

# Genomic Workflows

Jason Walker

BFX Workshop

April 22<sup>nd</sup>, 2024





Tim Malone, [www.timmalone.id.au](http://www.timmalone.id.au) [CC BY-SA 2.5 (<https://creativecommons.org/licenses/by-sa/2.5/>)]



Sarbjit Bahga [CC BY-SA 4.0 (<https://creativecommons.org/licenses/by-sa/4.0/>)]



Image by [ElasticComputeFarm](#) from [Pixabay](#)



[Ingredients for cucumber kimchi](#) by [David Davies](#) licensed under CC BY-SA 2.0



[Bengali cooking tools](#) by Amartyabag licensed under [CC BY-SA 3.0](#)



30 Years of Our Best Recipes

ANNUAL RECIPES

ANNUAL RECIPES

ANNUAL RECIPES

ANNUAL RECIPES



[Juicer with fruit](#) by [mathiasbaert](#) licensed under [CC BY 2.0](#)





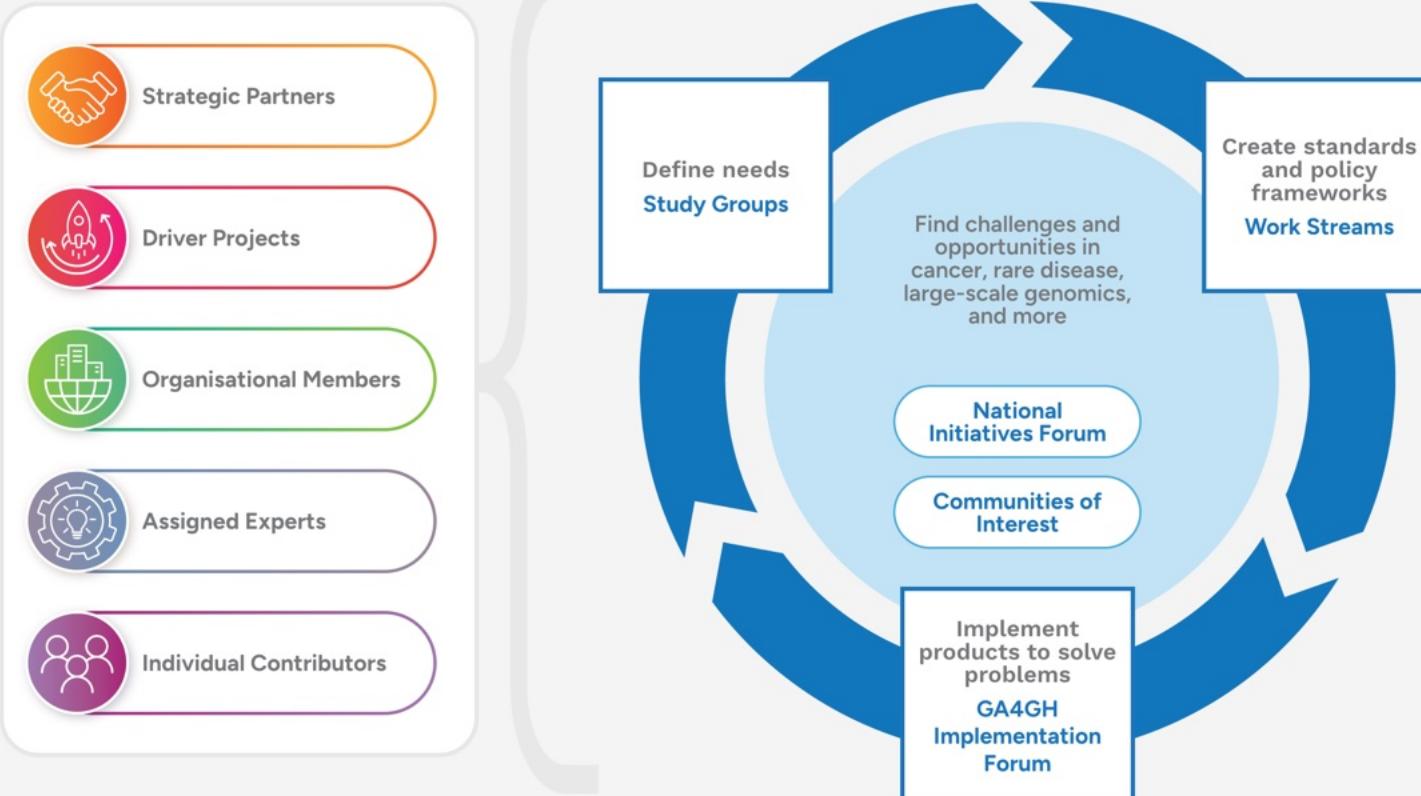
By Paul Keller (chefs at work Uploaded by Yarl) [CC BY 2.0 (<http://creativecommons.org/licenses/by/2.0>)], via Wikimedia Commons



