# NLP with MLlib:
# Global Empire-Building for Fun & Profit

Michelle Casbon

June 7, 2017

Spark Summit

San Francisco

**Qordoba**

# whoami  Michelle Casbon

- Where I work
  - Qordoba, Director of Data Science
- What I've done
- What I love
  - Natural language processing
  - Distributed systems
  - Emoji One

# Data Science Engineer

What my friends think I do
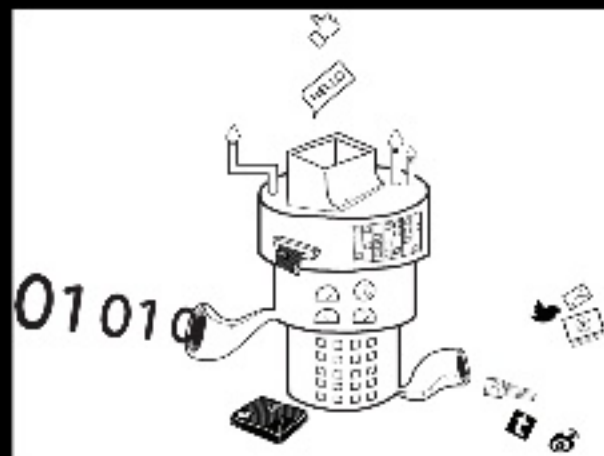
What my parents think I do

What society thinks I do

custom emoji

What my boss thinks I do

What I think I do

What I actually do

# TL;DR

- Importance of localization
    - Why it's hard
    - How it can be made easier
    - Examples
- The role of NLP in solving this problem
    - Brief introduction to NLP
- Spark's role in NLP
- Infrastructure
- Deployment
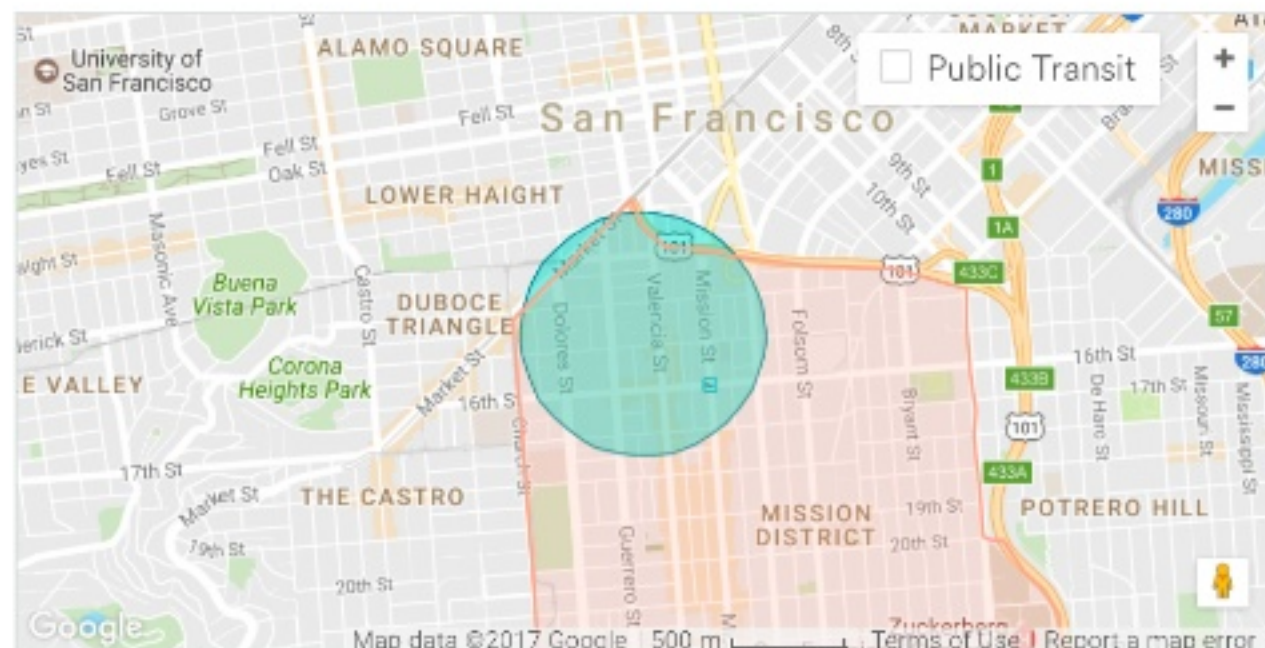
This apartment is in a killer location

Diese Wohnung befindet sich an einem mörderischen Standort

Exact location information is provided after a booking is confirmed.

What is  Qordoba?

# Every product in every language, by default

Fast, scalable way to go from one market to many

Together, our team has lived & worked in 34 countries

# Localization is hard

- People
  - Product managers
  - Marketers
  - Designers
  - Linguists
  - Engineers
- Things
  - Copies of copy
  - String files
  - Emails
  - Pull requests
- Wash, rinse, repeat

