# Human Centric

# Machine Learning Infrastructure

# @

**NETFLIX**

**Ville Tuulos**

QCon SF, November 2018

# Caveman Cupcakes

VS

the human is the bottleneck
VS

VS

the human is the bottleneck | the human is the bottleneck

**Build**

Build

Data Warehouse

Build

Compute Resources

Data Warehouse

**Build**

Job Scheduler

Compute Resources

Data Warehouse

Build

Versioning

Job Scheduler

Compute Resources

Data Warehouse

Build

Collaboration Tools

Versioning

Job Scheduler

Compute Resources

Data Warehouse

**Build**

Model Deployment

Collaboration Tools

Versioning

Job Scheduler

Compute Resources

Data Warehouse

**Build**

Feature Engineering

Model Deployment

Collaboration Tools

Versioning

Job Scheduler

Compute Resources

Data Warehouse

Build

ML Libraries

Feature Engineering

Model Deployment

Collaboration Tools

Versioning

Job Scheduler

Compute Resources

Data Warehouse

**Build**

How much data scientist cares

| |
|---|
| ML Libraries |
| Feature Engineering |
| Model Deployment |
| Collaboration Tools |
| Versioning |
| Job Scheduler |
| Compute Resources |
| Data Warehouse |

Build

How much data scientist cares

How much infrastructure is needed

ML Libraries

Feature Engineering

Model Deployment

Collaboration Tools

Versioning

Job Scheduler

Compute Resources

Data Warehouse

Deploy

# Deploy

## No plan survives contact with enemy

# Deploy

## No plan survives contact with enemy

# No model survives contact with reality

our ML infra supports
two **human** activities:
**building** and **deploying**
data science workflows.

# Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?

*Custom Titus executor.*

Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?

*Custom Titus executor.*

**How to** access data at scale?

*Slow!*



Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?
*Custom Titus executor.*

**How to** access data at scale?
*Slow!*

**How to** schedule the model to update daily?
*Learn about the job scheduler.*

Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?
*Custom Titus executor.*

**How to** access data at scale?
*Slow!*

**How to** schedule the model to update daily?
*Learn about the job scheduler.*

**How to** expose the model to a custom UI?
*Custom web backend.*



Bad Old Days

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?
*Custom Titus executor.*

**How to** access data at scale?

*Slow!*

**How to** schedule the model to update daily?
*Learn about the job scheduler.*

**How to** expose the model to a custom UI?
*Custom web backend.*

Bad Old Days



Time to production:
**4 months**

# Data Scientist built an NLP model in Python. Easy and fun!

**How to** run at scale?

*Custom Titus executor.*

**How to** access data at scale?

*Si...*

**How to** schedule the model to update daily?

*...ut the job scheduler.*

**How to** ... to a custom UI?

*Custom web ba...*

**How to** monitor models in production?

**How to** iterate on a new version without breaking the production version?

**How to** let another data scientist iterate on her version of the model safely?

**How to** debug yesterday's failed production run?

## Bad Old Days



...oduction:

**4 months**

**How to** backfill historical data?

**How to** make this faster?

ML Wrapping: **Metaflow**

models — ML Libraries: R, XGBoost, TF etc.

prototyping — Notebooks: **Nteract**

Job Scheduler: **Meson**

compute — Compute Resources: **Titus**

Query Engine: **Spark**

data — Data Lake: **S3**

Airflow

kubernetes

Metaflow

# Build

# How to get started?

input → **compute** → output

```python
def compute(input):
    output = my_model(input)
    return output
```

# How to get started?



**input** → **compute** → **output**

```python
def compute(input):
    output = my_model(input)
    return output
```

```
# python myscript.py
```

# How to structure my code?



```python
from metaflow import FlowSpec, step

class MyFlow(FlowSpec):

    @step
    def start(self):
        self.next(self.a, self.b)

    @step
    def a(self):
        self.next(self.join)

    @step
    def b(self):
        self.next(self.join)

    @step
    def join(self, inputs):
        self.next(self.end)

MyFlow()
```

# How to structure my code?



```
# python myscript.py run
```

```python
from metaflow import FlowSpec, step

class MyFlow(FlowSpec):

    @step
    def start(self):
        self.next(self.a, self.b)

    @step
    def a(self):
        self.next(self.join)

    @step
    def b(self):
        self.next(self.join)

    @step
    def join(self, inputs):
        self.next(self.end)

MyFlow()
```

