

# BIOINFORMATICS & HPC CORE

PRINCESS MARGARET CANCER CENTRE, UHN

Who are we and what do we do?

Natalie Stickle & Zhibin Lu

# OUTLINE

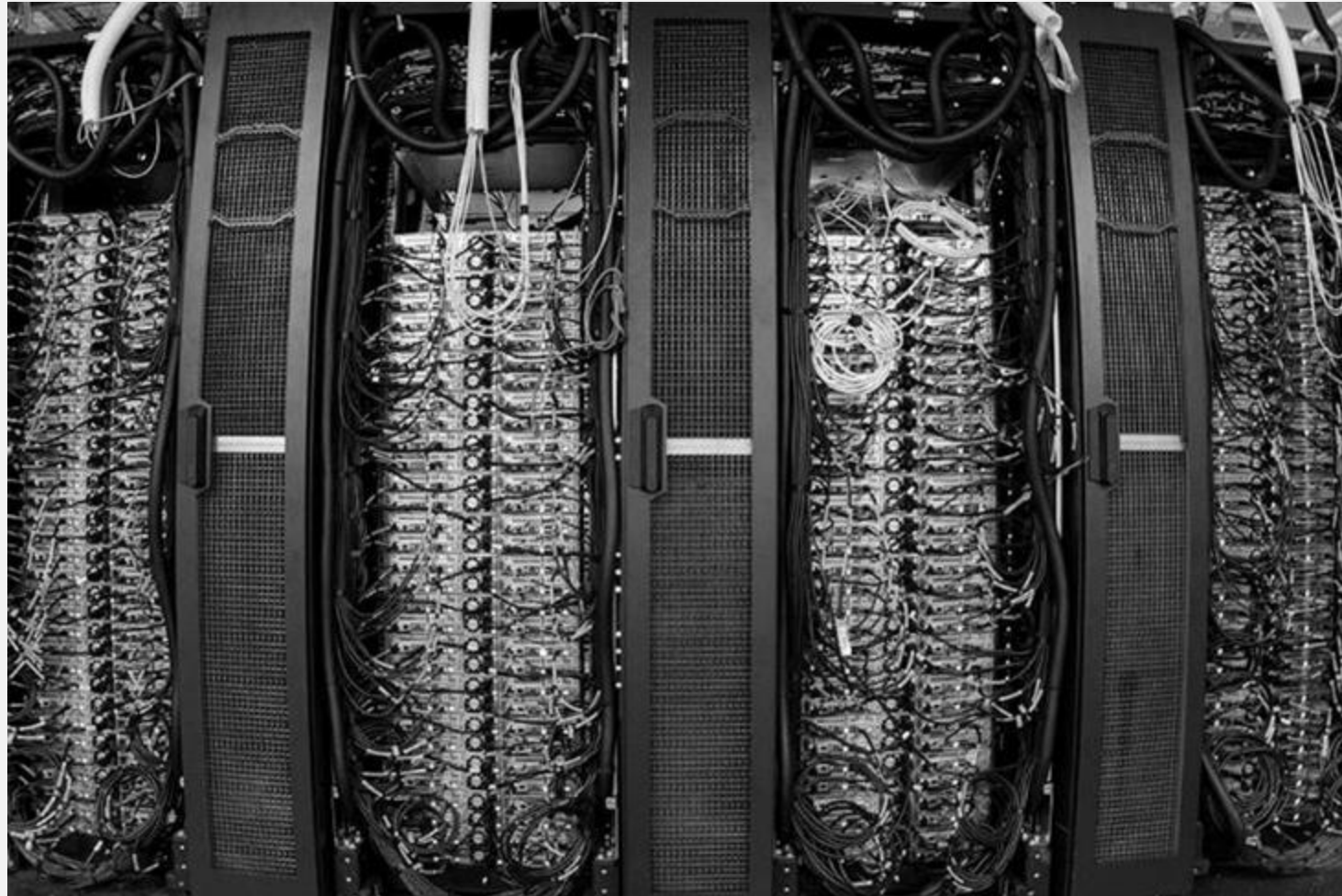
- What is Bioinformatics?
- What is HPC?
- The data problem.
- Bioinformatics & HPC Core mandate
- Internal HPC capabilities
- HPC4Health
- Core Services
  - Bioinformatic Analyses
  - Single Cell Genomics
  - Web sites, wikis and CMS
- Alignment within Canada and beyond

