

Domain Drivers – Lecture 3-4

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
 - 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...
 - High Energy Physics
 - Astrophysics
 - Electronics
 - Bioinformatics
 - Biomedical and bioinformatics domain
- BREAK
 - US Intelligence
 - 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

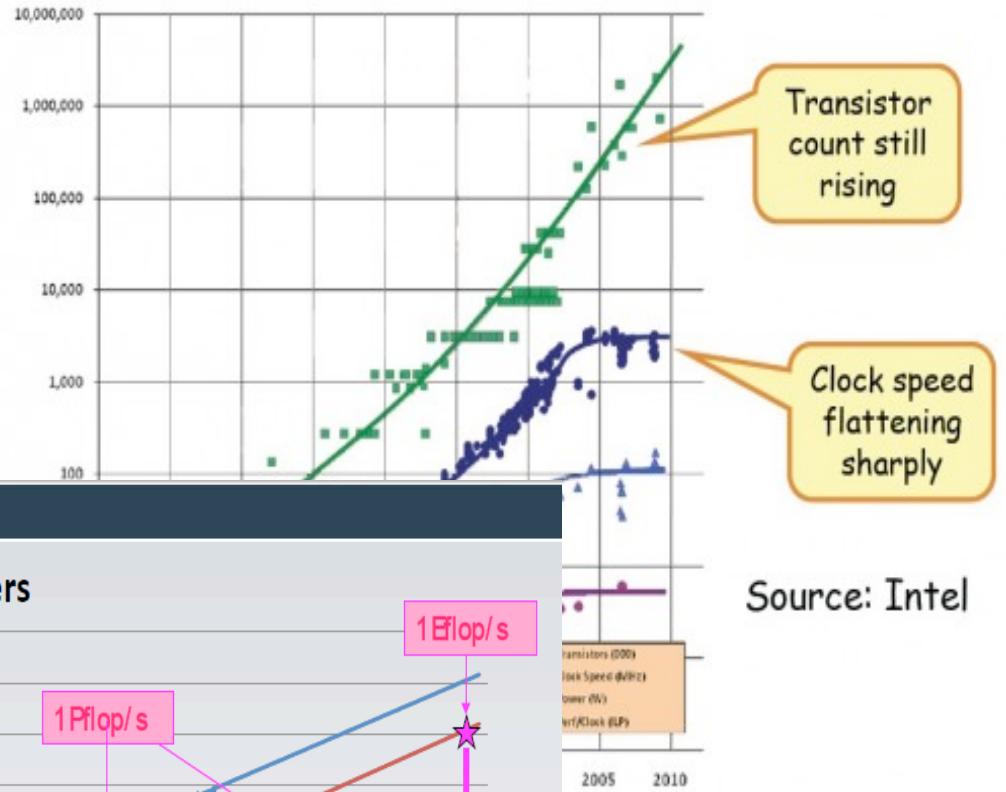
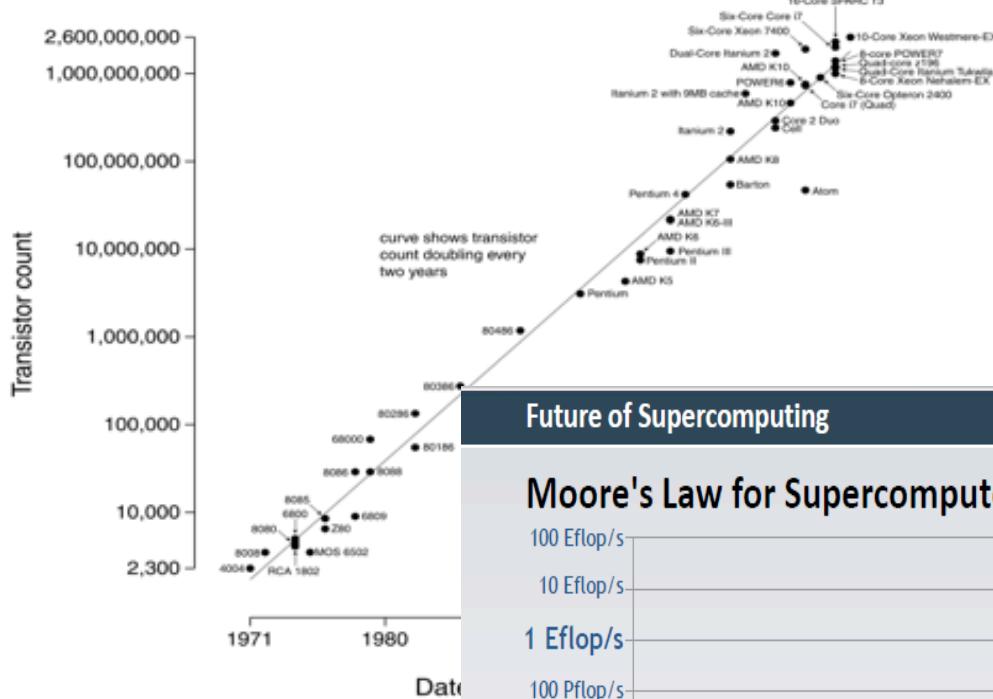
Project Portfolio

Subset of On-Going

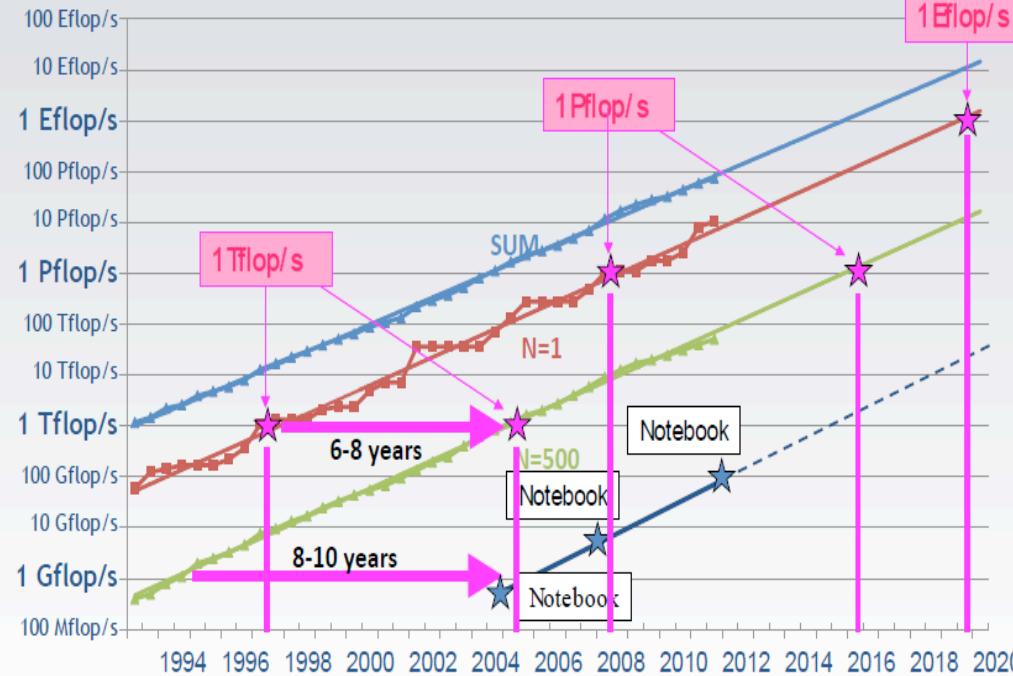
- EU European Platform for Study of Wolfram, Alstrom, Bardet Biedl (EuroWABB)
- Multicenter prospective study of biochemical profiles of monoamine-producing tumors (PMT Study)
- European Society of Hypertension Study on Pheo/PGL
- International DSD
- EU FW7 European Network for Study of Adrenal Tumors Cancer Research Platform (ENSAT-CANCER)
- VicHealth Health Indicators and Spatial Objective Data
- National Spinal Injury Research Platform
- Australian Urban Research Infrastructure Network (AURIN)
- 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
- 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

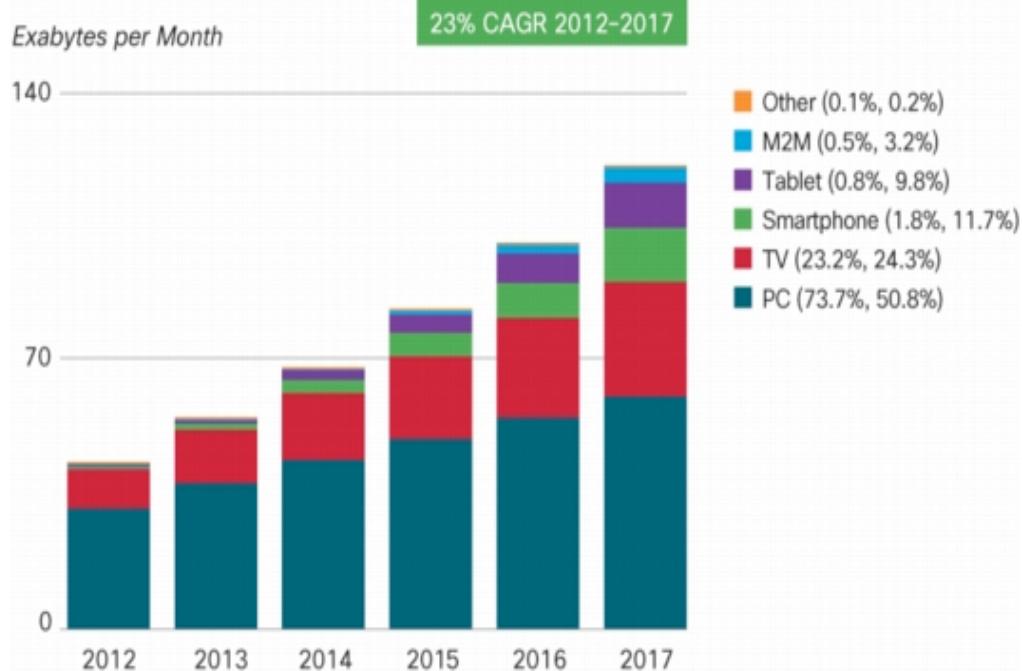
Microprocessor Transistor Counts 1971-2011 & Moore's Law



Moore's Law for Supercomputers



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

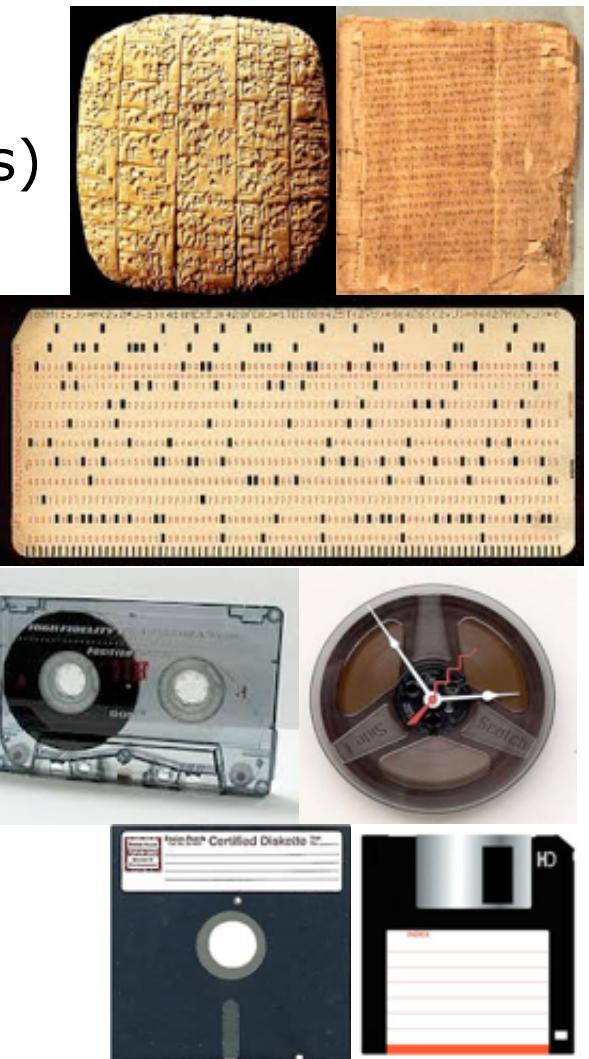
Table 6. Table A-1 Global IP Traffic, 2012-2017

	2012	2013	2014	2015	2016	2017	CAGR 2012-2017
IP Traffic, 2011-2016							
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%
By Type (PB per Month)							
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%
By Segment (PB per Month)							
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%
Total (PB per Month)							
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 circa 1960
 - punch cards,
 - punched streamer tape,
 - magnetic tape,
 - floppies,
 - ...
- ~RIP!

