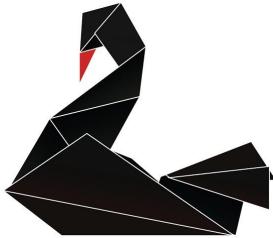


# How to do ML if you have lots of Google's GPUs

Vladimir “Vlejd” Macko

Powered by



CEAi

Google AI  
Residency Program



# Outline

## Not ML

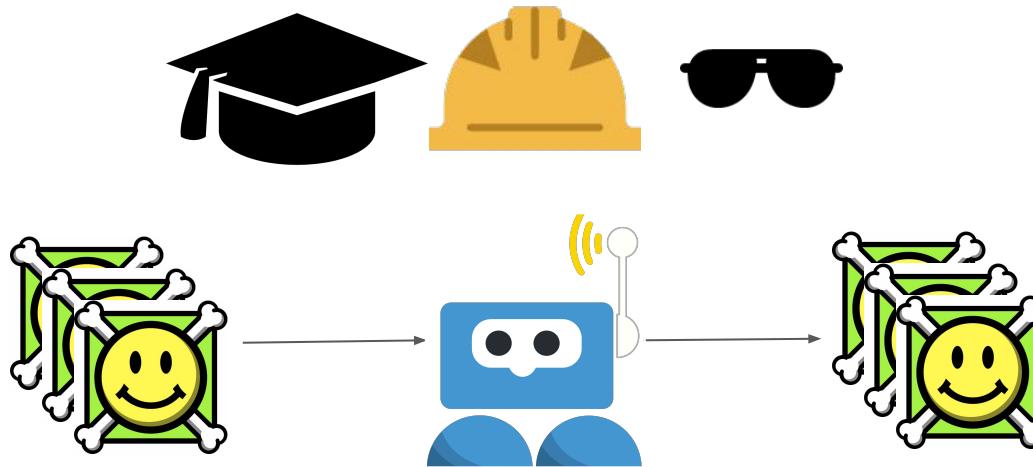
1. What is “Google” AI residency
2. How to get there

## ML

1. AutoML in theory
2. AutoML in practice



# Google AI Residency



Bc. Mgr. PhD. Pro.

# Process



# Indicators of getting in

- ✓ Google Internships
- ✓ Industry experience
- ✓ PhD. (not) required
- ✓ Luck!

# Indicators of getting in

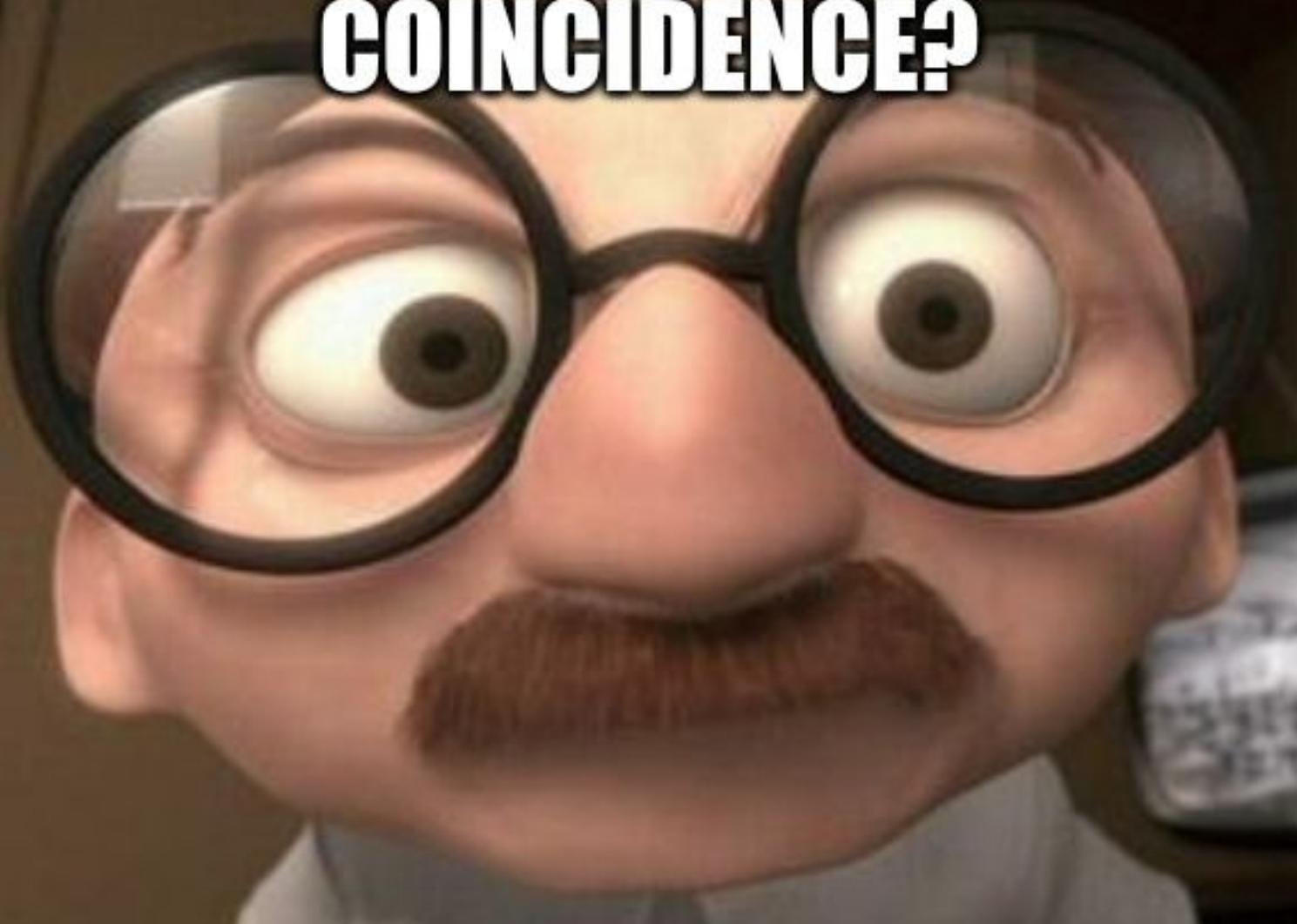
- ✓ Google Internships
- ✓ Industry experience
- ✓ PhD. (not) required
- ✓ Luck!
- ✓ Photo with Jeff Dean



