# Vincent Lesierse

**Technical Product Manager**

- **Vamp.io**

- **@vlesierse**

**Jörg Schad**
ML Platform Engineer

@joerg_schad

joerg@suki.ai

www.suki.ai

- **Suki.ai**
- **Mesosphere**
- **PhD Distributed DB Systems**

- **@joerg_schad**



| Monitoring & Operations | DATADOG | TensorBoard | | |
|---|---|---|---|---|
| Data & Streaming | Model Engineering | Model Training | Model Management | Model Serving |
| Distributed Data Storage and Streaming | Data Preparation and Analysis | Distributed Training using Machine Learning Frameworks | Storage of trained Models and Metadata | Use trained Model for Inference |

| Feature Catalogue | Jenkins | Continuous Integration |
|---|---|---|
| binder | Notebook Library | |
| TensorFlow Hub | Model Library | |

| DC/OS | Resource and Service Management | kubernetes | MESOS |
|---|---|---|---|

**KubeCon  CloudNativeCon**
China 2018

**Operating Deep Learning Pipelines Anywhere Using Kubeflow**
Jörg Schad & Gilbert Song, Mesosphere

**WEBCAST**

**AVOIDING PITFALLS WHEN CONTAINERIZING JAVA AND JVM APPLICATIONS**

O'REILLY

# A tale of two worlds

It was the best of times -where Microservices allowed broke our code in manageable pieces,
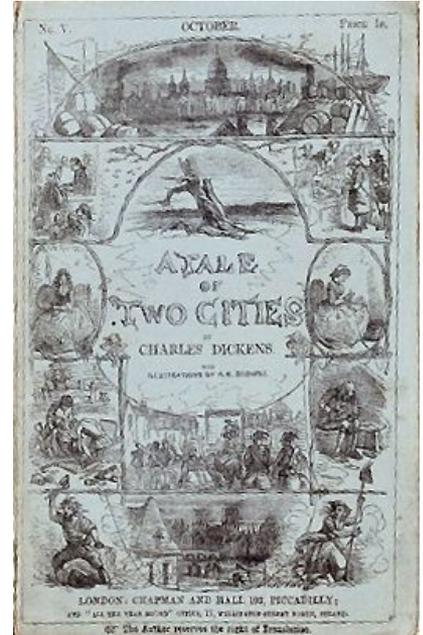
it was the worst of times -where we actually had to deploy all those pieces,

it was the age of wisdom - where Machine Learning allowed for cool Models,

it was the age of foolishness - where we suffered to productionize these models,

it was the epoch of belief, it was the epoch of incredulity,

it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way—in short, the period was so far like the present period, that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.[2]



4

# A tale of two worlds



Deployment of Microservices



Deployment of Machine Learning Models

# Microservice Deployments

Challenges

Blue/green….

Canary Releases

A/B Testing

Solutions

Service Meshes (Istio)

Vamp

# Vamp's history

**Our mission: Provide an easy-to-use and powerful solution to transform software from ideas into value as efficiently as possible, using cloud-native technologies.**

Important: Releasing is (much) more than technical deployment. Continuous validation requires observability and "golden metrics". Processes are crucial. Tools and technologies are means to an end.

2014: canary testing & releasing and SLA-based auto-scaling on top of Mesos+Marathon, using HAProxy for networking (service-discovery/load-balancing/mesh/ingress)

2016: added Kubernetes support

2019: added Istio support

# Vamp



Core — Open Source
Enterprise
Cloud

DC/OS  Kubernetes  HAProxy  Istio

Amazon Web Services  Microsoft Azure  Google Cloud

# Machine Learning….

# What you want to be doing

```
[ Get Data ] → [ Write intelligent machine learning code ] → [ Train Model ] → [ Run Model ]
```

Repeat

# What you're actually doing



*Sculley, D., Holt, G., Golovin, D. et al. Hidden Technical Debt in Machine Learning Systems*

# Challenge: Persona(s)

# The Rise of the *DataOps Engineer*

Combines two key skills:

- Data science
- Distributed systems engineering

The equivalent of *DevOps* for *Data Science*

- **Build** automation software to run machine learning systems
- **Operate** systems so they're available, scalable, and performant
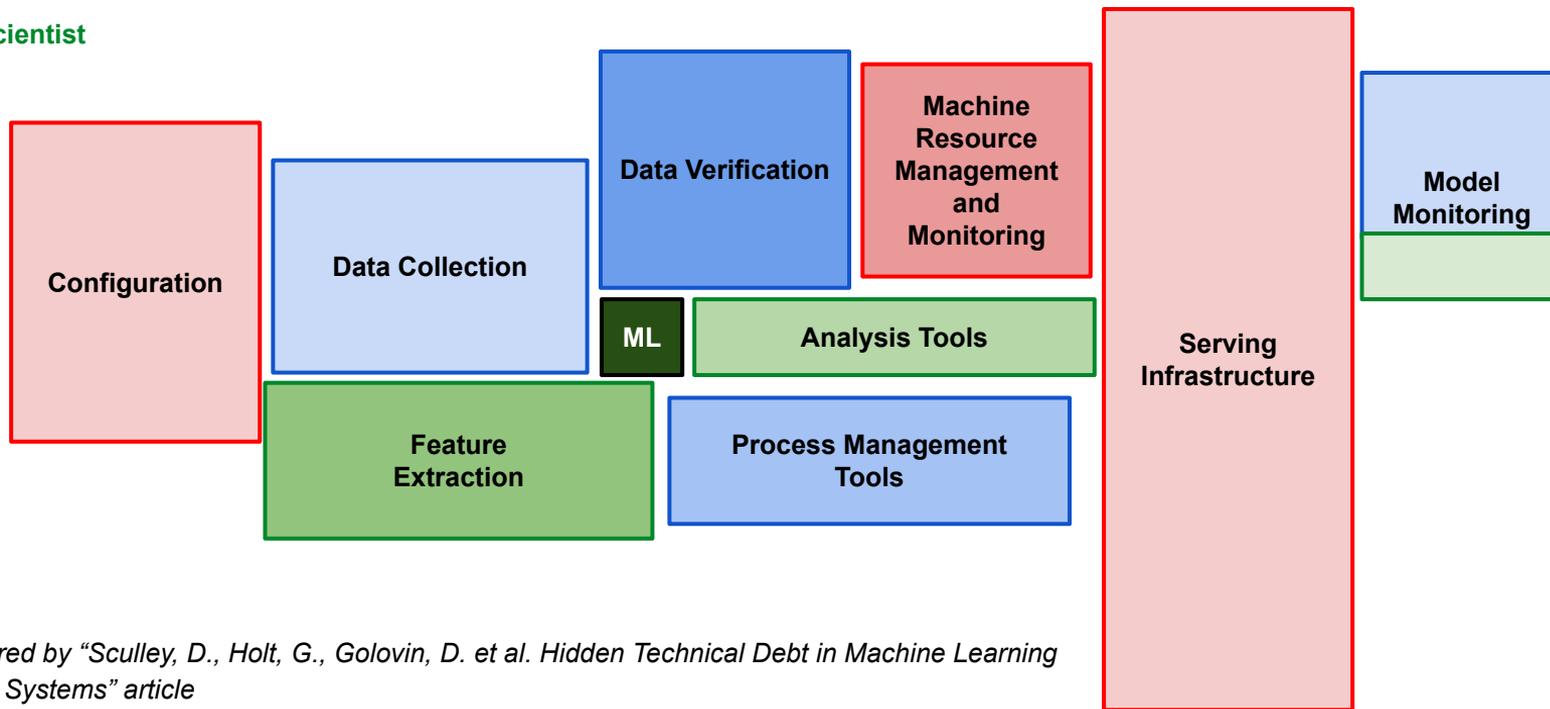- **Evangelize** tools and best practices among data scientists
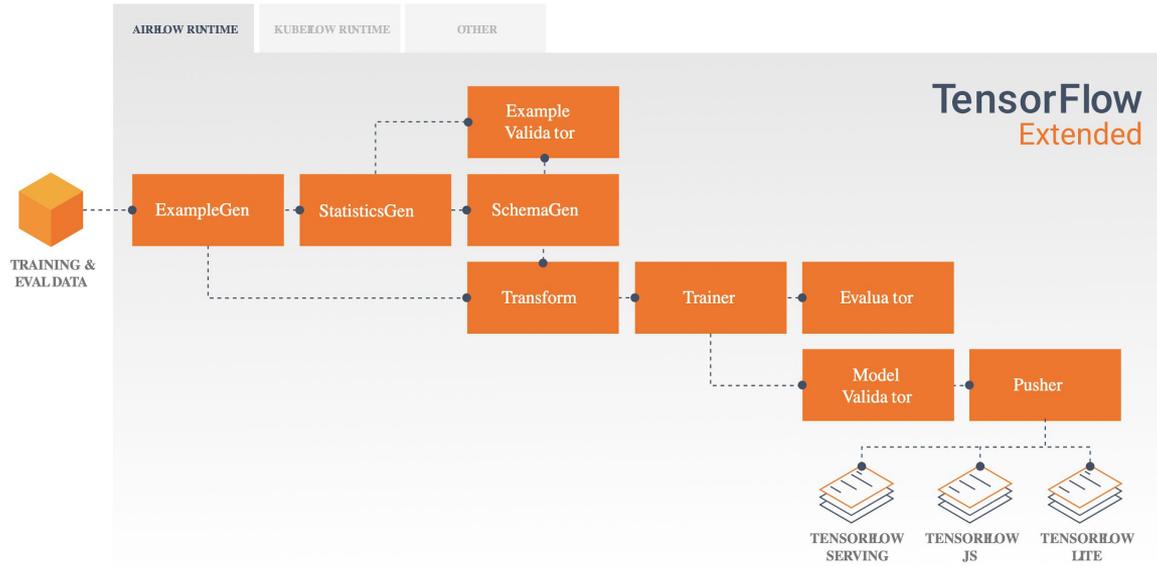
# Division of Labor

