

# docker for Data Science



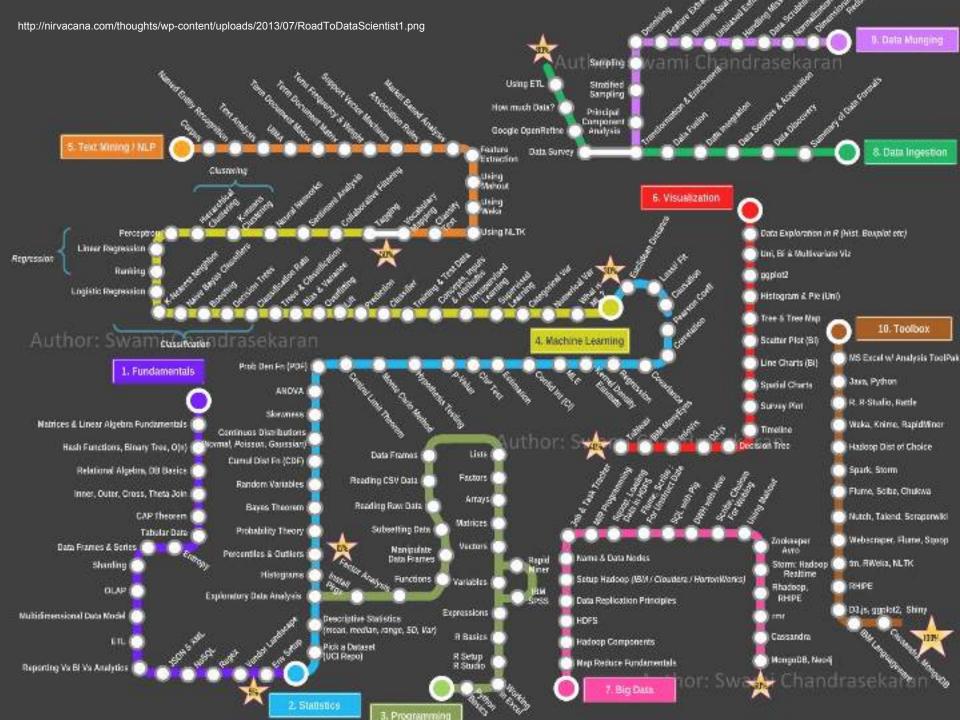
https://www.meetup.com/Data-Science-Meetup-Muenster



# **Daniel Nüst**

Institute for Geoinformatics
University of Münster

@nordholmen | <a href="http://nordholmen.net">http://nordholmen.net</a>



# **Docker for Data Science**

http://blog.kaggle.com/2016/02/05/how-to-get-started-with-data-science-in-containers/

<u>https://github.com/wiseio/datascience-docker</u> (Hackday container - nice!)

http://www.datadan.io/containerized-data-science-and-engineeringpart-2-dockerized-data-science/

http://www.slideshare.net/CalvinGiles/docker-for-data-science

https://civisanalytics.com/blog/data-science/2016/05/11/strata-2016 -talk/

https://www.quora.com/What-are-use-cases-for-Docker-in-Data-Science-and-Machine-Learning

"Isolation! Portability! Repeatability!"



# Agenda

What is Docker? Why?

What can it be used for?

Live Demo

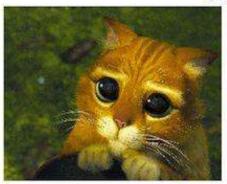


# Why containerization? Why Docker?

# **Motivation**

Pets vs. Cattle

# Service Model



- Pets are given names like pussinboots.cern.ch
- They are unique, lovingly hand raised and cared for
- When they get ill, you nurse them back to health



- Cattle are given numbers like vm0042.cern.ch
- They are almost identical to other cattle
- When they get ill, you get another one

 Future application architectures should use Cattle but Pets with strong configuration management are viable and still needed

# Motivations for Docker in mainstream IT

https://www.docker.com/use-cases



#### CI/CD

Enable developers to develop and test applications more quickly and within any environment



#### DEVOPS

Break down barriers
between Dev and Ops teams
to improve the app
development process



