# Big data, machine learning and maps: lessons learned on aerial imagery

Olivia Wilson
Research Scientist

Olivia.Wilson@os.uk
@oew1v07

Ordnance Survey

# OS is known for

▶ **Accurate, authoritative data**

When it needs to be right

▶ **Trusted and respected worldwide**

A domestic focus with international reach

▶ **225 years experience**
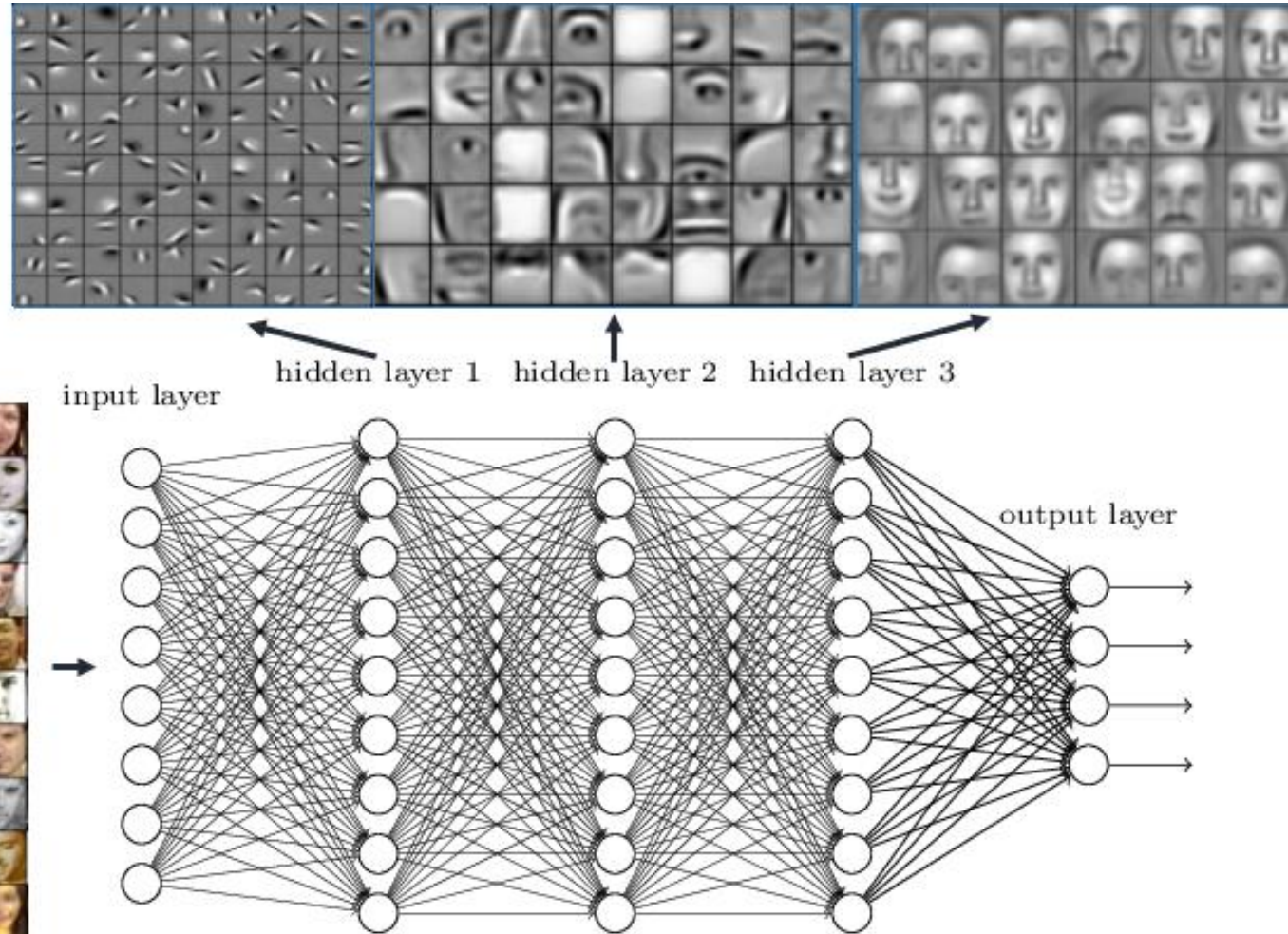The world's most experienced geospatial intelligence organisation.