# TensorFlow Extended

# Challenge: Serving Environment

- How to Deploy Models?
  - Zero Downtime
  - Canary
- Multiple Models?
  - Testing
  - Different Scenarios

  - Updating models

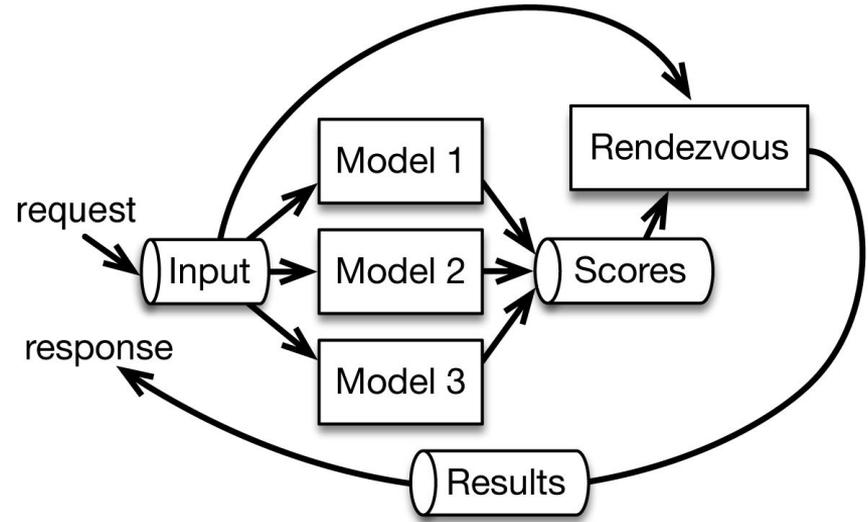  - Shadow models

  - A/B Canary testing

https://ai.googleblog.com/2016/02/running-your-models-in-production-with.html