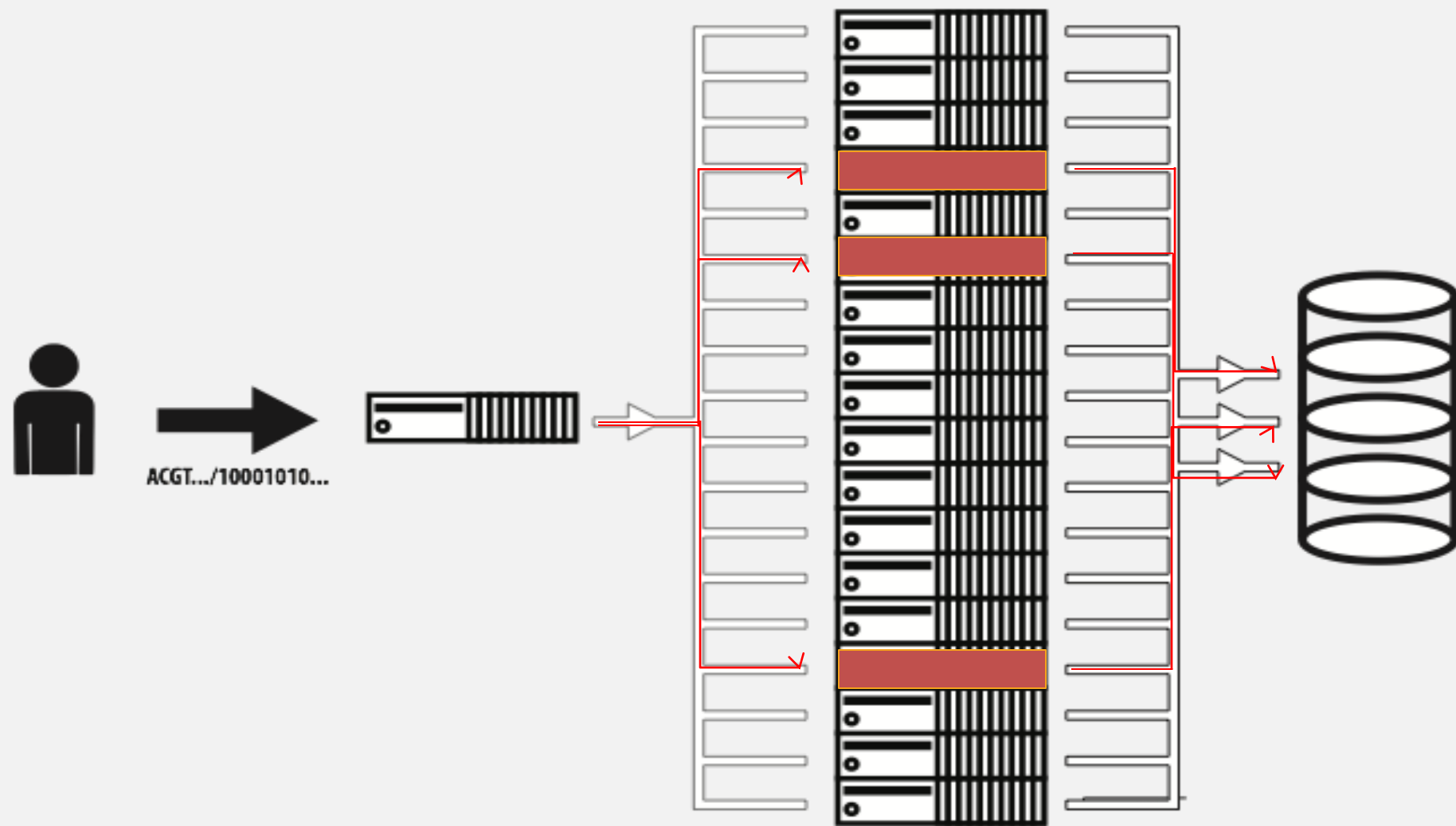
# WHAT IS BIOINFORMATICS?

- A “catch-all” phrase that means different things to different people
- Intersection of **computer science**, **mathematics** and **statistics** - application of these to problems in areas of basic biology such as **genomics**, **proteomics**, molecular biology, etc.
- Often employs state-of-the-art techniques in **artificial intelligence** and **machine learning** to create meaning and extract useful information from extremely large datasets
- Usually needs access to **High-Performance Computing (HPC)** resources due of size of data to be analyzed and complexity of the processing

# A SHORT INTRODUCTION TO HPC

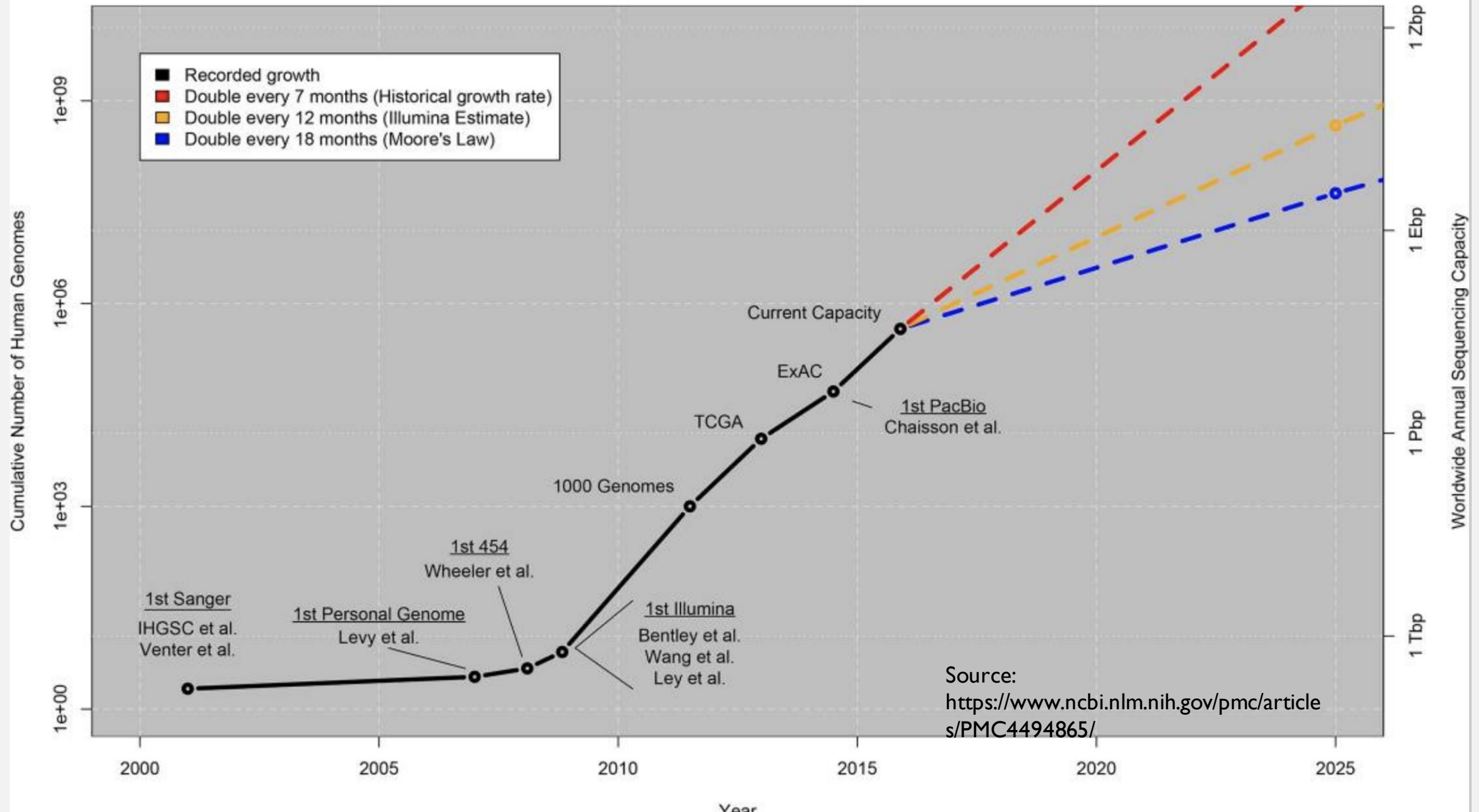
- Multiple servers connected together and integrated (called a cluster)
  - Each server in the cluster is called a node
- High speed connections
- Storage is optimized for fast I/O
  - eg. 80 Gb/s server vs. MAX 6 Gb/s desktop
- Scheduler software responsible for access to nodes



