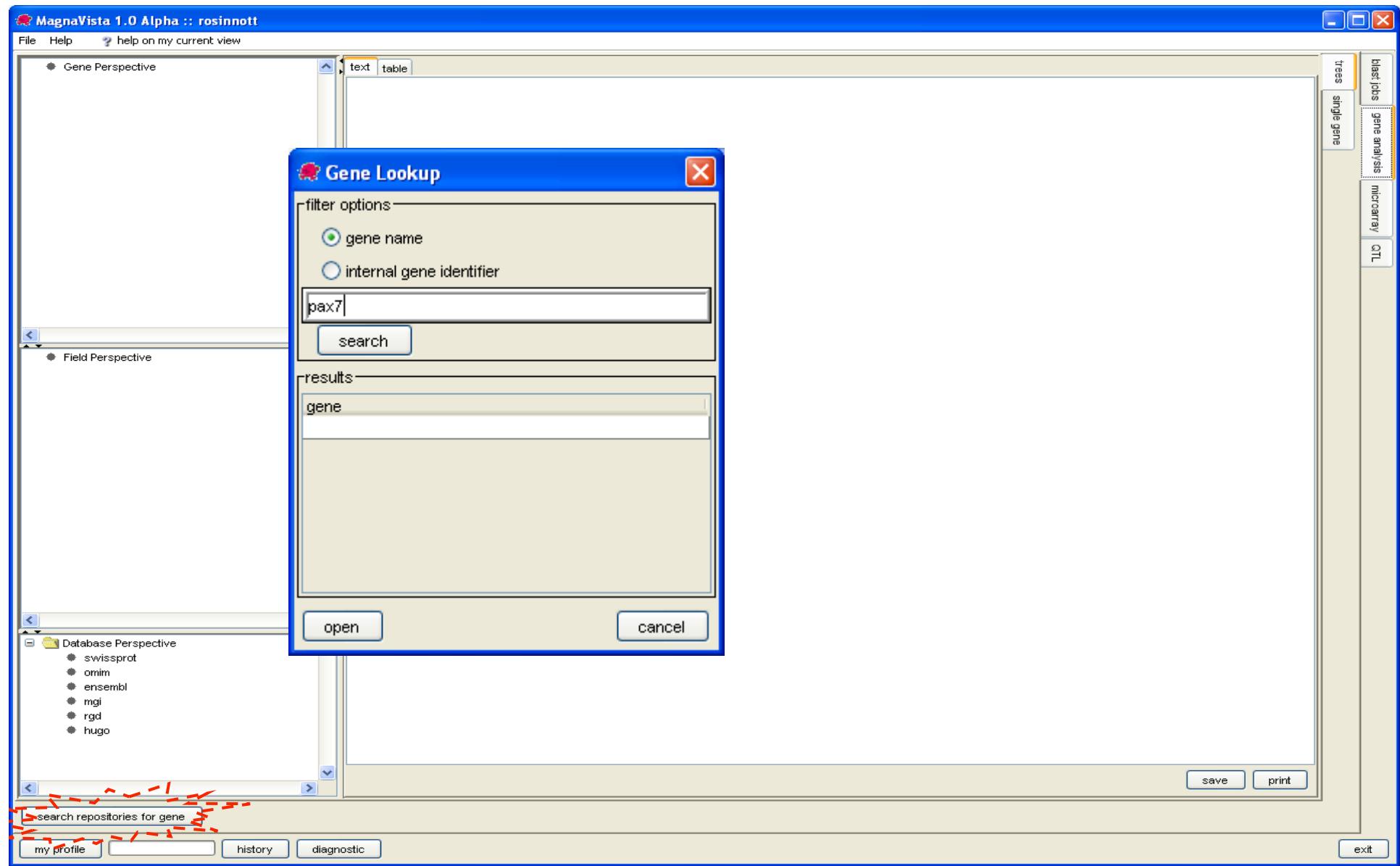
+



+



# 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	genes	databases	perspectives	field cross references	fields	Probe Sets	qtl	preferences			
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"/>

