



Data Engineering 101

Building Data Pipelines

March 10th, 2015

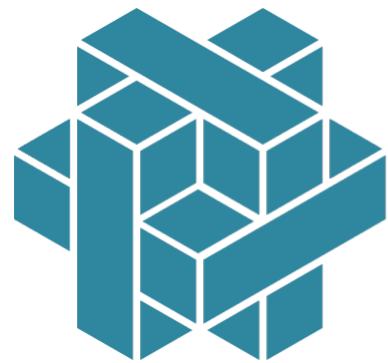


galvanize



Jonathan Dinu
VP of Academic Excellence, Galvanize
jonathan@galvanize.com
[@clearspandex](https://twitter.com/clearspandex)

Currently



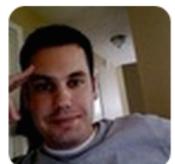
Zipfian
Academy



galvanize



- What is Data Engineering?
- Why is Data Engineering?!
- How is Data Engineering?!?
- Data Architectures
- Building a Pipeline (w/ Luigi)
- Q&A



Josh Wills

@josh_wills



+
Follow

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

Reply Retweet Favorite More



Josh Wills

@josh_wills

Engineer

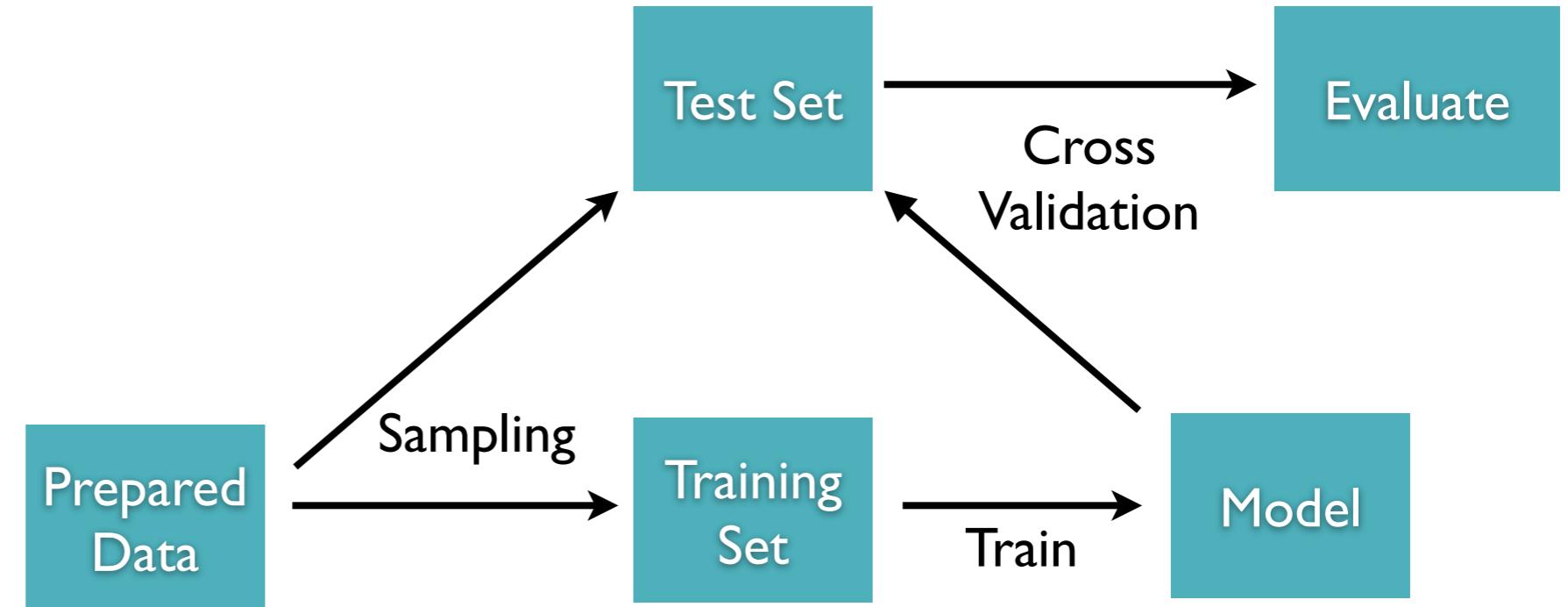


+
Follow

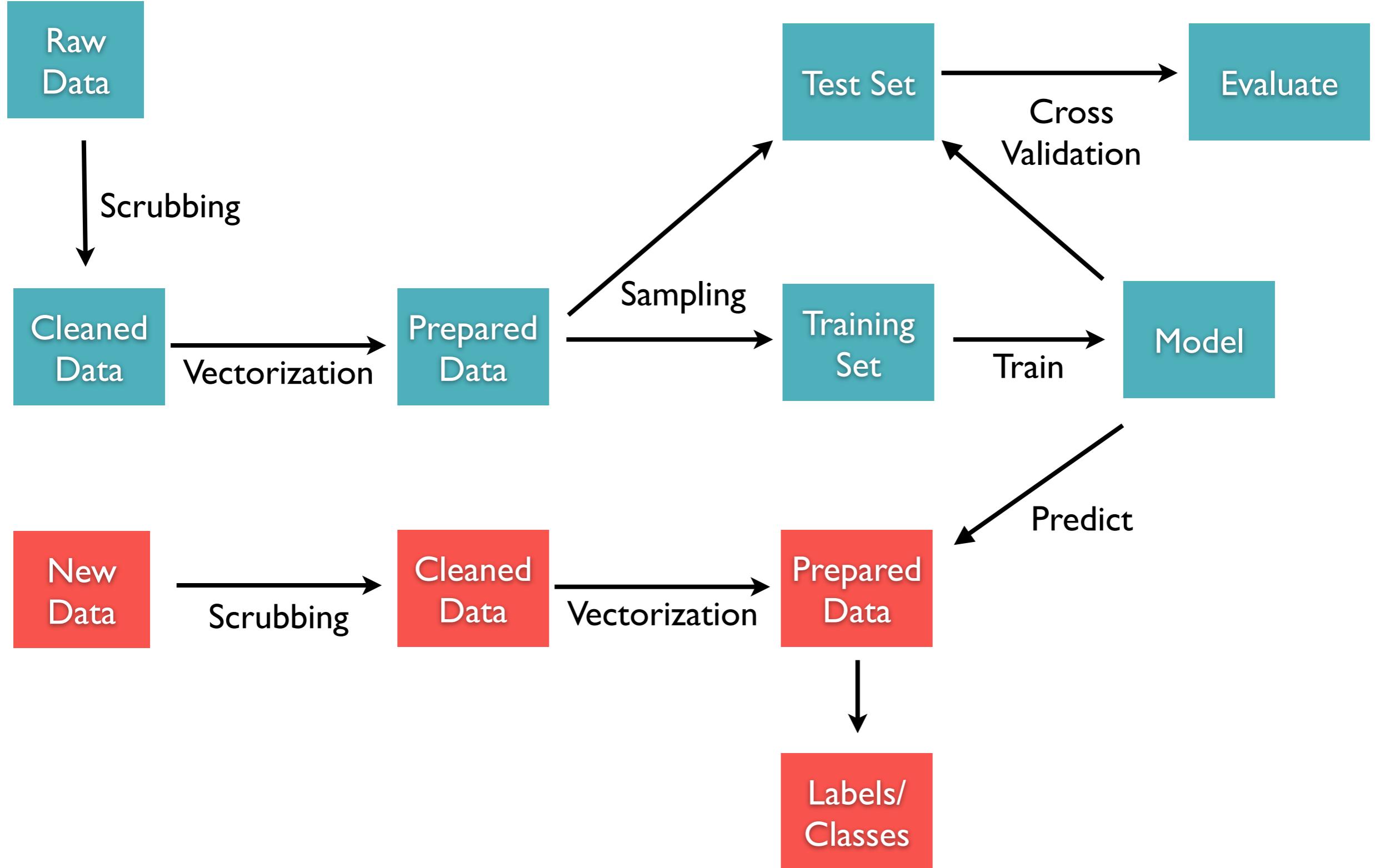
Data ~~Scientist~~ (n.): Person who is better at statistics than any software engineer and better at software engineering than any ~~statistician~~. Data Scientist

