

MOBILE

DEEP BUT SMALL

[martin @ reddragon.ai](mailto:martin@reddragon.ai)

[sam @ reddragon.ai](mailto:sam@reddragon.ai)

15 November 2017

WiFi : SG-Guest

Problems with Installation? **ASK!**

PLAN OF ACTION

TODAY

- Mobile : General & Specific
- Personal Project work

PLAN OF ACTION

27-NOV

- Reinforcement Learning(?)
- Finalize Projects

PERSONAL PROJECTS

- Form to fill in
- Plan to finish : **27**-Nov (last session)
- WSG DEADLINE : 30-Nov, including write-ups :
 - "Punchy Headline"
 - Minimum : README .md on GitHub page
 - Hosted slides / demo
 - Lightning Talk

SGINNOVATE



Looking for opportunities?
We find and match you to jobs based on
your skills and interests.

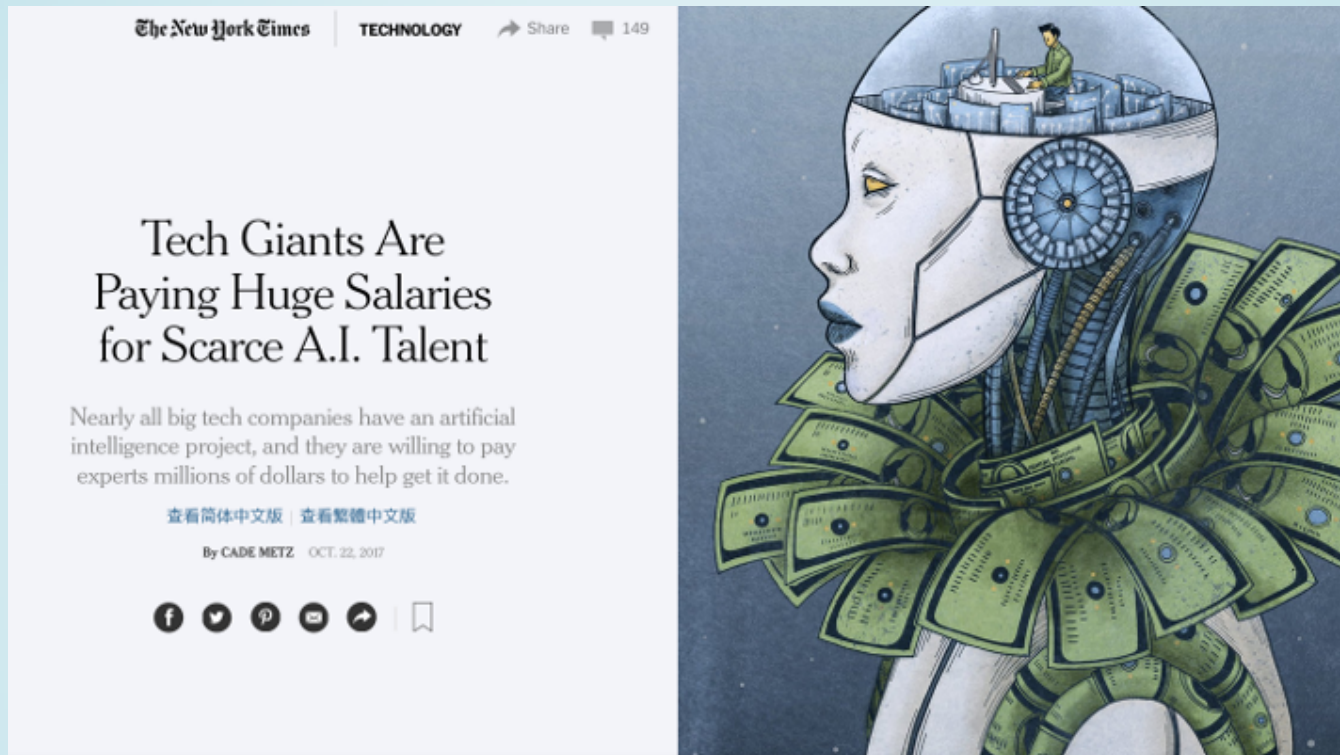
+65

PHONE NUMBER

APPLY

<http://bit.ly/2ialg60>
nicolette @ sginnovate.com

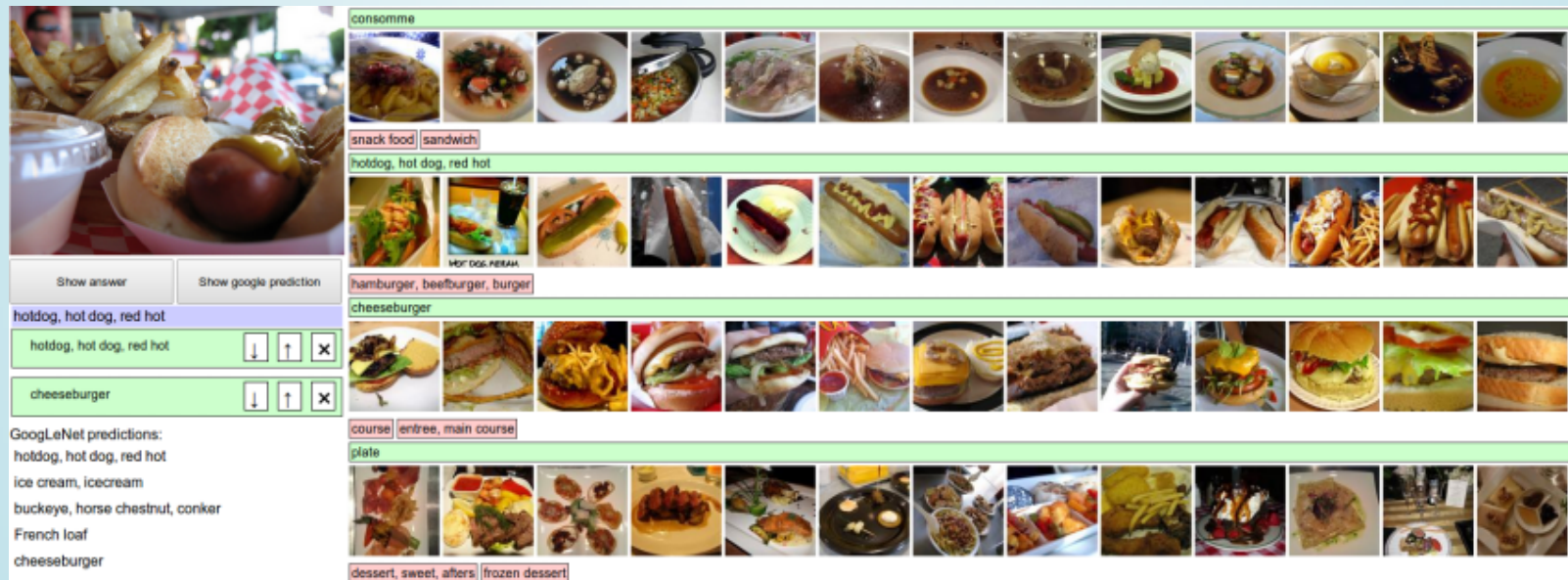
NEW YORK TIMES



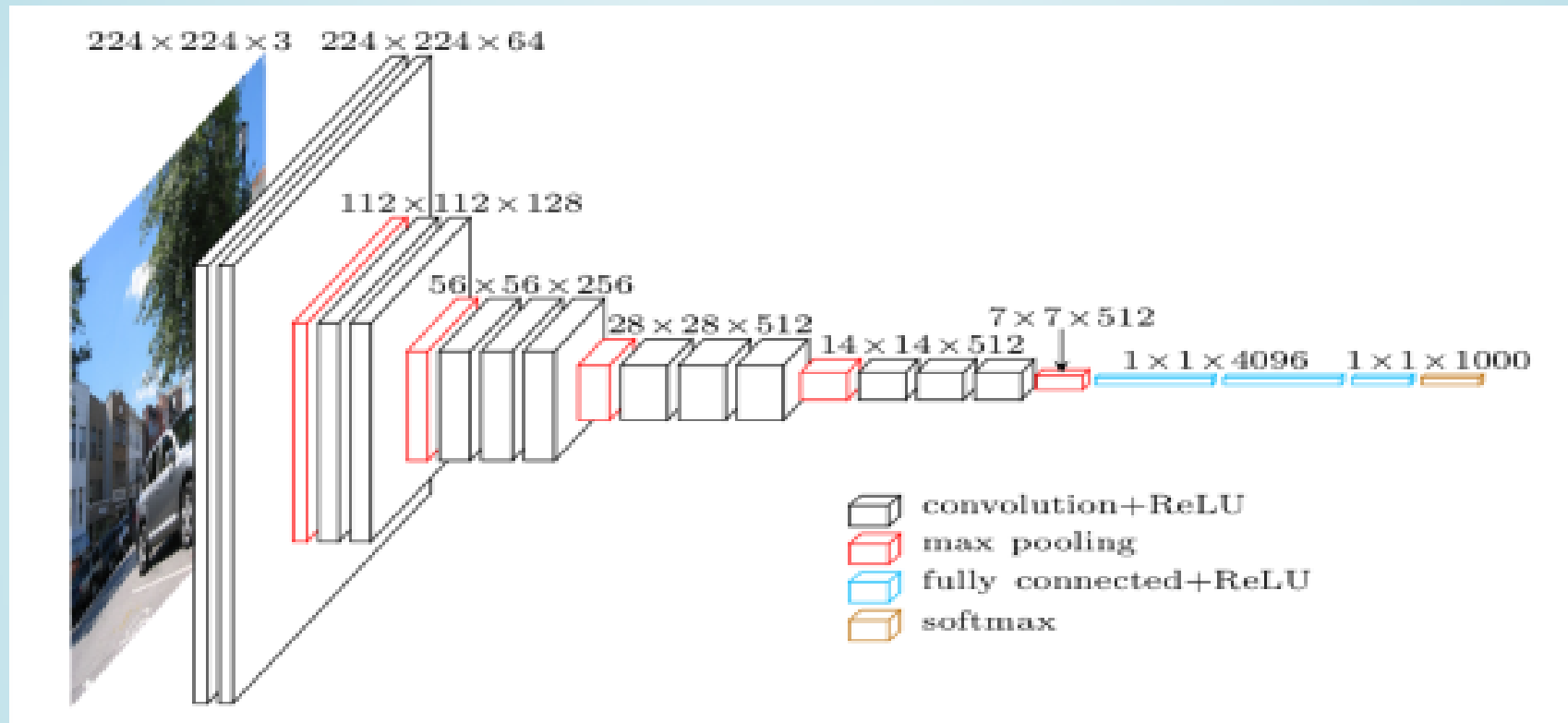
[Article Link](#)