Ordnance Survey

# What OS actually does

► Paper maps now a legacy – less than 6% of our business but still hugely significant

► One of the UK's largest stores of geographic digital data

► Making **100,000+** updates every day to over **650+ million** physical features.



Ordnance Survey

# Convolutional Neural Networks (CNNs)



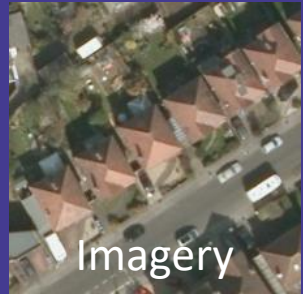Deep neural networks learn hierarchical feature representations

Learning deep representations

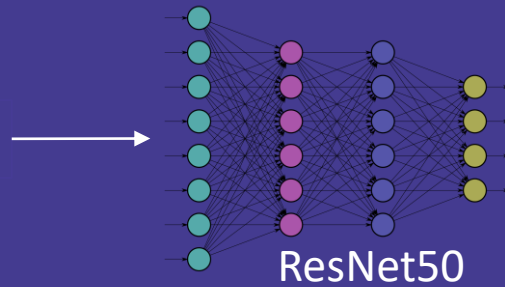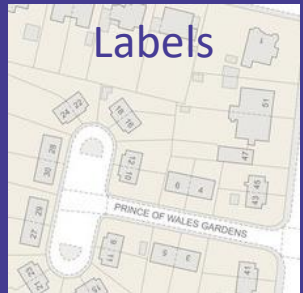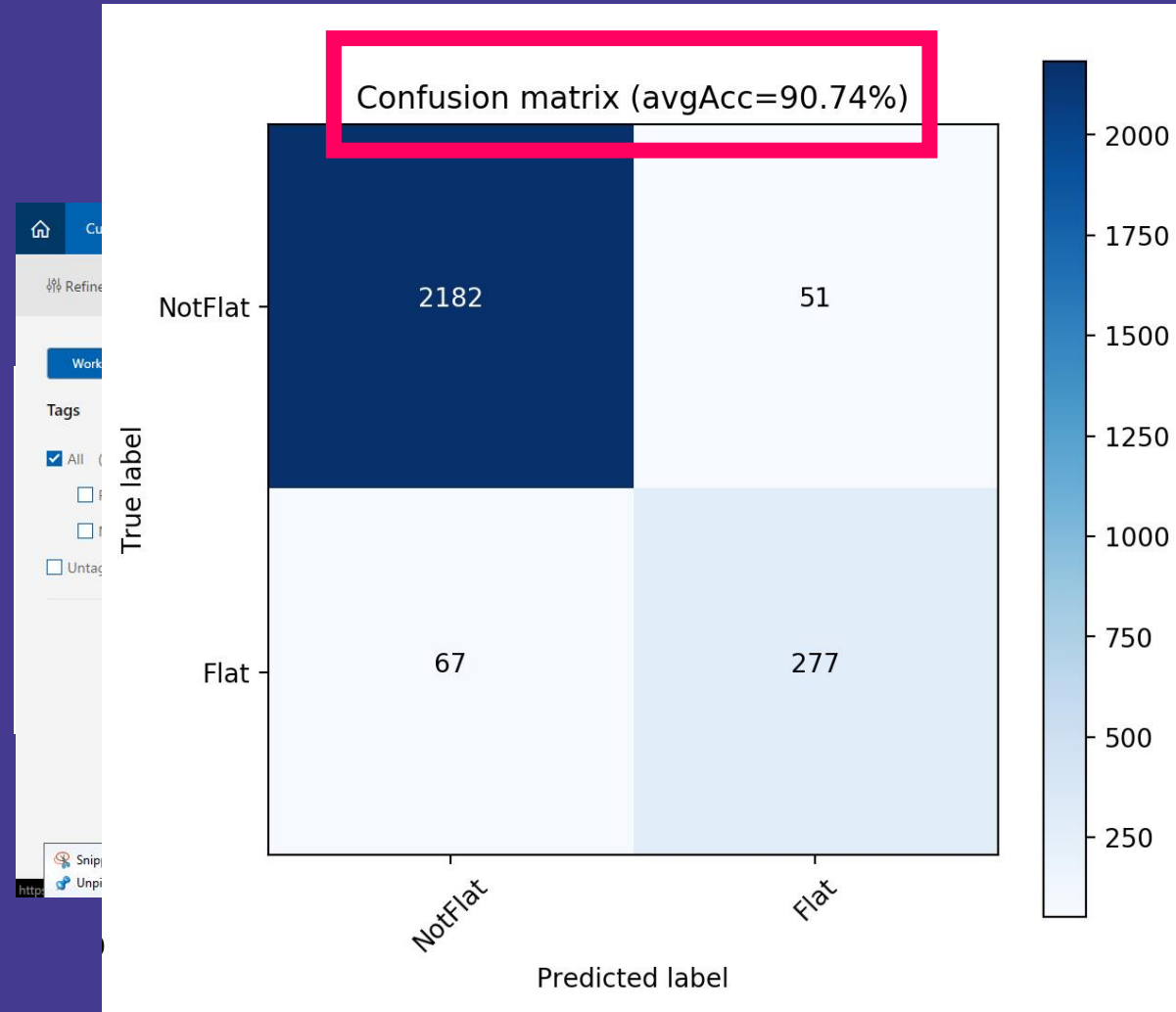Convolutional Neural Networks (CNNs)

