# Serverless ML Deployment on AWS

# Šimon Soták

- Studied Computer Graphics at Charles University in Prague
- Started out in Game Development (Keen SWH, Warhorse)
- Joined Represent in 2014

# Andrej Hoos

- Studied Computer Science and Physics at University of Glasgow
- Joined Represent in 2014

# Represent

# Arnold
# Schwarzenegger

## $1M
*revenue*

## 69K
*products sold*

CustomInk kupuje za několik miliard korun startup se slovenským spoluzakladatelem

**Represent**

## PRG

Prague, Czech Republic
Šaldova 12, Karlín

## LA

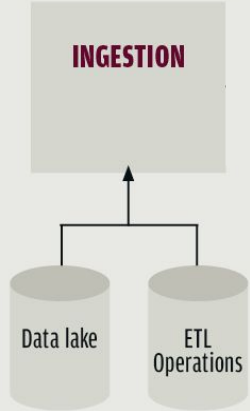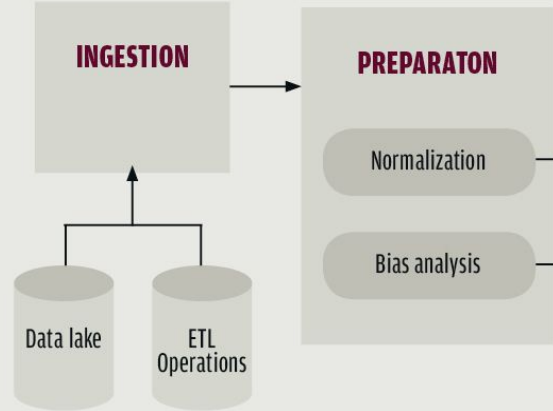Los Angeles, California
1680 Vine St, Hollywood

## DC

Washington metropolitan area
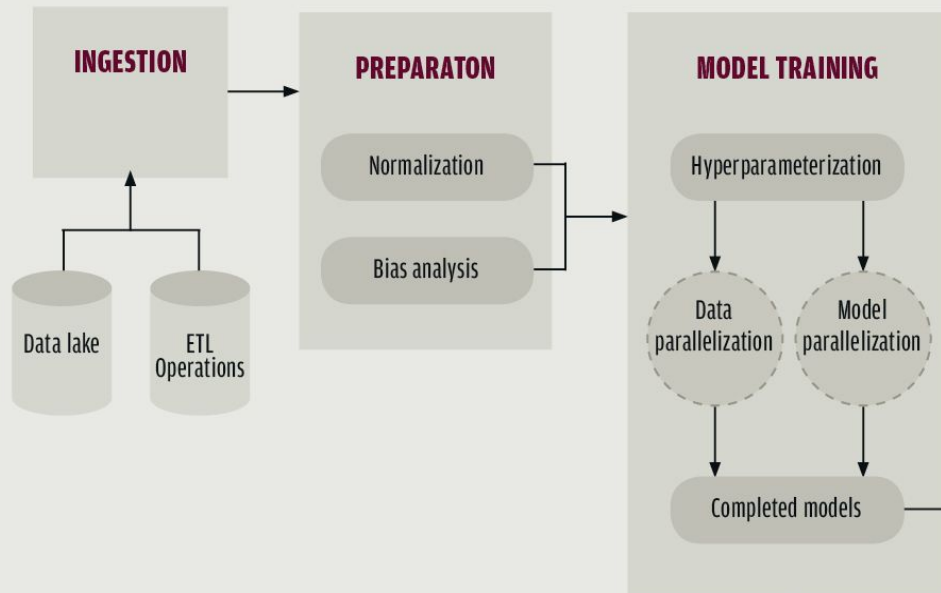2910 District Ave, Fairfax

