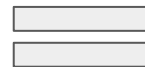
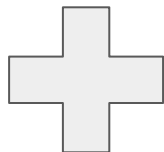




Tensor in the Sky with CloudML





Sylvain Lequeux



@slequeux



Romain Sagean



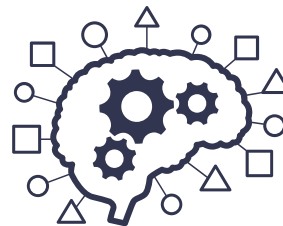
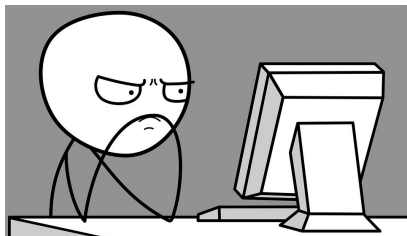
@scauglog

**Xebia** *Data Factory*

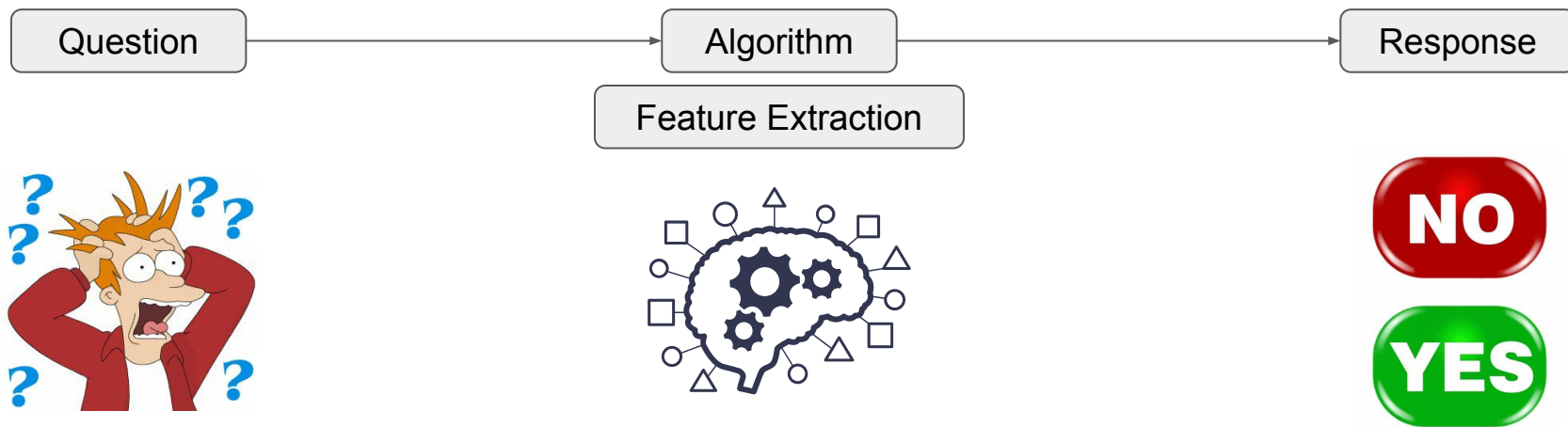
# WHY ?

