



Learning Transferable Architectures for ImageNet

Barret Zoph, Quoc Le

Thanks:

Vijay Vasudevan,
Irwan Bello, Jon Shlens,
Google Brain team

Current:

Solution =

ML Expertise + Data + Computation

Current:

Solution =

ML Expertise + Data + Computation

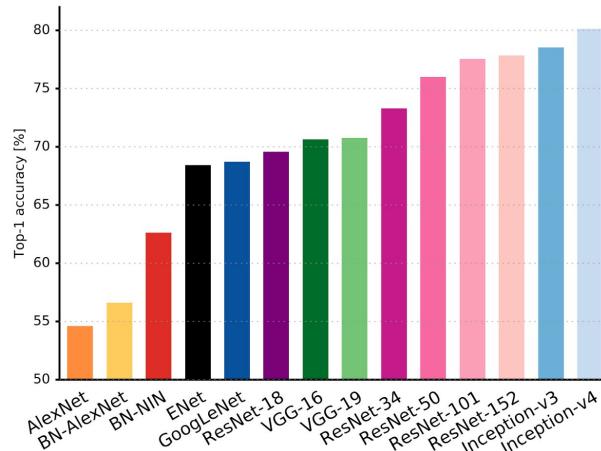
But can we turn this into:

Solution =

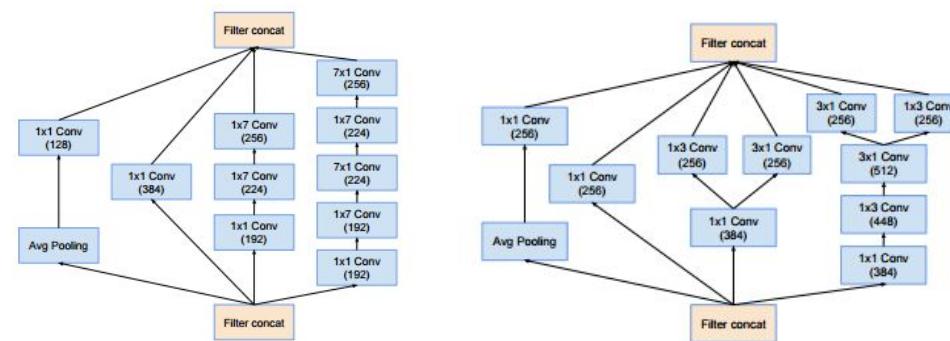
Data + 100X Computation

Importance of architectures for Vision

- Designing neural network architectures is hard
- Lots of human efforts go into tuning them
- There is not a lot of intuition into how to design them well
- Can we try and learn good architectures automatically?

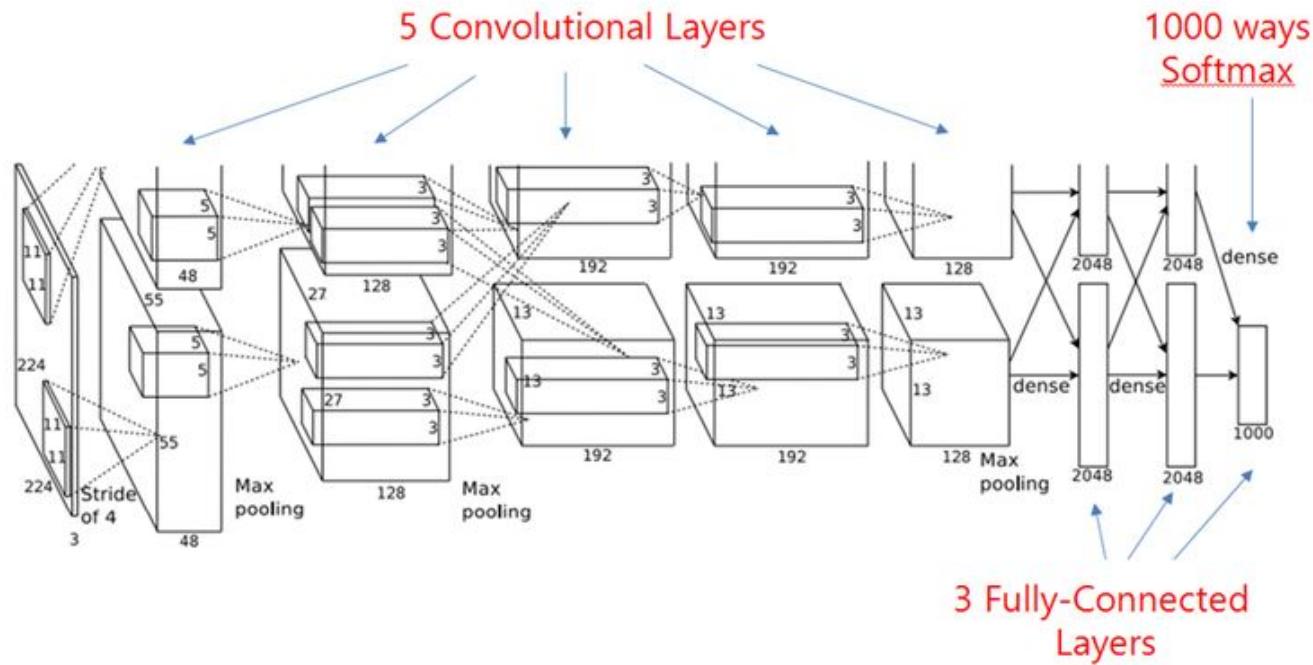


Canziani et al, 2017



Two layers from the famous Inception V4 computer vision model.
Szegedy et al, 2017

