VORO

# Sharing ML is harder!

# Clear Processes Help

STREAMBA



Data     Analysis     Documentation     Live

Data Scientist     Engineer

Exploratory     Production

https://www.oreilly.com/ideas/what-is-hardcore-data-science-in-practice

# Clear Processes Help
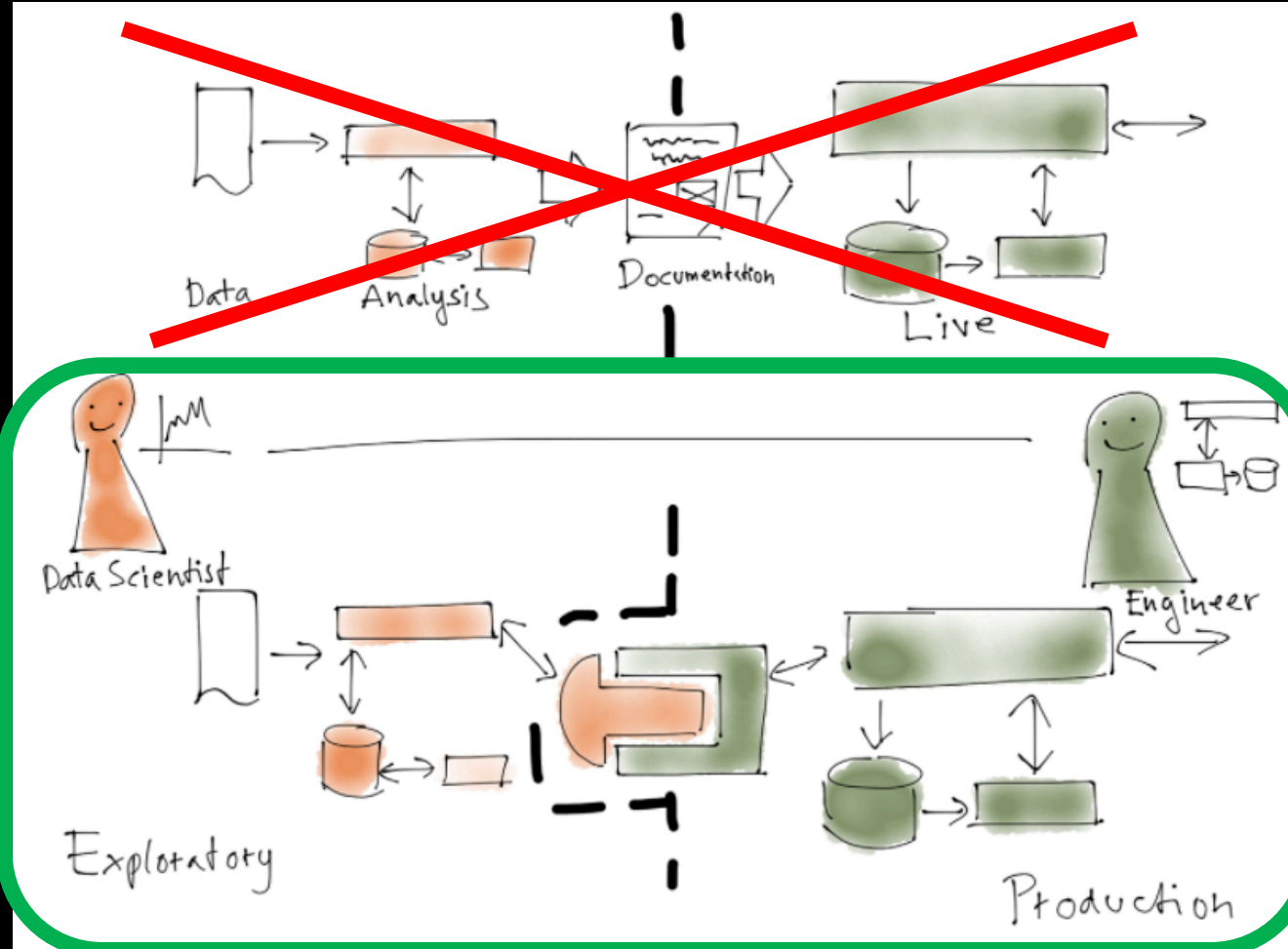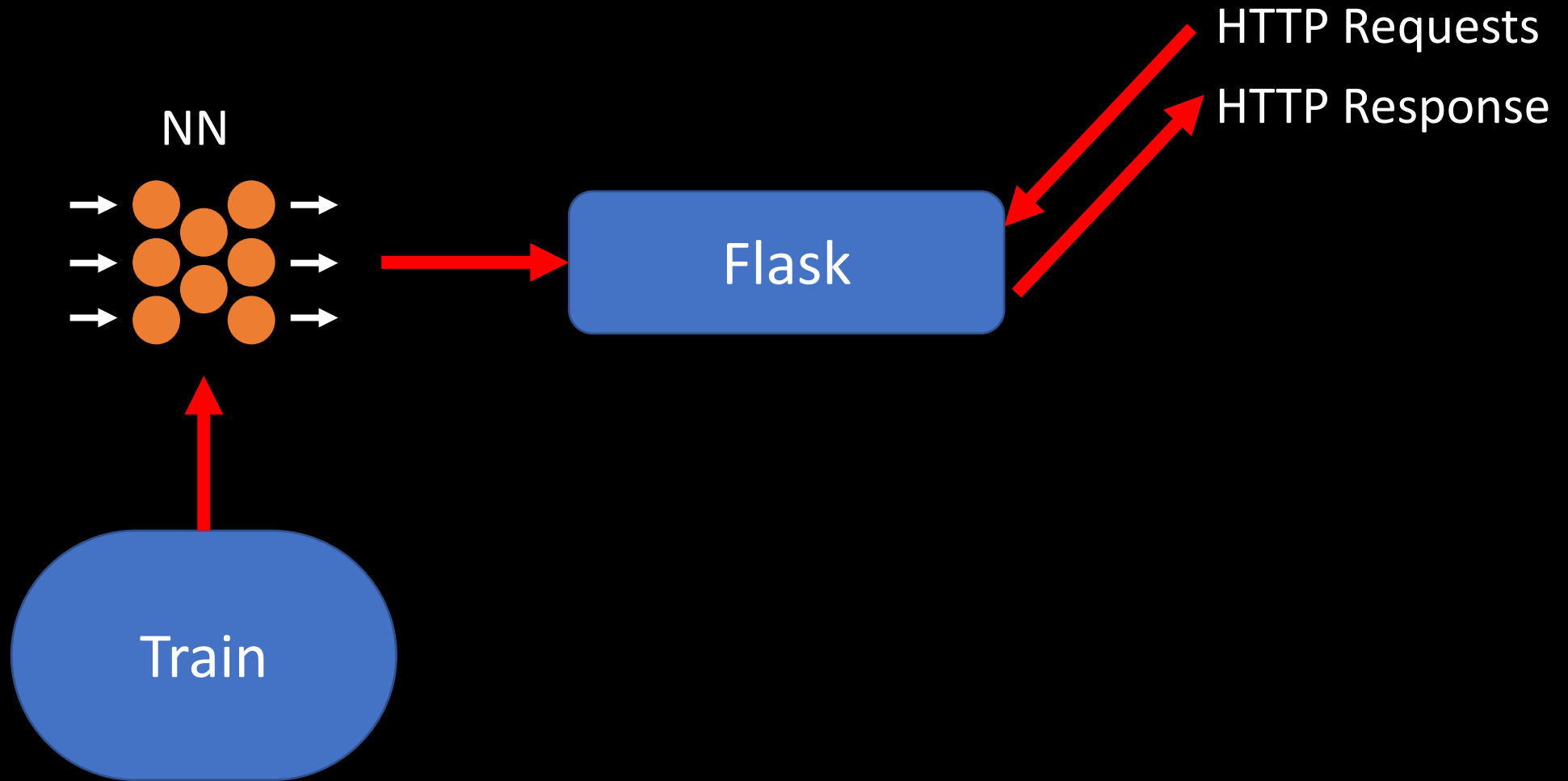
*Cross-collaboration, no siloes.*
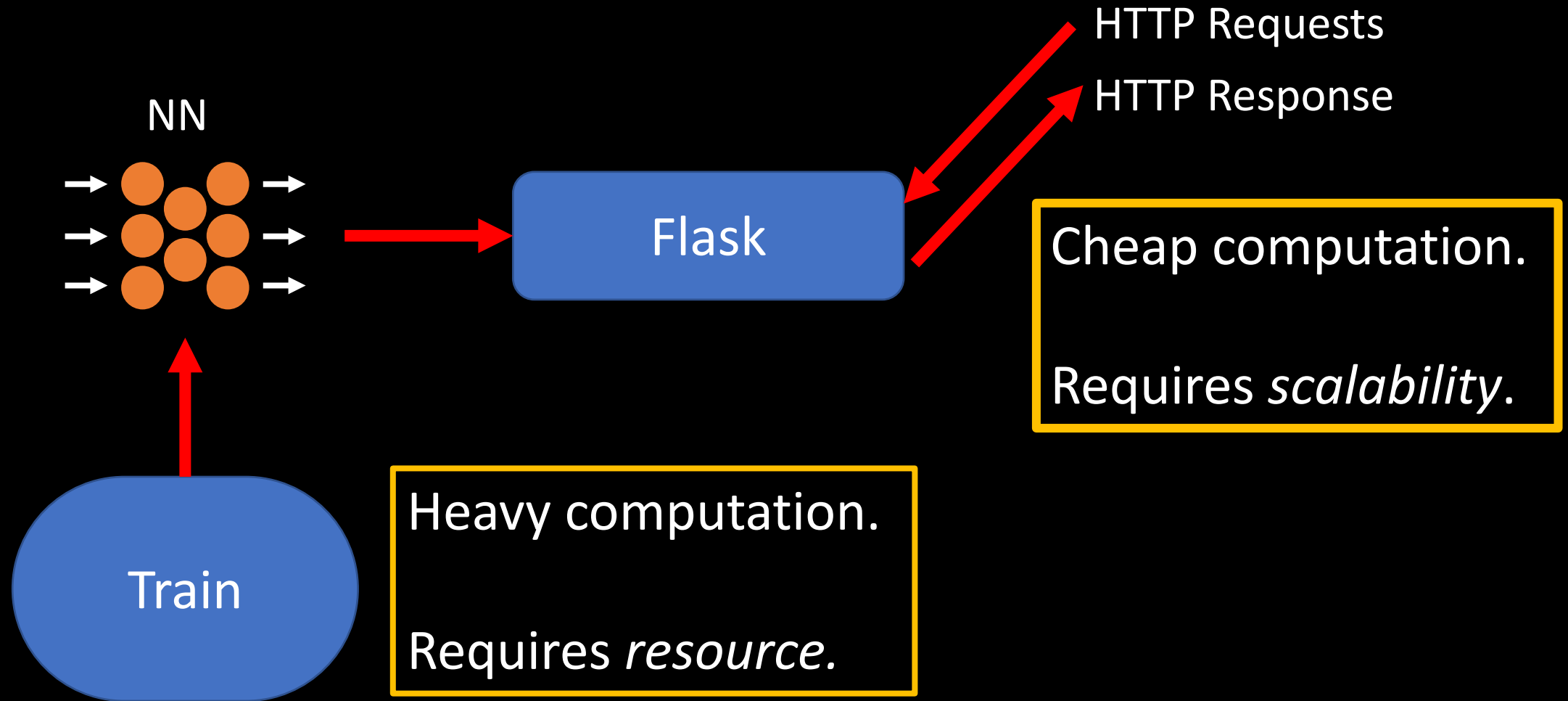
*Data team own their products.*