ok

my profile

history diagnostic

blast jobs gene analysis microarray CTL

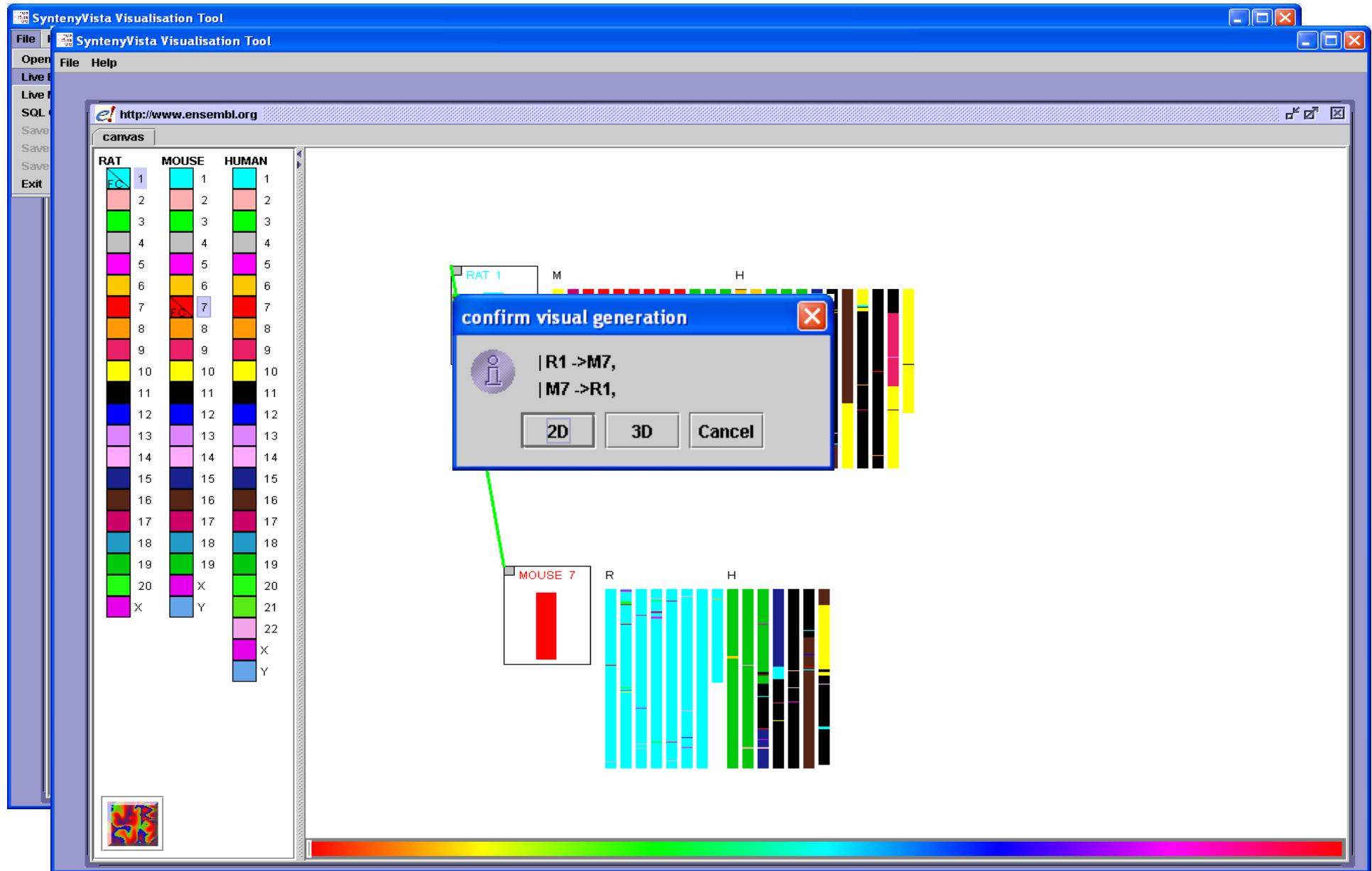
trees single gene

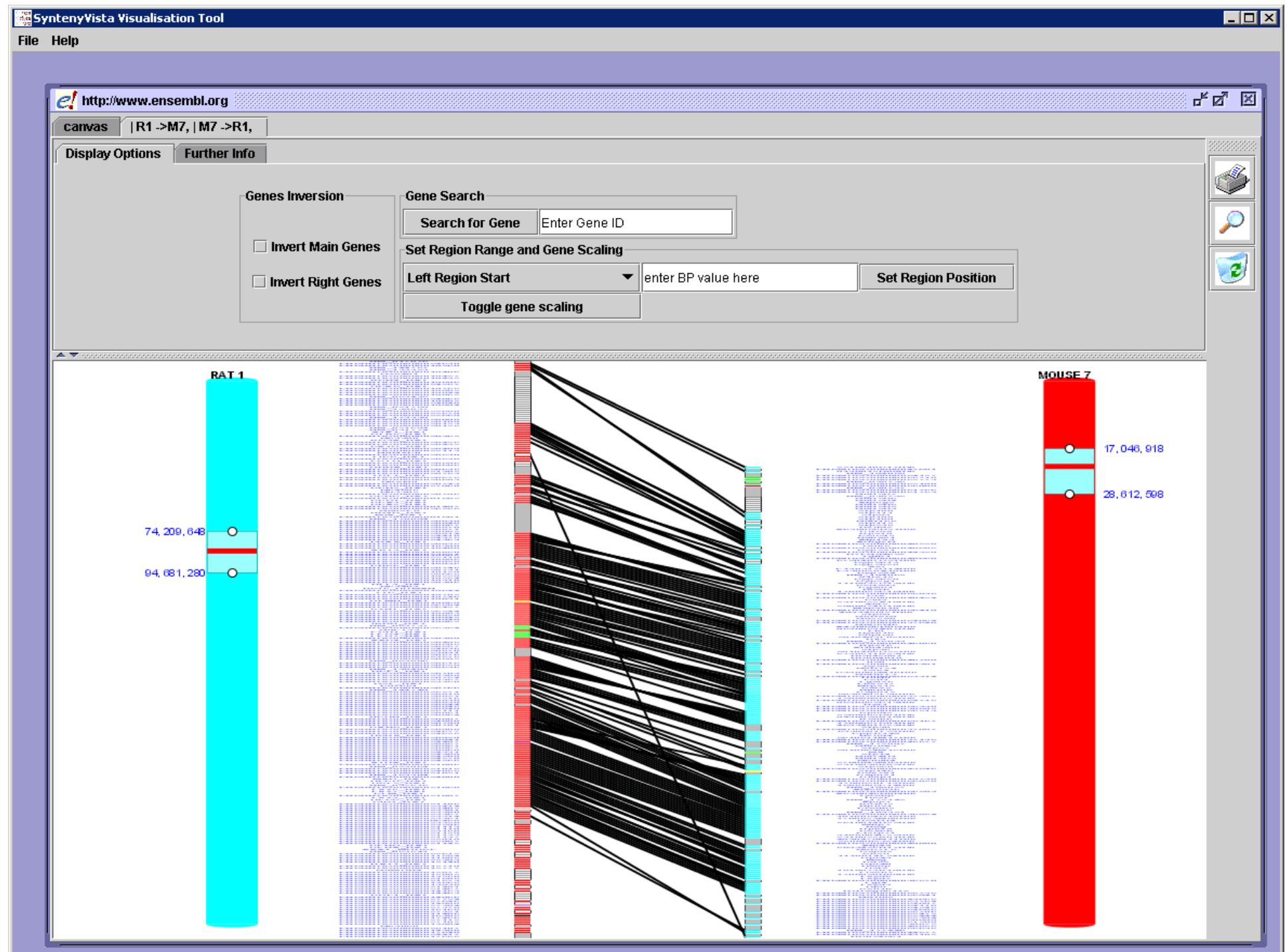
sh

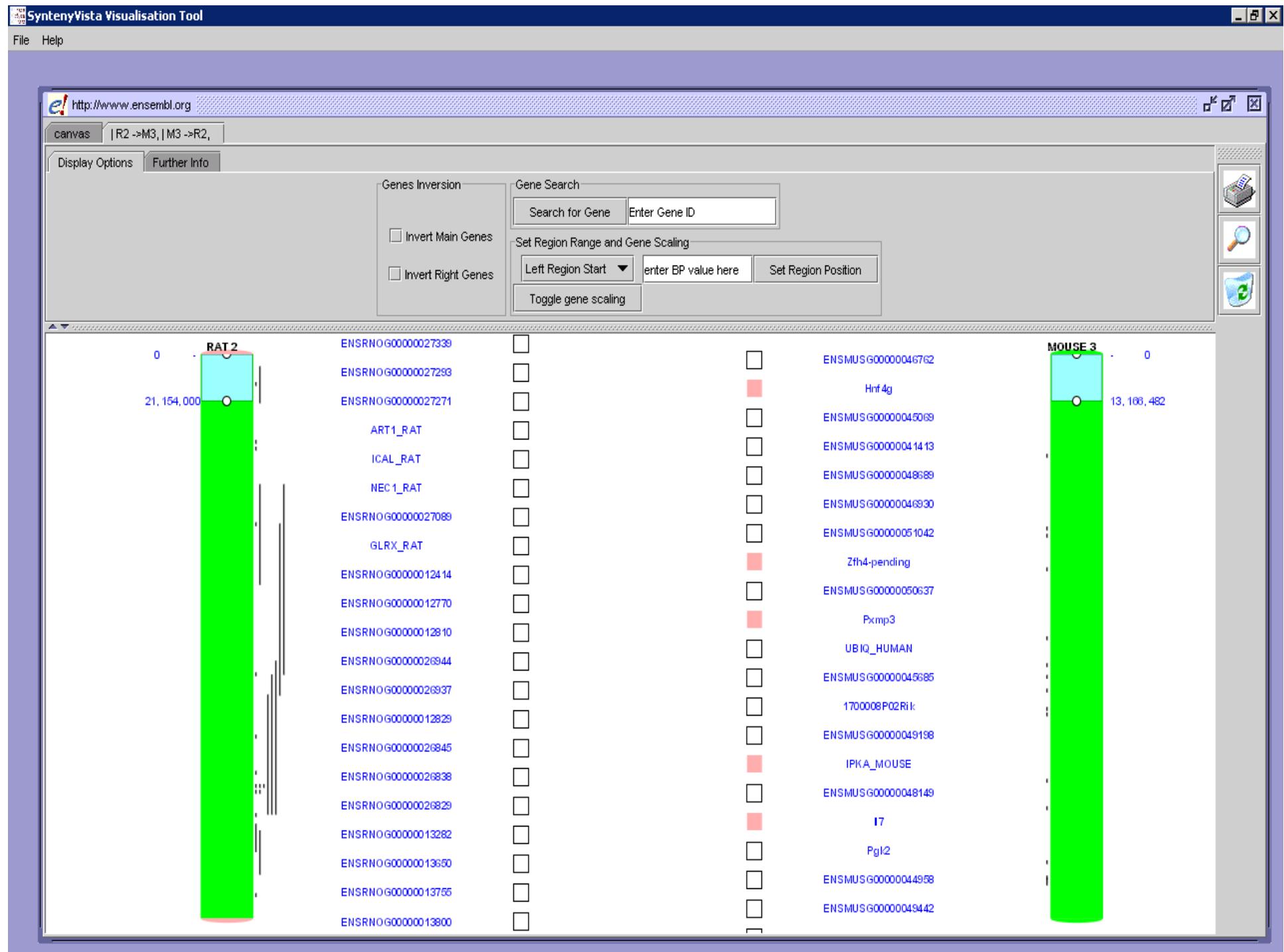
print

exit

# Importance of Data Visualisation







# BREAK

# Social Sciences

- Data, data everywhere
  - Researchers exploring all kinds of areas with societal impact
  - Many major resources distributed around the globe
    - Data Archives
      - literally thousands of surveys/studies crossing society
      - Most countries have their own similar archives
    - Office of National Statistics, Australian Bureau of Statistics, ... and the Census
      - ONS : 1971/81/91/01/11/...
      - ABS : .../01/06/11/16...
    - Geospatial information resources
      - Country profiles, regional profiles, city profiles, street42



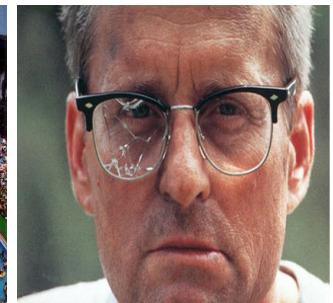