## How we work

Here's our community...

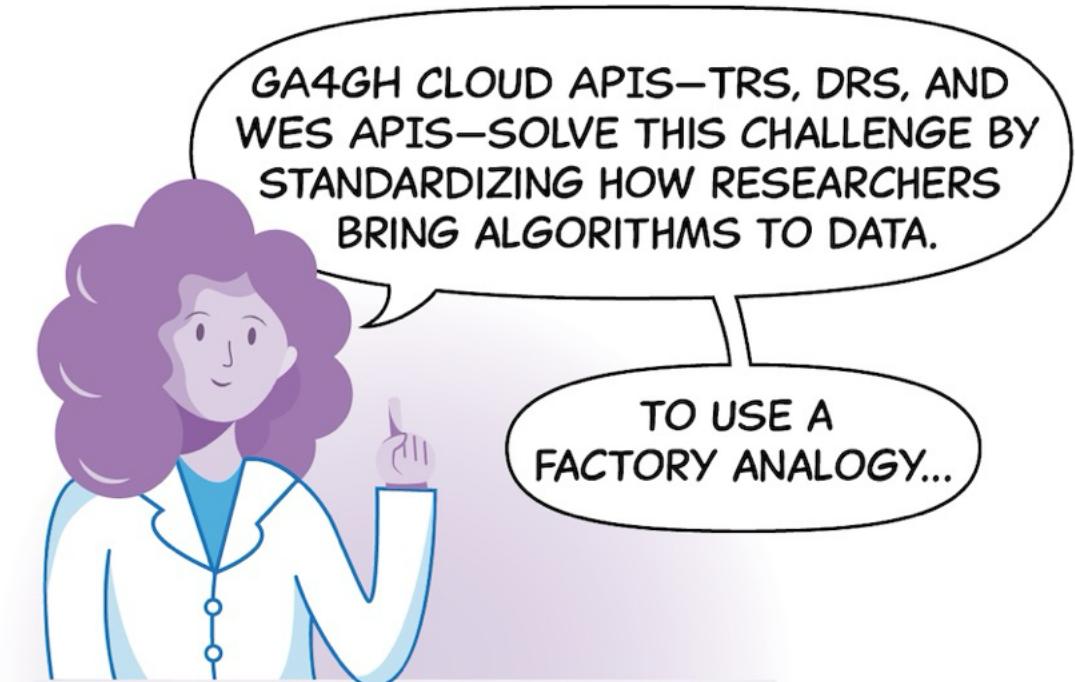
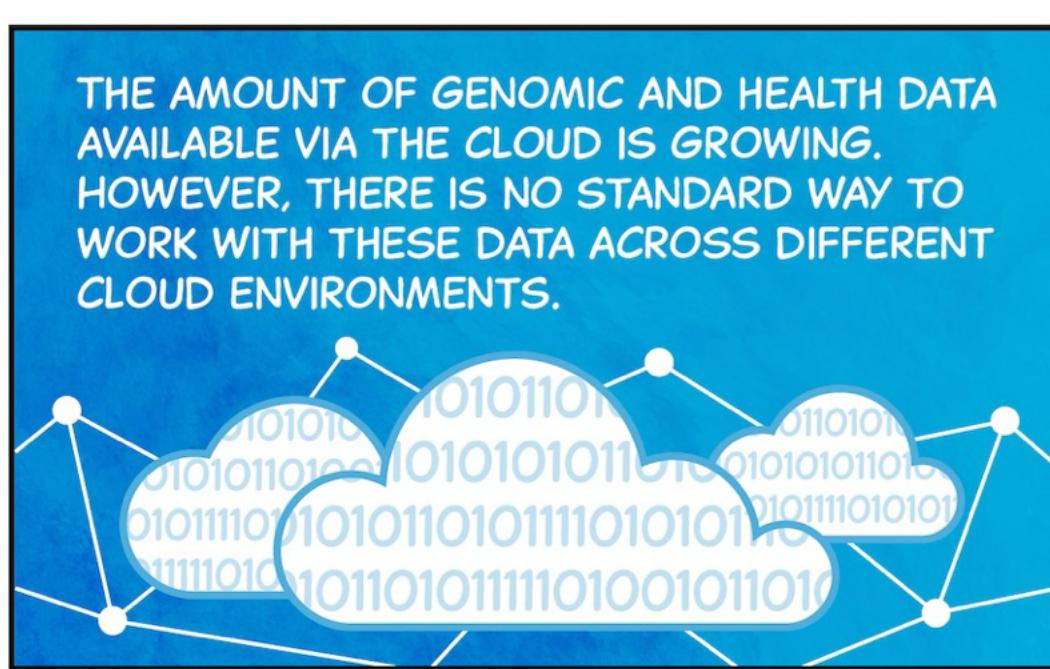
...and here's what we do.