- Common Metadata

# Challenge: Serving Environment

- How to Deploy Models?
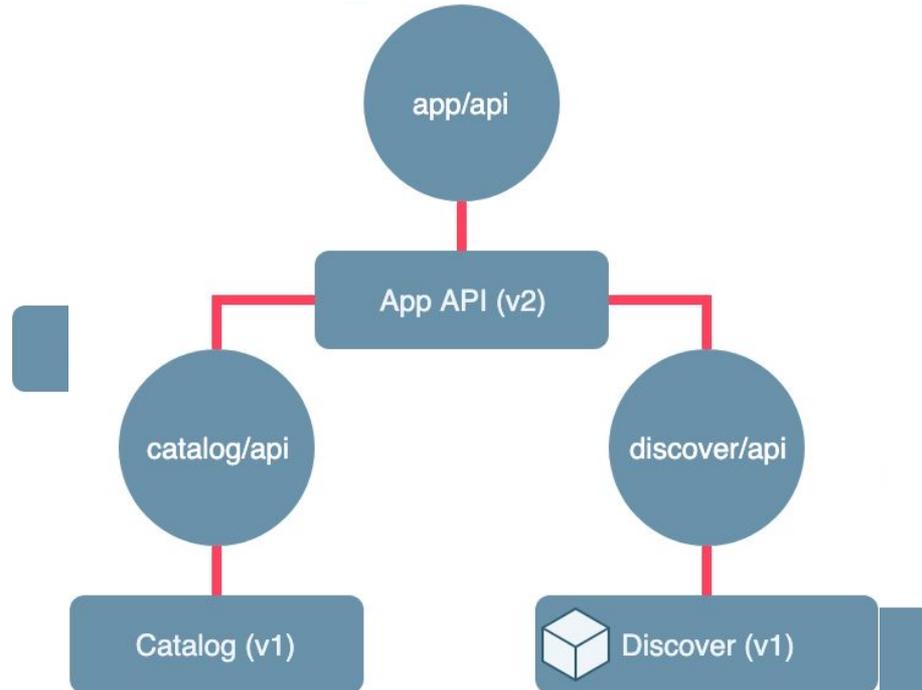  - Zero Downtime
  - Canary
- Multiple Models?
  - Testing



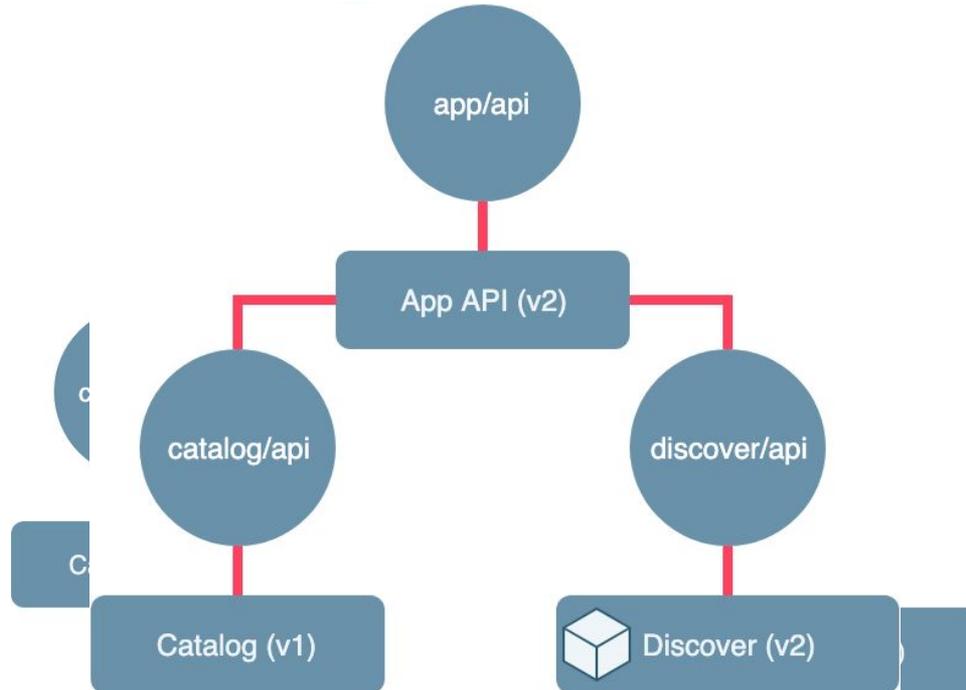https://mapr.com/ebooks/machine-learning-logistics/ch03.html