IMAGE COMPETITION

- ImageNet aka ILSVRC
- over 15 million labeled high-resolution images...
... in over 22,000 categories

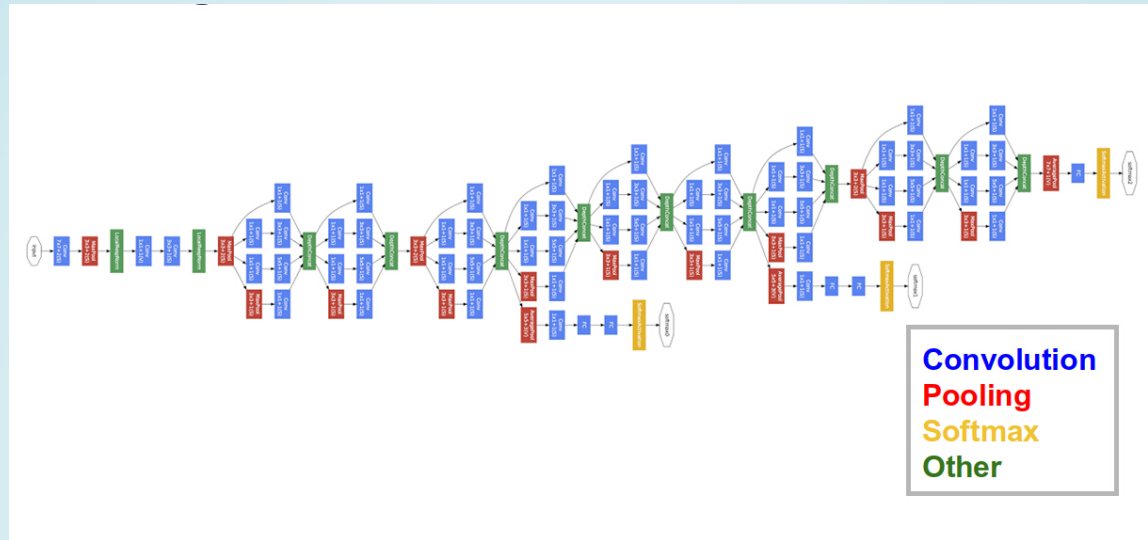