#### BIG DATA

Empower your enterprise to leverage big data analytics



# INFRASTRUCTURE OPTIMIZATION

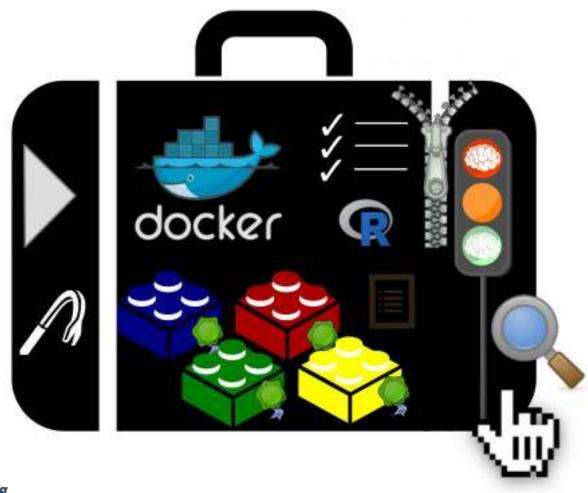
Decrease infrastructure costs while increasing its efficiency

# Motivation for Reproducible Research





# Executable Research Compendium





# Motivation for Data Scientists - Proof by quote

"There was only one problem — all of my work was done in my local machine in R. People appreciate my efforts but they don't know how to consume my model because it was not "productionized" and the infrastructure cannot talk to my local model. Hard lesson learned!"

-Robert Chang, data scientist at Twitter

"Data engineers are often frustrated that data scientists produce inefficient and poorly written code, have little consideration for the maintenance cost of productionizing ideas, demand unrealistic features that skew implementation effort for little gain... The list goes on, but you get the point". -Jeff Magnusson, director of data platform Stitchfix

But don't worry! There is a better way: **Dockerize your data** science applications for ease of deployment, portability, and integration within your infrastructure.

via Daniel Whitenack, <a href="http://www.datadan.io/containerized-data-science-and-engineering-part-2-dockerized-data-science/">http://www.datadan.io/containerized-data-science-and-engineering-part-2-dockerized-data-science/</a>

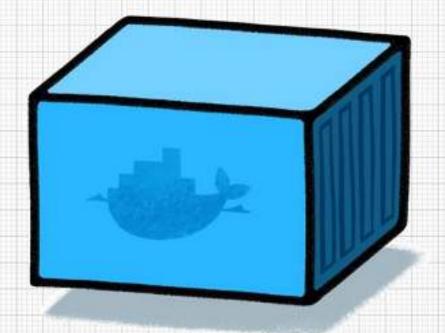


Home/Posts
About
Contact
Taiks/Schedule

(#) github

III linkedin

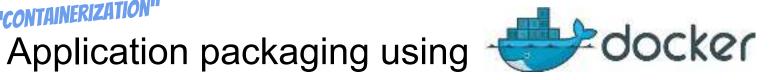
The real value of Docker is not technology



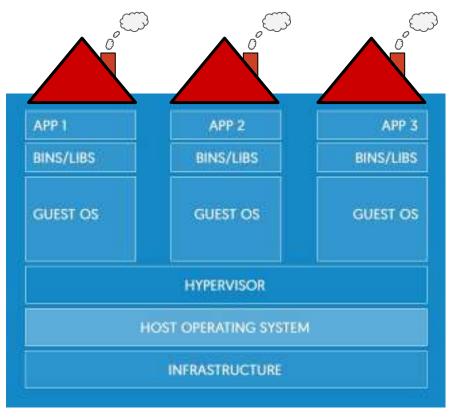
It's getting people to agree on something

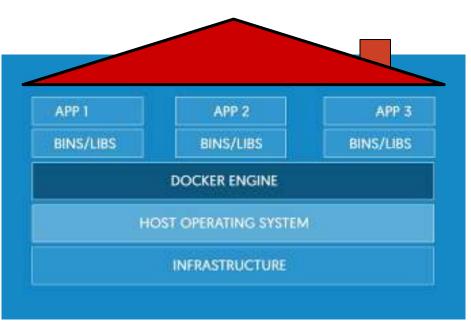
Slide by Docker inventor & Docker, Inc. CTO Solomon Hykes, DockerCon 2014

