# Pipelining solutions tend to proliferate



Randall Munroe, XKCD
https://www.xkcd.com/927/

# Use scripted workflows to automate processing

- Automate repetitive tasks

- Increase auditability and reproducibility

- Reduce human error

- Reduce time spent re-implementing the wheel

# Use a workflow language that humans can understand

```
workflow myWorkflowName {

    File my_ref
    File my_input
    String name

    call task_A {
        input: ref= my_ref, in= my_input, id= name
    }
    call task_B {
        input: ref= my_ref, in= task_A.out
    }

}
task task_A {      ...      }
task task_B {      ...      }
```

```
task task_A {

    File ref
    File in
    String id

    command {
        do_stuff -R ${ref} -I ${in} -O ${id}.ext
    }
    runtime {
        docker: "my_project/do_stuff:1.2.0"
    }
    output {
        File out= "${id}.ext"
    }

}
```
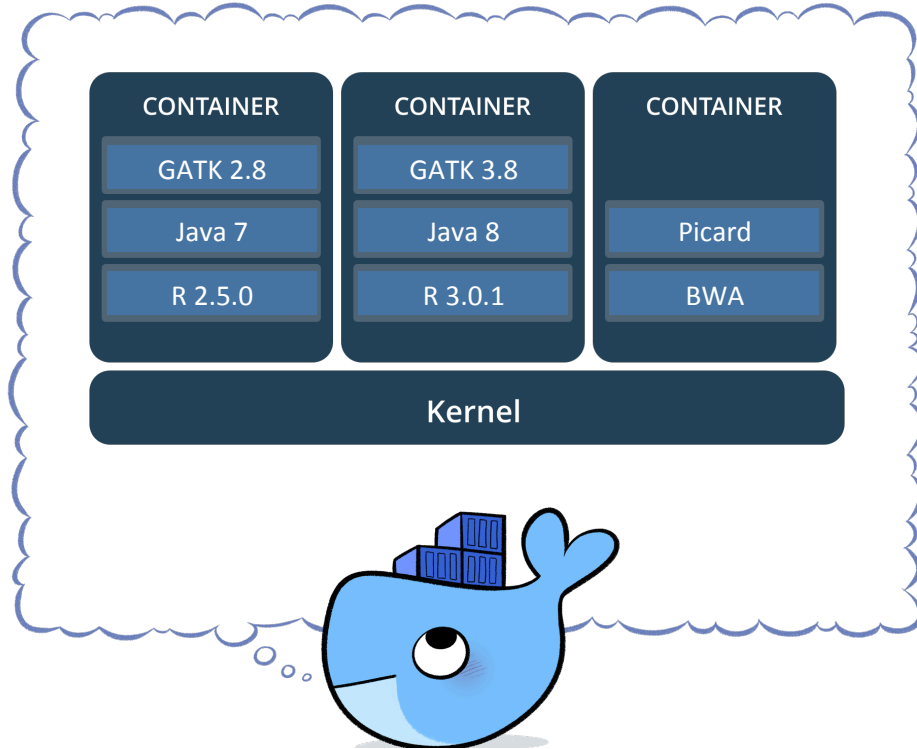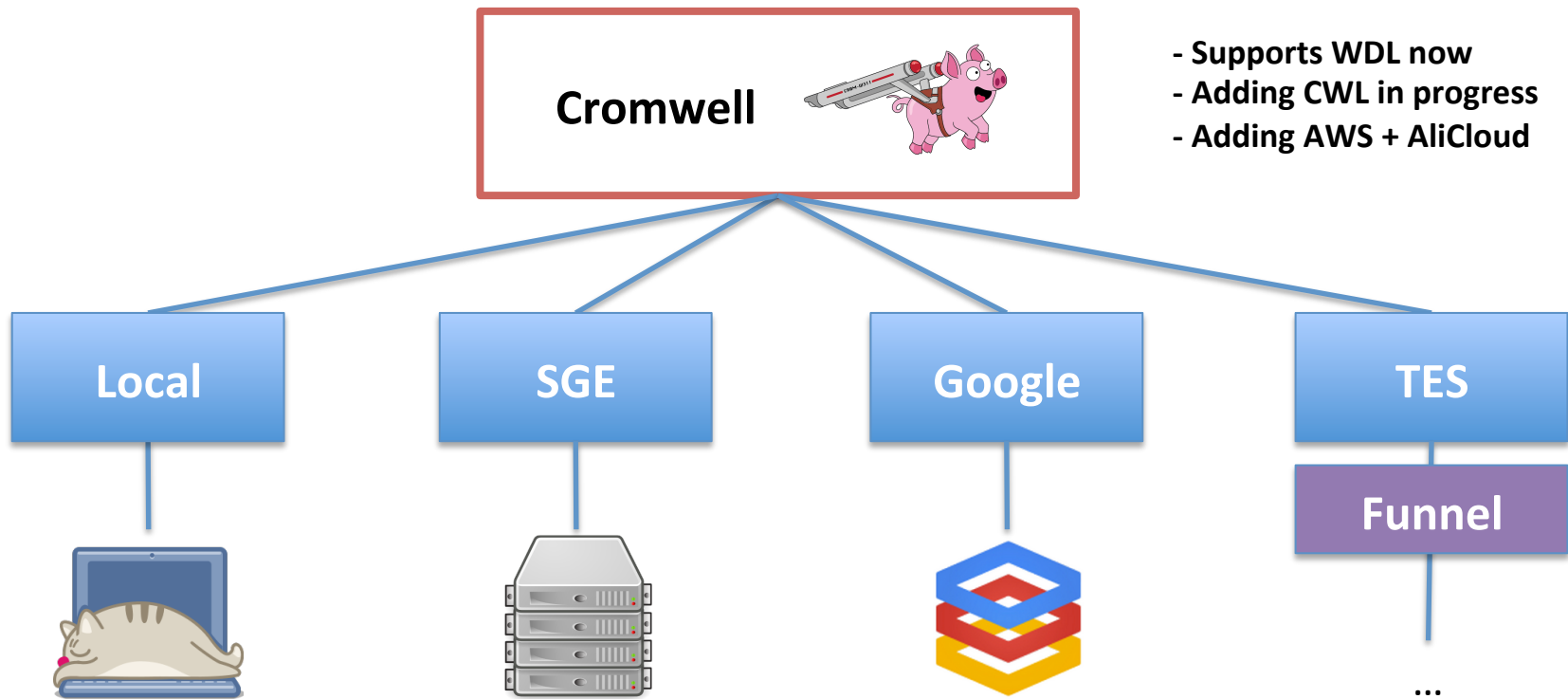
# Use containers for portability & reproducibility