Convolutional Architectures

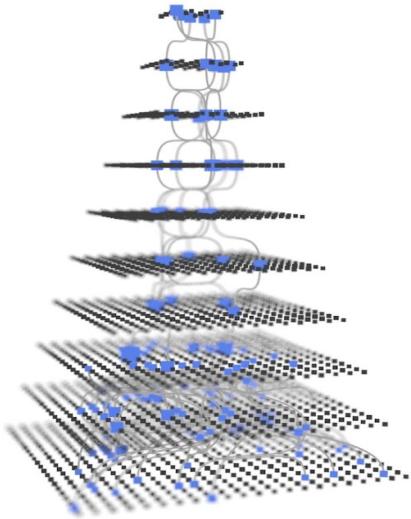


Krizhevsky et al, 2012

Neural Architecture Search

- Key idea is that we can specify the structure and connectivity of a neural network by using a configuration string
 - [“Filter Width: 5”, “Filter Height: 3”, “Num Filters: 24”]
- Our idea is to use a RNN (“Controller”) to generate this string that specifies a neural network architecture
- Train this architecture (“Child Network”) to see how well it performs on a validation set
- Use reinforcement learning to update the parameters of the Controller model based on the accuracy of the child model

Controller: proposes ML models

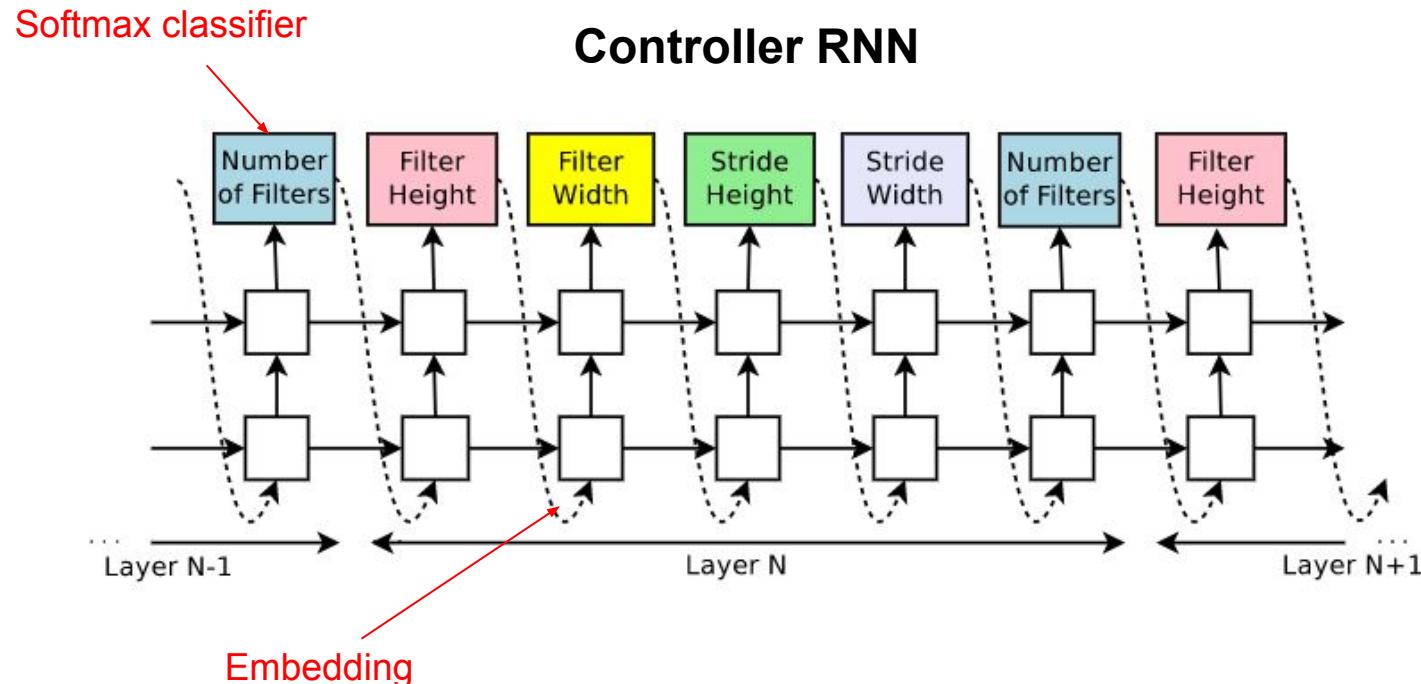


Iterate to
find the
most
accurate
model

Train & evaluate models



Neural Architecture Search for Convolutional Networks



Training with REINFORCE

$$J(\theta_c) = E_{P(a_{1:T};\theta_c)}[R]$$

Training with REINFORCE

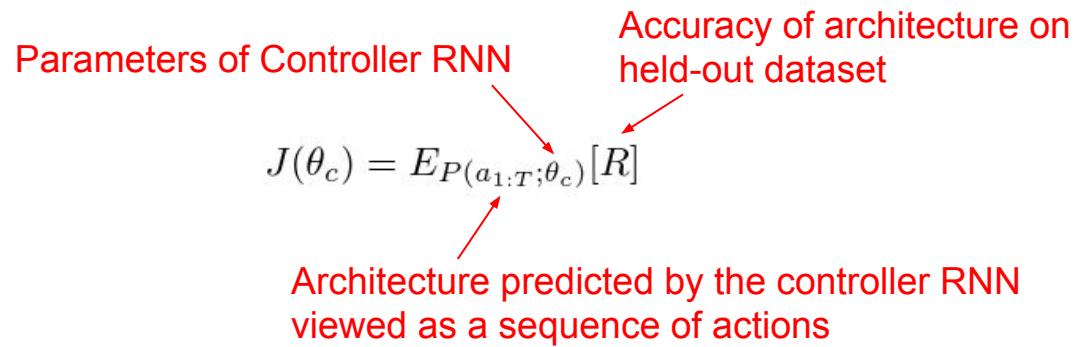
Parameters of Controller RNN

Accuracy of architecture on
held-out dataset

$$J(\theta_c) = E_{P(a_{1:T}; \theta_c)}[R]$$

Architecture predicted by the controller RNN
viewed as a sequence of actions