# "CONTAINERIZATION"



#### Houses vs. Appartments | "binary" vs. OS



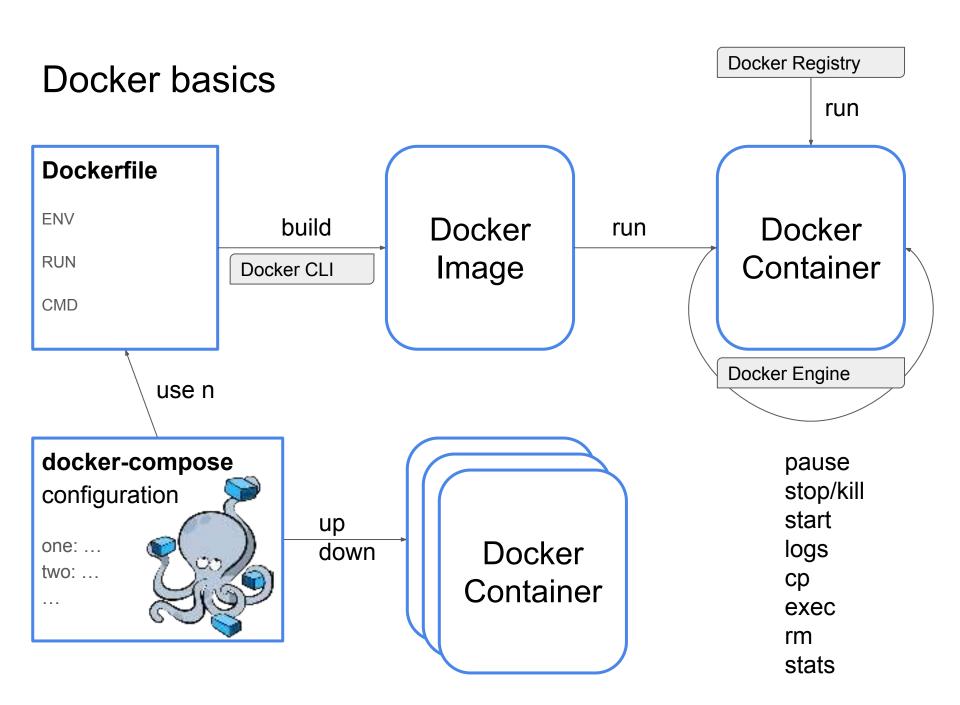


https://www.docker.com/what-docker

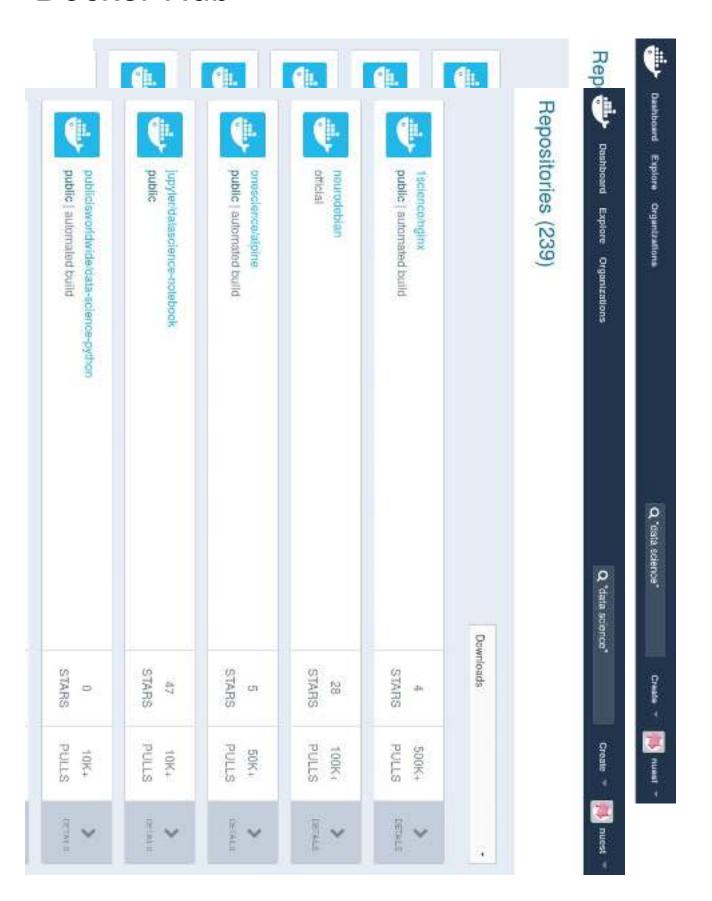
https://en.wikipedia.org/wiki/Operating-system-level virtualization

https://youtu.be/ki8CZkutoxQ

kernel features namespaces libcontainer, LXC cgroups resources



# **Docker Hub**







#### https://hub.docker.com/r/rocker/rstudio/

docker run --rm -it -p 8787:8787 rocker/rstudio

http://localhost:8787/ (rstudio/rstudio)

Great example: <a href="https://github.com/benmarwick/1989-excavation-report-Madjebebe">https://github.com/benmarwick/1989-excavation-report-Madjebebe</a>

docker run --rm -it -p 8787:8787 benmarwick/mjb1989excavationpaper <a href="http://localhost:8787/">http://localhost:8787/</a> (rstudio/rstudio)