# ImageNet

# What are we doing?

# Microsoft Roof Hack

▶ 20,000 labels

▶ Only in Hull and Southampton

▶ Webservice creating an SVM using imagenet features

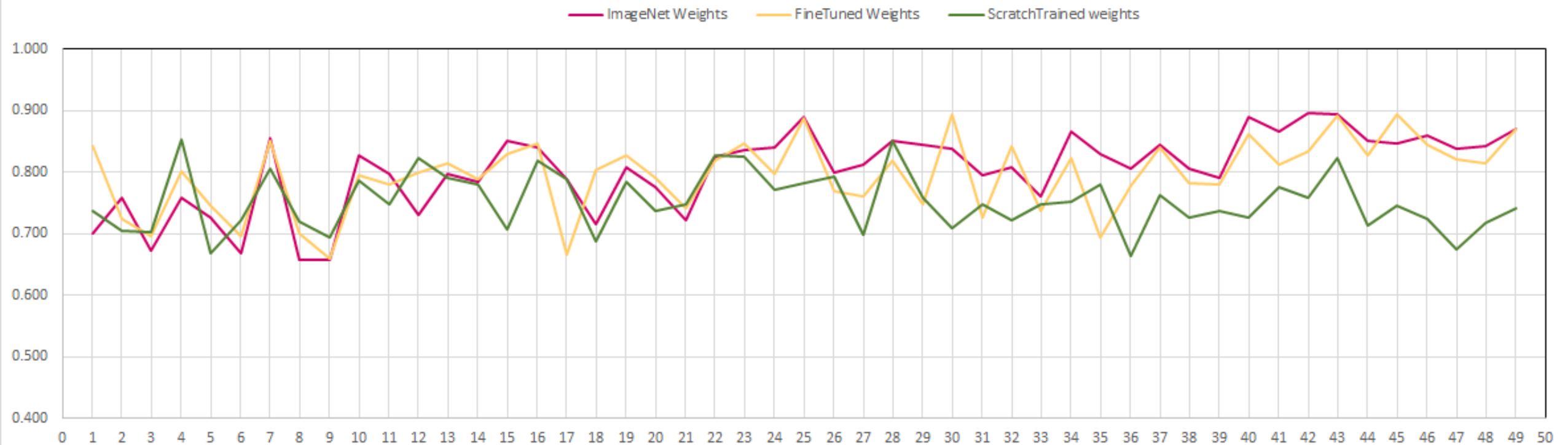▶ … but we weren't using equal classes

# Training a CNN

- 1 GPU
- 50 layers ResNet50 architecture
- 1.2 million image-class pairs (Southampton area)