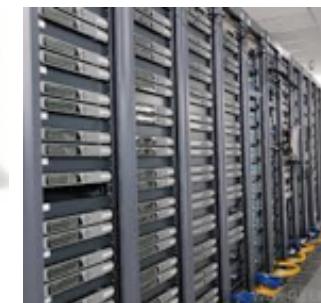


Data Present

- Data Storage today

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

- Dropbox
- Google
- Clouds
- ...



Data Deluge

- The combination of mobile devices and the Internet of Things will result in an estimated 100 billion connected devices by 2020.



at

(March 2007)
digital Universe

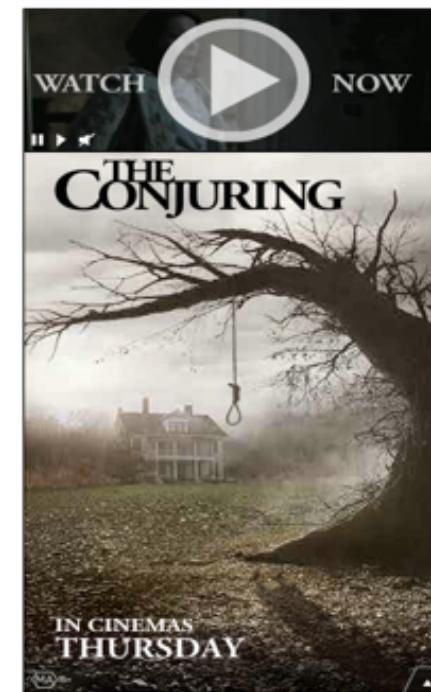
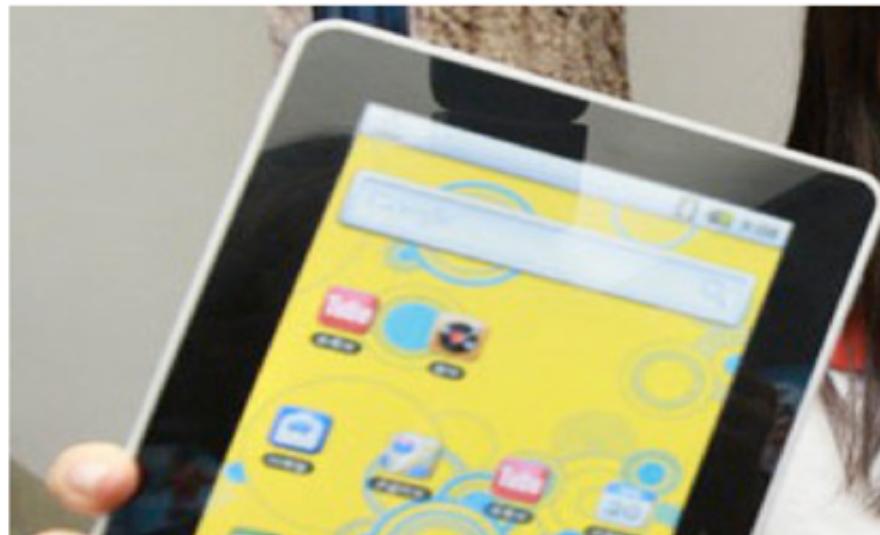
- The total amount of data created in 2012 is approximately 2.8 zettabytes.



2012.

ta Corporation

- By 2015, the amount of data created daily will be 40 times greater than it was in 2010.



YouTube,
is 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.

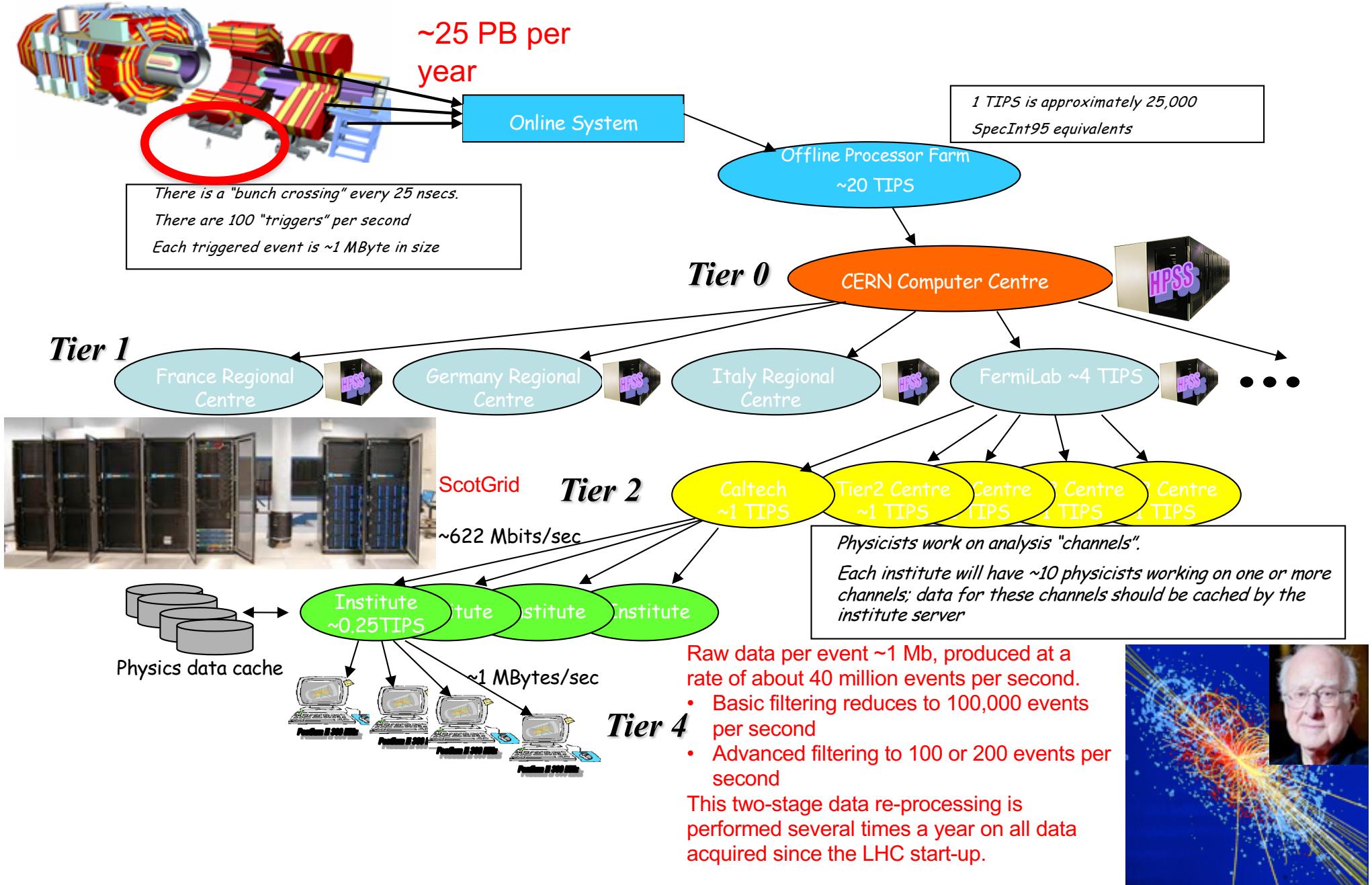
"I spend most of my time assuming the world is not ready for the technology revolution that will be happening to them soon," Schmidt said.

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.