# How to deal with models written in R?



```r
metaflow("MyFlow") %>%
  step(
    step = "start",
    next_step = c("a", "b")
  ) %>%
  step(
    step = "A",
    r_function = r_function(a_func),
    next_step = "join"
  ) %>%
  step(
    step = "B",
    r_function = r_function(b_func),
    next_step = "join"
  ) %>%
  step(
    step = "Join",
    r_function = r_function(join,
                  join_step = TRUE),
```

# How to deal with models written in R?



```
# RScript myscript.R
```

```r
metaflow("MyFlow") %>%
  step(
    step = "start",
    next_step = c("a", "b")
  ) %>%
  step(
    step = "A",
    r_function = r_function(a_func),
    next_step = "join"
  ) %>%
  step(
    step = "B",
    r_function = r_function(b_func),
    next_step = "join"
  ) %>%
  step(
    step = "Join",
    r_function = r_function(join,
                 join_step = TRUE),
```

**Metaflow** adoption at Netflix

**134** projects on Metaflow as of November 2018

# How to prototype and test my code locally?



```python
@step
def start(self):
    self.x = 0
    self.next(self.a, self.b)

@step
def a(self):
    self.x += 2
    self.next(self.join)

@step
def b(self):
    self.x += 3
    self.next(self.join)

@step
def join(self, inputs):
    self.out = max(i.x for i in inputs)
    self.next(self.end)
```

# How to prototype and test my code locally?



```
# python myscript.py resume B
```

```python
@step
def start(self):
    self.x = 0
    self.next(self.a, self.b)

@step
def a(self):
    self.x += 2
    self.next(self.join)

@step
def b(self):
    self.x += 3
    self.next(self.join)

@step
def join(self, inputs):
    self.out = max(i.x for i in inputs)
    self.next(self.end)
```

# How to get access to more CPUs, GPUs, or memory?

**16 cores, 1GPU**

start → A → join → end

**200GB RAM**

start → B → join

```python
@titus(cpu=16, gpu=1)
@step
def a(self):
    tensorflow.train()
    self.next(self.join)

@titus(memory=200000)
@step
def b(self):
    massive_dataframe_operation()
    self.next(self.join)
```

# How to get access to more CPUs, GPUs, or memory?

**16 cores, 1GPU**

```
start → A
      → B
A → join
B → join
join → end
```

**200GB RAM**

```
# python myscript.py run
```

```python
@titus(cpu=16, gpu=1)
@step
def a(self):
    tensorflow.train()
    self.next(self.join)

@titus(memory=200000)
@step
def b(self):
    massive_dataframe_operation()
    self.next(self.join)
```

# How to distribute work over many parallel jobs?



```python
@step
def start(self):
    self.grid = ['x','y','z']
    self.next(self.a, foreach='grid')

@titus(memory=10000)
@step
def a(self):
    self.x = ord(self.input)
    self.next(self.join)

@step
def join(self, inputs):
    self.out = max(i.x for i in inputs)
    self.next(self.end)
```

# How to access large amounts of input data?



```python
from metaflow import Table

@titus(memory=200000, network=20000)
@step
def b(self):
    # Load data from S3 to a dataframe
    # at 10Gbps
    df = Table('vtuulos', 'input_table')
    self.next(self.end)
```

# Case Study: Marketing Cost per Incremental Watcher

1. Build a separate model for every new title with marketing spend.
   **Parallel** foreach.

2. Load and prepare input data for each model.
   Download Parquet **directly from S3**.
   Total amount of model input data: **890GB**.

3. Fit a model.
   Train each model on an instance with **400GB of RAM, 16 cores**.
   The model is **written in R.**

4. Share updated results.
   Collect results of individual models, write to a table.
   Results shown on **a Tableau dashboard.**

# Deploy

# How to version my results and access results by others?

savin: unsampled_model

david: sampled_model



```python
# Access Savin's runs
namespace('user:savin')
run = Flow('MyFlow').latest_run
print(run.id) # = 234
print(run.tags) # = ['unsampled_model']

# Access David's runs
namespace('user:david')
run = Flow('MyFlow').latest_run
print(run.id) # = 184
print(run.tags) # = ['sampled_model']

# Access everyone's runs
namespace(None)
run = Flow('MyFlow').latest_run
print(run.id) # = 184
```