## Deep Learning

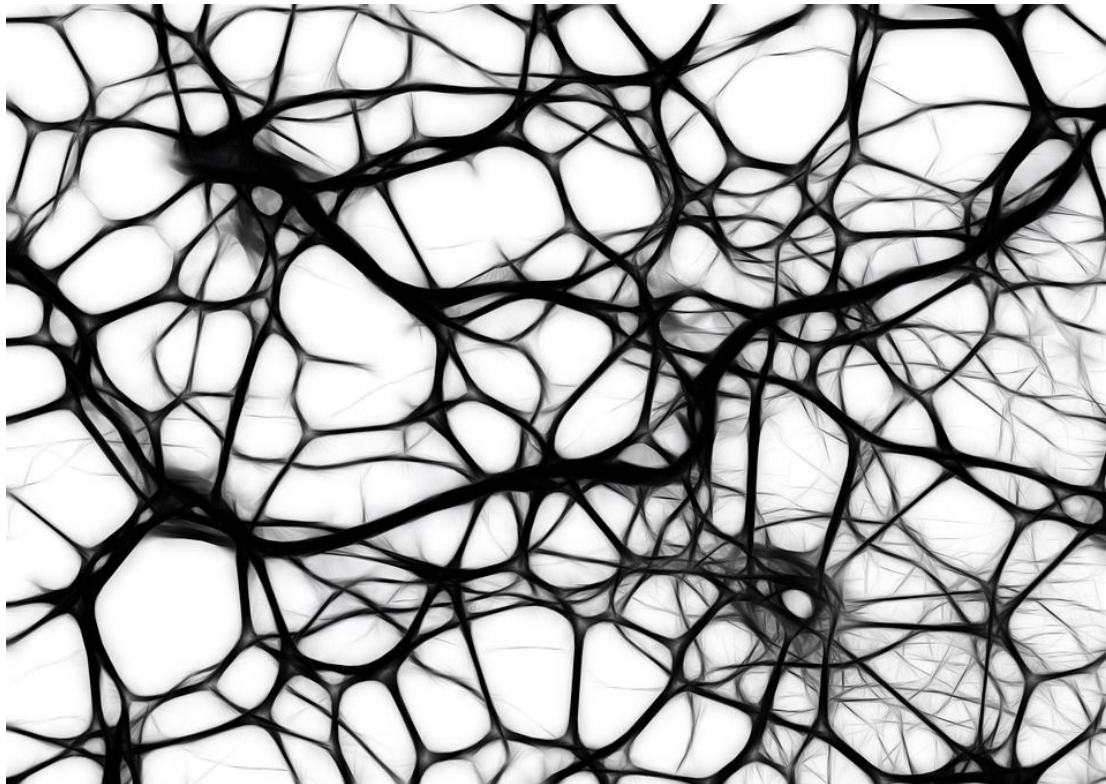
# Classic Machine Learning



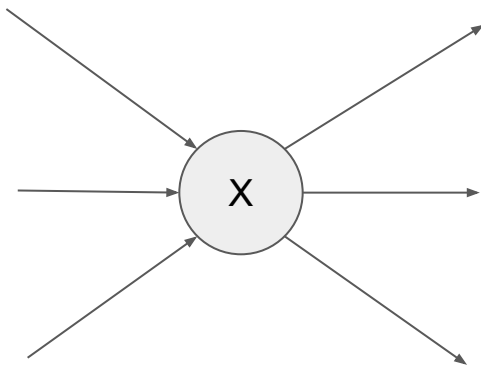
# Deep Learning



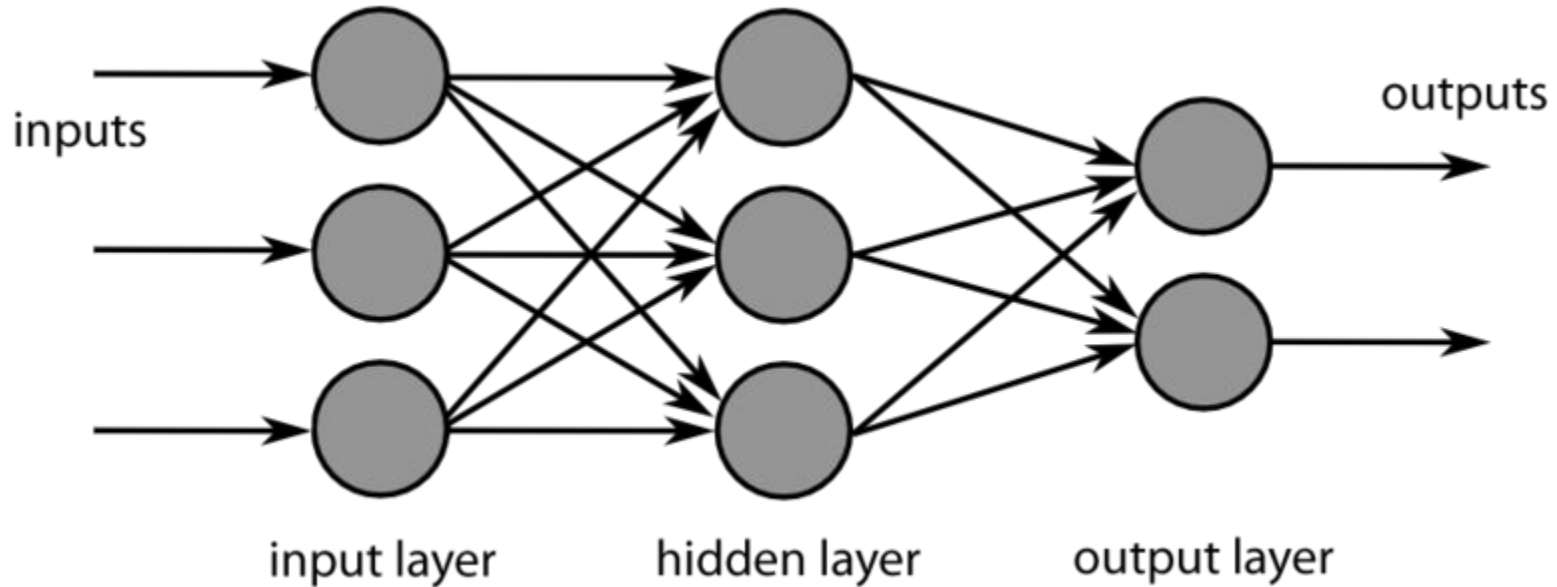
# | Neural Networks



# | Artificial Neural