

Improving Image Classifiers for Small Datasets by Learning Rate Adaptations

S Mishra, T Yamasaki, H Imaizumi

MVA 2019

Oral Session 3

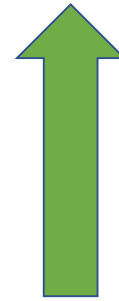
May 29, 2019



Scope

Rapid improvements in image classification tasks

- Larger better & detailed datasets
- Faster hardware resources
- Better architectures



However (the ugly truth)!

- More iterations to SOTA
- Longer train time
- Higher costs
- Small dataset reliability low



Scope

Deployment costs can adversely impact individuals or smaller groups.

SOLUTION?

- Organic combination of proven techniques, field tested on benchmark datasets.
- Optimization by learning rate (ν) adaptations.
- Transfer modus-operandi to smaller, untested data.
- Ensure repeatability.

Baseline

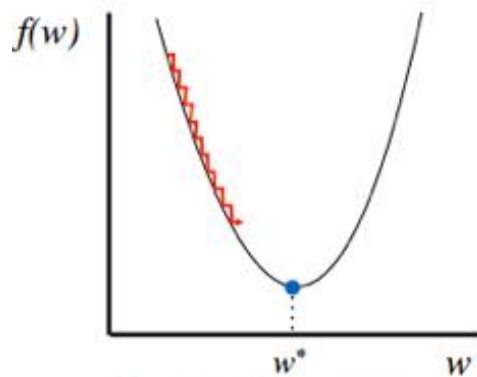
- Multi-class classification on CIFAR-10
- Test candidate architectures of increasing size/complexity
 - Resnet-34, ResNet-50, ResNet-101, ResNet-152
 - DenseNet161
- Baseline Performance
 - 5:1 split, Early stopping, lower LR restarts
 - BCE with logits loss
 - Train to 90%+ validation accuracy mark

Baseline

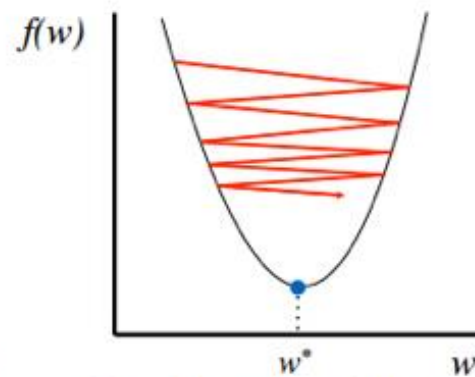
Architecture	Accuracy (Top-1)	Time (s)
ResNet 34	90.36%	17,757
ResNet-50	90.54%	34,039
ResNet-101	90.71%	60,639
ResNet-152	90.68%	91,888
DenseNet-161	93.02%	54,628

Learning Rate range-test

Learning rate under-exploited by Monotonic change.



Too small: converge
very slowly



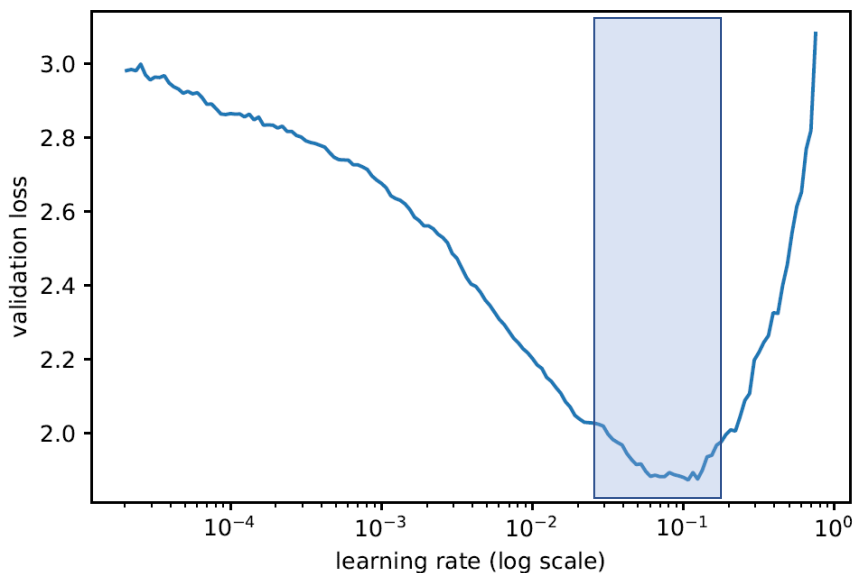
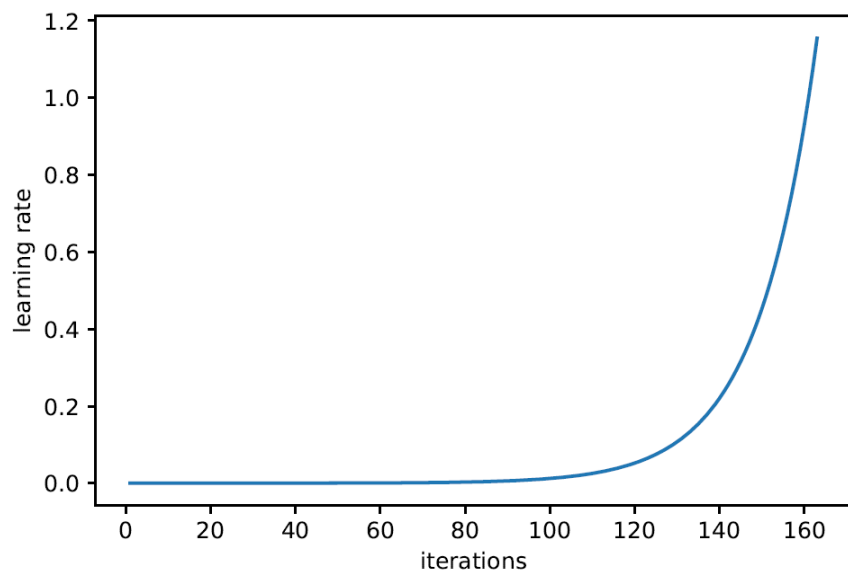
Too big: overshoot and
even diverge