# COINCIDENCE?

Indic

- ✓ G
- ✓ In
- ✓ P
- ✓ L
- ✓ P



# What to work on?

Specific project in mind (or not)

Lightning talks

AdaNet, graphnets for CO with Deepmind



# Not ML: DONE

# Machine learning

Have you tried to reproduce some  
State-Of-The-Art ML result recently?

Have you wrote your own code from scratch?

Have you managed to get **the same**  
performance?

# Why is it hard?

Architecture

Compiler flags

Hyperparameters

Hardware tricks

Data processing

Same hardware or ...

Training schedules

Same library or ...

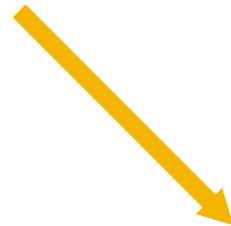
aux\_scaling  
batch\_norm\_decay  
batch\_norm\_epsilon  
cell\_name  
data\_format  
dense\_dropout\_keep\_prob  
distributed\_group\_size  
drop\_connect\_keep\_prob  
drop\_connect\_version  
drop\_path\_burn\_in\_steps  
enable\_hostcall  
gradient\_clipping  
image\_size  
iterations\_per\_loop  
l1\_decay\_rate  
label\_smoothing  
lr  
lr\_decay\_method  
lr\_decay\_value  
lr\_num\_epochs\_per\_decay

lr\_warmup\_epochs  
momentum\_rate  
moving\_average\_decay  
normal\_cell\_hiddenstate\_indices  
normal\_cell\_operations  
normal\_cell\_used\_hiddenstates  
num\_cells  
num\_epochs  
num\_label\_classes  
num\_reduction\_layers  
num\_shards  
num\_stem\_cells  
num\_train\_images  
optimizer  
per\_host\_input\_for\_training  
reduction\_cell\_hiddenstate\_indices  
reduction\_cell\_operations  
reduction\_cell\_used\_hiddenstates  
reduction\_size  
rmsprop\_decay