Networks



# Artificial Neural Networks





# WHY ?

Cloud

# | Google CloudML

Machine Learning Managed Service

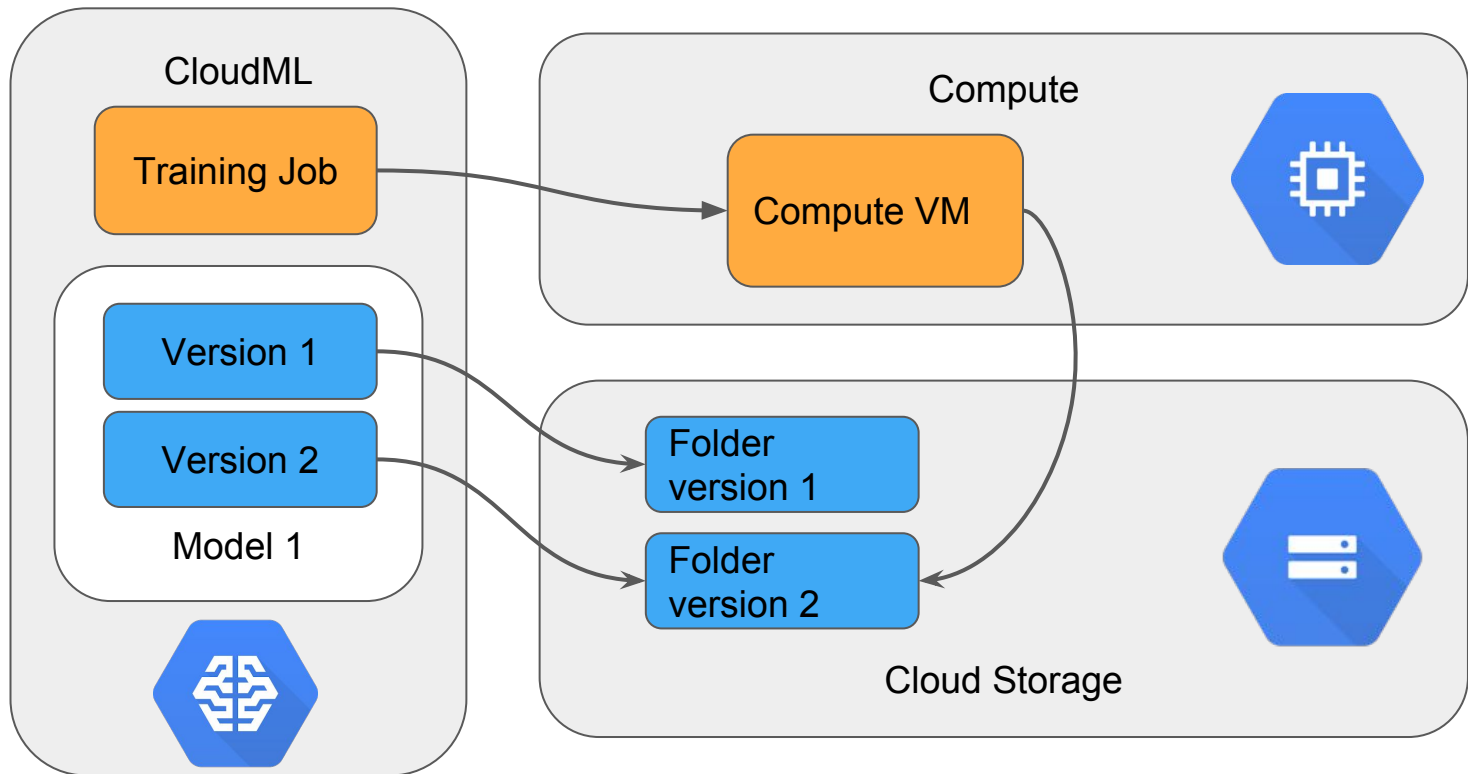
- ▼ ML Frameworks (beta) : scikit-learn, XGBoost
- ▼ DL Frameworks : Tensorflow, Keras

