



Data Engineering 101

Building Data Pipelines





Jonathan Dinu
VP of Academic Excellence, Galvanize
jonathan@galvanize.com
@clearspandex



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

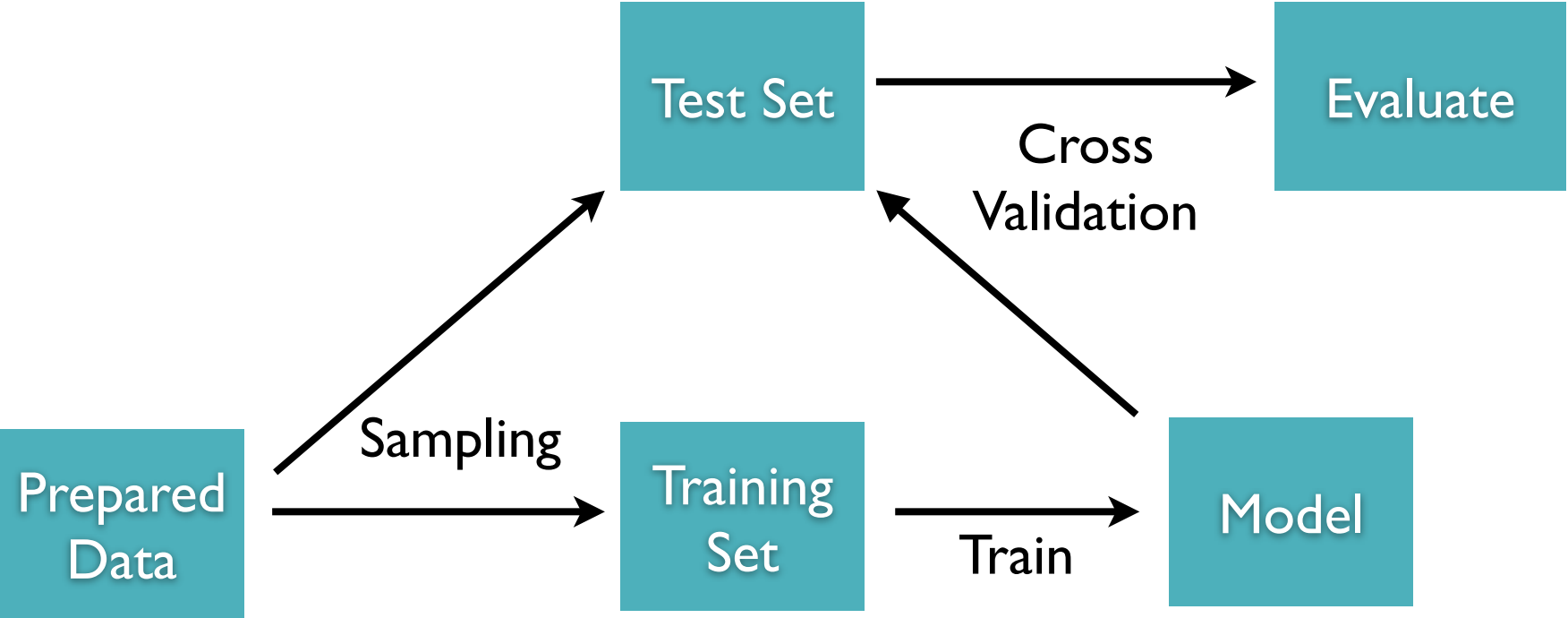
Data ~~Scientist~~ (n.): Person who is better at statistics than any software engineer and better at software engineering than any ~~statistician~~. Data Scientist