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

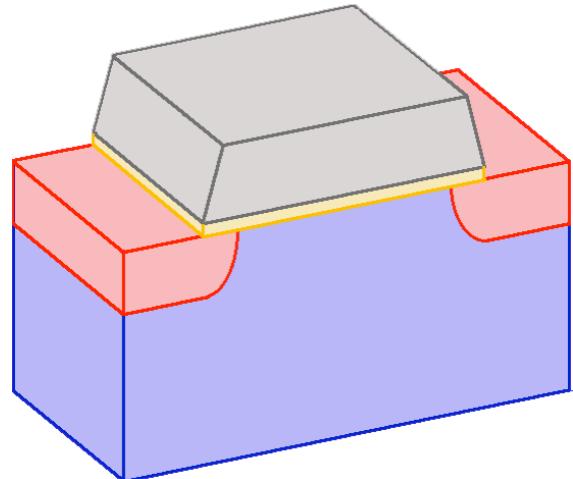
11 PDRAs (7 + 4)
7 PhD

Industrial partners

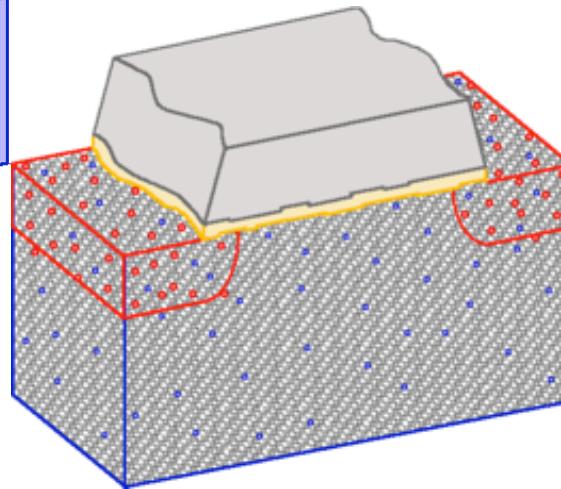


University partners

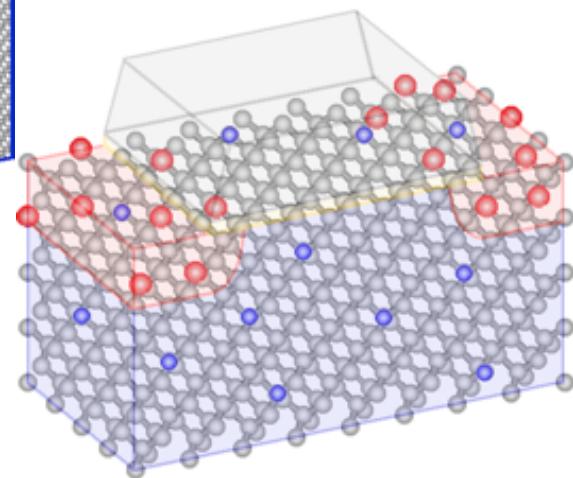
Semiconductor device variability



Historic simulation paradigm

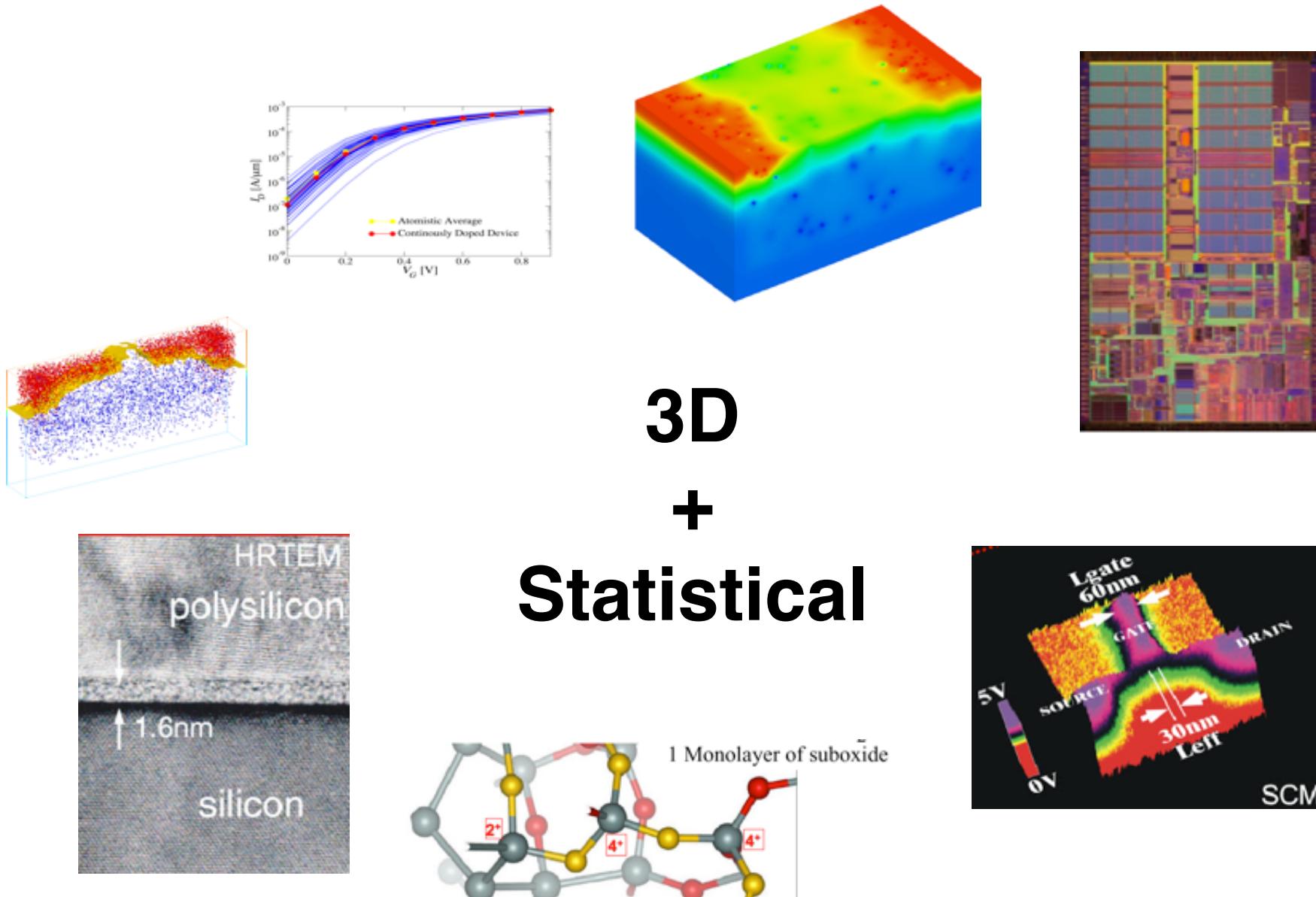


A 22 nm MOSFET
In production 2011

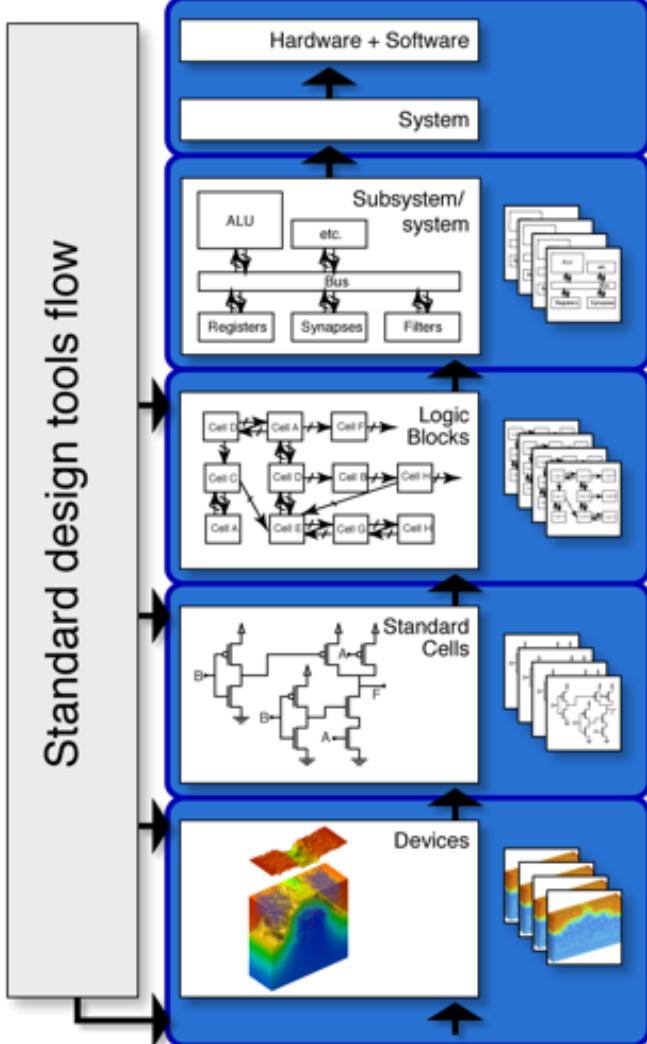


A 4.2 nm MOSFET
In production 2023

Challenges of NanoCMOS Design



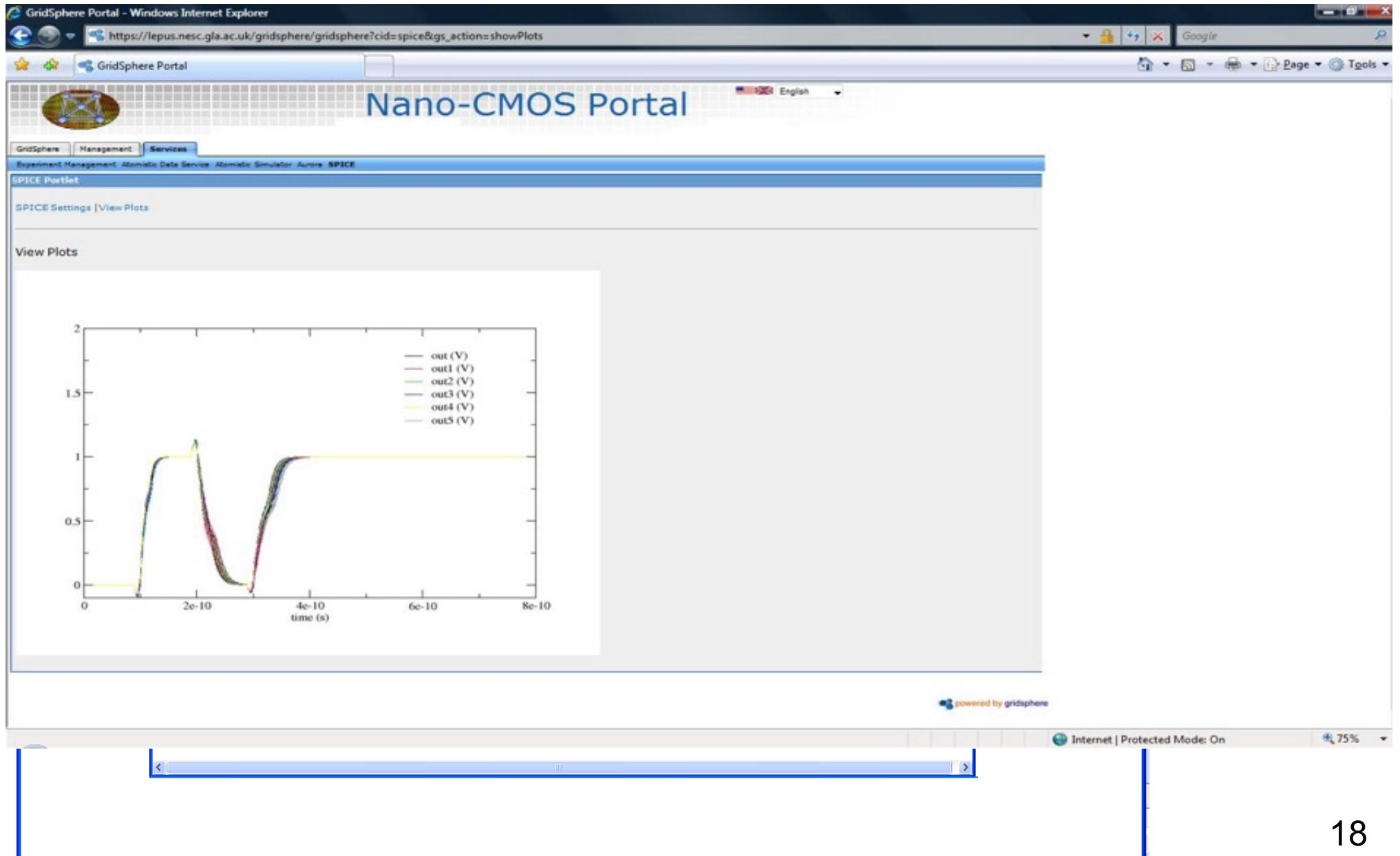
Challenges Hierarchical statistical system simulations