- Train from scratch with aerial imagery ("scratch train")
- Steal the ImageNet weights ("imagenet")
- Or, start with imagenet and then train using aerial imagery ("fine-tune")

Ordnance
Survey

# First Results of Transfer Learning



Accuracy of classifications of inference test data using features from different layers within trained ResNet 50
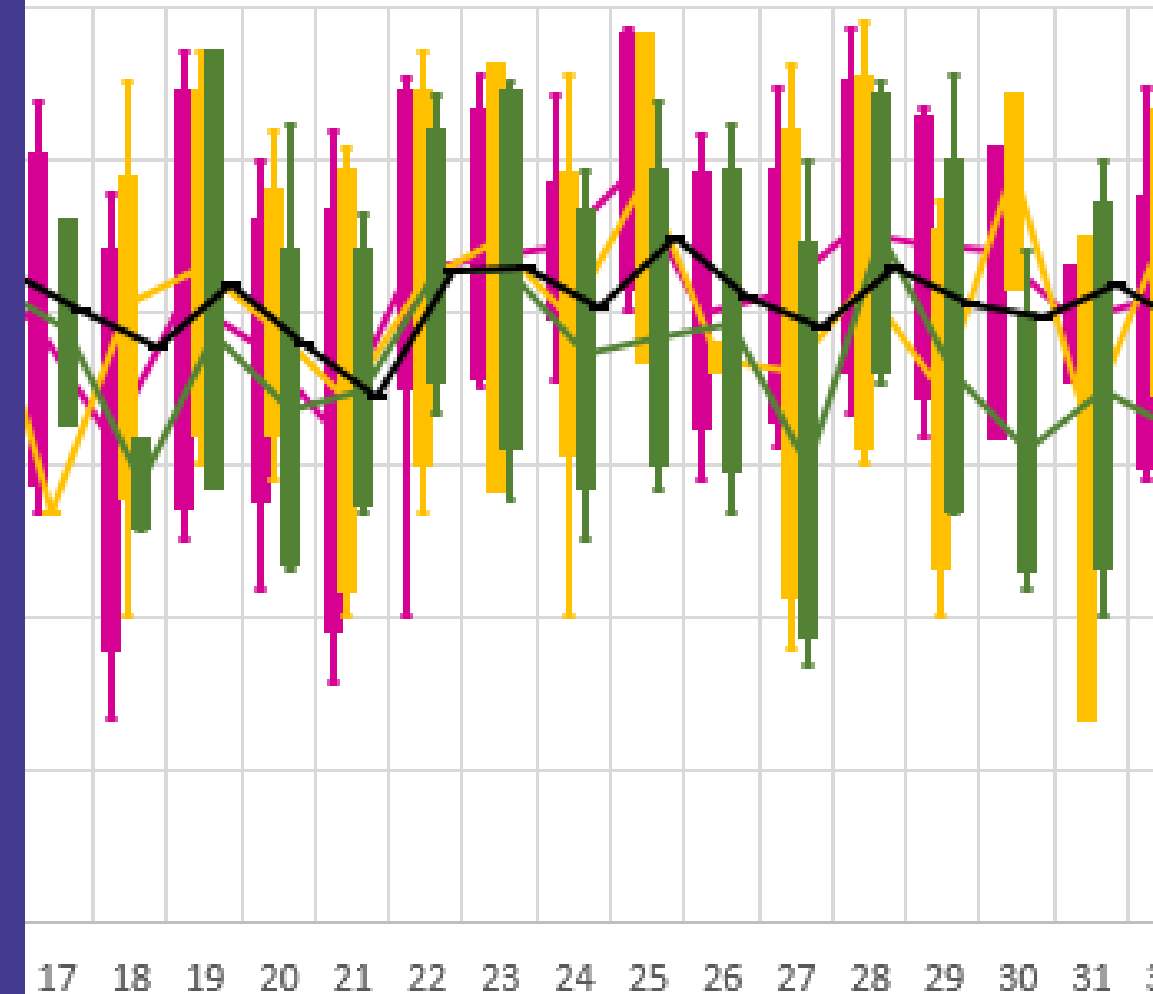
— ImageNet Weights — FineTuned Weights — ScratchTrained weights