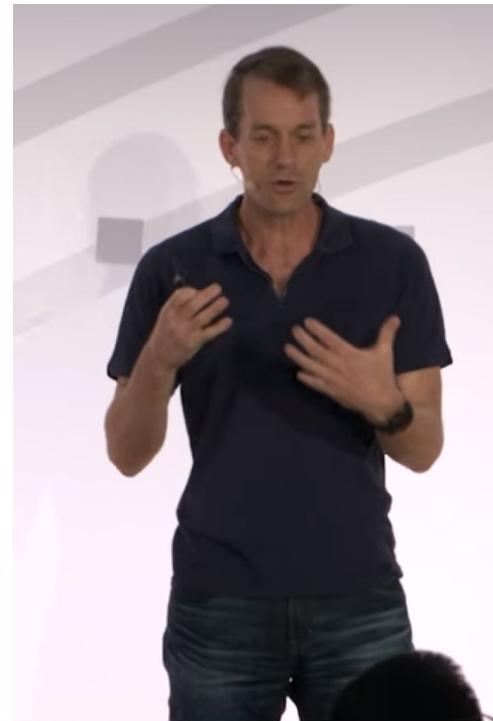
rmsprop\_epsilon  
rmsprop\_momentum\_rate  
shuffle\_buffer  
stem\_reduction\_size  
stem\_type  
train\_batch\_size  
train\_steps  
use\_aux\_head  
use\_bp16  
use\_nesterov  
use\_tpu  
version  
weight\_decay  
xla\_enable\_hlo\_trace  
xla\_jf\_use\_rotated\_pincer\_emitter2d

# Wise man once said

Solution = ML expertise + data + computation



Solution = ~~ML expertise~~ + data + **100X** computation



# ML expertise

|                        |   |                            |
|------------------------|---|----------------------------|
| Data preprocessing     | → | Autoaugment                |
| Architecture selection | → | <b>Architecture search</b> |
| Hyperparameter tuning  | → | Bayesian optimization      |
| Model selection        | → | Cross Validation           |

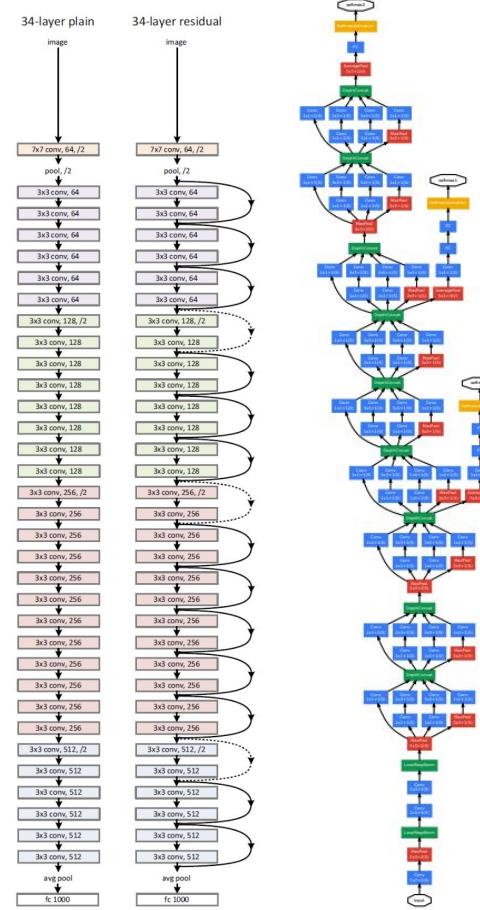
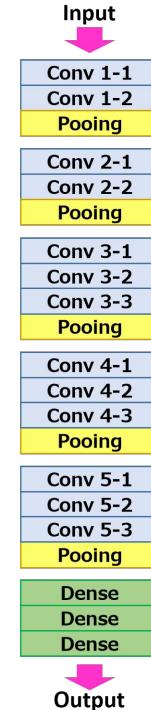


