



Scalable methods for data analysis and sharing in genomics

Guillaume Bourque

Dept. Human Genetics, McGill University

McGill Genome Center

Canadian Center for Computational Genomics (C3G)



@guilbourque

RDC Webinar Series

26 Novembre, 2020

Canadian Center for Computational Genomics (C3G)



*C3G is a Genome Canada Genomics Technology Platform (GTP) launched in 2015 that provides **bioinformatics analysis, HPC services and software solutions** to the life sciences research community.*



*Guillaume Bourque Ph.D.
Director, C3G Montreal*



*Michael Brudno Ph.D.
Director, C3G Toronto*



GenomeCanada

Large genomic projects at C3G



Cancer:

PROFYLE/Signature

MCC/QCC

DHDP (w/ Imagia)

Marathon of Hope

...

Infection/Immunity:

CoVSeQ/CanCOGeN

BQC19

MI4

Recodid

...

Epigenomics/genomics:

IHEC

4DN (NIH)

LSARP (Jabado)

CGEn

BRIDGET

Secure Cloud (w/ CQ)...

Software/Others:

CanDIG

GenAP

EpiShare

ClinDIG

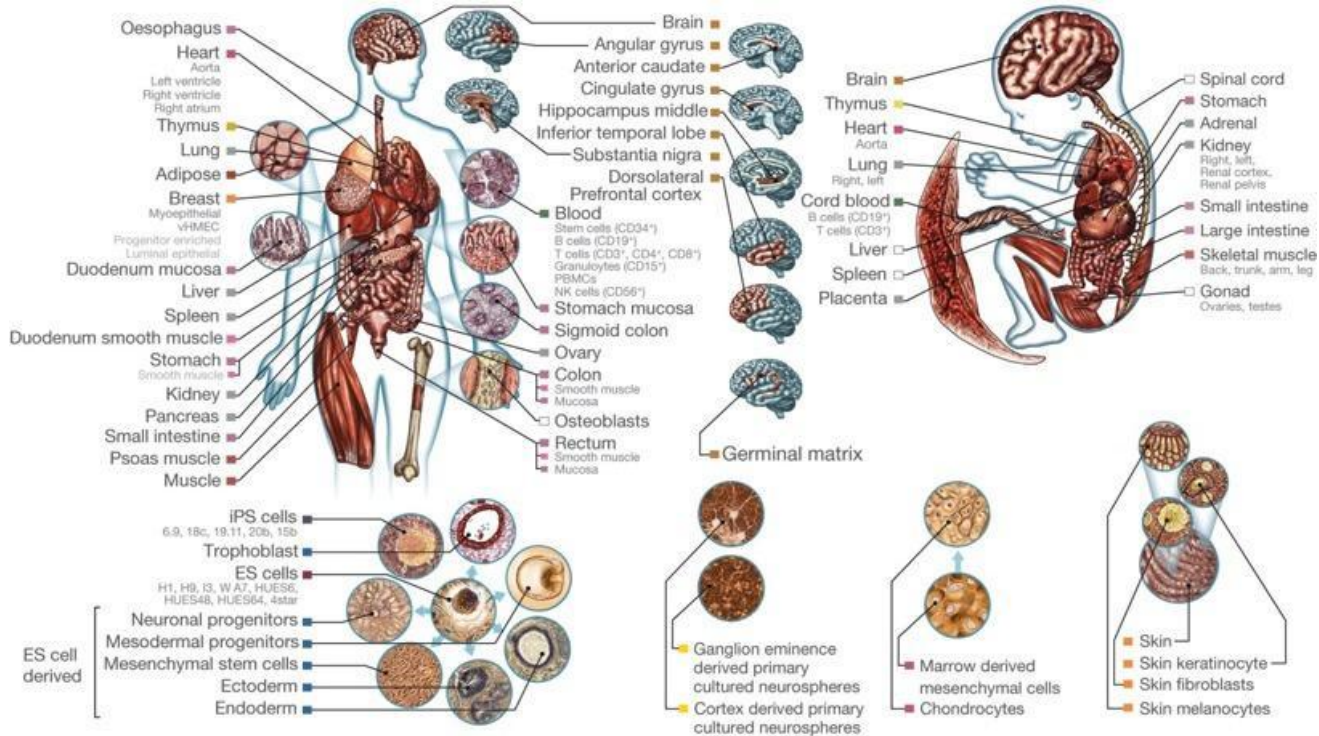
CHORD

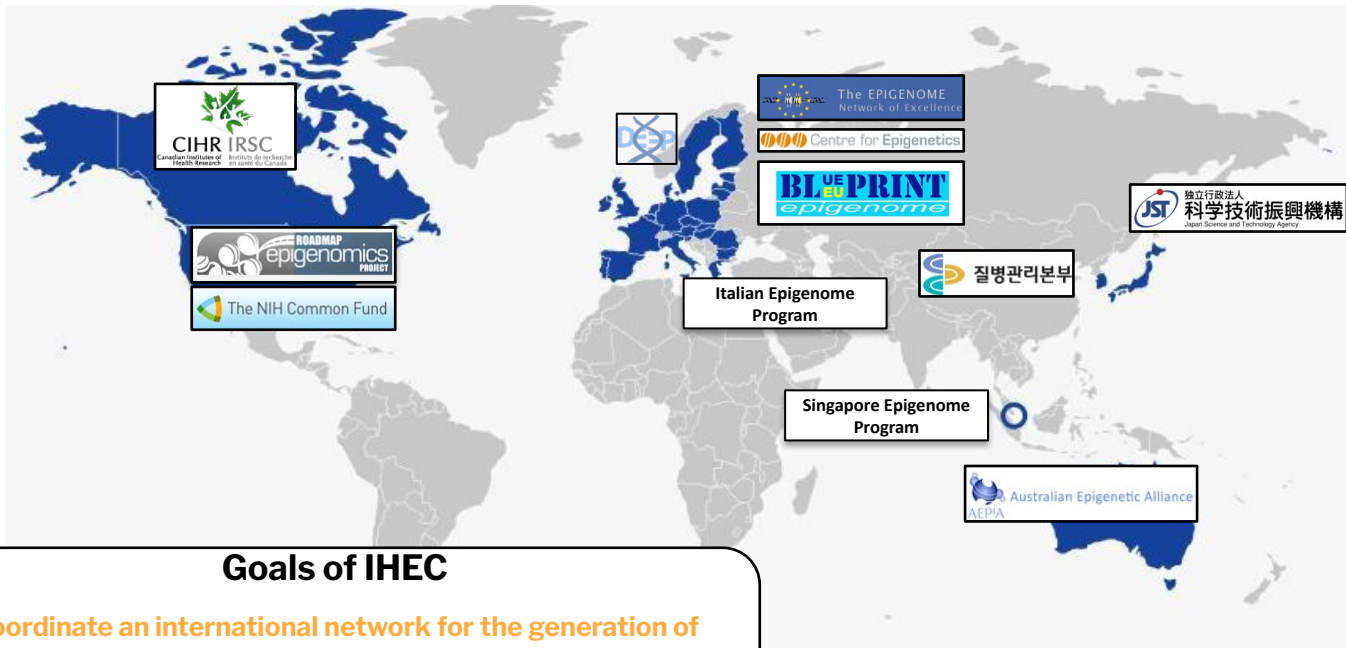
CINECA

Covid19 Resource Canada

...

One genome... Many epigenomes





Goals of IHEC

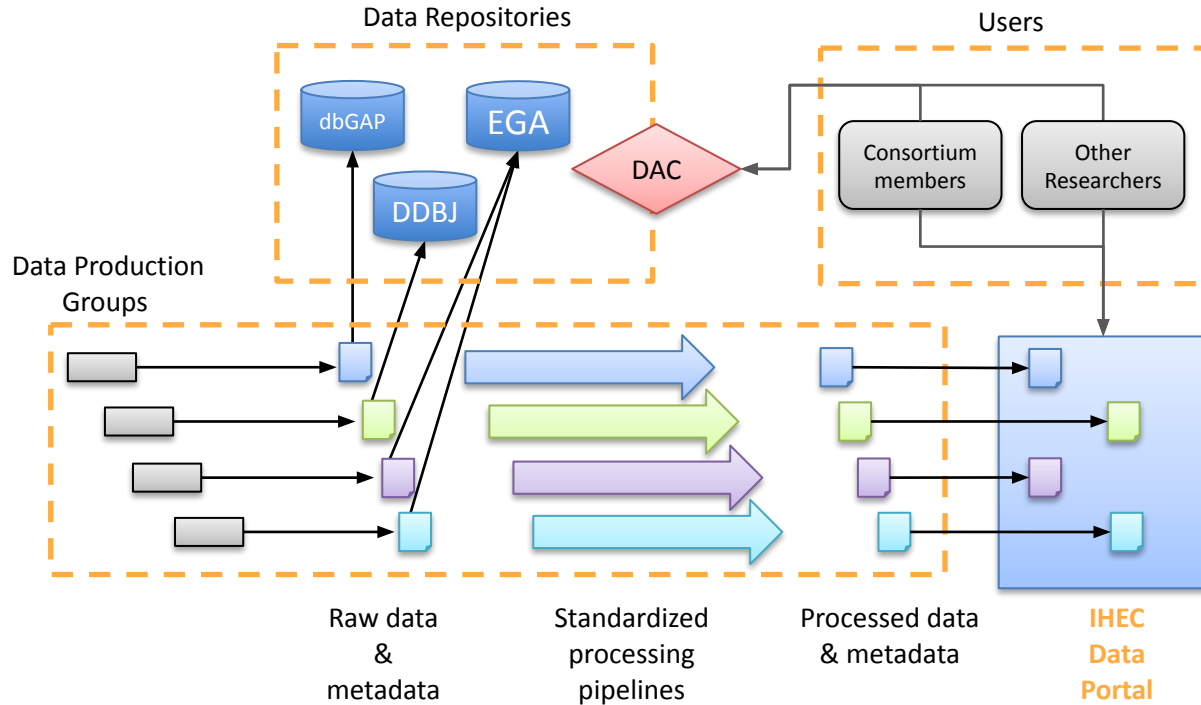
Coordinate an international network for the generation of 1000 or more reference epigenomes for a broad spectrum of human cell types and a wide range of developmental stages, laying the foundation to study the epigenetic mechanisms of human diseases.



