



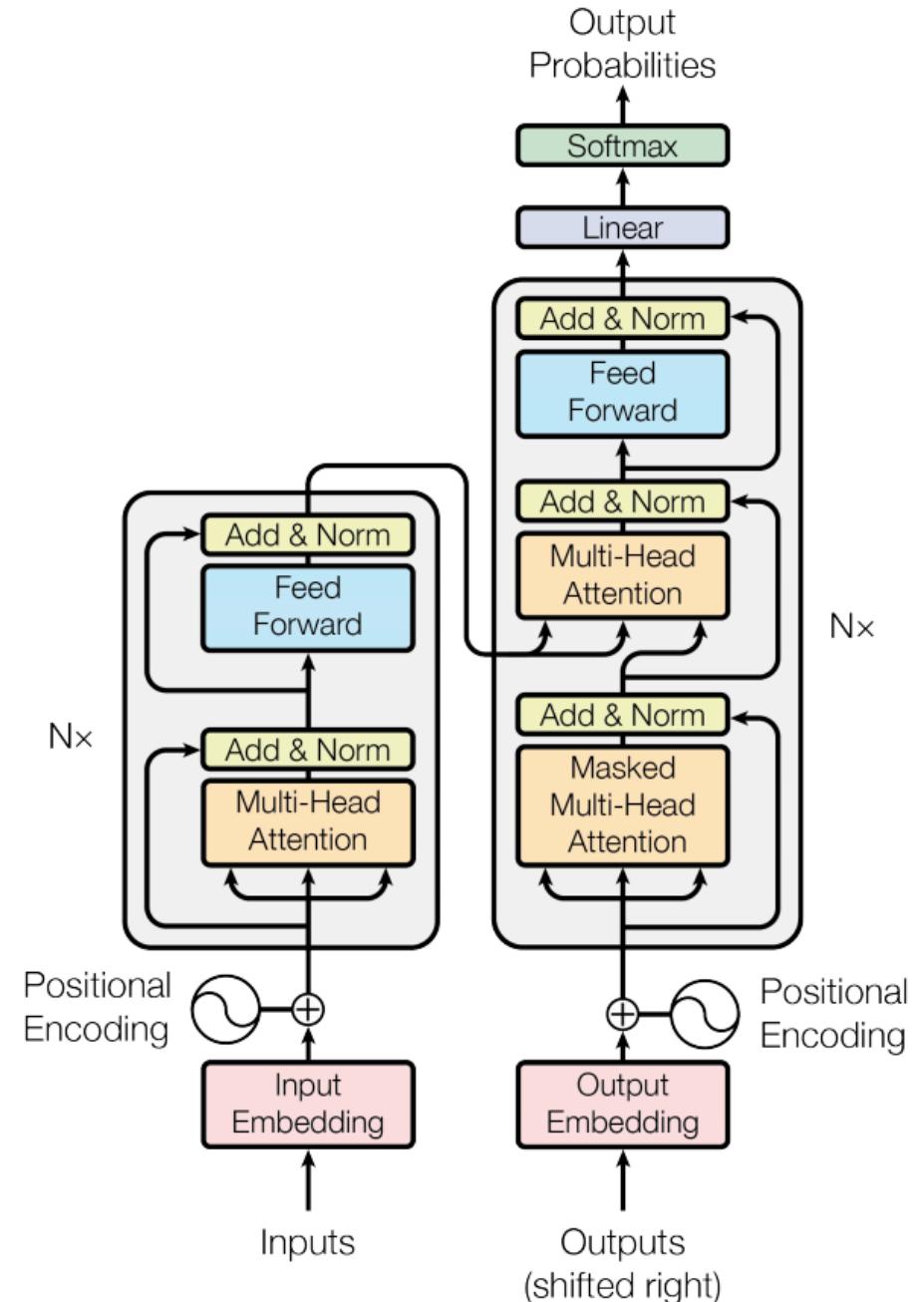
Automating machine learning pipelines with AWS Lambda