App/Git    Product

Design    Email    Vendor

Engineering    Marketing    Documentation

@texasmichelle

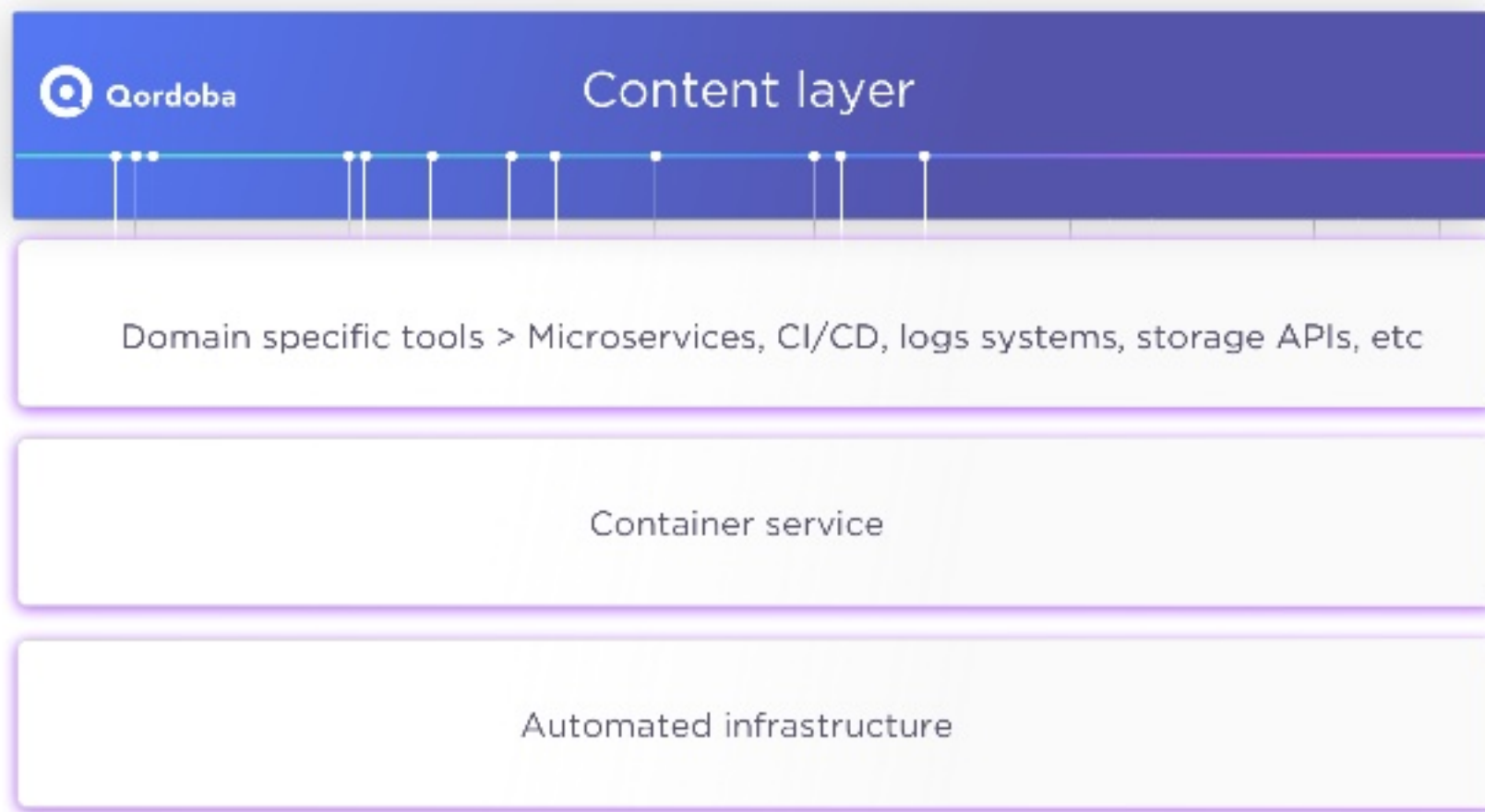# Localization is hard

- People
  - Product managers
  - Marketers
  - Designers
  - Linguists
  - Engineers
- Things
  - Copies of copy
  - String files
  - Emails
  - Pull requests
- Results
  - Got milk?
- Wash, rinse, repeat

App/Git

Product

Design

Email

Vendor

Engineering

Marketing

Documentation

You're doing it wrong

@texasmichelle

# The 🐳 of content

**Qordoba** — Content layer

Domain specific tools > Microservices, CI/CD, logs systems, storage APIs, etc

Container service

Automated infrastructure

@texasmichelle

# Problem

Content is disorganized & deeply embedded

Similar strings are translated over & over

A translation is wrong & needs to be updated everywhere

Adding new content takes time

Translations don't reflect the correct context

Making a change to content takes forever

# Solution

Content stored in one central location, separate from the application logic

Translation Memory & machine learning

The text is corrected once & changes are reflected everywhere, without a deployment

Real-time translation of new strings

Live editor

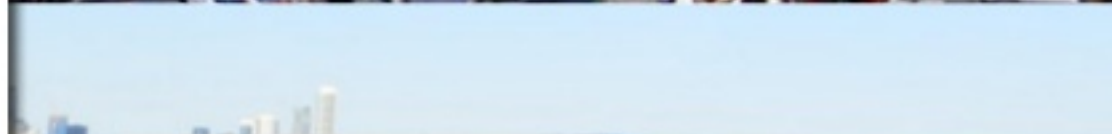A github integration propagates changes automatically

Japantown

すべて購入する

BUY
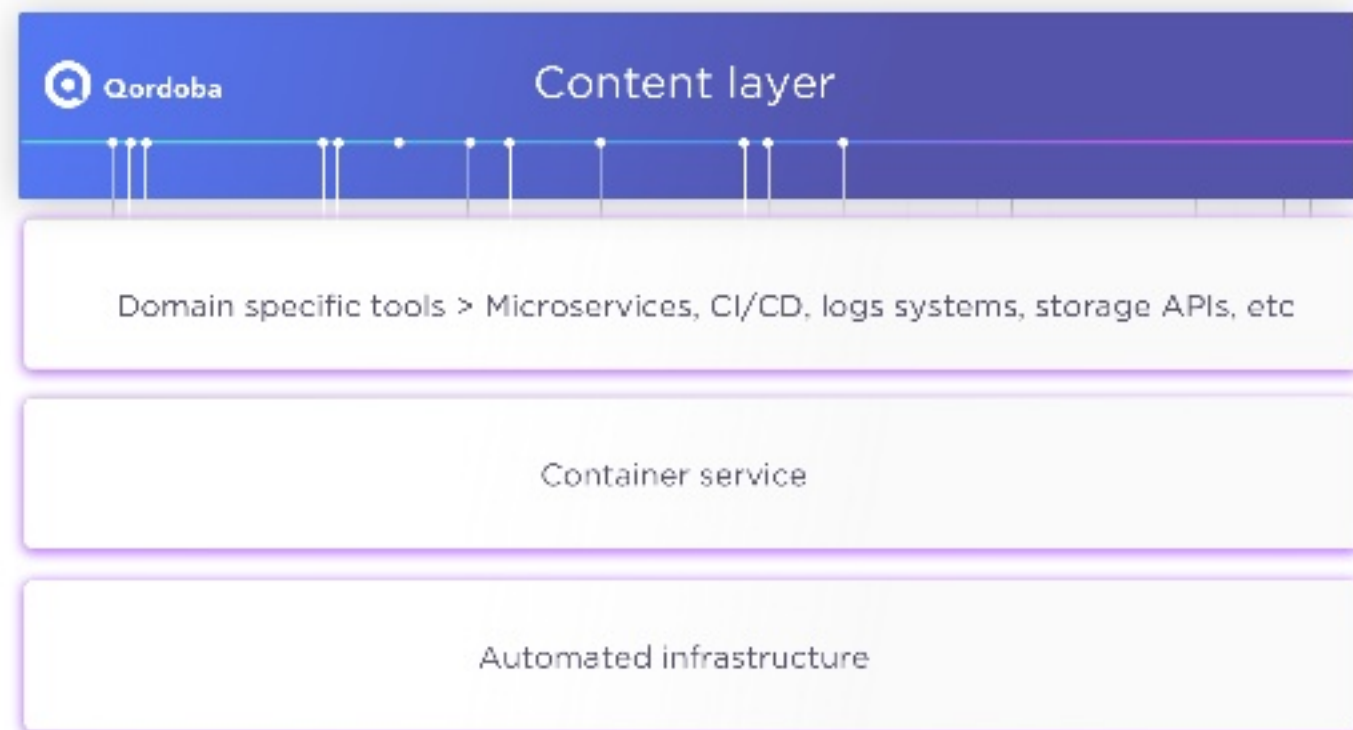ALL THE THINGS

Chinatown

買完全的東西

COMPRAR
TODAS LAS COSAS

Mission Dolores

Image: Christina Chung

# Containers

- Powered by 🐳 🔷 ☸️
- All the way down 🐢
  - Legal 🇫🇷
  - Informal 🇫🇷
  - Somber 🇫🇷

---

**Qordoba** — Content layer

Domain specific tools > Microservices, CI/CD, logs systems, storage APIs, etc

Container service

Automated infrastructure

@texasmichelle

# Affect Detection

- What is it?