# WHAT?

3 weeks to train (x 2)

Yet, at best AS GOOD AS ImageNet weights



erence test data using features from different

ImageNet Weights ▮  FineTuned Weights ▮  ScratchTrained we

17  18  19  20  21  22  23  24  25  26  27  28  29  30  31  3

Ordnance Survey

# But…

Aerial imagery is different to ImageNet data because:

❖ No foreground/background relationship in the scene

❖ Not 'composed', e.g. with objects arranged in the middle

❖ Different colour and texture
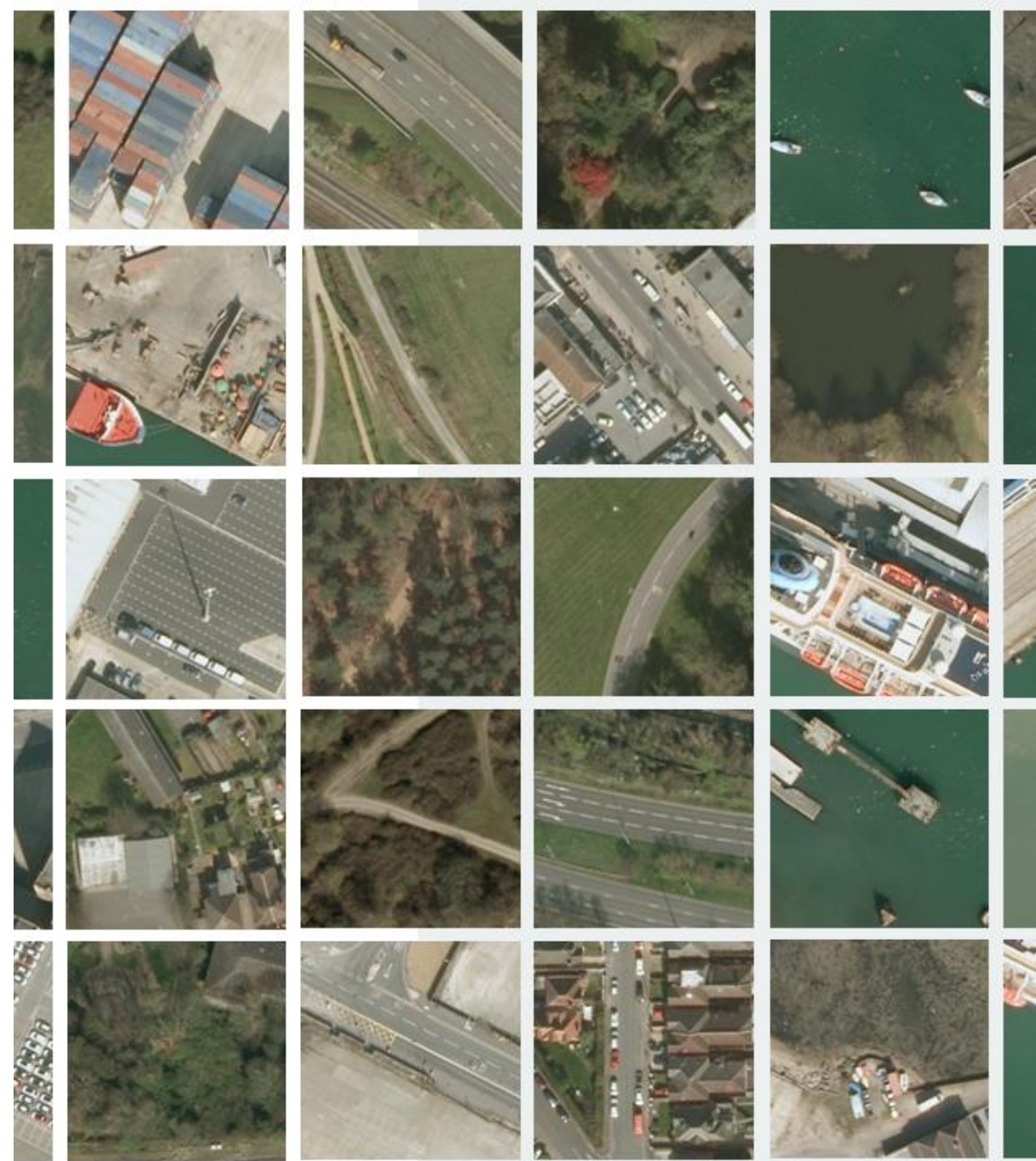
So training with aerial imagery MUST be valid