Reply Retweet Favorite More





Data Engineering





The Challenge



@goGIFGIF for updates.

GIFGIF

Search | Results | Data | About

Which better expresses **excitement?** [Change Question](#)

Info Share

NEITHER

Info Share

Achievements [Import | Export | Delete](#)

Your votes: 44

Analysis: Excitement

Best Worst

Global votes: 2,530,975

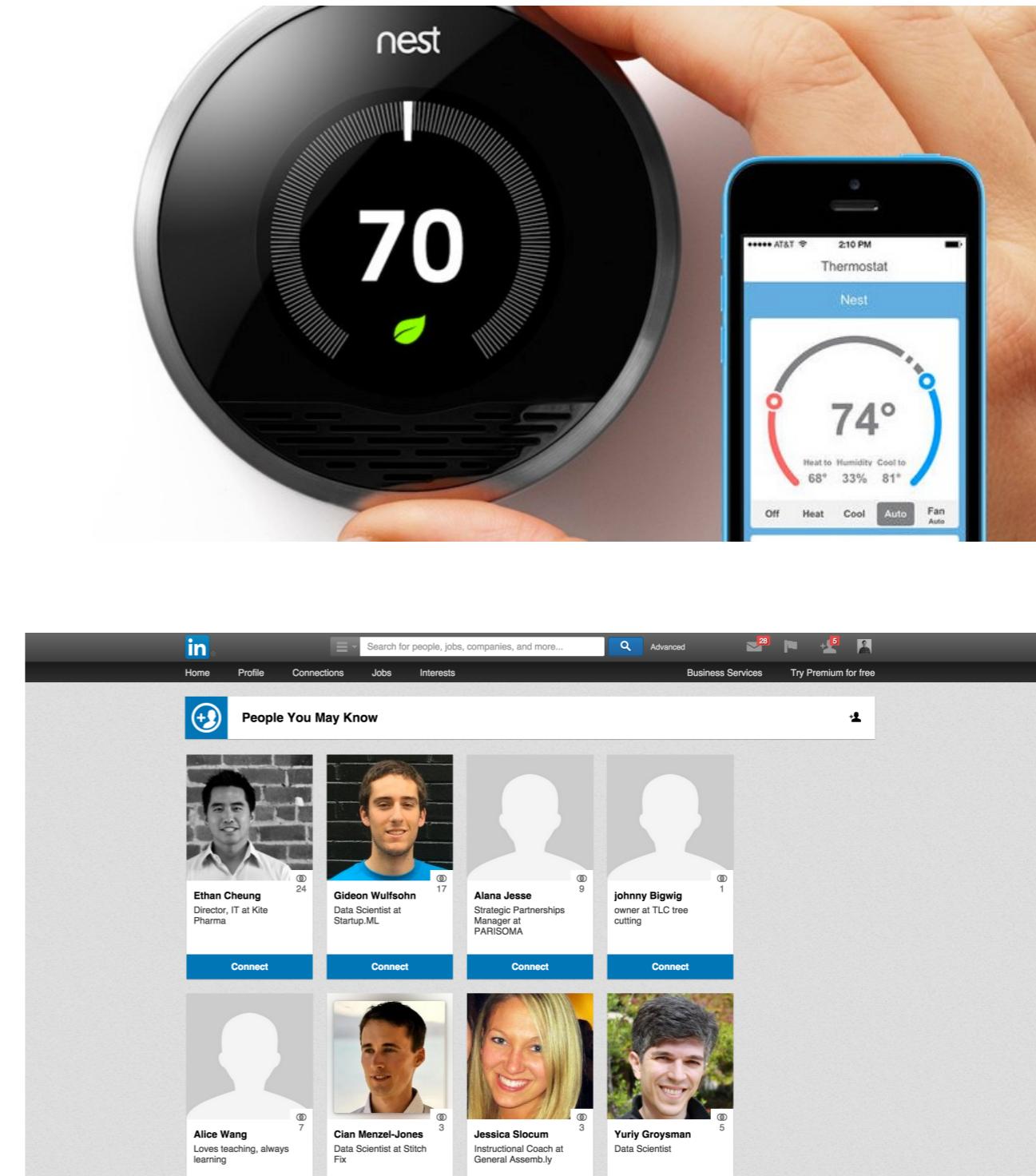
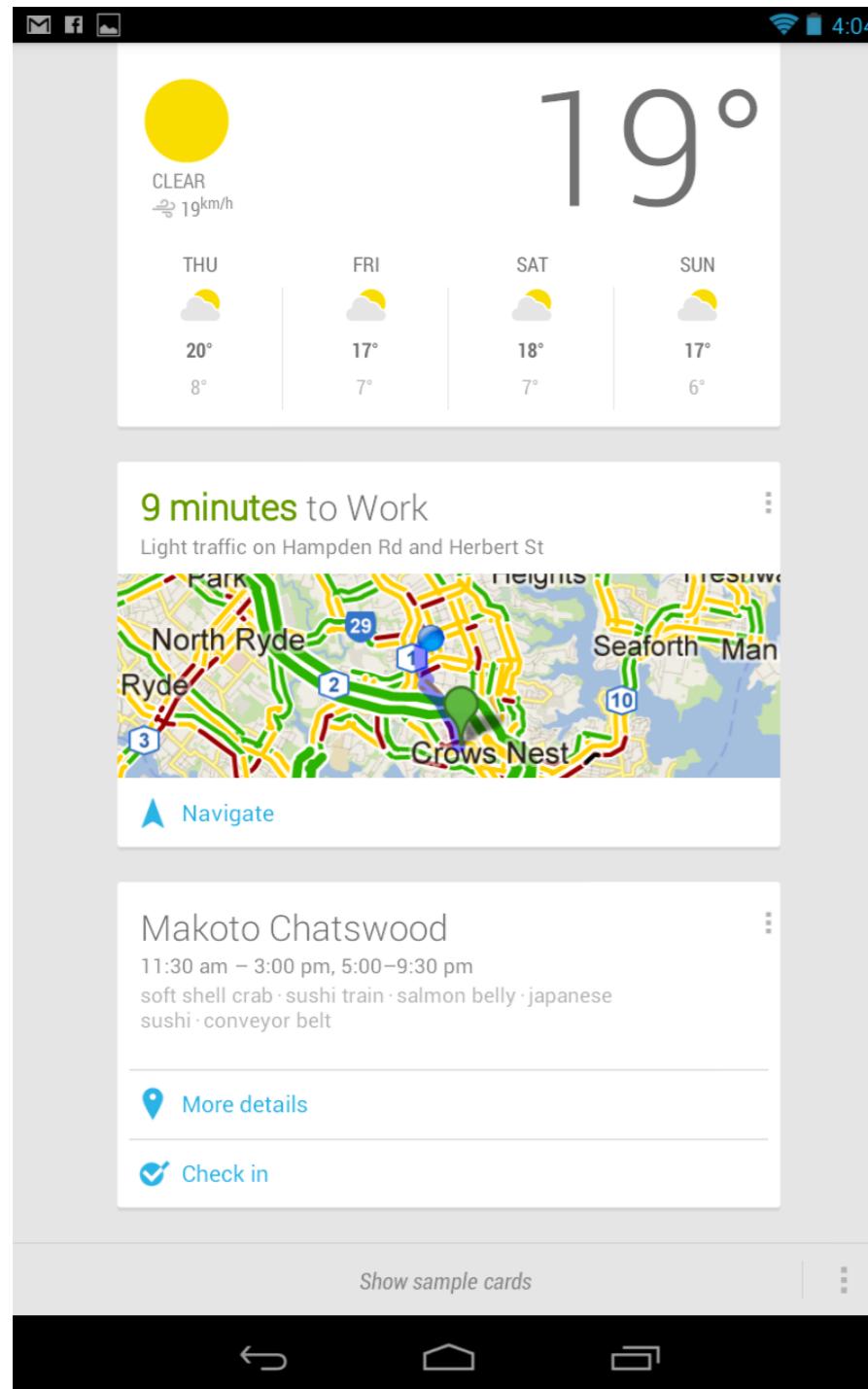