# Australian Urban Research Infrastructure Network (AURIN)

- EIF/NCRIS federally funded project
  - DIISRTE -> Innovation -> Education
  - \$35m+ project
  - [www.aurin.org.au](http://www.aurin.org.au)
    - University of Melbourne are lead agent
- Establishing an e-Infrastructure for Urban and Built Environment Researchers
  - Distributed, (completely!) heterogeneous datasets
  - Data interrogation services
  - Security (unit level data, health data, commercial data!)
  - Online analysis tools
  - Collaboration!!!



# AURIN Context

- Urban and built environment is extremely broad
  - health,
  - transport,
  - future population,
  - liveability,
  - crime,
  - housing,
  - design,
  - ...



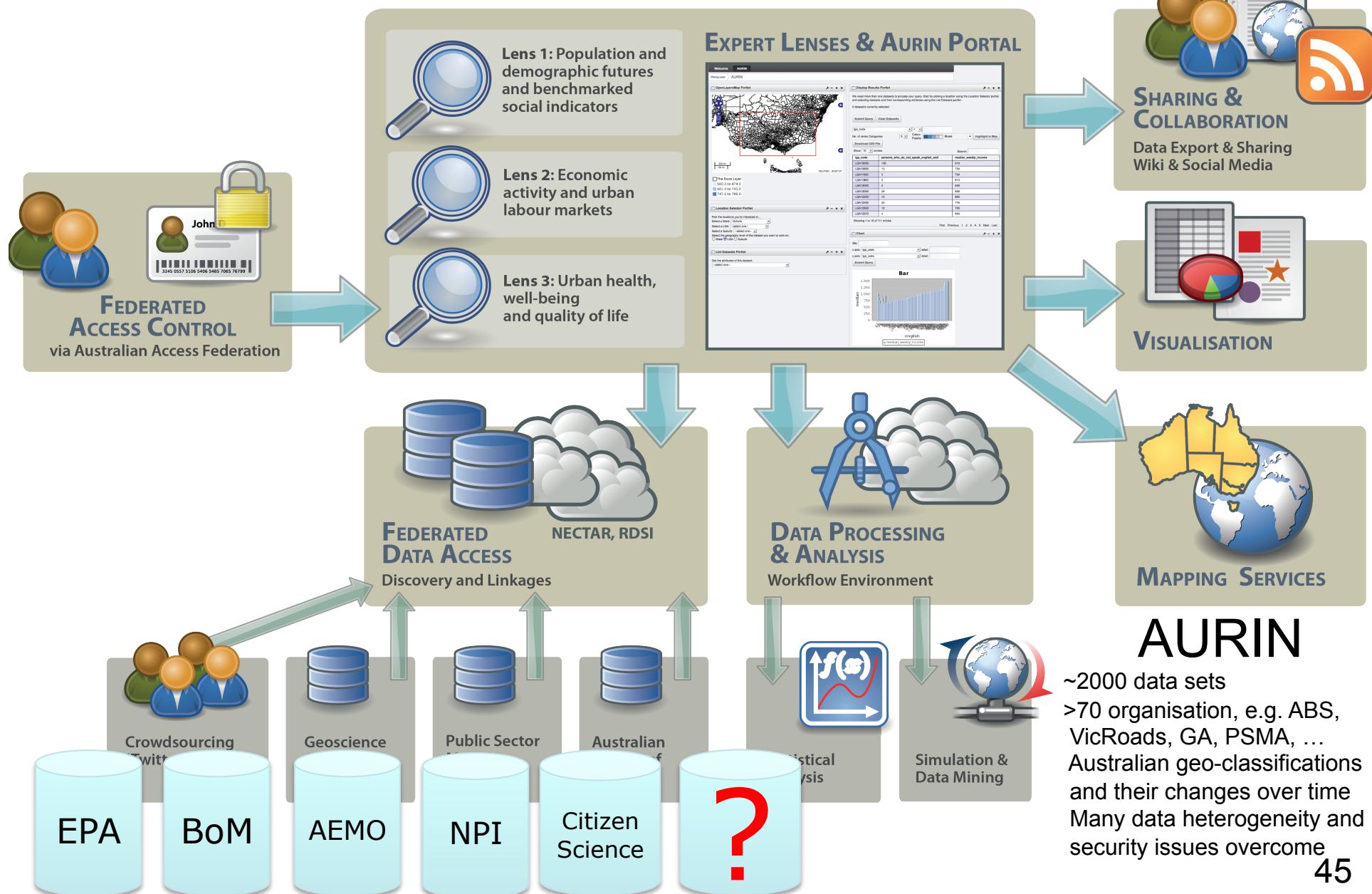
- Much research depends on access to data
  - There is LOTS (and LOTS) of data out there
  - Completely heterogeneous, e.g. geographically
  - ...
  - Data is more often than not silo'd



