



# CONTAINERIZING CONDA

Experiences Building 2000 Applications As  
Portable Docker And Singularity Images

Rion Dooley, John Fonner

@agaveapi

dooley@tacc.utexas.edu, jfonner@tacc.utexas.edu

# HELLO, MY NAME IS...

- ▶ Rion Dooley
- ▶ Manager, Web and Cloud Services Group,
- ▶ Texas Advanced Computing Center (TACC)
- ▶ PI of the Agave project (<https://agaveapi.co>)



# I LIKE LONG WALKS ON THE BEACH AND...

- ▶ That's taken me up and down the stack.
- ▶ Desktop -> Gateway -> Middleware -> SaaS -> PaaS
- ▶ TeraGrid File Manager -> GridChem/XUP -> XSEDE IIS/MPG -> GatewayDNA -> Agave

*Helping people work smarter, collaborate better, and dream bigger*

The logo features a stylized agave plant icon on the left, composed of several white, pointed leaves radiating from a central point. To the right of the icon, the word "Agave" is written in a large, bold, white sans-serif font. Below "Agave", the word "Platform" is written in a smaller, white sans-serif font.

# Agave Platform



AGAVE IS A MULTI-TENANT PAAS  
DELIVERING SCIENCE-AS-A-SERVICE SOLUTIONS  
IN *HYBRID* CLOUD ENVIRONMENTS

Think of it like



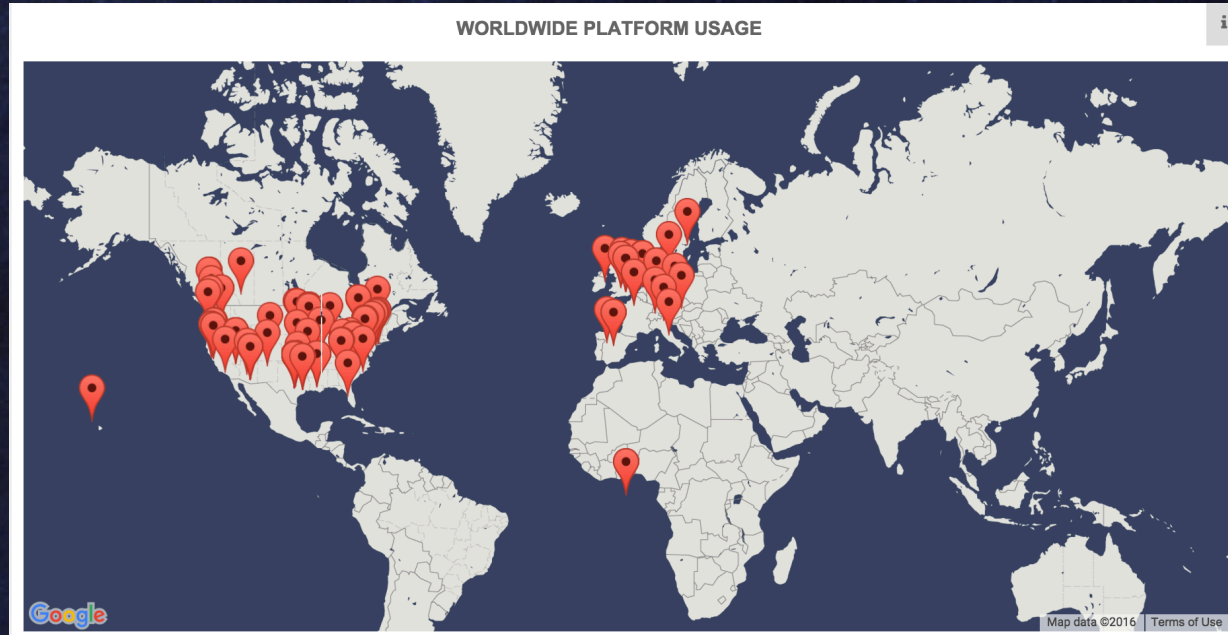
for Science



# WHAT IS AGAVE?

- ▶ Manage it like a cloud platform
- ▶ Scale it like cloud infrastructure
- ▶ Use it like a cloud service
- ▶ Works with your new and legacy infrastructure
- ▶ FOSS - On premise, hybrid, or hosted deployment
- ▶ Multi-tenant
- ▶ Secure by default