# Machine Learning

# Machine Learning Process

# Machine Learning Process

INGESTION

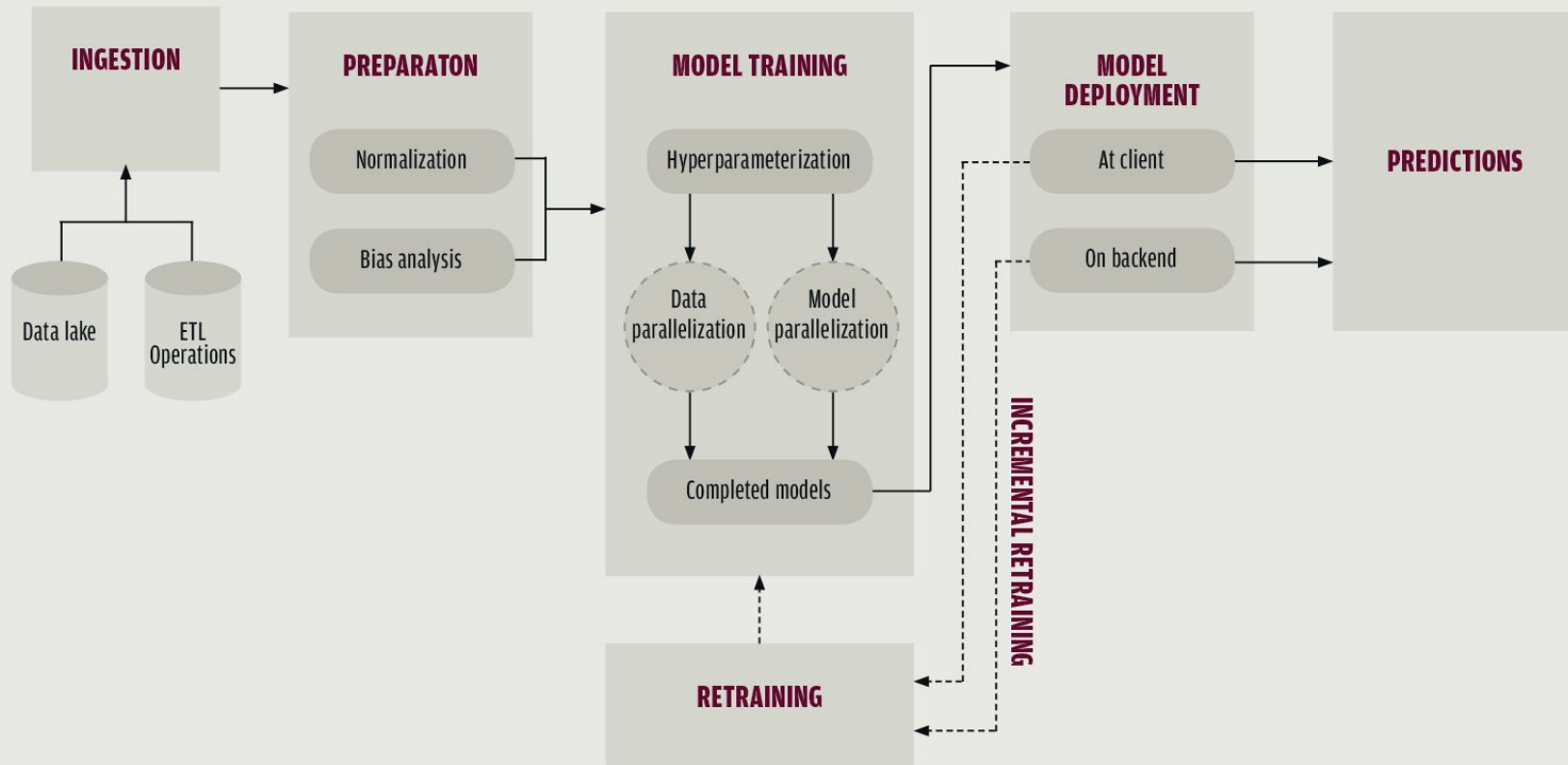PREPARATON

Normalization

Bias analysis

Data lake

ETL Operations

# Machine Learning Process

# Machine Learning Process

**Represent**

**INGESTION**

Data lake    ETL Operations

**PREPARATON**

Normalization

Bias analysis

**MODEL TRAINING**

Hyperparameterization

Data parallelization    Model parallelization

Completed models

**MODEL DEPLOYMENT**

At client

On backend

**PREDICTIONS**

**RETRAINING**

**INCREMENTAL RETRAINING**

# Serverless ML Deployment on AWS

# How to deploy ML models

| Approach | AWS Service |
|---|---|
| Hosted servers | EC2 |
| Docker / Kubernetes containers | ECS / EKS |
| **Serverless** | **Lambda** / Fargate |
| ML-tailored solutions | Sagemaker |

# What is serverless?

- Function as a Service
- Fully managed
- You provide only the code to be run
- Advantages
  - No maintenance required
  - Scaling "for free"
  - Easier development
  - Lower infrastructure cost

# Python ML Deployment on Serverless

- AWS Lambda
- Why Amazon?
- Why Python?

