

**LINUX**

**AND**

**GENOMICS**



# TWO REVOLUTIONS

CANADA'S MICHAEL SMITH GENOME SCIENCES CENTRE

MARTIN KRZYWINSKI

УДЯФН BUTTERFIELD



# Acknowledgements

asano, jennifer : astakhov, vadim : astakhova, tamara : astell, caroline : babakaiff, ryan : bala, miruna : barbazuk, stephen : barber, sarah : baross, agnes : bilenky, mikhail : bosdet, ian : brooks-wilson, angela : brown-john, mabel : butterfield, yaron : chan, susanna : chan, andy : chand, steve : chang, elbert : charest, david : charters, anita : chittaranjan, suganthi : chiu, gordon : chiu, readman : chow, william : chuah, eric : chun, hye-jung elizabeth : collins, jennifer : coughlin, shaun : crisostomo, lulu : del rio, luis : delaney, allen : despres, chantale : dhalla, noreen : farnoud, noushin : featherstone, ruth : field, matthew : fjell, chris : flibotte, stephane : freeman, doug : go, anne : gorski, sharon : griffith, obi : griffith, malachi : guan, jun : guin, ranabir (ran) : halaschek-wiener, julius : hanson, robyn : harrison, isabel : hassel, maik : hirst, martin : holt, robert : hou, claire : huang, peter : hume, lynn : jang, carrie : jendo, aga : jones, steven : khattra, jaswinder : kirkpatrick, robert : krzywinski, martin : kwong, leticia : leach, stephen : lee, darlene : lee, stephanie : lee, lisa : leung, amy : leung, derek : li, bernard : li, yvonne : liao, nancy : lin, keven : liu, jerry : ma, kevin : marcadier, julien : marra, marco : masson, amara : mathewson, carrie : matsuo, corey : mayo, mark : mayo, michael : mconechy, melissa : mcdonald, helen : mckay, sheldon : melnyk, brianna : messervier, steve : mills, courtney : missirlis, perseus : moksa, michelle : montgomery, stephen : moore, richard : moran, josh : morin, gregg : morin, ryan : narain, abhishek : nayar, tarun : novik, karen : o'connor, katie : oliveira, lisa : olson, teika : oveisi, mehrdad : palmquist, diana : pandoh, pawan : peloso, irene : persaud, deryck : petrescu, anca : pleasance, erin : prabhu, anna-liisa : pugh, trevor : qadir, mohammed : quayle, adrian : robertson, neil : robertson, gordon : roger, jennifer : rogers, sean : ruschkowski, sharon : ruzanov, peter : saeedi, parvaneh : santos, joseph roy : schein, jacquie : schnerch, angelique : schoeffel, kirk : shin, heesun : siddiqui, asim : sipahimalani, payal : sleumer, monica : smailus, duane : sohaib ali, mohammad : stott, jeff : tai, isabella : teague, kevin : tremblay, ashley : tsai, miranda : tsang, eddy : varabei, dmitry : vardy, jill : varhol, richard : warren, rené : wilson, gary : with, sheila : wong, david : wong, kim : wright, charlotte : wye, natasja : yang, george : yuen, wendy : zeng, thomas : zhao, yongjun

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USE LINUX SIG**

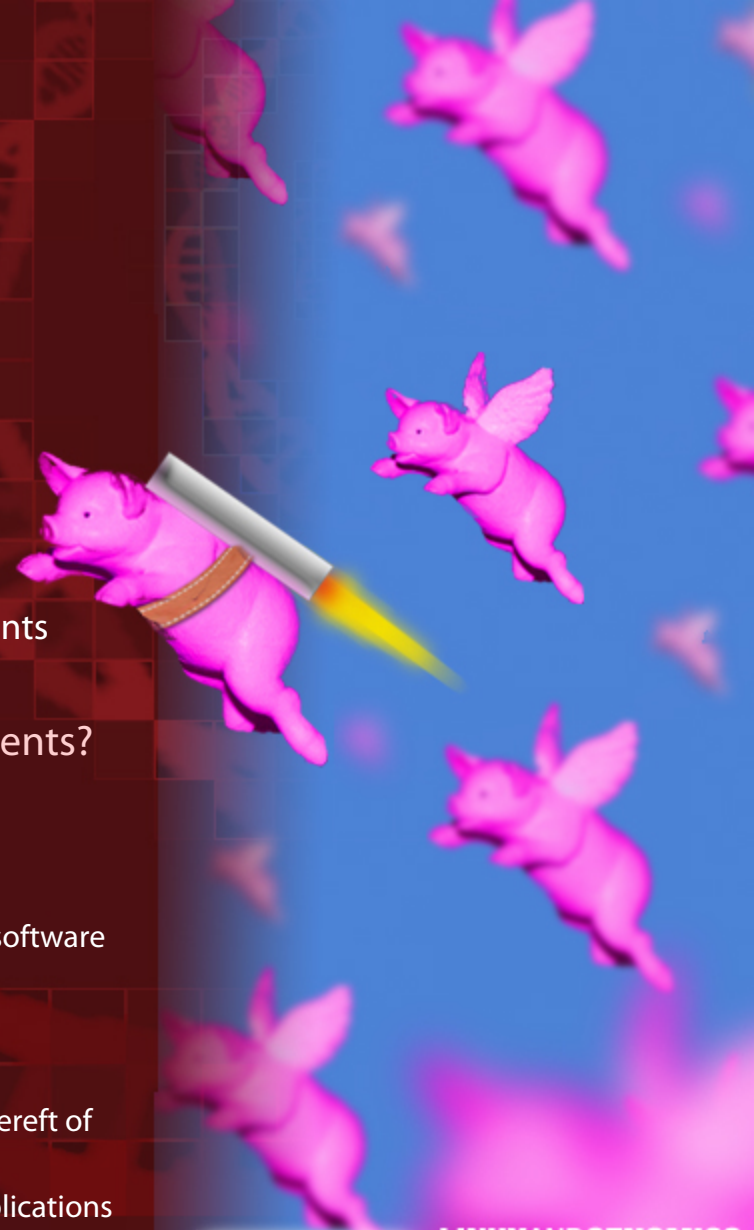
IN EVERY  
REVOLUTION  
THERE'S ONE  
MAN WITH  
A VISION