# Architecture selection by expert

Take conv/maxpool/dense layers

Constrained by memory & time

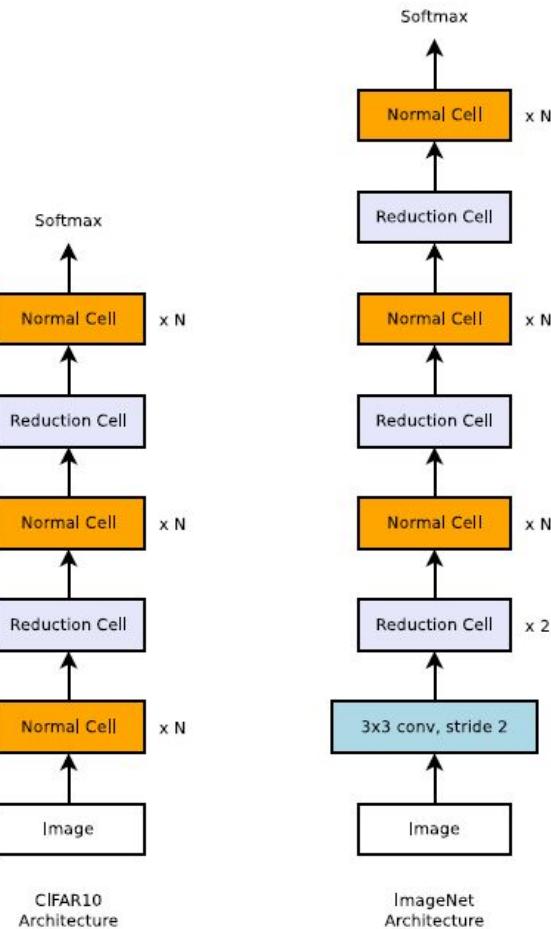
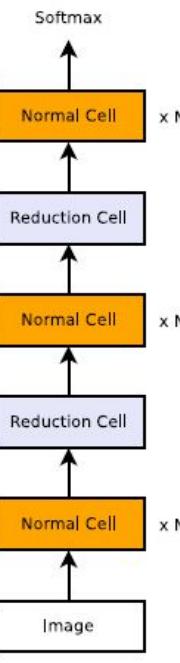
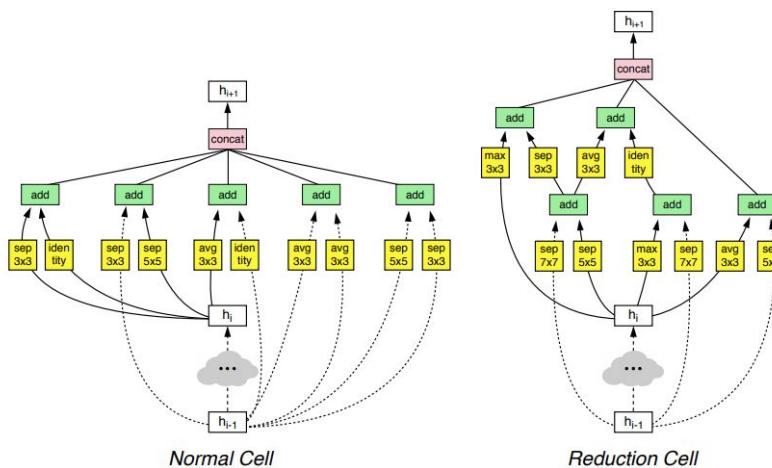
Fill parameter budget



Is this optimal?

# What blocks to stack?

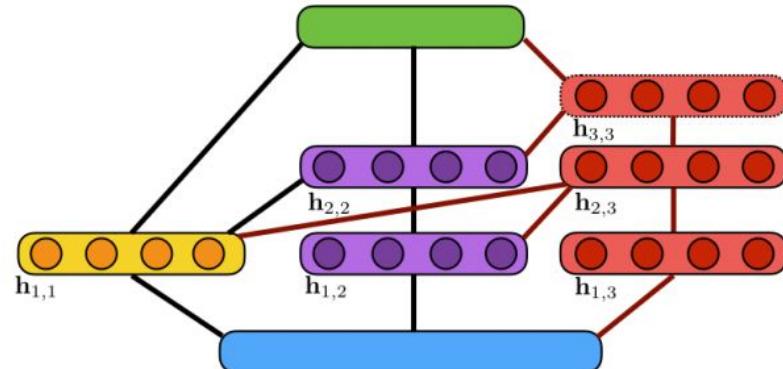
Zoph, B., & Le, Q. V. (2016). **Neural architecture search with reinforcement learning**. arXiv preprint arXiv:1611.01578.



# How to stack them?

Cortes, C., Gonzalvo, X., Kuznetsov, V., Mohri, M., & Yang, S. (2017, August). **Adanet: Adaptive structural learning of artificial neural networks**. In Proceedings of the 34th International Conference on Machine Learning-Volume 70 (pp. 874-883). JMLR.org.

Macko, V., Weill, C., Mazzawi, H., & Gonzalvo, J. (2019). **Improving neural architecture search image classifiers via ensemble learning**. arXiv preprint arXiv:1903.06236.



