



**Global Alliance**  
for Genomics & Health  
Collaborate. Innovate. Accelerate.

## Data standards and global variant databases

Thomas Keane (@drtkeane)  
Head of European Genome-phenome Archive  
EMBL-EBI



# Genomic standards

## Fast Healthcare Interoperability Resource



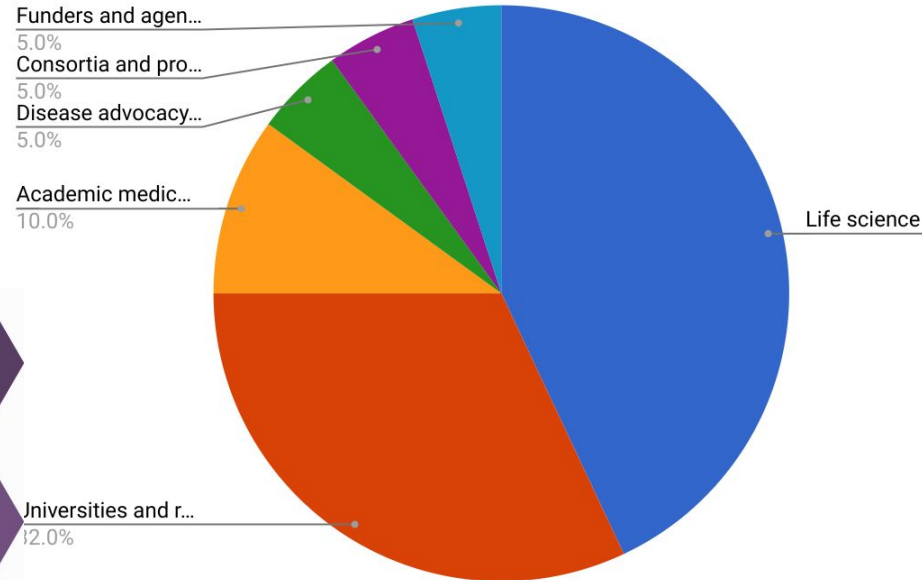
- Interoperability standard for exchange of healthcare information

## MPEG



- MPEG-G: standard for compression of genomic information

## Global Alliance for Genomics and Health (GA4GH)



# GA4GH mission

The Global Alliance for Genomics and Health aims to accelerate progress in genomic science and human health by developing standards and framing policy for responsible genomic and health-related data sharing.



**Ewan Birney**  
EMBL-European Bioinformatics Institute  
Hinxton, United Kingdom  
Chair, GA4GH  
Member, Steering Committee



**Peter Goodhand**  
Ontario Institute for Cancer Research  
Toronto, Canada  
Chief Executive Officer, GA4GH  
Member, Steering Committee



**David Haussler**  
University of California Santa Cruz  
Santa Cruz, United States  
Vice-Chair, GA4GH  
Member, Steering Committee