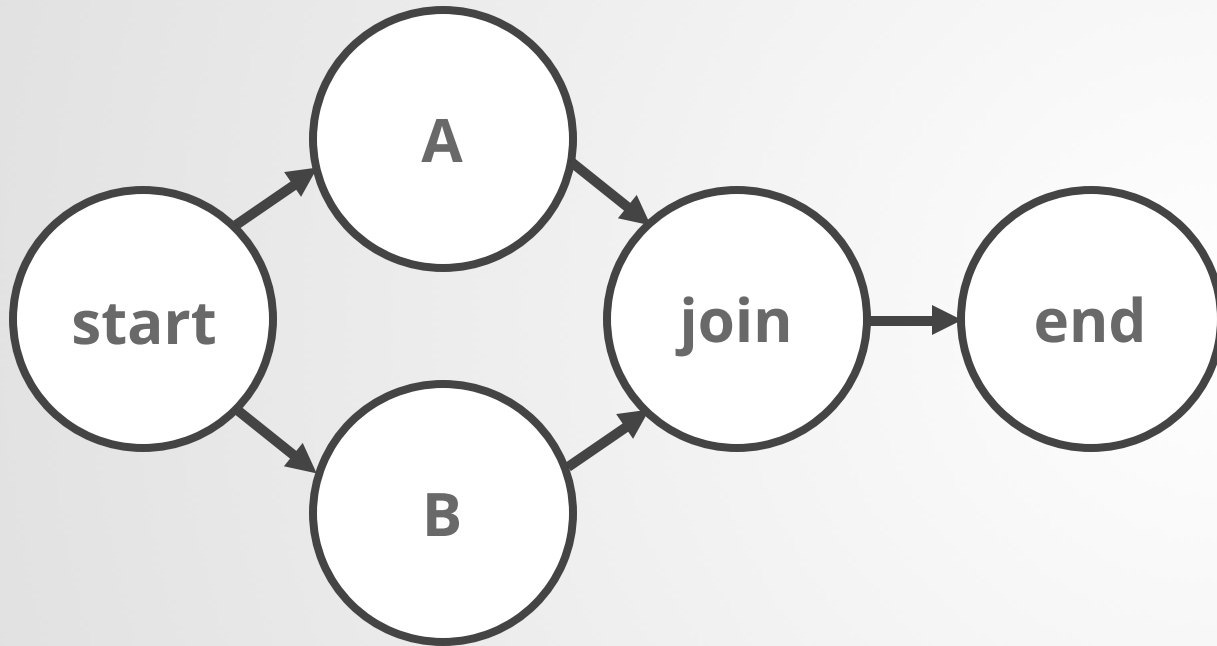# How to deploy my workflow to production?

# How to deploy my workflow to production?



```
#python myscript.py meson create
```

**26%** of projects get deployed to the production scheduler.

How quickly the first deployment happens?

# How to monitor models and examine results?



x=0

x+=2

A

x+=3

B

max(A.x, B.x)

start

join

end

```
In [145]:    1    from metaflow import Flow
             2    run = Flow('MyFlow').latest_run
             3    run

Out[145]:    Run('MyFlow/3')

In [152]:    1    list(run)

Out[152]:    [Step('MyFlow/3/end'),
              Step('MyFlow/3/join'),
              Step('MyFlow/3/b'),
              Step('MyFlow/3/a'),
              Step('MyFlow/3/start')]

In [150]:    1    run['start'].task.data.x

Out[150]:    0

In [151]:    1    run['a'].task.data.x

Out[151]:    2
```

# How to deploy results as a microservice?



```python
from metaflow import WebServiceSpec
from metaflow import endpoint

class MyWebService(WebServiceSpec):

    @endpoint
    def show_data(self, request_dict):
        # TODO: real-time predict here
        result = self.artifacts.flow.x
        return {'result': result}
```

# How to deploy results as a microservice?



```python
from metaflow import WebServiceSpec
from metaflow import endpoint

class MyWebService(WebServiceSpec):

    @endpoint
    def show_data(self, request_dict):
        # TODO: real-time predict here
        result = self.artifacts.flow.x
        return {'result': result}
```

```
# curl http://host/show_data
{"result": 3}
```