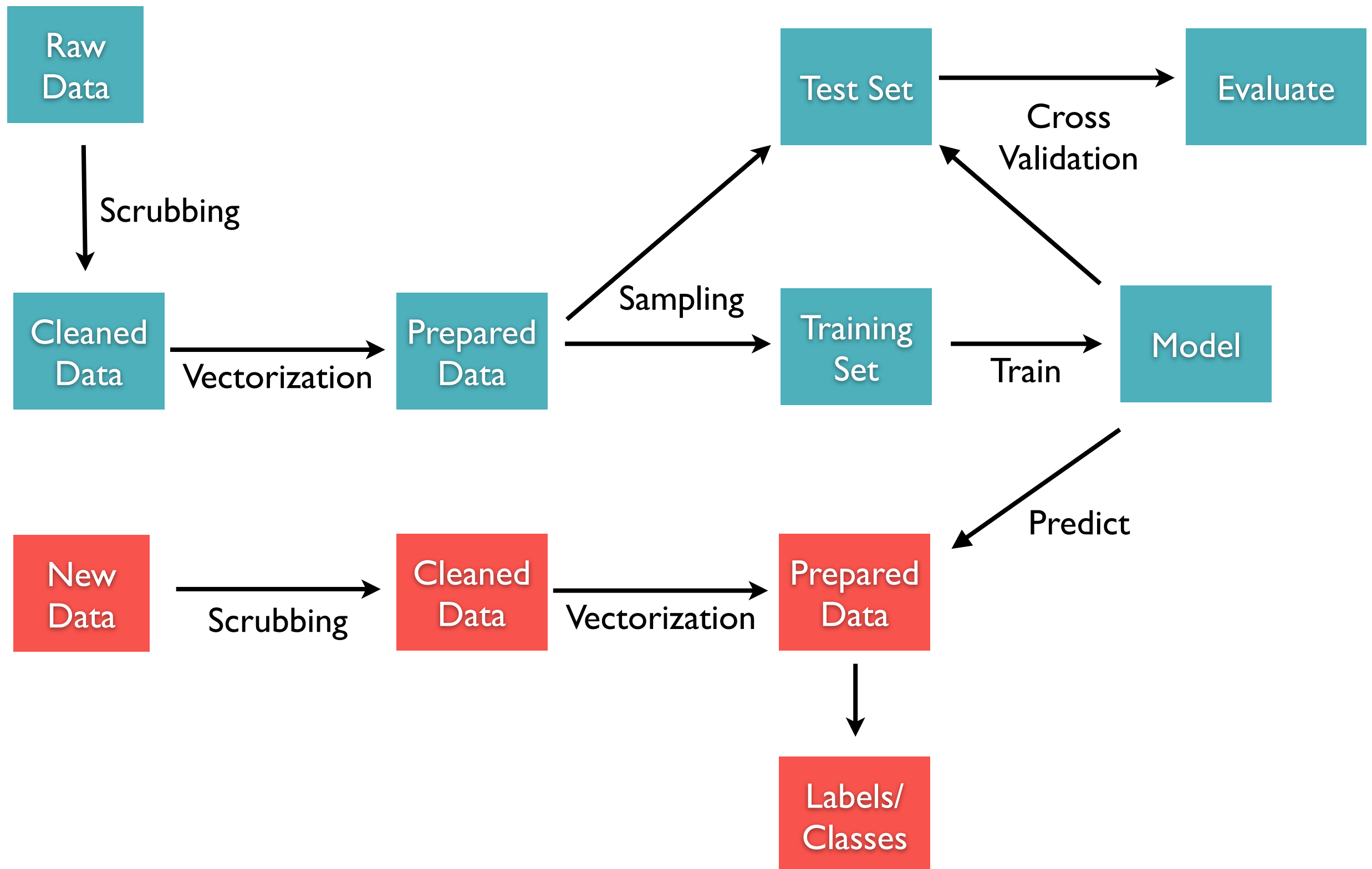
↩ Reply ↻ Retweet ★ Favorite ⋮ More



(In Reality)









The Challenge



At Scale

scrapy

Hadoop Streaming
(w/ BeautifulSoup4)

Snakebite (HDFS)

mrjob or luigi

Spark ML (pySpark)



Flask

Locally

requests

BeautifulSoup4

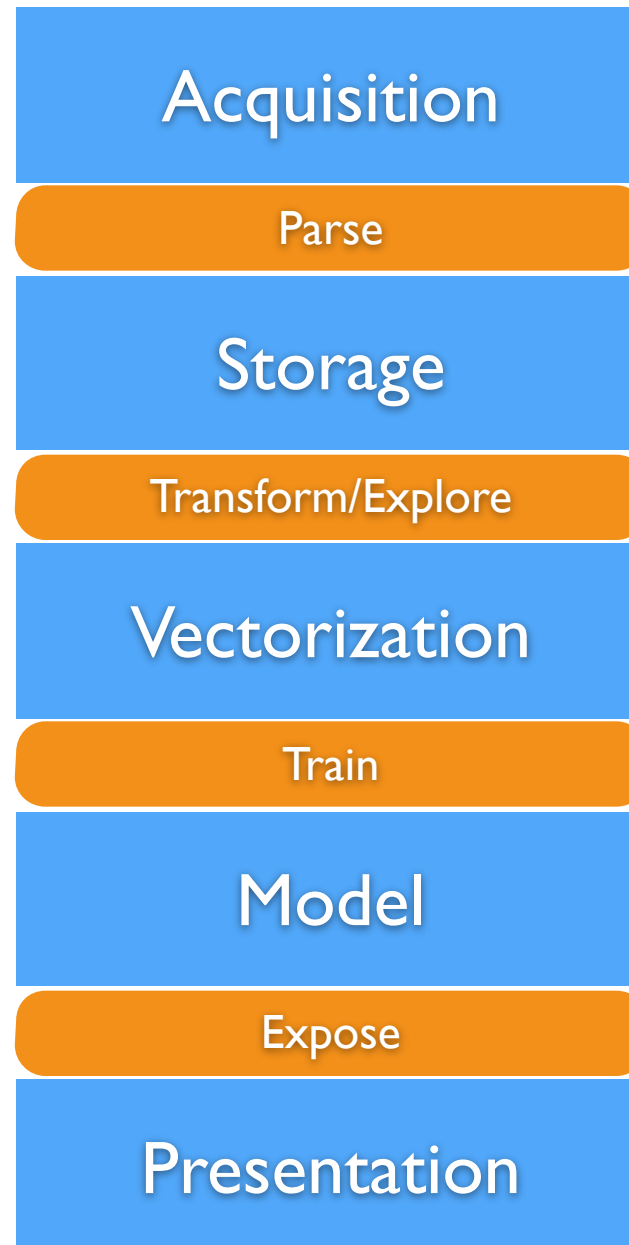
pymongo

pandas

scikit-learn/NLTK



Flask



Wt9t



(it's the pipes!)



Why Pipelines?