JAMES T. KIRK  
IN "MIRYARA, MIRYARA"



# two revolutions

- **revolution [n.]**
  - drastic and far-reaching way of thinking and behaving
  - paradigm shift
- promise beneficial change
  - increased quality of life for under-represented
  - independence, broader rights, protection from predatory elements
- who are the (a) under-represented, and (b) predatory elements?
  - in computers
    - under-represented: **the OS user**
    - predators: corporate bodies with monopolies, marketing, FUD
    - my ability to meaningfully participate in development of close-source software is negligible
  - in science
    - under-represented: **individuals**
    - predators: restrictive patents/licenses, large privately-funded centers bereft of public accountability, academic cliques isolated from social context
    - my ability to grasp new findings, evaluate their impact and foresee implications is negligible



**USENIX**

LINUXANDGENOMICS  
TWO REVOLUTIONS

USELINUXSIG



NO, THE WHOLE POINT  
OF LINUX ISN'T THAT  
IT'S FREE, AND NEVER  
HAS BEEN.

THE POINT WITH LINUX IS  
THAT I DIDN'T HAVE A  
GOOD OS ON MY MACHINE  
... SO I WROTE ONE.

LINUS TORVALDS  
COMP.OS.MINIX, 20 DEC 1992

# revolution 1 | linux

- UNIX on Intel CPU
  - free OS on cheap hardware
    - performance per unit cost is steadily dropping
    - “free” vs “absolutely free” – what is the real cost?
    - in our lab, computer hardware costs are negligible compared to personnel and laboratory equipment
    - my workstation costs ~ 1/50<sup>th</sup> of my salary
  - line between user and developer base blurred
  - “Beowulf” clusters are ubiquitous
    - cheap, limitless potential, free toolkits
- community encourages users to experiment
  - communities, not cliques
  - anyone can contribute
    - code, document, test, mirror
- rewards innovation and sharing
  - does not rely on scarcity tokens to drive its economy
  - skill and meaningful contribution is rewarded independent of socioeconomic status, reputation, or charisma

Ten years out, in terms of actual hardware costs you can **almost** think of hardware as being free – I’m **not saying it will be absolutely free.**

*B. Gates*

emphasis is mine

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USELINUXSIG**

# revolution 2 | genomics

- genomics, [n.]
  - study of organisms in terms of their genomes
- offers insight to fundamental building blocks of living systems
  - structure, biochemistry, evolution
  - organism as network – system biology
  - understand basic processes in a living cell
- genomics is the branch of golden era of biology
  - less classification, more integrated descriptive and predictive models
  - periodic table of elements, quantum mechanics
- applications are the fruits
  - accelerate drug development by identifying target pathways
  - recognize variants contributing to health and resistance to disease
  - generate diagnostics for early detection of cancer
  - understand reasons for individual drug resistance

It should be possible to understand the difference between a “bag of molecules” and a biological system.

*F. Collins et al.*

Nature (2003) 422: 835-847

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USELINUXSIG**

# common philosophy

## programming

kernel

patches

code repository

hacker

gatekeeper

software

## science

knowledge

papers

literature

scientist

editor

mindware

improvement of the "human condition"



# common process

## programming

new code should not crash system

add features or elegance

robust and logical

properly documented

major rewrites require damn good code

## science

new conclusions backwards-compatible  
with existing knowledge

increase range power of explanation and  
prediction; reduce complexity of  
fundamental laws

in accordance to observed data

reproducible

extraordinary claims require  
extraordinary evidence

## advancement

# embracing openness

- success in science requires open source principles
  - examine, verify and extend knowledge
  - knowledge must be publicly accessible to promote discovery
    - many eyes make all bugs shallow
  - scientists have created many open source tools
    - communication and collaboration
    - manipulate, munge, analyze large data sets
    - agree, adopt and improve
- integration of science into society requires open source principles
  - scientific process is not suitable for handling emergent ethical, legal and social issues (ELSI)
  - it is impossible to “test” wide-scale social models
  - social harm caused by misguided policies cannot be reversed
  - keeping science and its products open allows everyone to participate in ELSI discussions
    - experts from social sciences, law, humanities, cultural anthropology
    - public forums

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USELINUXSIG**

# public effort triumphs

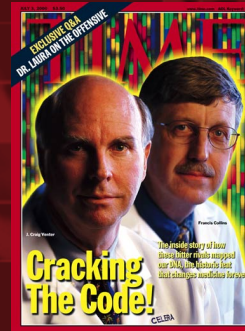
- public assembly of human genome done on UCSC centicluster
  - June 2000
  - 100 800-MHz P3 Linux boxes



UCSC kilocluster  
1024 Linux nodes

I thought it would help to get as much information about genes and the genome into the public domain to help discourage people from patenting it wholesale.