# Case Study: Launch Date Schedule Optimization

1. Batch optimize launch date schedules for new titles daily.
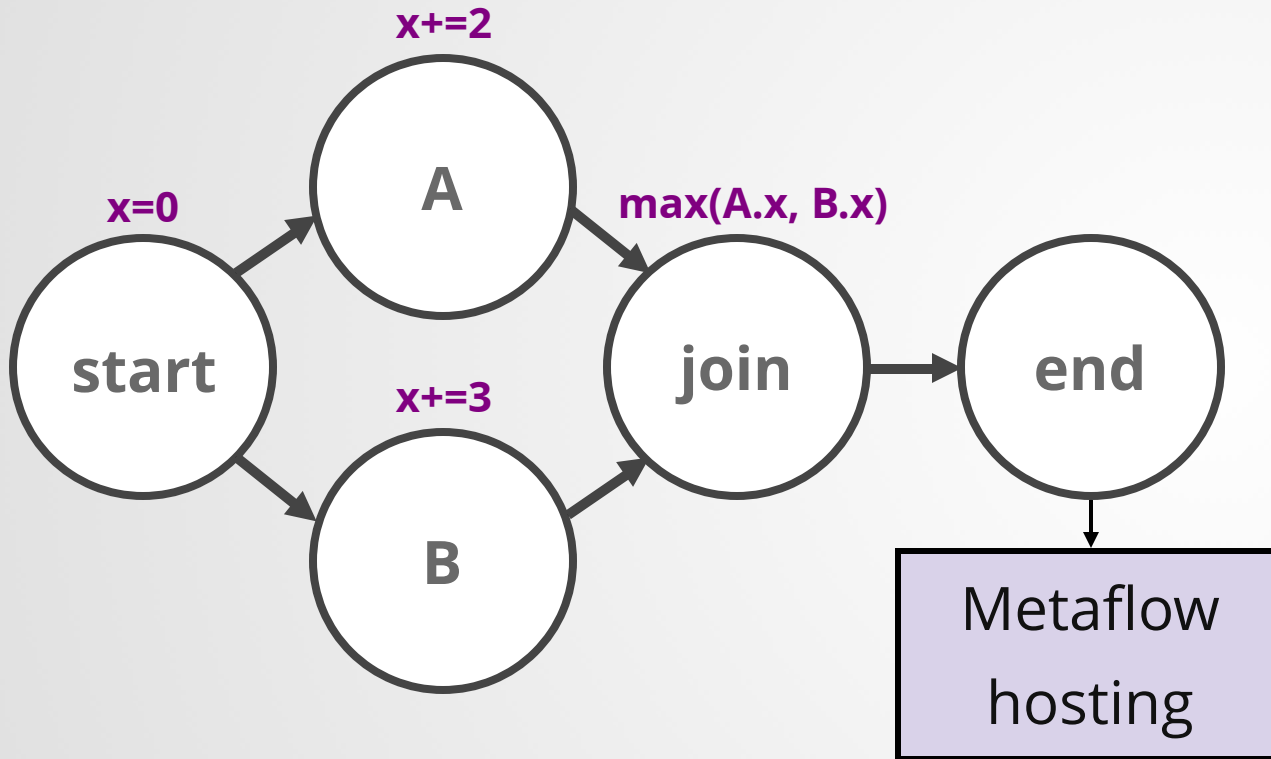   Batch optimization deployed on **Meson.**

2. Serve results through a custom UI.

   Results deployed on **Metaflow Hosting**.

3. Support arbitrary what-if scenarios in the custom UI.

   Run optimizer **in real-time in a custom web endpoint.**

Metaflow

# diverse problems

diverse problems

↓

diverse people

**diverse** problems

$\downarrow$

**diverse people**

$\downarrow$

**diverse models**

diverse problems

↓

diverse people

↓

diverse models

↓

help people build

**diverse** problems

↓

**diverse people**

↓

**diverse models**

↓

**help people build**

↓

**help people deploy**

diverse problems

↓

diverse people

↓

diverse models

↓

help people build

↓

help people deploy

↓

happy people, healthy business

# thank you!

@vtuulos

vtuulos@netflix.com

# Photo Credits

https://www.maxpixel.net/Isolated-Animal-Hundeportrait-Dog-Nature-3234285

Bruno Coldiori

https://www.flickr.com/photos/br1dotcom/8900102170/