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

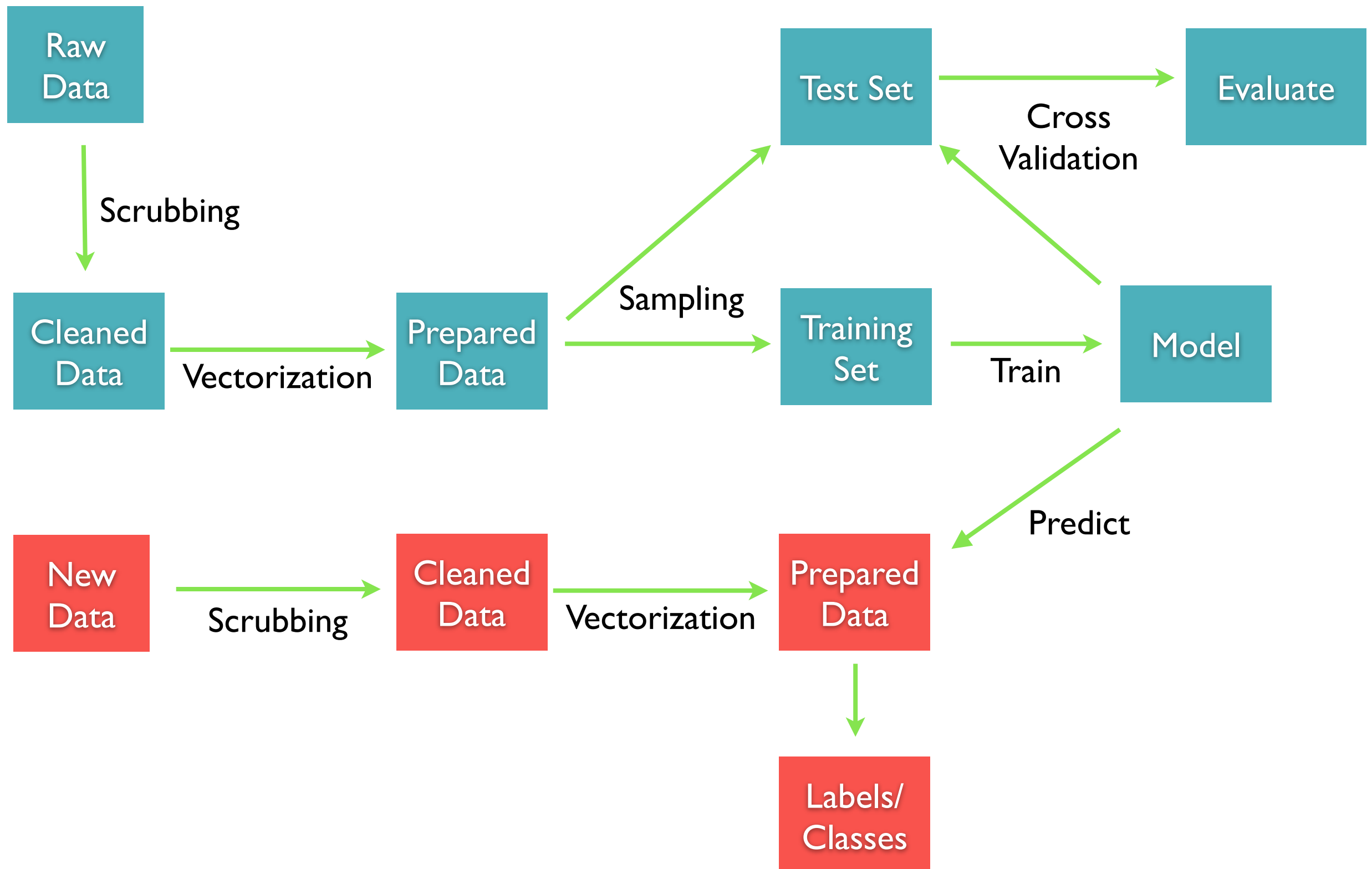
Why

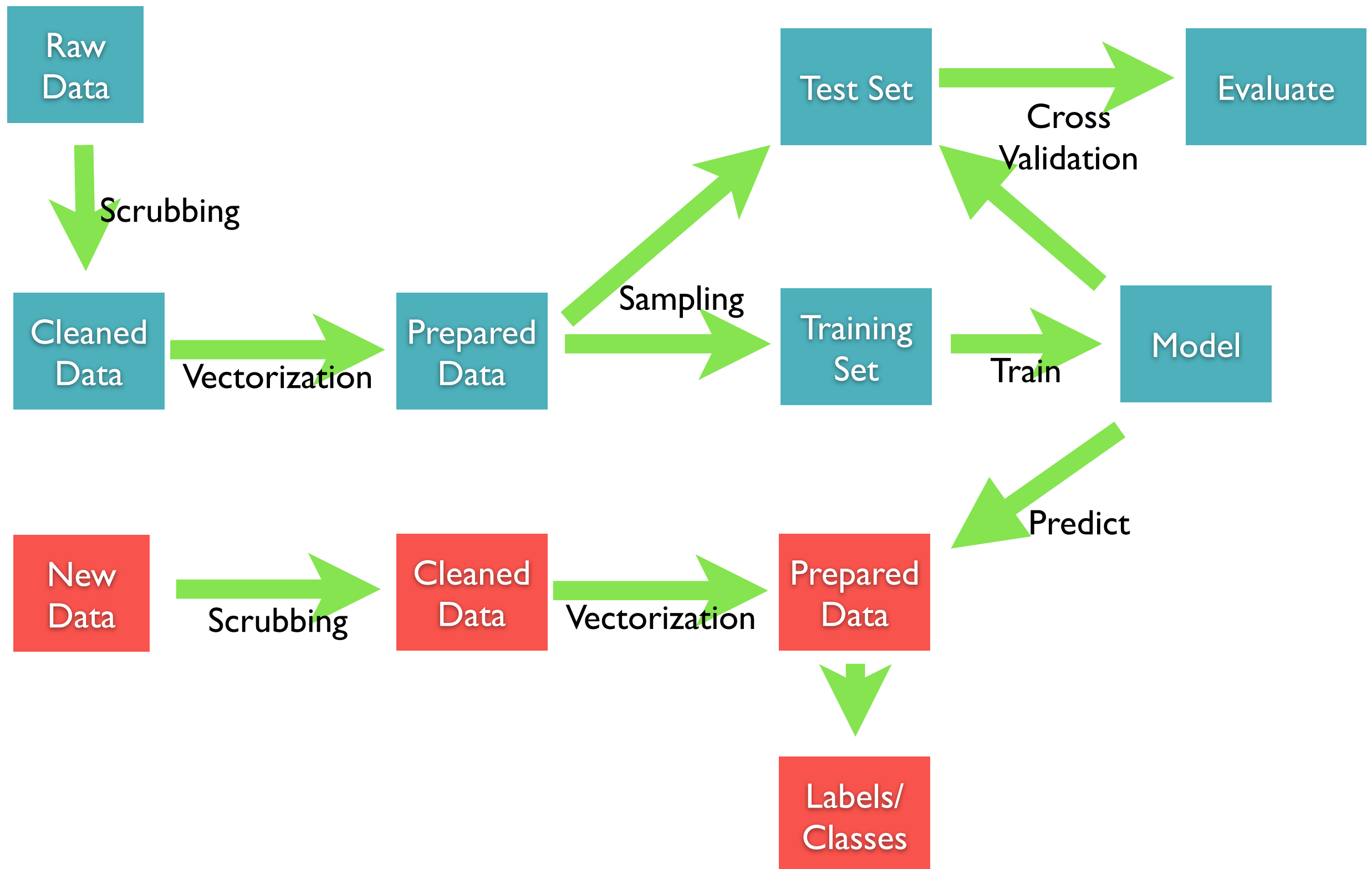
- Idempotence
- Checkpointing
- Native Hadoop Support
- But Works for arbitrary scripts (like make!)