NETWORKS HAVE EVOLVED



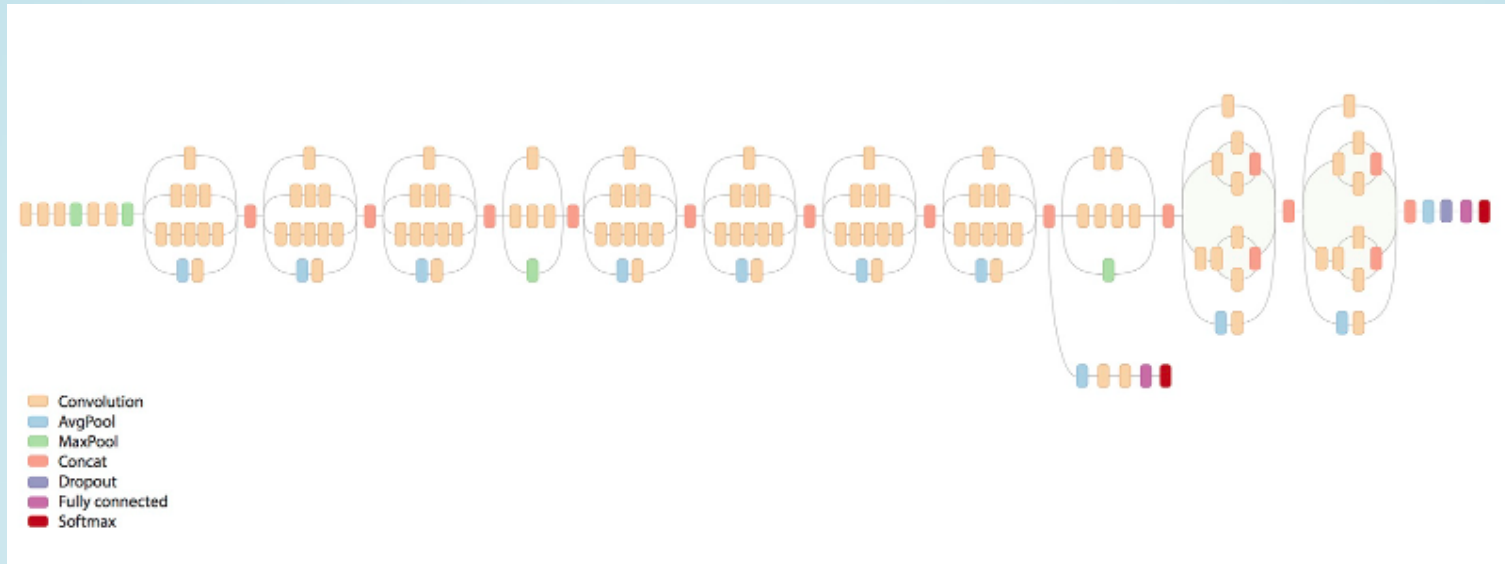
VGG 16 (2014)

GETTING DEEPER ...



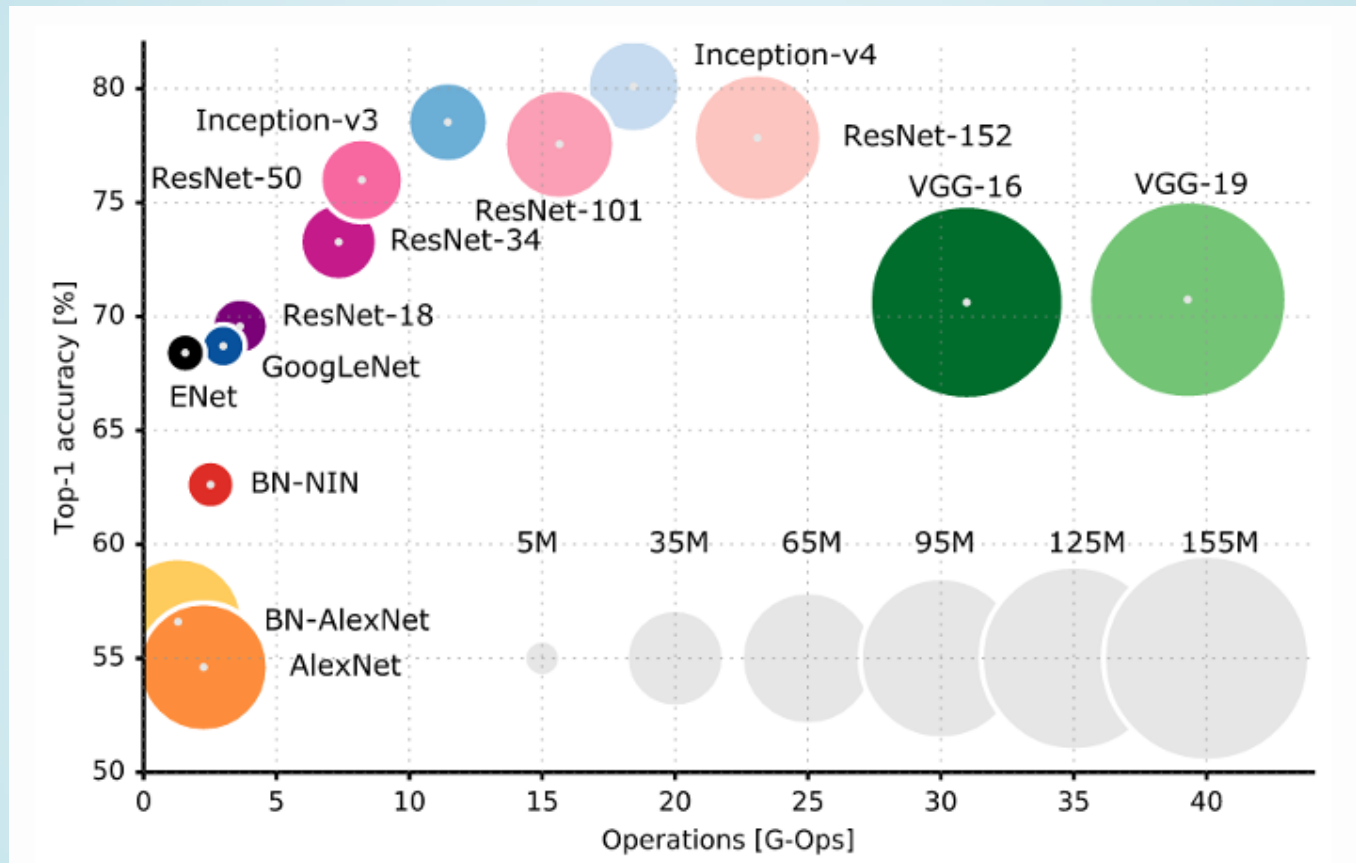
GoogLeNet (2014)