Christopher Grainger

February 2019

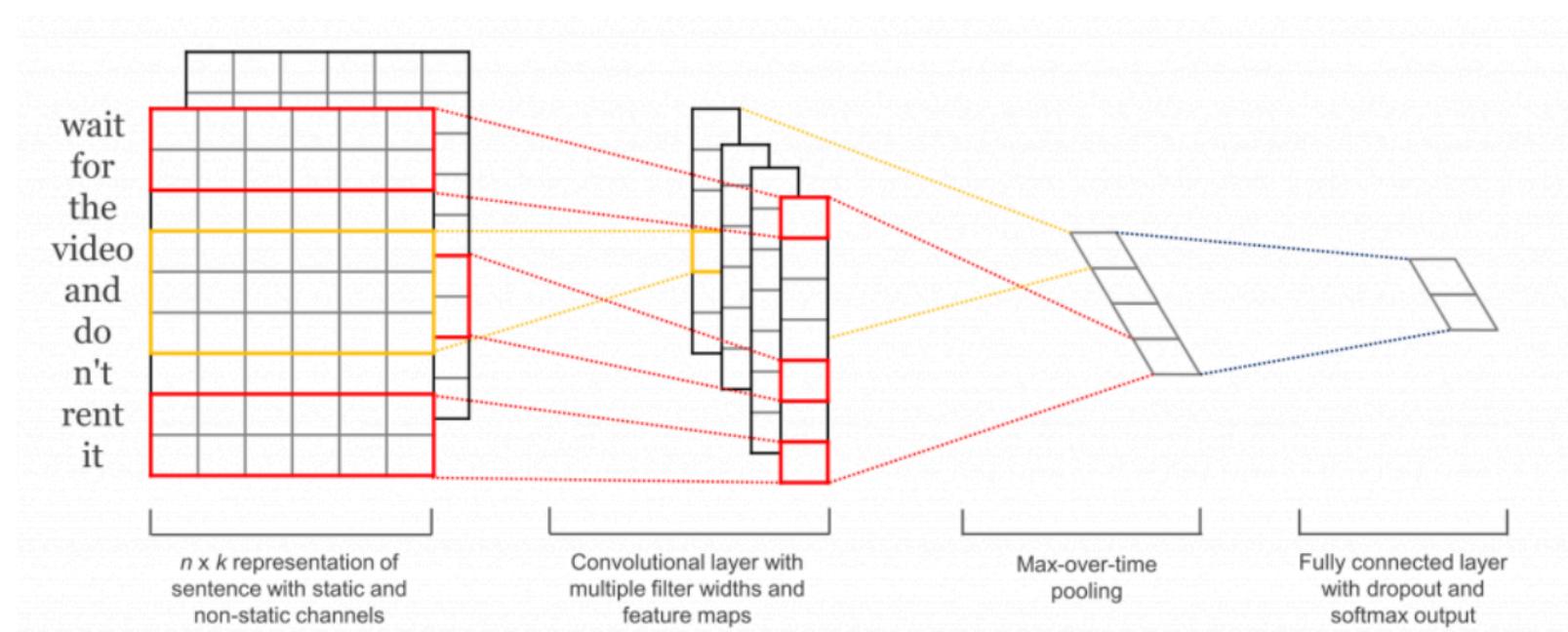
Today's Topics

- Going from research to production with ML
- Two problems: big models and big data
- Why Lambda?
- What we're doing now