Why

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

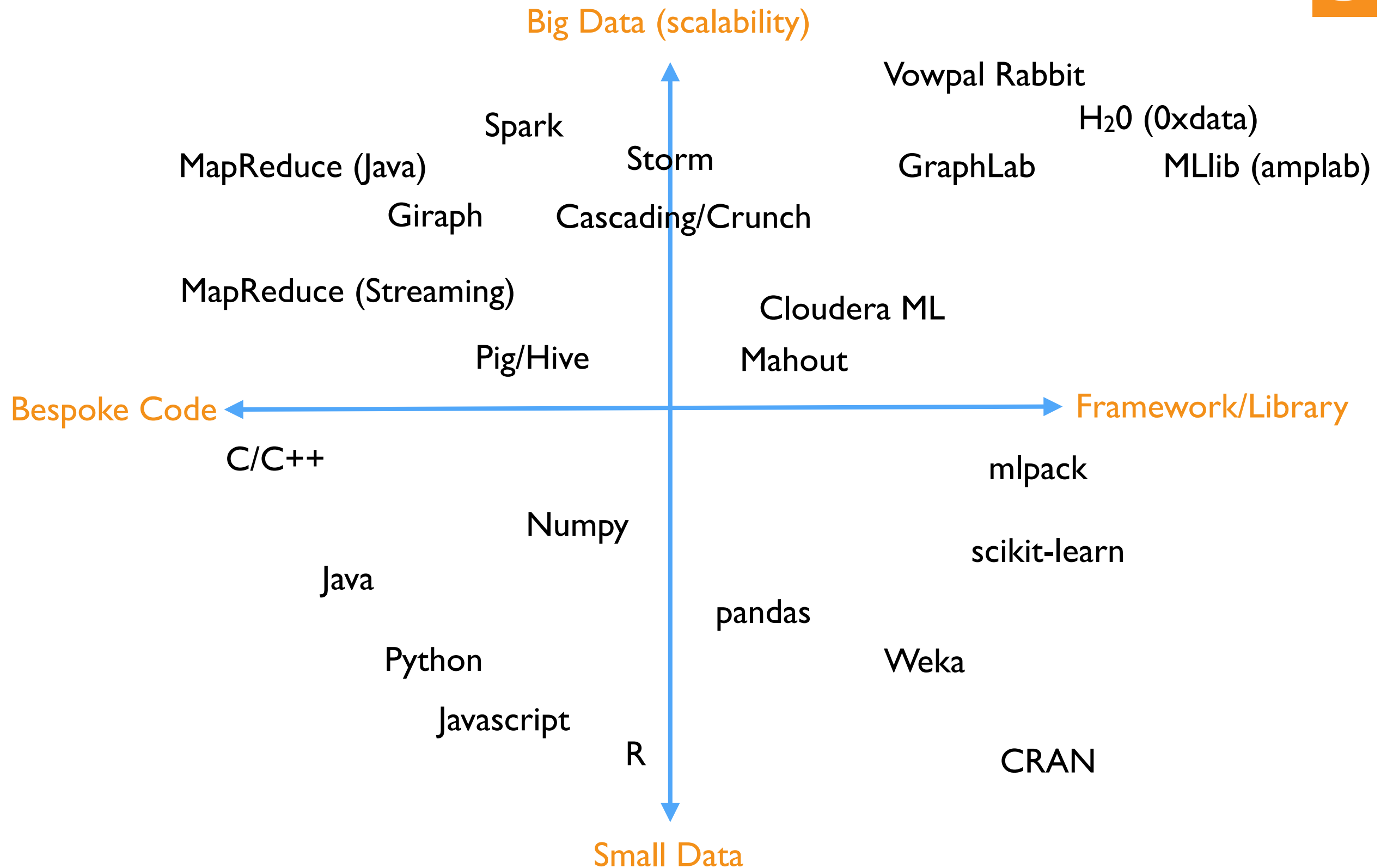
(also has a nice UI and sends emails :)