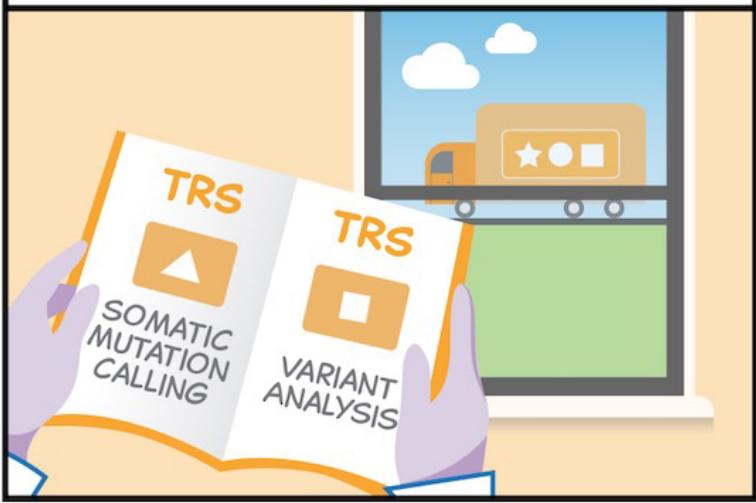
DISCOVER WHAT WE DO



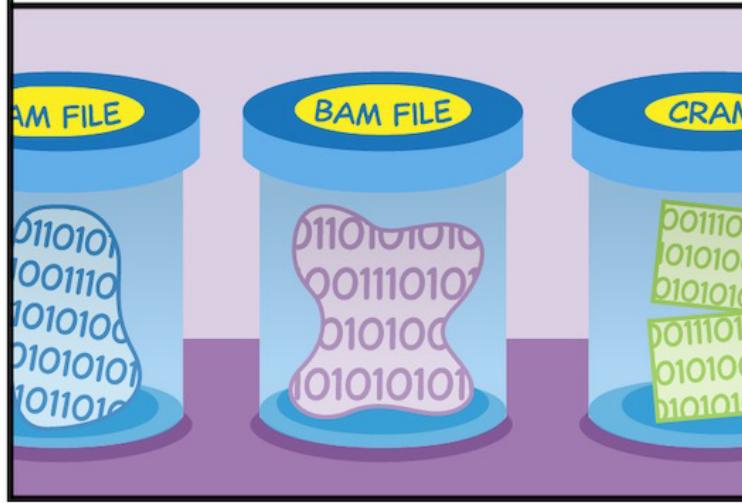
# GA4GH CLOUD APIs: BRINGING ANALYSES TO DATA



WITH TRS (TOOL REGISTRY SERVICE) API, RESEARCHERS CAN SHARE AND DISCOVER TOOLS AND WORKFLOWS FOR WORKING WITH DATA.



