

Large-Scale Evolution of Image Classifiers

Esteban Real¹ Sherry Moore¹ Andrew Selle¹ Saurabh Saxena¹
Yutaka Leon Suematsu² Jie Tan¹ Quoc V. Le¹ Alexey Kurakin¹

¹Google Brain

²Google Research

ICML, 2017

Presenter: Tianlu Wang

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Outline

1 Introduction

- Motivation
- Backgrounds

2 Related Work

- Neuro-evolution
- Non-evolutionary

3 Methods

- Algorithm Overview
- Encoding and Mutations
- More Details

4 Results

- Progress of experiments
- Comparisons
- Meta-parameters

5 Summary

Motivation

- AlexNet, GoogleNet, VGG, ResNet...
- Designing neural network architectures can be challenging
- Discover network architectures automatically

Outline

1 Introduction

- Motivation
- Backgrounds

2 Related Work

- Neuro-evolution
- Non-evolutionary

3 Methods

- Algorithm Overview
- Encoding and Mutations
- More Details

4 Results

- Progress of experiments
- Comparisons
- Meta-parameters

5 Summary

Backgrounds

- Achievements: evolution algorithm outputs a **fully-trained** model with no human participation
- Drawbacks: significant computation
- Image classification, CIFAR-10, CIFAR-100

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - **Neuro-evolution**
 - Non-evolutionary
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Neuro-evolution

- Weight evolution: back propagation
- Weight and architecture: *NEAT algorithm* (node and connection)

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - **Non-evolutionary**
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Non-evolutionary

- Bayesian optimization
- Reinforcement learning
- Q-learning

Outline

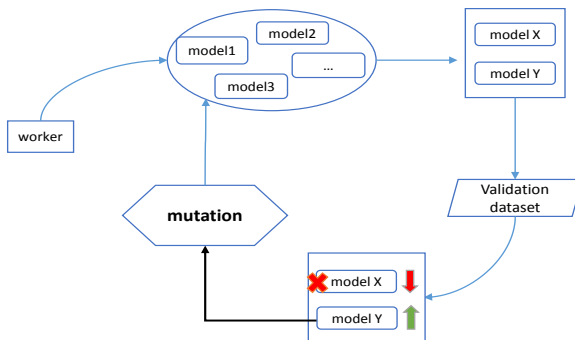
- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 **Methods**
 - **Algorithm Overview**
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Algorithm Overview

- Input: a population of models, each model is a **trained single-layer nonconvolutional model** with `learning_rate = 0.1`
- Measurement: accuracy on validation dataset

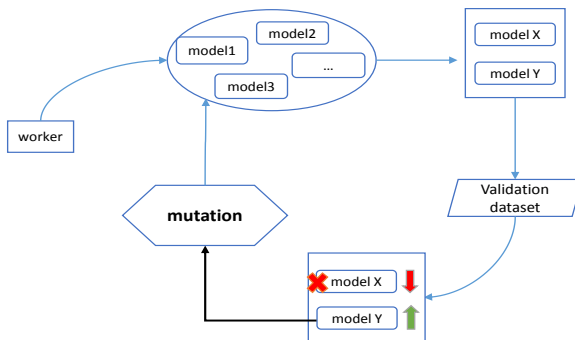
Algorithm Overview

- Input: a population of models, each model is a **trained single-layer nonconvolutional model** with `learning_rate = 0.1`
- Measurement: accuracy on validation dataset



Algorithm Overview

- Input: a population of models, each model is a **trained single-layer nonconvolutional model** with $\text{learning_rate} = 0.1$
- Measurement: accuracy on validation dataset



- When to stop?

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 **Methods**
 - Algorithm Overview
 - **Encoding and Mutations**
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Model Encoding

Individual model is encoded as a graph:

- Vertices
 - rank-3 tensor($\text{image_width} * \text{image_height} * \text{channels}$)
 - activations(batch normalization with ReLU or plain linear layer)
- Edges
 - Identity connections
 - Convolutions

Model Encoding

Individual model is encoded as a graph:

- Vertices
 - rank-3 tensor(image_width * image_height * channels)
 - activations(batch normalization with ReLU or plain linear layer)
- Edges
 - Identity connections
 - Convolutions

Inconsistent input:

- pick and keep primary one
- reshape(interpolation/truncation/padding) non-primary ones

Mutations

The worker picks a mutation at random from a set:

- ALTER-LEARNING-RATE
- IDENTITY (effectively means keep training)
- RESET-WEIGHTS
- INSERT/REMOVE CONVOLUTION
- ALTER-STRIDE
- ALTER-NUMBER-OF-CHANNELS
- FILTER-SIZE
- INSERT-ONE-TO-ONE
- INSERT/REMOVE SKIP

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 **Methods**
 - Algorithm Overview
 - Encoding and Mutations
 - **More Details**
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