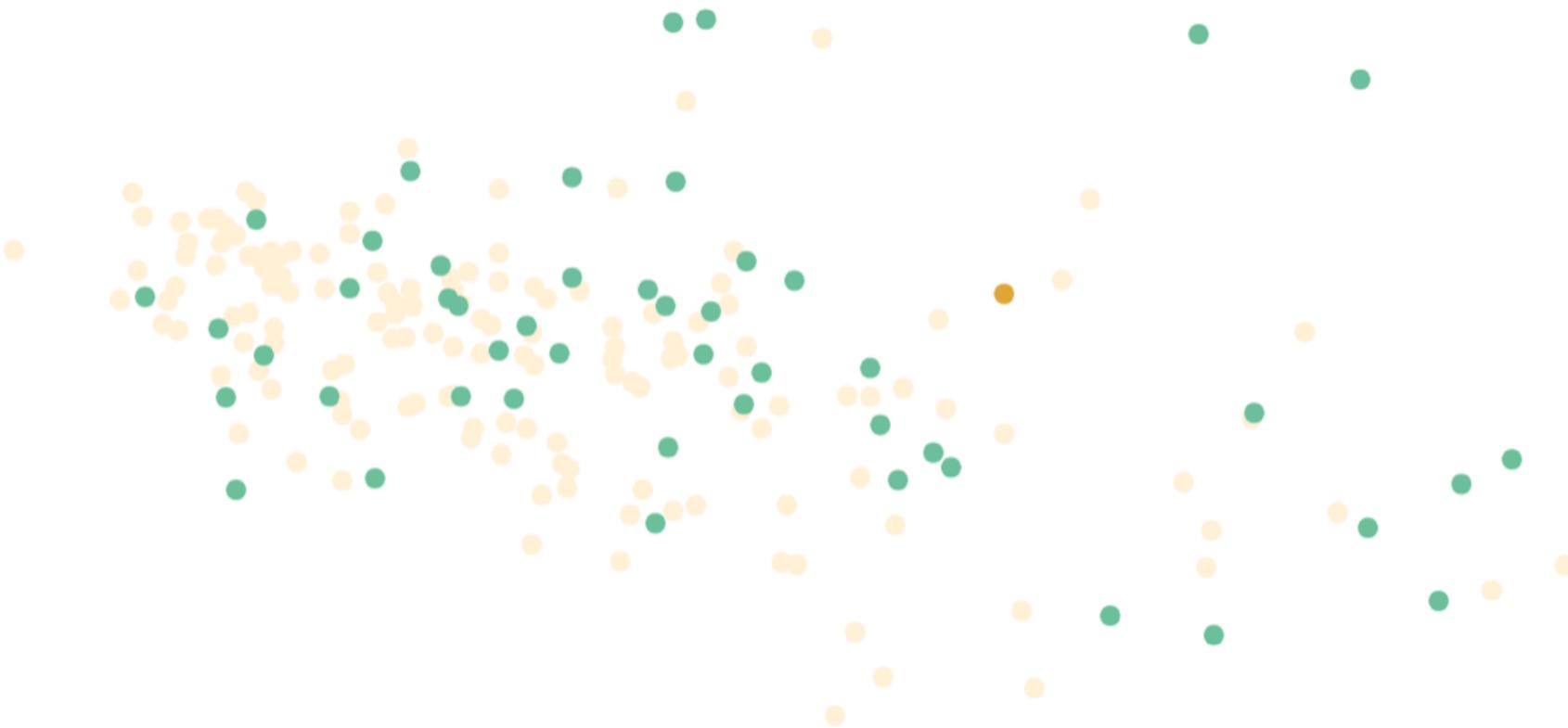


Okay, so what are we researching?