# So, can I host my Python ML model on AWS Lambda?

Yes*…

# Cold Start

- First request – "bootup" & init
- Instance "shuts down" after a while if no requests arrive
- Parallel cold starts
- Pre-warming
- ML Model loading

# Issue to overcome: Package size

- AWS Lambda – 250 MB max package size
- Scikit-learn – 200+ MB
- Tensorflow – over 250 MB

# Did you say "Easier development" earlier?

# Introducing: `lambdipy`

- Tool for building python packages for AWS Lambda
- [github.com/customink/lambdipy](github.com/customink/lambdipy)
- Builds "minified" packages for AWS Lambda
- Contains pre-built popular packages
- Active development – contributions welcome

# What does `lambdipy` do

- Builds any python ML package
    - Uses Docker with Lambda-like environment
    - Downloads pre-built popular packages
- **strip** binaries
    - Binaries take the most space in ML packages
    - This reduces their file size significantly
- Remove some unneeded files
    - Tests, caches, package metadata
- scipy: 140 MB → 33 MB

# Example

# MNIST dataset

# Simple image classifier

- Handwritten digit classification
- Inspired by blog post by Bikramjot Singh Hanzra[1]
- **sam** for deployment
- **pipenv** for dependency management

1. http://hanzratech.in/2015/02/24/handwritten-digit-recognition-using-opencv-sklearn-and-python.html

# Example: `train.py`

```python
dataset = datasets.fetch_mldata("MNIST Original")
```

# Example: `train.py`

```python
dataset = datasets.fetch_mldata("MNIST Original")

# Train a ML model
model = LinearClassifer()
model.fit(features, labels)
```

# Example: `train.py`

```python
dataset = datasets.fetch_mldata("MNIST Original")

# Train a ML model
model = LinearClassifer()
model.fit(features, labels)

joblib.dump(model, "model.pkl")
```

# Example: `predict.py`

```
model = joblib.load("model.pkl")
```

# Example: `predict.py`

```python
model = joblib.load("model.pkl")

def predict(context, _):
```

# Example: `predict.py`

```python
model = joblib.load("model.pkl")

def predict(context, _):
    request_body = json.loads(context["body"])
```

# Example: `predict.py`

```python
model = joblib.load("model.pkl")

def predict(context, _):
    request_body = json.loads(context["body"])

    # Load image from url and extract its features
    prediction = model.predict(image_features)
```

# Example: `predict.py`

```python
model = joblib.load("model.pkl")

def predict(context, _):
    request_body = json.loads(context["body"])

    # Load image from url and extract its features
    prediction = model.predict(image_features)

    return {
        "statusCode": 200,
        "body": json.dumps({"digit": int(prediction[0])})
    }
```

# Example: template.yml

```yaml
…

Resources:
  PredictFunction:
    Type: AWS::Serverless::Function
    Properties:
      CodeUri: ./build
      Handler: predict.predict
…
```

# Example: deployment

- Install dependencies
  - `pipenv install scikit-learn scikit-image numpy`
  - `pipenv install lambdipy aws-sam-cli --dev`
- Train the classifier
  - `pipenv run python train.py`
- Build the bundle
  - `pipenv run lambdipy build --from-pipenv -i predict.py -i model.pkl`
- Deploy
  - `pipenv run sam package --s3-bucket ...`
  - `pipenv run sam deploy --stack-name ...`

🤞🤞🤞 **Live Demo** 🤞🤞🤞

# Example: testing the predictions

**POST**

https://r7my4vk3t0.execute-api.us-east-1.amazonaws.com/Prod/predict
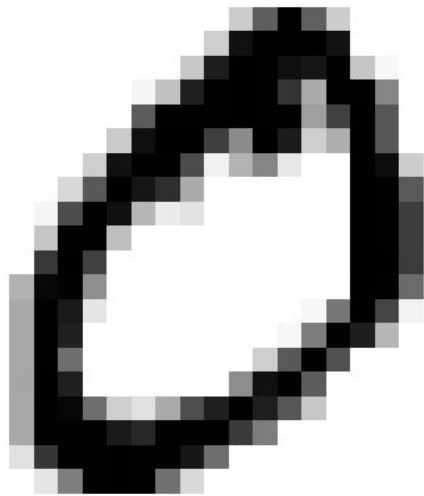```
{
    "url": "https://i.stack.imgur.com/FK0FB.png "
}
```

`{ "digit": 0 }`



```
curl -d '{"url": "https://i.stack.imgur.com/FK0FB.png"}' -X POST
https://r7my4vk3t0.execute-api.us-east-1.amazonaws.com/Prod/predict
```
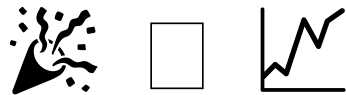
# Example: testing the predictions

**POST**

https://r7my4vk3t0.execute-api.us-east-1.amazonaws.com/Prod/predict
```
{
    "url": "https://goo.gl/K9ditV"
}
```

{ "digit": 2 }

```
curl -d '{"url": "https://goo.gl/K9ditV"}' -X POST
https://r7my4vk3t0.execute-api.us-east-1.amazonaws.com/Prod/predict
```

🎉 □ 📈

# Production ML system complete!

# Alternative approaches

- Zappa
  - Deploys Python web apps to serverless environments
- Google Cloud Functions
  - Python support currently in Beta

# Thank you!

This presentation: tiny.cc/aws-lambdipy

We are hiring and always looking to have good conversations

represent.com/join

 Represent