... EVEN MORE COMPLEX



Google Inception-v3 (2015)

MODEL SIZE / PERFORMANCE



[github.com/mdda/deep-learning-workshop/
notebooks/2-CNN/4-ImageNet/0-modelzoo-tf-
keras.ipynb](https://github.com/mdda/deep-learning-workshop/notebooks/2-CNN/4-ImageNet/0-modelzoo-tf-keras.ipynb)

BUT WHAT ABOUT MOBILE?

- Better performance \Rightarrow larger network
- But mobile needs us to :
 - Compress
 - Downgrade
 - Restructure

ENERGY USAGE

Table 1: Energy table for 45nm CMOS process [20]. DRAM access uses three orders of magnitude more energy than simple arithmetic and 128x more than SRAM.

| Operation | Energy [pJ] | Relative Cost |
|----------------------|-------------|---------------|
| 32 bit int ADD | 0.1 | 1 |
| 32 bit float ADD | 0.9 | 9 |
| 32 bit Register File | 1 | 10 |
| 32 bit int MULT | 3.1 | 31 |
| 32 bit float MULT | 3.7 | 37 |
| 32 bit 32KB SRAM | 5 | 50 |
| 32 bit DRAM | 640 | 6400 |

EIE: Efficient Inference Engine on Compressed Deep Neural Network (ISCA'16)

COMPRESS / DOWNGRADE

- Sparsity
- High precisions are not required
- Quantisation
- Compressibility

SPARSITY

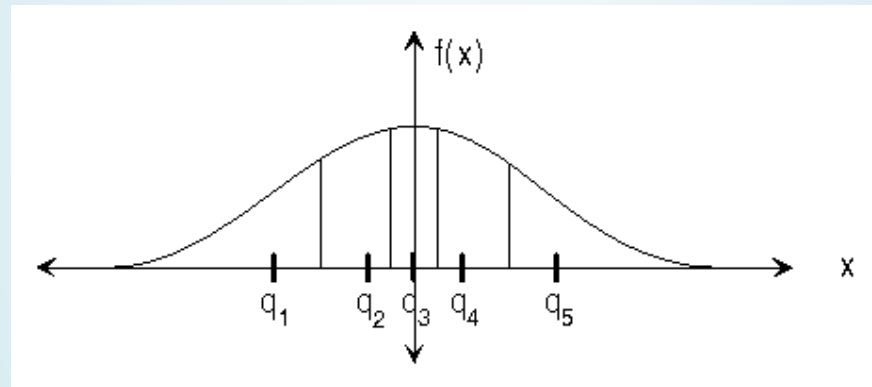
- Weights near zero \rightarrow Zero
- Clamp these weights during training

LOW-PRECISION TRAINING

- Quantise weights in forward pass
- Use 'full resolution' derivative to do backprop
- 6-bits per parameter seems to work

QUANTISATION

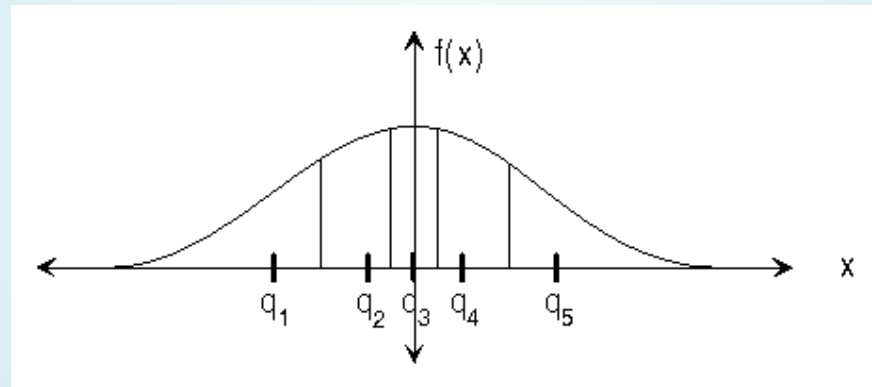
- Bucket weights into a few levels
- Store bucket positions and bucket indexes
- eg: 4 buckets (2 bits per weight index)



Quantisation in TensorFlow

COMPRESSIBILITY

- Bucket weights into a few levels
- Store bucket positions and bucket indexes
- But make bucket positions ZIPable