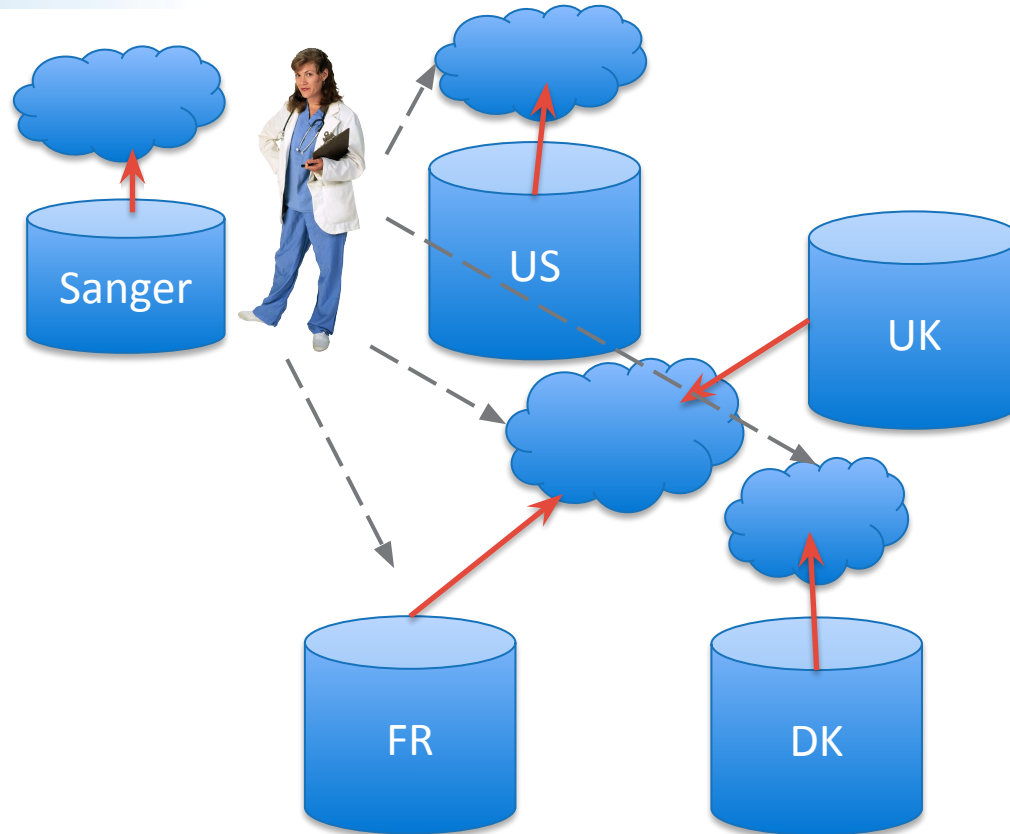


**Kathryn North**  
Murdoch Childrens Research Institute  
Melbourne, Australia  
Vice-Chair, GA4GH  
Member, Steering Committee  
Lead, Partner Engagement Initiative

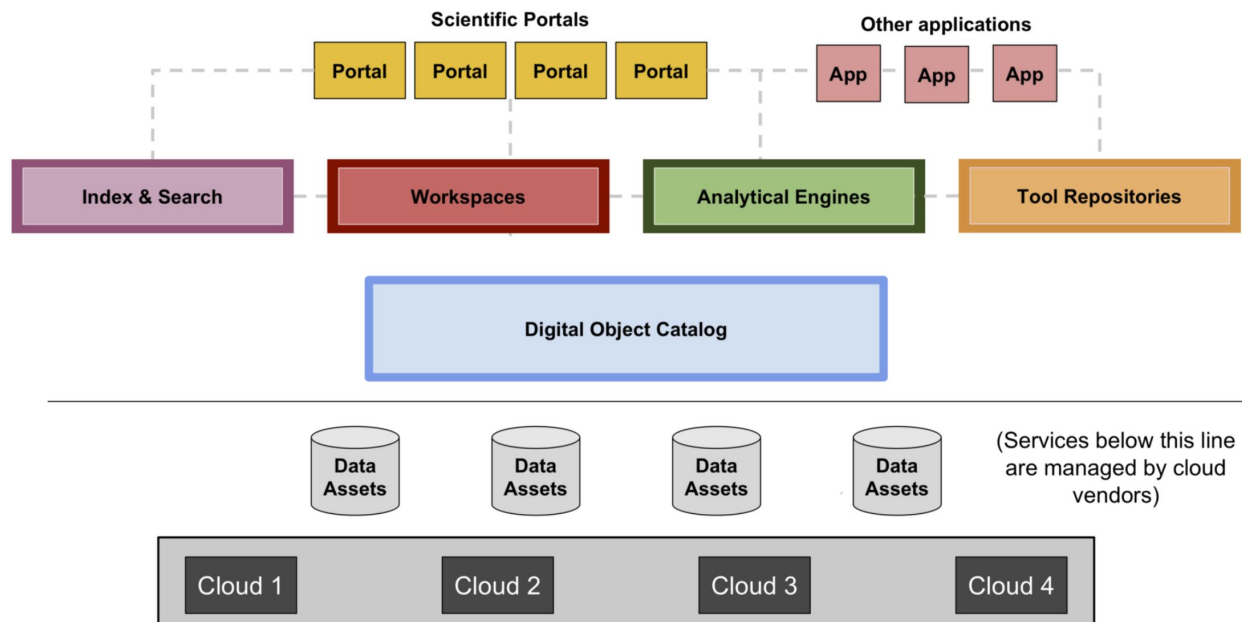
# Where do we want to get to?