- Dense representations of documents

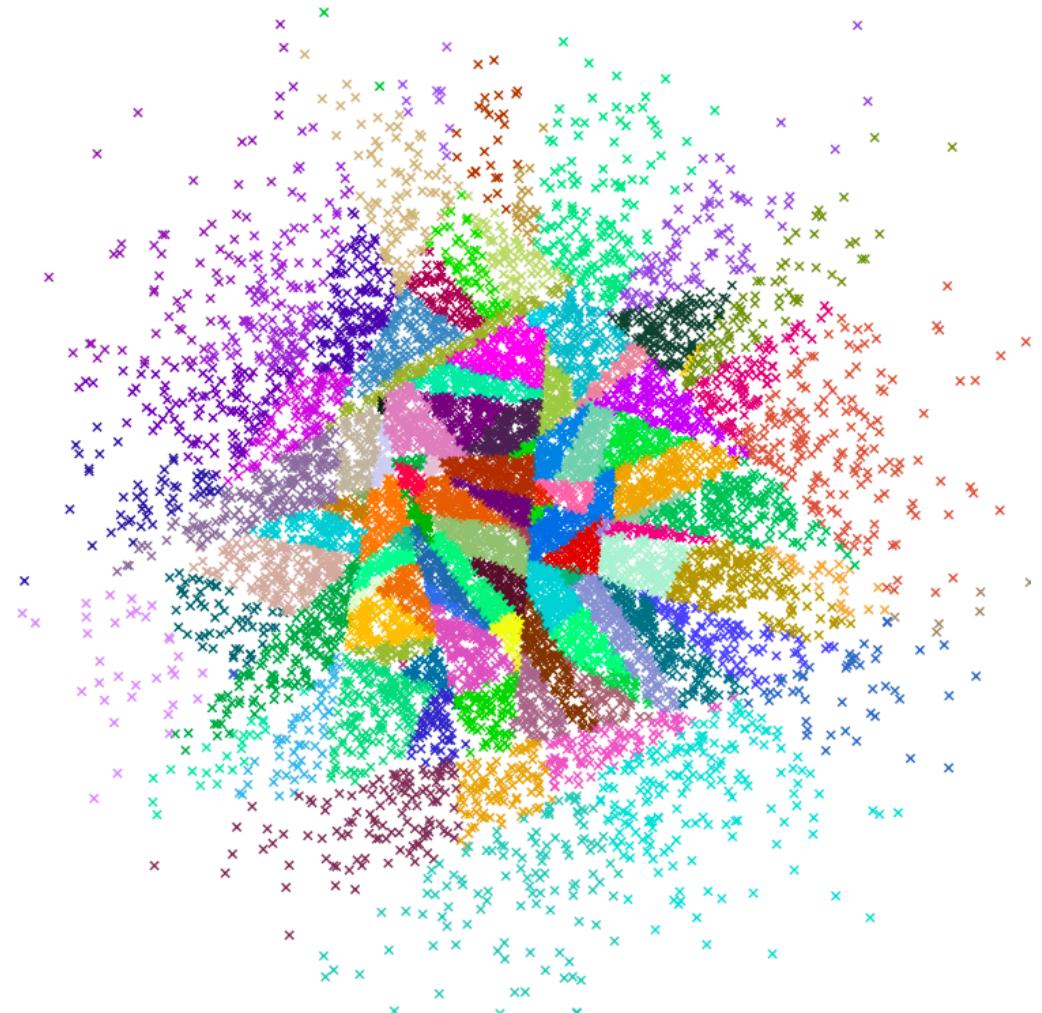


Okay, so what are we researching?

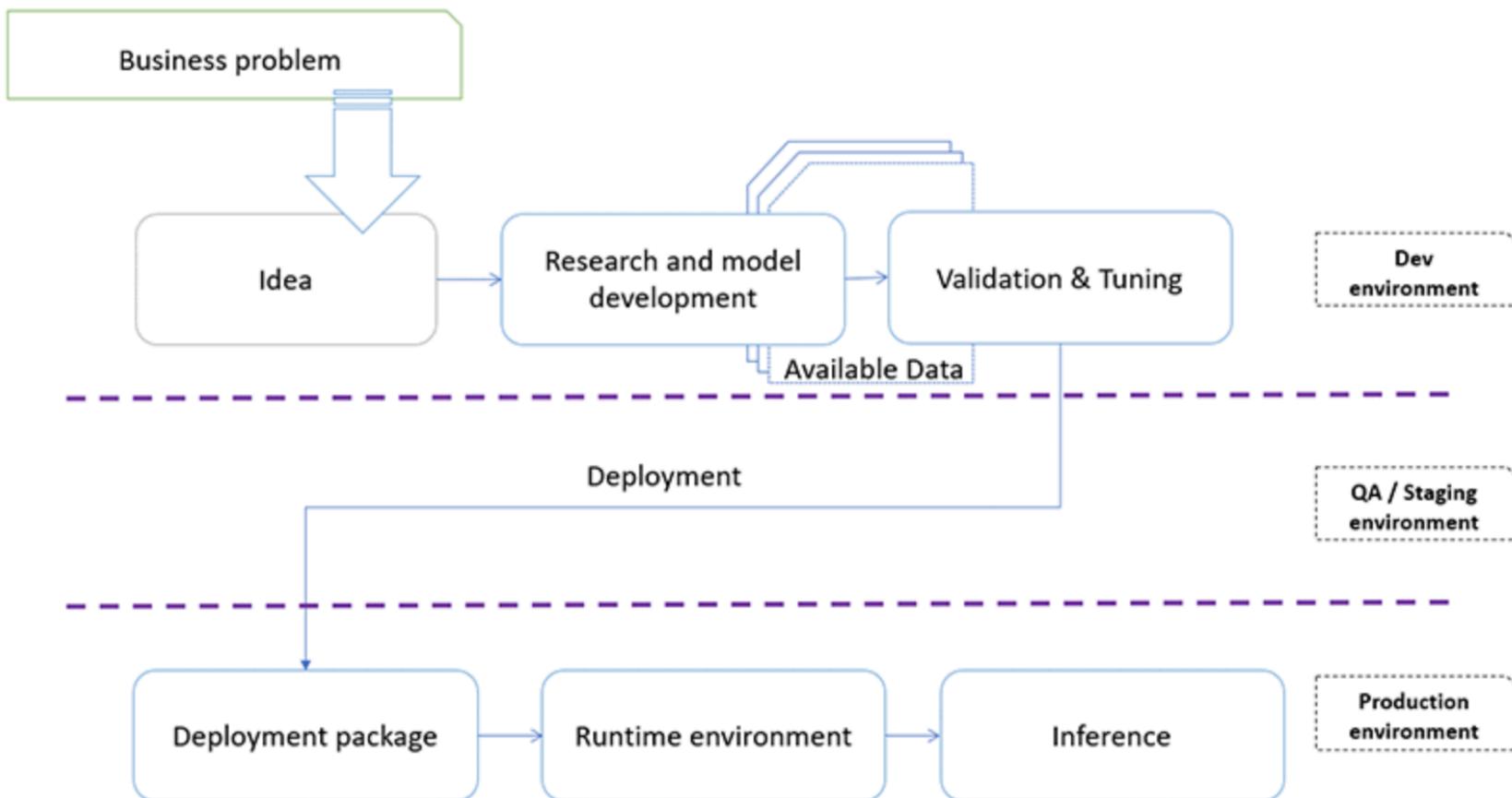


And how are we using it?