# WORLDWIDE USAGE



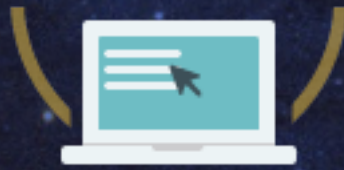


# WHAT DOES IT DO?

MANAGE  
DATA



RUN  
CODE

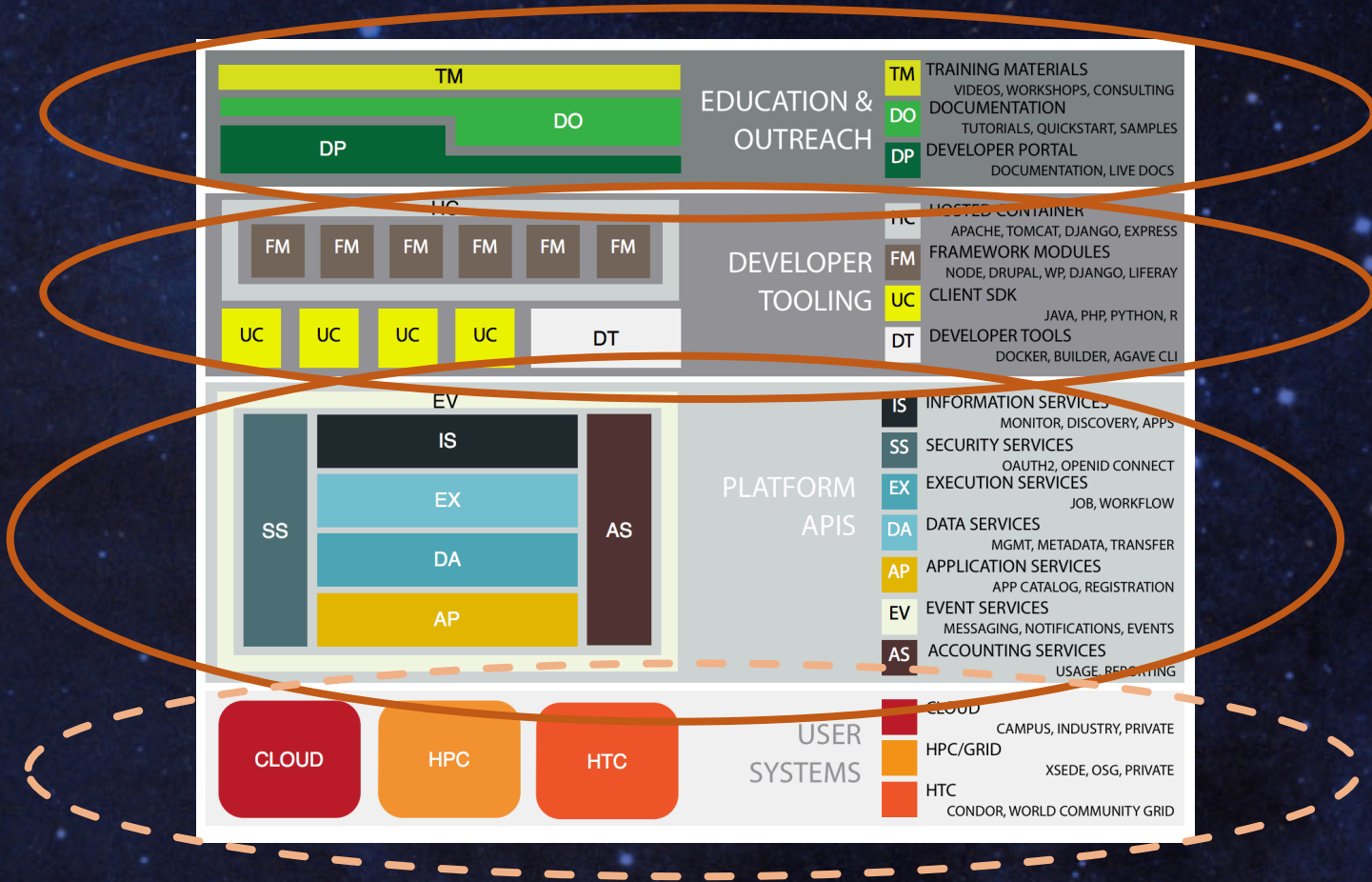


COLLABORATE  
ANYWHERE



CONNECT  
ANYTHING







# WHAT ABOUT CONTAINERS AND CONDA?

- ▶ About 4 years ago, we were in ops hell
- ▶ Build an test, deployments, scaling, configuration management, polyglot services, networking, multi-cloud, stateful everything, app installs, versioning, dependency management
- ▶ Enter a tech talk about Docker...

# WE THOUGHT CONTAINERS WERE COOL

- ▶ Docker was still at 0.4, but we were hooked
- ▶ Encapsulation and the ecosystem solved so many of the problems we had
- ▶ But we weren't just looking at our hosting, we were looking at a bigger picture.

*The juice was worth the squeeze!*



# DOCKER FOR COMPUTATION

Compute containers are the Magic 8 Ball of science...

- ▶ Compartmentalize code
- ▶ Eliminate build and run complexities
- ▶ Introduce portability, reuse, & versioning
- ▶ Widgetize the creation of a scientific pipeline

...but better because results are reproducible.



*Compute containers enable reproducible science by composition.*

# DOCKER FOR DATA

Data containers serve as universal adapters between compute containers

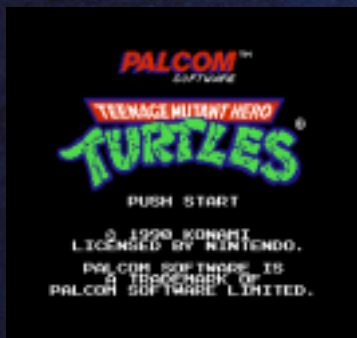
- ▶ Transform data
- ▶ Bridge file systems
- ▶ Enable distributed data access
- ▶ Virtualize interfaces



*Data containers* enable clean integration between containers and *standardize* how we interact with *distributed data*.



# WHAT DO THESE HAVE IN COMMON?



# KONAMI CODE





# DOCKER FOR RESEARCH

- ▶ Compute and data containers *are* cheat codes when it comes to science
- ▶ Building a library of “Science Apps” in the Docker Hub
- ▶ Built from a few value-added trusted base images
  - ▶ Attribution
  - ▶ Extended metadata
  - ▶ Documentation
  - ▶ Extension points