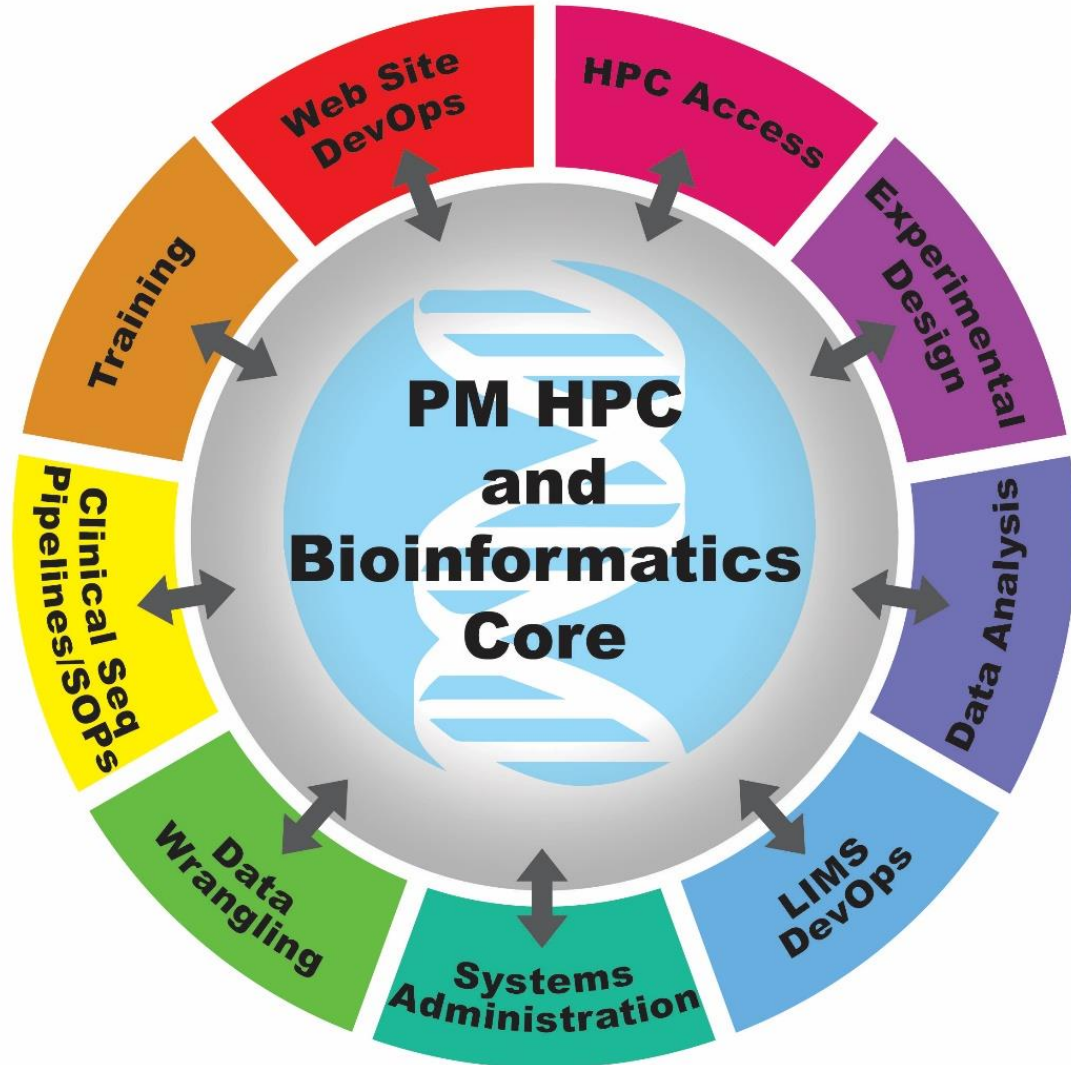


# THE DATA PROBLEM

## Growth of DNA Sequencing



## OUR CORE SERVICES

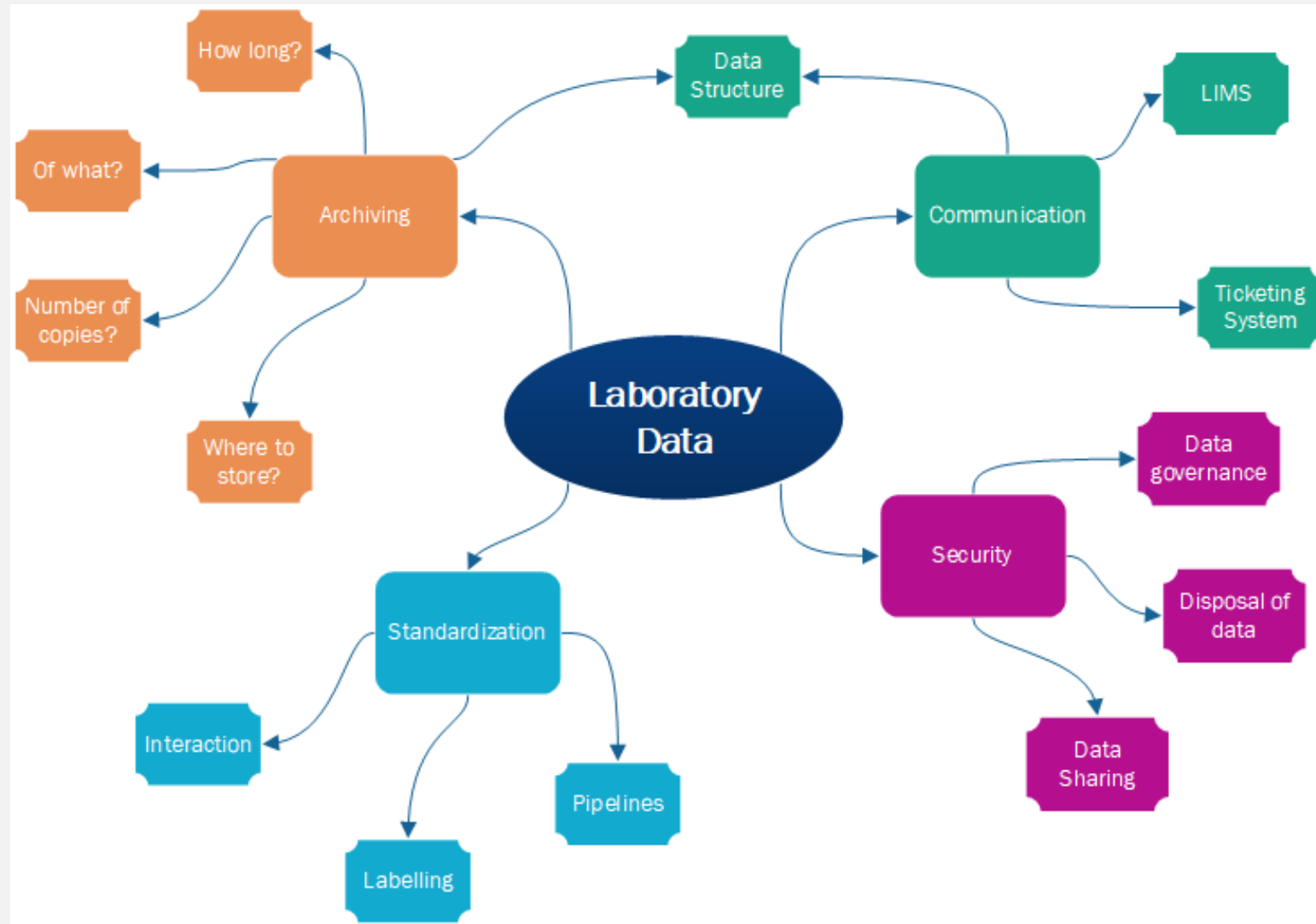




# PURPOSE OF BIOINFORMATICS CORE

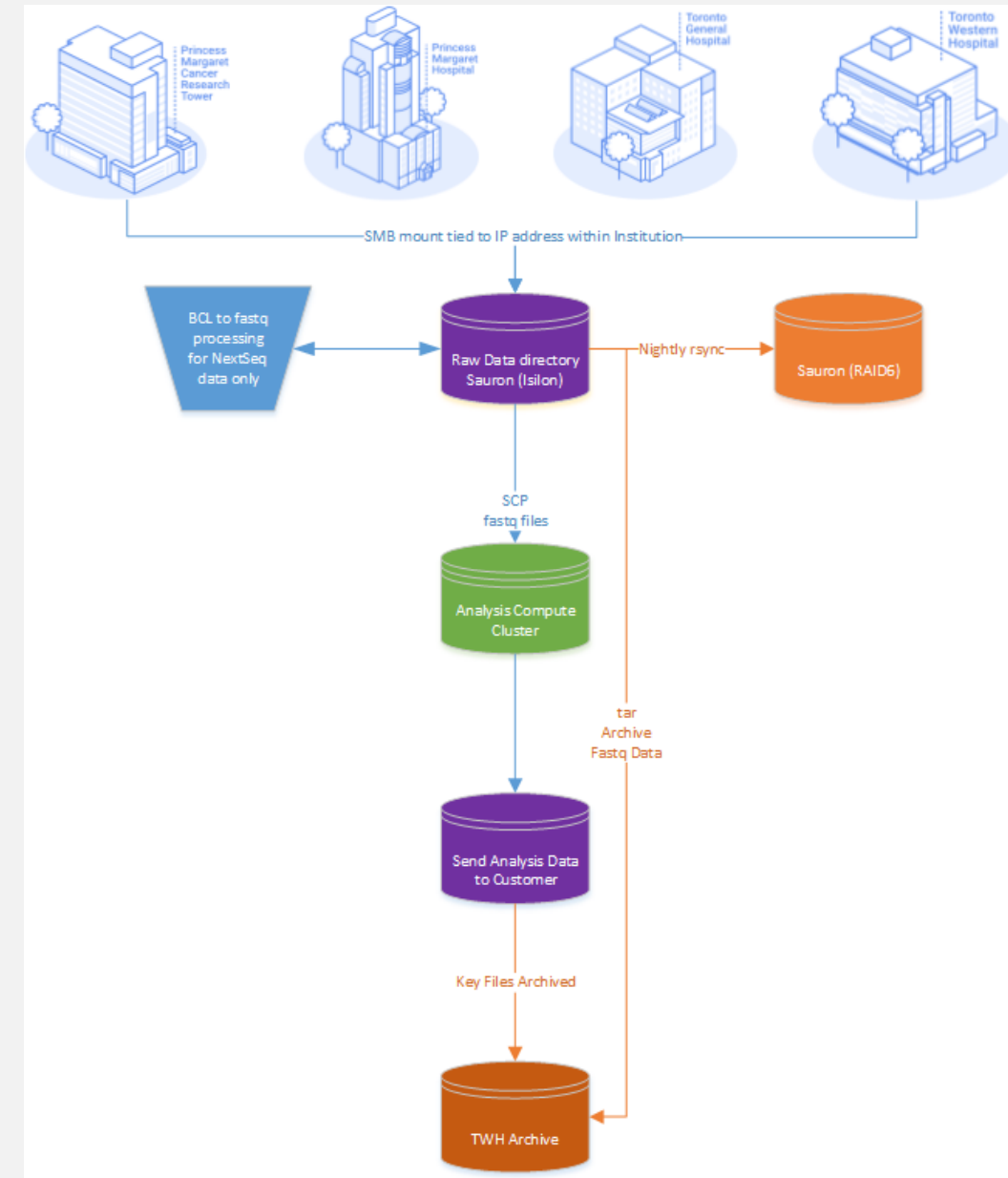
Factors to consider:

- Experimental Design
- Best practices for analysis from literature
- Need for HPC for large scale calculations
- Vast array of possible software options
- Implementation of governance for data (storage, security, access, sharing, etc.)