Ordnance Survey

# So…

We need to experiment to improve our model.

It's a many dimensional search space:
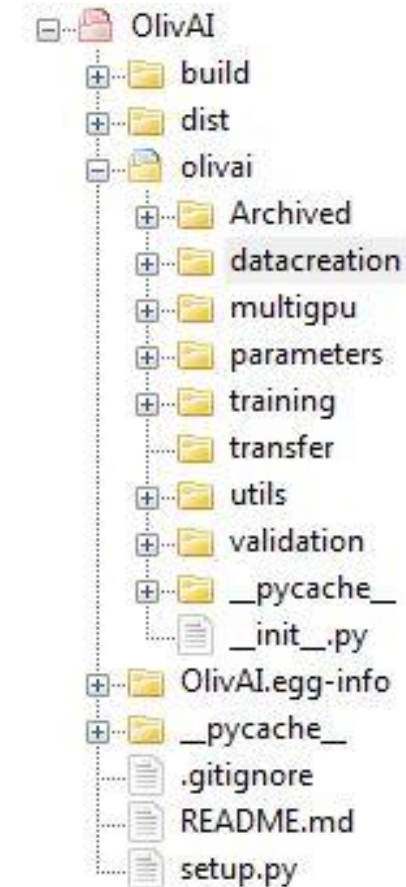
- Architecture
- Training time
- Hyperparameters

# But…

3 weeks is too long

# olivai

► 79 files

► ~ 3000 lines of code

► ~ 1200 lines of commenting (!)

► 'Olivia' → OlivAI

```
OlivAI
    build
    dist
    olivai
        Archived
        datacreation
        multigpu
        parameters
        training
        transfer
        utils
        validation
        __pycache__
        __init__.py
    OlivAI.egg-info
    __pycache__
    .gitignore
    README.md
    setup.py
```

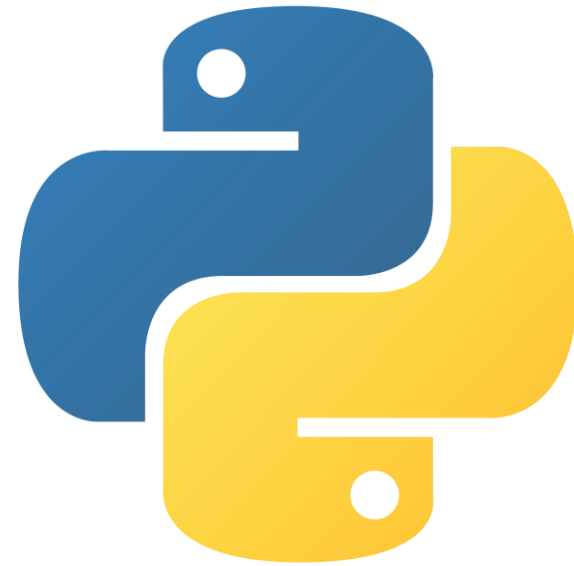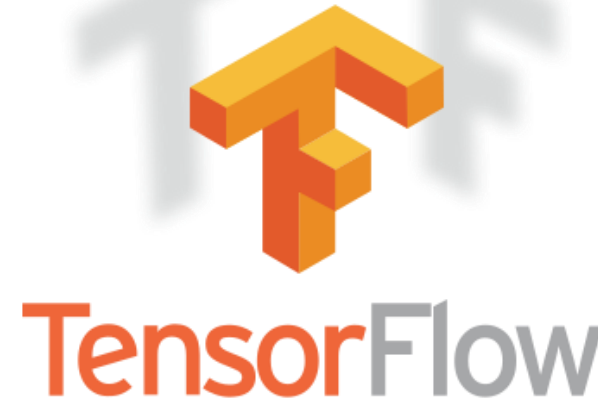## olivai

▶ Built on Python, TensorFlow and Keras

Why Keras?