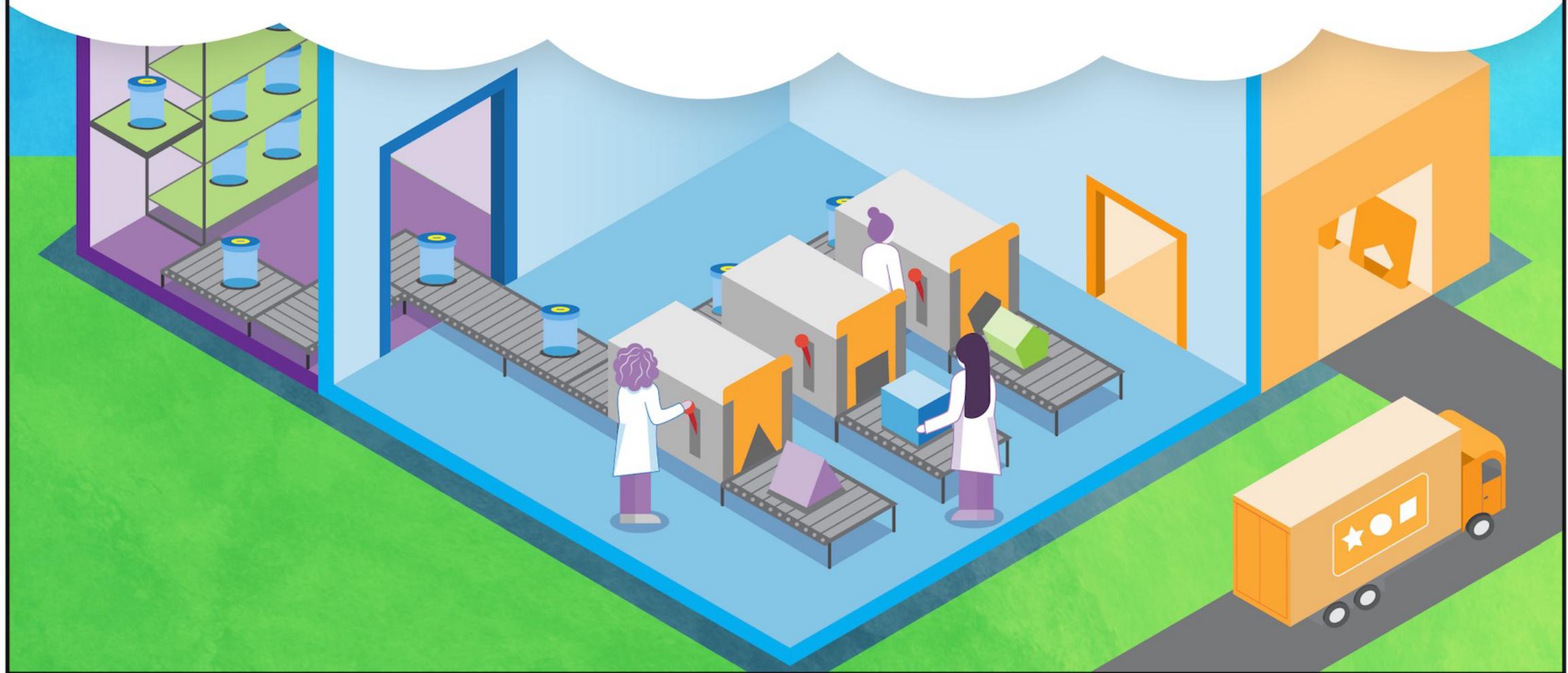
DRS (DATA REPOSITORY SERVICE) API PACKAGES DATA IN A UNIVERSAL WAY, DESPITE HOW THE DATA WAS STORED OR MANAGED ORIGINALLY.



