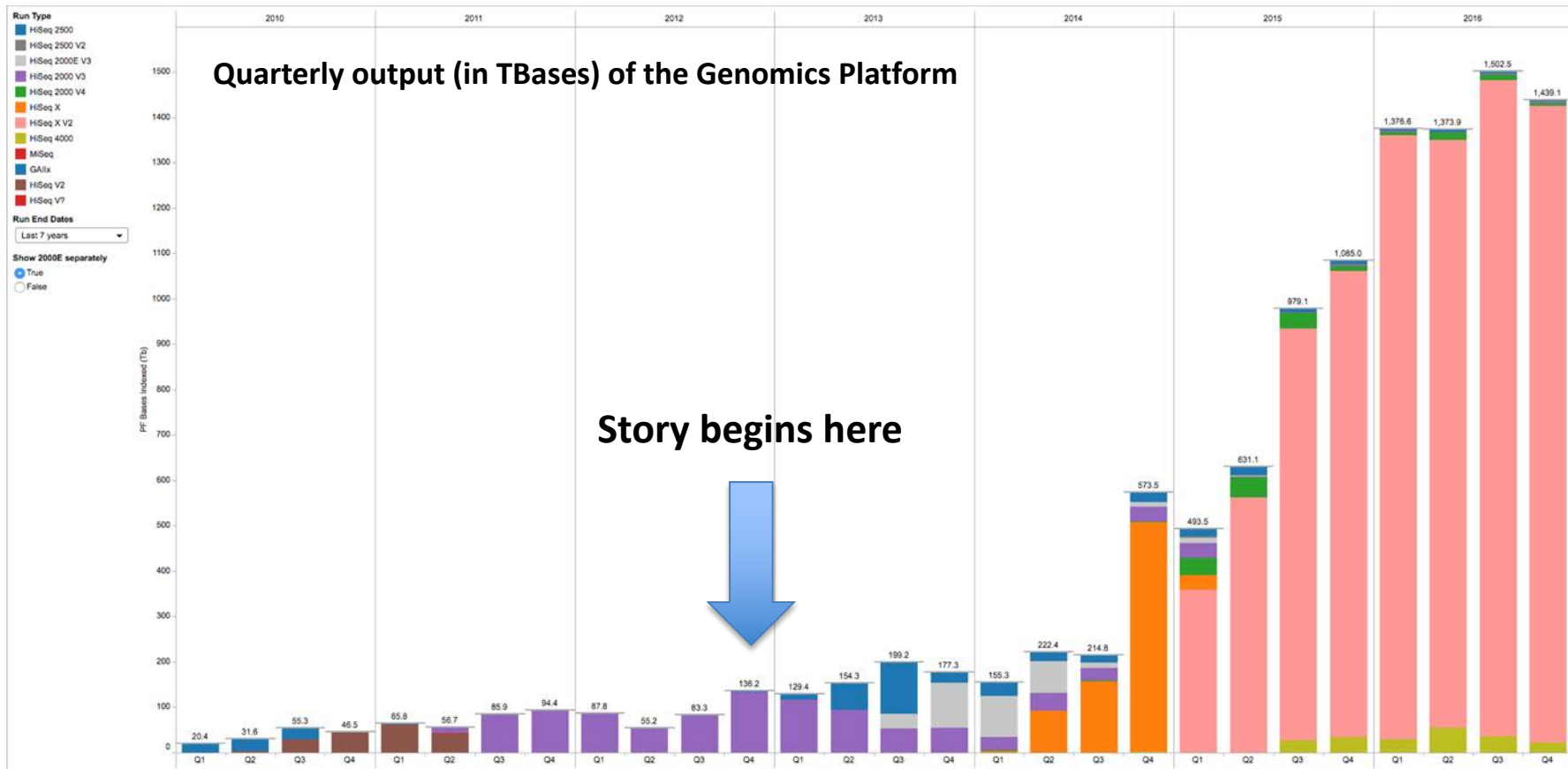# Getting started with WDL & Cromwell

## Bioinformatics workflows at any scale

Ruchi Munshi

Data Sciences Platform

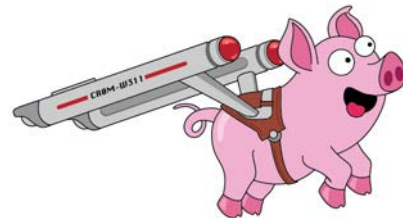Broad Institute

# The backdrop: data generation set to explode



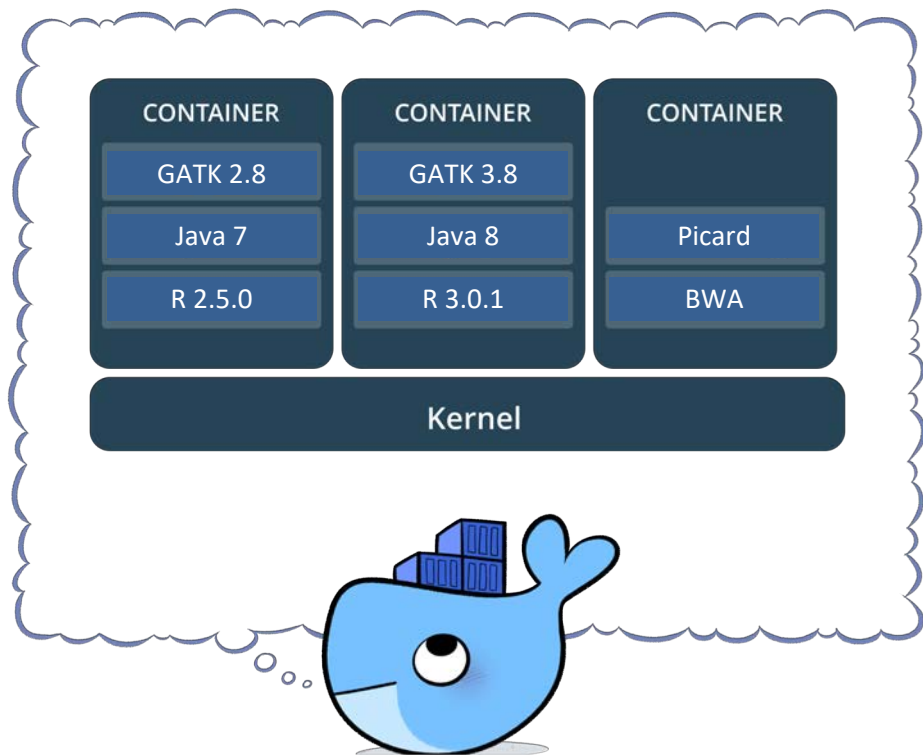Quarterly output (in TBases) of the Genomics Platform

Story begins here