Training with REINFORCE



$$\nabla_{\theta_c} J(\theta_c) = \sum_{t=1}^T E_{P(a_{1:T}; \theta_c)} \left[\nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R \right]$$

Training with REINFORCE

$$J(\theta_c) = E_{P(a_{1:T};\theta_c)}[R]$$

Parameters of Controller RNN

Accuracy of architecture on held-out dataset

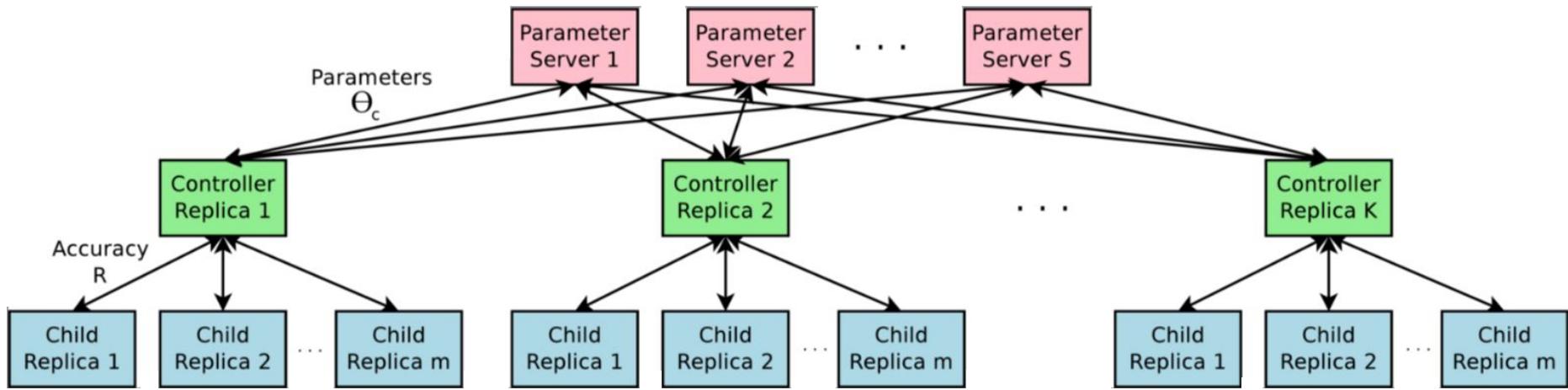
Architecture predicted by the controller RNN viewed as a sequence of actions

$$\nabla_{\theta_c} J(\theta_c) = \sum_{t=1}^T E_{P(a_{1:T};\theta_c)} [\nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R]$$

Number of models in minibatch →

$$\frac{1}{m} \sum_{k=1}^m \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R_k$$

Distributed Training



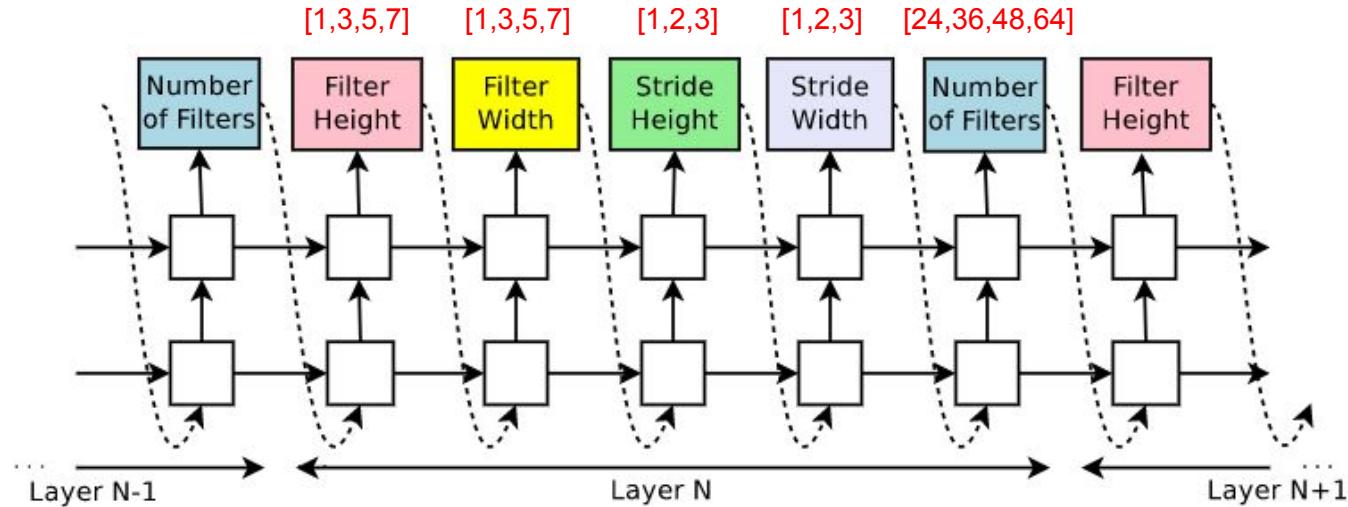
Overview of Experiments

- Apply this approach to Penn Treebank and CIFAR-10
- Evolve a convolutional neural network on CIFAR-10 and a recurrent neural network cell on Penn Treebank
- Achieve SOTA on the Penn Treebank dataset and almost SOTA on CIFAR-10 with a smaller and faster network
- Cell found on Penn Treebank beats LSTM baselines on other language modeling datasets and on machine translation