# NGS DATA STORAGE & SECURITY

- All sequencers write directly to the core's storage
- Run is complete – data is moved to directory where initial processing can be completed – then sent to customer or further analyzed by the bioinformaticians in the core
- All raw data is backed up every day while in the initial storage directories
- Raw data is archived after initial run and analysis is completed
- If analysis is completed by the core's bioinformaticians key analysis files will also be archived.



## INTERNAL HPC CAPABILITIES

- 1.1 PB of high speed Infiniband distributed storage (Isilon)
- 200+TB of general purpose scratch (used for processing raw sequence data)
- ~1.4 PB of spinning disk backup at TWH (triple mirrored)
- 656 cores of CPU (384 Ivy Bridge/272 Sandy Bridge @ between 10 and 16GB of RAM/core)
- 10 Gigabit ethernet
- GPGPU node (nVidia Tesla M2090)
- Isolated web servers for web apps
- >120 users

## INTERNAL HPC CAPABILITIES CONT'D

- >200 TB of public data (some require approval eg.TCGA)
- Dozens of OS applications (Centos 6/7 series, module system)
- SGE, Slurm scheduler with multiple queues
- Custom scripts written for general use (pipelines for RNA-seq and DNA-seq)
- Command Line ONLY!

# HPC4HEALTH

- ~2500 CPU cores latest Intel Xeon chips (>5000 total)
- 10 Gigabit ethernet
- 256 GB RAM/node
- 1.2 PB high-speed Isilon Storage
- ~20 Gb/s to data centre
- GPGPU nodes (4 nVidia Tesla P100 12GB each node)
- Torque/moab, Slurm scheduler

# WHAT IS THE HPC4HEALTH?

[Home](#)[Overview](#)[About ▾](#)[Contact Us](#)

## MISSION STATEMENT

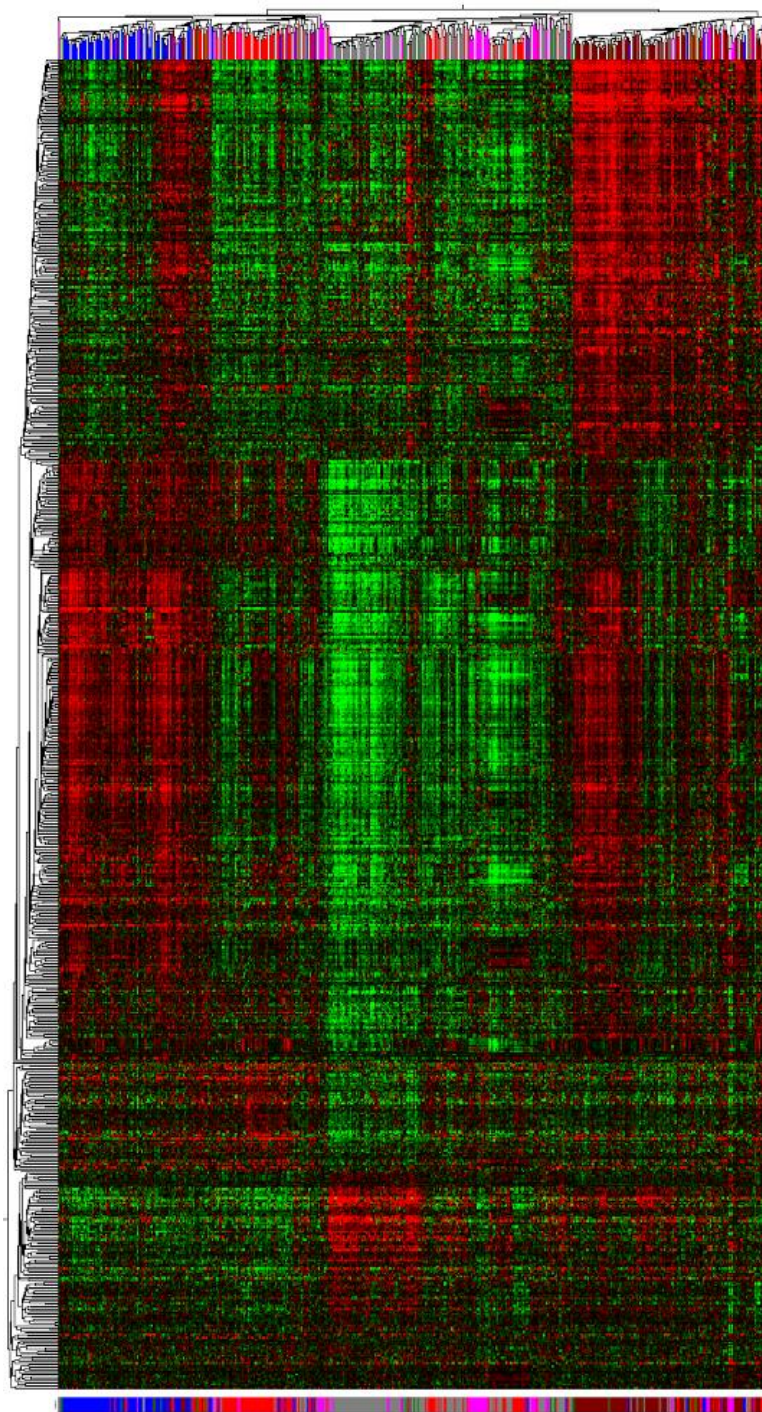
[Computing for Health](#)

From genomics to medical imaging, almost every discipline in health-care is dealing with a "Data Deluge" of information. Translating this into something that will ultimately benefit patients requires massive amounts of computation and storage in an environment that is fast, secure, and run with efficiency. HPC4Health is a consortium of health providers who are working together to build the next-generation of compute engine for clinical research.