- ❑ 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 t

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

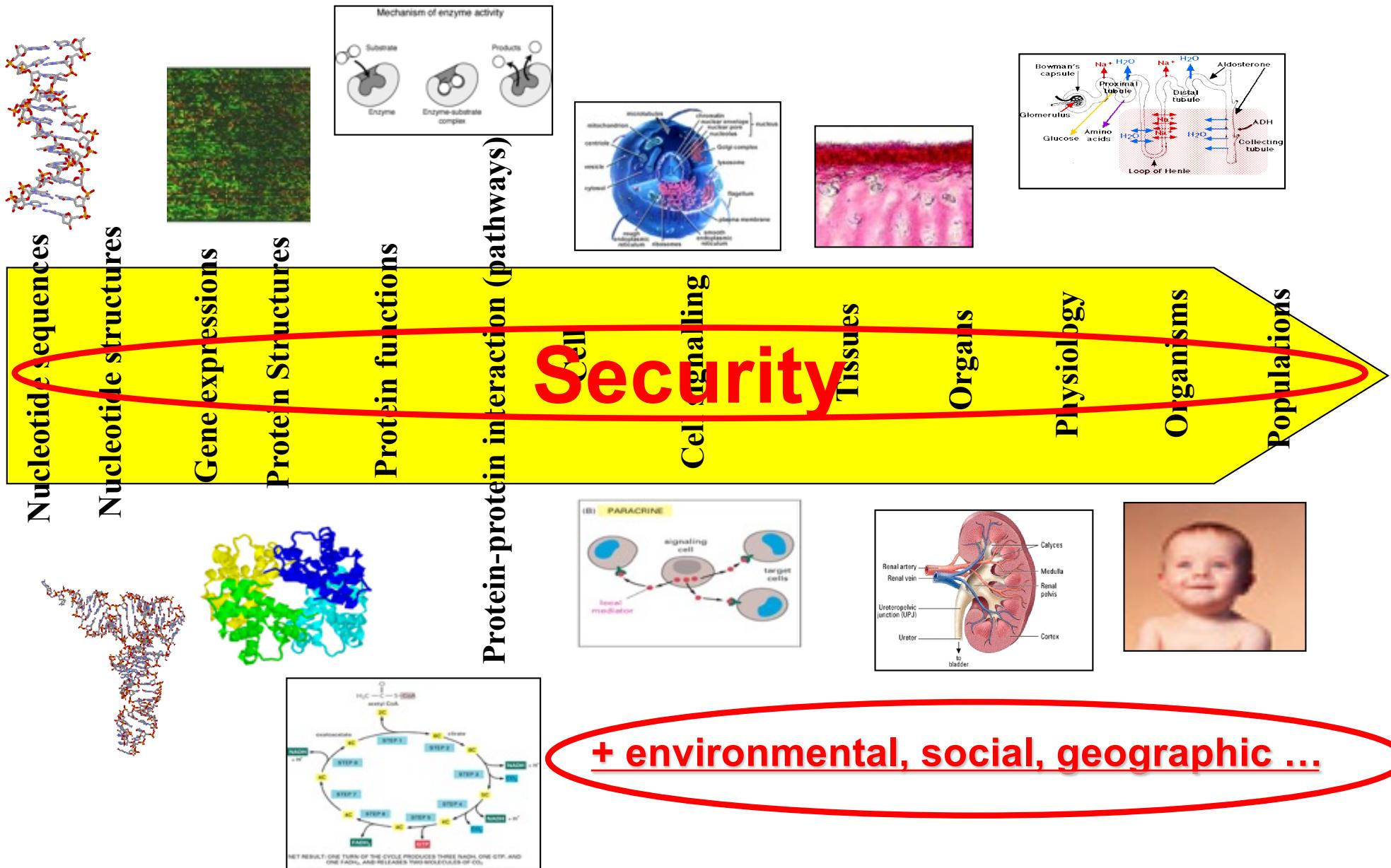
The -g flag!!!

The screenshot illustrates a hybrid approach to managing distributed data. On the left, a Mozilla Firefox window displays a Flash-based application for managing file records, showing a sidebar with 'asks', 'Jobs', and 'File Records' options. The main content area shows 'File Record Details' with 'Core Data' and a table of name-value pairs. A large, stylized 3D watermark reading 'Data vs Metadata' is overlaid across the center of the application. On the right, a terminal window shows the XML representation of the metadata captured by the '-g' flag, including details like created time, date, host, and various identifiers for the job submission.

```
<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/nanocomstest/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/nanocomstest/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/nanocomstest/task/items/16d68fffc-1027-45fa-a7c9-8627d3695a75/" />
```

Done nanodata.vidar.ngs.manchester.ac.uk

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 microbiome data. The data is presented in a JSON-like structure with many nested objects, representing a complex and distributed dataset. The middle section shows a woman holding a newborn baby. A small, purple, handheld device with a screen and buttons is attached to the baby's chest, likely a glucose monitor. The bottom section is a screenshot of the ENDIA (Environmental Determinants in Islet Autoimmunity) website. It features a logo of a stylized flower inside a circle. The page includes navigation links for Home, About Us, What's Involved, Regional Program, Health Professionals, Researcher Resources, News, and Contact. A prominent blue header banner asks, "Why are more children getting Type 1 Diabetes?". Below this, a message states, "Please note that recruitment has closed". A "Contact the Team" button is located at the bottom left of the banner.

Participant Login Staff Login

environmental determinants of islet autoimmunity

Why are more children getting Type 1 Diabetes?

Please note that recruitment has closed

Contact the Team

Home About Us What's Involved Regional Program Health Professionals Researcher Resources News Contact

Participant Login Staff Login

environmental determinants of islet autoimmunity

Why are more children getting Type 1 Diabetes?

Please note that recruitment has closed

Contact the Team

Messy

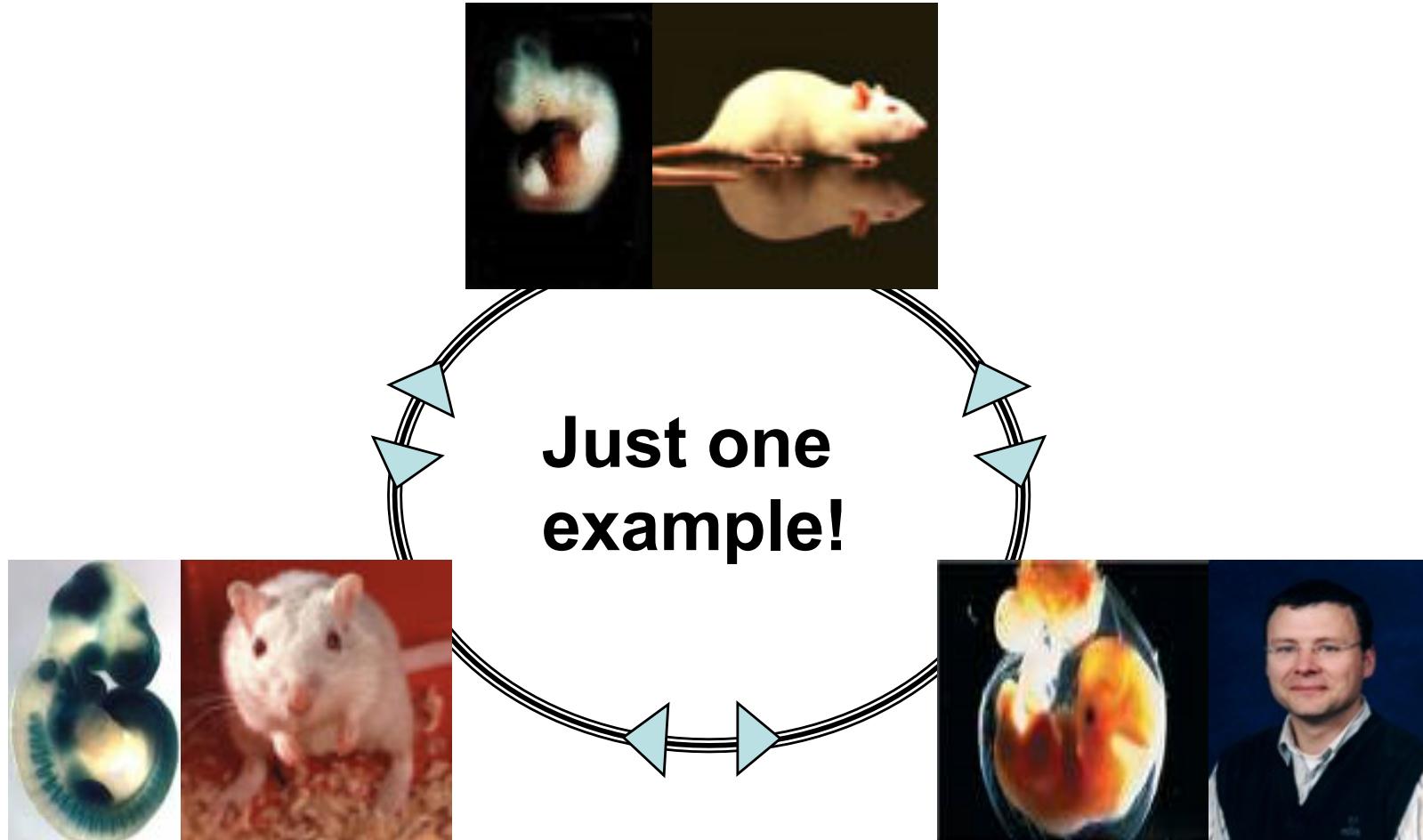
GPL96 - Notepad

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 NM_0009
ing factor 2, 45kDa (ILF2), mRNA 24234746 NM_004515 ILF2 0200053_at nucleophila, member A, mRNA (c
SCC3L1 0200059_s_at ras homolog gene family, member A Homo sapiens nucleophila, member A, mRNA (c
2288 AF275719 HSP90AB1 0200065_s_at ADP-ribosyltransferase 1, clone 24537 ADP-ri
eterogeneous nuclear ribonucleoprotein M Homo sapiens nucleophila, member A, mRNA (c
RNA (cDNA clone MGC:4498 IMAGE:2964510), complete cds 33873259 0200079_s_at lysyl-tRNA synt
014267 C1orf58 0200085_s_at alpha (FNTA), transcript variant 1 Homo sapiens nucleophila, member A, mRNA (c
rase, CAAX box, alpha (FNTA), transcript variant 1 Homo sapiens proteasome (prosome, macropain) subunit, beta type
omo sapiens ATPase, H⁺ transporting, alpha (FNTA), transcript variant 1 NM_003945 ATP6VOE
U 0200594_x_at nucleophila, member A, transcript variant 1 Homo sapiens he
, mRNA 4507676 NM_003299 HSP90AA1 0200594_x_at nucleophila, member A, transcript variant 1 Homo sapiens he
ein kinase, cAMP-dependent, regulat Homo sapiens nucleophila, member A, transcript variant 3, mRNA 47132579//4713
006265 RAD21 0200600_s_at SET domain repeat domain 1 (WDR1), transcript variant 1, m
rotein complex 2, beta 1 subunit (AP2B1), transcript variant 1 NM_001030006//NM_001282 AP2B1
n 3 (phosphorylase kinase, delta) Homo sapiens nucleophila, member A, transcript variant 1 NM_0051
mRNA 47419913 NM_004184 WARS 0200605_s_at SET translocation (myeloid leukemia-associated) NM_002840 PTPRF
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
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 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
nre 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

