# Kubeflow

Wojciech Barczyński [Head of Engineering]
SMACC.io | Hypatos.ai

## Wojtek Barczynski

- Software Developer
- System Engineer
- Head of Engineering at hypatos.ai and SMACC.io

# Hypatos / SMACC.io

- Fintech Machine Learing
- Data capturing from document
- Validation
- Automation
- Deep learning

# Hypatos / SMACC.io
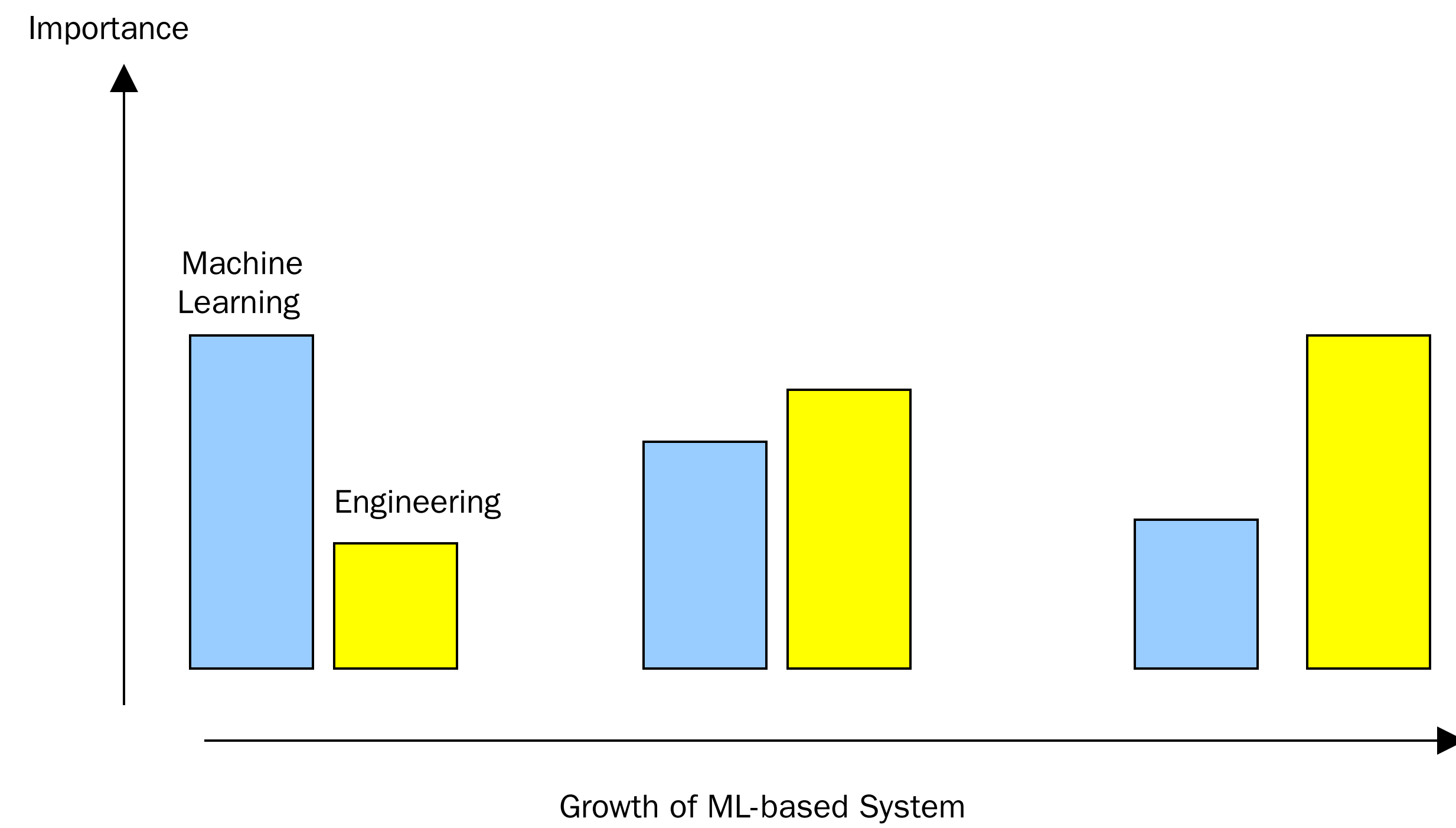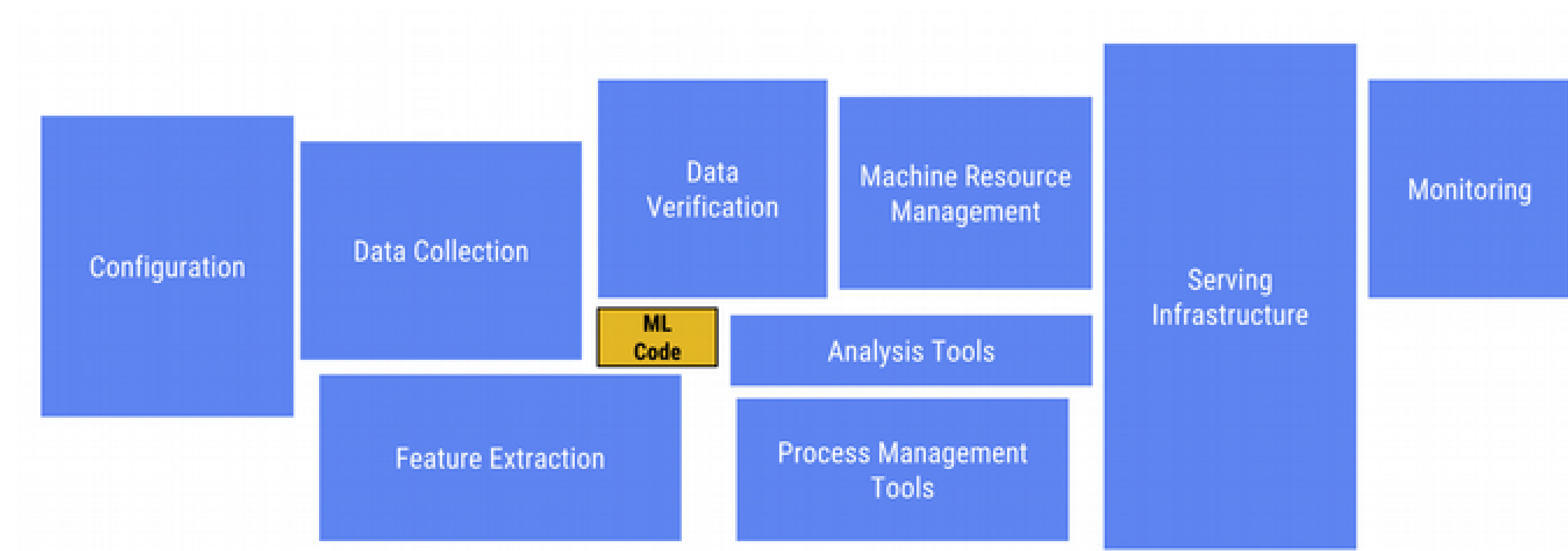


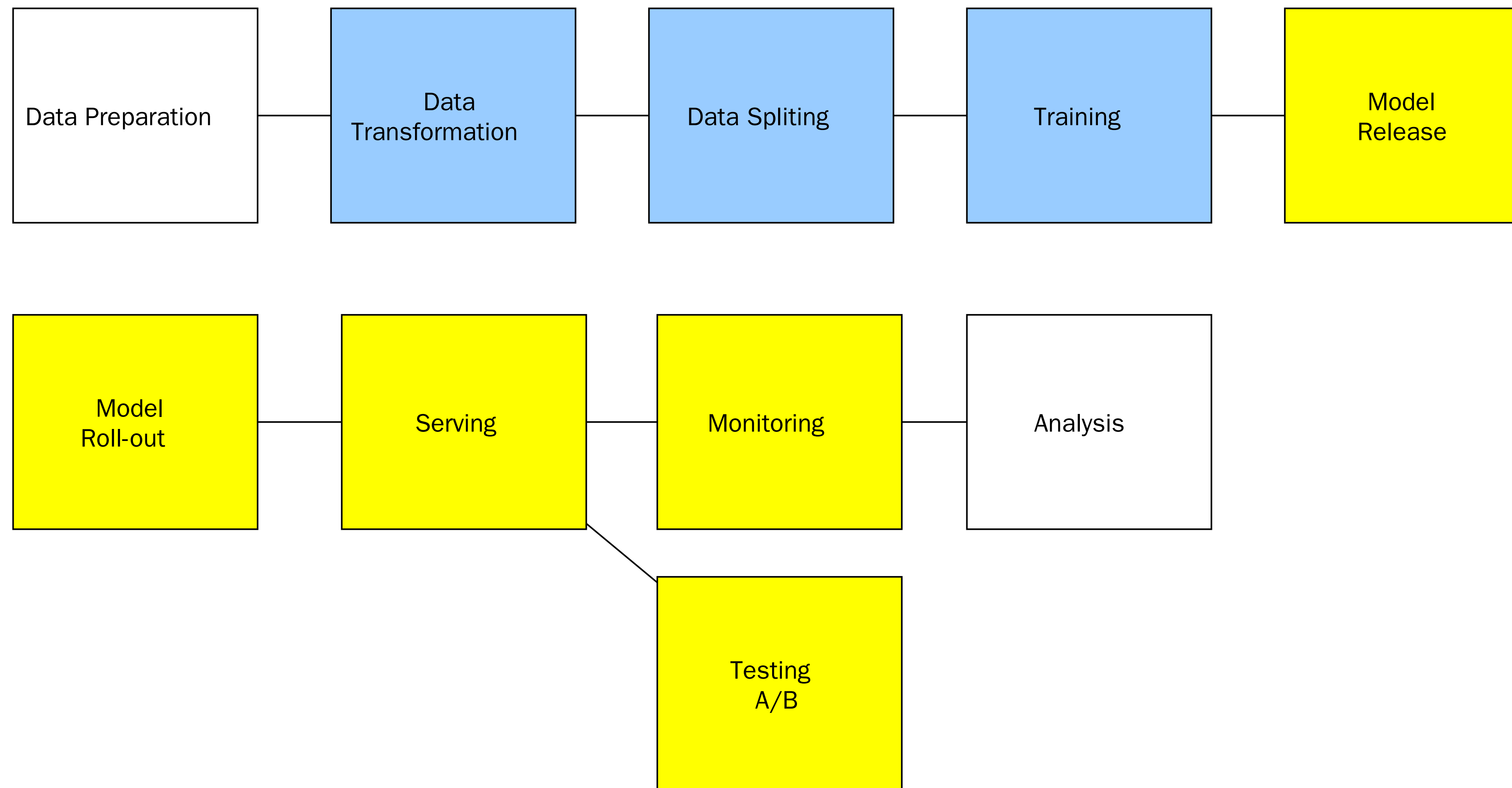CloudNative: Prometheus, Grafana, Fluentd, Grafana Loki

# Machine Learning



Importance

Machine
Learning

Engineering

Growth of ML-based System

# Big Companies

# Machine Learning

| Data Preparation | Data Transformation | Data Spliting | Training | Model Release |
|---|---|---|---|---|

| Model Roll-out | Serving | Monitoring | Analysis |
|---|---|---|---|

| Testing A/B |
|---|

# ML pipeline

Look a lot like
Continuous Integration / Deployment

# ML pipeline

What did we lean from XX years of CD/CI?