What you can do

- ▼ Train
- ▼ Predict



# Google CloudML



# WHY ?

## Use Case

# Heat the Neurons of Your Smartphone with Deep Learning

Qian Jin | @qianjin | qianjin@xebicon.fr  
Yohann Benoit | @yohannbenoit | yohannbenoit@xebicon.fr  
Sylvain Lequesux | @sylvainlequesux | sylvainlequesux@xebicon.fr

Xebicon

Xebicon TV

Xebicon'17



# Local Mode

At least, it'll work

Included into Tensorflow & Keras

There are pre-trained models

- ▼ Save time for training

Param count approximation

- ▼ InceptionV3 : 23M
- ▼ VGG16 : 138M

Adapted architecture

- ▼ Convolutional
- ▼ MaxPooling
- ▼ Recurrent
- ▼ ...



# TIME TO DEMO



# LOCAL PREDICTIONS

# Export model

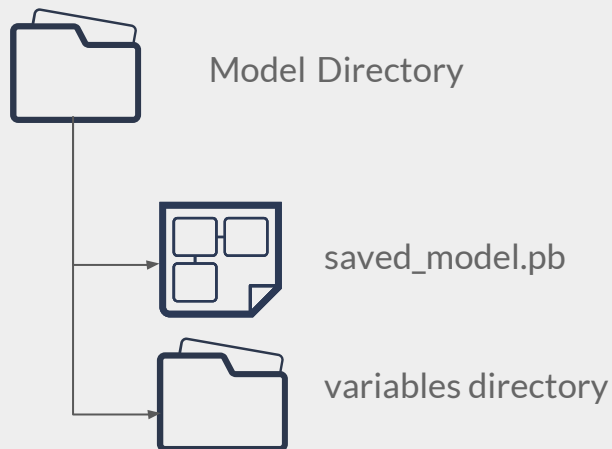
## Or How to Reuse

Weight stored in-memory

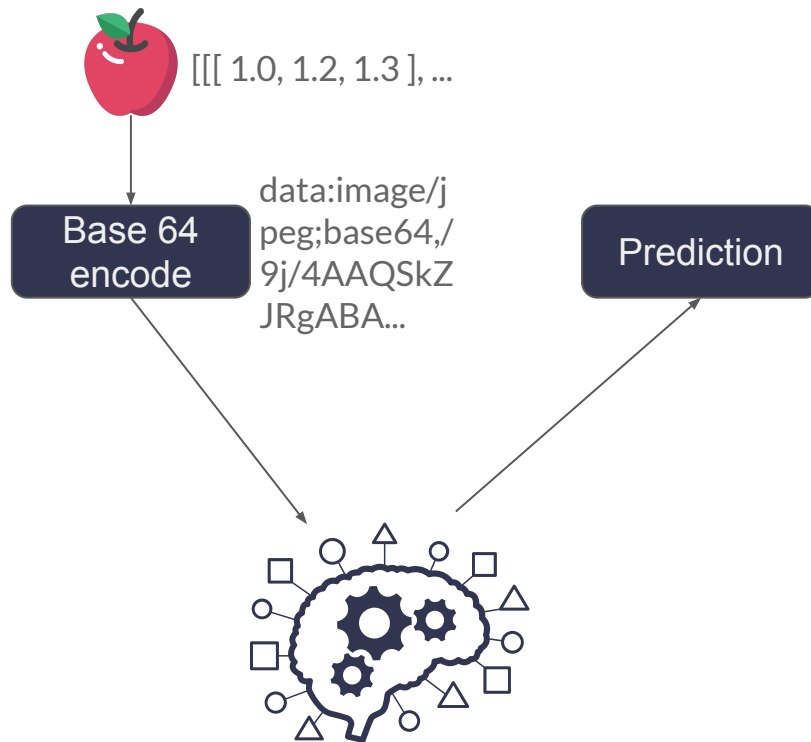
- ▼ Need to retrain each time you start
- ▼ Non-deterministic behaviour