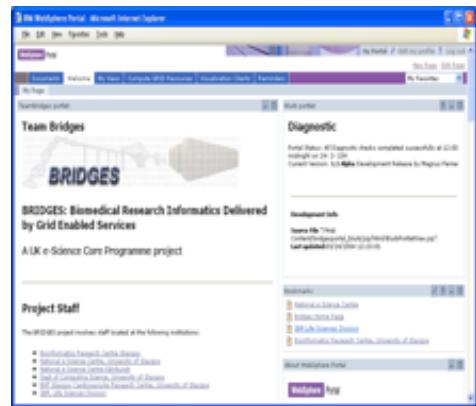
**Next Generation Sequencers
+ 100 TB data**

Translational Research

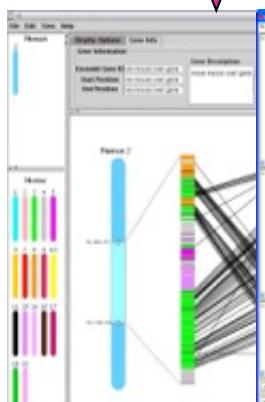


BRIDGES Project

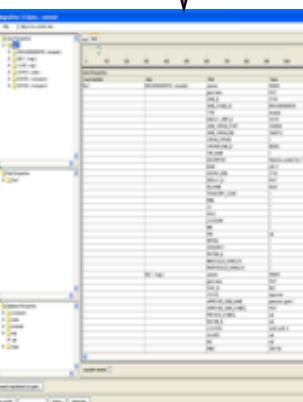
VO Authorisation



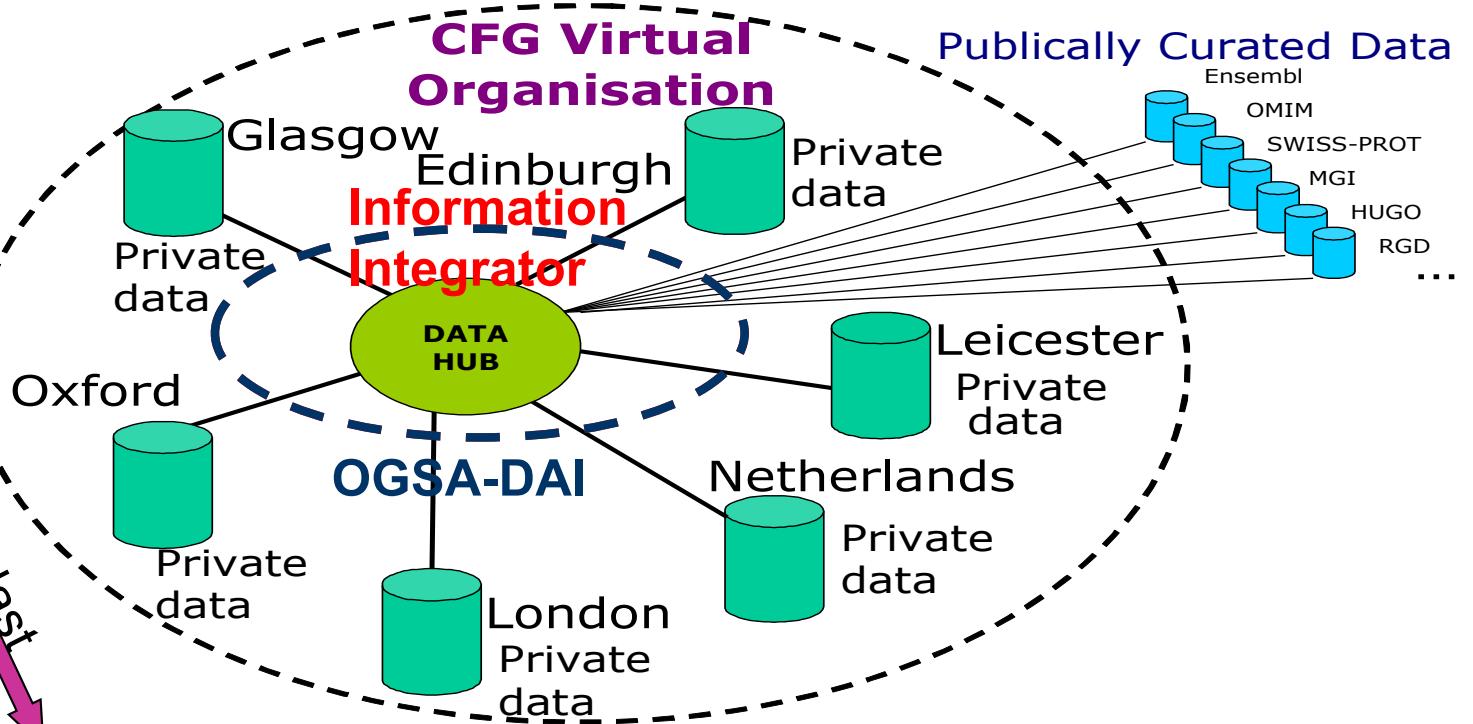
Synteny Service



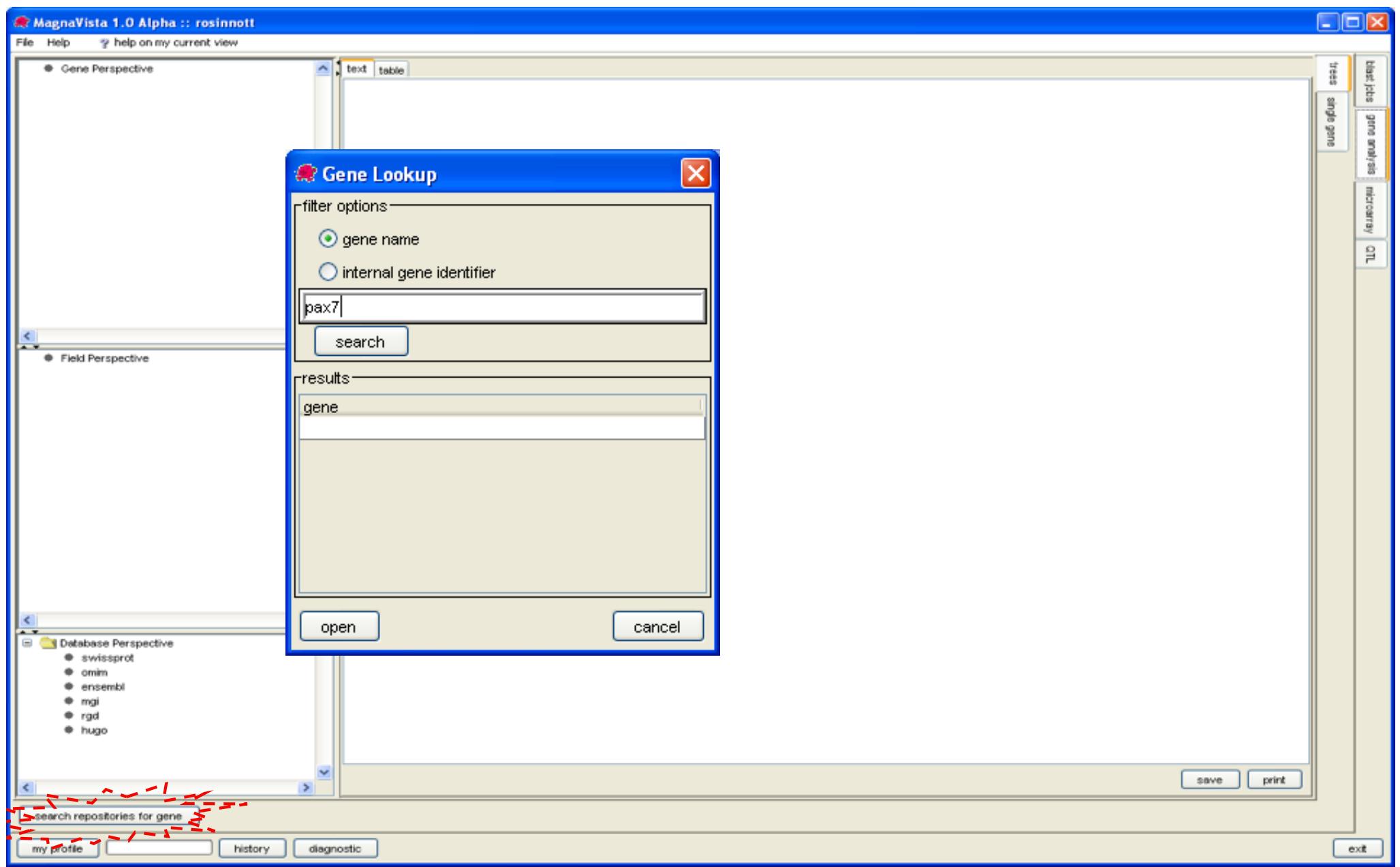
Magna Vista Service



blast



MagnaVista



MagnaVista

MagnaVista 1.0 Alpha :: rosinnott

File Help ? help on my current view

Gene Perspective
Gene Identifier
Field Perspective
Database Perspective

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"/>					

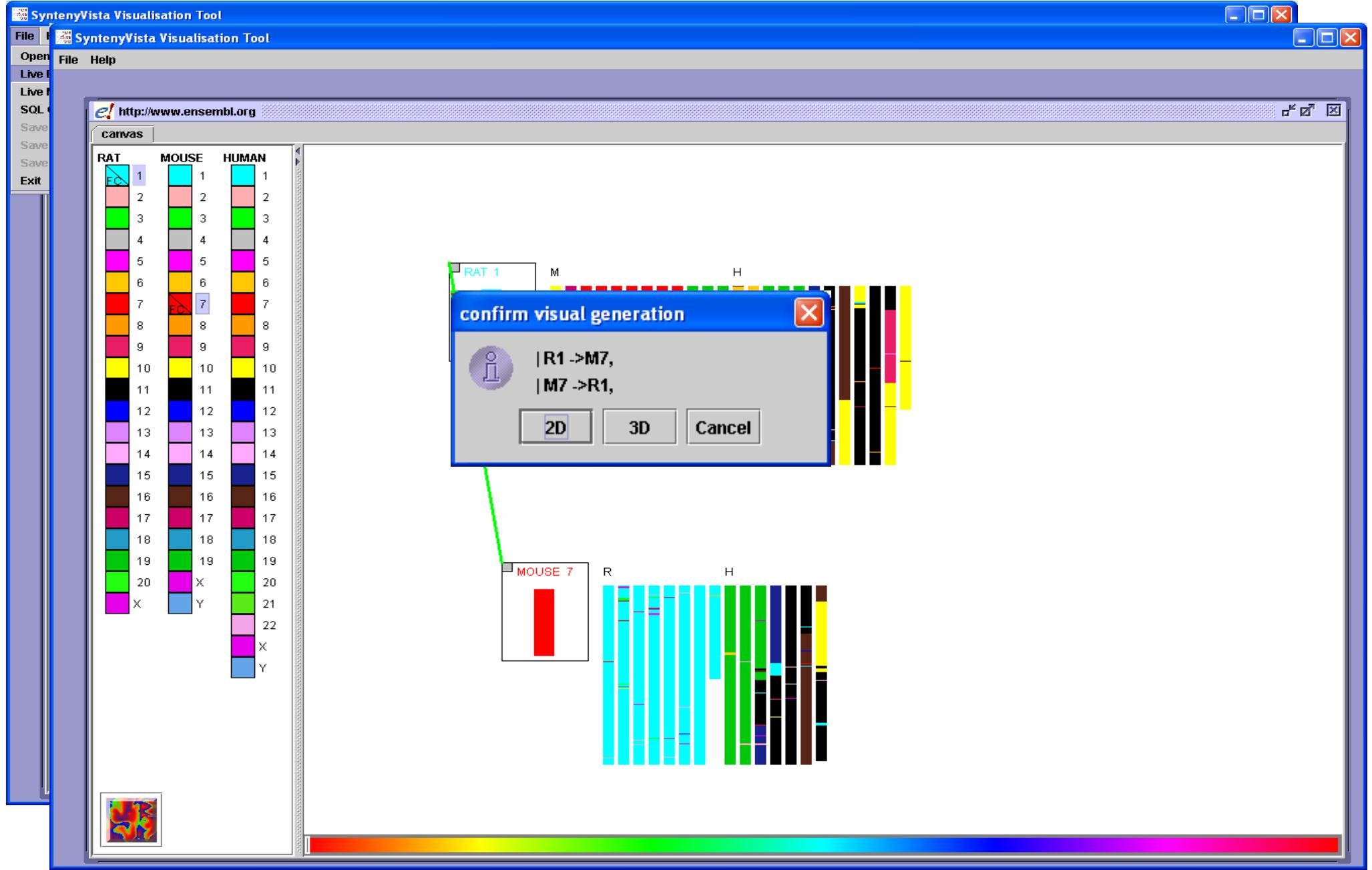
genes databases perspectives field cross references fields Probe Sets qtl preferences