- Simple approximate nearest neighbours
- Embedding documents on the fly
- Visualising in two-dimensional space



An ideal world



But I can train ImageNet on 512 GPUs in 2.5 minutes!

Not so fast

- Efficient use of GPUs for NLP is getting much better, but still way slow
- Models take days to run a single epoch on 8xV100s.
- Concept is still the same: we want to cut models and put them in S3 manually, then trigger inference to deployment



What do we need?

A database

Somewhere to hold the documents, metadata, representations, and dim-reduced representations.

Training instance

We need somewhere to train our ML models.

API

Other applications need to interact with the models and data.

Model storage

Trained models need to go somewhere.

Cutting-edge ML in six easy steps (it's a piece of cake)

- Train
- Test (and benchmark)
- Version
- Deploy
- Publish new model
- Repeat

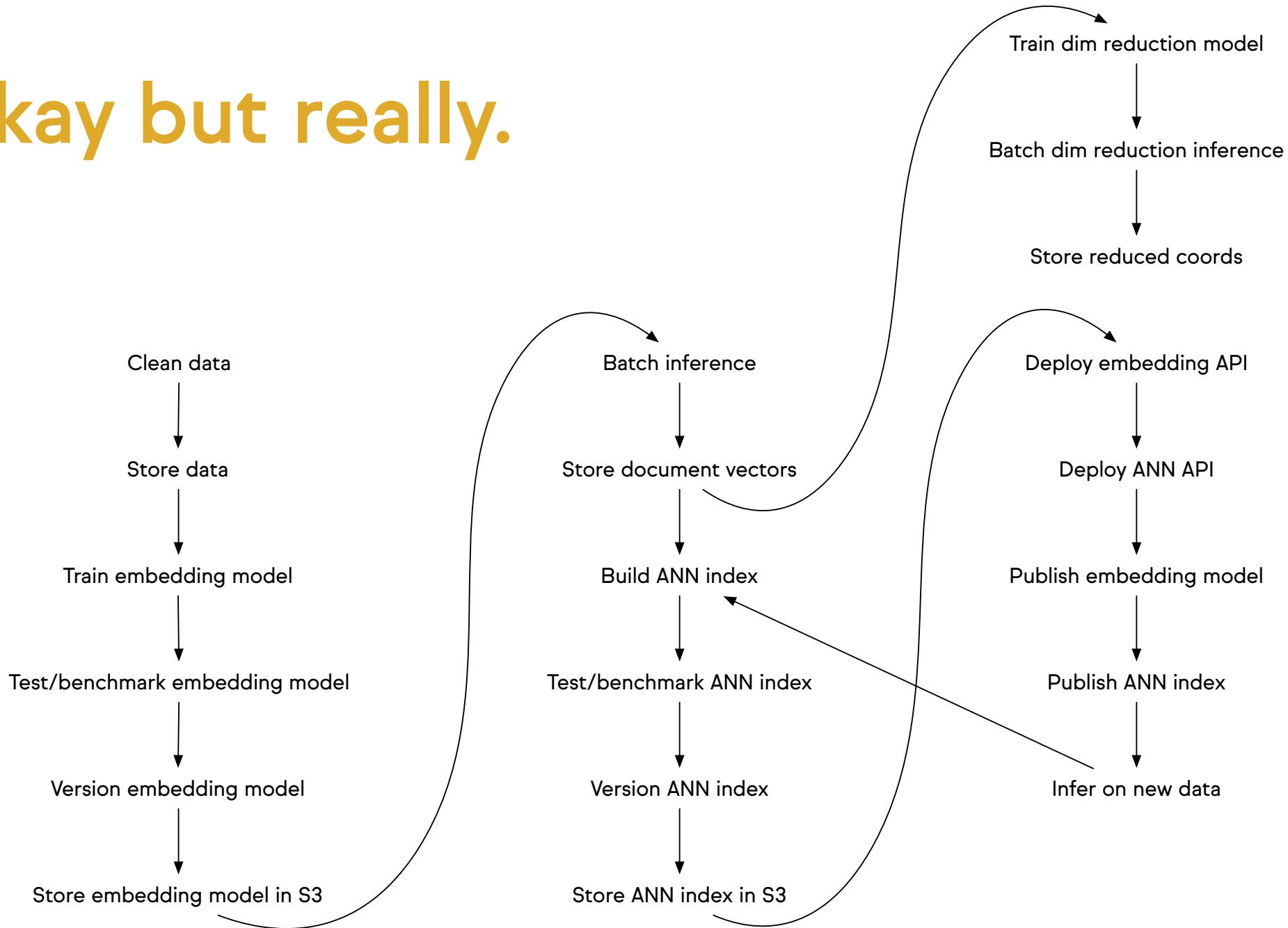


What do we really need?

- A data warehouse
- A production database
- A data lake
- ETL for data prep
- Deep learning training instances
- Deep learning inference instances for batch/streaming
- ANN index building instance
- PCA training instance
- Benchmarking
- RESTful API for document embedding (with GPU)
- RESTful API for nearest neighbours (without GPU)
- Secrets management
- Autoscaling
- CI/CD for models and code

What do we actually really need?

Okay but really.



Why Lambda?

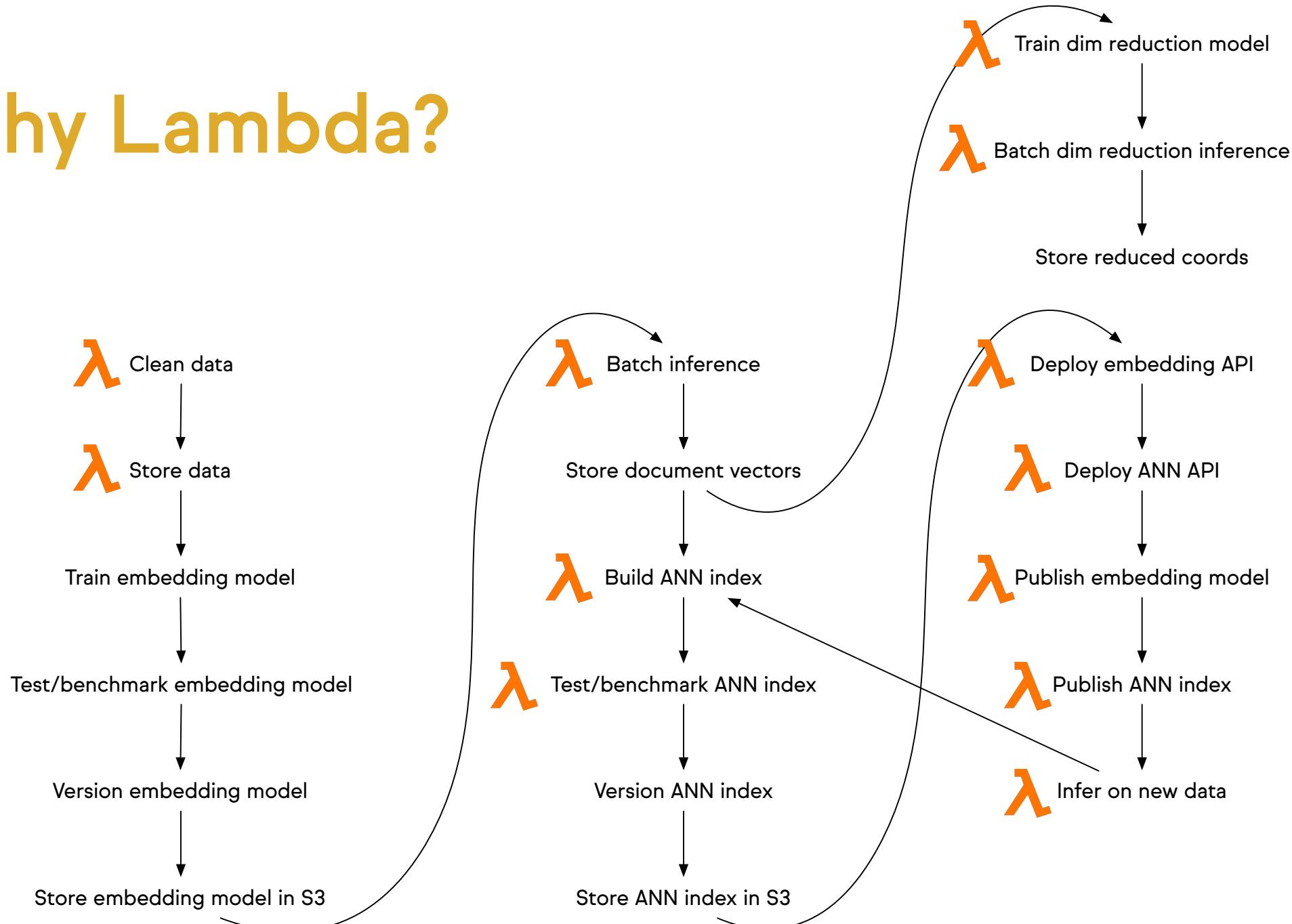
- Why not:
 - Ansible?
 - CloudFormation?
 - Jenkins?
 - Airflow?
 - Whatever else?