- Why do we want it?
  - Hands-off translations 👐
  - Workflow transitions 🔄

- What does this have to do with ⭐ ?

# Affect Detection

- I had dinner with my wife 🍽️     `fear`

- I had dinner with my girlfriend 🍽️     `joy`

- Help Wanted 🐺     `sadness`

- Busca empleo 🐺     `anger`

- This apartment is in a killer location 🏙️     `joy`

- Diese Wohnung befindet sich an einem mörderischen Standort 🏙️     `fear`

@texasmichelle

# Requirements

- REST interface 💨
  - Response time
- Scalability 🇩🇪 🇩🇰 🇩🇴
  - Languages
  - Deployment
  - Models
- Accuracy 💯
- Open source 🪶

PredictionIO

@texasmichelle

# Initial Model

Wow, your first day at the new school! Lisa, have fun. Bart, don't!

Oh boy, sleep! That's where I'm a viking.

When she sees you'll do anything she says, she's bound to respect you!

It takes two to lie. One to lie and one to listen.

Credit: Matt Groening, Fox

Not like this....

1　2　3　4
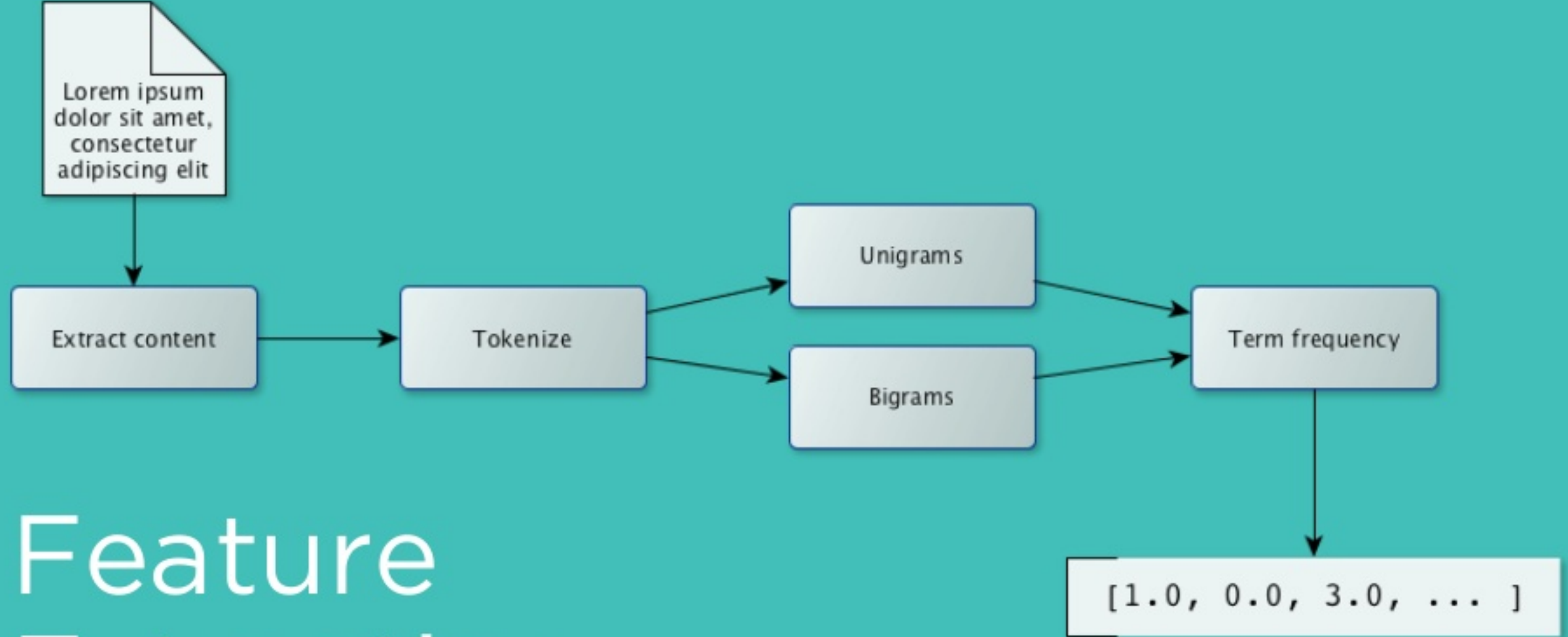
Like this!

1　2　3　4　5

Henrik Kniberg

What does the 🚲 look like?
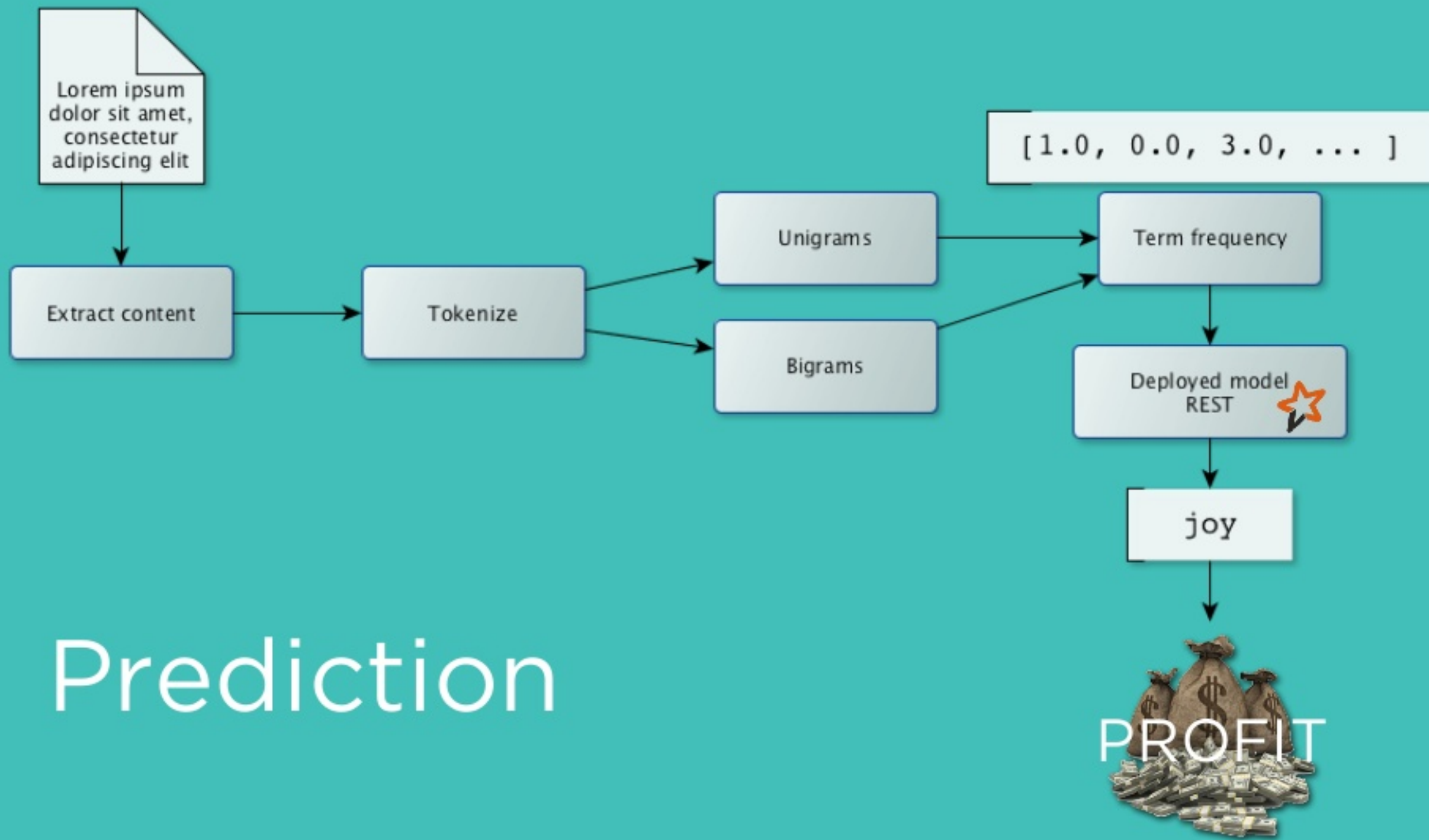
Feature
Extraction