Neural Architecture Search for CIFAR-10

- We apply Neural Architecture Search to predicting convolutional networks on CIFAR-10
- Predict the following for a fixed number of layers (15, 20, 13):
 - Filter width/height
 - Stride width/height
 - Number of filters

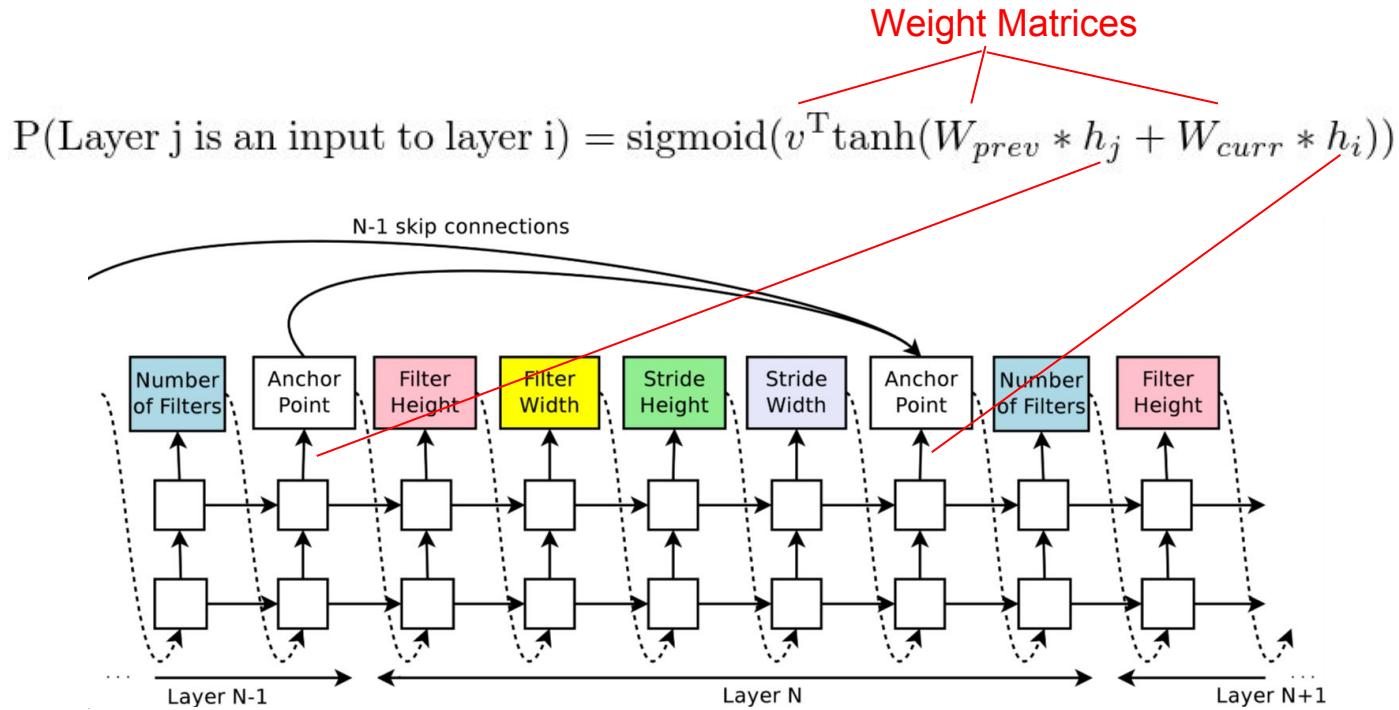
Neural Architecture Search for CIFAR-10



CIFAR-10 Prediction Method

- Expand search space to include branching and residual connections
- Propose the prediction of skip connections to expand the search space
- At layer N, we sample from N-1 sigmoids to determine what layers should be fed into layer N
- If no layers are sampled, then we feed in the minibatch of images
- At final layer take all layer outputs that have not been connected and concatenate them

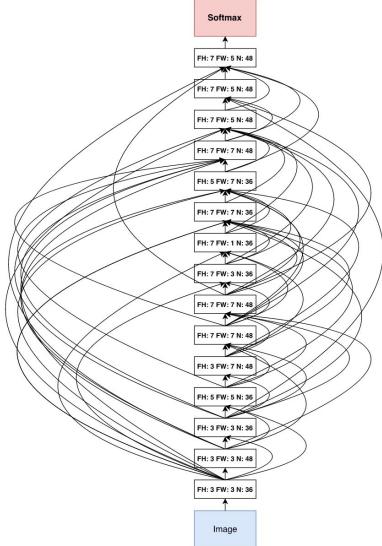
Neural Architecture Search for CIFAR-10



CIFAR-10 Experiment Details

- Use 100 Controller Replicas each training 8 child networks concurrently
- Method uses 800 GPUs concurrently at one time
- Reward given to the Controller is the maximum validation accuracy of the last 5 epochs squared
- Split the 50,000 Training examples to use 45,000 for training and 5,000 for validation
- Each child model was trained for 50 epochs
- Run for a total of **12,800** child models
- Used curriculum training for the Controller by gradually increasing the number of layers sampled