# Almost AdaNet

**Input:**  $S, I, G$

$$f_0 \leftarrow 0$$

**for**  $i = 1$  **to**  $I$  **do**

$$H_i \leftarrow G(f_{i-1})$$

**for**  $h_i^{(j)}$  **in**  $H_i$  **do**

$$\theta^{(j)} \leftarrow \operatorname{argmin}_{\theta} \mathcal{L}_h(S, h_i^{(j)}(\theta))$$

**end for**

$$j^* \leftarrow \operatorname{argmin}_j \mathcal{L}_f(S, \frac{1}{i} \cdot h_i^{(j)}(\theta^{(j)}) + \sum_{k=1}^{i-1} \frac{1}{i} \cdot h_k)$$

$$h_i \leftarrow h_i^{(j^*)}(\theta^{(j^*)}); f_i \leftarrow \sum_{k=1}^i \frac{1}{i} h_k$$

**end for**

**return**  $f_I$



# What about overfitting?

$Pr(\text{test error} < \text{training error} + f(\text{complexity})) = \text{big}$

Optimize upper bound on **test error**

Ensemble generalizes better than single tower

# Let's classify

# Image classification

Used to be SOTAs 2018

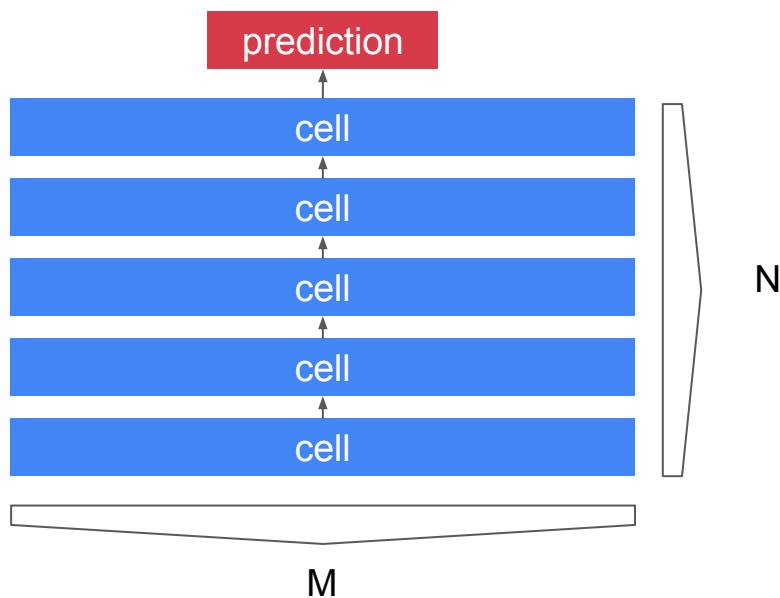
| MODEL                 | #PARAMETERS | CIFAR-10        | CIFAR-100 |
|-----------------------|-------------|-----------------|-----------|
| NASNET-A(6@768)       | 3.3M        | 2.65            | 16.58     |
| NAONET (F=36)         | 10.6M       | 3.18            | 15.67     |
| NAONET (F=64)         | 28.6M       | 2.98            | N/A       |
| NASNET-A(7@2304)      | 32.6M       | 2.40            | 16.03     |
| AMOEBA.NET-B (6, 128) | 34.9M       | $2.13 \pm 0.04$ | 15.80     |
| NAONET (F=128)        | 128M        | 2.11            | 14.75     |

Classification error  
lower is better

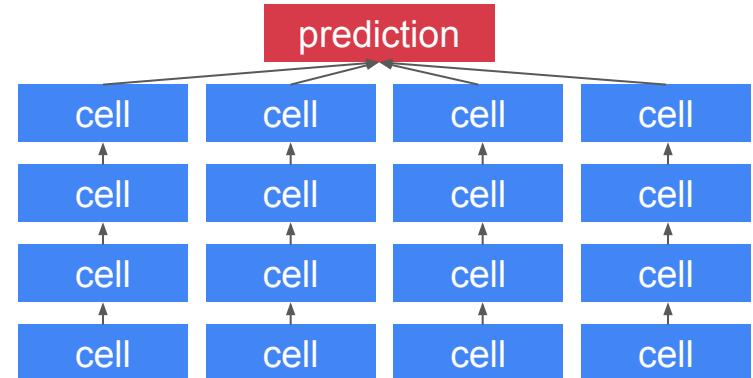
# Single tower vs. Ensemble

NASNet(7@2304)

NASNet(N@M)



10x NASNet(6@768)