`[1.0, 0.0, 3.0, ... ]`

Add classification

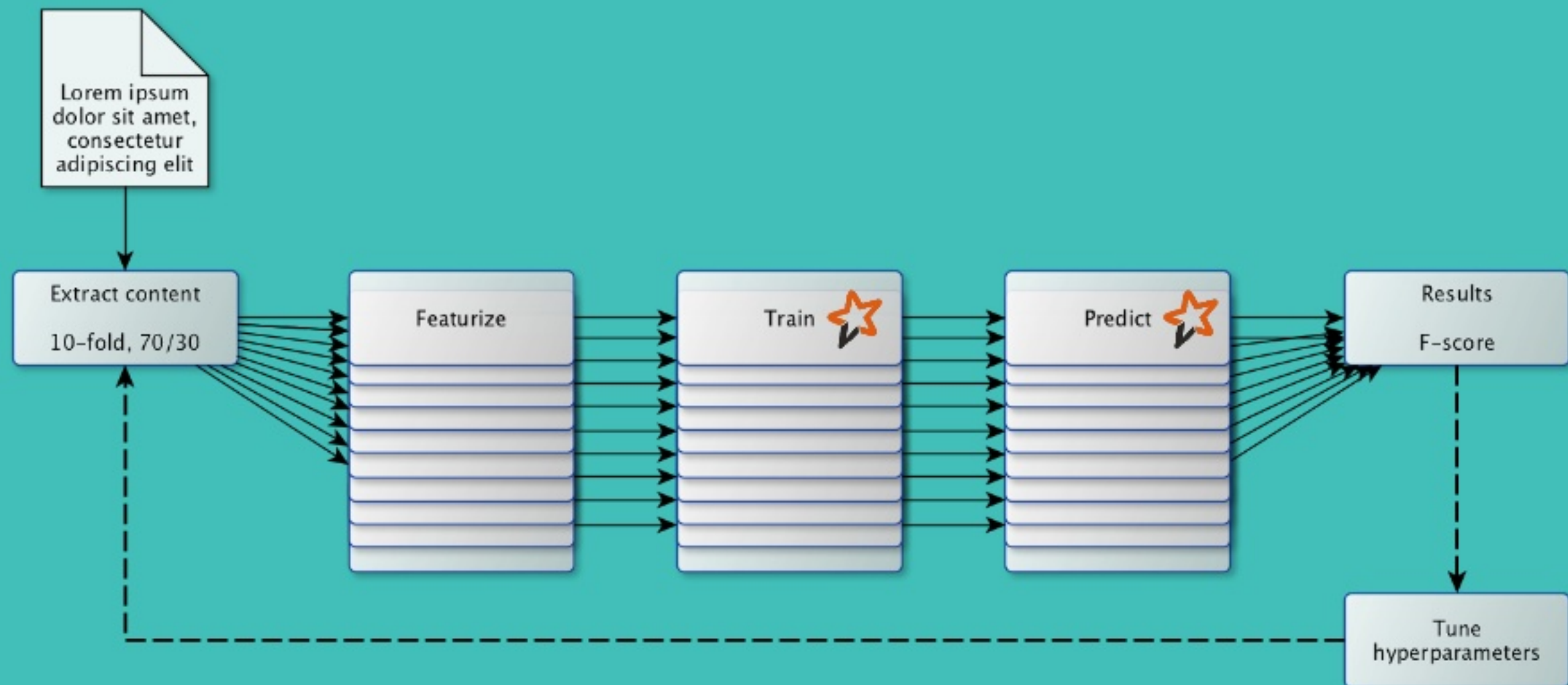LogisticRegression
WithLBFGS

Deploy
LogisticRegressionModel

`[1.0, [1.0, 0.0, 3.0, ... ]]`

Training

Lorem ipsum dolor sit amet, consectetur adipiscing elit

Extract content → Tokenize → Unigrams / Bigrams → Term frequency

`[1.0, 0.0, 3.0, ... ]`

Deployed model REST

joy

Prediction

PROFIT

Cross-validation

# Logging & Reporting

logstash

kibana

elastic

# Microservices

Scala

akka

APACHE Spark

PredictionIO

# Event Bus

Cloud Pub/Sub