Neural Architecture Search for CIFAR-10



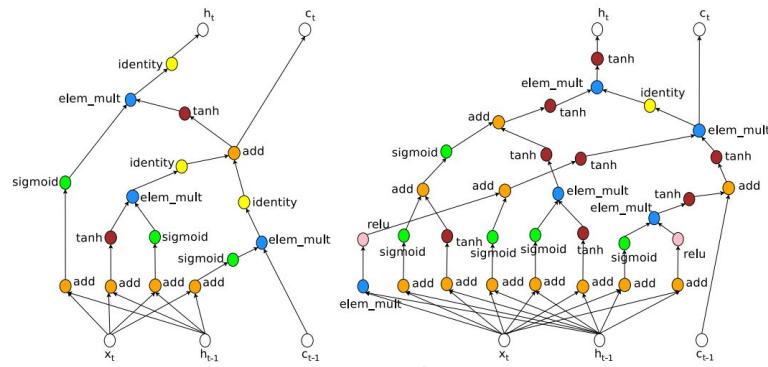
Model	Depth	Parameters	Error rate (%)
Network in Network (Lin et al., 2013)	-	-	8.81
All-CNN (Springenberg et al., 2014)	-	-	7.25
Deeply Supervised Net (Lee et al., 2015)	-	-	7.97
Highway Network (Srivastava et al., 2015)	-	-	7.72
Scalable Bayesian Optimization (Snoek et al., 2015)	-	-	6.37
FractalNet (Larsson et al., 2016) with Dropout/Drop-path	21 21	38.6M 38.6M	5.22 4.60
ResNet (He et al., 2016a)	110	1.7M	6.61
ResNet (reported by Huang et al. (2016c))	110	1.7M	6.41
ResNet with Stochastic Depth (Huang et al., 2016c)	110 1202	1.7M 10.2M	5.23 4.91
Wide ResNet (Zagoruyko & Komodakis, 2016)	16 28	11.0M 36.5M	4.81 4.17
ResNet (pre-activation) (He et al., 2016b)	164 1001	1.7M 10.2M	5.46 4.62
DenseNet ($L = 40, k = 12$) Huang et al. (2016a) DenseNet ($L = 100, k = 12$) Huang et al. (2016a) DenseNet ($L = 100, k = 24$) Huang et al. (2016a) DenseNet-BC ($L = 100, k = 40$) Huang et al. (2016b)	40 100 100 190	1.0M 7.0M 27.2M 25.6M	5.24 4.10 3.74 3.46
Neural Architecture Search v1 no stride or pooling Neural Architecture Search v2 predicting strides Neural Architecture Search v3 max pooling Neural Architecture Search v3 max pooling + more filters	15 20 39 39	4.2M 2.5M 7.1M 37.4M	5.50 6.01 4.47 3.65

5% faster

Best result of evolution (Real et al, 2017): 5.4%

Best result of Q-learning (Baker et al, 2017): 6.92%

Neural Architecture Search for PTB



Model	Parameters	Test Perplexity
Mikolov & Zweig (2012) - KN-5	2M [‡]	141.2
Mikolov & Zweig (2012) - KN5 + cache	2M [‡]	125.7
Mikolov & Zweig (2012) - RNN	6M [‡]	124.7
Mikolov & Zweig (2012) - RNN-LDA	7M [‡]	113.7
Mikolov & Zweig (2012) - RNN-LDA + KN-5 + cache	9M [‡]	92.0
Pascanu et al. (2013) - Deep RNN	6M	107.5
Cheng et al. (2014) - Sum-Prod Net	5M [‡]	100.0
Zaremba et al. (2014) - LSTM (medium)	20M	82.7
Zaremba et al. (2014) - LSTM (large)	66M	78.4
Gal (2015) - Variational LSTM (medium, untied)	20M	79.7
Gal (2015) - Variational LSTM (medium, untied, MC)	20M	78.6
Gal (2015) - Variational LSTM (large, untied)	66M	75.2
Gal (2015) - Variational LSTM (large, untied, MC)	66M	73.4
Kim et al. (2015) - CharCNN	19M	78.9
Press & Wolf (2016) - Variational LSTM, shared embeddings	51M	73.2
Merity et al. (2016) - Zoneout + Variational LSTM (medium)	20M	80.6
Merity et al. (2016) - Pointer Sentinel-LSTM (medium)	21M	70.9
Inan et al. (2016) - VD-LSTM + REAL (large)	51M	68.5
Zilly et al. (2016) - Variational RHN, shared embeddings	24M	66.0
Neural Architecture Search with base 8	32M	67.9
Neural Architecture Search with base 8 and shared embeddings	25M	64.0
Neural Architecture Search with base 8 and shared embeddings	54M	62.4

RL vs random search

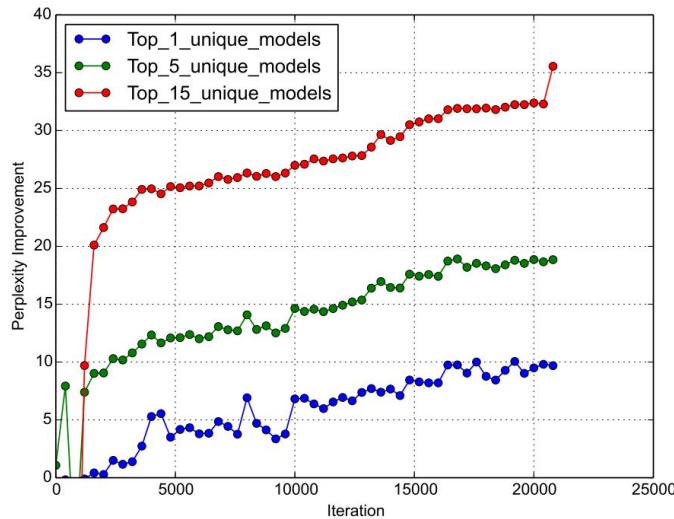


Figure 6: Improvement of Neural Architecture Search over random search over time. We plot the difference between the average of the top k models our controller finds vs. random search every 400 models run.