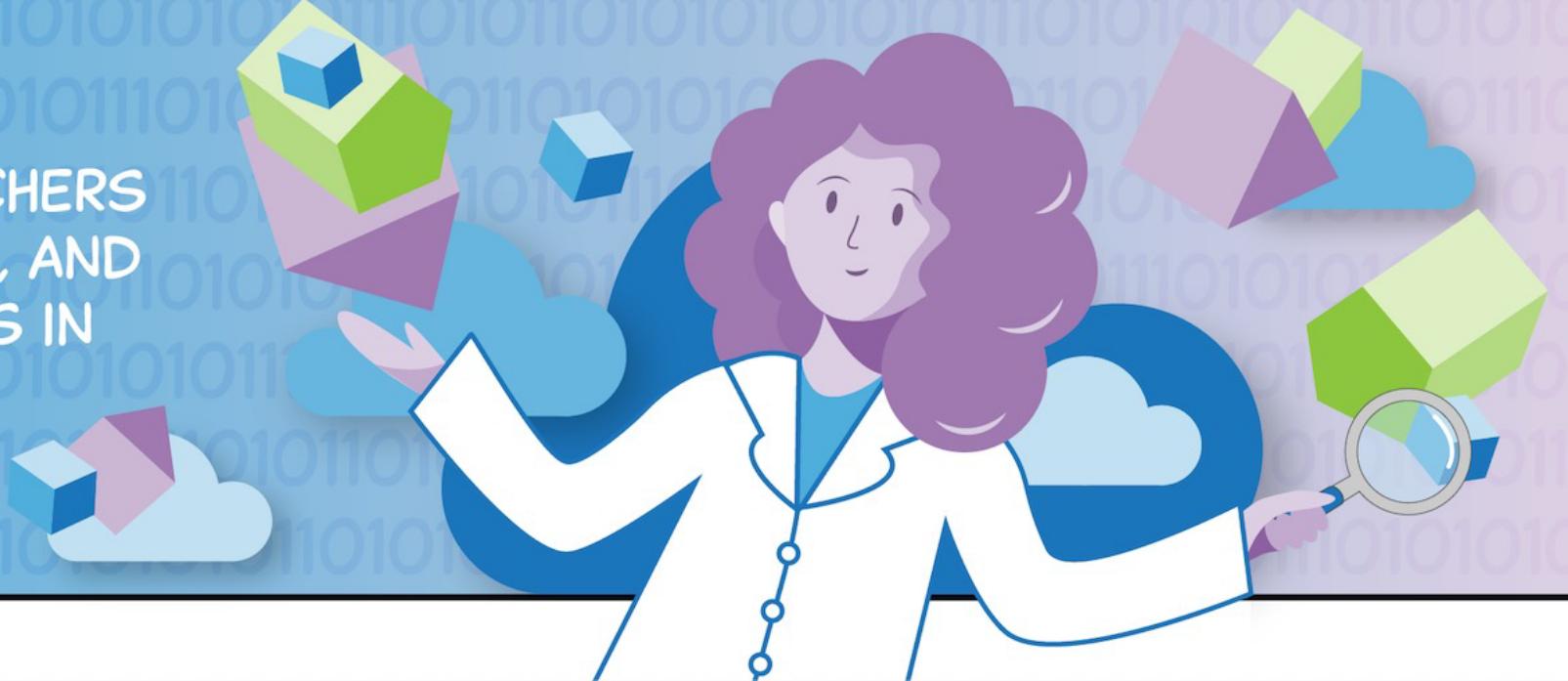
WES (WORKFLOW EXECUTION SERVICE) API ALLOWS RESEARCHERS TO EXECUTE ANALYSES ACROSS DIFFERENT ENVIRONMENTS AND ACHIEVE CONSISTENT RESULTS.