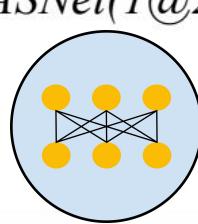


# Single tower vs. Ensemble

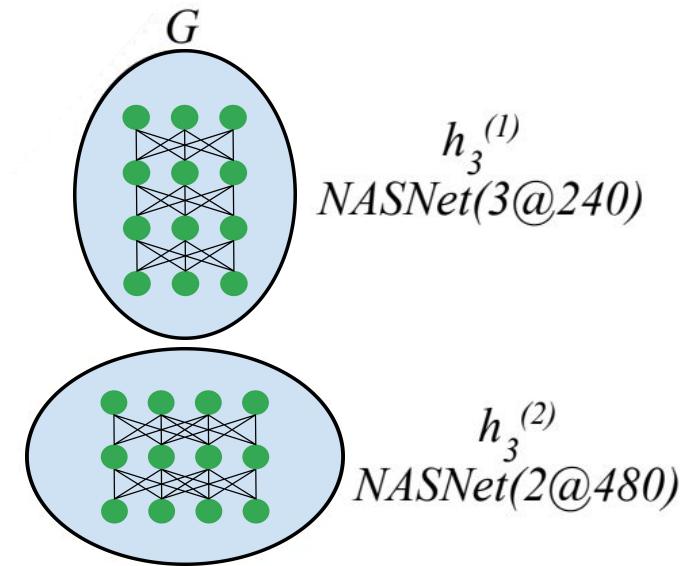
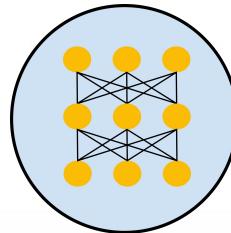
| MODEL                 | #PARAMETERS | CIFAR-10           | CIFAR-100           |
|-----------------------|-------------|--------------------|---------------------|
| NASNET-A(6@768)       | 3.3M        | 2.65               | 16.58               |
| NAONET (F=36)         | 10.6M       | 3.18               | 15.67               |
| NAONET (F=64)         | 28.6M       | 2.98               | N/A                 |
| NASNET-A(7@2304)      | 32.6M       | 2.40               | 16.03               |
| AMOEBA.NET-B (6, 128) | 34.9M       | $2.13 \pm 0.04$    | 15.80               |
| NAONET (F=128)        | 128M        | 2.11               | 14.75               |
| 10× NASNET(6@768)     | $\leq 33M$  | <b>2.29 ± 0.07</b> | <b>14.58 ± 0.05</b> |