# What we leant from XX years of CD/CI?

Nobody likes it

# Kubeflow

- Easy the pain
- Unified experience
- One to rule them all
- Low bar; High ceiling

# Focus

- Scalability
- Composition
- Portability

# Focus
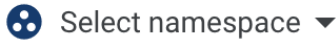
- Enable Dev(ML)ops culture

# DEMO

# Heart: pipelines

# Gears: component

# User Experience: notebook

# User Experience: notebook

# Python SDK

- Let data scientist and engineers work together

Kubeflow User story

# YAML vs Python SDK

```yaml
apiVersion: argoproj.io/v1alpha1
kind: Workflow
metadata:
  generateName: charts-of-accounts-
spec:
  arguments:
    parameters:
      - name: aws-cli-image
        value: "pbsmacc/aws-cli:latest"
      - name: prepare-dataset-image
        value: "smaccio/accounting-charts-prepare-dataset:v0.2
      - name: trainer-image
        value: "smaccio/accounting-charts-classifier:v0.2.0"
      - name: minio-client-image
```

# Python SDK

```python
@dsl.pipeline(
    name='Taxi Cab on-prem',
    description='Exar.'
)
def taxi_cab_classification(
    training = tf_train_op(
        preprocess.output,
        validation.outputs['schema'],
        learning_rate,
        hidden_layer_size,
        steps,
        'tips',
        '/mnt/%s' % preprocess_module,
        '/mnt',
```

# Python SDK

```python
def kubeflow_deploy_op(model: 'TensorFlow model', tf_server_na
                       pvolumes, step_name='deploy'):
    return dsl.ContainerOp(
        name=step_name,
        image='gcr.io/ml-pipeline/ml-pipeline-kubeflow-deploye
        arguments=[
            '--cluster-name', 'tfx-taxi-pipeline-on-prem',
            '--model-export-path', model,
            '--server-name', tf_server_name,
            '--pvc-name', pvc_name,
        ],
        pvolumes=pvolumes
    )
```

# More Python: Fairing SDK

- All power of kubeflow
  from your local jupyter notebook
- For hybird cloud

https://github.com/kubeflow/fairing

# User Experience: Tensorboard

# Tracking: Artifacts

- Emitted by steps as metadata

# Tracking: Artifacts

- Emitted by steps as metadata
  - `mlpipeline-ui-metadata.json`

## Focus on Data Scientist

- Self-service
- Provide familiar user experience and tools
- Share the knowledge
- Hide the engineering complexity

## Batteries Included

### Scale trainings

- TFJobs
- MPI Training
- PyTorch Training
- MXNet Training

# Machine Learning

| | | | | |
|---|---|---|---|---|
| Data Preparation | Data Transformation | Data Spliting | Training | Model Release |

| | | | |
|---|---|---|---|
| Model Roll-out | Serving | Monitoring | Analysis |

Testing A/B

## Kubeflow

- **How to serve the model**
- Operating service
- Observability, e.g., metric collection

# **Kubeflow**

- Operating service
- Observability, e.g., metric collection
- Deployment strategies
- ...

One CloudNative project comes to mind - Istio.

**Kubeflow**

Serving

- ML Model servers
- seldon.io
- kfserving

## ML Model Servers

- TFserving
- PyTorch Serving
- ...

## ML Model Servers

- Minimum configuration
- Serve a given trained ML model

  with Istio integration if needed

# Seldon.io

- More complex use cases

seldon.io

## Istio

- Observability Grafana
- Deployment strategies

# Istio

- All served models are available with Istio

# kfserving

# kfserving

```yaml
apiVersion: "serving.kubeflow.org/v1alpha2"
kind: "KFService"
metadata:
  name: "xgboost-iris"
spec:
  default:
    predictor:
      xgboost:
        storageUri: "gs://kfserving-samples/models/xgboost/iri
```

# Architecture



Laptop

Docker Images | IDE

Kubectl

my-model-dir
  ./train.py
  ./model.pb
  ./Dockerfile
  ./k8s_manifests/
      ./tfjob_train.yaml
      ./deploy_model.yaml

Auth proxy

ISTIO GateWay

K8s Cluster

kubeflow-alice namespace

Jupyter

Flask App

TFJob

TFServing

kubeflow-system namespace

Jupyter Spawner UI | Jupyter Controller

TFJob Dashboard UI | TFJob  Controller

StudyJob Controller | Katib UI

Pipelines servers | Profiles Controller

Alice Profile | Seldon Controller

Cloud Resources

Object Store | Identity Provider | NFS Share | Image Registry | Metadata

# Architecture

**Libraries and CLIs - Focus on end users**

Arena | kfctl | kubectl | fairing

**Systems - Combine multiple services**

katib | pipelines | Model DB

kube bench | notebooks | TFX

**Low Level APIs / Services (single function)**

TFJob | PyTorch Job | Pipelines CR

Argo | Jupyter CR | MPI CR

Seldon CR | Study Job | Spark Job

Metadata

IAM

Orchestration

Scheduling

Developed By Kubeflow

Developed Outside Kubeflow

\* Not all components shown

# What did we decided?

- Large project with many moving parts
- Take bits that we need and keep delivering
- Invest more time into observability