*Jim Kent*



Venter, Collins  
TIME July 3 2000



Venter, Clinton, Collins

# ACTGs of genomics

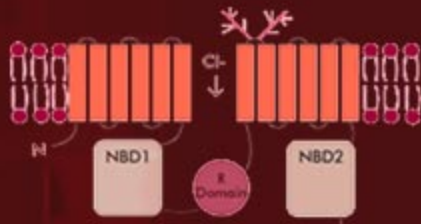


transcription  
translation

proteins

> chr 7, 159 Mb

...  
ATTATGCCTGGCACCA  
TTAAAGAAAATATCAT  
**CTTTGGTGTTCCTAT**  
GATGAATATA...



phe  
508

normal

K	E	N	I	I	F	G	V	S	Y	D					
A	A	G	A	A	A	T	A	T	C	A	T	C	A	T	G

cystic fibrosis

K	E	N	I	I	G	V	S	Y	D						
A	A	G	A	A	A	T	A	T	C	A	T	C	A	T	G

membrane transporter

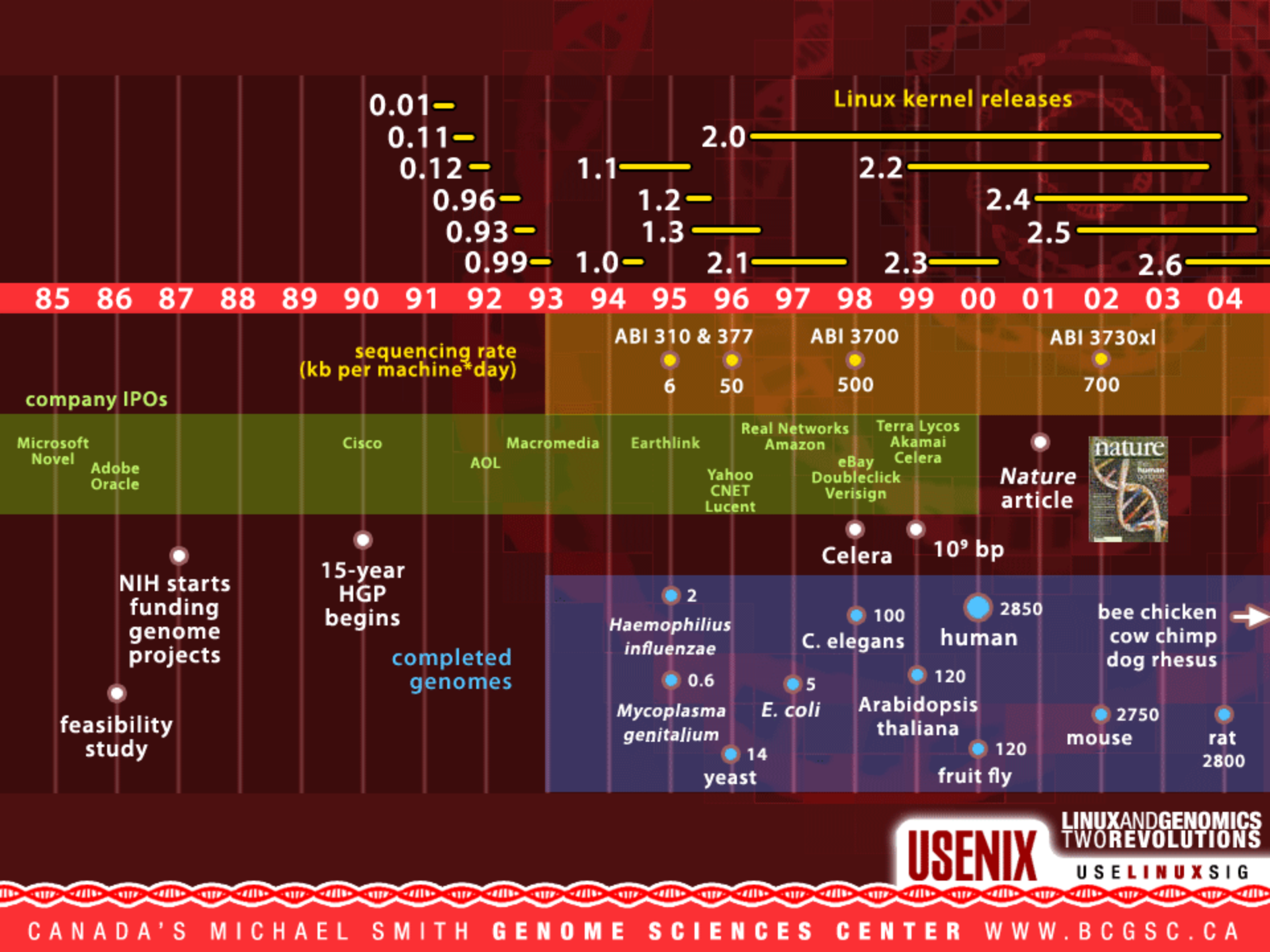
proteins participate  
in biochemical pathways  
and mediate processes