# Grow subnetworks

$h_1$   
*NASNet(1@240)*



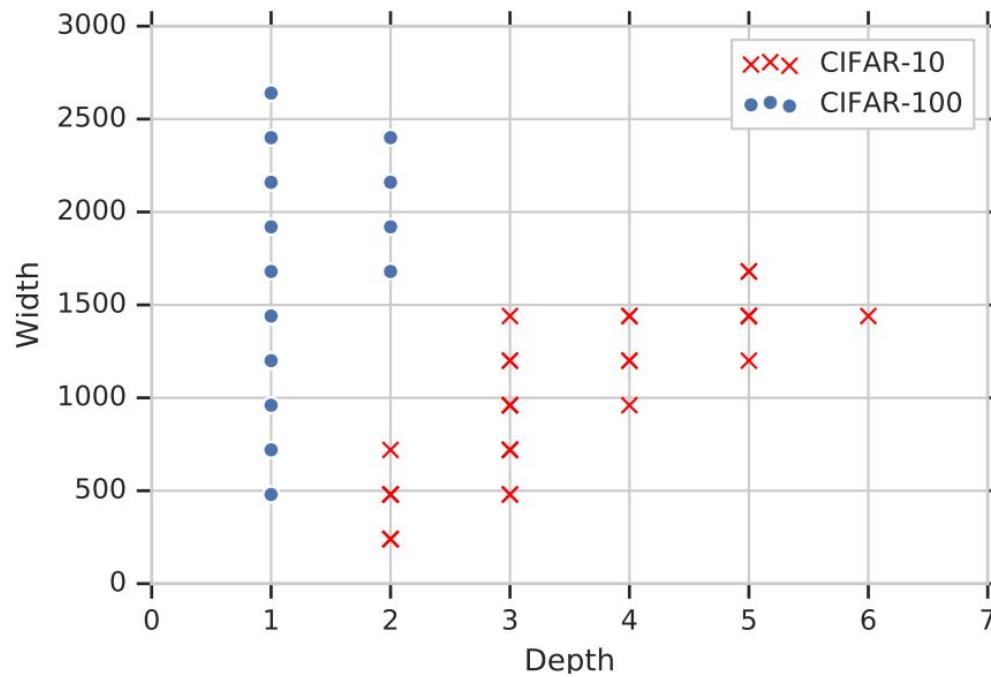
$h_2$   
*NASNet(2@240)*



# Grow subnetworks

| MODEL                 | #PARAMETERS | CIFAR-10                          | CIFAR-100                          |
|-----------------------|-------------|-----------------------------------|------------------------------------|
| NASNET-A(6@768)       | 3.3M        | 2.65                              | 16.58                              |
| NAONET (F=36)         | 10.6M       | 3.18                              | 15.67                              |
| NAONET (F=64)         | 28.6M       | 2.98                              | N/A                                |
| NASNET-A(7@2304)      | 32.6M       | 2.40                              | 16.03                              |
| AMOEBA.NET-B (6, 128) | 34.9M       | $2.13 \pm 0.04$                   | 15.80                              |
| NAONET (F=128)        | 128M        | 2.11                              | 14.75                              |
| 10× NASNET(6@768)     | $\leq 33M$  | <b><math>2.29 \pm 0.07</math></b> | <b><math>14.58 \pm 0.05</math></b> |
| ADANAS ( $G_D$ )      | $\leq 33M$  | $2.38 \pm 0.06$                   | $15.55 \pm 0.07$                   |

# Architectures for CIFAR-10 vs CIFAR-100



# Ensembling tricks

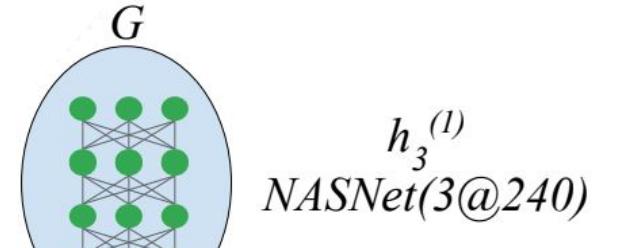
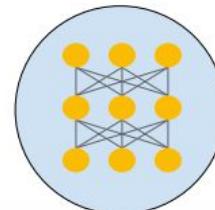
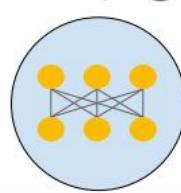
$$\mathcal{L}(h_i(x), y)$$

$$\mathcal{L}(h_i(x), y) + L_{KD}(f_{i-1}, h_i(x))$$

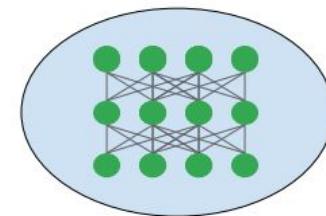
$$CrossEntropy(h_{i-1}(x), h_i(x))$$

$$CrossEntropy(f_{i-1}(x), h_i(x))$$

$$h_1 \quad h_2 \\ NASNet(1@240) \quad NASNet(2@240)$$

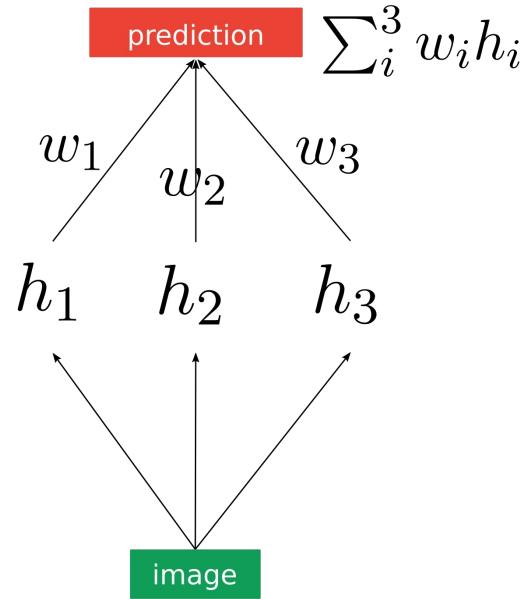
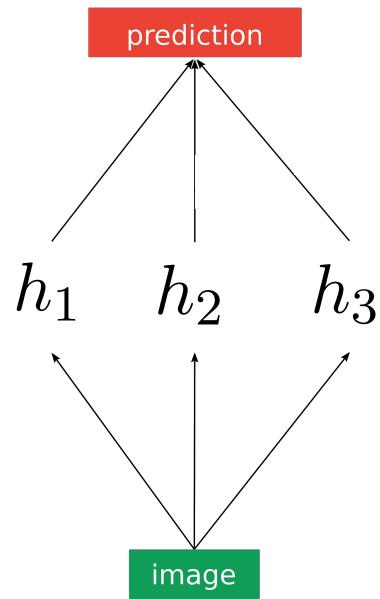