# Jupyter Notebook



https://github.com/jupyter/docker-stacks/tree/master/datascience-notebook

#### wget

https://raw.githubusercontent.com/datascience-meetup-muenster/talks/master/meetup-01/what a re data scientists and tools.ipynb

docker run -it -p 8888:8888 -v \$(pwd):/home/jovyan/work jupyter/datascience-notebook

http://localhost:8888/

See also <a href="https://www.dataquest.io/blog/docker-data-science/">https://www.dataquest.io/blog/docker-data-science/</a>



# **ELK** stack

```
git clone https://github.com/deviantony/docker-elk.git
cd docker-elk
# add filter to logstash/config/logstash.conf:
# filter {
# grok { match => { "message" => " %{COMBINEDAPACHELOG}"}}
# }
#}
docker-compose up
http://localhost:5601/app/kibana
http://localhost:9200/
```

Example data: http://ita.ee.lbl.gov/html/contrib/NASA-HTTP.html
nc localhost 5000 > access\_log\_Aug95
docker-compose down (-v)



#### https://hub.docker.com/r/sverhoeven/cartodb/

```
docker run --rm -it -p 3000:3000 -p 8080:8080 -p 8181:8181
--name carto sverhoeven/cartodb
sudo sh -c 'echo 127.0.1.1 cartodb.localhost >> /etc/hosts'
docker run --rm -it -p 80:80 --link carto:cartodb.localhost
spawnthink/cartodb-nginx
```

http://cartodb.localhost dev/pass1234

# Deep convolutional network in Amazon Cloud



https://civisanalytics.com/blog/data-science/2016/05/11/strata-2016-talk/

```
https://github.com/mdagost/pug_classifier
git clone https://github.com/mdagost/pug_classifier.git
docker run -d -p 8888:8888 -v /home/ubuntu/pug_classifier:/home/jovyan/work
mdagost/pug_classifier_notebook
```

# Interested in "geo"? Go to OSGeo wiki +

https://wiki.osgeo.org/wiki/DockerImages

https://wiki.osgeo.org/wiki/DockerImagesMeta

http://geocontainers.org

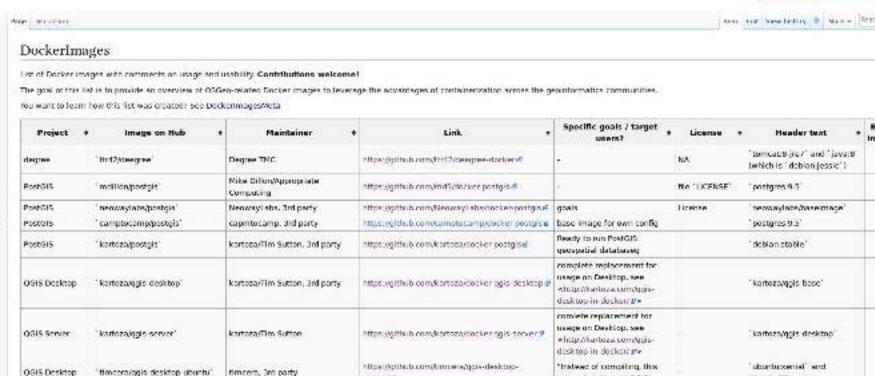


A Danie Plant Talk Preferences V

BhuntuGIS repos

image is j. I latest QSIS





thurtus.

# Core arguments for Data Scientists

(all the Docker advantages... write once, biz ops, cloud, etc.)

# Reproducibility

Project separation + don't clutter dev machine

Environment (re)creation, documentation

Adopt good practices on the way (dev cred)

Easy collaboration

Easy transition from testing to production

# More from the Dockerverse

Docker **Machine** (provision remote host or clusters)

Docker Cloud (hosting of Dockerized apps)

Docker **Toolbox** (for older Mac and Windows OS)



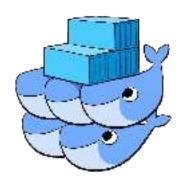


Docker Universal Control Pane (custer management and monitoring UI)

Docker **Swarm mode** (container orchestration)

Docker Trusted Registry (own enterprise image storage)





# Thanks for your attention!

# What are your questions?

https://github.com/nuest

http://www.slideshare.net/nuest/

http://nördholmen.net

http://o2r.info daniel.nuest@wwu.de







Want more Docker?
Watch Dockercon Keynote!



http://bit.ly/2cjrqQl