We do not have such a large team

# What did we pick?

1. Mostly Model + Code as Software Components
2. Automation project:
   - Argo in YAML
   - tfserving
   - tensorflow_transform.beam

# Deployment?

- Git-driven deployment
- Version the model and the code

we might use argo here as well

# Observability

- Prometheus + Grafana
- Dedicated metric collector

## Keep an eye on kubeflow

- Enterprise client projects on-prem
- Growing team

## Summary

- Enable Dev(ML)Ops culture
- Hide the engineering complexity

## Summary

- The learning part is the most compelling
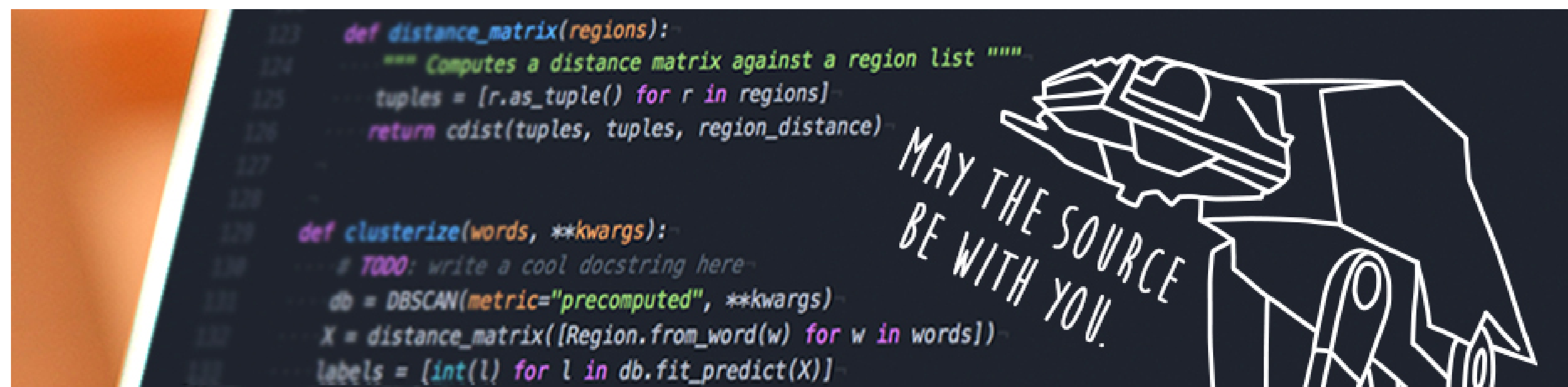- and self-service

# QUESTIONS?

```python
123  def distance_matrix(regions):
124      """ Computes a distance matrix against a region list """
125      tuples = [r.as_tuple() for r in regions]
126      return cdist(tuples, tuples, region_distance)
127
128
129  def clusterize(words, **kwargs):
130      # TODO: write a cool docstring here
131      db = DBSCAN(metric="precomputed", **kwargs)
132      X = distance_matrix([Region.from_word(w) for w in words])
133      labels = [int(l) for l in db.fit_predict(X)]
```

MAY THE SOURCE
BE WITH YOU.

# Big thanks to Piotr Brzostowski and whole BER+WAW team.

BACKUP

# Development

- How to handover to engineering?
- How did I trained the model X?
- Lineage and Metadata

# Operation

- How the model performs in production
- Is it better?
- Which data should I add to the next training?
- Low performance → roll back
- Keep the TCO reasonable

# Operation

- Observability: Prometheus, Grafana
- A/B testing: service mesh
- Serving model?