# Building A Fun-Sized MLOps Stack From Scratch

**Mikiko**

Head of MLOps

Powered By:

{eature{orm

# Goals

- What are the main problems MLOps tries to solve.

- What are the most common tools being used & their drawbacks.

- What are some OSS projects & tools that have been developed in the past 2-3 years and how do they solve some of the pain points of the prior tools.

- What is the realistic roadmap for companies that are forever "not-Google" scale but want to continue improving their data and ML maturity.
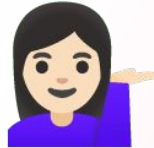
# How?

🤷🏻‍♀️

**By**
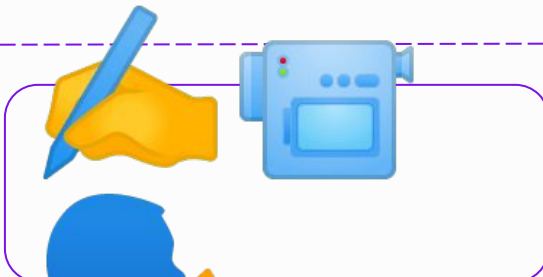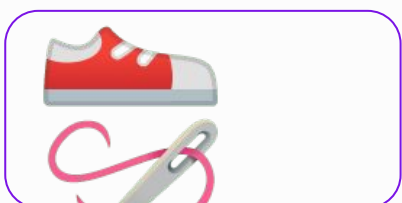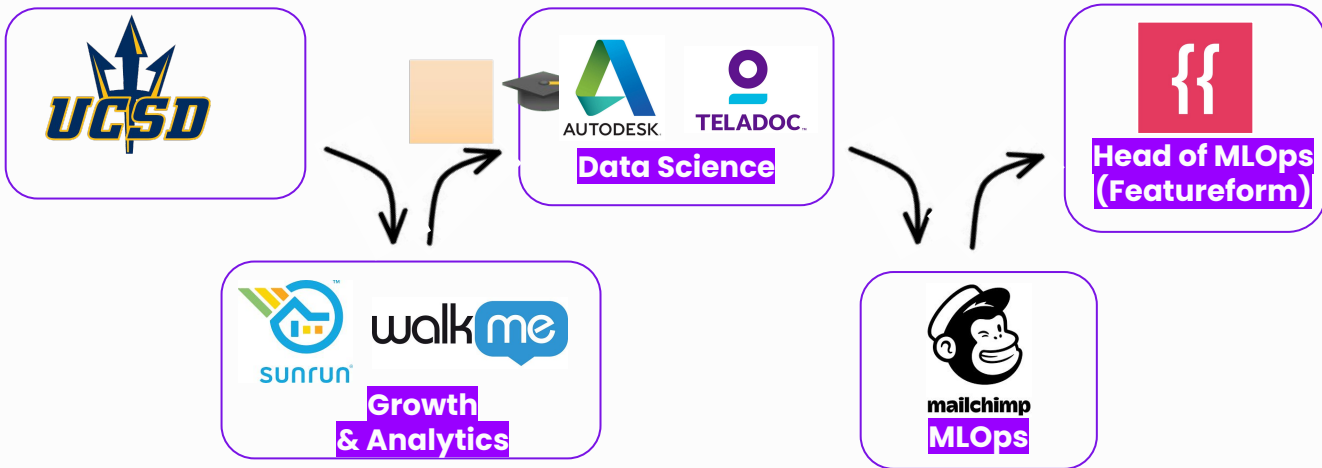
- Describing the original promises of MLOps (& the current shortfalls).

- Understanding the "Jobs-To-Be-Done" of Data Scientists (& how the current ecosystem supports them).

- Describe the pain-points of the Solo Data Scientist, the SMB Data Science Team, & the areas of opportunity for Enterprises.

- Propose stacks that can be easily implemented in a relatively short period (sometimes even a day!).

# Who Am I?

💁🏻‍♀️

# Mikiko Bazeley (But Call Me Mickey 🐭)

**UCSD**

**Data Science** — AUTODESK · TELADOC

**Head of MLOps (Featureform)**

**Growth & Analytics** — sunrun · walk me

**MLOps** — mailchimp

# What's The Current Landscape Look Like?

🗺️

THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE

# 2012



Big Data Landscape

© Matt Turck (@mattturck) and ShivonZilis (@shivonz)

{eature{orm

THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE

— This Isn't Even All Of Them —

# What Was The Original Promise of MLOps?

🤲

# Design => Train => Run

# The Three Dimensions: **Velocity, Throughput, Risk**

# We Failed

# The OG Users of

# MLOps:

🧑🏻‍💻 **Data Scientists**

🧑🏻‍🔬

# Trying To Mimic The Heavyweights



Figure 4. End-to-end MLOps architecture and workflow with functional components and roles

## Find database objects

Starting with...