RESTRUCTURE 5X5

- Normal CNN '5x5' layer operation :
 - Has $(5 \times 5 + 1)$ parameters (per input channel)
- Convert to 2 stacked 3x3 layers :
 - Has $2 \times (3 \times 3 + 1)$ parameters (per input channel)
- Parameter count : $26 \times 50 = 1300 \rightarrow 20 \times 50 = 1000$

RESTRUCTURE 3X3

- Normal '3x3' CNN layer operation :
 - For each output channel:
 - Run separate '3x3' kernels over all input channels
 - Allows anywhere-to-anywhere interactions
- Parameter count : $(3 \times 3 + 1) \times 50 = 500$

SEPARABLE CONVOLUTIONS

- Separable CNN layer operation :
 - For each output channel:
 - Run one '3x3' kernel over all input channels
 - Do a weight sum (a '1x1' convolution) over results
 - Separate *texture* vs *layer* operations
- Parameter count : $(3 \times 3 + 1) \times 1 + 1 \times 1 \times (50 + 1) = 61$

+VARIATIONS

- Need to be careful that factorisation doesn't destroy performance
- Lots of scope for experimentation :
 - Xception
 - Depthwise Separable Convolutions for Neural Machine Translation

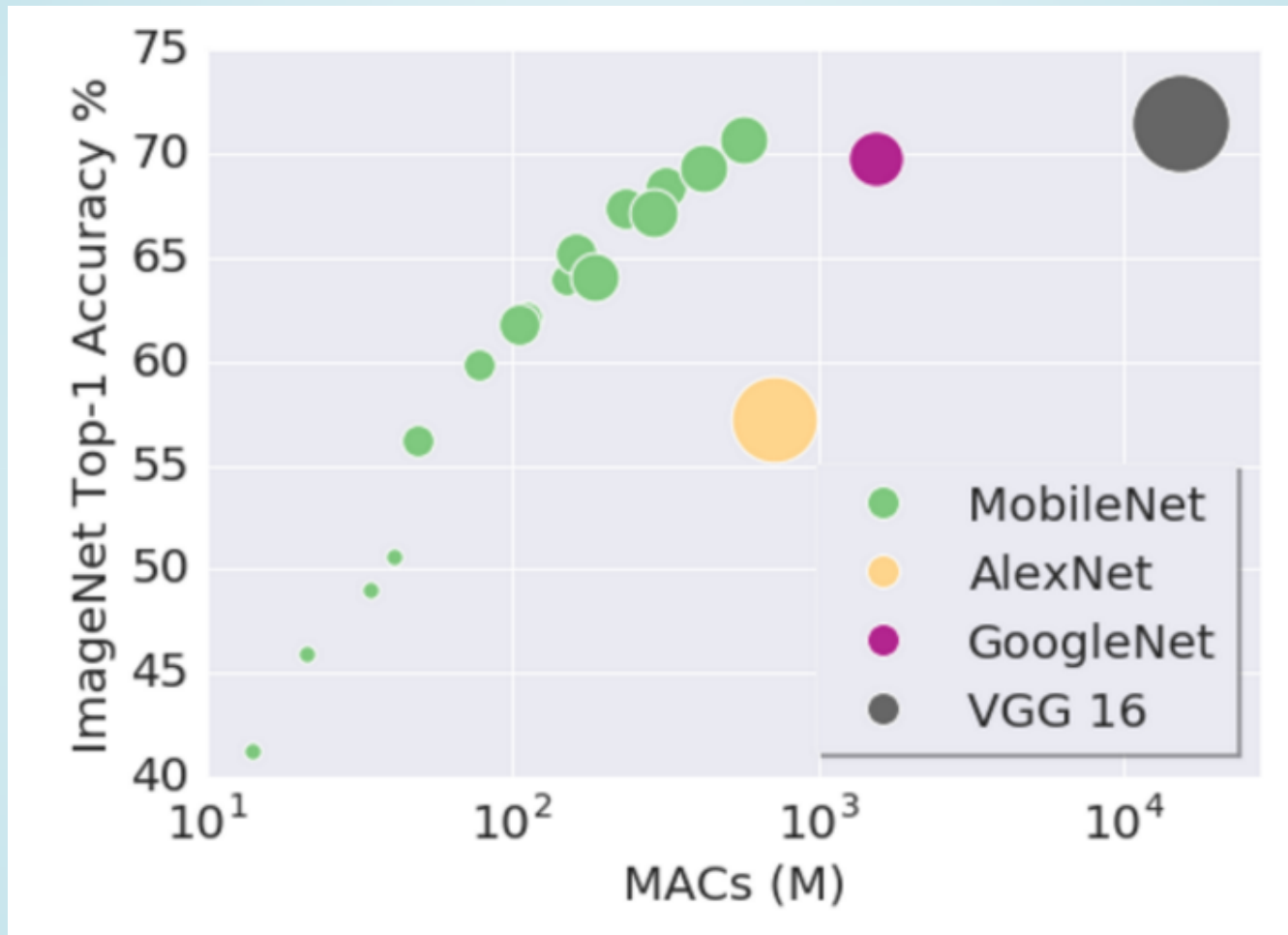
PRACTICALITIES

- Understand tradeoffs
- Use pre-defined models
- Hardware should start to arrive soon