Global Alliance  
for Genomics & Health



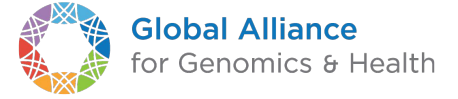
# A Data Biosphere for Biomedical Research



## Data Biosphere principles

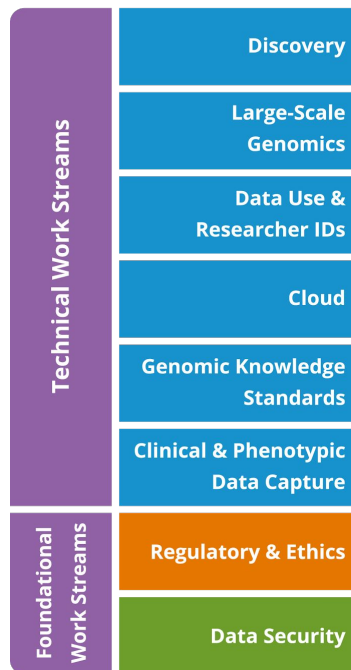
- (1) modular, functional components + interfaces
- (2) community-driven
- (3) open, open-source licenses
- (4) standards-based, to ensure interoperability

# Core Mantra for GA4GH



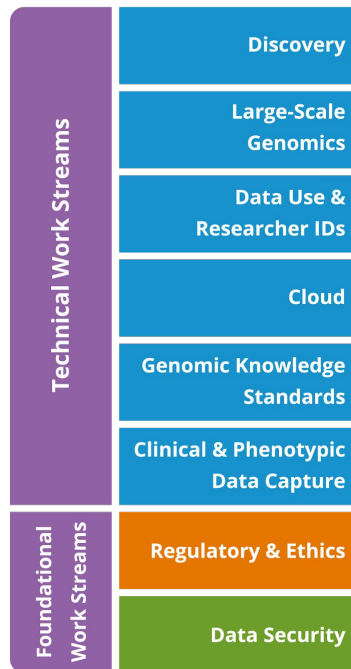
***Standardise on interfaces, compete on implementations***

# GA4GH Connect





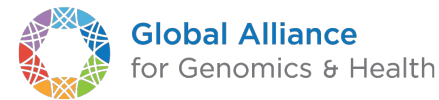
Real-World Driver Projects



		Real-World Driver Projects									
Technical Work Streams	Discovery	✓		✓		✓		✓			
	Large-Scale Genomics		✓		✓		✓		✓		
	Data Use & Researcher IDs	✓		✓		✓	✓				✓
	Cloud		✓	✓						✓	
	Genomic Knowledge Standards		✓				✓	✓	✓		
	Clinical & Phenotypic Data Capture	✓			✓	✓	✓				✓
Foundational Work Streams	Regulatory & Ethics										
	Data Security										

Partner Engagement

# 2017 Driver Projects



**All of Us Research Program**  
United States



**Australian Genomics**  
Australia



**BRCA Challenge**  
International



**CanDIG**  
Canada



**ClinGen**  
United States



**ELIXIR Beacon**  
Europe



**ENA / EVA / EGA**  
Europe



**Genomics England**  
United Kingdom



**Human Cell Atlas**  
International



**ICGC-ARGO**  
International



**Matchmaker Exchange**  
International



**Monarch Initiative**  
International



**National Cancer Institute  
Genomic Data Commons**  
United States



**TOPMed**  
United States



**Variant Interpretation  
for Cancer Consortium**  
International

# Data Use and Researcher ID Workstream

*Goal: Define the standards, both regulatory and technical, to facilitate both axes of access control (AAI & Data Use)*

## **Data Use Restrictions: What are you doing with the data?**

“The donor wants her data used only for non-commercial cancer research”



		Satisfies Data Use Restrictions	
		Yes	No
Appropriate Permissions	Yes	Access	No access
	No	No access	No access