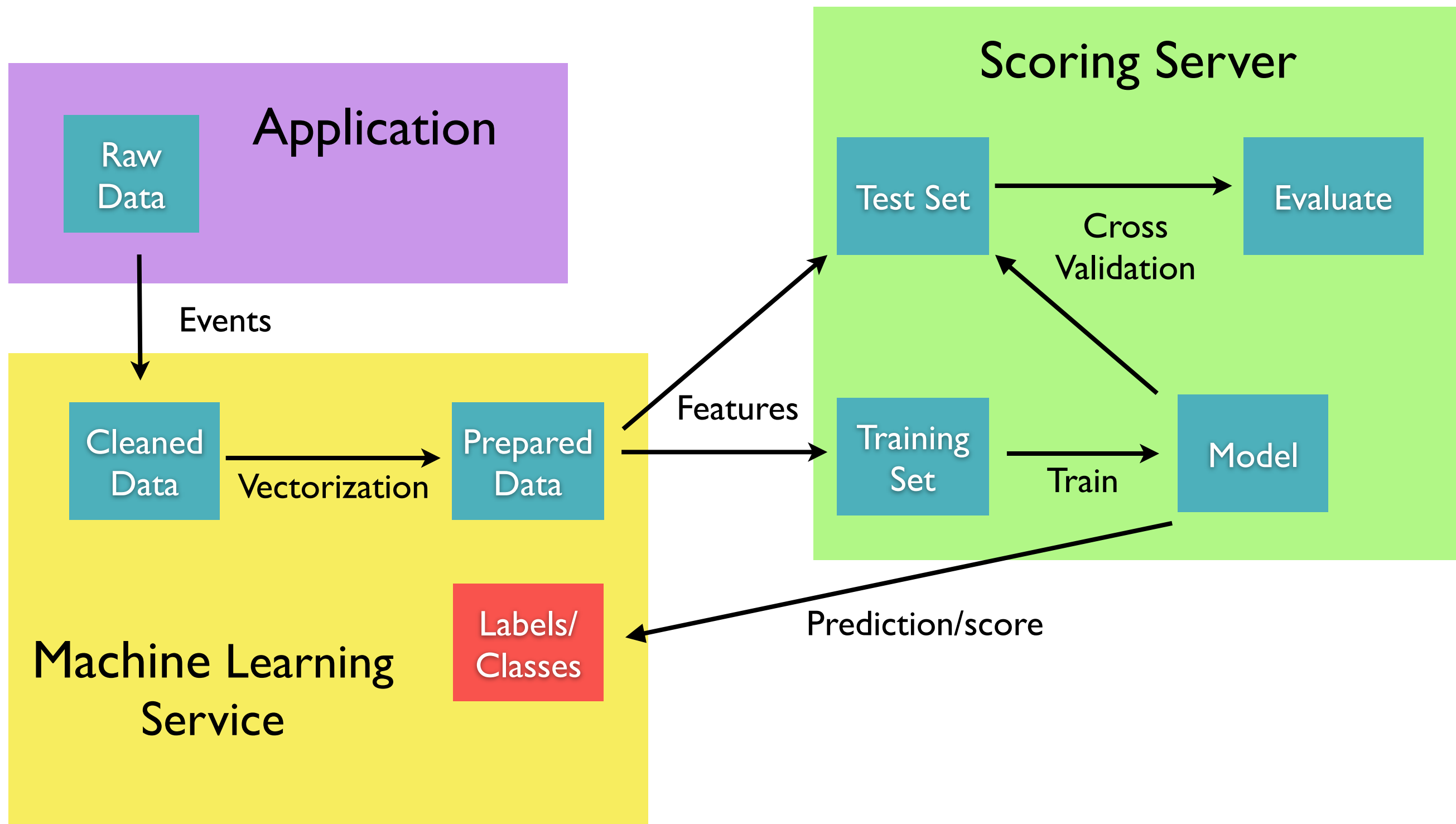




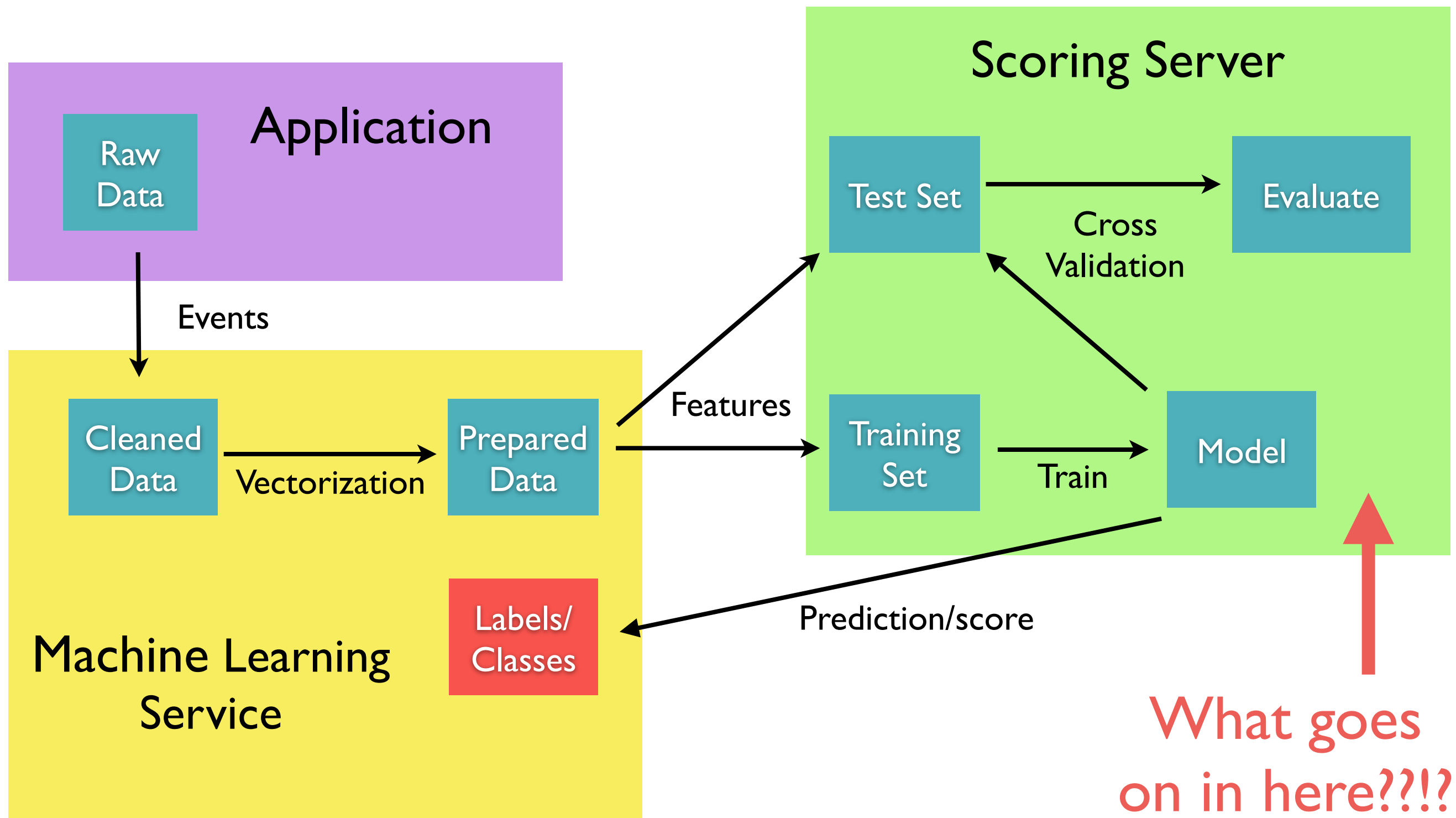
LIVE CODE