# A tale of ~~two~~ one world
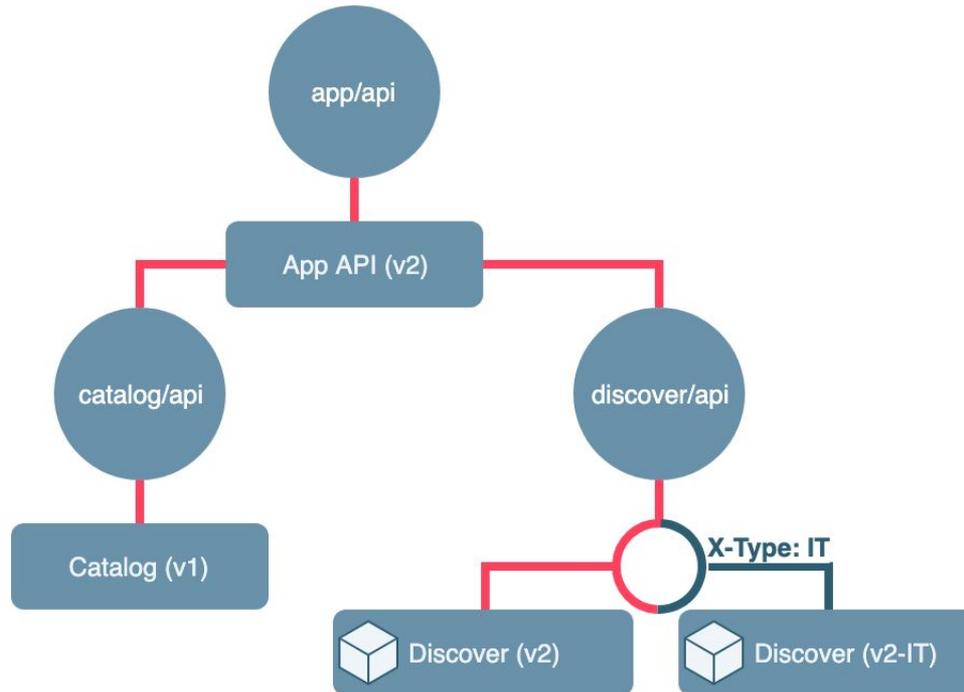




Deployment of Services
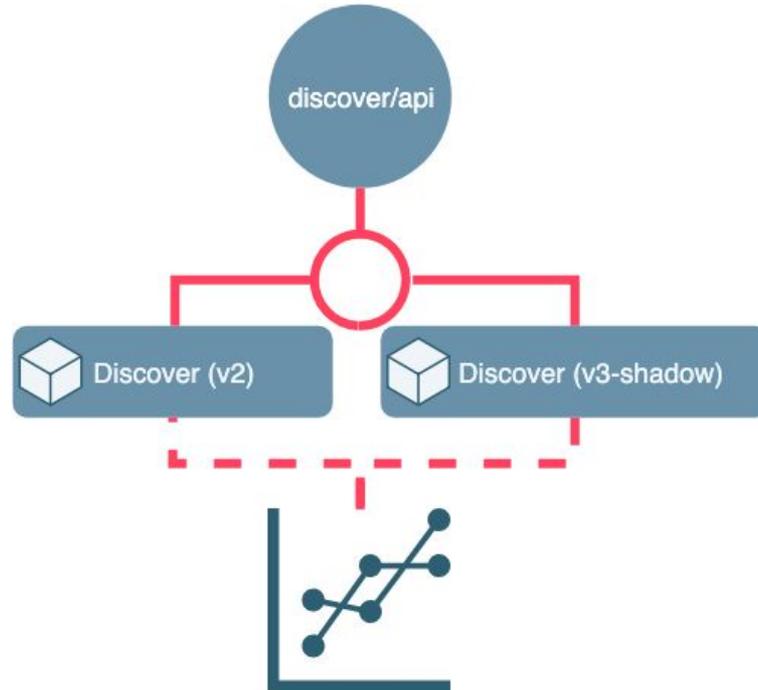
# Demo: Canary Release Service

# Demo: Canary Release ML Model

# Scenario: Model Segmentation

# Scenario: Shadow Traffic

# Talk with us

Join use for a chat the at the Mesosphere stand ….

# What's next for Vamp?

- Integrating Istio for smart networking
- AI/ML based workflows & policies for self-learning and optimisation
- More CI and APM integrations
- Additional release, validation, experimentation  and optimisation workflows
- Hybrid containers & serverless support
- Cloud-based version with new & improved core "engine"