- > alter molecules
- > transport molecules
- > catalyze reactions

USENIX

LINUXANDGENOMICS  
TWO REVOLUTIONS

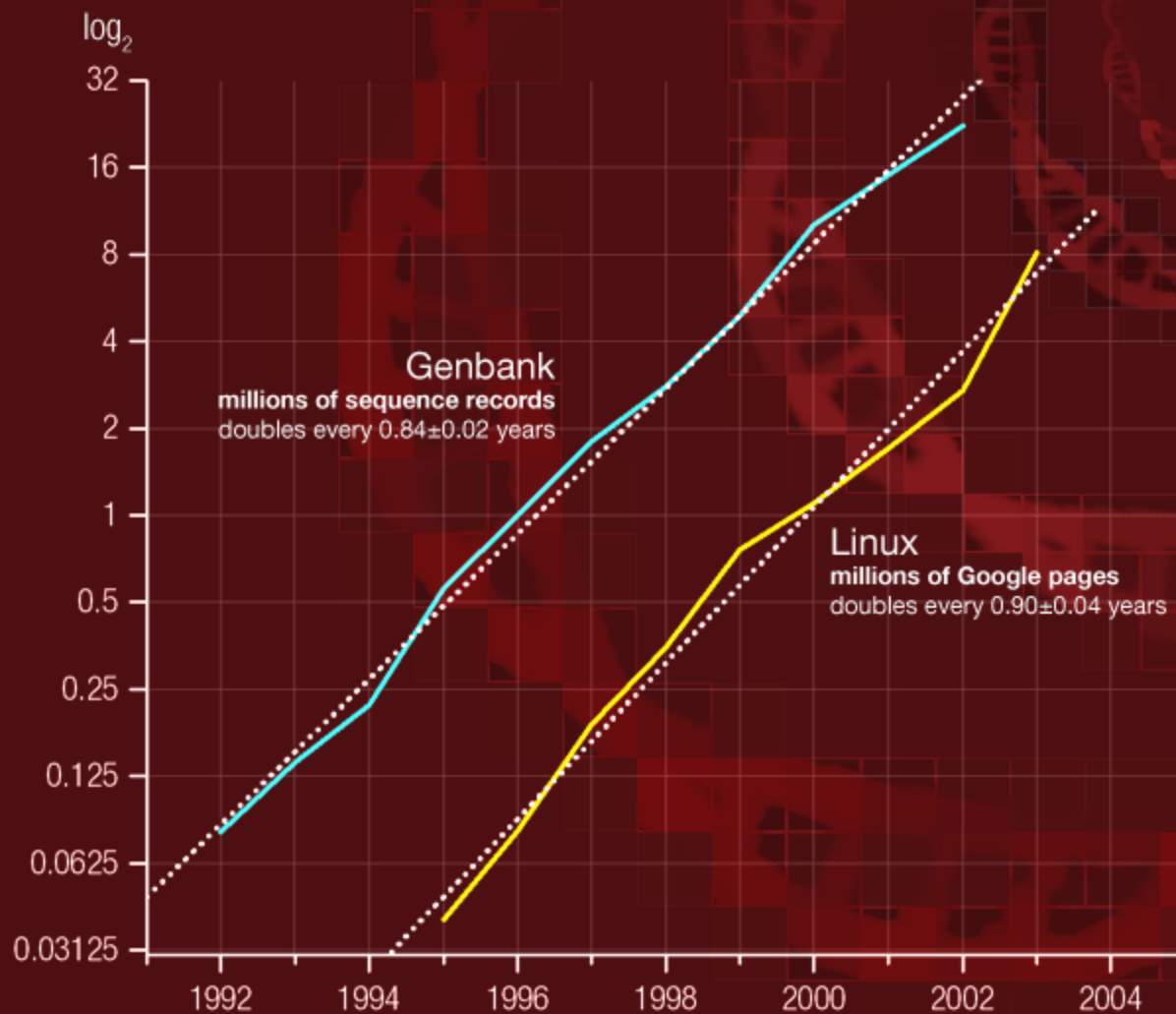
USELINUXSIG



**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USELINUXSIG**



SCIENTIFIC PROGRESS  
WILL BE MAXIMIZED  
BY EARLY, OPEN AND  
CONTINUING ACCESS  
TO LARGE DATA SETS.

...REWARDING AND PROTECTING  
THE INTERESTS OF SCIENTISTS  
WHO WISH TO SHARE THEIR  
DATA WITH THE COMMUNITY  
IN SUCH A GENEROUS MANNER.



FRANCIS COLLINS

A VISION FOR THE FUTURE OF GENOMICS RESEARCH, NATURE (2003) 422: 835-847

# genomics data is public

- the net started at CERN
  - physicists needed to share documents and large datasets
  - genome centers and scientists have been networked from the start
- public data release is a mandate of sequencing centers
  - data is submitted to Genbank, a public repository of sequence information, as soon as it is collected
  - "Bermuda Principles" ratified during a 1996 sequencing conference, call for **automatic and rapid** release of primary sequence data to the public domain
  - internet is the natural forum

The top screenshot shows the NCBI 'Submit to GenBank' page. It features a search bar, navigation tabs (PubMed, Entrez, BLAST, OMM, Books, TaxBrowser, Structure), and a sidebar with links to various NCBI resources. The main content area is titled 'Submitting Sequence Data to GenBank' and includes sections for 'Submitting Sequence Data to GenBank' and 'Receiving an Accession Number for your Manuscript'. The bottom screenshot shows the EMBL-EBI 'EMBL Nucleotide Sequence Database' page. It features a search bar, navigation tabs (Home, About EBI, Research, Services, Toolbox, Databases, Downloads, Submissions), and a sidebar with links to various EBI resources. The main content area is titled 'EMBL Nucleotide Sequence Database' and includes sections for 'EMBL Nucleotide Sequence Database' and 'EMBL Fetch'.