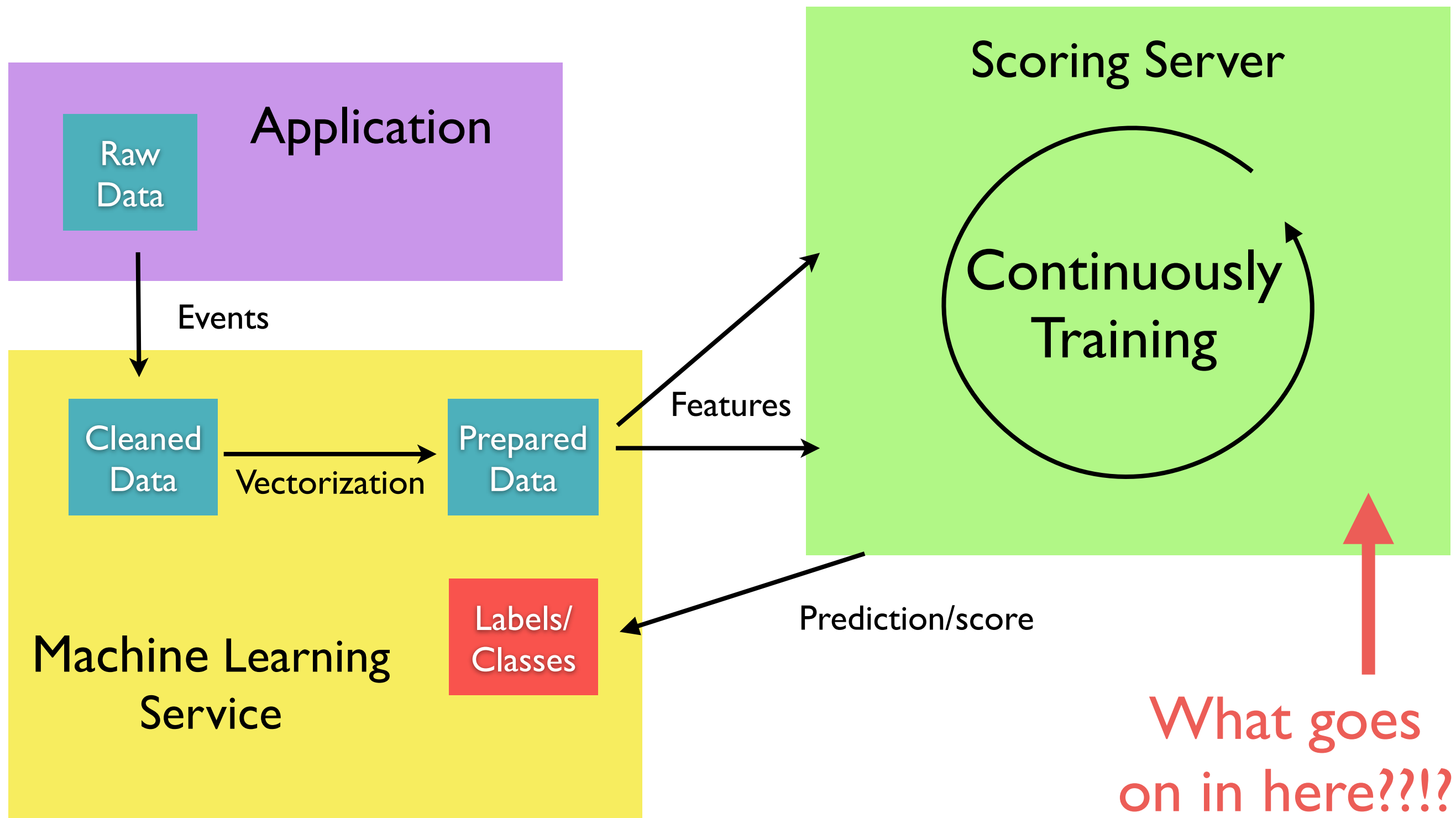
Machine Learning



Machine Learning



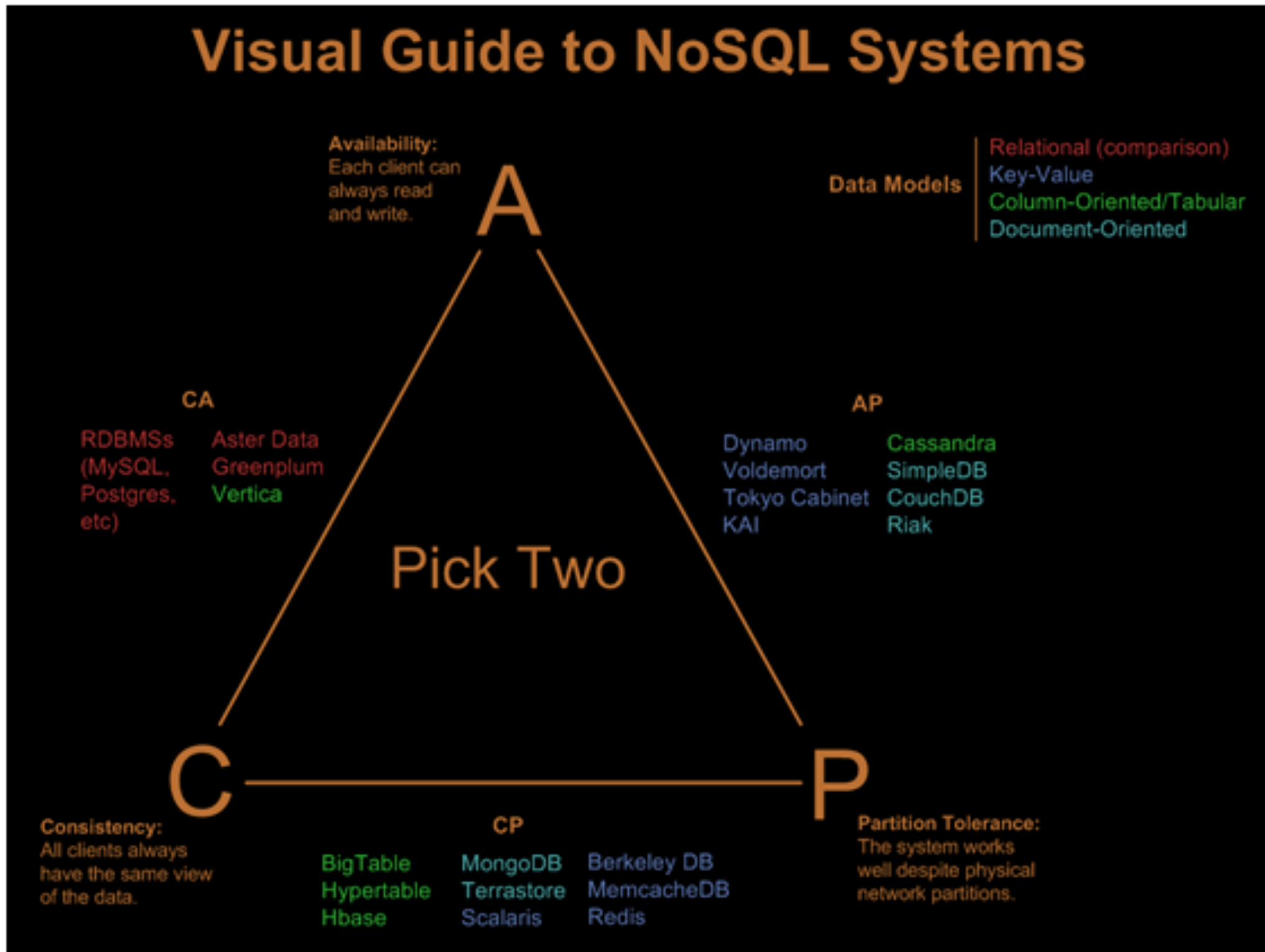