# Meet Cromwell & WDL

- Execution engine that can
  - Run on any platform (HPC *and* on Cloud)
  - Seamlessly scale based on workflow needs
  - Provide maximal flexibility for all use cases
  - https://github.com/broadinstitute/cromwell


- Workflow language that humans can read/write
  - Methods developers and biomedical scientists at large
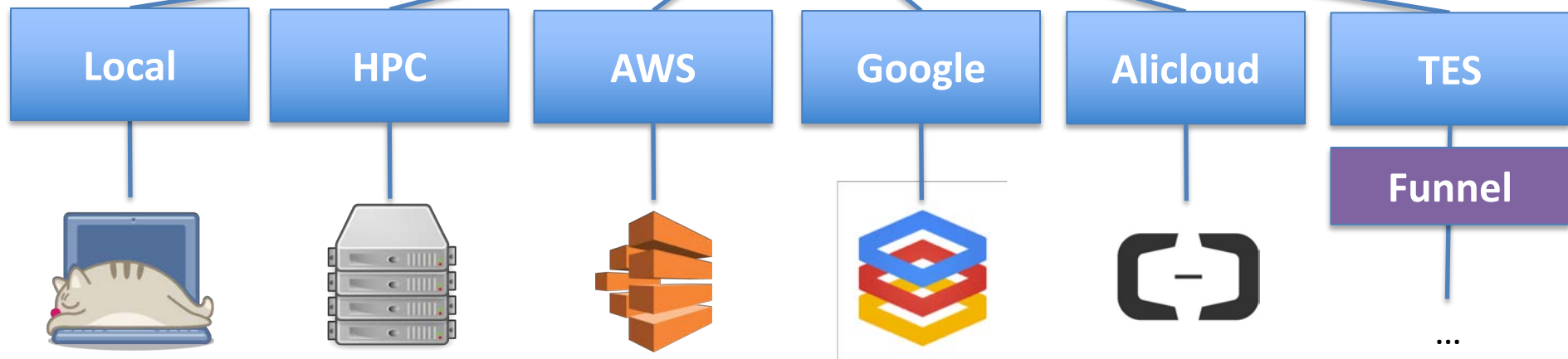  - https://github.com/openwdl/wdl/