  - ▶ Common CLI

Science apps are *cheat codes* when it comes to discovery.

# DOCKER FOR RESEARCH

The screenshot shows the Bitbucket overview page for the repository 'taccaci/sci-apps'. The page includes a sidebar with navigation icons, a main content area with repository statistics, and a description of the repository's purpose.

**Overview**

SSH: git@bitbucket.org:taccaci/sci-apps.git

Last updated: 2014-10-24	1 Branch	0 Tags
Language: —	0 Forks	2 Watchers
Access level: Admin [revoke]		

**TACC Scientific Apps Repository**

The goal of this project is to create a repository of **Docker** images for scientific applications. The main repository of images is the **Docker hub**. The work here does not try to replicate such content. We want to add some features and conveniences for developers to the subset of applications most commonly used by scientists in the TACC environment.

The features in the current images are:

- Included documentation:** we define a standard way to add documentation and examples to the image. An containerized application in this framework automatically accepts a set of common options, such as `--help`, `--help` for examples of use, etc. (see [ref:using](#)).
- Attribution metadata:** an uniform way to specify version, authors, license, and other metadata. This can be retrieved through command line options, same as the documentation.
- Analytics:** the containers can automatically send information about the application being run to a central server, to collect usage and analytics. This feature can be disabled, if necessary.
- Developer friendliness:** we want to make it easy for developers to extend the images. The images provide simple ways to translate the above features to new images through inheritance. It is easy to add background processes and to set the environment for an application.