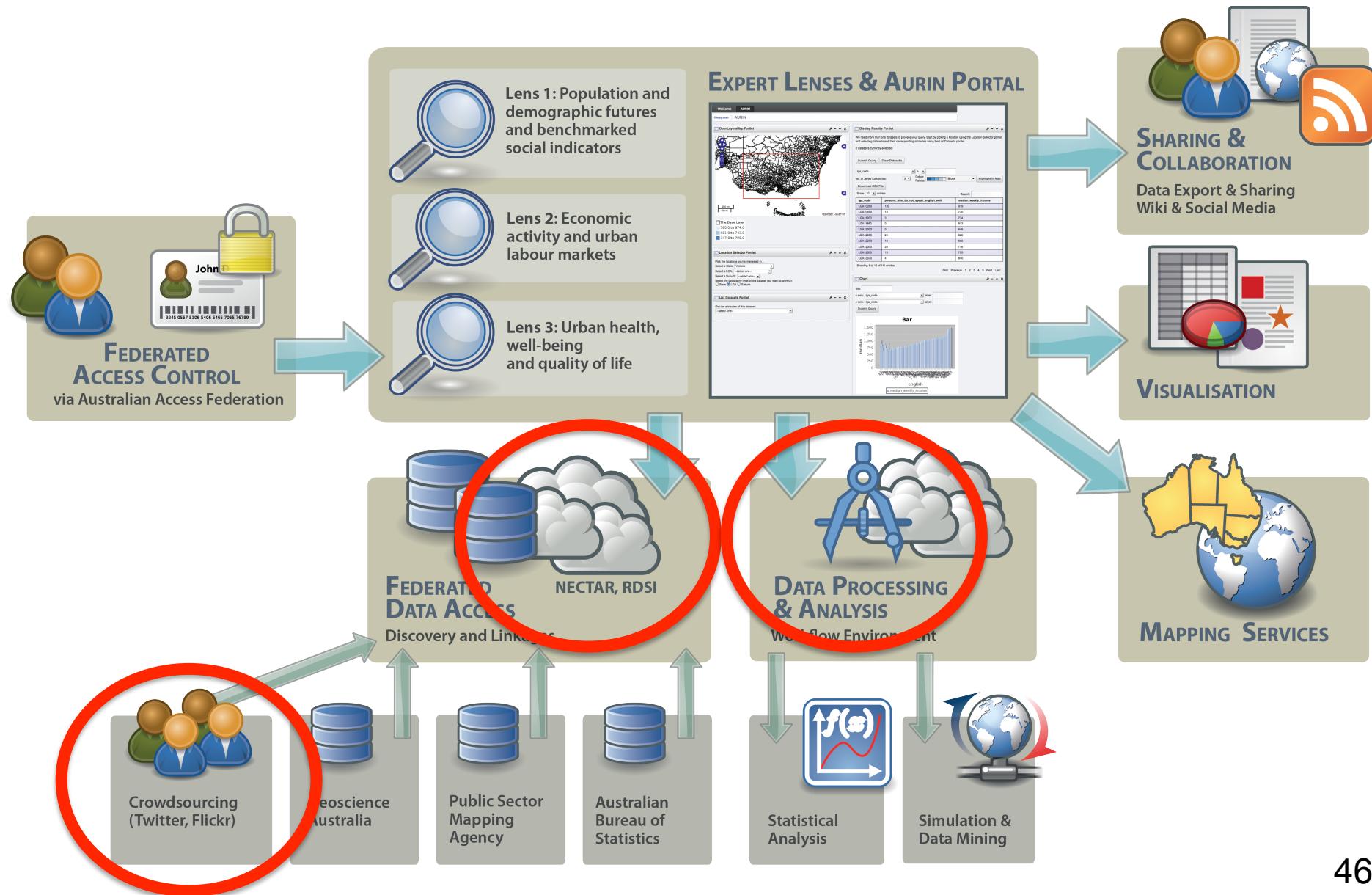
- Requires tools to find, interrogate, analyze and visualize data and enforce good research methodologies  
巩固
  - Consolidate tools and best practice/community know-how!
  - Allow researchers to share results, interact and collaborate
    - No single expert!
- Allow data providers to keep control of their data and its use
  - Authentication and authorisation (and auditing/accounting)



# AURIN Simplified



# AURIN Architectural Components



# Demonstration

(note – Assignment II)

# Clinical Sciences

- Research into:
  - cancer
    - breast, bone, pancreatic, liver...
  - paediatric endocrinology
    - disorders of sex development, adrenal tumours, hormonal imbalances, ...
  - brain
    - neurological disorders, brain trauma, ...
  - ...
- Pretty much every hospital, GP, clinical system I've come across is different

# Data Sharing and Ethics

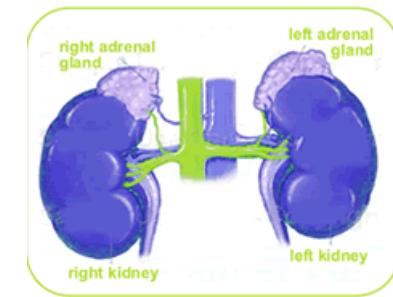
- Clinical data sharing differences globally
  - Ethics
    - Committees / bodies
    - Structures and processes
  - Consent
    - Opt-in
    - Opt-out
  - Data sharing
    - Local only
    - National only
    - Project collaborators only
    - All international researchers
- Need for completely discretionary data sharing models
  - Driven by patient/parental consent
  - Ethics, oversight, ...

Designed with worst case scenarios in mind!