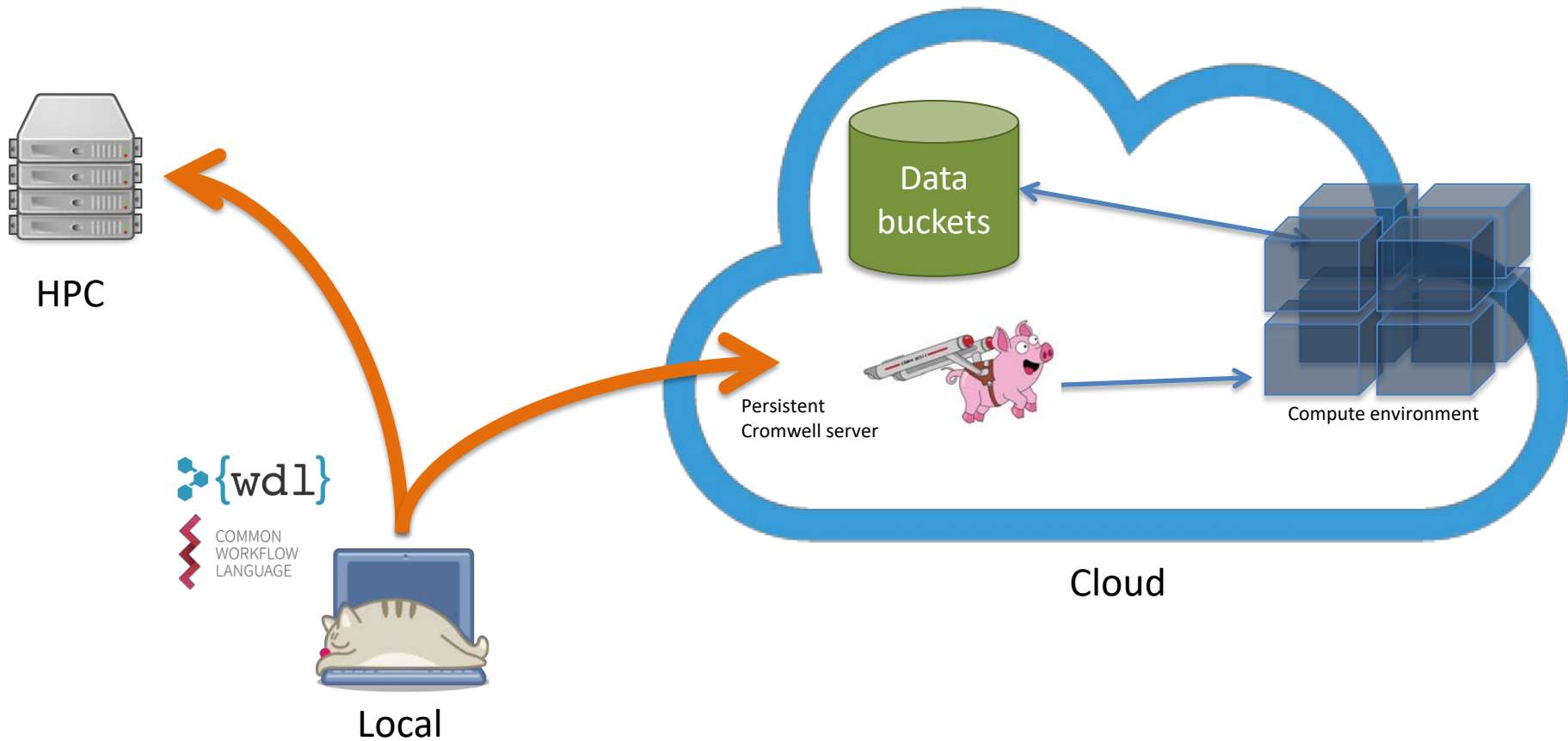# Use containers for portability & reproducibility



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

Takes the guesswork out of running workflows on different platforms!