SQUEEZENET

- SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size
- 1x1 and 3x3 layers
- No fully-connected layer

MOBILENETS ~ TF



[github.com/tensorflow/research/slim/
nets/mobilenet_v1.md](https://github.com/tensorflow/research/slim/nets/mobilenet_v1.md)

MOBILENETS ~ KERAS

[github.com/keras/applications/
mobilenet.py](https://github.com/keras/applications/mobilenet.py)

```
from keras.applications.mobilenet import MobileNet
from keras.applications.mobilenet import preprocess_input, decode_predictions

img = keras_preprocessing_image.load_img(img_path)
#...
x = preprocess_input(img)

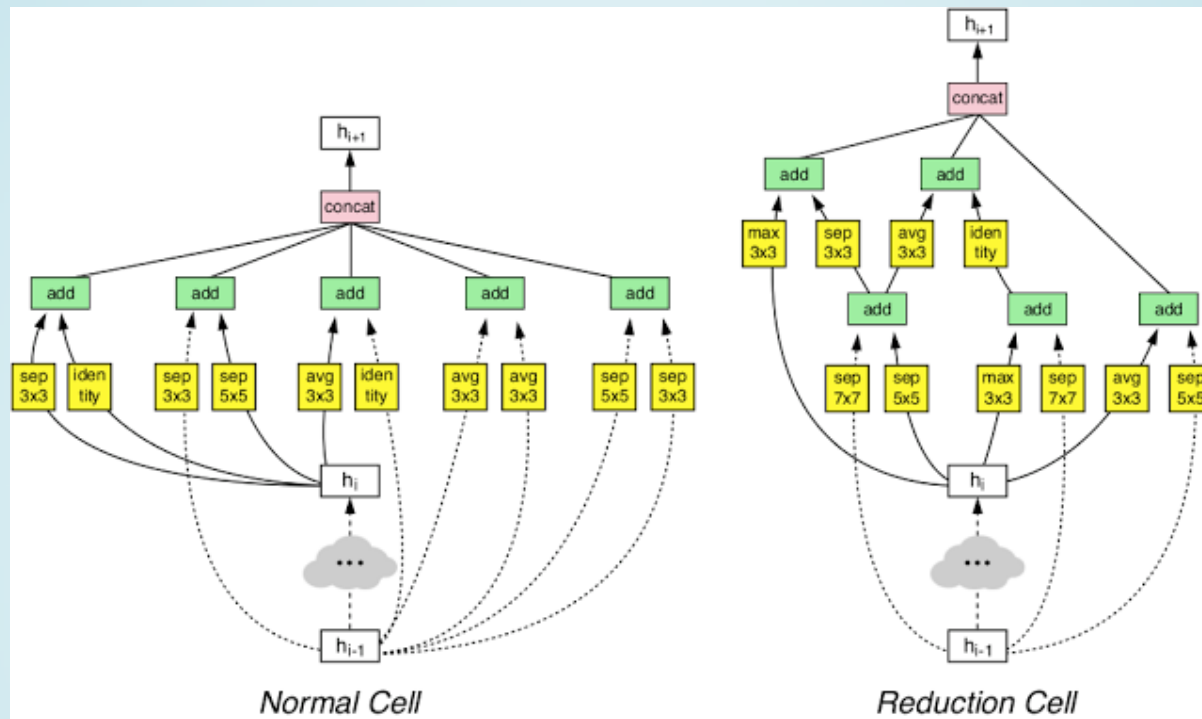
model = MobileNet(weights='imagenet')
preds = model.predict(x)

predictions = decode_predictions(preds, top=1)
```