▶ Well used and documented

▶ More seriously, multiple backends

▶ Easier than raw TensorFlow

▶ PyTorch was in its infancy when we started

▶ Better error handling than Neon

# Running

```python
from olivai.training import data_gen, train_on_gen
args = parse_args()

# Create the generator for the data for the model to train on
gen = data_gen(args.folder, target_size=(224,224),
               batch_size=args.batch_size)


serial_model = train_on_gen(args.arch, gen, args.batch_size,
                            args.epochs, args.gpus,
                            args.parameter_server,
                            args.method, args.num_class)
```

Ordnance
Survey

# Pipeline

```python
def precompile(model_name, num_classes, gpus, parallel, epochs,
               run_test, batch_size, method,
               validation_data=None):
    (parallel_model, serial_model) = create_model(model_name,
                                                   num_classes,
                                                   gpus, method
                                                   parallel)
    callbacks = default_callbacks(parallel, serial_model,
                                  batch_size)
    fit_dict = create_fit_dict(callbacks, epochs, run_test,
                               validation_data, parallel)
    return (parallel_model, serial_model, fit_dict)
```

Ordnance
Survey

# Pipeline

```python
def train_on_gen(model_name, gen, batch_size, epochs, gpu_count,
                 num_classes=12):
    out = precompile(model_name, num_classes, gpu_count,
                     epochs, batch_size, validation_data,
                     parallel=True)
    (parallel_model, serial_model, fit_dict) = out
    parallel_model.compile(**compile_dict)
    parallel_model.fit_generator(gen, **fit_dict)
    return serial_model
```

# Parallelisation

▶ Parallelisation is hard!

▶ Data vs model parallelism?

**What actually happens**

▶ Create serial model (on CPU)

▶ Pass serial model to `make_parallel` ( rossumai/keras-multi-gpu)

▶ `make_parallel` distributes model onto GPUs

▶ When fitting the data batch (128 in my case) is divided between the GPUs (32 on each)