# Use a workflow execution engine that runs anywhere

# Run using HPC **and** Cloud resources!



HPC

{wdl}

COMMON
WORKFLOW
LANGUAGE

Local

Data buckets

Persistent
Cromwell server

Compute environment

Cloud

## Command Line

- Simple, self-contained command

- Appropriate for independent analysts

- Call Caching

```
java -jar cromwell.jar \
    run hello.wdl \
    hello_inputs.json
```

## Server

- API endpoints

- More scalable, appropriate for production environments

- Call caching

# HPC Backend Configuration

## Managing data

- Root directory

- Data handling strategies

- Support for object stores

```
backend.providers.MyHPCBackend {
    config.root = "/path/to/execution_directory"
    filesystems {
        local {
            localization: [
                "hard-link", "soft-link", "copy"
            ]
        },
        gcs { ... },
        s3 { ... }
    }
}
```

https://cromwell.readthedocs.io/en/stable/tutorials/HPCIntro/

# HPC Backend Configuration

## Managing resources

- CPU

- Memory

- Custom attributes

```
backend.providers.SGE.config {
  runtime-attributes = """
Int cpu = 1
Float? memory_gb
String? sge_queue
String? sge_project
"""

}
```

https://cromwell.readthedocs.io/en/stable/tutorials/HPCIntro/

# HPC Backend Configuration

## Run command

- Built-in variables

- Full flexibility

```
backend.providers.SGE.config {
  check-alive = "qstat -j ${job_id}"
}
```

```
backend.providers.SGE.config {
  kill = "qdel ${job_id}"
}
```

```
backend.providers.SGE.config {
  submit = """
  qsub \
  -terse \
  -V \
  -b y \
  -N ${job_name} \
  -wd ${cwd} \
  -o ${out}.qsub \
  -e ${err}.qsub \
  -pe smp ${cpu} \
  ${"-l mem_free=" + memory_gb + "g"} \
  ${"-q " + sge_queue} \
  ${"-P " + sge_project} \
  /usr/bin/env bash ${script}
  """
}
```

https://cromwell.readthedocs.io/en/stable/tutorials/HPCIntro/

**So of course we decided to create a new one.**

# Workflow description Language (WDL)
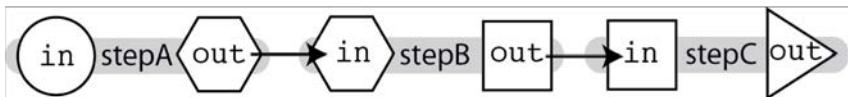
```
task echoHelloWorld {
    command {
        echo 'Hello, World!'
    }
    runtime {
        docker: "phusion/baseimage"
        disks: "local-disk 10 HDD"
        memory: "1 GB"
        preemptible: 3
    }
}

workflow printHelloAndGoodbye {
    call echoHelloWorld
}
```
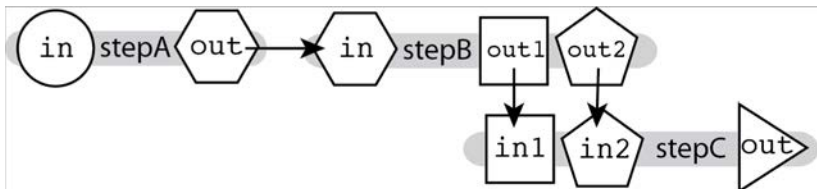
containers

resourcing

cost savings!