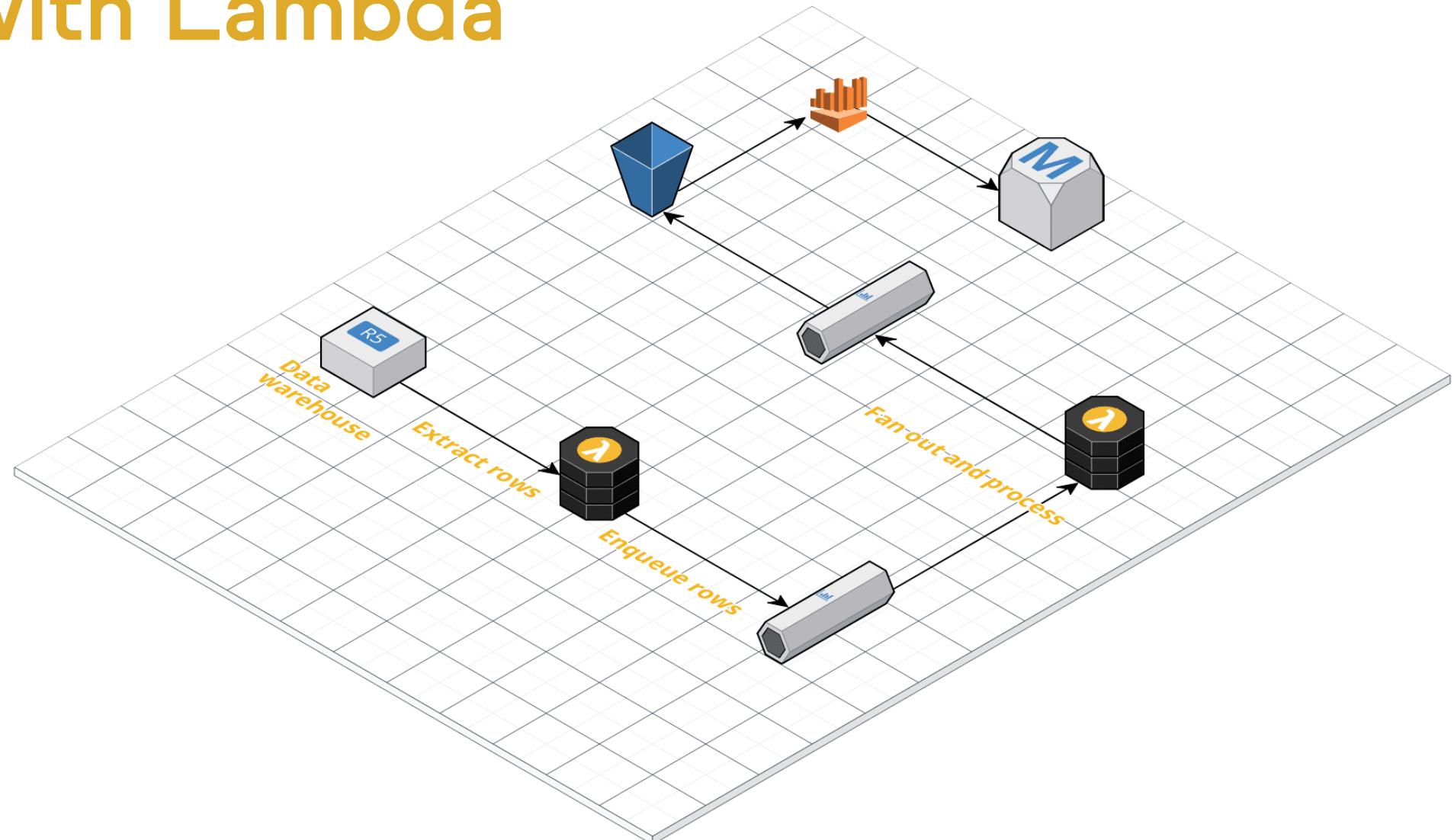
CloudFormation



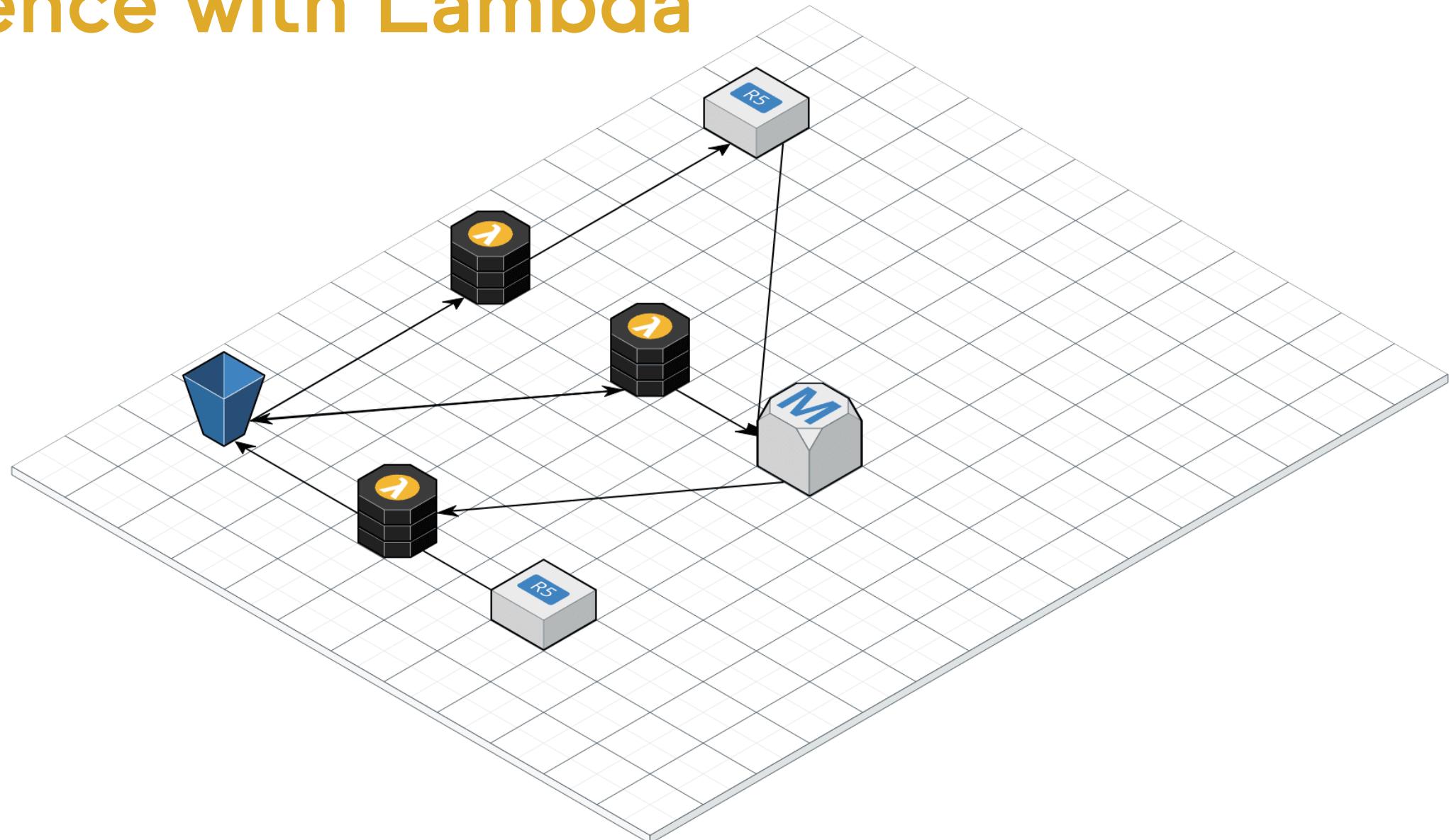
Why Lambda?



ETL with Lambda



Inference with Lambda





Thank you! Questions?

Also: we're hiring!

[careers@amplified.ai](mailto:ccareers@amplified.ai)