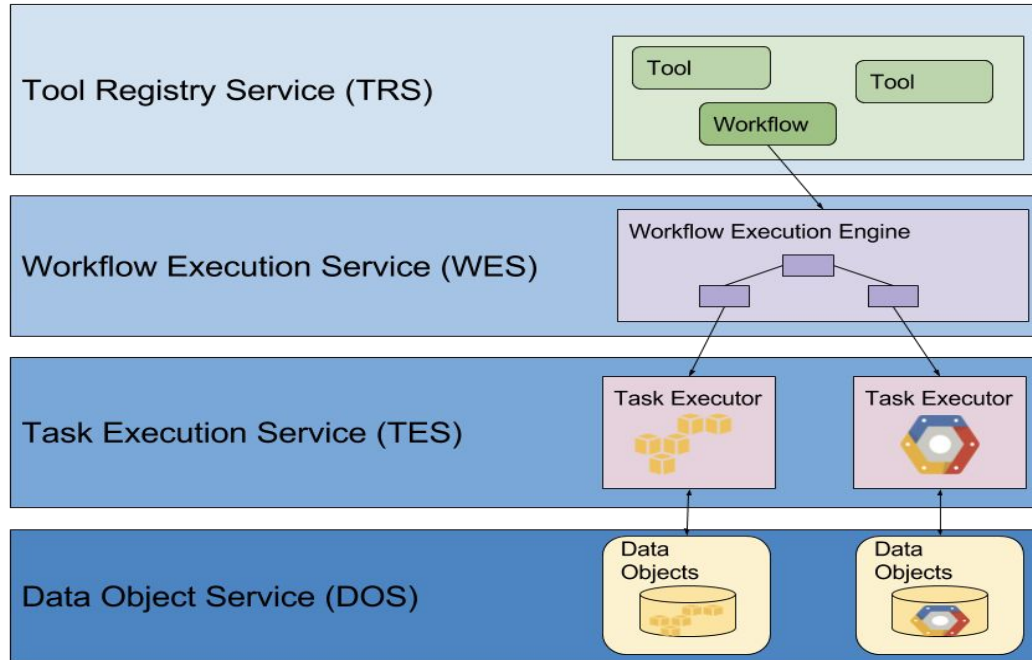


## **Authentication and Authorization Infrastructure (AAI): Who are you?**

“Only consortium members can access this data.”

Ravi Pandya (Microsoft)  
Anthony Philippakis (Broad)

# Cloud Workstream



**v1.0.0 -  
Dockstore.org**

**v0.1 - CWL  
workflow-service**

**v0.2 - Funnel**

**Pre-release - inspired  
by HCA, ICGC, GDC,  
etc.**

David Glazer (Verily)  
Brian O'Connor (UCSC)

# Discovery Workstream



Global Alliance  
for Genomics & Health



Beacon Network

Matchmaker Exchange  
Genomic discovery through the exchange of phenotypic & genotypic profiles

Our focus is on building standards for **federated, secured networks** of data and services, forming an “Internet of Genomics” (IoG), and **asking meaningful questions** across it.

Mark Fiume (DNASTack)  
Harindra Arachchi (Broad)

- Focus Areas and Driver Project feedback
  - Variant Representation: create a standard model for computer readable variant representation
  - Variant Annotation: represent and link annotations, including their evidence and provenance, to variants
- **June 2018:** Continue current VMC spec development and start to extend into structural variants, engage DPs as pilot projects.
- **Dec 2018:** Refine due to pilot feedback, expand structural variant capability, expand approach to allele equivalence
- **June 2019:** Review and expand VMC based on feedback, publish manuscripts with both structural & allele equivalence foci.

# Clinical and Phenotypic Data Capture Workstream



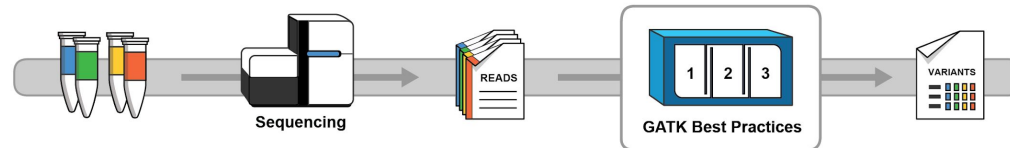
- Focus Areas: 3 sub-groups
  - Data representation & interoperability (Melanie Courtot, Tudor Groza, Allison Heath)
  - Data exchange & interoperability (Grant Wood, Michael Baudis, Gil Alterovitz)
  - Implementation/Education/Engagement (Andrea Ramirez, John Mattison)
- How were the following deliverables prioritized based on engagement with Driver Projects and other Work Streams
  - Leads and sub-group leads directly linked to Driver Projects
    - Melissa Haendel & Tudor Groza → *Monarch Initiative*
    - David Hansen → *Australian Genomics*
    - Michael Baudis → *ELIXIR Beacon*
    - Andrea Ramirez → *All of Us*

David Hansen (Australian Genomics)  
Melissa Haendel (Monarch Initiative)