A container encapsulates
**all the software dependencies**
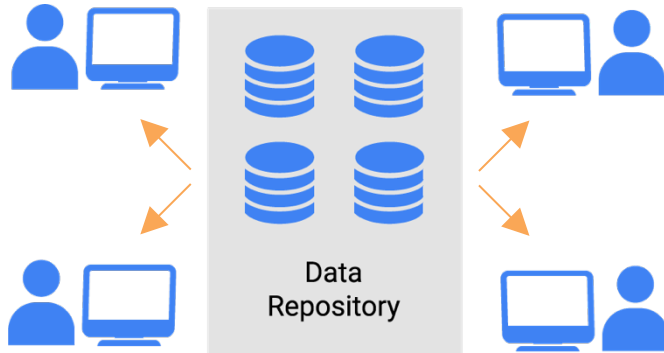associated with running a program

Takes the guesswork out of running
pipelines on different platforms!

*Modified from https://www.docker.com/what-container*

# Use a workflow execution engine that runs anywhere

**Cromwell**

- Supports WDL now
- Adding CWL in progress
- Adding AWS + AliCloud

**Local**

**SGE**

**Google**

**TES**

**Funnel**

...

*https://github.com/broadinstitute/cromwell*

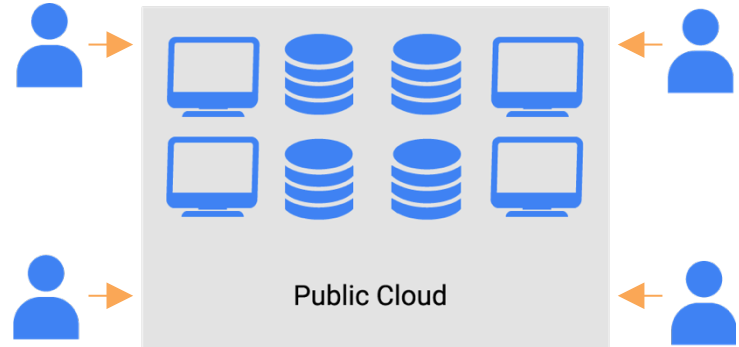# Use the Cloud to share methods & data



**Traditional Way:** Bring data to the researchers

Data Repository

**Problems**
Data sharing = data copying
Requires big infrastructure at each site
Largely fixed compute
Individual security implementations

**Cloud Way:** Bring researchers to the data

Public Cloud

**Solutions**
True data sharing
Cloud provides the infrastructure
Elastic compute and storage
Centralized security implementation

# Clouds are elastic!



**Genome processing requests per day over the last several months in the cloud**

Genome processing requested by the Genomics Platform has been "spiky"

- We haven't needed to pay for compute power when we weren't using it
- We can easily tolerate the spikes without being forced to maintain a backlog of "things to process once everything calms down"

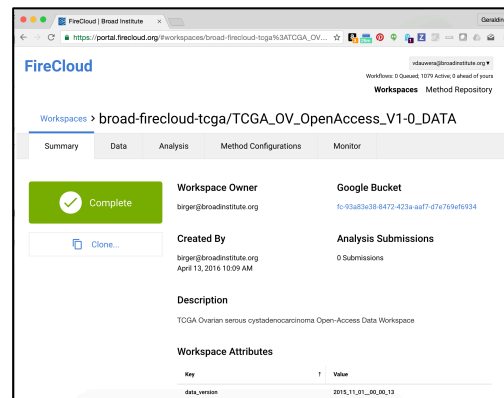# Several paths to running GATK pipelines on Cloud

Operators submit workflows written in WDL to Cromwell, which parses the script, generates individual jobs and dispatches them for execution via the specified backend.

**LOCAL**

Cluster

localize / delocalize

NFS — Storage

SGE

Command-line tools like GATK

Script + list of inputs

{wdl}

**Cromwell Execution Service**

Other cloud backends in development

**GOOGLE CLOUD**

Storage

gs://

localize

delocalize

Compute Engine

pull container image

docker

Genomics Pipelines API

# Self-service analysis for all: Workbench + FireCloud

- Collaborative **cloud-based** analysis platform built on top of **Google Cloud Platform**
- **Free to access** / compute & storage charged by Google
- Access published **data** and/or add your own
- Access existing **methods** and/or add your own
- **Execute** analyses in auditable pipelines
- **Share** data, methods and results with collaborators

**FireCloud Portal**

**Workbench**

**API**

| Workspaces |
| Data Library |

**Workflows Repository**

**Execution (Cromwell)**

**IAM & Groups**

FireCloud puts GATK4 workflows in everyone's hands

# GATK4 workflows preloaded in FireCloud workspaces

# Coming next: Interactive analysis with Notebooks

**SEE YOU ON THE CLOUD!**