[www.ncbi.nlm.nih.gov/](http://www.ncbi.nlm.nih.gov/) [www.ebi.ac.uk/emb/](http://www.ebi.ac.uk/emb/)

[www.gene.ucl.ac.uk/hugo/bermuda.htm](http://www.gene.ucl.ac.uk/hugo/bermuda.htm)

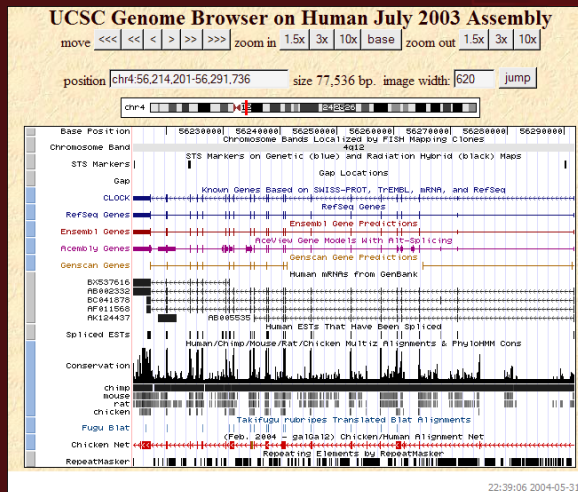
USENIX  
LINUXANDGENOMICS  
TWO REVOLUTIONS  
USELINUXSIG

CANADA'S MICHAEL SMITH GENOME SCIENCES CENTER WWW.BCGSC.CA



# data viewers are public

- openly available genome browsers
  - permit public mining of information
  - visualization of annotations
- Ensembl and UCSC browsers run on MySQL
  - entire data set can be downloaded freely
  - you can become an Ensembl mirror
    - ~100 Gb, MySQL/Perl



genome.ucsc.edu

www.ensembl.org

**USENIX** LINUXANDGENOMICS  
TWO REVOLUTIONS  
USE LINUX SIG

# genomics community is online

- Lincoln Stein, CSHL (NY, USA)
  - CGI.pm, GD.pm
  - generic genome browser
  - genome knowledge base
  - generic model organism database construction set
    - modular and extensible framework for storing biological information
  - distributed sequence annotation system (DAS)
    - XML web service for contributing sequence information
  - "How Perl Saved The Human Genome Project"
- Jim Kent, UCSF (CA, USA)
  - Gigassembler algorithm assembly algorithm
    - open source, Linux cluster
  - public effort beat Celera by 3 days
- Ewan Birney, Wellcome Trust (UK)
  - Ensembl browser
  - BioPerl bundle



**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

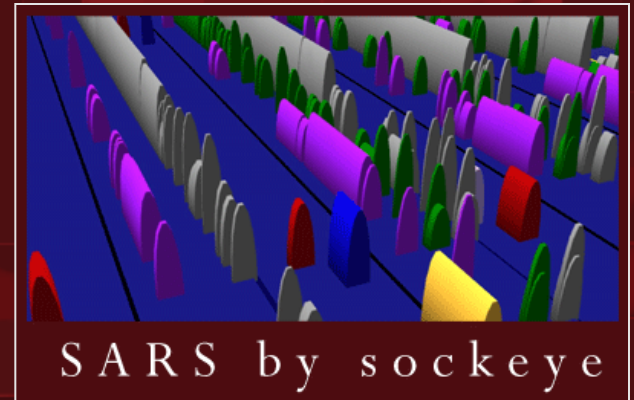
**USE LINUX SIG**

# SARS – linux called to action

- sudden acute respiratory syndrome
  - first case in 1999
  - 8,100 cases, 800 deaths since Nov 2002
- our center was the first to publicly release the full coronavirus sequence assembly
  - 30kb ~ 1/100,000 size of human genome
  - assembled on 8-way Linux machine
  - published on Zope/Apache web server
  - sequence data used to characterize virus and search for potential therapies



SARS affects everyone  
[www.stileproject.com](http://www.stileproject.com)



sockeye 3D genome browser

[www.who.int/csr/sars/country/table2004\\_04\\_21/en/](http://www.who.int/csr/sars/country/table2004_04_21/en/)

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS  
USE LINUX SIG**

# open data - SARS genomes in Genbank

- complete genomic sequence was immediately submitted to Genbank
- 293 sequences from 151 different strains of SARS in Genbank database
- currently >1,600 SARS-related publications indexed by Pubmed

2: [Weinstein RA](#)

**Planning for epidemics--the lessons of SARS.**  
N Engl J Med. 2004 Jun 3;350(23):2332-4. No abstract available.  
PMID: 15175434 [PubMed - in process]

139: [Shurtleff AC](#) [Related Articles](#), [Links](#)

**Bioterrorism and emerging infectious disease - antimicrobials, therapeutics and immune-modulators. SARS coronavirus.**  
IDrugs. 2004 Feb;7(2):91-5.  
PMID: 15057645 [PubMed - indexed for MEDLINE]

The screenshot shows the NCBI Entrez Nucleotide search results for the SARS coronavirus complete genome. The search criteria are 'Nucleotide' for 'SARS coronavirus, complete genome'. The results show a single entry with accession number NC\_004718.3. The sequence is displayed in FASTA format, starting with ATATTAGGTTTACCTACCCAGGAAAAGCCCAACCACTCGATCTCTGTAGATCTGTTCTCTAAAACGA. A box highlights the entry details:  1: [NC\\_004718](#) SARS coronavirus, complete genome gi|30271926|ref|NC\_004718.3|[30271926].

www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?db=nucleotide&val=30271926  
www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=227859



- [faq](#)
- [code](#)
- [awards](#)
- [journals](#)
- [subscribe](#)
- [older stuff](#)
- [rob's page](#)
- [preferences](#)
- [submit story](#)

## Canadian Lab Unravels SARS With A Beowulf Cluster

Posted by [timothy](#) on Sunday April 13, @03:21PM from the you-knew-they'd-come-in-handy dept.

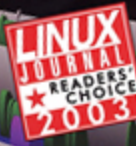
*Amad writes "A Canadian Genetics Research Lab in BC, Canada has used a Linux sequence the genetic code of the virus linked to SARS. This lab is the first to crack the public. You can read an article about the discovery, or check out the lab."*

# B.C. team unravels SARS

By Charlie Anderson  
Staff Reporter

Scientists at the B.C. Cancer Agency have made a major breakthrough in solving the puzzle of the killer virus known as SARS. The Vancouver scientists are the first to crack the genetic code of the Severe Acute Respiratory Syndrome virus — which will speed the diagnosing of victims of the lethal disease and help with the work of finding a vaccine.