# Large Scale Genomics Workstream

*Goal: Create standardized methods for accessing large-scale genomic data by file-based, API-based, cloud-based, and distributed access.*

- Key stakeholders
  - Platform vendors, clinical genomics platforms, genome sequencing providers, bioinformatics developers, academic sequencing centers



Thomas Keane, EMBL-EBI  
Oliver Hofmann, University of Melbourne

# NGS data formats (2009-date)



## Legacy of the 1000 Genomes Project

- Standardisation at the file level
- Few population scale projects

## Sequencing reads

- SAM/BAM (2009-)
- CRAM (2016-)

## Genetic variation

- VCF/BCF (2010-)

## Enabled many applications of NGS

From Richard Durbin <rd@SANGER.AC.UK>★  
Subject **Re: File format for SNP calls and individual genotypes** 15/07/09 19:32  
To 1000ANALYSIS@LIST.NIH.GOV★

Hello Gabor,  
I started looking at using some version of the new proposed format last night for all my current activities on comparing call sets, and immediately hit problems. There is no doubt that having a common format would be good. But it needs to be flexible enough to contain additional information that we may be carrying around for filtering etc.

After thinking about this, and talking to Heng this morning, I propose the following. First an example:

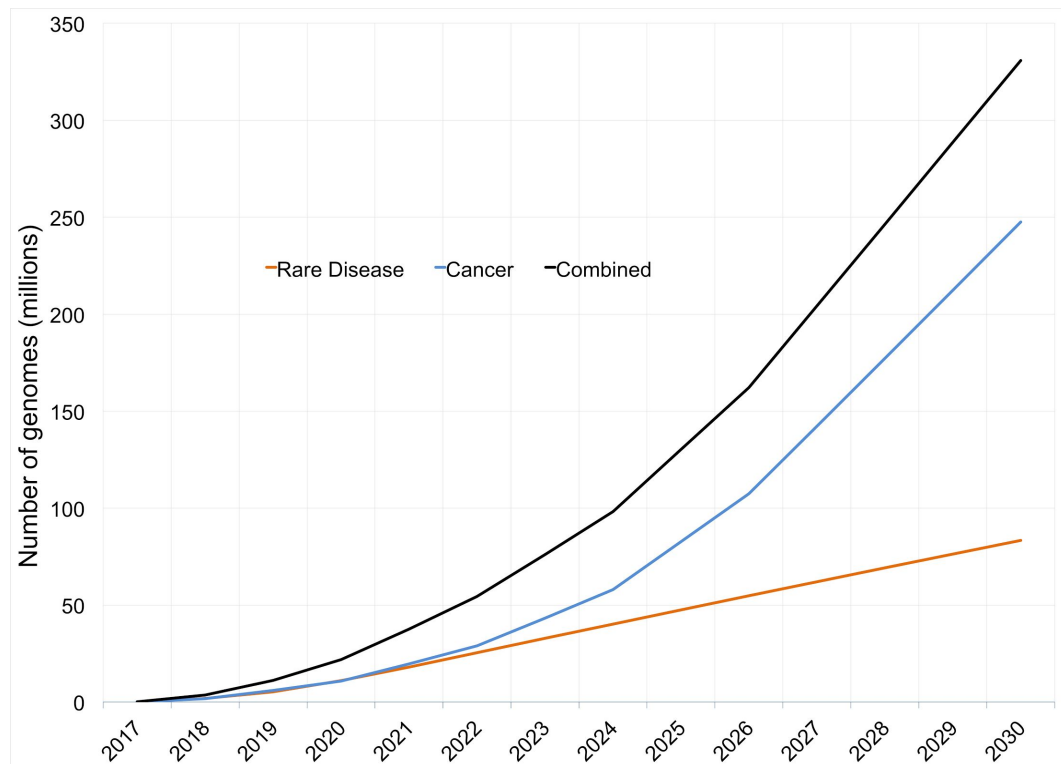
```
1 193 192 0 MQ:99 C T C/T:1f 77 C/T:115:19  
1 25 0 MQ:72 A T T/T:- T/T:7:5  
- 77 0 G D2 G/D '02:17:11  
- 0 . G .
```