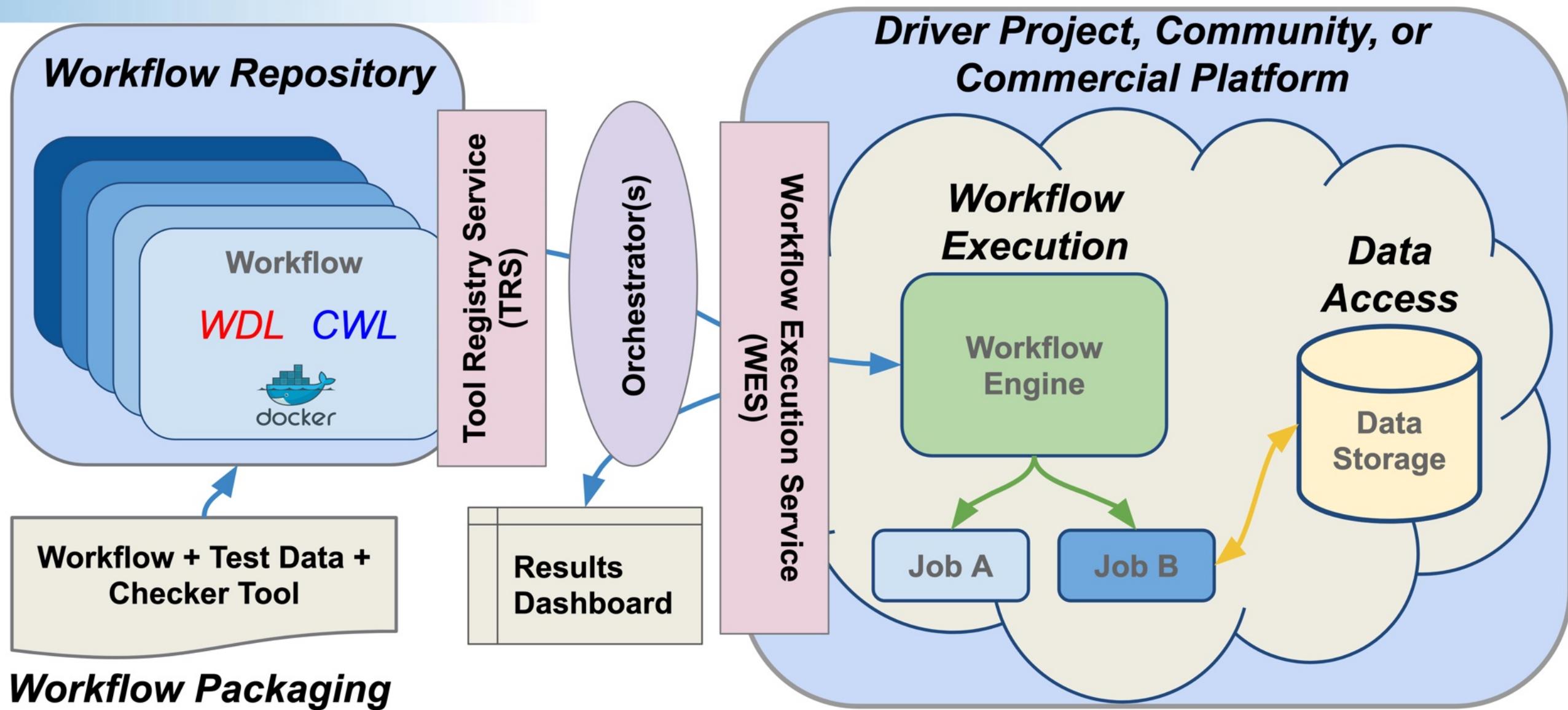
TOGETHER, THE GA4GH CLOUD APIS ENABLE A SEAMLESS ECOSYSTEM...



...ALLOWING RESEARCHERS  
TO ACCESS, ANALYZE, AND  
WORK WITH DATASETS IN  
THE CLOUD.



# Leveraging Our Existing Cloud Work Stream Standards





COMMON  
WORKFLOW  
LANGUAGE

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



<https://openwdl.org/>



<https://www.nextflow.io/>



<https://snakemake.github.io/>



<https://galaxyproject.github.io/>

# *Common Workflow Definition Languages*

- Defining a process, including the software tools, settings, and configuration, as a workflow results in portable and sharable best practices.
- CWL, WDL, Nextflow, Galaxy, and Snakemake are the most *common* ways to represent workflow definitions in Genomics.
- Standard ways to represent inputs, outputs and tools used in workflows.
- Combined with Docker these definitions are portable to many High Performance Computing environments including public Clouds.
- The definitions themselves adhere to FAIR Principles for Research Software (FAIR4RS), making the workflows digital data objects as well as the results.

# FAIR4RS - Findable

**F: Software, and its associated metadata, is easy for both humans and machines to find.**

F1. Software is assigned a globally unique and persistent identifier.