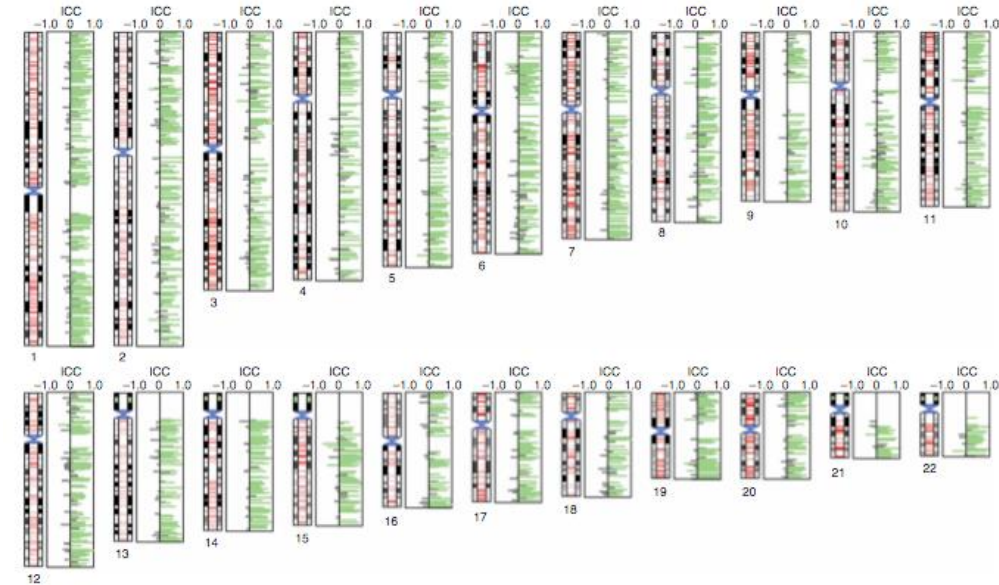


# BIOINFORMATIC ANALYSES

Expression Signatures in  
489 Ovarian Cancer  
Patients from TCGA



## Genome-wide Methylation Patterns



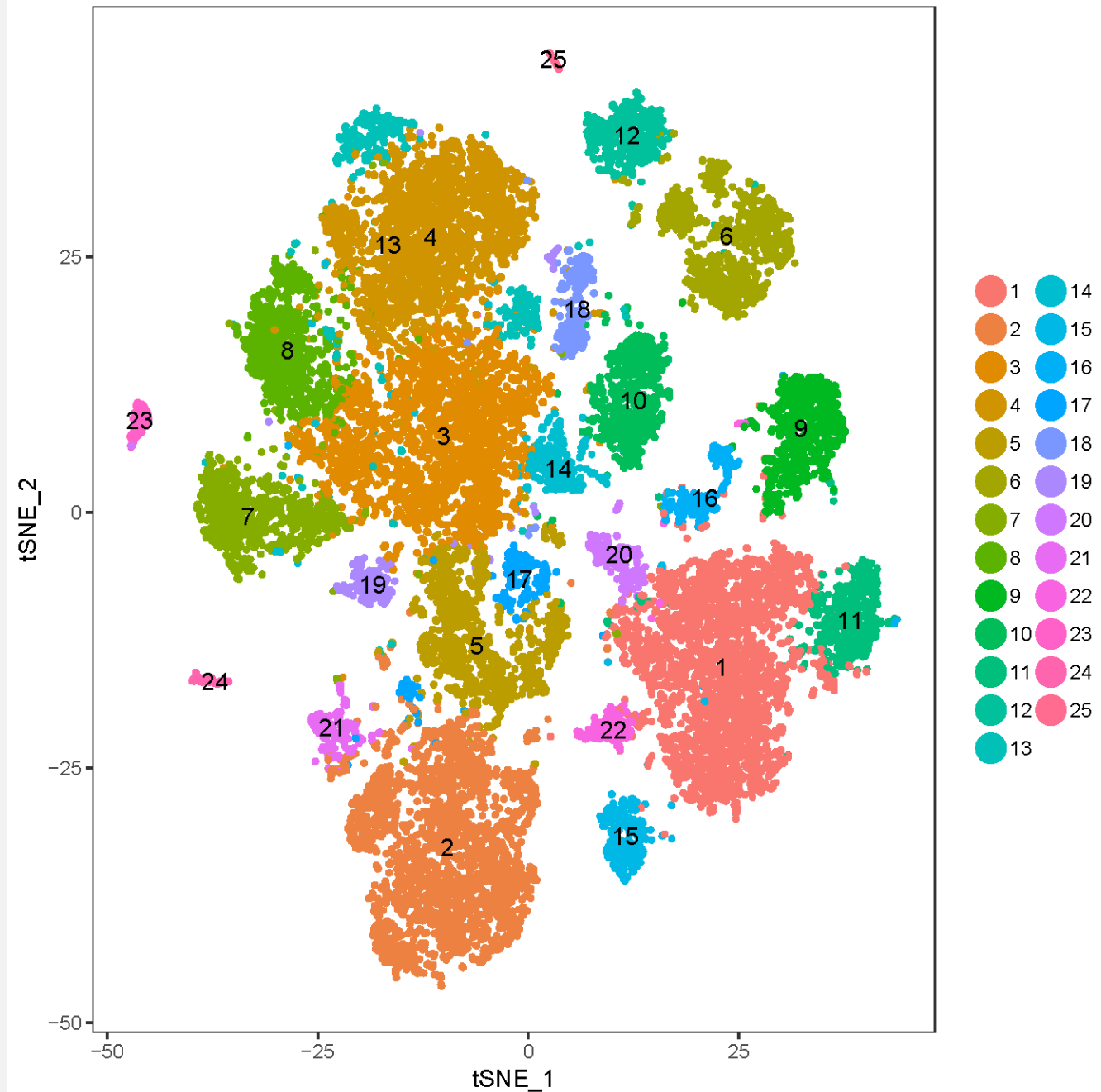
## ANNOTATED SEQUENCING RESULTS

Sequence Report Table for Patient74nf0.hardfiltered.noFail.avinput.hg19_multianno																							
Files included:Patient74nf0.hardfiltered.noFail.avinput.hg19_multianno.txt																							
Column Headings																							
Chr	Start	End	Ref	Alt	Func.refGene	Gene.refG	ExonicFunc.refGene	AAChange.refGene		snp137	cosmic64	genomi	1000g2012apr_a	avsisf	LJB_PhyloP	LJB_Pi	LJB_SiF	LJB_Si	LJB_PolyPhe	LJB_Pr	LJB_GERP++	LJB_Mut	LJB_MutationTaster_Pre
Patient74nf0.hardfiltered.noFail.avinput.hg19_multianno.txt																							
chr1	205073034	205073034	G	C	exonic	RBBP5	nonsynonymous SNV	RBBP5:NM_001193272:exon5:c.C473G;p.S158C,RBBP5:NM_001193272:exon5:c.C473G;p.S158C						0.13	0.999794	C	0.84	T	0.899	D	5.7		0.90917 D
chr1	114438951	114438951	A	G	exonic	AP4B1	nonsynonymous SNV	AP4B1:NM_001253852:exon8:c.T1439C;p.L480S,AP4B1:NM_001253852:exon8:c.T1439C;p.L480S	rs1217401			0.35		0.13	0.999183	C	0.82	T	0.0010	B	5.69		0.73428 P
chr1	44877932	44877932	G	C	exonic	RNF220	nonsynonymous SNV	RNF220:NM_018150:exon2:c.G163C;p.D55H	-						0.999769	C	0.99	D	0.529	P	5.62		0.96010 D
chr1	103468308	103468308	T	G	exonic	COL11A1	nonsynonymous SNV	COL11A1:NM_0080630:exon20:c.A1690C;p.N564H,COL11A1:NM_0080630:exon20:c.A1690C;p.N564H							0.998469	C	1.0	D	0.788621	NA	5.59		0.99853 D
chr1	113153597	113153597	T	C	exonic	ST7L	nonsynonymous SNV	ST7L:NM_017744:exon3:c.A317G;p.H106R,ST7L:NM_017744:exon3:c.A317G;p.H106R	rs143476452						0.998408	C	1.0	D	0.981	D	5.55		0.0 N
chr1	98348885	98348885	G	A	exonic	DPYD	nonsynonymous SNV	DPYD:NM_000110:exon2:c.C85T;p.R29C,DPYD:NM_000110:exon2:c.C85T;p.R29C	rs1801265				0.77	0.17	0.999726	C	0.98	D	0.713714	NA	5.49		0.71975 NA
chr1	84956143	84956143	C	T	exonic	RPF1	nonsynonymous SNV	RPF1:NM_025065:exon5:c.C604T;p.R202C		0.0014					0.999121	C	1.0	D	1.0	D	5.48		0.99999 D