**Open-Source Desktop Publishing: Scribus**



# LINUX JOURNAL

## GENOME SCIENTISTS REPORT: How we sequenced the SARS virus in five days

- Secure IMAP mail servers
- Disaster plans for your web site
- Design C++ programs for security
- Highly available cluster management with OSCAR
- Virtual security zones
- Secure Web servers
- Disaster plans for your web site
- Design C++ programs for security
- Highly available cluster management with OSCAR
- Virtual security zones



**SARS cases**

Following is a breakdown of suspected SARS cases worldwide:

- Canada: 274 cases, 13 deaths.
- France: 2,336 cases, 60 deaths.
- Hong Kong: 1,109 cases, 30 deaths.
- Italy: 51 cases, 0 deaths.
- Singapore: 140 cases, nine deaths.
- Taiwan: 23 cases, 0 deaths.
- U.S.: 100 cases, 0 deaths.



溫哥華解碼專家  
**SARS 變種複雜**  
難醫似愛滋

張國榮  
溫市4大瀑布地點  
布殊·肺炎帶擊  
全球旅遊平價情報

**TOP STORY**  
**SARS BUSTERS**  
How the virus crackers got started.  
By Lila MacLellan

**USENIX**  
LINUX AND GENOMICS  
TWO REVOLUTIONS  
USE LINUX SIG

# public feedback

**Subject: transcription starting position of genes in TOR2**

**Dear BCGSC:**

**I downloaded your sequences of TOR2. Thank you very much for your great work about this.**

**JZ, Palo Alto, CA**

**Subject: SARS draft sequence**

**To all involved, CONGRATULATIONS! We will have a PCR-based SARS diagnostic up and running by next week thanks to YOU.**

**MJM, Madison, WI**

**Subject: You have to be NUTS!**

**My daughter doesn't think its such a good idea to have the gene sequencing for the new coronavirus on the internet. I don't either! There should have been a better way! You must be crazy!**