AUTOML

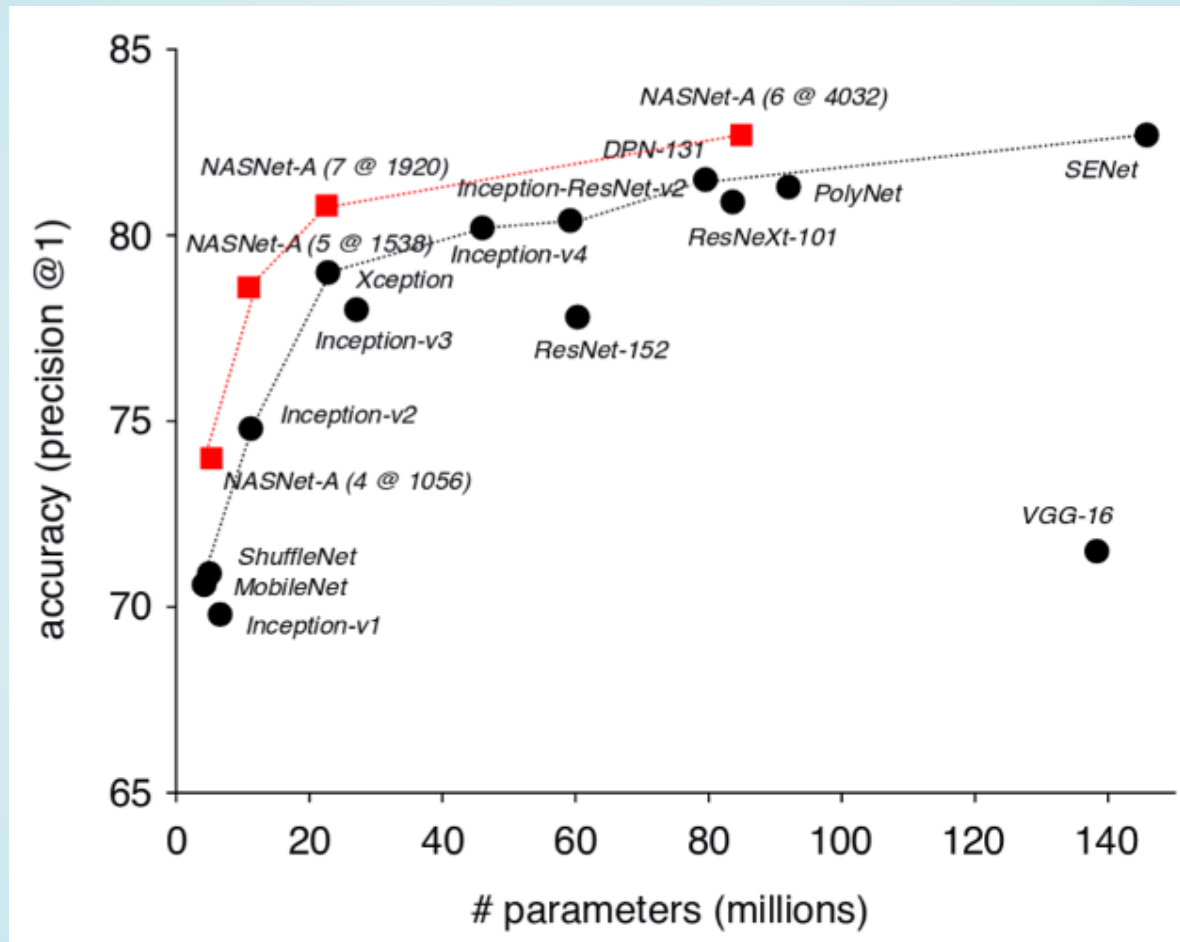
- Recent idea : Search for architecture
 - Optimise for desired performance trade-off
 - Makes use of many GPUs/TPUs
 - [Google Blog post](#)
- NB: Not so flexible for 'mash-ups'

NASNET DESIGN



Learning Transferable Architectures for Scalable Image Recognition (Nov-2017)

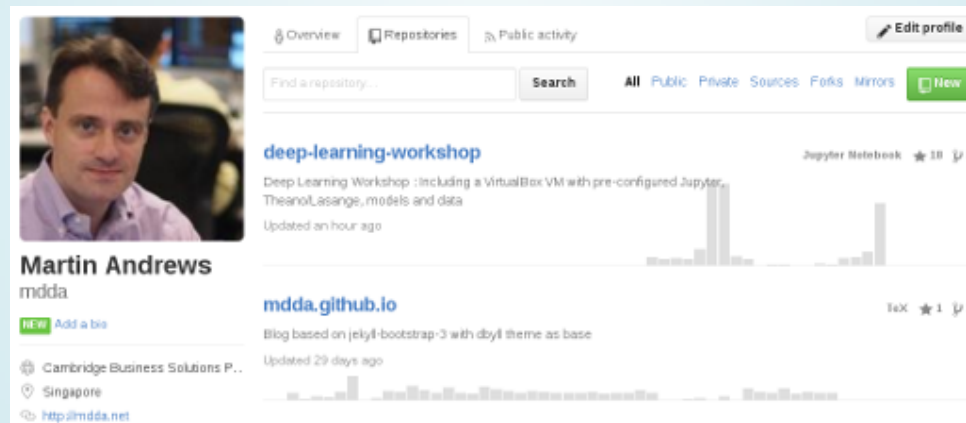
NASNET ~ TF



[github.com/tensorflow/research/
slim/nets/nasnet](https://github.com/tensorflow/research/slim/nets/nasnet)

WRAP-UP

- Explore structure vs accuracy tradeoffs
- Even tiny models work 'well enough'
- Lots more behind all this



* Please add a star... *

- QUESTIONS -

MARTIN.ANDREWS @
REDDRAGON.AI

My blog : <http://mda.net/>

GitHub : [mda](#)