**The variant call format and VCFtools**  
Sequence analysis  
The variant call format and VCFtools  
Sequence analysis  
Sequence Alignment/Map format and SAMtools  
for the CEU track  
Heng Li<sup>1,2</sup>, Bob Handsaker<sup>2,3</sup>, Alec Wysoker<sup>2,3</sup>, Tim Fennell<sup>2,3</sup>, Jue Ruan<sup>3</sup>, Nils Homer<sup>4</sup>, Gabor Marth<sup>5</sup>, Goncalo Abecasis<sup>6</sup>, Richard Durbin<sup>1,\*</sup>, and 1000 Genome Project Data Processing Subgroup<sup>1</sup>  
Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge CB20 1SA, UK; <sup>2</sup>Broad Institute of MIT and Harvard, Department of Computer Science, University of California Los Angeles, Los Angeles, CA 90045; <sup>3</sup>Department of Biostatistics, Boston College, Chestnut Hill, MA 02467; <sup>4</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA; and <sup>5</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA; and <sup>6</sup>Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA; and <sup>7</sup>Department of Biotechnology Information, MD 20894, USA and <sup>8</sup>Department of Statistics, University of Oxford, Oxford OX1 3TG, UK  
\*Correspondence: richard.durbin@sanger.ac.uk  
doi:10.1093/bioinformatics/btp330  
Advance Access publication June 7, 2011  
Vel. 27, no. 15, 2011, pages 2715-2718  
doi:10.1093/bioinformatics/btp330

# How Many Genomes?



Global Alliance  
for Genomics & Health



Genomics in healthcare: GA4GH looks to 2022

Ewan Birney, Jessica Vamathevan, Peter Goodhand  
doi: <https://doi.org/10.1101/203554>

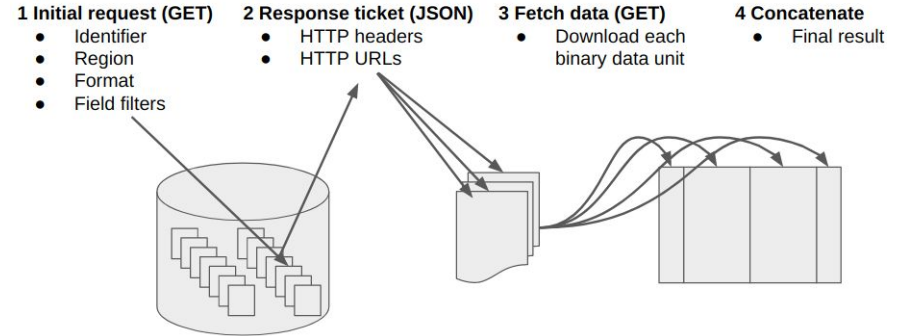
# Htsget Streaming API

Securely accessing or streaming genomic data

- htsget: A standardised non-file based API for securely streaming genomic data
- Spec maintainer: Mike Lin (DNAnexus)

## Milestones

- Launched 1.0 (Oct 17)
- VCF support (Feb 18)
- POST support (Mar 18)
- htsget paper (Jun 18)



	Servers									
	WSI		DAWS		DAZ		GCP		EGA	
Clients	BAM	CRAM	BAM	CRAM	BAM	CRAM	BAM	CRAM	BAM	CRAM
htsget	Green	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow
WSI	Green	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow
EGA	Green	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow
Samtools	Green	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow

**Table 1:** Interoperability status for clients (side) and servers (top). Green boxes indicate fully passing interoperability (see 2 above for details), yellow indicates combination not supported by the server. WSI=Wellcome Sanger Institute, DAWS=DNAnexus on Amazon AWS, DAZ=DNAnexus on Azure, GCP=Google Cloud Platform, EGA=European Genome-phenome Archive. Server URLs are given in section 1 above.

# Htsget Streaming API

Who has implemented it so far?

DNAnexus



Samtools  
BCFtools  
HTSlib

**verily**



Genome Maps

\*coming soon

# Genomics enters healthcare

In 2017 active genomic medicine programmes are already underway in many countries. Finland, the UK, the US, and Australia are a few examples.



# Thank you

<https://www.ga4gh.org/>



**Global Alliance**  
for Genomics & Health  
Collaborate. Innovate. Accelerate.

## 1 MILLION PEOPLE



## 18,000 people by 2021

To improve patient outcomes and support research, the Australian healthcare system is building a Federation of clinical and genomic data.

10% of Finland's population expected to have some genomic data in healthcare by 2020.

## 10%

## 100,000 PATIENTS

The UK National Health System plans to sequence 100,000 individuals by 2020.