*Faster and more robust to change.*

*Contracts between ML services and main product(s).*



https://www.oreilly.com/ideas/what-is-hardcore-data-science-in-practice

# ML as a Service

# ML as a Service

STREAMBA

**Train model**

Locally,
Cloud,
On Premises
data centre.

**Save model**

Cloud (S3,
Google Cloud
Storage, Azure
Storage).

**Load model**

**Serve model**

Cloud (Google
App Engine,
Cloud ML
Engine, AWS
Lambda, EC2).

# Cloud ML

- Managed clusters set up for Tensorflow jobs.

- Can train model on cluster / locally.

- Can serve predictions via an endpoint managed by CloudML. Model loaded into Google Cloud Storage to be read.

# Cloud ML



- Managed clusters set up for Tensorflow jobs.

- Can train model on cluster / locally.

- Can serve predictions via an endpoint managed by CloudML. Model loaded into Google Cloud Storage to be read.

# STREAMBA

# Google App Engine



- Managed VM infrastructure.

- Supported languages and Docker containers.

- Built to scale.

- Very easy to use but you may want even more control (Compute Engine / EC2).

# Google App Engine



- Managed VM infrastructure.

- Supported languages and Docker containers.

- Built to scale.

- Very easy to use but you may want even more control (Compute Engine / EC2).

# Conclusions

STREAMBA

- It is important you know how to share your ML work.

- Learn the tools – just pick one and start playing.

- GCP has options specifically for Tensorflow/Scikit-learn and more general options. AWS and Azure have great tools too.

- <u>Don't fall into the trap of thinking "that's the engineer's job", it's *your job!*</u>