Use Protobuf format to store model

- ▼ Google's data exchange format
- ▼ Stores weights and architecture



## Export model

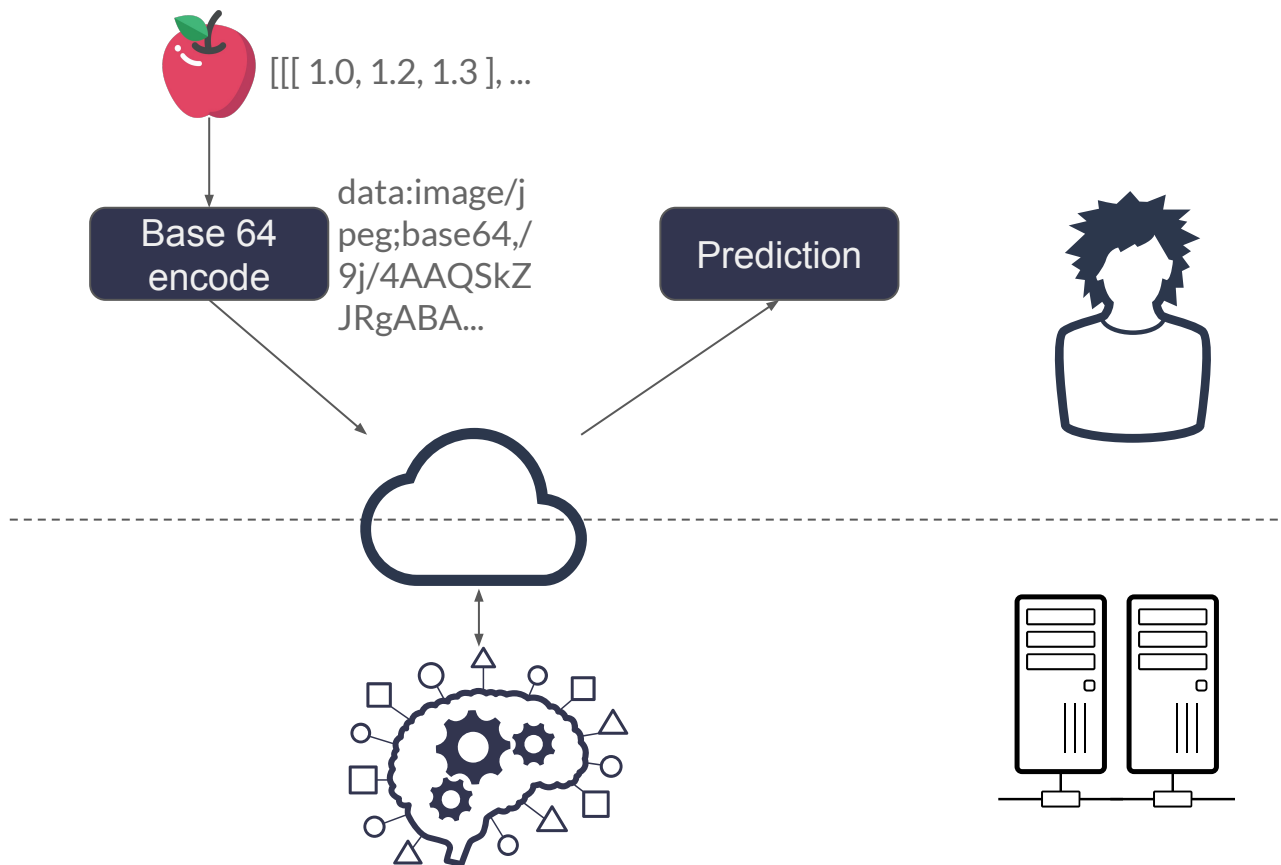




# Serving In the Cloud

No diamonds, sorry

# Serving





# Google functions

## Proxify HTTP requests

Perform controls :

- ▼ Authentication
- ▼ Bandwidth controls
- ▼ ...