Learn the Optimization Update Rule

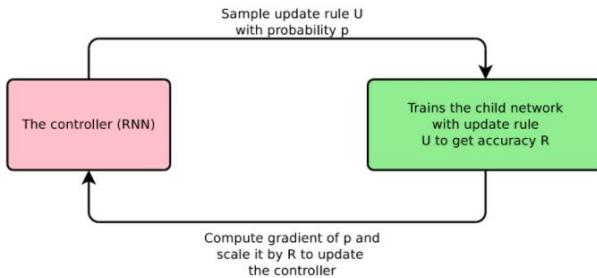


Figure 1. Overview of Neural Optimizer Search.

Optimizer	Final Val	Final Test	Best Val	Best Test
SGD	92.0	91.8	92.9	91.9
Momentum	92.7	92.1	93.1	92.3
ADAM	90.4	90.1	91.8	90.7
RMSProp	90.7	90.3	91.4	90.3
$[e^{\text{sign}(g)*\text{sign}(m)} + \text{clip}(g, 10^{-4})] * g$	92.5	92.4	93.8	93.1
$\text{clip}(\hat{m}, 10^{-4}) * e^{\hat{v}}$	93.5	92.5	93.8	92.7
$\hat{m} * e^{\hat{v}}$	93.1	92.4	93.8	92.6
$g * e^{\text{sign}(g)*\text{sign}(m)}$	93.1	92.8	93.8	92.8
$\text{drop}(g, 0.3) * e^{\text{sign}(g)*\text{sign}(m)}$	92.7	92.2	93.6	92.7
$\hat{m} * e^{g^2}$	93.1	92.5	93.6	92.4
$\text{drop}(\hat{m}, 0.1)/(e^{g^2} + \epsilon)$	92.6	92.4	93.5	93.0
$\text{drop}(g, 0.1) * e^{\text{sign}(g)*\text{sign}(m)}$	92.8	92.4	93.5	92.2
$\text{clip}(\text{RMSProp}, 10^{-5}) + \text{drop}(\hat{m}, 0.3)$	90.8	90.8	91.4	90.9
$\text{ADAM} * e^{\text{sign}(g)*\text{sign}(m)}$	92.6	92.0	93.4	92.0
$\text{ADAM} * e^{\hat{m}}$	92.9	92.8	93.3	92.7
$g + \text{drop}(\hat{m}, 0.3)$	93.4	92.9	93.7	92.9
$\text{drop}(\hat{m}, 0.1) * e^{g^3}$	92.8	92.7	93.7	92.8
$g - \text{clip}(g^2, 10^{-4})$	93.4	92.8	93.7	92.8
$e^g - e^{\hat{m}}$	93.2	92.5	93.5	93.1
$\text{drop}(\hat{m}, 0.3) * e^w$	93.2	93.0	93.5	93.2

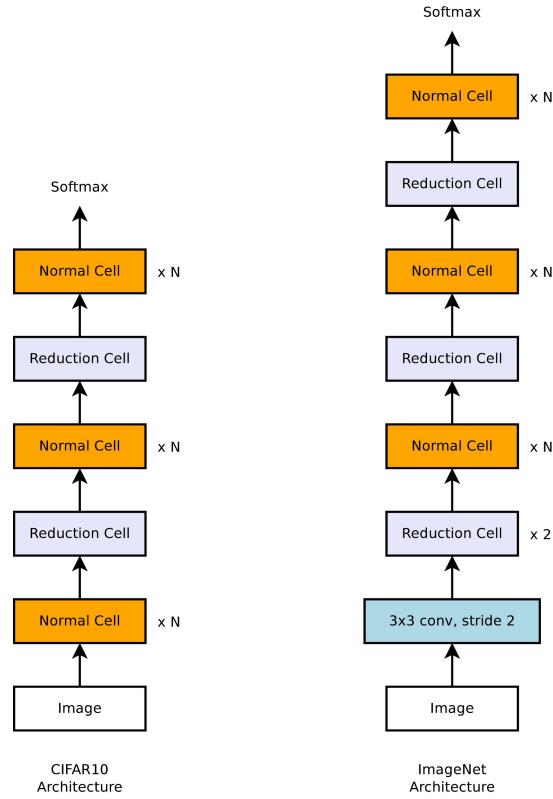
Table 1. Performance of Neural Optimizer Search and standard optimizers on the Wide-ResNet architecture (Zagoruyko & Komodakis, 2016) on CIFAR-10. Final Val and Final Test refer to the final validation and test accuracy after for training for 300 epochs. Best Val corresponds to the best validation accuracy over the 300 epochs and Best Test is the test accuracy at the epoch where the validation accuracy was the highest.