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USELINUXSIG**

# linux at the Genome Sciences Centre



HP 718 DLT autoloader

1999

36 Gb RAID-5  
3 x 18 GB SCSI  
DPT IV RAID card

O'Reilly  
bunch of books

52 GB RAID-0  
3 x 18 GB SCSI  
software raid

dual 400-MHz P2  
512 Mb RAM  
RH 5.2 2.0.36

# linux at the Genome Sciences Centre

Raidion RAID-5  
2 x 8 x 36 GB  
400 GB

VA VAR server 900  
VA 6.0.3 2.0.36

RAIDION.u2w  
RAID controllers

2000

4:41pm up 393 days, 6:48, 9 users, load average: 0.37, 0.29, 0.63

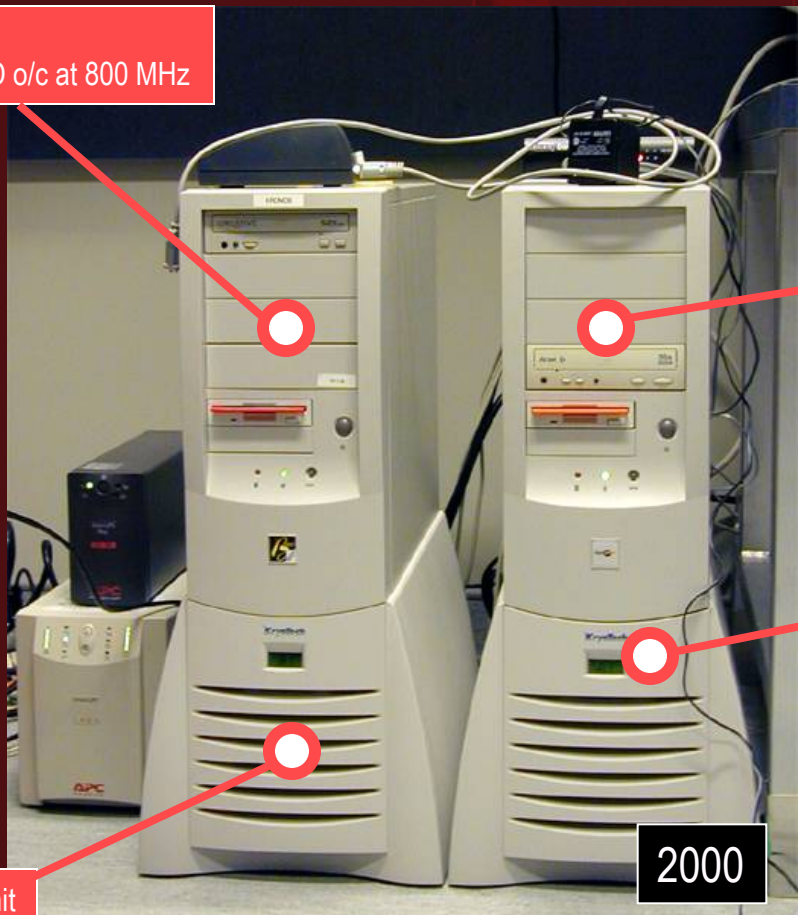




# linux at the Genome Sciences Centre

Kryotech  
700 MHz AMD o/c at 800 MHz

Kryotech "SuperG"  
800 MHz AMD o/c at 1 GHz

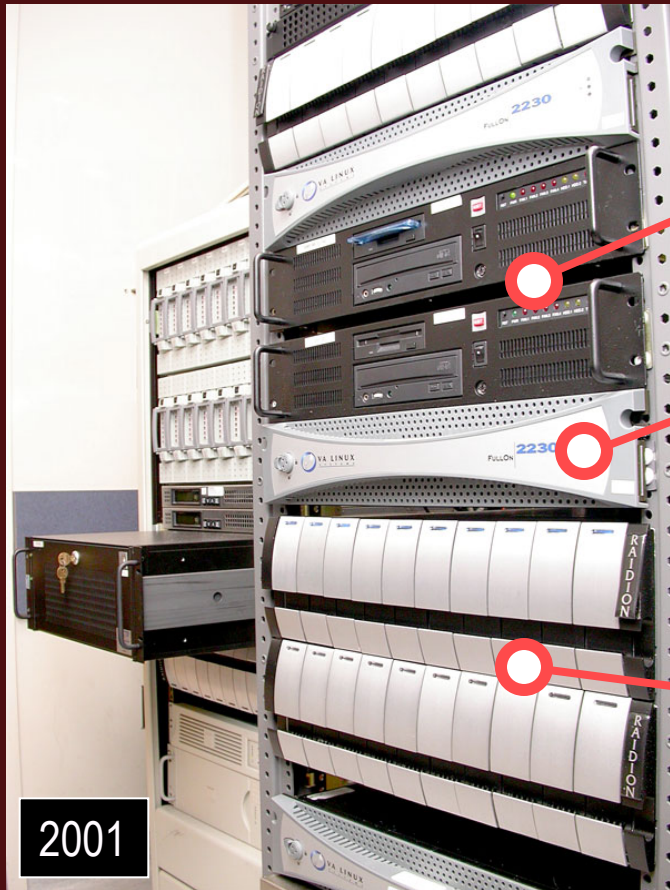


cooling unit  
- 40 C



2000

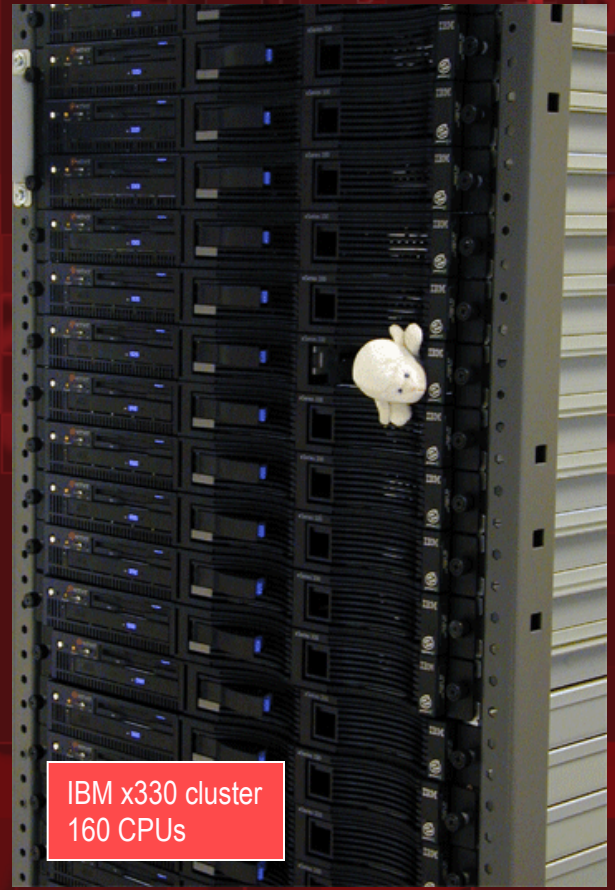
# linux at the Genome Sciences Centre



dual 1 GHz P3  
4 GB RAM

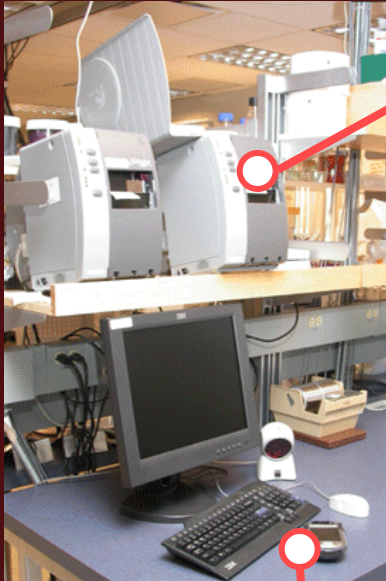
VA Linux 2230  
dual 800 MHz P3  
2 Gb RAM

Raidion RAID-5  
3 x 10 x 72 Gb SCSI  
Raidion.ha RAID controller  
2 TB



IBM x330 cluster  
160 CPUs

# linux at the Genome Sciences Centre



Zebra S600  
ZPL printer

- wireless interface to ZPL printers for on-demand barcodes
  - Perl/Apache, Compaq IPAQ, Zebra S600
- MySQL database stores all events, objects and data