kafka

# Storage

MariaDB

cassandra

DATASTAX

# Continuous Delivery

Jenkins

GitHub

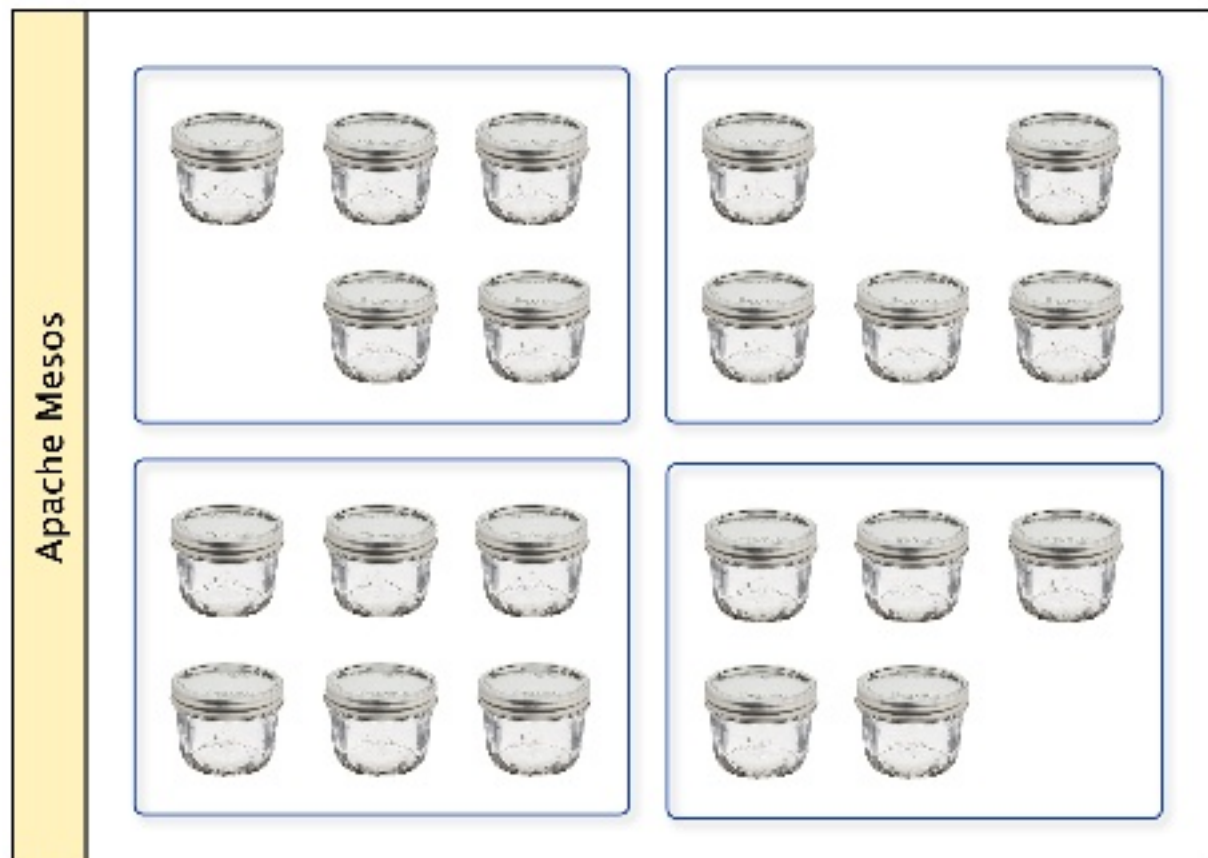JIRA

JFrog Artifactory

docker

MESOS
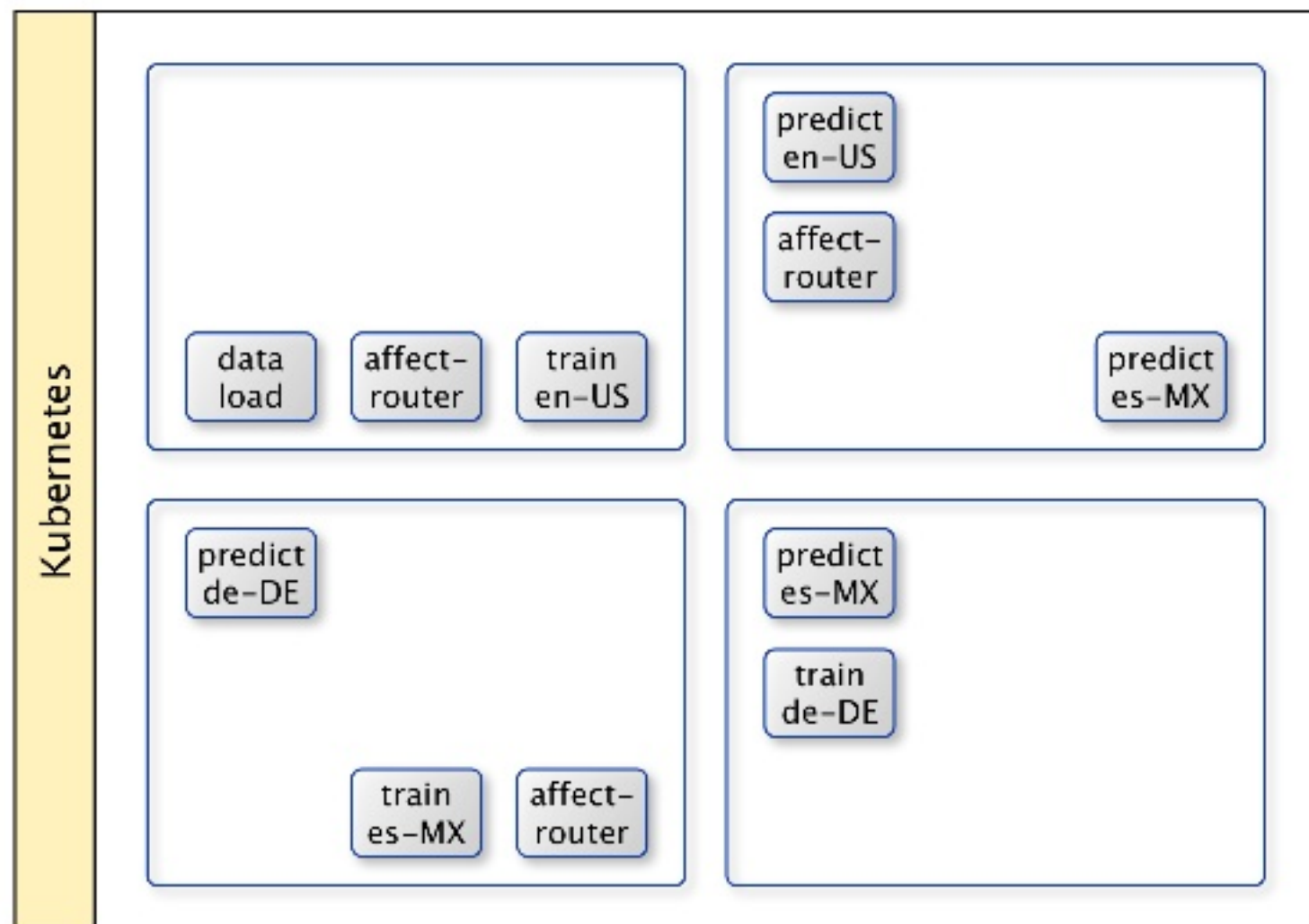
kubernetes

# System Components

# Before

# How to prevent conflict between data science & engineering
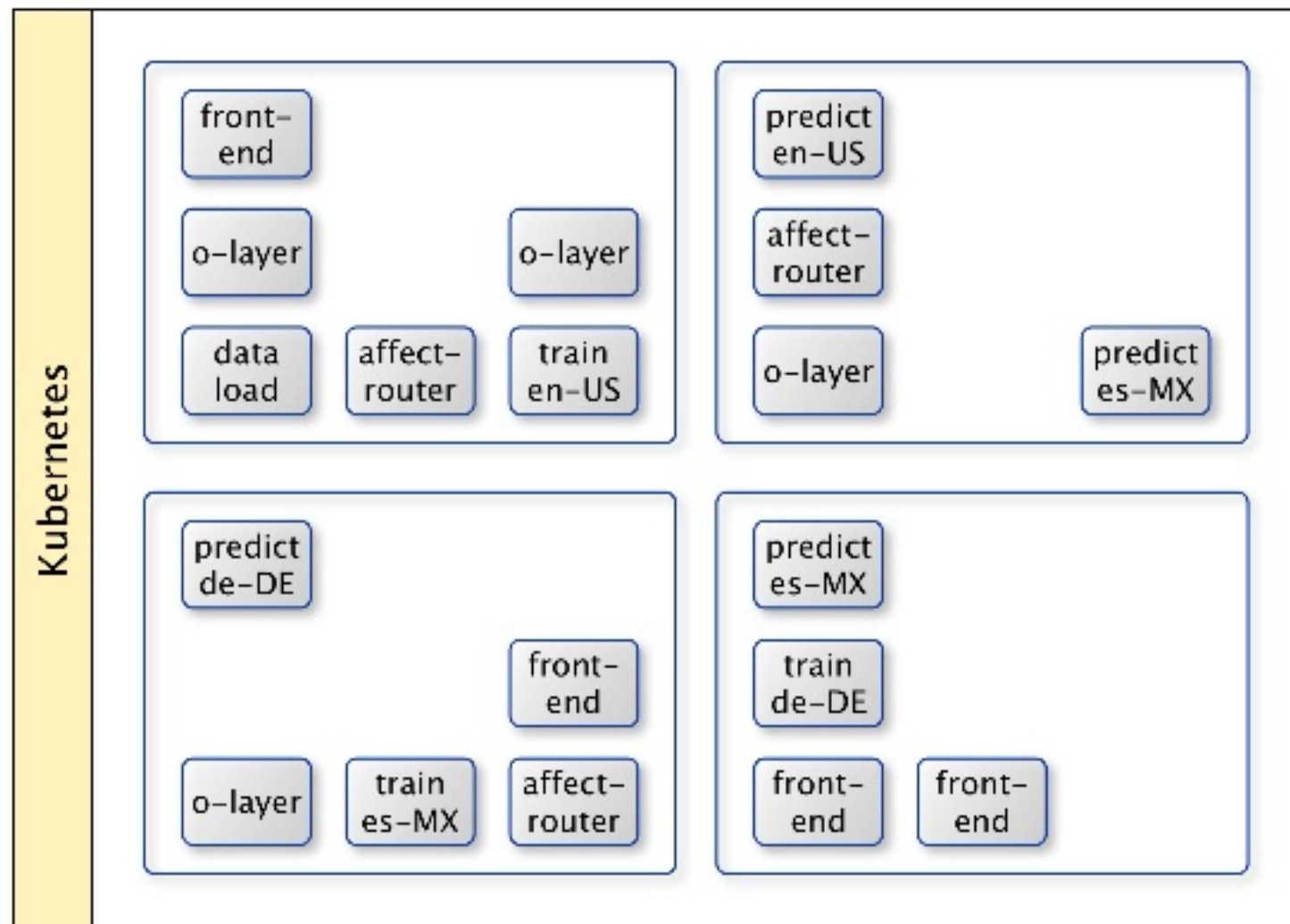


@texasmichelle

# Affect detection deployment
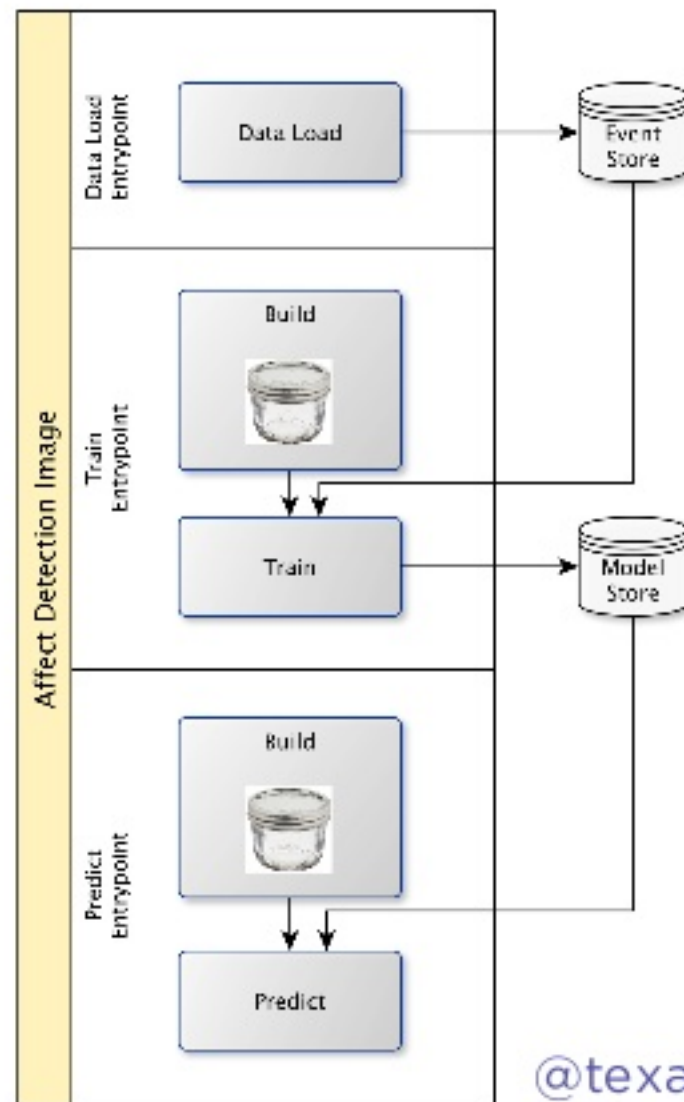
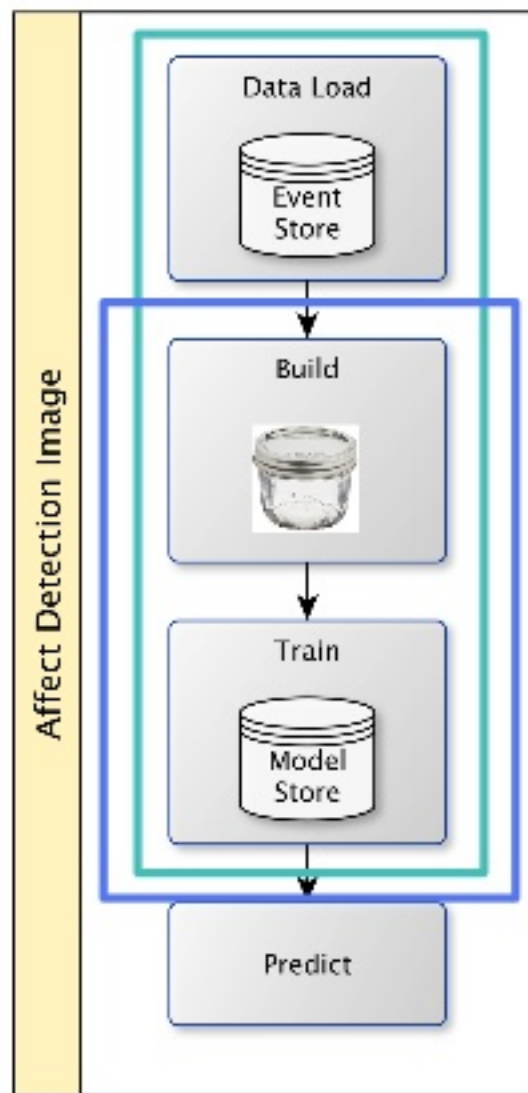

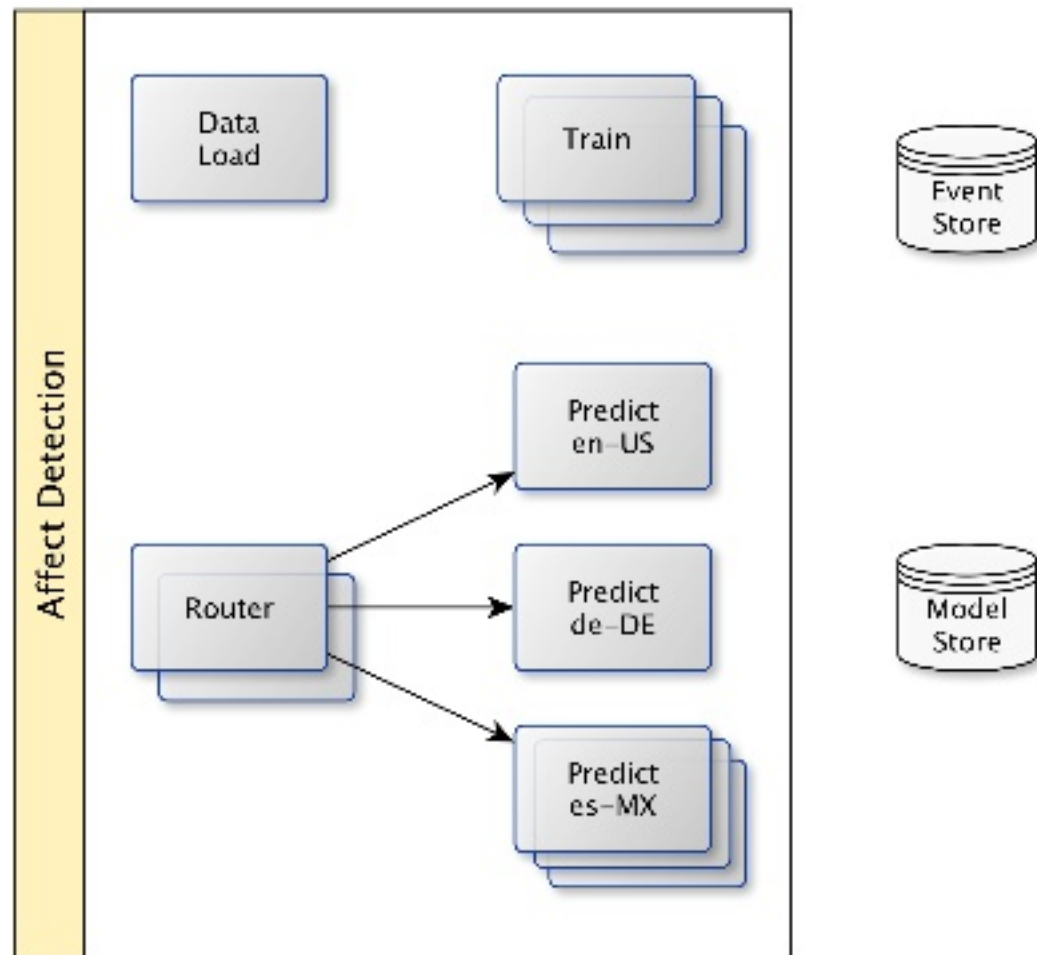@texasmichelle