IHEC-generated data

- Data generated by IHEC is aimed to be available for everyone's own research, following the consortium's Open Science mission.
- Human data produced by most IHEC members has different levels of access:
 - Controlled access data
 - Potentially personally identifiable
 - Access is limited to protect study participants
 - Archived at repositories such as EGA and DDBJ
 - Obtaining the data involves applying for access
 - Fully public data
 - Meta-information considered non-personally identifiable
 - Processed data (analysis pipelines output)

IHEC data integration and sharing strategy



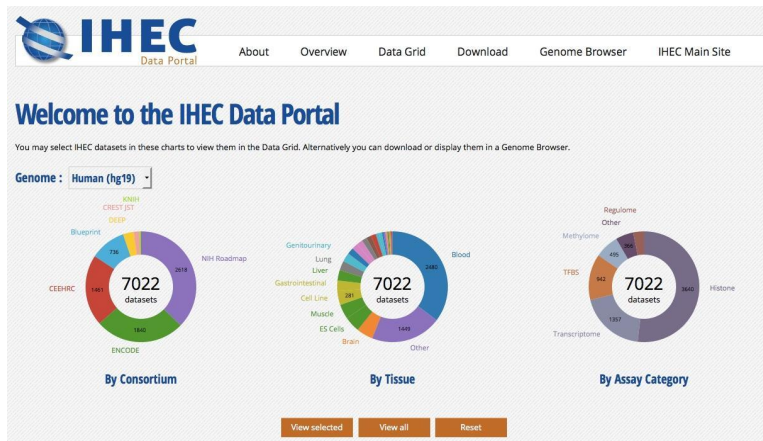


IHEC Data Portal

- Launched in June 2014, <http://epigenomesportal.ca/ihec>
- Includes:
 - Over **10,800** human epigenomic datasets (hg19 and hg38)
 - Over **280** mouse and primate datasets
 - Over **>290** full reference epigenomes
 - Data from: Blueprint, CEEHRC, CREST, DEEP, ENCODE, KNIH, NIH Roadmap



David
Bujold



Bujold et al.
Cell Systems
2016



The epiATLAS project

- IHEC Integrative Analysis workgroup initiative
- Initial goals:
 - a. establish an analyst-friendly and widely sharable compendium of **quality-controlled, consistently processed, reference epigenomic maps** from all areas of IHEC.
 - Prepare a gold standard dataset
 - Develop ways to ease access to the raw data
 - Improve the overall experience of accessing and analysing IHEC data
 - b. initiate and support numerous hypothesis-driven as well as exploratory analysis projects based on the IHEC epigenome compendium.
 - c. coordinate the publication of the resulting analyses in the form of a flagship and companion papers.



Challenges of current model

There are multiple challenges bound to using controlled access data, even before getting to the bulk of the analysis!