Neural Optimizer Search using Reinforcement Learning,
 Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc Le. ICML 2017

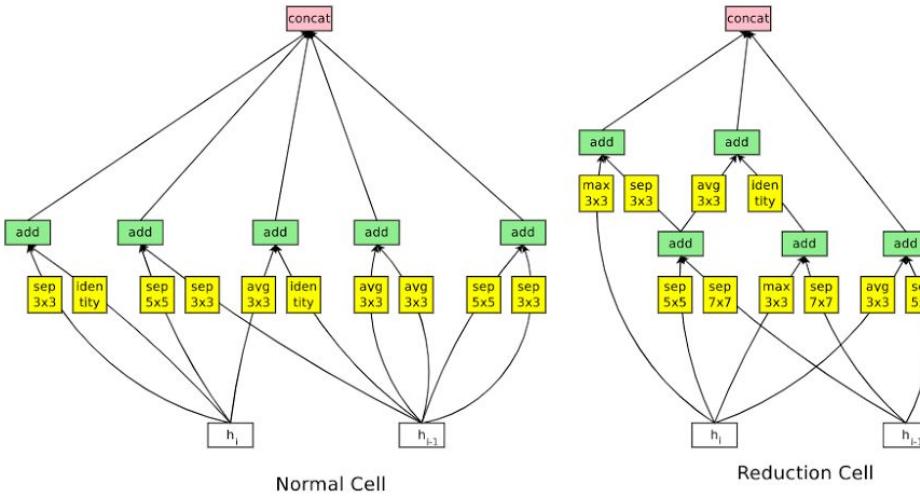
Neural Architecture Search for ImageNet

- Neural Architecture Search directly on ImageNet is expensive
- Key idea is to run Neural Architecture Search on CIFAR-10 to find a “cell”
- Construct a bigger net from the “cell” and train the net on ImageNet

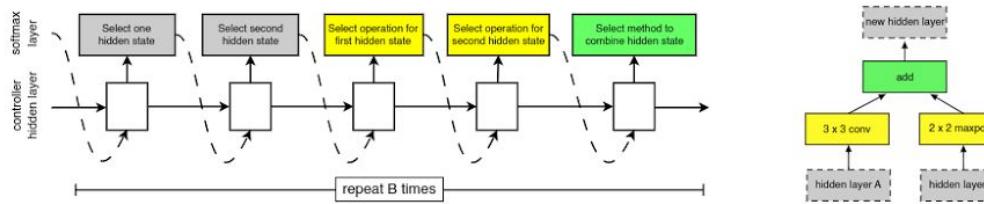
Neural Architecture Search for ImageNet



Neural Architecture Search for ImageNet

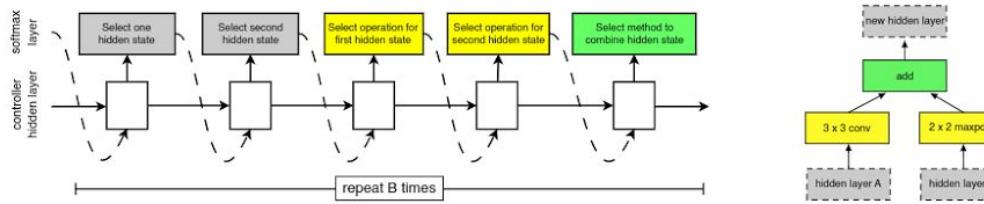


How the cell was found



- Step 1.** Select a hidden state from h or from the set of hidden states created in previous blocks.
- Step 2.** Select a second hidden state from the same options as in Step 1.
- Step 3.** Select an operation to apply to the hidden state selected in Step 1.
- Step 4.** Select an operation to apply to the hidden state selected in Step 2.
- Step 5.** Select a method to combine the outputs of Step 3 and 4 to create a new hidden state.

How the cell was found



Step 1. Select a hidden state from h or from the set of hidden states created in previous blocks.

Step 2. Select a second hidden state from the same options as in Step 1.

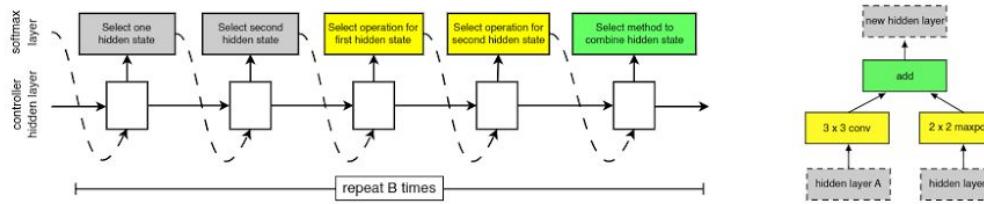
Step 3. Select an operation to apply to the hidden state selected in Step 1.

Step 4. Select an operation to apply to the hidden state selected in Step 2.

Step 5. Select a method to combine the outputs of Step 3 and 4 to create a new hidden state.

- identity
- 1x7 then 7x1 convolution
- 3x3 average pooling
- 5x5 max pooling
- 1x1 convolution
- 3x3 depthwise separable convolution
- 7x7 depthwise separable convolution
- 1x3 then 3x1 convolution
- 3x3 dilated convolution
- 3x3 max pooling
- 7x7 max pooling
- 3x3 convolution
- 5x5 depthwise separable convolution