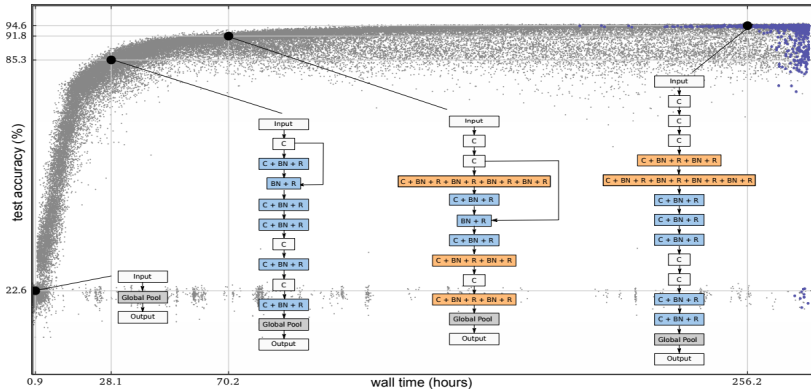
More Details

- Poor initial conditions(12th slide)
- 45,000 training; 5,000 validation; 10000 test
- SGD with momentum of 0.9, batch size 50, weight decay 0.0001
- Computation cost: floating-point operations
- Inherit parameters' weights whenever possible

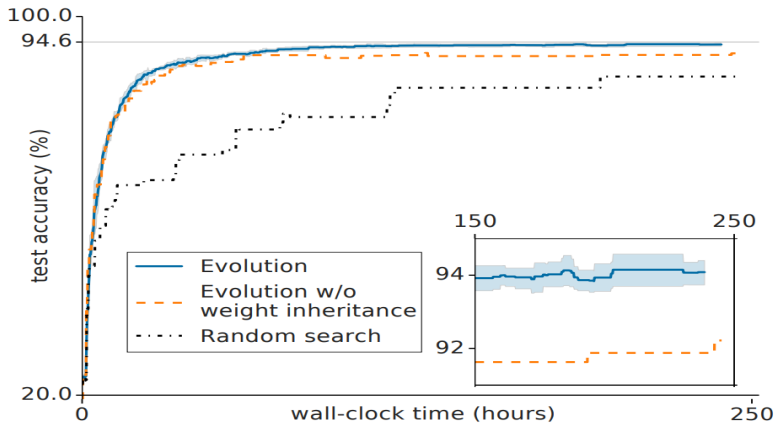
Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - Meta-parameters
- 5 Summary

Progress of an evolution experiment



Repeatability of results and controls



Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - **Comparisons**
 - Meta-parameters
- 5 Summary

Compared to hand-designed networks

STUDY	PARAMS.	C10+	C100+	REACHABLE?
MAXOUT (GOODFELLOW ET AL., 2013)	–	90.7%	61.4%	No
NETWORK IN NETWORK (LIN ET AL., 2013)	–	91.2%	–	No
ALL-CNN (SPRINGENBERG ET AL., 2014)	1.3 M	92.8%	66.3%	YES
DEEPLY SUPERVISED (LEE ET AL., 2015)	–	92.0%	65.4%	No
HIGHWAY (SRIVASTAVA ET AL., 2015)	2.3 M	92.3%	67.6%	No
RESNET (HE ET AL., 2016)	1.7 M	93.4%	72.8% [†]	YES
EVOLUTION (OURS)	5.4 M 40.4 M	94.6%	77.0%	N/A
WIDE RESNET 28-10 (ZAGORUYKO & KOMODAKIS, 2016)	36.5 M	96.0%	80.0%	YES
WIDE RESNET 40-10+D/O (ZAGORUYKO & KOMODAKIS, 2016)	50.7 M	96.2%	81.7%	No
DENSENET (HUANG ET AL., 2016A)	25.6 M	96.7%	82.8%	No

Compared to auto-discovered networks

STUDY	STARTING POINT	CONSTRAINTS	POST-PROCESSING	PARAMS.	C10+	C100+
BAYESIAN (SNOEK ET AL., 2012)	3 LAYERS	FIXED ARCHITECTURE, NO SKIPS	NONE	–	90.5%	–
Q-LEARNING (BAKER ET AL., 2016)	–	DISCRETE PARAMS., MAX. NUM. LAYERS, NO SKIPS	TUNE, RETRAIN	11.2 M	93.1%	72.9%
RL (ZOPH & LE, 2016)	20 LAYERS, 50% SKIPS	DISCRETE PARAMS., EXACTLY 20 LAYERS	SMALL GRID SEARCH, RETRAIN	2.5 M	94.0%	–
RL (ZOPH & LE, 2016)	39 LAYERS, 2 POOL LAYERS AT 13 AND 26, 50% SKIPS	DISCRETE PARAMS., EXACTLY 39 LAYERS, 2 POOL LAYERS AT 13 AND 26	ADD MORE FILTERS, SMALL GRID SEARCH, RETRAIN	37.0 M	96.4%	–
EVOLUTION (OURS)	SINGLE LAYER, ZERO CONVS.	POWER-OF-2 STRIDES	NONE	5.4 M 40.4 M ENSEMB.	94.6% 95.6%	77.0%

Outline

- 1 Introduction
 - Motivation
 - Backgrounds
- 2 Related Work
 - Neuro-evolution
 - Non-evolutionary
- 3 Methods
 - Algorithm Overview
 - Encoding and Mutations
 - More Details
- 4 Results
 - Progress of experiments
 - Comparisons
 - **Meta-parameters**
- 5 Summary

Improve the method

- Large population size
- More training steps
- Increase mutation rate
- Reset all weights

Summary

- Neuro-evolution starts from trivial initial conditions and yields fully trained models
- Construct large, accurate networks for two challenging and popular image classification benchmarks
- Large search space and high computation cost