**PRODUCTION**

▼ Tables

- NEW_V2
- NEW_V2_JAKE_FINAL
- NEW_V2_DO_NOT_TOUCH
- NEW_V3_TEST
- NEW_V53_NOT_SURE
- NEW_V2_MODIFIED
- NEW_V4__MODIFIED_TODAY
- NEW_IM_NOT_SURE
- NEW_MIGHT_BE_OLD
- NEW_NOT_THAT_OLD
- NEW

Where have you faced the biggest challenges in productionizing models?

| Collecting, curating and cleaning data | Data transformation and data versioning and lineage | Deploying and serving | AI monitoring, observability and explainability | Labeling | Synthetic data | Training |

# Ouch!

Need to handle evolving data & Data sources can suddenly change without announcement

Insufficient data documentation

Creating Meaningful Data Alerts is Challenging

"A Hotbed of Bias." Efforts to Assess, Prevent, & Mitigate Bias

Industry-Classroom Mismatch.

Training data is often insufficient and incomplete

"Experiment, Iterate, See We're Getting Closer." A Model Is Never Finished

"Data Literacy Is Not a Silver Bullet." On Communication & Collaboration

Dev-Prod Env Mismatch

Industry-Classroom Mismatch.

Undocumented Tribal Knowledge

Product requirements require input from the model team & Lack of ML knowledge in managers

# What We

# Should Be Doing

💪

📌 **Best Practice #1**

# Treating The ML Platform As

# A Product

# Characteristics of a Product

Does What It Needs ⚙️🛠️

Users/Customers 🙋🙋🙋🙋🙋🙋

Products

Named 🎀

Long-lived ⏳🏞️🏛️🗽

Owned 👩🏻

Evolve 🦖➡️🐓🐊

# Applied to Platforms

Composable ⚙️ 🧰

Quick, Easy Start 🟢 🏁 🏎️

Internal Community 🙋🙋🙋🙋🙋🙋

Platforms

Secure & Compliant 🧧

Up-to-Date ⏳ 🗻 🏛️ 🗽

Self-Service 👩🏻‍💼

Attraction, Not Coercion 🦖 ➡️ 🐓 🐊

{eatureform

# Measuring Data Science DevExp

| What We Should Care About | What We Should Measure (Examples) |
|---|---|
| Product<br>(Activation, Engagement, Adoption) | Weekly Active Users, Engagement, Adoption Rate, etc |
| User Satisfaction | NPS, WAU Retention, Sessions Per User, # Tickets, etc |
| Platform Performance<br>(Reliability, Availability, Scalability) | SLOs, Latency, # Incidents, etc |

# Empire State of Mind



Different Types of Platform Teams

APIs, tools, services, knowledge & support

Unified Developer Experience Characterized By Ease of Use + High Bar

BRICKSET

📌 **Best Practice #2**

**Prioritizing Enablement,**

**The Last Mile of True**

**Platform Adoption**

# The Stereotypes & Tropes

# Being User Centric

Do you have anyway to get feedback?

How do we know if it meets their needs?

Are we prioritizing the right work?

If the platform is working for them?

If the platform is working for them?

# Do You Know Who Your Customer Is?



**Degree Field of Study**

Legend:
- Computer Science
- Engineering
- Business/Economics
- Math/Stats
- Natural Sciences
- Data Science
- Soc Sci and Lib Arts
- Other

Y-axis: % (0, 20, 40, 60, 80, 100)

X-axis (Job Title): Data Scientist, Machine Learning Engineer, Software Engineer, Data Analyst, Data Engineer