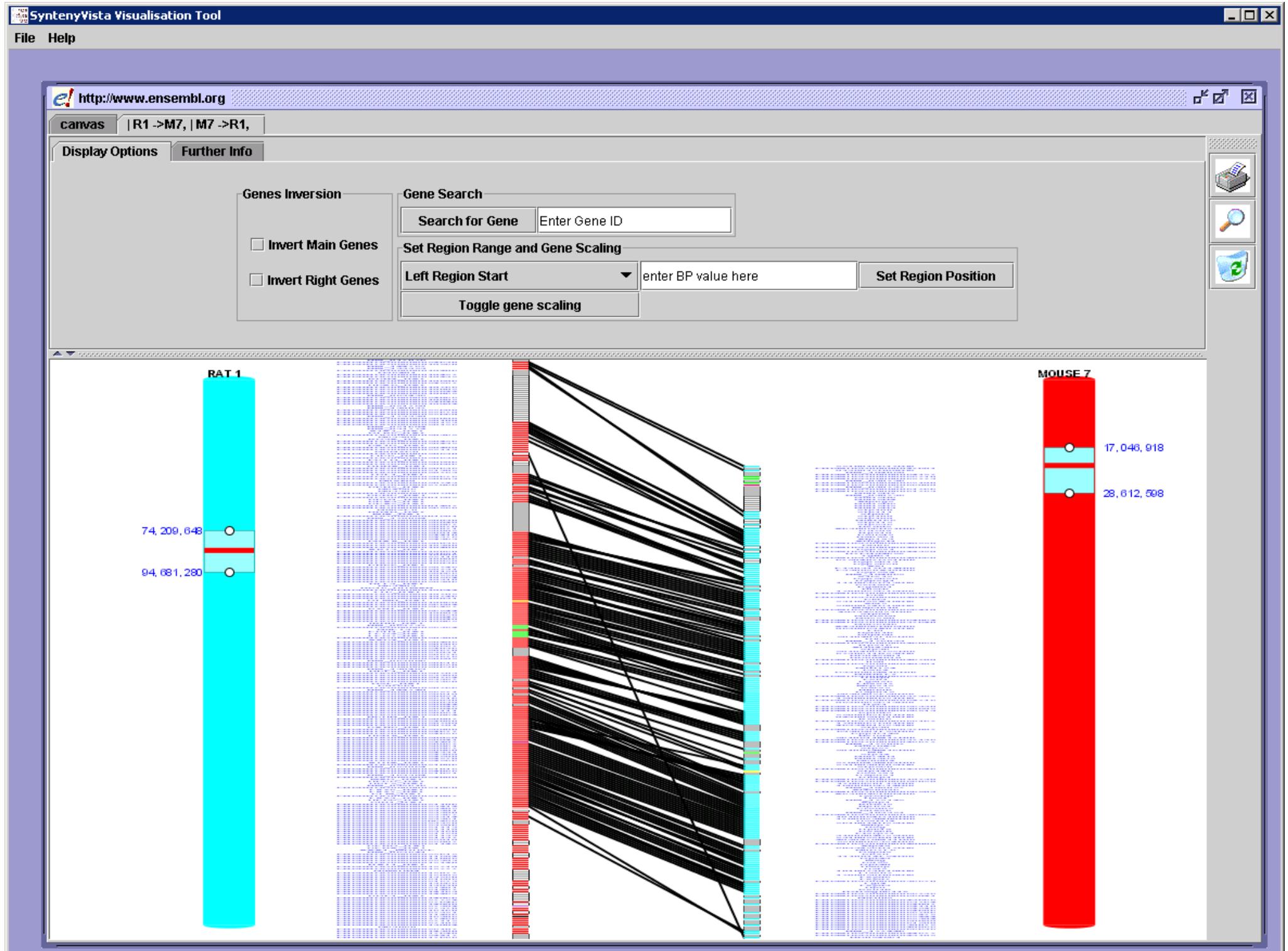
swissprot omim ensembl mgi rgd hugo

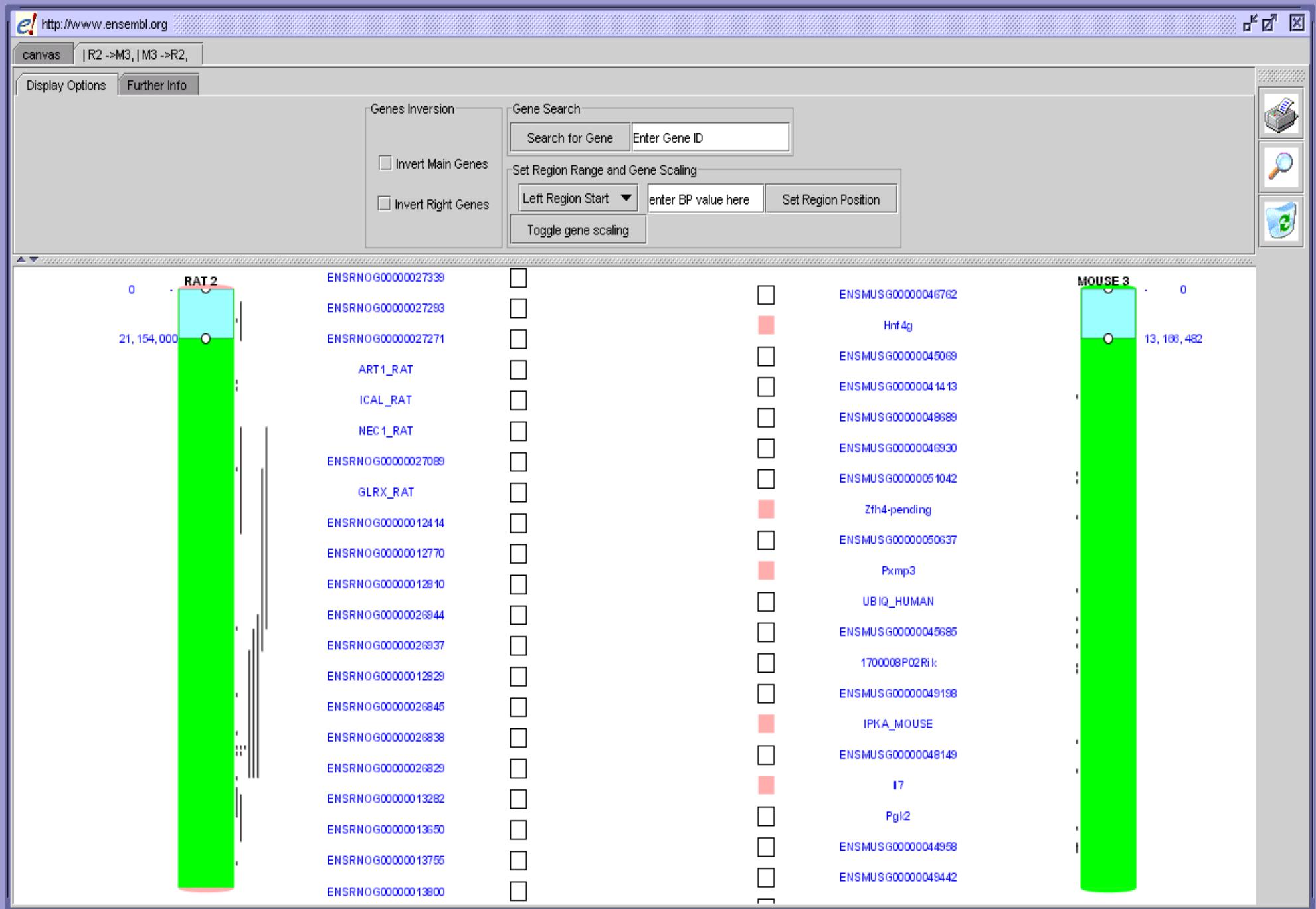
print exit

my profile history diagnostic

Importance of Data Visualisation



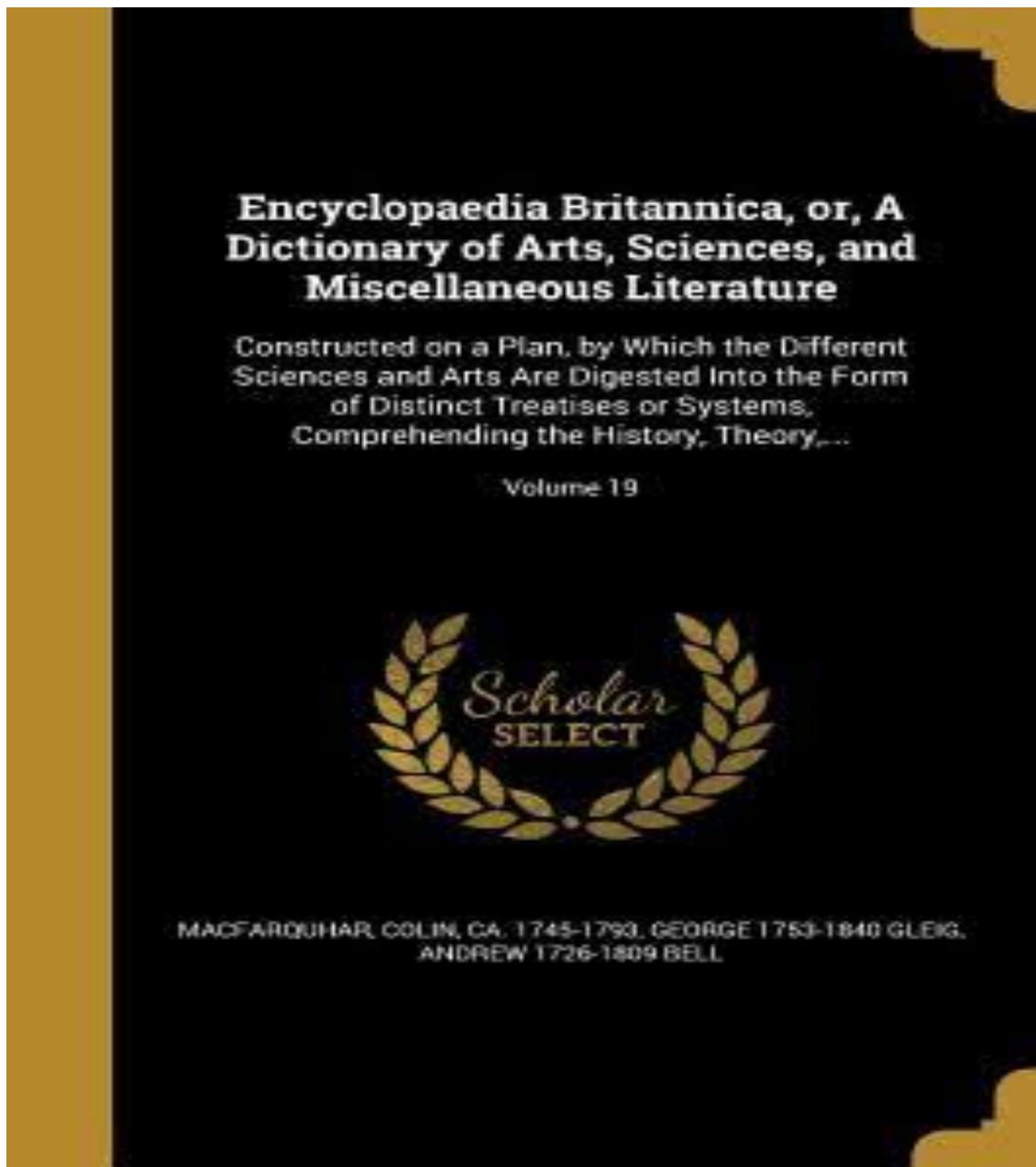




BREAK

Data -> Knowledge?

- Once upon a time...



Crowdsourcing Knowledge & Reasoning

- Many approaches that *work* (?)...

The screenshot shows the Stack Overflow homepage with the URL <https://stackoverflow.com/questions?sort=votes>. The main content area displays a list of questions, each with a title, a green box indicating the number of answers, and a snippet of the question. The sidebar on the left includes links for Home, PUBLIC, Stack Overflow, Tags, Users, Jobs, Teams, Q&A for work, and Learn More. The right sidebar features advertisements for Stack Overflow Jobs and a section titled "Jobs near you" with job listings for Junior Web Developer and Full Stack Engineer.

Home

PUBLIC

Stack Overflow

Tags

Users

Jobs

Teams

Q&A for work

Learn More

All Questions

Ask Question

16,632,793 questions

Newest 369 Featured Frequent Votes Active Unanswered

22009 Why is it faster to process a sorted array than an unsorted array?
votes
22 answers
1.3m views

Here is a piece of C++ code that seems very peculiar. For some strange reason, sorting the data miraculously makes the code almost six times faster. #include <algorithm> #include <ctime> #...

java c++ performance optimization branch-prediction

asked Jun 27 '12 at 13:51
GManNickG 293k ● 39 ● 406 ● 505

18288 How to undo the most recent commits in Git?
votes
75 answers
7.0m views

I accidentally committed wrong files to Git, but I haven't pushed the commit to the server yet. How can I undo those commits from the local repository?

git git-commit git-reset git-revert

community wiki
73 revs, 49 users 14%
Peter Mortensen

14125 How do I delete a Git branch both locally and remotely?
votes
10 answers

I want to delete a branch both locally and on my remote project fork on GitHub. Failed Attempts to Delete Remote Branch \$ git branch -d remotes/origin/bugfix error: branch 'remotes/origin/bugfix' ...

Make your next move
with a career site
that's by developers

Get started

stack overflow JOBS

Jobs near you

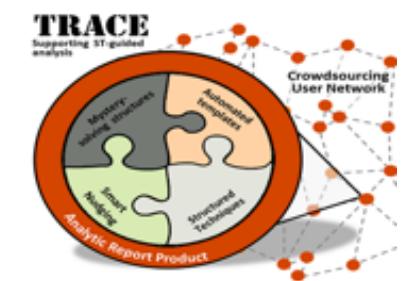
Junior Web Developer
Hays plc ● Melbourne, Australia
A\$75K - A\$85K
javascript cf#

Full Stack Engineer: Impact over 4,000
charities and social enterprises
EthicalJobs ● Collingwood, Australia
A\$80K - A\$110K
react php

CREATE Program



- CREATE
 - Commenced in 2017
 - Involve(d) four teams*
 - TRACE – Trackable Reasoning and Analysis for Collaboration and Evaluation (Syracuse)
 - Co-Arg – Cogent Argumentation System with Crowd Elicitation (George Mason University)
 - BARD – Bayesian Argumentation via Delphi (Monash)
 - SWARM – Smartly-assembled Wiki-style Argument Marshalling (UniMelb)
 - www.swarmproject.info



* Interesting paradigm of funding...

SWARM Overview



• Solution overview



Kalukistan Bomb Problem

Instructions