Machine Learning

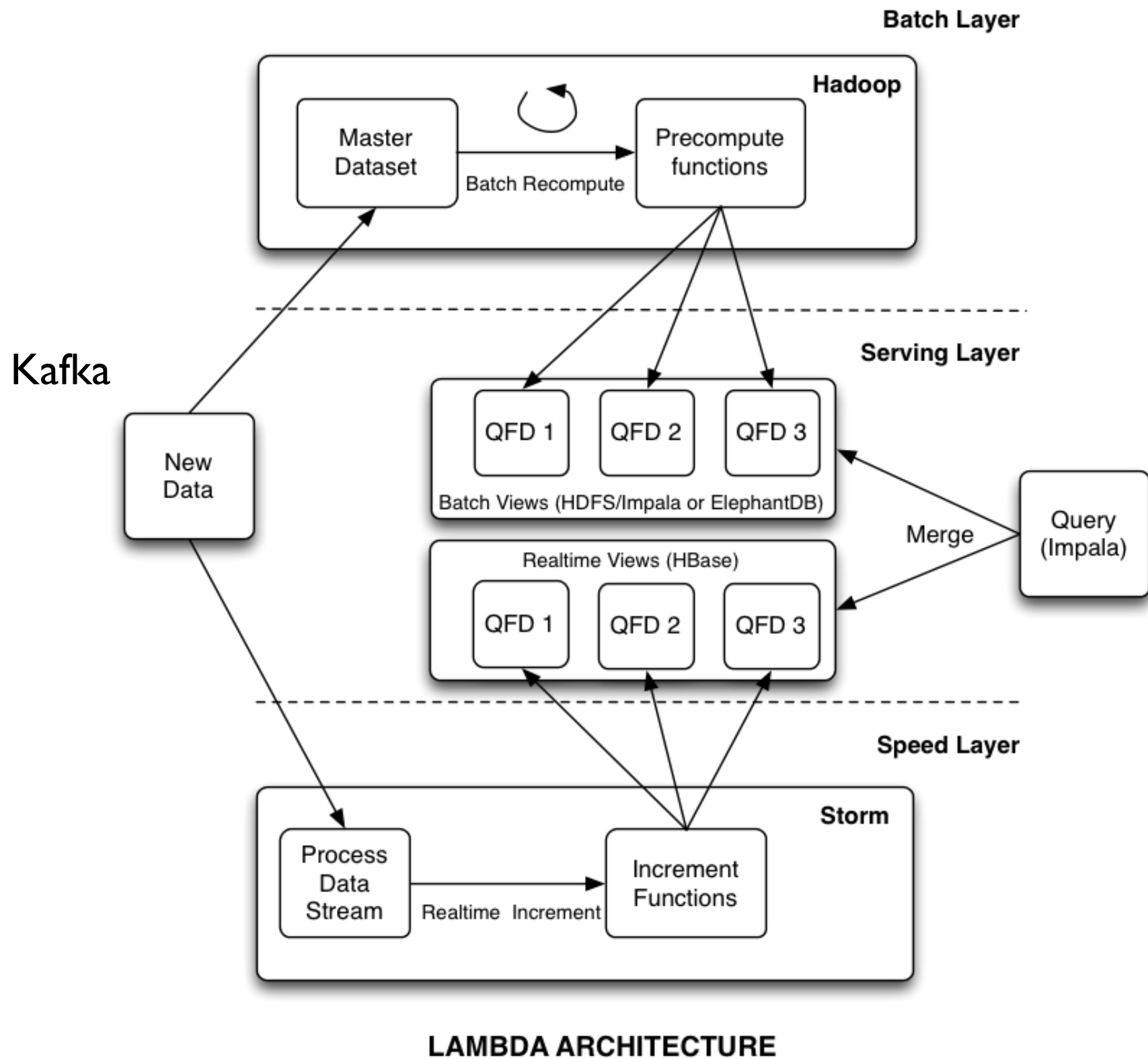


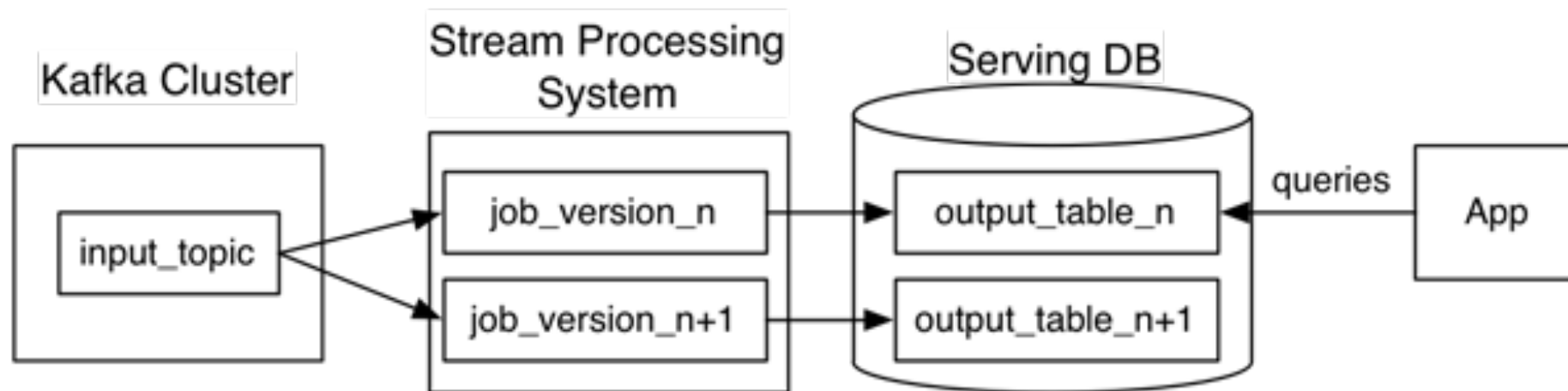