# After



@texasmichelle

# Single Model



@texasmichelle

# Many Models

# Affect Detection Image

Affect Detection

Mysql Client

Apache PredictionIO (incubating) 0.10.0

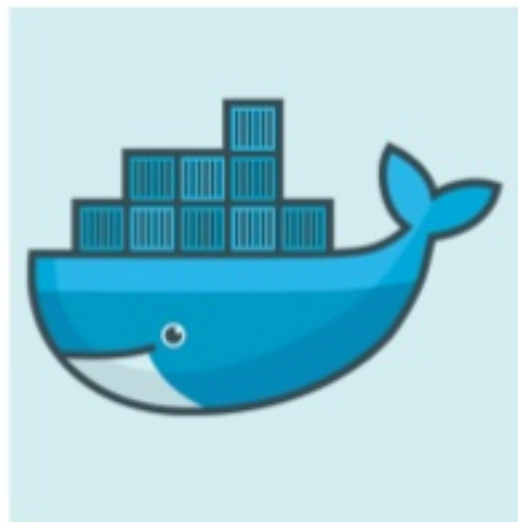Apache Spark 1.6.3 Native BLAS

Java 8

Ubuntu 16.04

Affect Detection

Mysql Client

Apache PredictionIO (incubating) 0.10.0

Apache Spark 1.6.3 Native BLAS

Java 8

Ubuntu 16.04

# Router Image

affect-router

qordoba-builder

affect-router

qordoba-builder

# Qordoba Builder Image

sbt

git

Google Cloud SDK