**Recent activity**

- 1 commit: Pushed to taccaci/sci-apps by [Walter Meneses](#) · 2014-10-24
- 1 commit: Pushed to taccaci/sci-apps by [Walter Meneses](#) · 2014-10-24
- 1 commit: Pushed to taccaci/sci-apps by [Walter Meneses](#) · 2014-10-24
- 1 commit: Pushed to taccaci/sci-apps by [Walter Meneses](#) · 2014-10-24

The screenshot shows the Docker Hub page for the organization 'taccsciapps'. The page lists several repositories, including 'samtools', 'fastqc', 'bwa', 'ipython', and 'bioperl'.

**taccsciapps**

Repositories

Repository	Updated	Stars	Forks
taccsciapps/samtools	3 months ago	0	0
taccsciapps/fastqc	3 months ago	0	0
taccsciapps/bwa	3 months ago	0	0
taccsciapps/ipython	3 months ago	0	0
taccsciapps/bioperl	3 months ago	0	0

<https://bitbucket.org/taccaci/sci-apps/overview>



# AND THEN NOTHING HAPPENED

- ▶ Well, we got tired
- ▶ And frustrated
- ▶ And lost a lot of sleep
- ▶ And wrote a bunch of code that will never see the light of day...except the previous slide.

# IT'S THE ECOSYSTEM STUPID

- ▶ We forgot the very reason that made us love containers in the first place, the ECOSYSTEM!
- ▶ We returned to our roots and started looking for places we could fill in the gap between IT and research
- ▶ Better base images and Dockerfiles
- ▶ Integration with autobuilds
- ▶ Private repos as part of your personal storage
- ▶ Transparent build and run on demand



# BUILDING TOWARDS SOMETHING BIGGER

- ▶ That went pretty well. We added several hundred docker images to the Cyverse infrastructure.
- ▶ Ran 100k or containers
- ▶ Taught users how to contribute their own images
- ▶ Then we hit the Peter Covney problem
- ▶ Just not enough cycles and/or hours in the day for the demand

# THEN THE WHEELS CAME OFF

- ▶ Containers are portable. TACC has lots of cores. Problem solved
- ▶ Except Docker won't or can't run on many HPC systems
- ▶ Security issues
- ▶ Scheduler issues
- ▶ Data issues
- ▶ Repository issues
- ▶ Networking issues
- ▶ Accounting issues



# ENTER SINGULARITY

- ▶ Singularity is “Containers for HPC”
- ▶ Singularity is NOT Docker
- ▶ Technology aside, there is no ecosystem, which made onboarding users nearly impossible
- ▶ But the potential is significantly larger to serve the research community since it's a portable, secure runtime that can swap out 1 for 1 with existing system calls.

# CONDA GETTING THE HANG OF IT

- ▶ Built our own automated build farm to convert Docker image to Singularity images
- ▶ 1 for 1 on published images
- ▶ Published the entire BioConda repository as Singularity images
- ▶ Added transparent support for one click execution of any image.
- ▶ Published images as publicly available downloads



# STRENGTH IN NUMBERS

- ▶ We're not alone...not even close
- ▶ Currently working with Vanessa Sochat from the Singularity team to make images available on the new singularity-hub.org and public registry.
- ▶ Working with Björn Grüning on the BioContainers side to upstream generating the Singularity images, something they already want to do.
- ▶ Working with Common Workflow Language team to ensure we have a standard for invoking containers, regardless of the runtime.

# DOING OUR PART

- ▶ Added first-class support for running Singularity containers in the Agave Platform.
- ▶ Working on REPL support and additions to the Agave CLI tools to help people go from daydream to discovery faster than ever.
- ▶ Incorporating the container ecosystem into all TACC resources.
- ▶ Clearing a smooth migration path for computation between local, cloud, htc, and HPC resources.
- ▶ Creating additional content to raise awareness, train, and support users as they begin looking at these technologies.





THANK YOU!

FOLLOW US

@agaveapi

agaveapi.co