Assume you (as a team) are the Identity Intelligence (ID2) Officer serving with Special Operations Command Asia (SOCA) and currently stationed in Kalukistan. Perform the duty specified in the Memorandum.

Note:

- Your report should consider multiple hypotheses and estimate their likelihood.
- You are not required to have a prior understanding of forensics or DNA science.

Background
As an Identity Intelligence (ID2) Officer serving with Special Operations Command Asia (SOCA) and currently stationed in Kalukistan, your role is liaising with and providing training and assistance to the Central Government's National Security Forces and the local police (gendarmerie) in the capital, Khyberadad.

Kalukistan is a Central Asian country with vast mineral wealth, a weak central government and a history of corruption and endemic inter-tribal/ethnic warfare. Historically and geographically at the crossroads of trade routes and imperial ambitions, it also has a long history of military incursions and occupations by numerous foreign powers over the last century.

Most recently, the Central Government has aligned itself with a

Quoll20

Type in your new response here ...

Hints

EMU52 created 2 months ago updated 2 months ago
CROCODILE8 created 2 months ago updated 2 months ago

Here is a timeline

RESOURCE

DINGO13 created 2 months ago updated 2 months ago

So far, I have 3 hypotheses:

Hypothesis 1
The owner is the bomb: money transfer business the shop during the day from the Customer.

Hypothesis 2
The owner and courier:

READINESS RATING simple detail

Reasoning

Completeness	Not at all (0-19%)	10
Correctness	Moderately (45-55%)	50
Logic	Moderate (45-55%)	50
Evidence	Well (56-80%)	68
Alternatives	Poorly (20-44%)	32

Communication

Clarity	With huge effort (20-44%)	32
Format	Moderately (45-55%)	50

second on key assumptions. As you will see from my comment, I am still unsure we need to argue that this was the bomb used in the ministry bombing.
15:11 Crocodile8:
and lastly I have a question about whether

SWARM Overview



• Key features

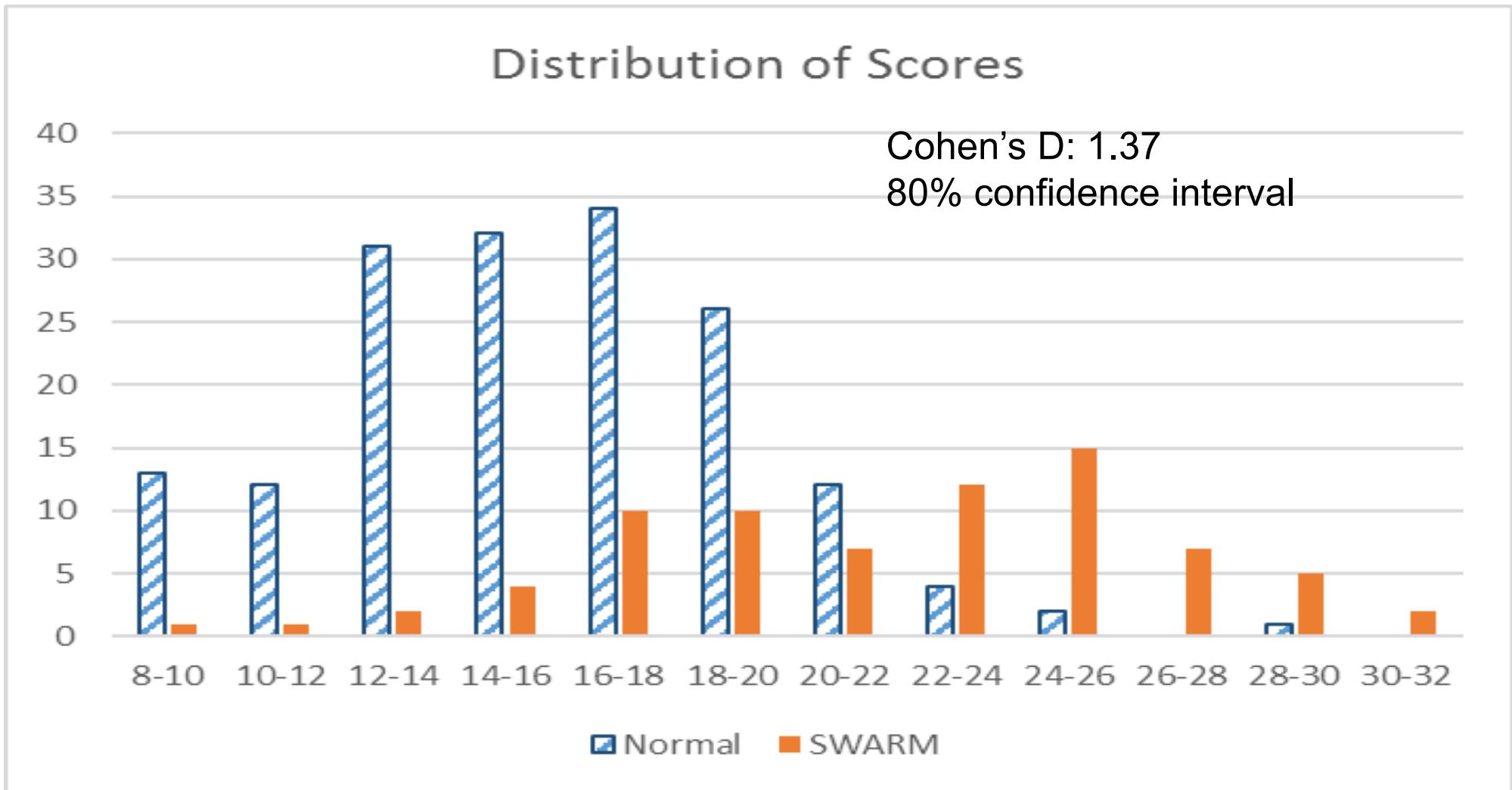
- Anonymous
 - Avatars and rating
- No leader
- Team size flexibility
 - Dockerized AWS/NeCTAR
 - Kubernetes/Docker SWARM
- Public vs Private
 - Off platform work supported
- Arbitrary contributions
 - No payments/obligations
- Social interactions (chat)
 - Social warmth
- Ad hoc use encouraged

A screenshot of the Swarm platform interface. The main title is "Kakulistan Bomb Problem". It includes sections for "Instructions", "Background", and "Hypotheses". The "Hypotheses" section shows contributions from users "Tweeps" and "Oscarouf", each with a rating of 78. A timeline section shows a post from "Oscarouf" dated April 13th, 2018, at 13:11. The "Resource" section shows a contribution from "Oscarouf" with a rating of 2.0. The "Hypothesis 1" section contains text about the owner being a bomb maker. The "Hypothesis 2" section contains text about the owner and courier working together. The bottom right corner has a "Message" button.