# ENSAT-CANCER



- European Network for Study of Adrenal Tumours (ENSAT – [www.ensat.org](http://www.ensat.org)) CANCER
  - 6m Euro; 5-year project;
    - Started 2011
    - Funds one techie (@UniMelb)
  - Focus on 4 types of 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.)
  - e-'ing
    - Secure web-based tumour databases
    - Biobanking
    - Imaging
    - Clinical trials
    - Seamless interactions between all of these



## Summary

	Associated Study and Registry Distribution							Total
	ACC	Pheo	NAPACA	APA	Total	Principal Investigator	Study Protocols	Study sites/eCRFs
EURINE-ACT	515	7	1678	229	2429	<a href="#">Wiebke Arlt</a>	<a href="#">EURINE-ACT</a>	
Ki-67	36	0	1	0	37	<a href="#">Martin Fassnacht</a>		
Stage III/IV ACC	35	0	0	0	35	<a href="#">Eric Baudin</a>		
PMT	0	229	133	2	364	<a href="#">Graeme Eisenhofer</a>	<a href="#">PMT</a>	<a href="#">PMT</a>
TMA	3	188	0	0	191	<a href="#">Ronald de Krijger</a>		
Long-term PHPGL	0	242	0	0	242	<a href="#">Pierre-Francois Plouin</a>		
AVIS-2	0	0	0	2	2			
PMT3	0	12	0	0	12	<a href="#">Graeme Eisenhofer</a>		
ADIUVO	46	0	0	0	46	<a href="#">Massimo Terzolo</a>	<a href="#">ADIUVO</a>	
ADIUVO Observational	61	0	0	0	61	<a href="#">Massimo Terzolo</a>		
HairCo-2	9	0	1	0	10	<a href="#">Marcus Quinkler</a>		
FIRST-MAPPP	0	41	0	0	41	<a href="#">Eric Baudin</a>	<a href="#">FIRST-MAPPP</a>	
German Cushing Registry	0	1	62	0	63	<a href="#">Martin Reincke</a>		
German Conn Registry	0	0	0	609	609	<a href="#">Martin Reincke</a>		<a href="#">German Conn Registry</a>
CHIRACIC	0	0	3	0	3			
ACC Molecular Marker	35	0	0	0	35			
UK Pheo Audit	0	669	0	0	669	<a href="#">Srirangalingam Umasuthan</a>		
Lysosafe	318	0	0	0	318	<a href="#">Felix Beuschlein</a>	<a href="#">Lysosafe</a>	
FIRMACT	5	0	0	0	5	<a href="#">Martin Fassnacht</a>	<a href="#">FIRMACT</a>	<a href="#">FIRMACT</a>
MAPP-Prono	0	127	0	0	127	<a href="#">Eric Baudin</a>	<a href="#">MAPP-Prono</a>	
MIBG Impact	0	115	0	0	115	<a href="#">Henri Timmers</a>		
Predict Ancillary FIRM-ACT	2	0	0	0	2	<a href="#">Eric Baudin</a>		
FAMIAN	0	0	3	0	3	<a href="#">Stefanie Hahner</a>		<a href="#">FAMIAN</a>

# ENSAT-CANCER

ENSAT - European Network for the study of adrenal Cancers

[https://registry.ensat.org/jsp/search/search\\_result.jsp?dbid=1&dbn=ACC&mainsearch=custom&showformsearch=0](https://registry.ensat.org/jsp/search/search_result.jsp?dbid=1&dbn=ACC&mainsearch=custom&showformsearch=0)



welcome to the ens@t registry

ENSAT Home | ACC Home | ACC Search | ACC Exported Data |  Search | [Select...] | Filter | Welcome, Richard | Sign (14)

## ACC Search Results

There are **16** records matching the following query:

Parameter	Condition
Identification.year_of_birth	Identification.year_of_birth >= 1970 AND Identification.year_of_birth <= 2012 <input checked="" type="checkbox"/>
AND	
Identification.sex	Identification.sex LIKE 'M' <input checked="" type="checkbox"/>
AND	
ACC_DiagnosticProcedures.cushings_syndrome	ACC_DiagnosticProcedures.cushings_syndrome LIKE 'Yes' <input checked="" type="checkbox"/>
AND	
ACC_DiagnosticProcedures.hypertension	ACC_DiagnosticProcedures.hypertension LIKE 'Yes' <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		
FRPA3-48	Eric Baudin (eric.baudin@igr.fr)	25 Feb 2011	05 May 2010	M	1979	Local	<a href="#">Detail</a>	<a href="#">Delete</a>
GYWU-	Martin Fassnacht	25 Oct				National		

### Associated Record Search

- [Biomaterial](#)
- [Chemoembolisation](#)
- [Chemotherapy](#)
- [Follow-up](#)
- [Mitotane](#)
- [Pathology](#)
- [Radiofrequency](#)
- [Radiotherapy](#)
- [Surgery](#)

[ACC Home](#)

# Prospective Monoamine-Producing Tumor (PMT) study

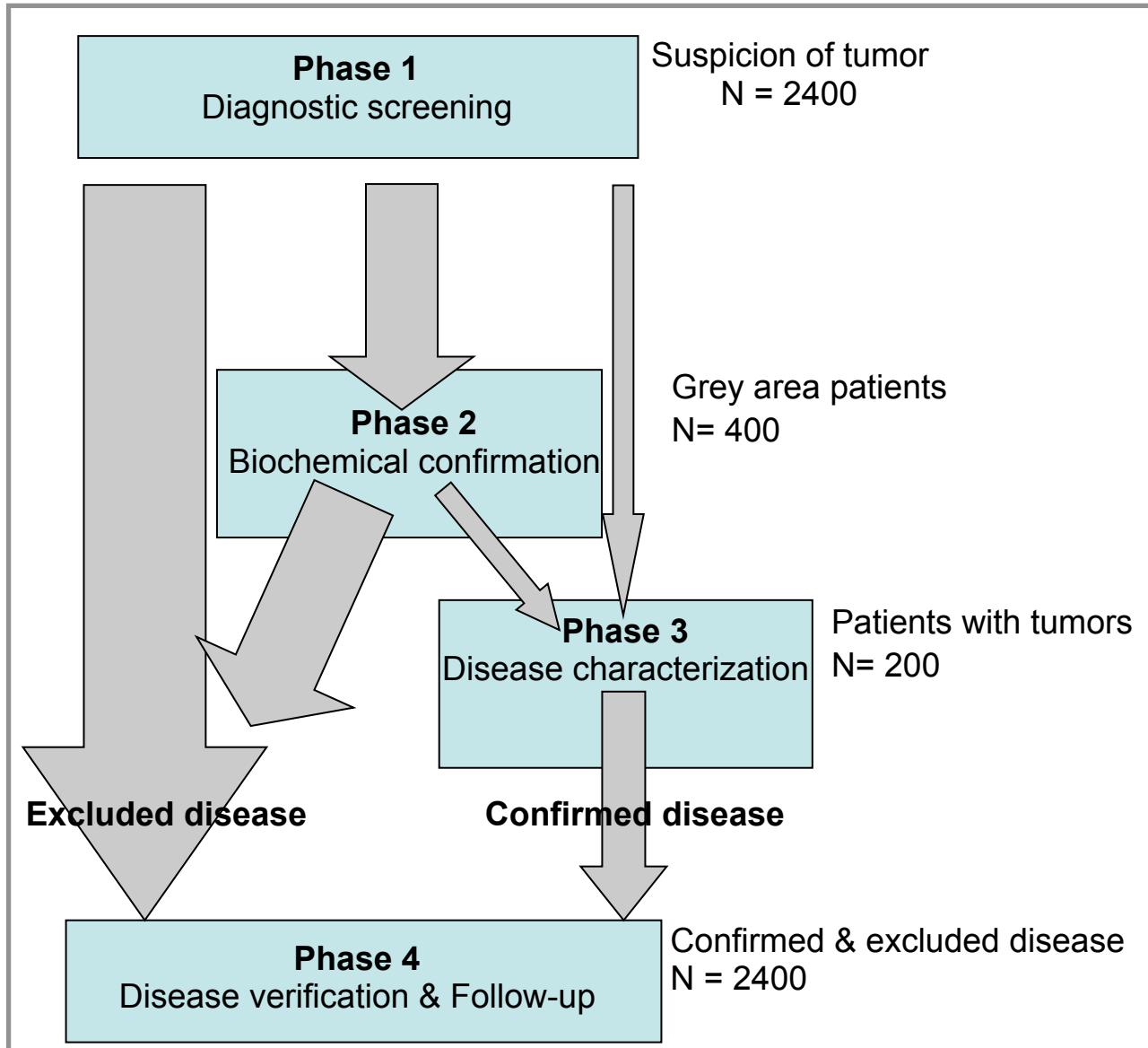
## Aims:

### 生物标记

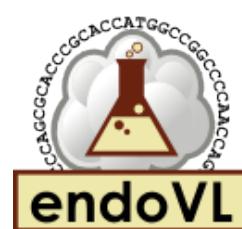
- Diagnostic biomarkers - miRNA & utility of plasma and urinary free metanephrenes
- Follow-up testing - effective strategy for follow-up testing
- Genotype-phenotype relationships
- Characterize catecholamine metabolomic & secretory phenotypes in relation to cardiovascular and other manifestations of disease
- LC-MS/MS based metabolomic profiling of bioenergetic signatures
- Novel biomarkers of malignancy

## Contact for patient enrolment:

Dr. Roland Därr  
Roland.Daerr@uniklinikum-dresden.de



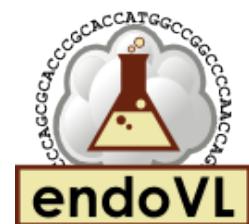
# Endocrine Genomics VL



- Recently funded NeCTAR VL project
  - \$1.03m (ran for 1 year!)
  - Platform for range of endocrine related disorders with associated - omics analytics
  - Includes
    - Type-1 & Type-2 Diabetes (paediatric and adult) and rare forms of diabetes related disorders
      - Working with APEG, JDRF & ADS
    - Thyroid-related and severe obesity-related disorders
      - Working with Australian Thyroid Foundation
    - Neuroendocrine and adrenal disorders
      - Supporting clinical oncology society of Australia
    - Bone related disorders
      - Supporting Australia and New Zealand Bone and Mineral Society and Osteoporosis Australia
    - Disorders of Sex Development
      - Supporting the Australasian DSD Network



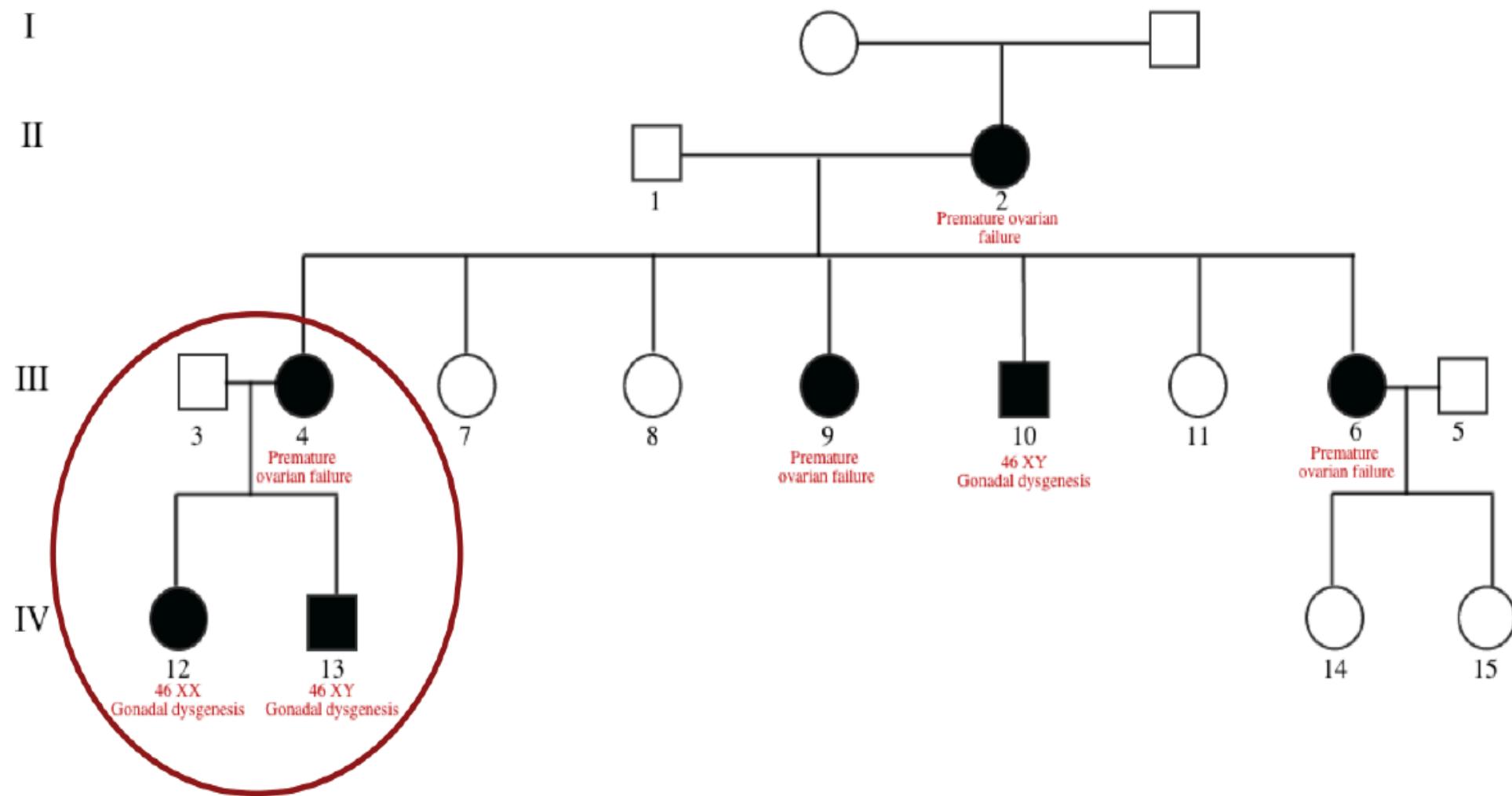
# Endocrine Genomics VL



- Beyond ENSAT (registry/trials)
  - Integration of workflows for –omics data analysis
    - Whole genome/exome sequencing
    - Sharing how analytics is undertaken (repeatable data analysis pipelines)
    - Leveraging Australia-wide expertise in this space
      - Life Science Computation Centre (Victoria)
      - Australia Genetics Research Foundation
      - Bioplatforms Australia
      - National ICT Australia
  - Multiple –omics groups offering multiple omics analysis
  - Comparison of resulting analysed data in integrated system
  - Use of Cloud for data and analytics
    - ...and overcoming security challenges this gives rise to!