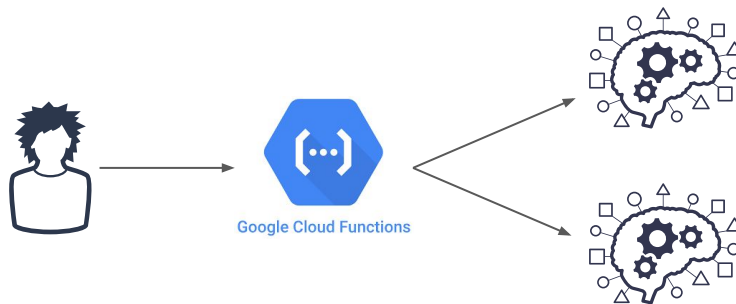
Transform data

- ▼ Input preprocessing
- ▼ Output formatting

# Google functions

## A/B test your models

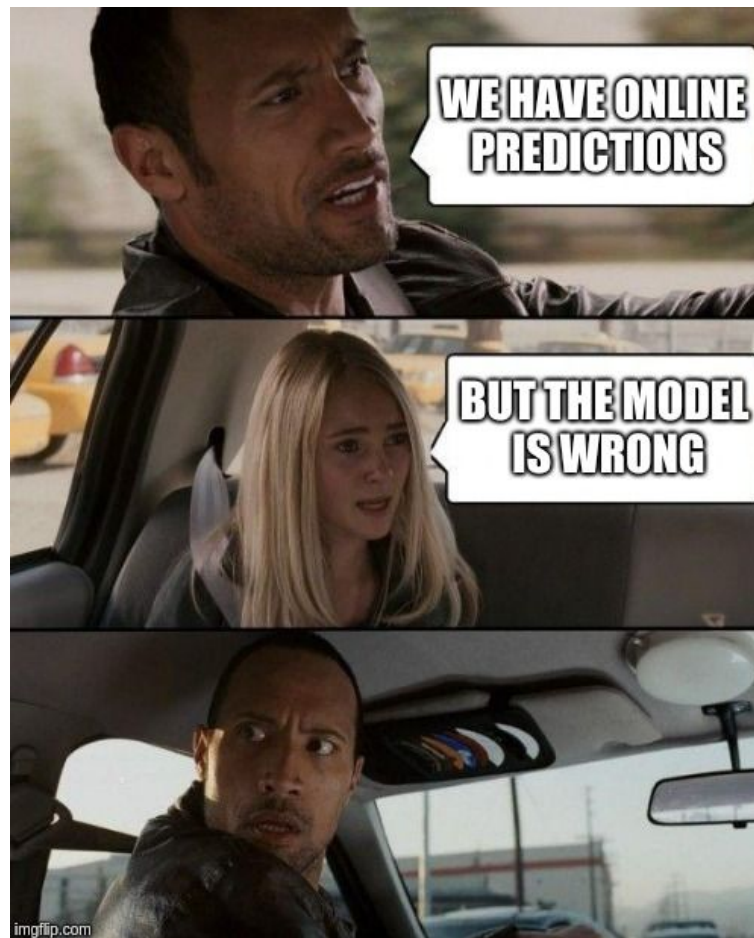
Implements A/B testing logic inside this function





# Learning In the Cloud

How about Magritte ?



**Transfer learning** [...] is a research problem in machine learning that focuses on storing knowledge gained while solving one problem and applying it to a different but related problem.



# Transfer learning

## Retraining without retraining

First layers learn basic features

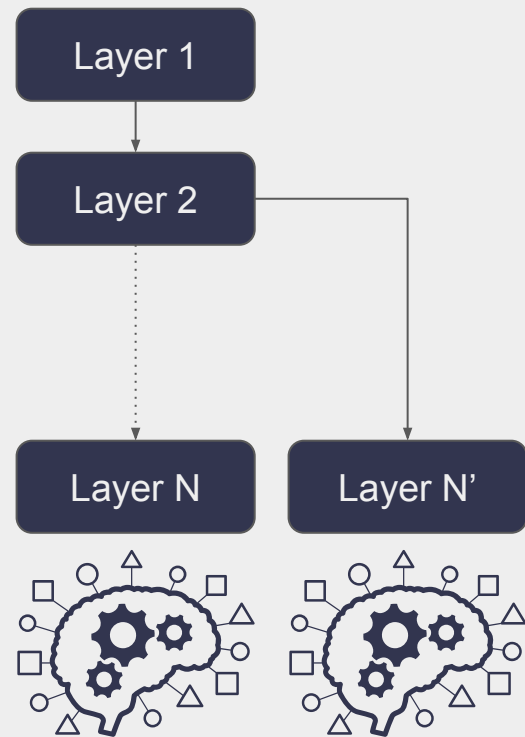
- ▼ Lines
- ▼ Colors
- ▼ ...