- Obtaining access
 - Application to a Data Access Committee (DAC)
- Downloading
 - Getting the data from a controlled access repository
- Comparing datasets across projects
 - Metadata is often hard to collate across projects
- Analysing the data
 - Heavy use of resources

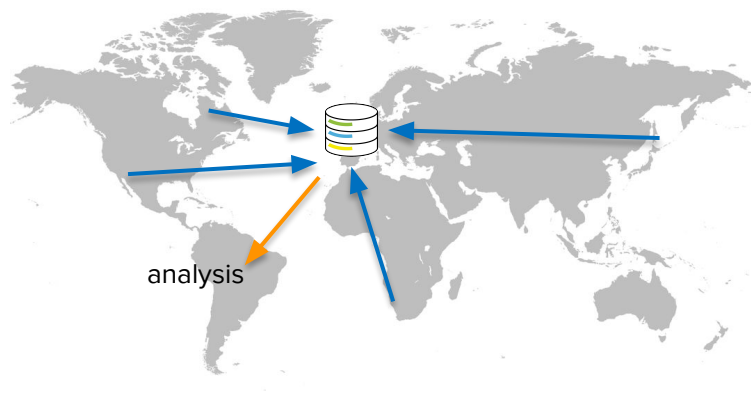
Accessing raw data

	Korea Epigenome Consortium	Blueprint Consortium	McGill EMC	BC Cancer Agency	DEEP (Germany)	Singapore Epigenome Project	CREST (Japan)
a) Acknowledgements							
b) Application renewal							
c) Confidentiality							
d) Data destruction							
e) Data usage							
f) Ethics review							
g) Evidence of PI's experience							
h) Intellectual property							
i) IT security requirements							
j) Naming of jurisdiction for disputes							
k) Liability exclusions & warranty limitations							
l) Process in the face of access requests							
m) Publication embargos							
n) Rapid publication							
o) Progress report							
p) References to laws							
q) References to policies							
r) Rules on student access to data							

Fig. 1 Terms and Conditions identified in IHEC DAAs.

Saulnier et al.
Scientific Data,
 2019

Open Research Data

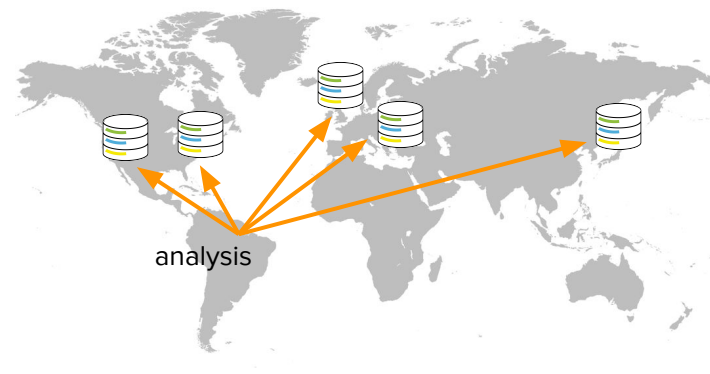


Aggregate data globally

Download, analyze locally

Continues for basic research

Healthcare Data with Research Use



Aggregate data locally (via VMs)