chr1	46195375	46195375	T	C	exonic	IPP	nonsynonymous SNV	IPP:NM_001145349:exon4:c.A791G;p.T264R,IPP:NM_000589:exon4:c.A791G;p.T264R	rs28375469	ID=cosmic		0.30		0.19	0.998274	C	0.59	T	0.0	B	5.46		0.97077 P
chr1	115168600	115168600	A	T	exonic	DENND2C	nonsynonymous SNV	DENND2C:NM_198459:exon2:c.T6A;p.D2E,DENND2C:NM_198459:exon2:c.T6A;p.D2E	rs7541738			0.13		0.2	0.998841	C	0.11	T	0.07	B	5.36		0.12126 P
chr1	39908506	39908506	G	A	exonic	MACF1	nonsynonymous SNV	MACF1:NM_012090:exon75:c.G13048A;p.A4350T	rs587404			0.35		0.63	0.981078	C	0.42	T	0.235978	NA	5.35		8.0E-6 P
chr1	38265496	38265496	C	A	exonic	MANEAL	nonsynonymous SNV	MANEAL:NM_152496:exon2:c.C329A;p.S110Y,MANEAL:NM_152496:exon2:c.C329A;p.S110Y							0.998899	C	1.0	D	0.895	D	5.27		0.99610 D
chr1	11150621	11150621	C	G	exonic	EXOSC10	nonsynonymous SNV	EXOSC10:NM_001001998:exon6:c.G748C;p.E250Q,EXOSC10:NM_001001998:exon6:c.G748C;p.E250Q						0.21	0.999348	C	0.74	T	0.024	B	5.23		0.93950 D
chr1	197087086	197087086	G	A	exonic	ASPM	stopgain SNV	ASPM:NM_001206846:exon17:c.C398T;p.Q1300X,ASPM:NM_001206846:exon17:c.C398T;p.Q1300X							0.999584	C	0.90383	NA	0.735479				



# SINGLE CELL GENOMICS

- Princess Margaret Genomics Centre is one of 2 in North America (only Canadian) 10x Genomics certified service provider
- Our core has extensive experience with this data



WEB SITES,  
WIKIS,  
CMS

cd13

Symbol  
Synonyms  
Description  
Species

www.cancergen



The Princess Margaret Cancer Centre  
Ovarian Cancer  
Prevention Initiative

Home Our Initiative Ovarian Cancer Genetics Resources Login / Sign up



#### Our Initiative

Facts about our ovarian cancer prevention initiative & how to participate if you are eligible

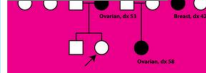
#### What is Ovarian Cancer?

Learn more about the biology and treatment of ovarian cancer



#### Genetic Counselling

Learn facts about genetic counselling and testing



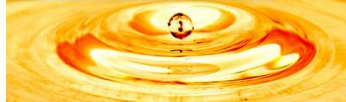
#### How can I help?

Volunteer, spread the word, or donate!