# Do You Know Who Your Customer Is?

{eature{orm

31

📌 **Best Practice 3**

**Seamless Iteration**

**& Making Every Data Scientist**

**A 10X Data Scientist**

# 1. Solving Cognitive Overload

# 2. Cut Friction & Increase Flow

AMAZON · GOOGLE · MICROSOFT · ORACLE

Manu Cornet: Conway's Law

🤯 So Many Tools & Services Out There, And Yet… 🤯

# … Raise Your Hand ✋

# If Your Platform Is Basically

# S3 +Spark + Redis

# (And Bash Scripts)

# And While Those Are Great Tools...

# ... There Are Certain Classes of Problems They Don't Solve

🙅🏻‍♀️

# MLOps problems fall into two categories

# Specifically,

## Problems Around

👉 **Workflows & The Data**

**Science "Jobs-To-Be-Done"**👈

# MLOps problems fall into two categories

# Let's Describe The

# Data Science

# ✅ "Jobs-To-Be-Done" ✅

**&**

**Propose A Framework**

**For How We Can Effectively Map Tools**

**To Create An Effective Stack**

# DS Jobs To Be Done By:
# {Level of Abstraction} VS {Stage of Lifecycle}



Dimension:
Level of Abstraction

Layer 4:
Platform

Layer 3:
Workflow

Layer 2:
Services

Layer 1:
Compute Frameworks

Layer 0:
Hardware

Dimension: Stage of ML Lifecycle

Dataset Eng
+
Data Analysis
+
Feature
Engineering
& Serving

Model
Experimentation
+
Training

Deployment
+
Serving

Monitoring
+
Evaluation

# {Level of Abstraction} VS {Stage of Lifecycle}



Dimension:
Level of Abstraction

Layer 4:
Platform

Layer 3:
Workflow

Layer 2:
Services

Layer 1:
Compute Frameworks

Layer 0:
Hardware

Dimension:
Stage of ML Lifecycle

Dataset Eng
+
Data Analysis
+
Feature
Engineering
& Serving

Model
Experimentation
+
Training

Deployment
+
Serving

Monitoring
+
Evaluation

# Choosing The Right Layer For the Right Job To be Done

# Who Is Doing

# What*

**\*(As can be determined by their docs)**

# Choosing The Right Layer For the Right Job To be Done

# Let's Apply The Framework For:

✅ **The Solo Data Scientist**

✅ **The SMB Data Science Team**

✅ **The Enterprise DS Org**

# The Needs of

# "The Solo Data Scientist"

📝

# User Story: The Solo Data Scientist

| I need to… | Pain-Point |
|---|---|
| Keep track of data, model, & code artifacts, including changes & experimentation runs. | **Versioning & Documentation** |
| Quickly iterate between features, algorithms, & hyperparameter tuning. | **Experimentation Tracking** |
| Train models on a "large enough" amount of data with access to GPUs. | **Serverless GPU** |
| Do everything with the least amount of overhead possible with the least amount of steps. | **Compatibility With Existing Product Stack** |

# Stack 1:

# The Duke Nukem Stack

# Stack 1: **The Duke Nukem Stack (Solo Data Scientist)**

# The Needs of

# "The SMB Data Science Team"

📝

# User Story: **The DS Team**

| We need to… | Pain-Point |
|---|---|
| Collaborate with other members of the DS team (and potentially even external partners) on projects, with visibility into progress or health of data science assets. | **Collaboration** |
| Share & distribute knowledge asynchronously, while getting ahead of human bottlenecks & the accumulation of tribal knowledge. | **Documentation & Discoverability** |
| Ensure we're not "reinventing the wheel" across the organization & repeating work. | **Reuse & Resource Sharing** |
| Be notified when model pipelines and prediction services aren't working as expected with insight into failure conditions. | **Fine-grained Monitoring & Evaluation** |

# Stack 2:

# The Serious Business Stack

Stack 2: **The Serious Business Stack (SMB)**

# The Needs of

# "The Enterprise Org"

📝

# User Story: **The Organization**

| We need to... | Pain-Point |
|---|---|
| Span and unify multiple infrastructure providers (including multi-cloud and on-prem), model deployment patterns, and model serving architectures as seamlessly as possible. | **Heterogeneous Infrastructure** |
| Handle a wide variety of regulation around data & models, log compliance related information & data, & streamline communication & visibility. | **Governance, Access Control, Audit Logs** |
| Interface with non-DS teams (including other engineering teams, as well as non-eng teams like legal & marketing). | **Cross-Functional Workflows** |

# Stack 3:

# The Olly Olly Oxen Free Stack

# Stack 3: **The Olly Olly Oxen Free Stack (Enterprise)**



Dimension:
Level of Abstraction

| | | | | |
|---|---|---|---|---|
| **Layer 4:** **Platform** | MLOps Platform Layer | | | |
| **Layer 3:** **Workflow** | Feature Orchestration & Management | Experiment Orchestration & Management | Model Orchestration & Management | Evaluation Orchestration & Management |
| **Layer 2:** **Services** | Feature Service | Training Service | Serving Service | Monitoring Service |
| **Layer 1:** **Compute Frameworks** | Across Layer 1 | | | |
| **Layer 0:** **Hardware** | Across Layer 0 | | | |

Dimension:
Stage of
ML Lifecycle

| Dataset Eng + Data Analysis + Feature Engineering & Serving | Model Experimentation + Training | Deployment & Serving | Monitoring & Evaluation |
|---|---|---|---|

# There Will Be No

# "Modern MLOps" Stack

🙅🏻‍♀️

# But Wait!

# There's Hope!

🙏🏼

ALWAYS HOPE, THERE IS
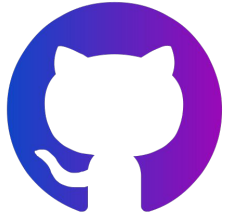
{eature{orm

# In Closing

📫

# Takeaways

- Describing the original promises of MLOps (& the current shortfalls).

- Understanding the "Jobs-To-Be-Done" of Data Scientists (& how the current ecosystem supports them).

- Describe the pain-points of the Solo Data Scientist, the SMB Data Science Team, & the areas of opportunity for Enterprises.

- Propose stacks that can be easily implemented in a relatively short period (sometimes even a day!).

# Feel Free To Chat With Me During Office Hours

👩🏻‍🏫

**Repository**

bit.ly/3Yz2G95

**Docs**

bit.ly/423SE2W

**LinkedIn**

Mikiko Bazeley

Head of MLOps