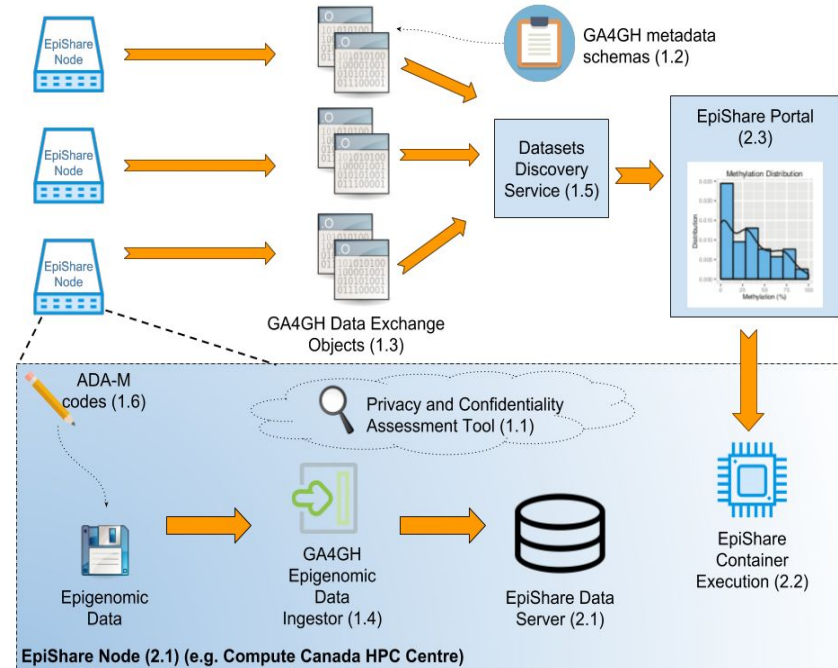
Collate analyses

New approach for research and healthcare

EpiShare



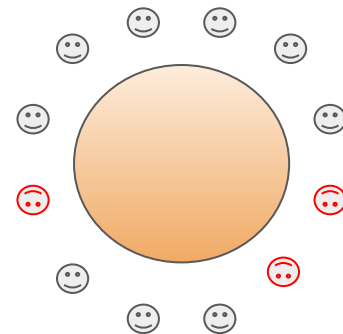
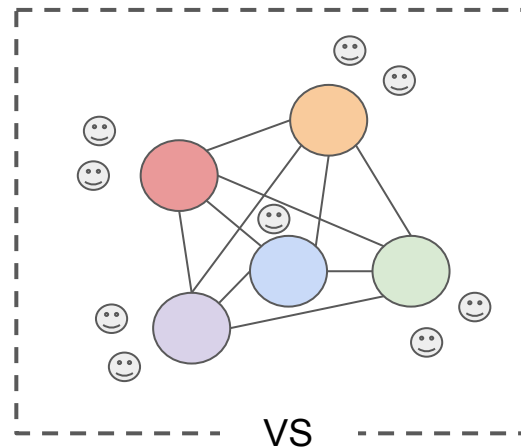
- Genome Canada funded project (2018-2021)
- Aims at extending the GA4GH APIs, etc. for epigenomic data
- Will create a resource to make data more easily discoverable
- Will enable the launch of multi-omics analyses on controlled-access datasets at their storage location
- Now a GA4GH Driver Project



Global Alliance
for Genomics & Health

Federated platform

- Enables **fine-grained control over data** while participating in data sharing
 - Shared normalization/QC pipelines
- **Resource allocations** from e.g. Compute Canada can be used **at a node level**
- Federation **maintain locality** while optionally **sharing in a network**
 - Some data may not be allowed to be stored in another jurisdiction





CanDIG/Bento platform

- Released as a self-hosted software infrastructure
 - SIGNATURE instance deployed at the Calcul Quebec Secure Cloud

- Data types-specific microservices include:
 - Clinical/phenotypic metadata
 - Genomic variants (VCF)
 - RNA-Seq expression data (in progress)

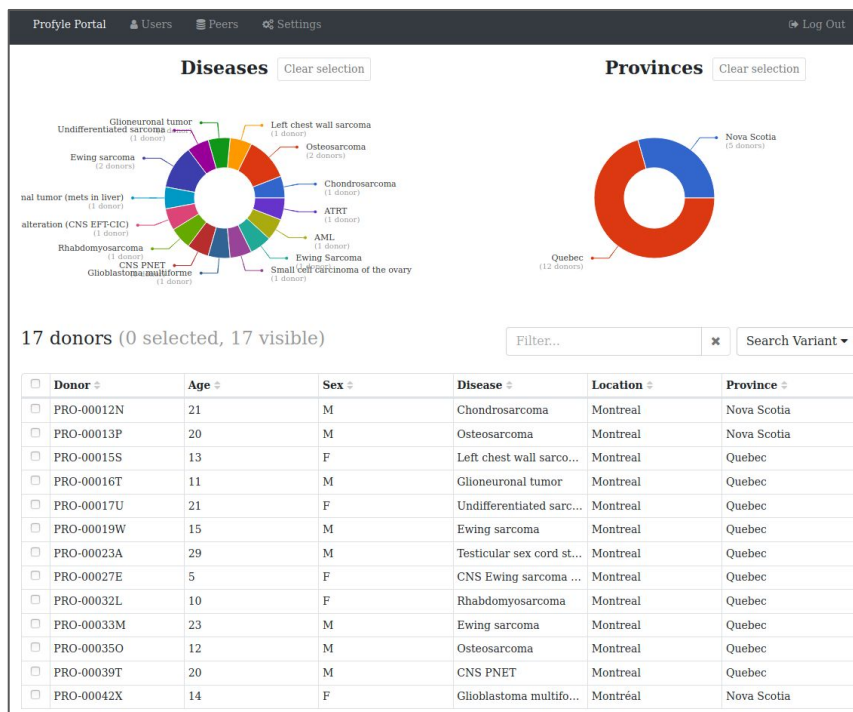
The screenshot shows the Bento dashboard interface. At the top, there is a navigation bar with the Bento logo and several menu items: SunBoard, Data Discovery, Data Explorer, Data Manager, and Peers. The main content area is titled 'Dashboard' and includes a subtitle 'Node status and health monitor'. Below this, there is an 'Overview' section with a table showing system metrics: Node URL (https://signature.c3g.calculquebec.ca/), Projects (1), Datasets (1), and Network Size (1). The 'Services' section contains a table with columns for Artifact, Name, Version, URL, Data Service?, and Status. Each row represents a different microservice, such as 'service-registry', 'drop-box', 'ves', 'Federation', 'notification', 'event-relay', 'metadata', 'variant', 'drs', and 'log-service'. The 'Data Service?' column has a checkmark for 'notification', 'metadata', and 'variant', and an 'X' for others. The 'Status' column for all services shows a green 'Active' indicator.