#### Impact of Prevention

Learn about how this initiative has the potential to save lives



READ MORE

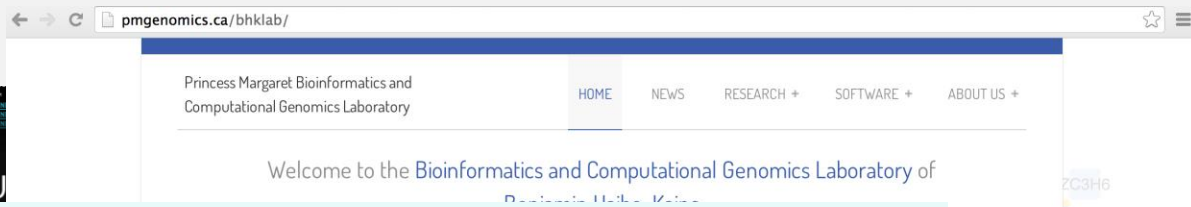
The **Cancer Genomics Program** at the **Princess Margaret Cancer Centre** is dedicated to advancing personalized cancer medicine through the identification of genetic mutations and molecular mechanisms that drive cancer.

Our ultimate vision is to have all patients who may benefit from this approach undergo comprehensive molecular characterization of their cancer by 2018.

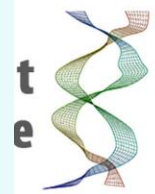
Collaboration

Participation

Communication



a poster



#### Started

reliable and cost-effective genomic solution?


help you with preliminary experimental and consultation.

#### Client Portal

Access your project and obtain results and data via our secured **User Portal!**



# WIKIS



PRINCESS  
MARGARET  
GENOMICS  
CENTRE

- ▼ Home
  - Main
  - Event Calendar
- ▼ Data Analysis
  - High Throughput Sequencing
  - Microarray
  - Software
- System
- ▼ Web
  - Web Sites
  - Web Servers
  - Software
  - Document
- Database
- Repository
- Backup
- Help
- Tools

Carl Talk Preferences Watchlist Contributions Log out

Pmgc Discussion Read Edit View history ☆ Search

## Pmgc:HTS RNASeq


**Contents** [hide]

- 1 Tuxedo Package
  - 1.1 Raw Data Structure
  - 1.2 Tools and Reference
  - 1.3 Wrapper
  - 1.4 Step by Step
    - 1.4.1 Alignment
    - 1.4.2 RNAQC
    - 1.4.3 Transcript Assembly
    - 1.4.4 Merged GTF
    - 1.4.5 quantify expression levels
    - 1.4.6 Differential expression
    - 1.4.7 Example

### Tuxedo Package

 [edit]

Workflow using the tuxedo package can be found at Cufflinks website.



- scripts and wrappers are located under /mnt/work1/software/ocigcscripsts/production/tophat2\_wrapper
- There is a "qsub" sub-directory under Output directory, and they are real scripts to perform the jobs. You can repeat the job by "qsub" each script.
- By default, there is a log file under output directory called 'process\_log.txt'.

### Raw Data Structure

 [edit]

- Fastq files for each sample should be in a directory named "Sample\_<sample name>"
- Fastq files from multiple lanes for the same sample should be put into the same directory.
- Fastq files can be gzipped or unzipped.
- Fastq file name format "<sample name>\_<barcode>\_<lane number>"

Example:

```
fastq
├── Sample_D001-Blood
│   └── D001-Blood_AGAGTC_L001_R1.fastq
```

# TRACK EVERYTHING!

## 141205\_SN1080\_0189\_BC6345ACXX [\[edit\]](#)

---

- 51 cycles + 7 indexing cycles
- [sample sheet](#)
- bcl2fastq with mismatch=1, barcode=6
- barcode error for sample 3088 on lane 6. Changed to CAGATC and re-ran the conversion.
- Project\_DeCarvalho\_Hung (Roxana)
  - lane 1-4
  - fastq files were copied into /mnt/work1/users/decarvalhogroup/141205\_SN1080\_0189\_BC6345ACXX\_Hung on Dec 10
- Project\_DeCarvalho (Tiago)
  - lane 5-8
  - fastq files were copied into /mnt/work1/users/lupiengroup/141205\_SN1080\_0189\_BC6345ACXX on Dec 10

## 141204\_SN1080\_0188\_AC5V7DACXX [\[edit\]](#)

---

- 51 cycles + 8 indexing cycles
- [sample sheet](#)
- bcl2fastq with mismatch=1, barcode=6
- Project\_Dick
  - lane 1-6,8
  - fastq files were copied into /mnt/work1/users/lupiengroup/141204\_SN1080\_0188\_AC5V7DACXX\_Naoya. Emailed Naoya [tnaoya19760517@gmail.com](mailto:tnaoya19760517@gmail.com) and Nadia [npenrod@uhnresearch.ca](mailto:npenrod@uhnresearch.ca) on Dec 8.
- Project\_DeCarvalho\_Hung
  - lane 7
  - fastq files were copied into /mnt/work1/users/decarvalhogroup/141204\_SN1080\_0188\_AC5V7DACXX\_Hung

## 141120\_SN1068\_0163\_AC5FTVACXX [\[edit\]](#)

---

- 51 cycles + 9 indexing cycles
- [sample sheet](#)
- bcl2fastq with mismatch=1, barcode=6
- Project\_Lupien
  - lane 1-4
  - fastq files were copied into /mnt/work1/users/lupiengroup/141120\_SN1068\_0163\_AC5FTVACXX\_ken on Nov 25.
- Project\_Dick
  - lane 5-8
  - fastq files were copied into /mnt/work1/users/lupiengroup/141120\_SN1068\_0163\_AC5FTVACXX\_Dick on Nov 25.

# ALIGNMENT - CANADIAN BIOINFORMATICS

- Canadian Centre for Computational Genomics (C3G)
- Canadian Bioinformatics Workshops (CBW)
- Compute Canada Bioinformatics Helpdesk
- CanDIG (<http://CanDIG.github.io>)
- AACR Project GENIE

# BIOINFORMATICS & HPC CORE

## Director

Carl Virtanen<sup>RC</sup>

## Manager

Natalie Stickle<sup>RC</sup>

## System Administration

### Manager

Zhibin Lu<sup>RC</sup>

### Analyst

Qun Jin<sup>RC</sup>

## Bioinformaticians

### Clinical

Irene Chae<sup>C</sup>

Roozbeh Dolatshahi<sup>C</sup>

Gregory Downs<sup>C</sup>

### Research

Richard De Borja<sup>RC</sup>

Abdellatif Daghrach<sup>RC</sup>

Neke Ibeh<sup>R</sup>