F1.1. Components of the software representing levels of granularity are assigned distinct identifiers.

F1.2. Different versions of the software are assigned distinct identifiers.

F2. Software is described with rich metadata.

F3. Metadata clearly and explicitly include the identifier of the software they describe.

F4. Metadata are FAIR, searchable and indexable.

# FAIR4RS - Accessible

## **A: Software, and its metadata, is retrievable via standardised protocols.**

A1. Software is retrievable by its identifier using a standardised communications protocol.

A1.1. The protocol is open, free, and universally implementable.

A1.2. The protocol allows for an authentication and authorization procedure, where necessary.

A2. Metadata are accessible, even when the software is no longer available.

# FAIR4RS - Interoperable

**I: Software interoperates with other software by exchanging data and/or metadata, and/or through interaction via application programming interfaces (APIs), described through standards.**

- I1. Software reads, writes and exchanges data in a way that meets domain-relevant community standards.
  
- I2. Software includes qualified references to other objects.

# FAIR4RS - Reuseable

**R: Software is both usable (can be executed) and reusable (can be understood, modified, built upon, or incorporated into other software).**

R1. Software is described with a plurality of accurate and relevant attributes.

R1.1. Software is given a clear and accessible license.

R1.2. Software is associated with detailed provenance.

R2. Software includes qualified references to other software.

R3. Software meets domain-relevant community standards.