Product Built on Data



Product that Generates Data

Products



Questions? tweet @clearspandex



Product that Generates Data (that you sell)

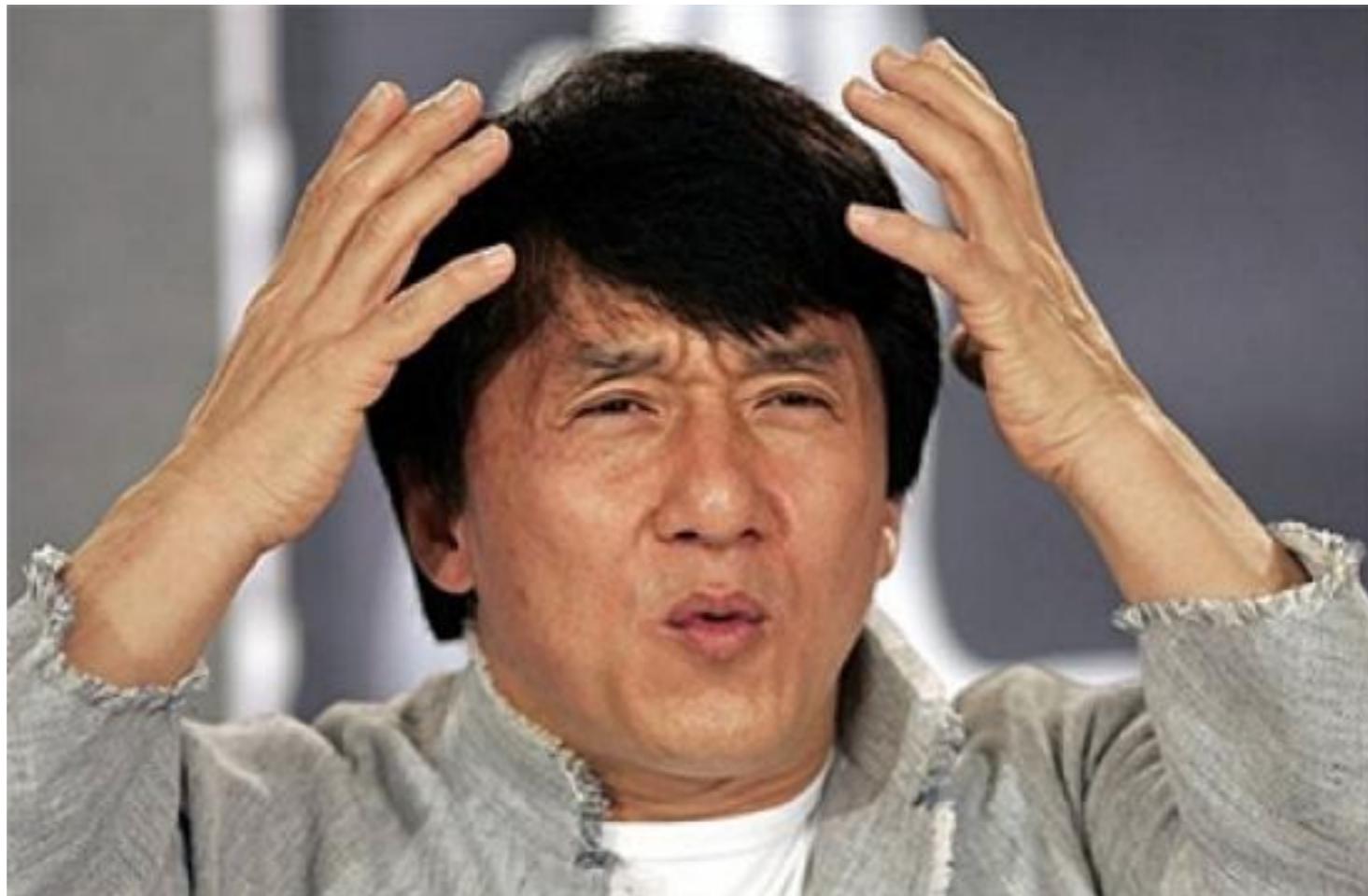


Product that Generates Data (that you sell)

i.e. Facebook



But.... How?!?!!?