# Basic WDL plumbing options

## LINEAR CHAINING


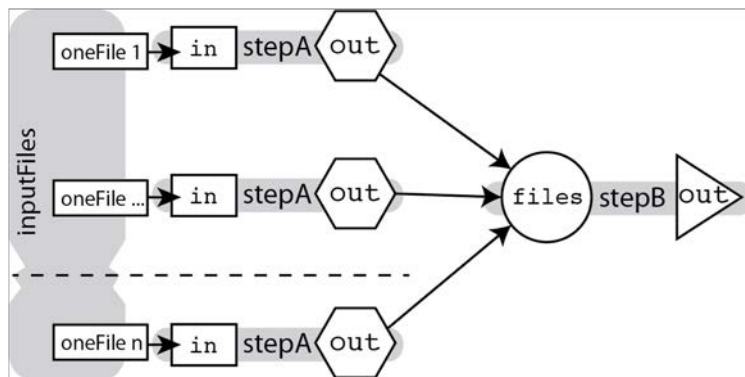
```
call stepA
call stepB { input: in=stepA.out }
call stepC { input: in=stepB.out }
```

## MULTI-IN/OUT



```
call stepC { input :
             in1=stepB.out1,
             in2=stepB.out2 }
```

## SCATTER-GATHER



```
Array[File] inputFiles

scatter(oneFile in inputFiles) {
    call stepA { input: in=oneFile }
}

call stepB { input: files=stepA.out }
```

# But what about CWL?



Multiple workflow languages coming to Cromwell, starting with CWL
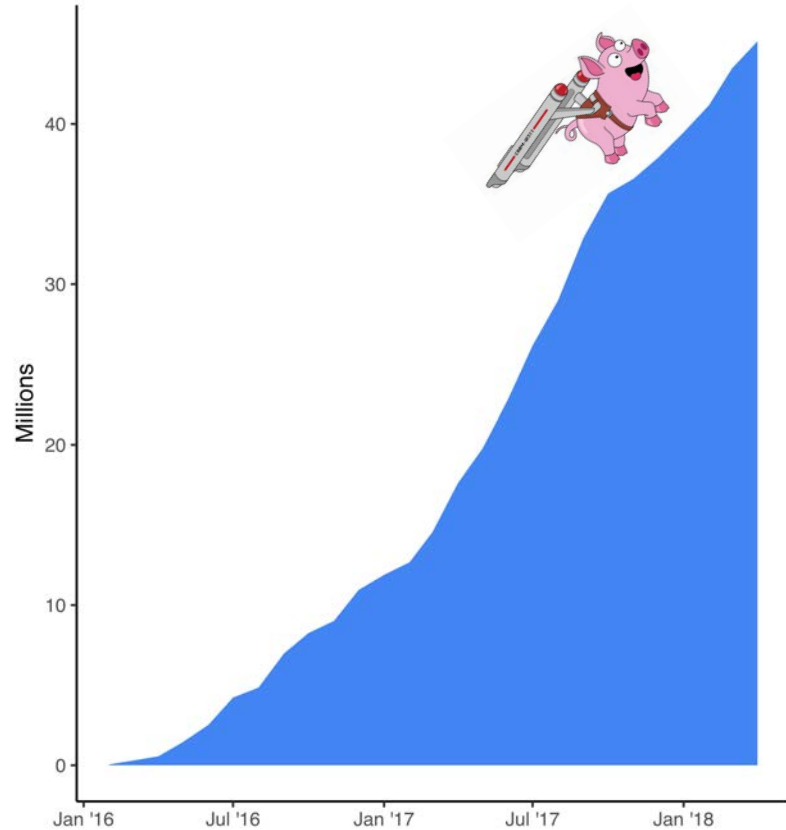
Posted by jgentry on 2 Jan 2018

Thanks to our **Workflow Object Model (WOM)**, Cromwell now supports multiple versions of WDL as well as CWL 1.0!
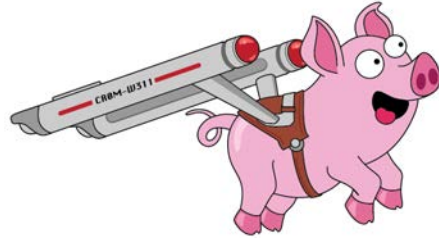
# Cromwell has been busy

**Cromwell in production at Broad:**

Processed **47.5 million jobs**
over the last two years

And this is just the tip of the iceberg!

# Want to discuss further?



**My Email:**
rmunshi@broadinstitute.org

**More Information:**

**Docs:** http://cromwell.readthedocs.io/en/develop/

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

**WDL:** http://www.openwdl.org

**Example Pipelines:** https://github.com/gatk-workflows