Docker 17.03

Java 8

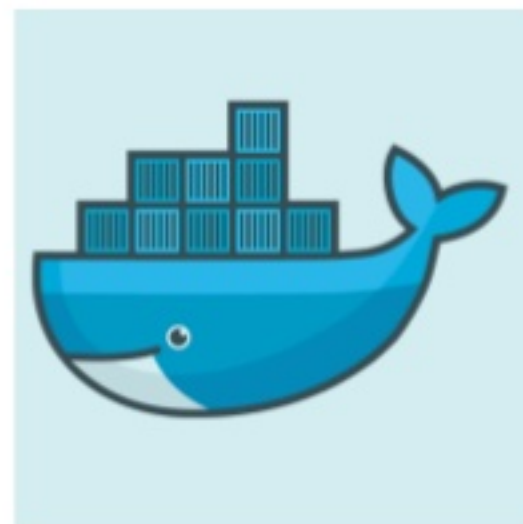Ubuntu 16.04

sbt

git

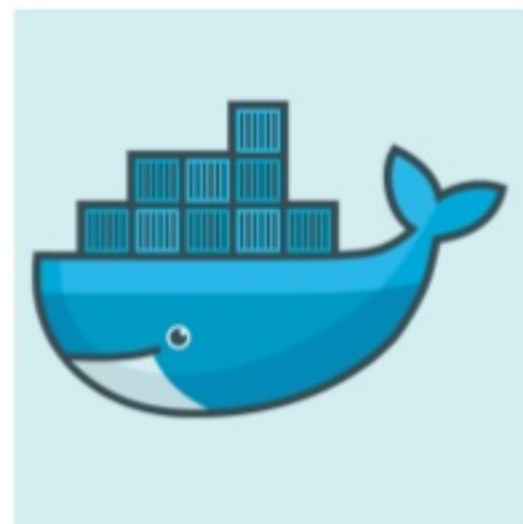Google Cloud SDK

Docker 17.03

Java 8

Ubuntu 16.04

# Qordoba Builder Image

# Qordoba Builder Image

# Affect Detection

- Featurization improvements
- Multi-label instead of mutual exclusivity
- Confidence scores
- Additional algorithms 🌲🌲🌲
- Expanded training set
- More evaluation metrics
- ~~Classification~~ Graph?
- Unsupervised models

OH, LOGISTIC REGRESSION ISN'T ENOUGH FOR YOU?