$$h_3^{(1)} \\ NASNet(3@240)$$



$$h_3^{(2)} \\ NASNet(2@480)$$

# Ensembling tricks



# Ensembling tricks

| MODEL                 | #PARAMETERS | CIFAR-10                          | CIFAR-100                          |
|-----------------------|-------------|-----------------------------------|------------------------------------|
| NASNET-A(6@768)       | 3.3M        | 2.65                              | 16.58                              |
| NAONET (F=36)         | 10.6M       | 3.18                              | 15.67                              |
| NAONET (F=64)         | 28.6M       | 2.98                              | N/A                                |
| NASNET-A(7@2304)      | 32.6M       | 2.40                              | 16.03                              |
| AMOEBA.NET-B (6, 128) | 34.9M       | $2.13 \pm 0.04$                   | 15.80                              |
| NAONET (F=128)        | 128M        | 2.11                              | 14.75                              |
| 10× NASNET(6@768)     | $\leq 33M$  | <b><math>2.29 \pm 0.07</math></b> | <b><math>14.58 \pm 0.05</math></b> |
| ADANAS ( $G_D$ )      | $\leq 33M$  | $2.38 \pm 0.06$                   | $15.55 \pm 0.07$                   |
| ADANAS +w+ AKD        | $\leq 33M$  | <b><math>2.26 \pm 0.05</math></b> | $14.83 \pm 0.14$                   |

# Other results with AdaNet

- Used in Google AutoML services
- Internal google applications
- Speech recognition, Text, Structured data
- GANs, AutoEncoders

**Ensembling rocks**

# Let's talk compute

|                           | Resources     | GPU-hours   |
|---------------------------|---------------|-------------|
| <b>NASNet-A(6@768)</b>    | 1 GPU, 1.5 d  | 36          |
| <b>NASNet-A(7@2304)</b>   | 1 GPU, 5 d    | 120         |
| <b>Ensembling</b>         | 10 GPU, 5 d   | 1.2k*       |
| <b>Growing</b>            | 10 GPU, 7 d   | 1.6k*       |
| <b>NAS search</b>         | 800 GPU, 28 d | <u>537k</u> |
| <b>NAS search 2</b>       | 500 GPU, 4 d  | 48k         |
| <b>NASNet on Imagenet</b> | 100 GPU, 8 d  | <u>19k</u>  |

# How to get that many GPUs?

A. Amazon

B. Beg for a grant

C. Cloud

D. Donate your soul to Big Tech Corporation™

# What will AutoML mean?

Privatization of ML?

Will AutoML take our jobs?

Will it help us?

Are we doomed?

Can you do this at home?





# Try this at home: [tensorflow/adanet](#)

```
import adanet
import tensorflow as tf
head = MultiClassHead(n_classes=10)
feature_columns = ...
estimator = adanet.AutoEnsembleEstimator(
    head=head,
    candidate_pool={
        "linear":
            tf.estimator.LinearEstimator(
                head=head,
                feature_columns=feature_columns,
                optimizer=...),
        "dnn":
            tf.estimator.DNNEstimator(
                head=head,
                feature_columns=feature_columns,
                optimizer=...,
                hidden_units=[1000, 500, 100])},
    max_iteration_steps=50)
estimator.train(input_fn=train_input_fn, steps=100)
metrics = estimator.evaluate(input_fn=eval_input_fn)
predictions = estimator.predict(input_fn=predict_input_fn)
```

- ✓ Multiple levels of usage
- ✓ Infrastructure
- ✓ Tests

# Final remarks

Fork me on GitHub

tensorflow / adanet

Used by 13 Watch 163 Star 2,826 Fork 403

Code Issues 45 Pull requests 2 Projects 0 Wiki Security Insights

Branch: master adanet / research / improve\_nas /

Create new file Upload files Find file History

cweill Run the 'Evaluator' before 'evaluate', 'predict', and 'export\_saved\_m... ... Latest commit d1667ed on Aug 8

..

|                  |   |              |
|------------------|---|--------------|
| images           | Add arxiv link and some images to the paper open source code.             | 6 months ago |
| trainer          | Run the 'Evaluator' before 'evaluate', 'predict', and 'export_saved_m...' | 2 months ago |
| README.md        | Remove accidental javascript from README.md.                              | 6 months ago |
| config.yaml      | Opensourcing code related to paper: Improving Neural Architecture Sea...  | 7 months ago |
| config_test.yaml | Opensourcing code related to paper: Improving Neural Architecture Sea...  | 7 months ago |
| setup.py         | Opensourcing code related to paper: Improving Neural Architecture Sea...  | 7 months ago |

# Final remarks



[to apply for AI residency 2020](#)

Fin