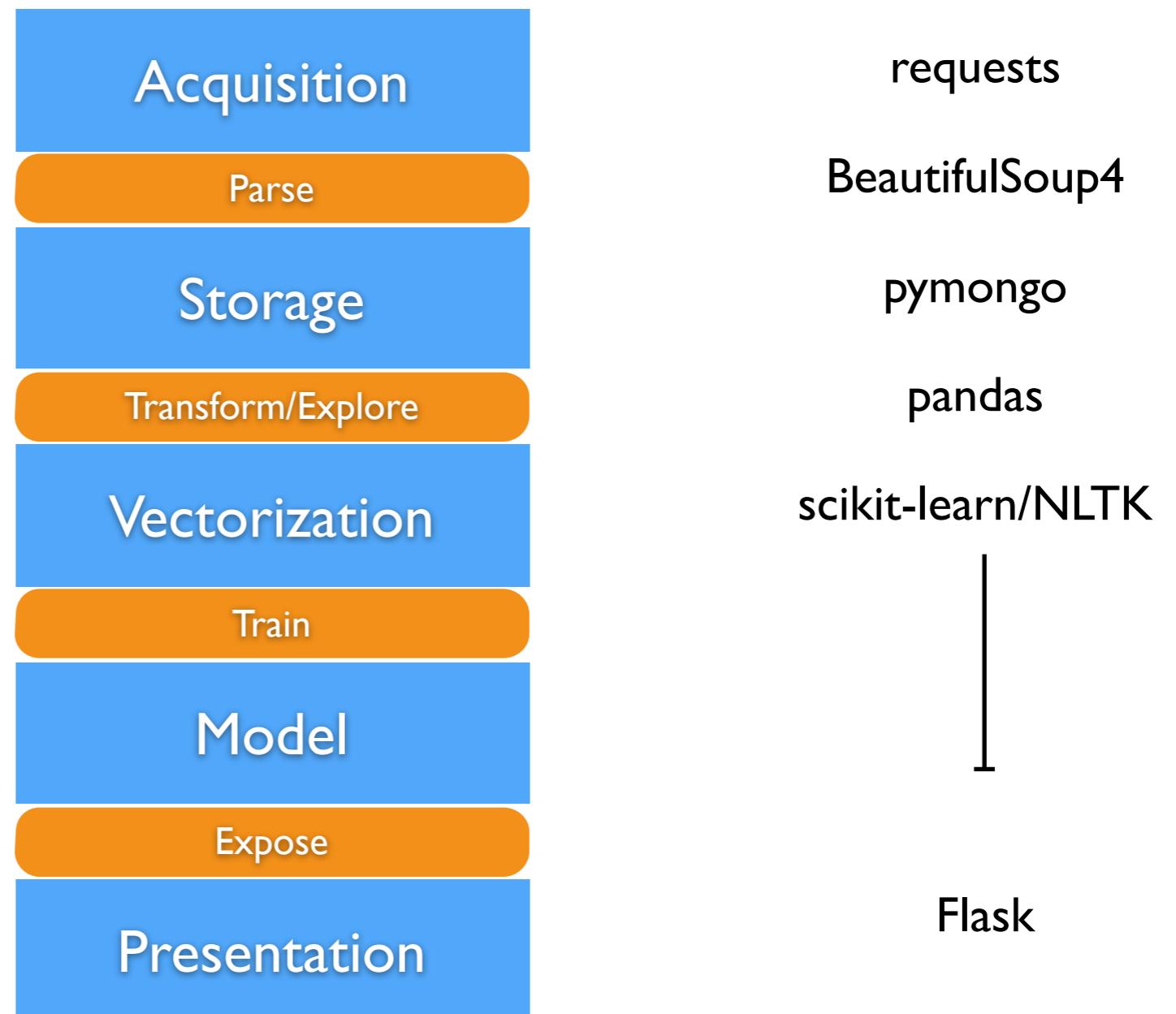


At Scale

scrapy
Hadoop Streaming
(w/ BeautifulSoup4)
Snakebite (HDFS)
mrjob or luigi
Spark ML (pySpark)