End layers learn complex features

- ▼ Objects
- ▼ Complex shapes
- ▼ ...

Transfer learning is 2-steps :

- ▼ Freeze first layers
- ▼ Train new end layers



# Online training

## Infinite power

- ▼ Define a setup.py file
- ▼ Define a config.yml file
  - ▼ This could be done directly in step 3 by passing args
  - ▼ Defines CPU/GPU/TPU usage
  - ▼ Defines runtime version
  - ▼ ...
- ▼ Run it
  - ▼ Using gcloud command line
  - ▼ Use -- to separate gcloud args from your job's args

```
(xke-cloudml) root@b4dc5ee1316c:/opt/demo# gcloud ml-engine jobs submit training
$JOB_NAME \
  --job-dir $OUTPUT_PATH \
  --runtime-version 1.5 \
  --module-name trainer.task \
  --package-path trainer/ \
  --region europe-west1 \
  --scale-tier basic-gpu \
  -- \
  --num-epochs 3 \
  --steps-per-epoch 3 \
  --data-project $PROJECT_ID \
  --data-bucket $BUCKET \
  --data-path $DATA_PATH
```



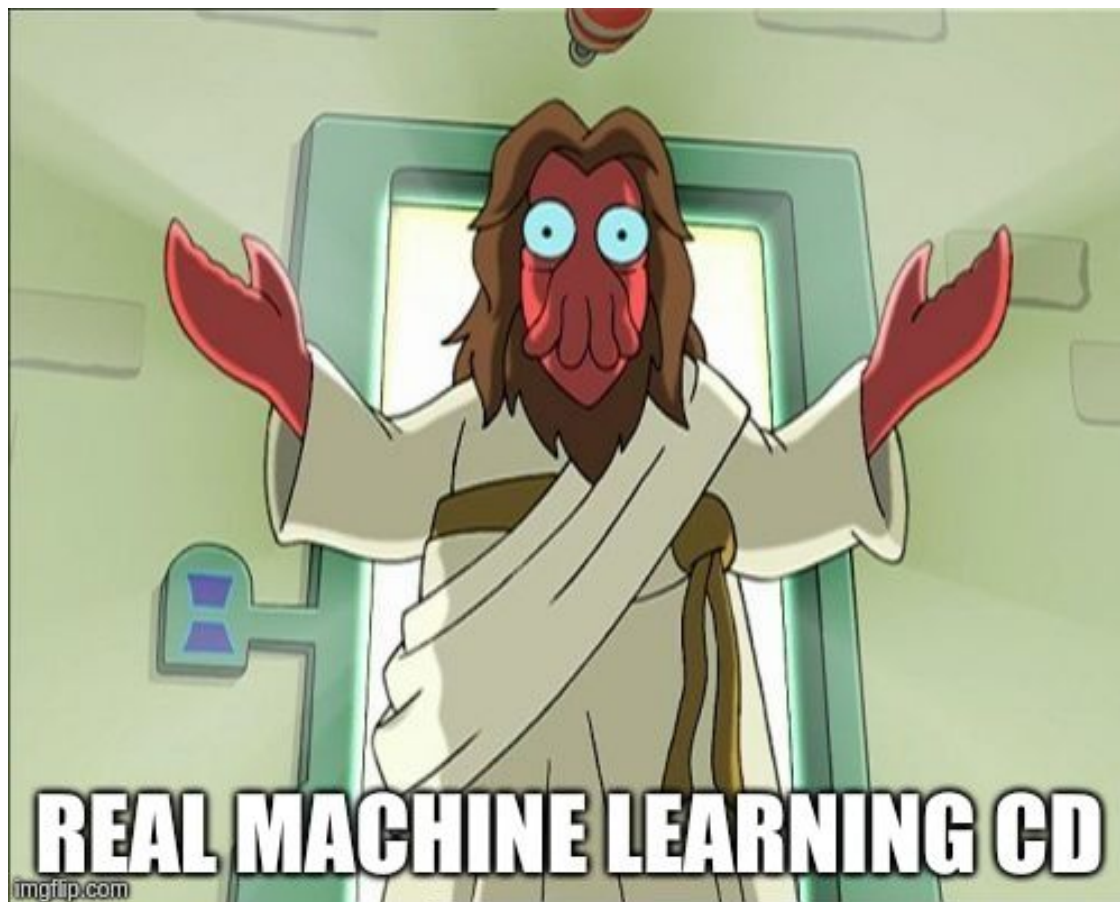
# So, what ?