# Case Study (Blind) DSD -omics analysis



# Cloud-based -omics Analytics Pipelines



# Integration of Phenotype & Genotype

EndoVL Home

EndoVL Home Portal DSD

### Assessments

Assessment Date	Phallus Length (mm)	Phallus Size	Urinary Meatus Site	Labioscrotal Fusion	Right Gonad	Left Gonad	Mullerian Structures	Wolffian Structures	EMS
19/08/2013	Within the reference range for male	Normal	Yes	Labioscrotal	Labioscrotal	No	Yes	Yes	12.0
19/08/2013	Within the reference range for male	Normal	Yes	Labioscrotal	Labioscrotal	No	Yes	Yes	12.0

### Gene Analysis

Screening data last updated on: 10/12/2013 14:50  
Template: DNA - Technique: ?

– select screening –

 View in GenomeBrowser

Chr hg19	Position Start/End	Allele	Type	Variant DNA	Variant DBID
NC_000001.10	12783..12783	Parent #2	subst	g.12783G>A	DDX11L1_000002
NC_000001.10	13116..13116	Parent #2	subst	g.13116T>G	DDX11L1_000003

# WGS Visualisation

Genomes    Genome Browser    Tools    Mirrors    Downloads    My Data    About Us    View    Help

## UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base Zoom out 1.5x 3x 10x

chr1:12,733-112,783 100,051 bp. enter position, gene symbol or search terms go

chr1 (p36.33) | p31.1 | 1q12 | q41 | 4344

Scale chr1: 20,000 30,000 40,000 50,000 60,000 70,000 80,000 90,000 100,000 110,000 hg19 variants29.bbged g.12783G>A | g.13116T>G | g.13118A>G | g.13868A>G | g.14610T>C | g.14653C>T | g.14717G>A | g.14932G>T | g.15190G>A | g.15274A>T | g.15820G>T | g.15985dup | g.16495G>C | g.16841G>T | g.16957G>T | g.16977G>A | g.17538C>A | g.17697G>C |

UCSC Genes (RefSeq, GenBank, CCDS, RFam, tRNAs & Comparative Genomics)