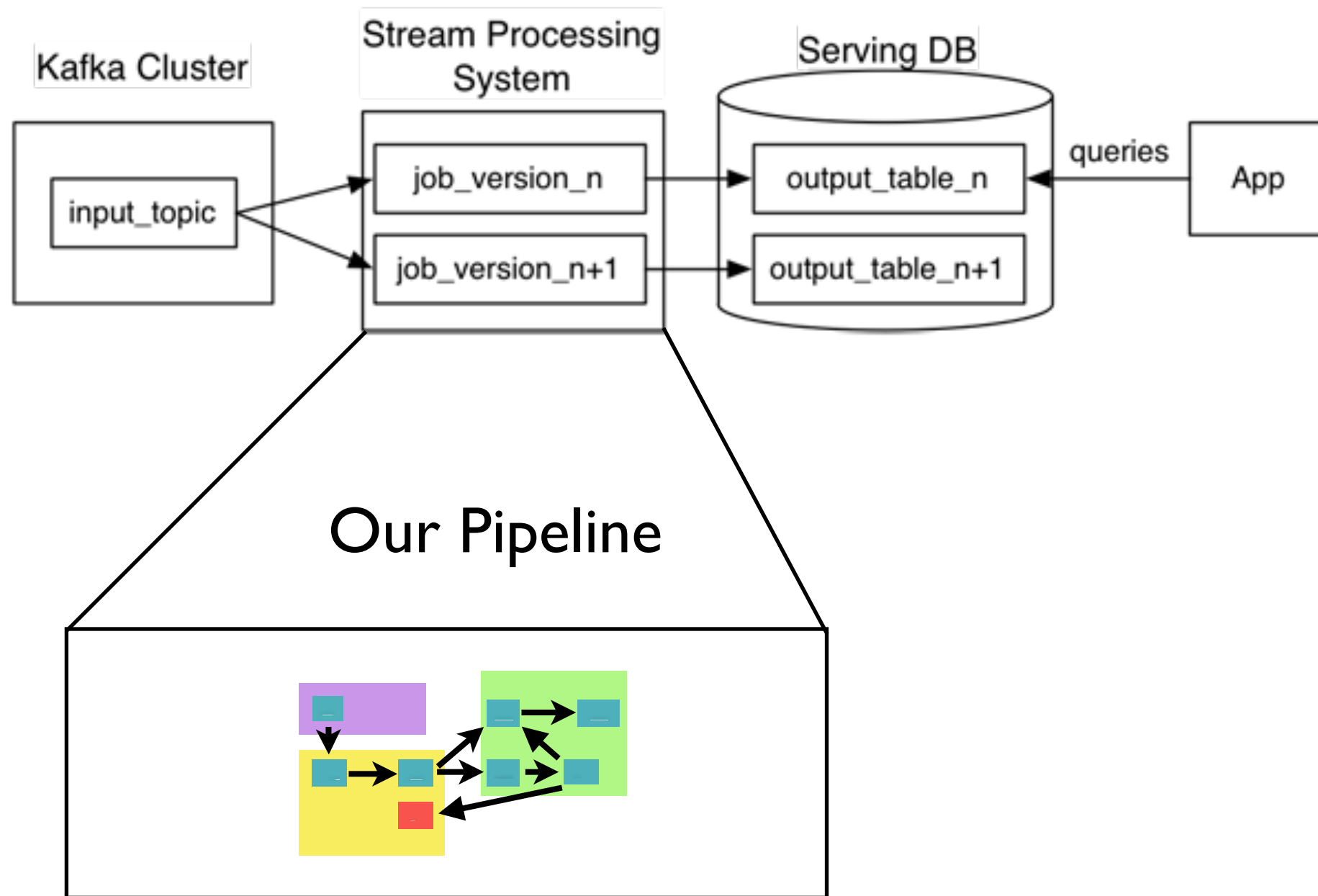
Modern Architectures













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