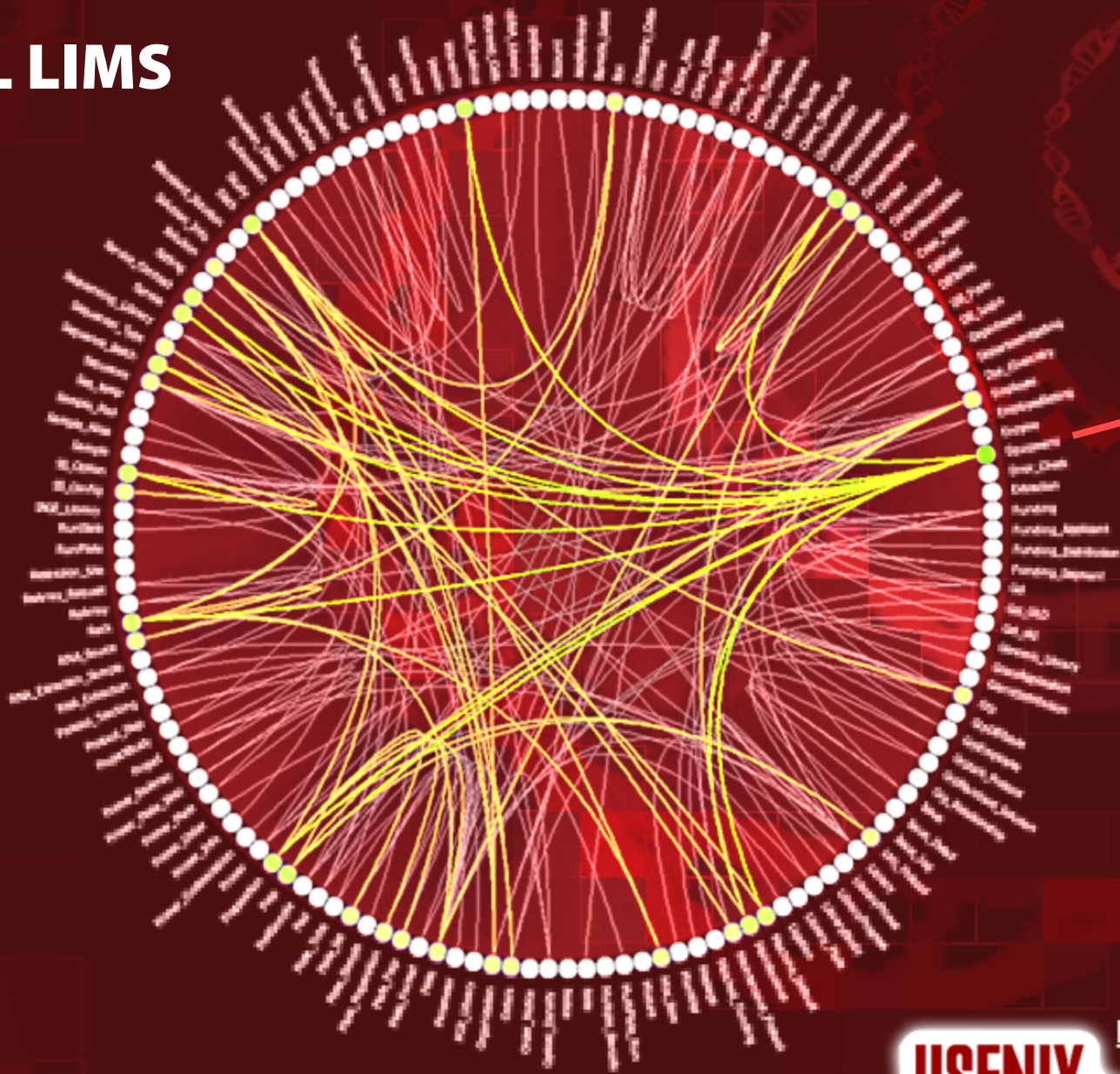
Zebra Z4Mplus  
ZPL printer



Compaq iPAQ



# MySQL LIMS



Equipment

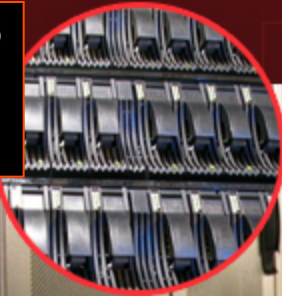
[mkweb.bcgsc.ca/schemaball](http://mkweb.bcgsc.ca/schemaball)

**USENIX**

**LINUXANDGENOMICS  
TWO REVOLUTIONS**

**USE LINUX SIG**

**IBM SAN RAID-5**  
1 x 14 x 72 GB  
2 x 14 x 146 GB  
5 TB



**Sun L700 LTO2 robot**  
396 slots 4 drives



**NetApp FAS960/R150**  
7 x 14 x 146 GB  
6 x 10 x 250GB  
26 TB



**IBM Bladecenter**  
14 x 2 x 2.4 GHz Xeon  
1.5 GB  
RH 9 2.4.20



**IBM x440**  
3 x 8 x 1.5 GHz Xeon  
8 GB  
SUSE Enterprise 8 2.4.21

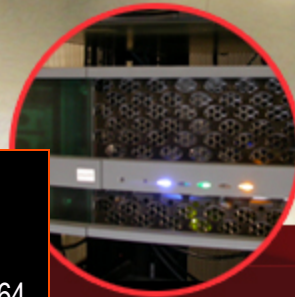


**IBM x330**  
90 x 2 x 1.0-1.4 GHz P3  
1 GB  
RH 9 2.4.20

**AMD Opteron**  
40 x 2 x 2.0 GHz  
2 GB  
RH 9 2.4.20



**NetApp C6100**  
DNFS NetCache  
14 x 36GB



**AMD Opteron**  
4 x 1.4GHz  
32 GB  
SUSE 9 2.4.21 x86\_64

**USENIX** LINUXANDGENOMICS  
TWO REVOLUTIONS  
USE LINUX SIG



MAKE OF LITTLE PLANS.  
THEY HAVE OF MAGIC TO  
STIR MEN'S BLUFFS...

...MAKE BIG PLANS.

DANIEL BURNHAM  
ARCHITECT



***\*nix | systems | db | programming | automation***

**[www.bcgsc.ca/about/employment](http://www.bcgsc.ca/about/employment)**

**talk slides [mkweb.bcgsc.ca/sars/usenix](http://mkweb.bcgsc.ca/sars/usenix)**



**CANADA'S MICHAEL SMITH GENOME SCIENCES CENTER [WWW.BCGSC.CA](http://WWW.BCGSC.CA)**