Flask



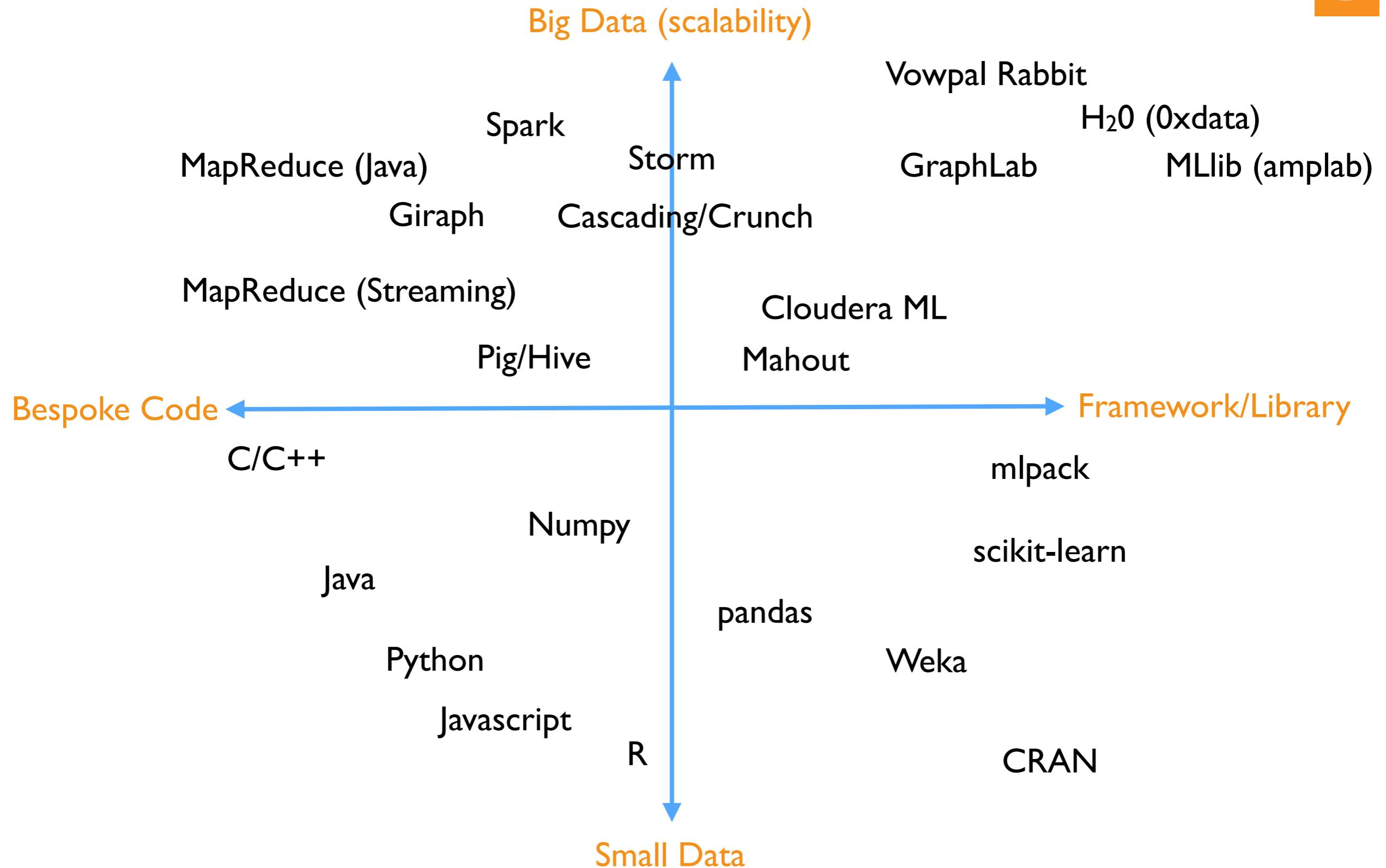
Locally

requests
BeautifulSoup4

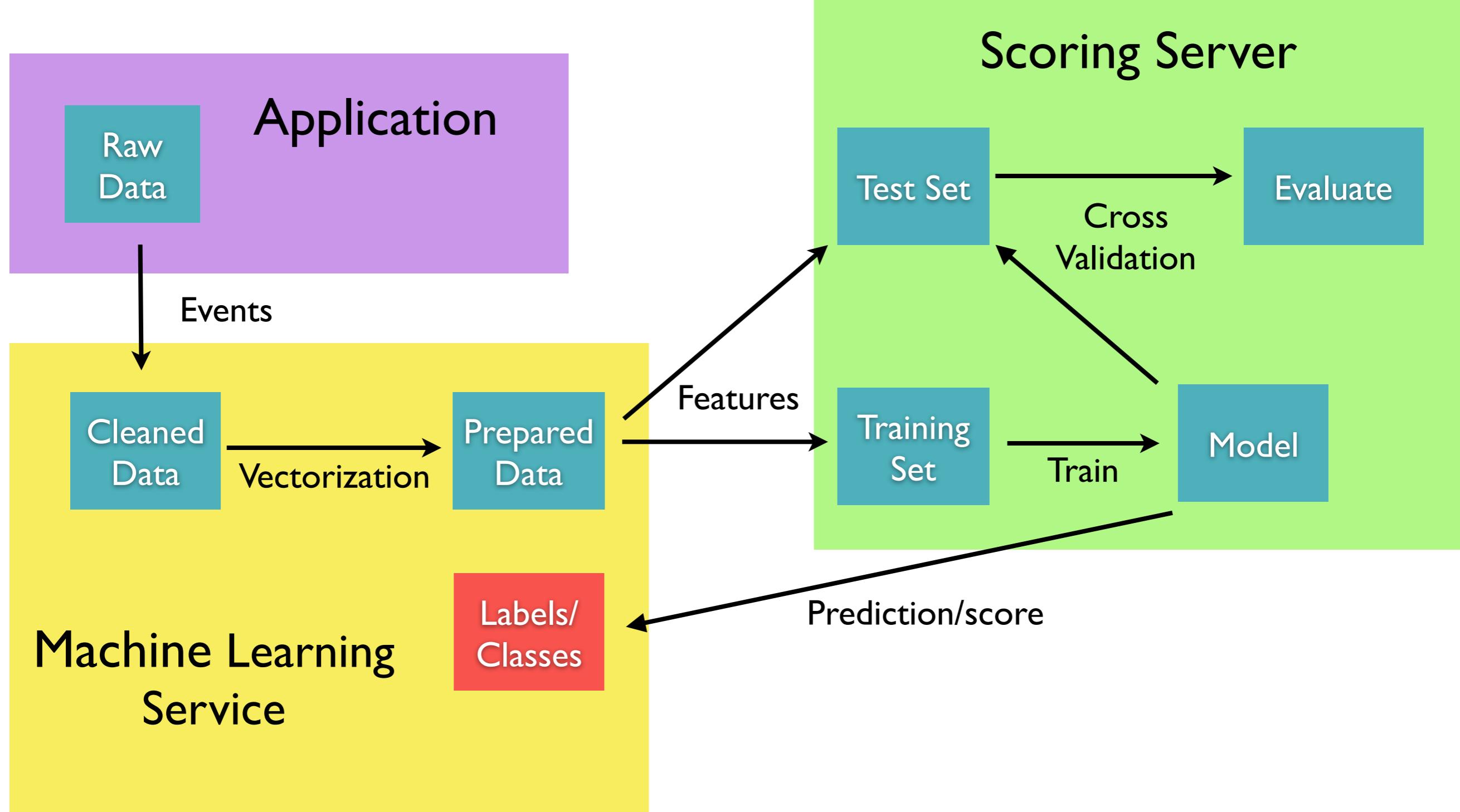
pymongo
pandas

scikit-learn/NLTK

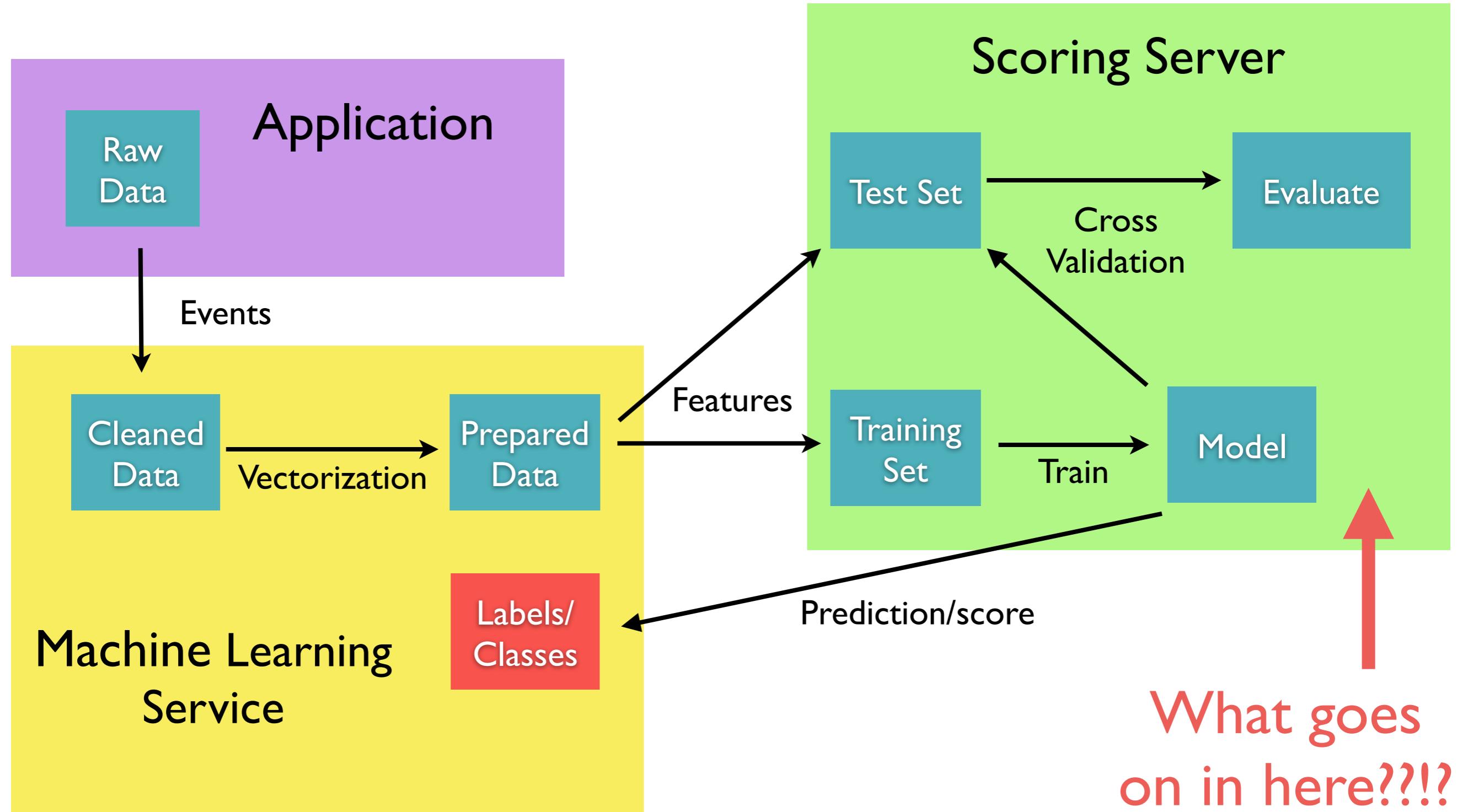
Flask



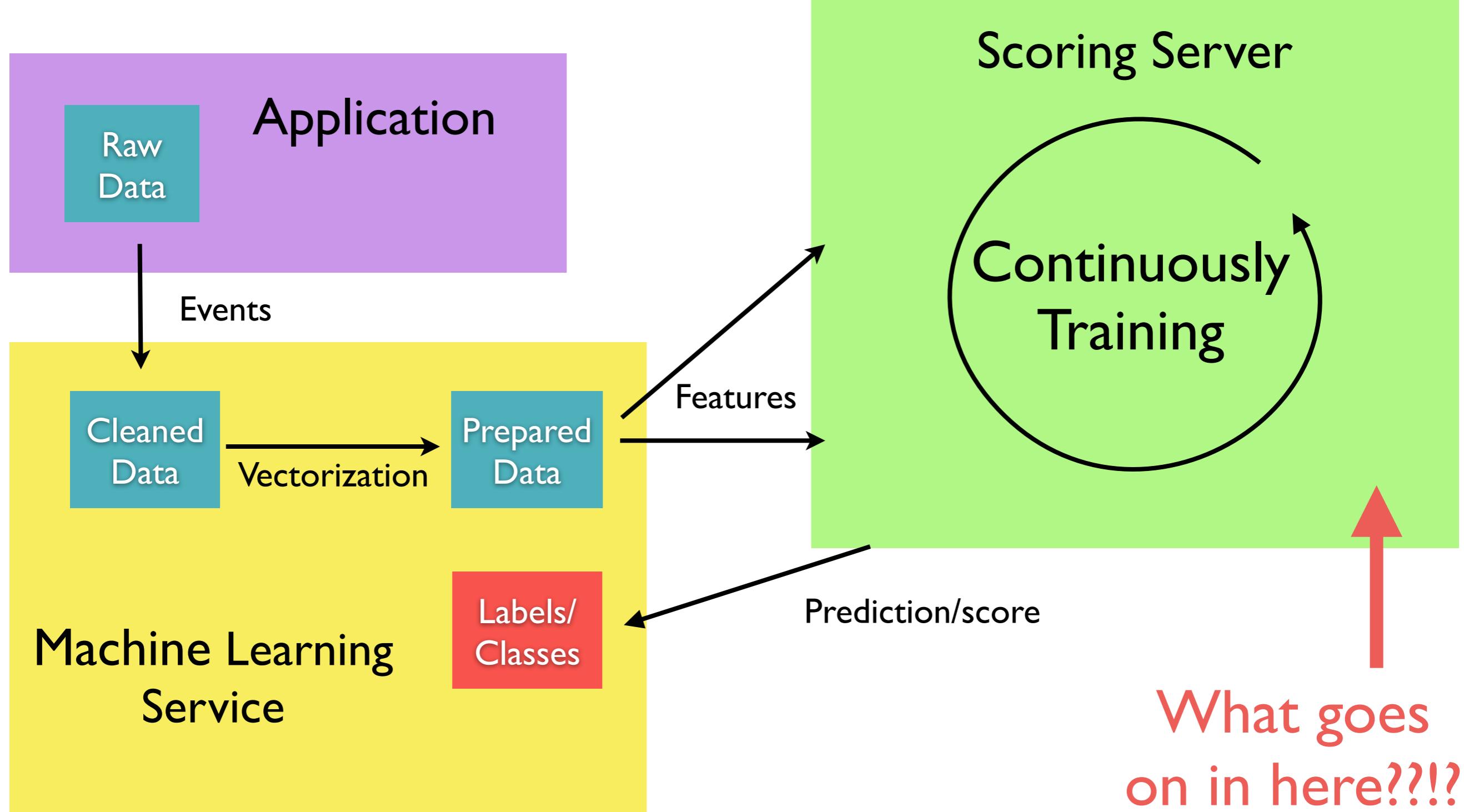
Machine Learning



Machine Learning



Machine Learning



Code



LIVE CODE

Code



u.9t



(it's the pipes!)



Why Pipelines?

- Always keep raw data
- Data Lineage
- Apply a series of transforms to data.
- Flexible, Modular, Extensible (and testable!)

Why Luigi

- Idempotence
- Checkpointing
- Native Hadoop Support
- But Works for arbitrary scripts (like make!)

Why Luigi

- Idempotence
- Checkpointing
- Native Hadoop Support
- But Works for arbitrary scripts (like make!)

(also has a nice UI and sends emails :)