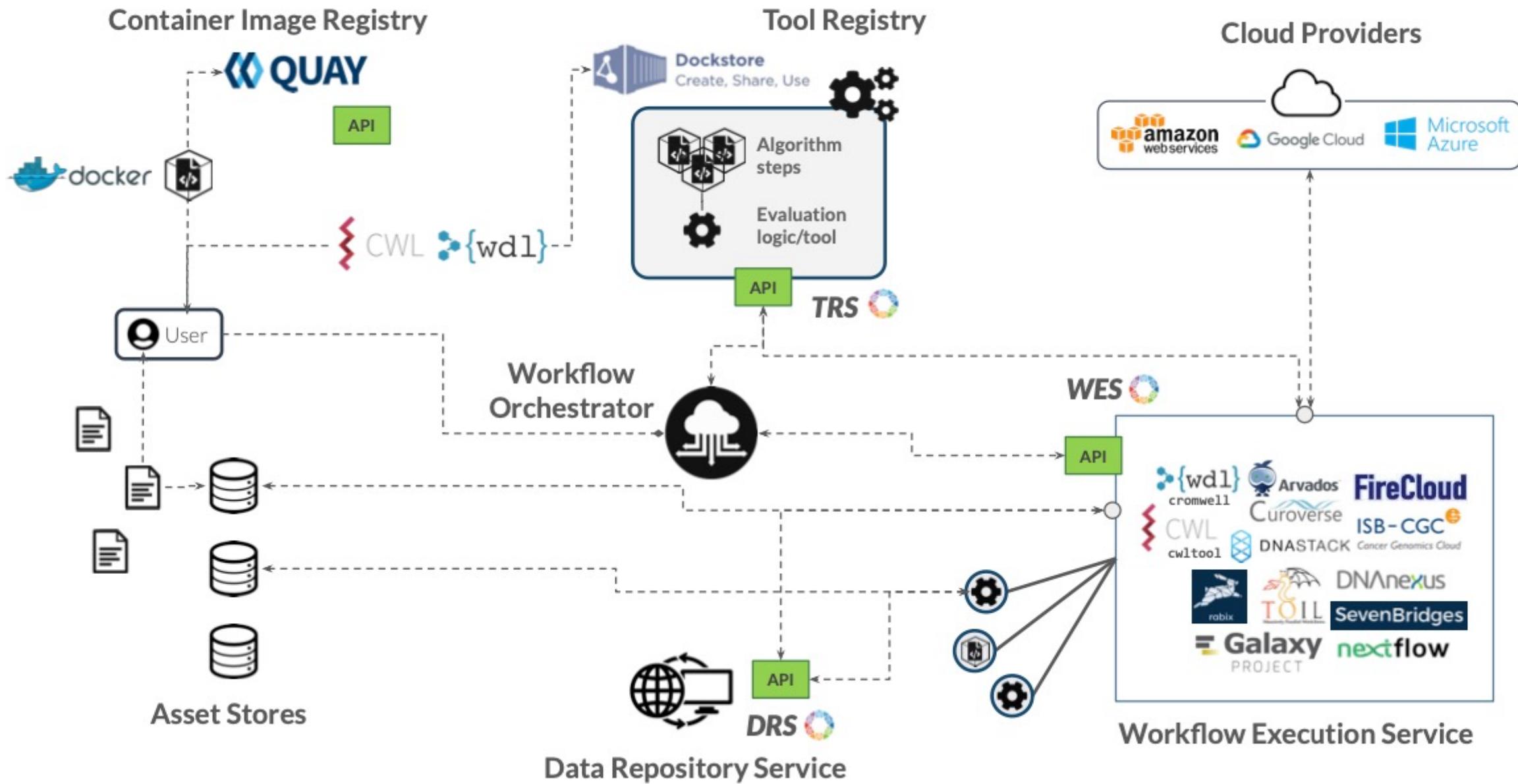
# Language-level Comparison

Aspect	Nextflow	CWL	WDL
Parent language	Ruby and Groovy	N/A	N/A
Compilation	Interpreted	Compiled	Compiled
GUILs	NextflowWorkbench <sup>27</sup> , DolphinNext <sup>28</sup>	Rabix composer	Pipeline Builder <sup>29</sup>
DSL features	complete, extensible in Groovy and Java	Limited standard library, extensible via javascript	Limited standard library
Variables	Qualified, unique within scope	Typed, unique identifiers	Typed, fully qualified names
Loops	Parallel queue channels	Parallel scatter via ScatterFeatureRequirement	Parallel scatter
Conditionals	Via when declaration within a process	When and pickValue fields proposed in CWLv1.2	If blocks producing optional output types
Enforcing good practices	nf-core ( <a href="https://nf-co.re/">https://nf-co.re/</a> )	CWL guide ( <a href="https://www.commonwl.org/user_guide/rec-practices/">https://www.commonwl.org/user_guide/rec-practices/</a> )	–

# Executor-level Comparison

WfMS		Remarks
Language	Execution engine	
<i>Nextflow</i> (DSL-1, DSL-2)		Complete WfMS, supporting conditionals, loops and nested logic
CWL	<i>cwltool</i> <sup>†</sup>	The official reference implementation of an execution engine for the complete CWL standard <sup>33</sup> ; no cluster or cloud support
	<i>arvados</i> <sup>†</sup> (1.0, 1.1, 1.2)	Most feature-rich CWL runner, albeit with tedious setup
	<i>toil-cwl-runner</i> <sup>†</sup> (1.0.1)	Optimized for cloud environments, less stable in batch environments (Section Scalability)
	<i>cwl airflow</i> <sup>†</sup> (1.1)	Works with <code>celery</code> and Kubernetes clusters, not readily with HPC CRMs
	<i>REANA</i> <sup>†</sup> (Documentation missing)	Cloud-optimized platform. For HPC, only CERN Slurm and HTcondor are supported
	<i>Cromwell</i> (1.0)	Supports CWL workflows via WOM, with comparable performance in both languages (Section Scalability)
	<i>cwl-tes</i> (1.0)	Partial implementation at present, with tedious setup. GA4GH TES API compatible
	<i>rabix executor</i> (sbg:draft-2, 1.0)	Single node local executor is no longer supported by the original developer team at Seven Bridges
	<i>Cromwell</i> (draft-2, 1.0)	De facto standard for executing WDL workflows. Support for nested loops is version-dependent
WDL	<i>toil-wdl-runner</i> (draft-2)	No support for modularity or nesting of loops and conditionals. Support for batch systems is also rudimentary
	<i>miniWDL</i> (draft-2, 1.0)	No cluster or cloud support. Includes Cromwell wrapper

# WES in the ecosystem of GA4GH APIs



# A Workflow Example:

[https://github.com/genome/bfx-workshop/tree/master/lectures/week\\_23/SingleSampleQc](https://github.com/genome/bfx-workshop/tree/master/lectures/week_23/SingleSampleQc)