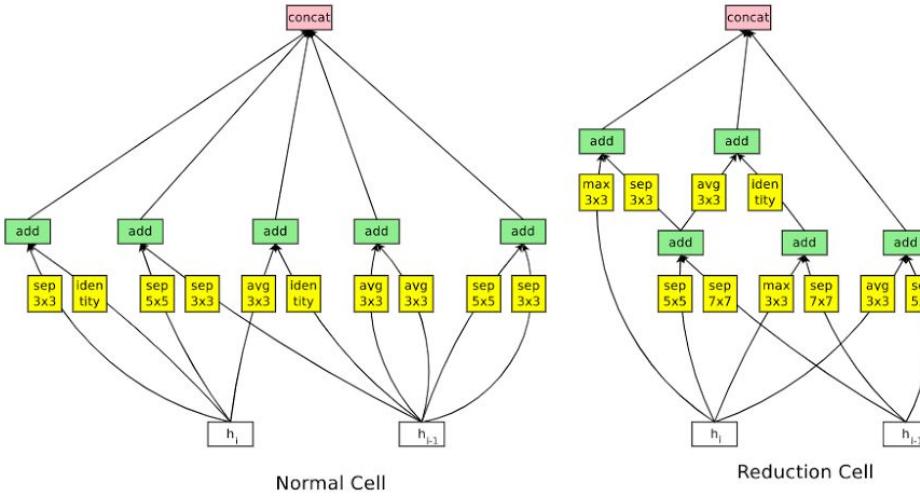
How the cell was found



- Step 1.** Select a hidden state from h or from the set of hidden states created in previous blocks.
- Step 2.** Select a second hidden state from the same options as in Step 1.
- Step 3.** Select an operation to apply to the hidden state selected in Step 1.
- Step 4.** Select an operation to apply to the hidden state selected in Step 2.
- Step 5.** Select a method to combine the outputs of Step 3 and 4 to create a new hidden state.

- 1. Elementwise addition
- 2. Concatenation along the filter dimension

The cell again



Performance of cell on CIFAR-10

Model	Depth	# parameters	Error rate (%)
DenseNet ($L = 40, k = 12$) [23]	40	1.0M	5.24
DenseNet ($L = 100, k = 12$) [23]	100	7.0M	4.10
DenseNet ($L = 100, k = 24$) [23]	100	27.2M	3.74
DenseNet-BC ($L = 100, k = 40$) [24]	190	25.6M	3.46
Shake-Shake 26 2x32d [17]	26	2.9M	3.55
Shake-Shake 26 2x96d [17]	26	26.2M	2.86
NAS v1 no stride or pooling [61]	15	4.2M	5.50
NAS v2 predicting strides [61]	20	2.5M	6.01
NAS v3 max pooling [61]	39	7.1M	4.47
NAS v3 max pooling + more filters [61]	39	37.4M	3.65
NASNet-A	N=6	3.3M	3.41

Performance of Neural Architecture Search and other state-of-the-art models on CIFAR-10.

Performance of cell on ImageNet

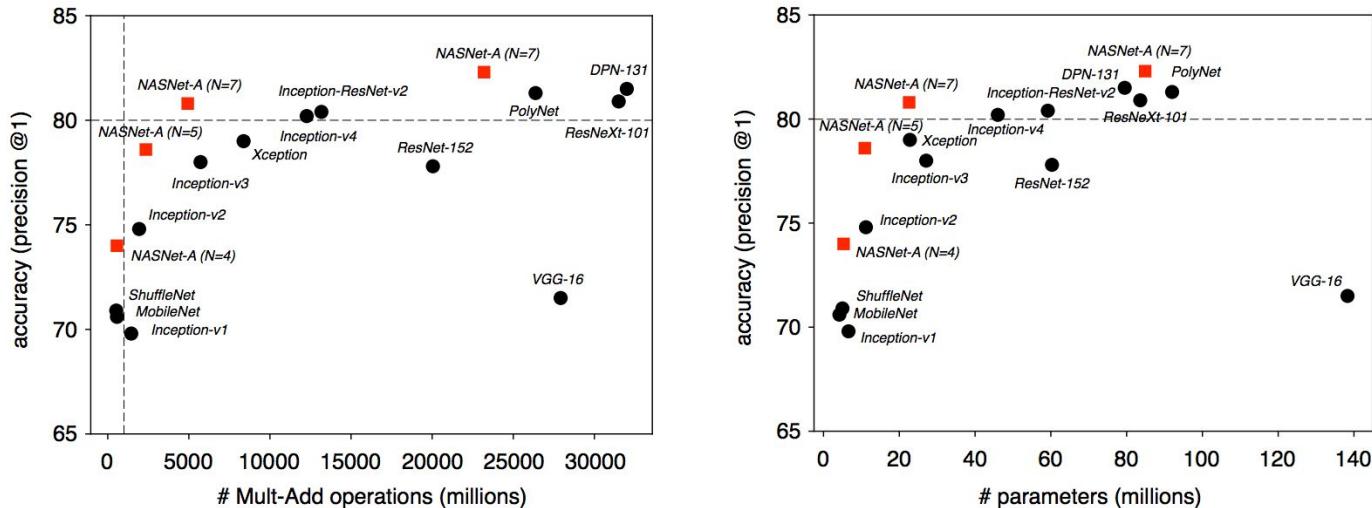


Figure 5: Accuracy versus computational demand (left) and number of parameters (right) across top performing CNN architectures on ImageNet 2012 ILSVRC challenge prediction task (compiled as of July 2017). Computational demand is measured in the number of floating-point multiply-add operations to process a single image. Black circles indicate previously published work and red squares highlight our proposed models. Vertical dashed line indicates 1 billion multiply-add operations. Horizontal dashed line indicates 80% precision@1 prediction accuracy.

References

- [Learning Transferrable Architectures for Scalable Image Recognition](#)
- [Neural Architecture Search with Reinforcement Learning](#)
- [Neural Optimizer Search with Reinforcement Learning](#)
- [Neural Combinatorial Optimization with Reinforcement Learning](#)