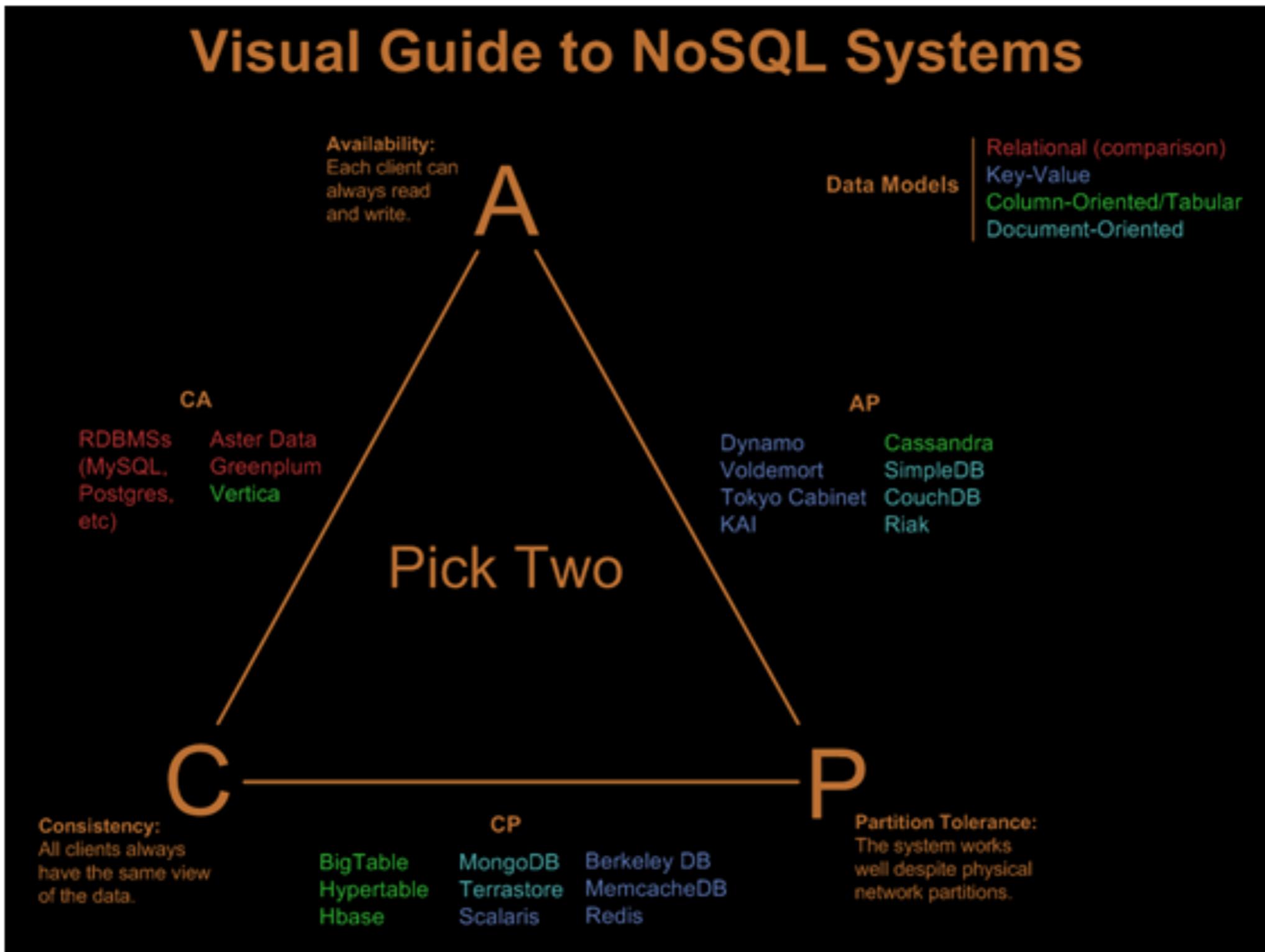
Modern Architectures

Architecture

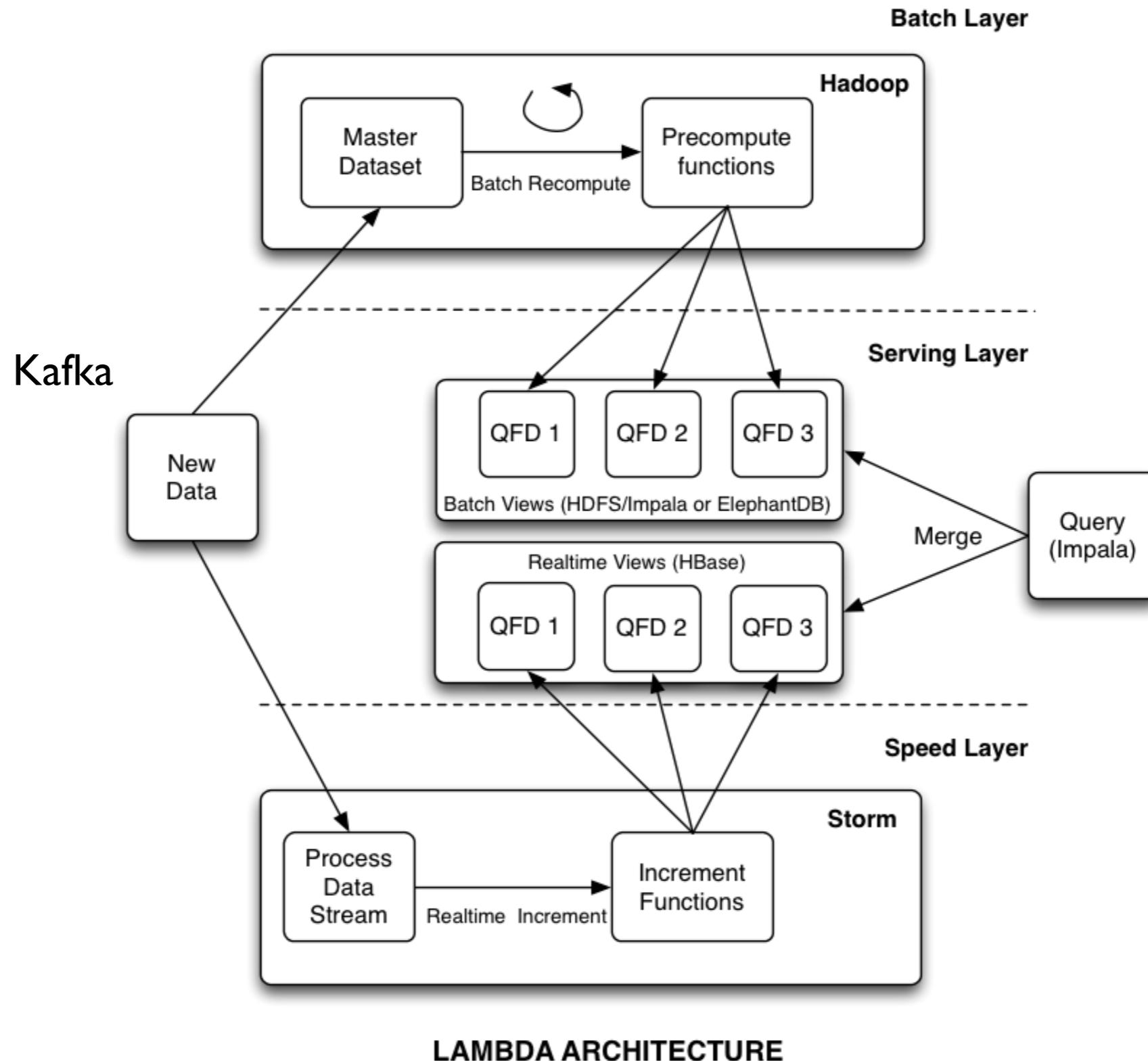




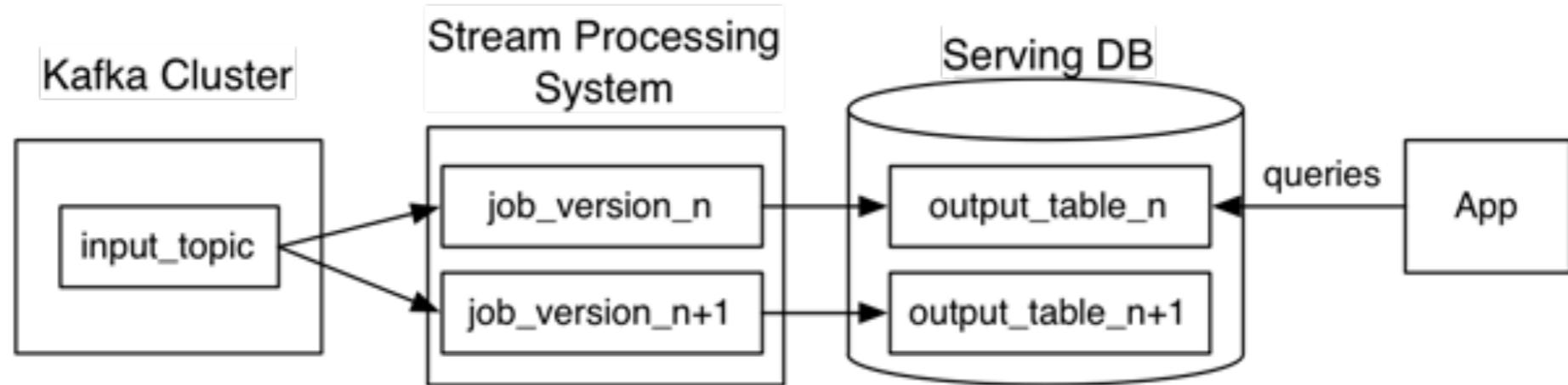
Visual Guide to NoSQL Systems



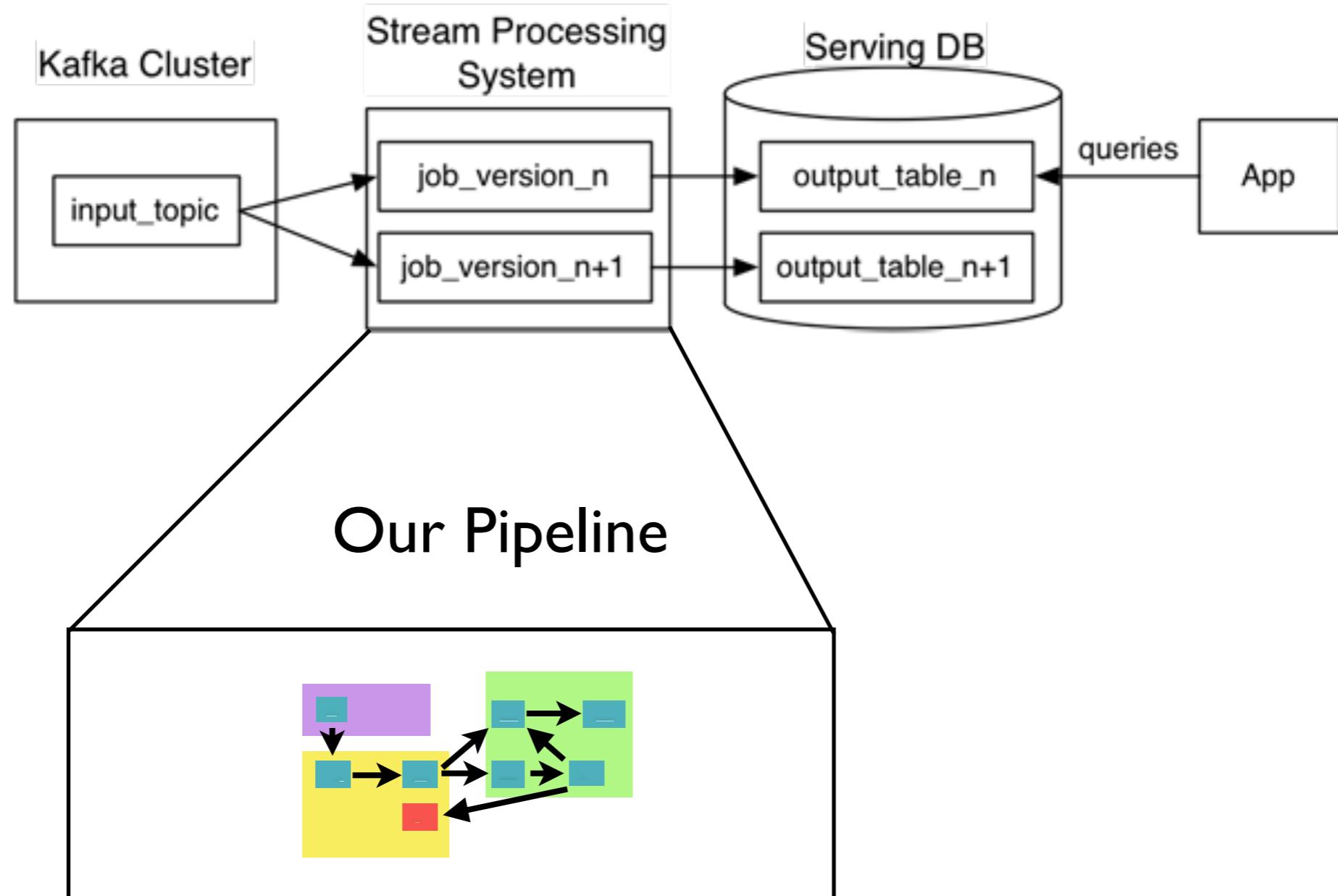
Architecture



Architecture



Architecture





“

So, why the excitement about the Lambda Architecture? I think the reason is because people increasingly need to build complex, **low-latency** processing systems. What they have at their disposal are two things that don't quite solve their problem: a scalable **high-latency batch system** that can process **historical data** and a **low-latency stream processing system** that **can't reprocess** results. By duct taping these two things together, they can actually build a working solution.

- Jay Kreps

METRICS

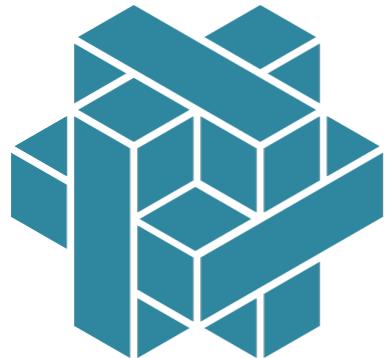


Coda Hale

@coda

github.com/codahale

METRICS EVERYWHERE



Zipfian
Academy



galvanize

Data Science
Immersive

Masters in Data
Science

Data Engineering
Immersive

Weekend
Workshops



We're Hiring!

- Full-time Instructors
- TAs
- Mentor (volunteer)

Questions?



galvanize

Thank You!

Jonathan Dinu
VP of Academic Excellence, Galvanize
jonathan@galvanize.com
[@clearspandex](https://twitter.com/clearspandex)