TELL ME MORE ABOUT RANDOM FOREST ENSEMBLES
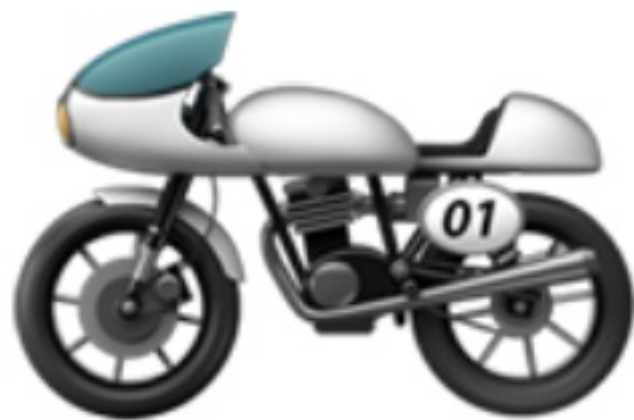
memegenerator.net

@texasmichelle

# Affect Detection

- Lighter images
- Lighter deployment
- Weaveworks

# TL;DR

- Importance of localization
  - Why it's hard
  - How it can be made easier
  - Examples
- The role of NLP in solving this problem
  - Brief introduction to NLP
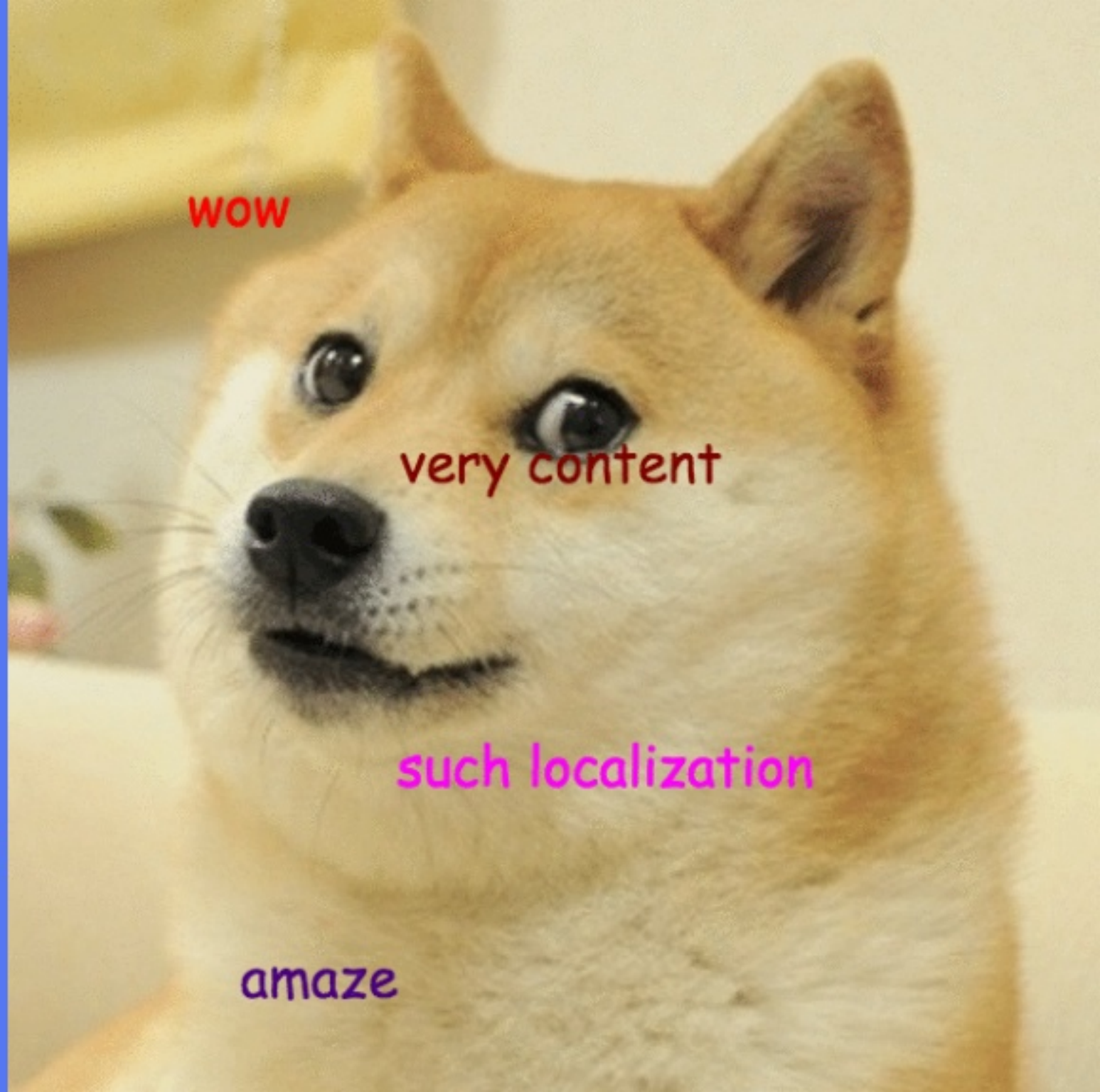- Spark's role in NLP
- Infrastructure
- Deployment

@texasmichelle

# Challenge

Don't just build products that scale. Build products that scale content.