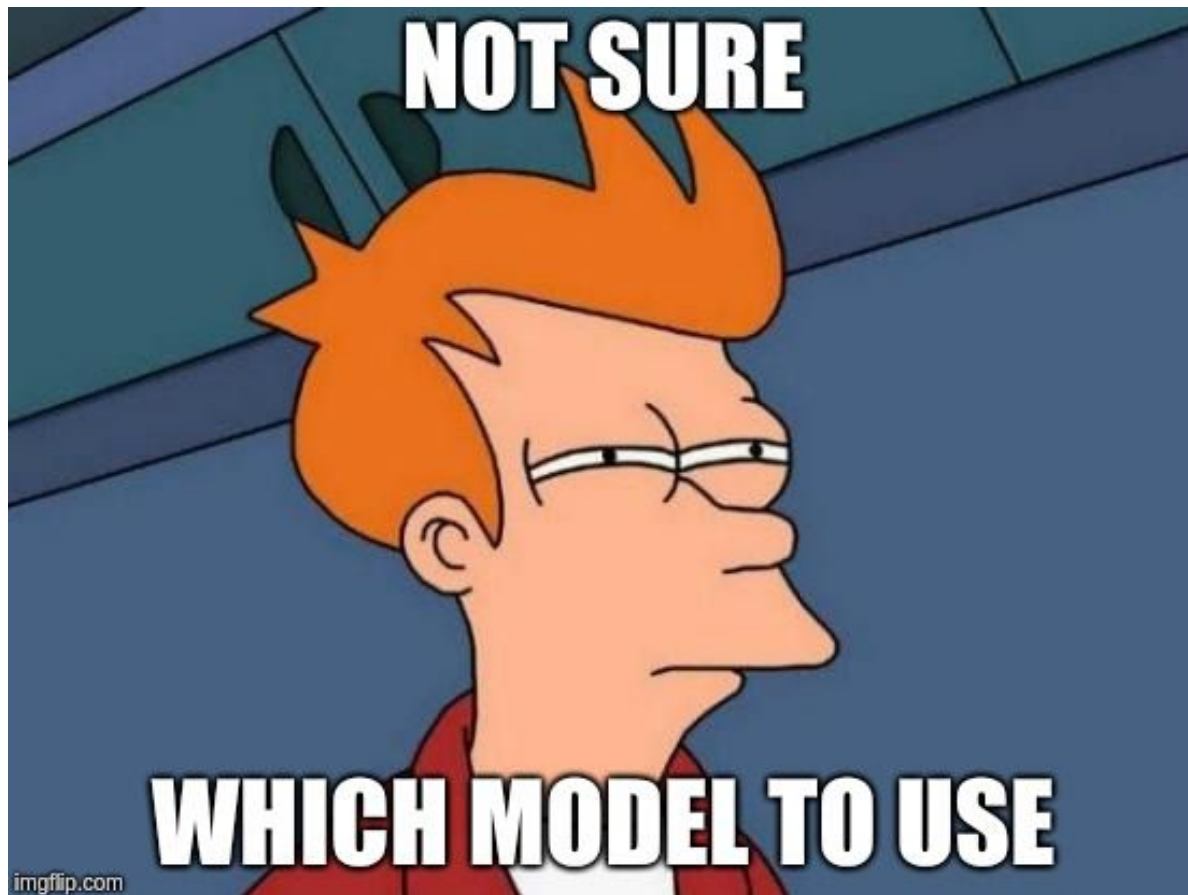
Conclusions : 2 for 1











# Thank you

<https://github.com/slequeux/xke-cloudml>