Proof of the Pudding



- Does SWARM help to reason better?
 - ASIO, VicPolice, ...
 - 81 reports on platform; 167 off platform (normal)



SWARM Clouds



- Developed on openStack (NeCTAR) The Nectar logo, which consists of a yellow hexagonal icon made of smaller hexagons followed by the word "nectar" in a bold, black, sans-serif font.
 - Scripted solution using Docker & Kubernetes
- Deployed to AWS (US)
 - \$1000+ / month for basic use
 - (costs ramp up a LOT!)
- Benefits of scripted solution
 - Developed/test/trialled on free Cloud (NeCTAR)
 - Deployed to AWS when ready
 - Not possible to do if would have used AWS Elastic Container Solution for Kubernetes (EKS)
 - (...or Azure AKS or Google GKE etc etc)



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, street41

Australian Urban Research Infrastructure Network (AURIN)

- EIF/NCRIS federally funded project
 - DIISRTE -> Innovation -> Education
 - \$40m+ project (+\$18.9m)
 - 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

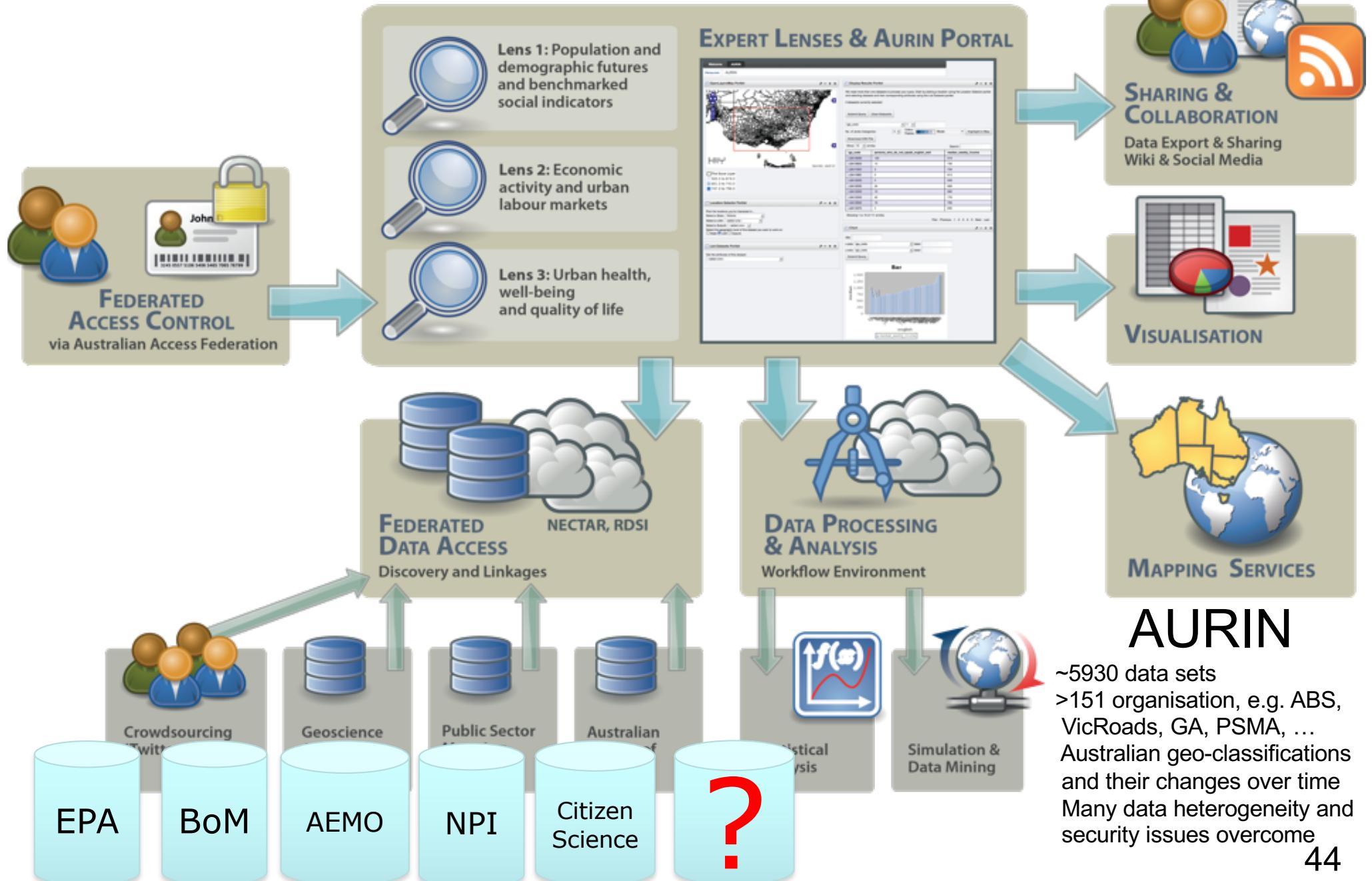
- Consolidate tools and best practice/community know-how!
- Allow researchers to share results, interact and collaborate
 - No single expert!

- Allow data providers to keep control of their data and its use

- Authentication and authorisation (and auditing/accounting)



AURIN Simplified



AURIN Example (in one slide!)

https://portal.aurin.org.au

Lung Cancer Patient Data

CSV JSON Title Name

PostCodes Current smoker Ex-smoker Never smoked Not stated (blank) Grand Total

PostCodes	Current smoker	Ex-smoker	Never smoked	Not stated	(blank)	Grand Total
1000	1					1
3000	2	5	2	1	1	11
3002		1			1	2
3003		1			1	2
3004		2			1	3
3006		1			1	2
3008			1			1
3011	1	1			1	3
3012	1	2			1	4
3015	1	1			1	3

Records Number : 382

Substance Quantity

The dashboard displays two main maps of Melbourne. The top map shows Lung Cancer Patient Numbers Per Postcode, with areas color-coded by patient count. The bottom map shows PM2.5 Readings (2014), with areas color-coded by reading range. The AURIN interface includes a sidebar with various data layers and analysis tools.

Logged in as Richard Sinnott | Project LungCancer2016 (All changes saved) Help MyAURIN

Area

- Area Selection Greater Melbourne (gccsa/2GMEL)
- Bounding-box Selection [144.8457,-37.8382,145.1588,-37.5711]

Data

- Lung Cancer Patient Data
- SA2 Health Risk Factors - Modelled Esti... Melbourne - North East (sa4/209)
- Victorian Road Traffic Volumes [144.8807,-37.7851,145.5800,-37.2629]
- National Pollutant Inventory - Site Emis... [144.3336,-38.5030,145.8784,-37.1751]
- National Pollutant Inventory - Site Emis... [144.3336,-38.5030,145.8784,-37.1751]

Visualise

- Maps, Charts and Graphs
- Victorian Road Traffic Volumes - GeoJ...
- Never Smoked
- Nitrogen in Water (2014) Victoria
- National Pollutant Inventory - Site E...
- PM2.5 emitters

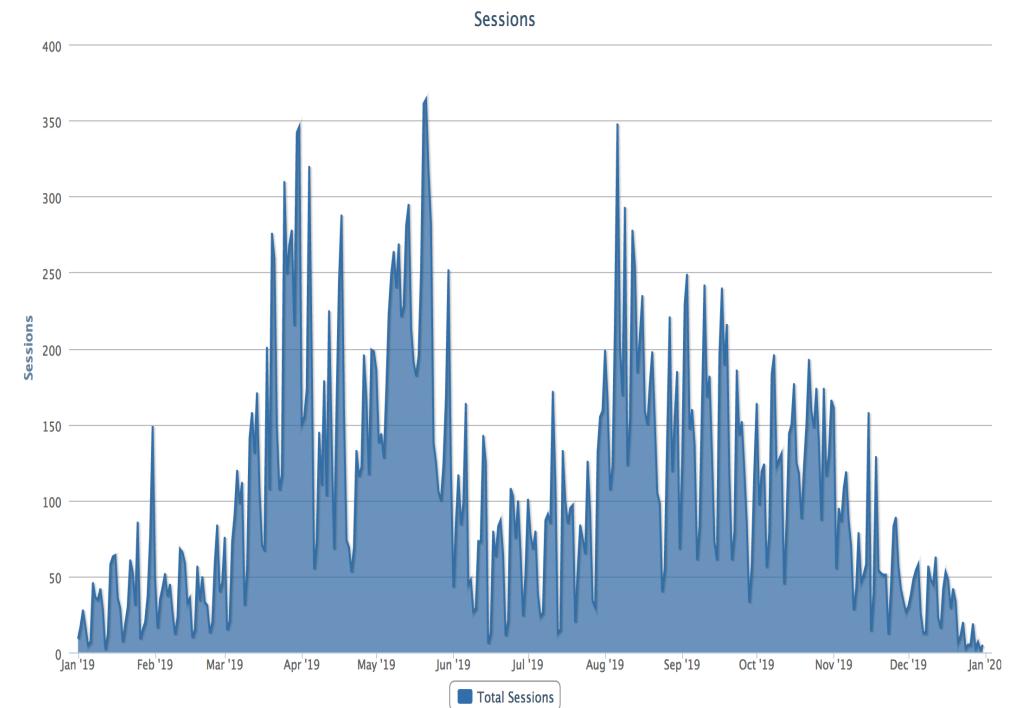
Analyse

- Tools
- Bar Chart - 1
- Scatter Plot - 1
- Scatter Plot - 2
- Bar Chart - 2
- Summary Statistics - 1

Disclaimer Terms of Use Copyright Report an Issue

AURIN Clouds

- Original plan to use NeCTAR
 - Early reliability issues
- Actual plan
 - Servers purchased (VMware)
 - Used for production system
 - 8 years old/partial refresh 2016
 - 3 days outage in 8 years!!!
 - Failover system on NeCTAR
- Dev, Staging, Production
 - 20,000+ users
 - ~4million lines of code
 - CouchDB, web services (ReST, ...), geoJSON, ...
 - Much of the code not written by my team!!!
 - Move to using Docker containers/Kubernetes ongoing
 - e.g. L. Chen, Y. Pan, [R.O. Sinnott](#), *Auto-Scaling a Cloud-based Walkability Tool through Kubernetes and Docker Swarm*, CLOSER 2020, Prague, Czech Republic, May 2020.



Demonstration in Workshops

(note – Assignment II)

AURIN Homework

(<https://portal.aurin.org.au>)

Find the suburb (SA2) in Greater Melbourne that had the highest number of Jobseeker recipients in June 2020.

(not assessed!)

Questions ... ?