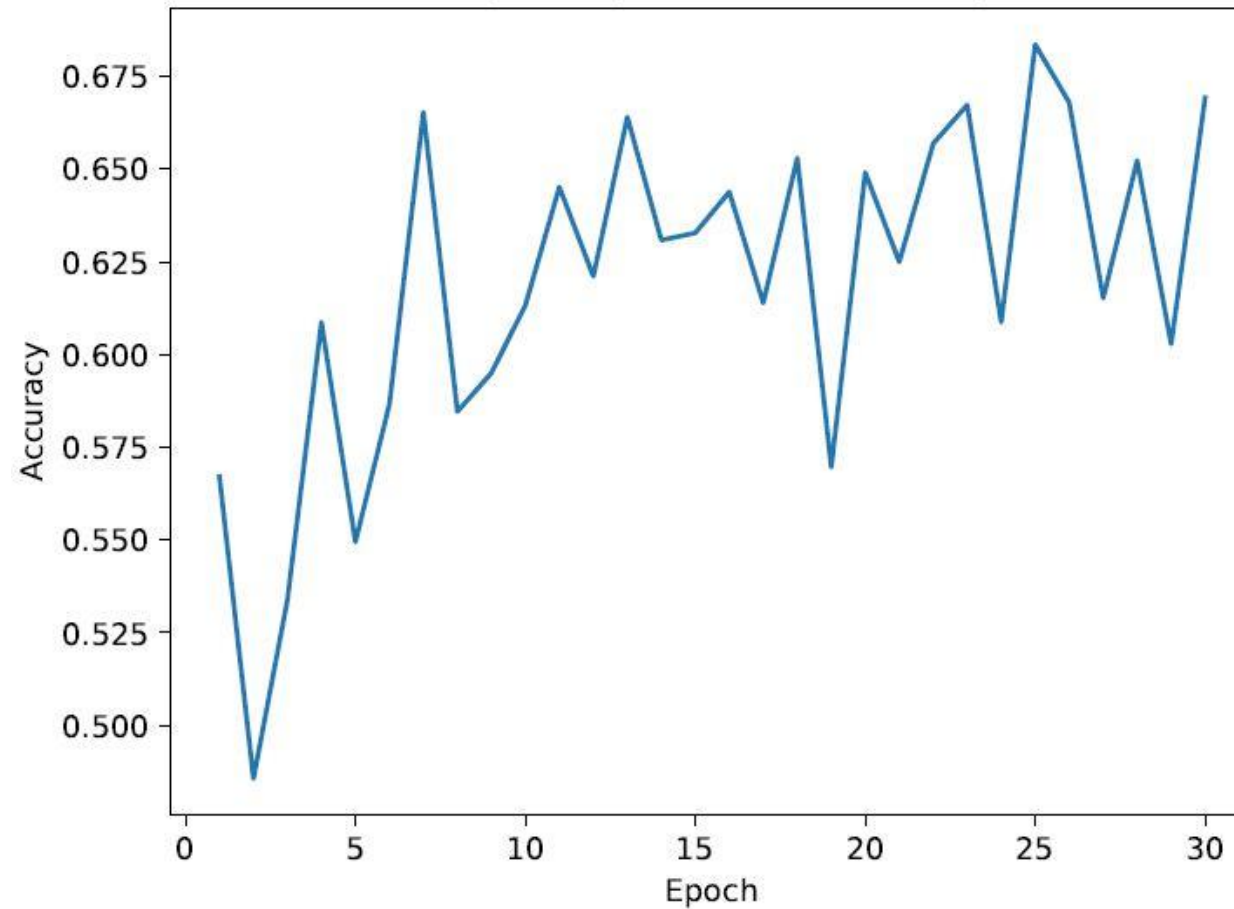▶ On batch end each model comes back with different weights and updates them centrally

Ordnance Survey

3 Weeks - > 2 days ☺

# What decreased training time

▶ Parallelisation not using buggy parallelisation

▶ Also conda installs of keras and tensorflow (on keras > 2.2.2)

# Results

## Test accuracy of ResNet at each epoch

# TopoNet V2

Speed ups mean we able to iterate quickly

TopoNet is fully trained

We can iterate quickly

We can quickly change parameters

We have many to choose from

So which TopoNet is best for what task?

# Transfer Learning

```python
def create_feature_model(model, trained_weights, num_acts=None,
                         upper_layer=None):

    model.load_weights(trained_weights)

    acts = [l for l in model.layers if isinstance(l, keras.layers.Activation)]


    if upper_layer is not None:

        layer_names = [l.name for l in acts[:upper_layer + 1]]

    else:

        layer_names = [l.name for l in acts]

    outputs = [model.get_layer(layer).output for layer in layer_names]
    inp = model.get_input_at(0)

    feature_model = keras.Model(inputs=inp, outputs=outputs)
    return feature_model
```

Ordnance
Survey

# Transfer Learning

```python
def create_features(feature_model, folder, batch_size,
                    target_size):
    gen = data_gen(folder, target_size=target_size,
                   batch_size=batch_size)

    features = feature_model.predict_generator(gen)

    return features
```

▶ This returns a list of arrays for each layer of the network (chosen using `create_feature_model`)

▶ Each array is of shape `(num_images, num_filters, size_of_filters)`

# So where are we now?

Complete
Class
Gabled
Hipped
Thatched
Tiled



Validation accuracy of Roof_Type_3M_Ponly in Wiltshire

0.86

Flat
Gabled
Hipped

# Testing

# Continuous Integration

# Conclusions

- Developed a robust python library (**olivai**) with documentation, testing and continuous integration with necessary functionality.

- Data remains '90%' of the work, and probably always will (long live the domain expert)

- We can now iterate quickly; essential to our work

- We've trained some deep networks and are now in a position to train in anger rather than mild frustration

# Any Questions?

Olivia Wilson
Research Scientist

Olivia.Wilson@os.uk
@oew1v07