WASH7P WASH7P WASH7P WASH7P WASH7P WASH7P WASH7P WASH7P

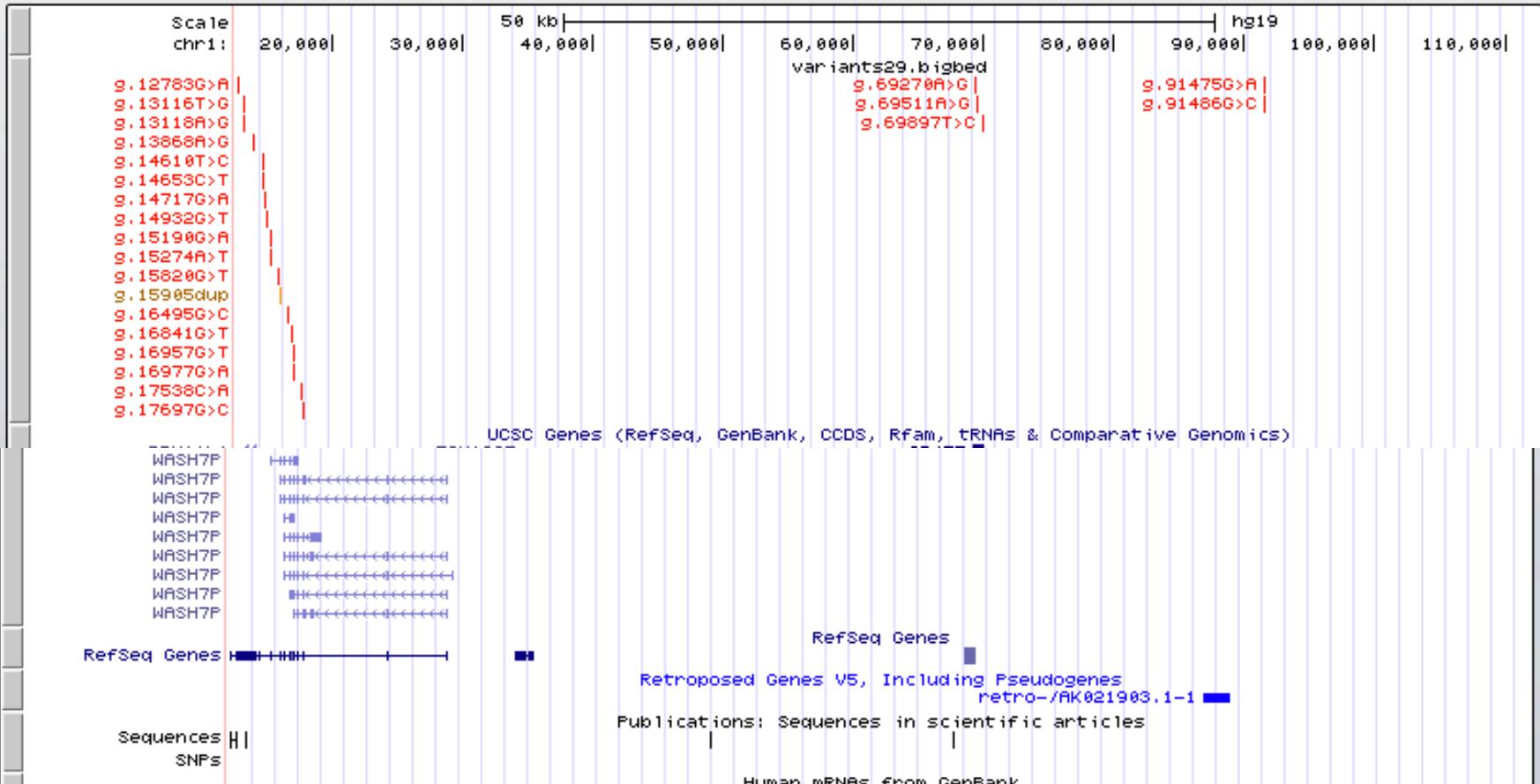
RefSeq Genes

Retroposed Genes V5, Including Pseudogenes retro-/AK021903.1-1

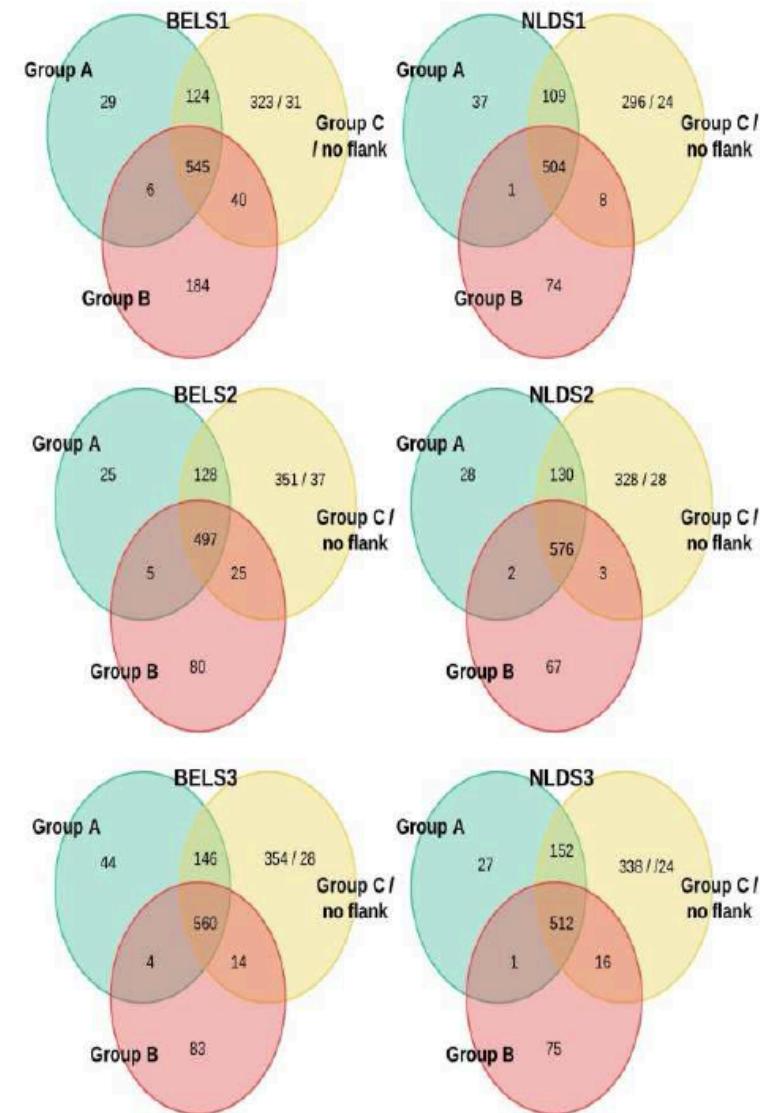
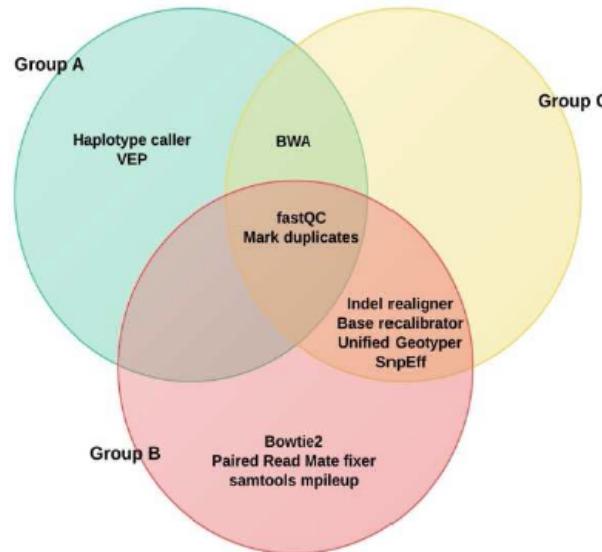
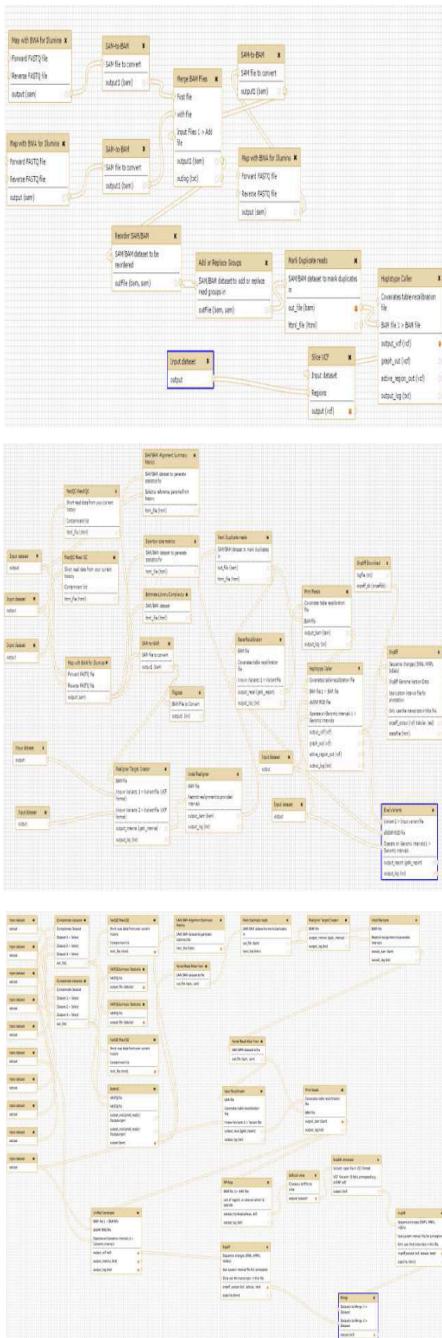
Publications: Sequences in scientific articles

Human mRNAs from GenBank

Sequences SNPs



# And the winner is...



# Conclusions

- All application domains facing data “issues”
- How do we build e-Infrastructures that support this?
- Specialised research specific Google-farms?
  - Yes but life sciences
    - EBI,
    - SANGER,
    - NCBI,
    - ...
  - The diversity
    - researchers interested in rat, mouse, human, barley, jellyfish, cancer, diabetes, Tasmanian devil, ...
    - Sonic Hedgehog, ARSE, 18wheeler, pray for elves, pokemon, werewolf, BRCA1, BRCA2, ...

# Conclusions...ctd

- Importance of community driven standards
  - But forever a moving target
    - Ontologies
    - Communities often not aligned/mature
    - ...
- The toothbrush culture!!!
  - Big Pharma/Grants...
- Security!!!
  - User-driven
  - Provider-driven
  - Virtual Organisations
    - Simple enough?
- Technology is still often specialised/hard
  - But it can be achieved ...if you know how! ;o)

# Questions ...?