Artifact	Name	Version	URL	Data Service?	Status
service-registry	Bento Service Registry	0.4.1	https://signature.c3g.calculquebec.ca/api/service-registry	X	Active
drop-box	Bento Drop Box Service	0.4.1	https://signature.c3g.calculquebec.ca/api/drop-box	X	Active
ves	Bento VES	0.3.2	https://signature.c3g.calculquebec.ca/api/ves	X	Active
Federation	Bento Federation Service	0.7.0	https://signature.c3g.calculquebec.ca/api/federation	X	Active
notification	Bento Notification Service	0.3.0	https://signature.c3g.calculquebec.ca/api/notification	✓	Active
event-relay	Bento Event Relay	0.4.0	https://signature.c3g.calculquebec.ca/api/event-relay	X	Active
metadata	Metadata Service	1.3.3	https://signature.c3g.calculquebec.ca/api/metadata	✓	Active
variant	Bento Variant Service	0.4.0	https://signature.c3g.calculquebec.ca/api/variant	✓	Active
drs	CHORD Data Repository Service	0.2.0	https://signature.c3g.calculquebec.ca/api/drs	X	Active
log-service	Bento Log Service	0.1.0	https://signature.c3g.calculquebec.ca/api/log-service	X	Active



CanDIG demonstration project

Visualization tools and an overview of the project in the dashboard

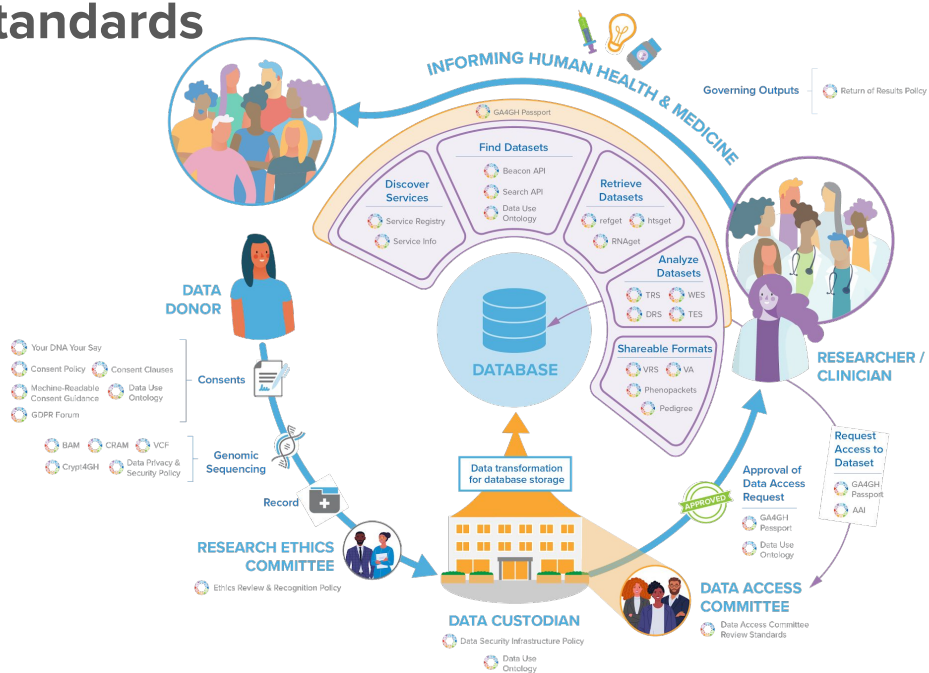




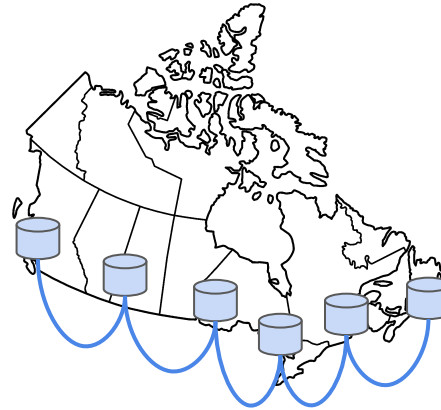
Global Alliance for Genomics & Health



GA4GH sets of standards

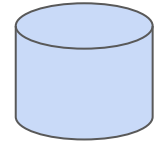


Data Federation using GA4GH in Canada



GA4GH
Standards

International Projects
and Consortia



SecureData4Health project



- A secure cloud infrastructure for analysis and sharing of genomic and health data
- 20M proposal submitted to CFI's Innovation fund 2020
- Team:

Guillaume Bourque	McGill University
Vincent Ferretti	Université de Montréal
Michael Brudno	University Health Network
Anne-Claude Gingras	Sinai Health System
Anna Goldenberg	The Hospital for Sick Children
Benjamin Haibe-Kains	University Health Network
Julie Hussin	Université de Montréal
Pierre-Étienne Jacques	Université de Sherbrooke
Bartha Knoppers	McGill University
Jacques Simard	Université Laval

SecureData4Health project



International Projects and Consortia

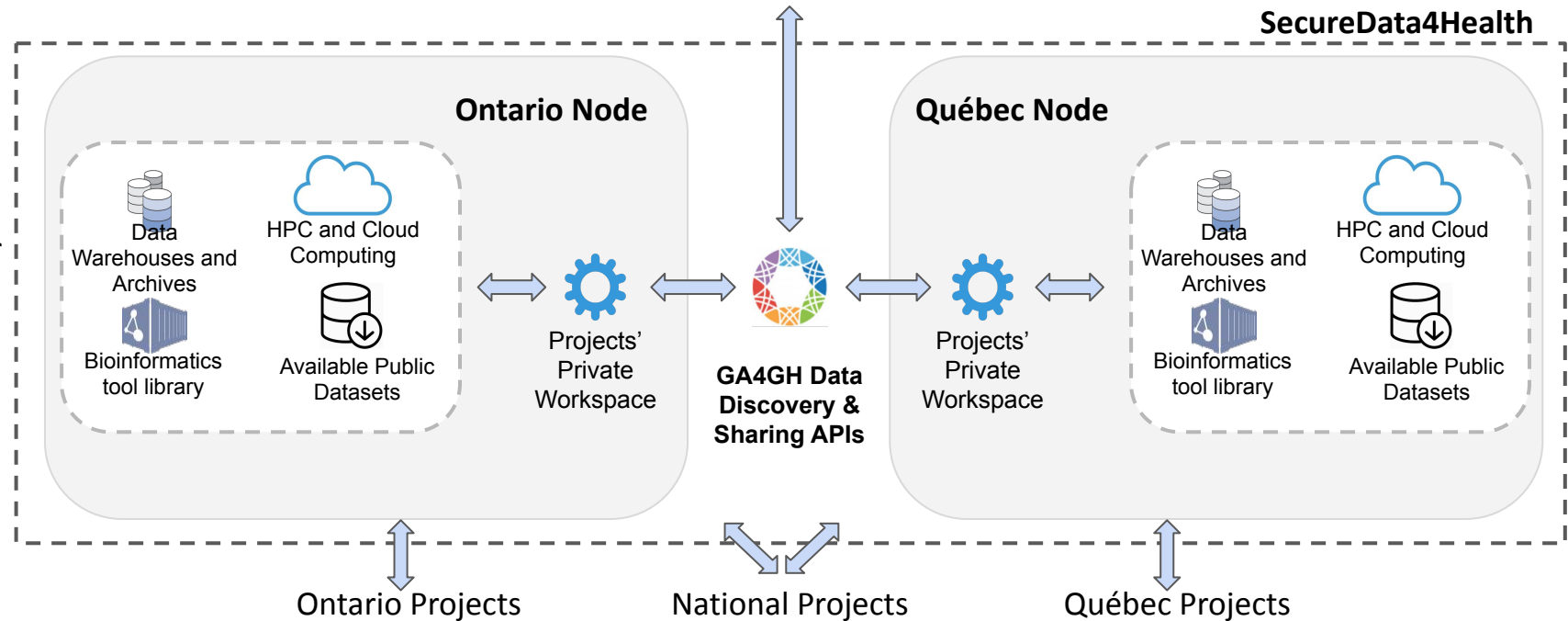
SecureData4Health

Ontario Node

Québec Node

Security

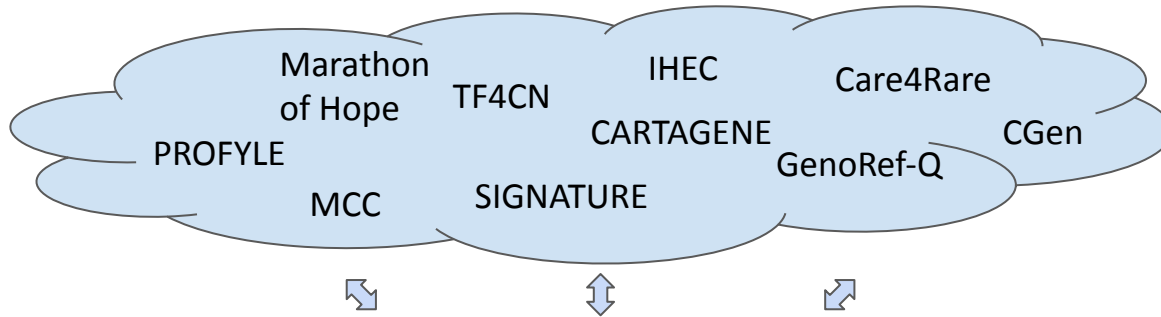
Privacy



Ontario Projects

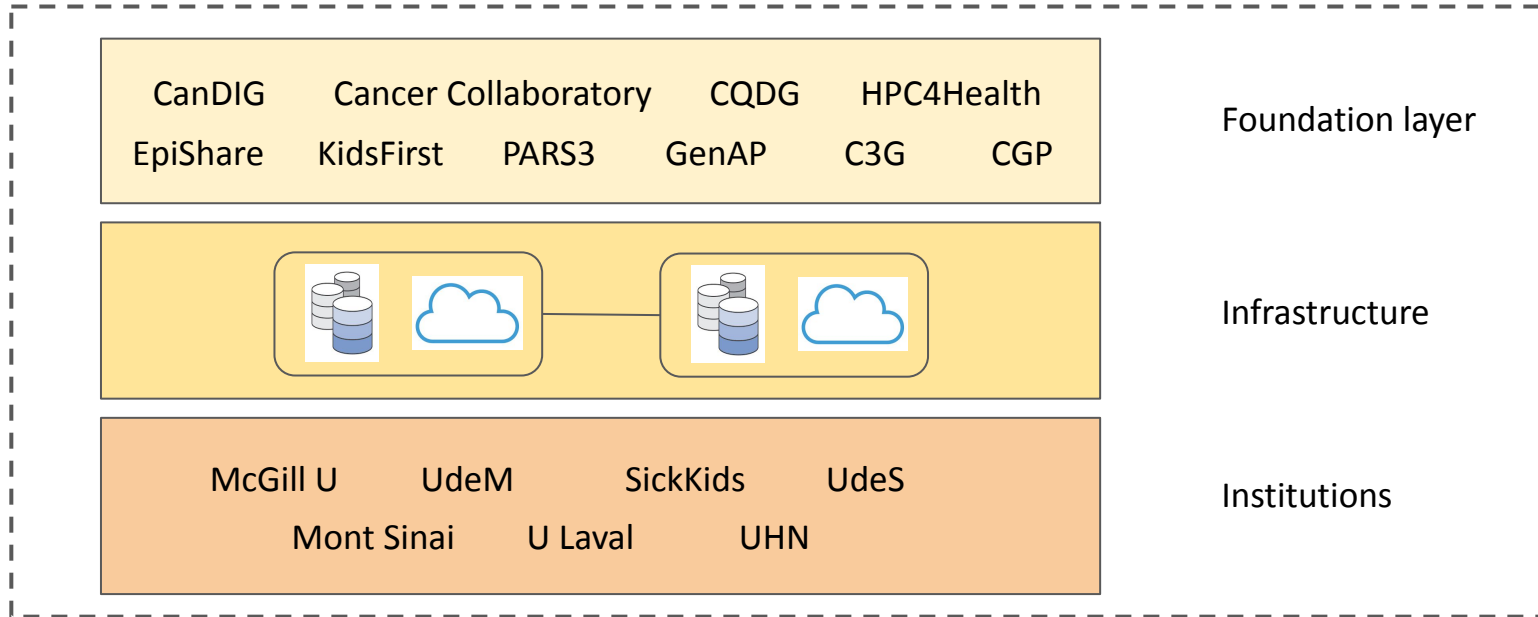
National Projects

Québec Projects



Projects

SecureData4Health



Foundation layer

Infrastructure

Institutions

Acknowledgements

Analysis team

Francois Lefebvre
Pascale Marquis
Gary Lévesque
Emmanuel Gonzalez
Senthil Duraikannu
Alain Pacis
Robert Syme

Development team

Mathieu Bourgey
Edouard Henrion
Robert Eveleigh
Rola Dali
Hector Galvez Lopez
Paul Strenenowich
Pierre-Olivier Quirion
Pubudu Nawarathna
Ulysse Fortier-Gauthier

Data team

David Bujold
Romain Gregoire
David Lougheed
David Anderson
Ksenia Zaytseva
David Brownlee
Simon Chenard
Adrielle Houweling

Lab

Patricia Goerner-Potvin
Maxime Caron
Jeffrey Hyacinthe
Cristian Groza
Qinwei Zhuang
Mio Shibata
Xun Chen
Lindsay Dayton
Ksenia Egorova
Arber Kacollja
Rachade Hmamouchi



The group is recruiting!
guil.bourque@mcgill.ca



@guilbourque