Courtesy:

Cyclical Learning rates for training NN, L. Smith [2017]

Deep Learning, S. Verma et al. 2018

Learning Rate range-test



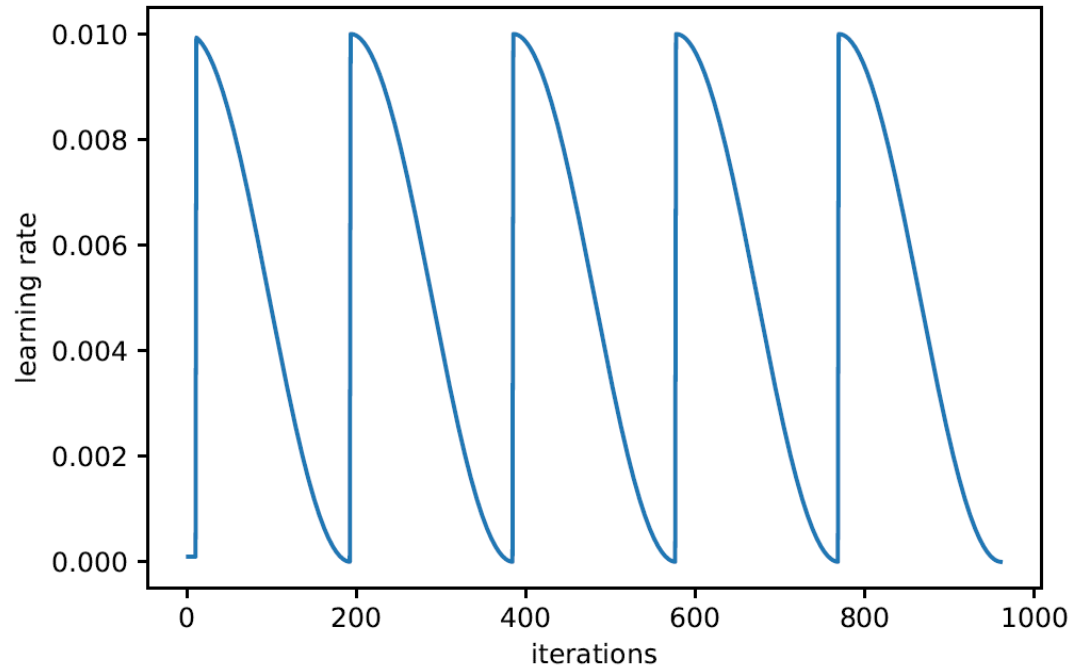
Steadily increase the LR and observe the Cross entropy loss
Test several mini-batches to see a point of inflexion

SGD-R

Avoid monotonicity

Cycle LR ν by Cosine scheduling function

$$\nu(t) = \frac{1}{2} \left(1 + \nu \cos \left(\frac{t\pi}{T} \right) \right) + \epsilon$$

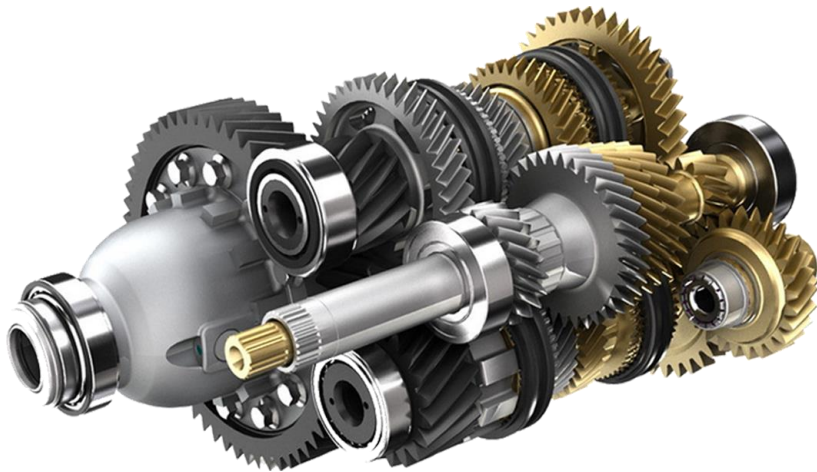


Perform initial coarse fit by tuning the last (or last few) FC layer

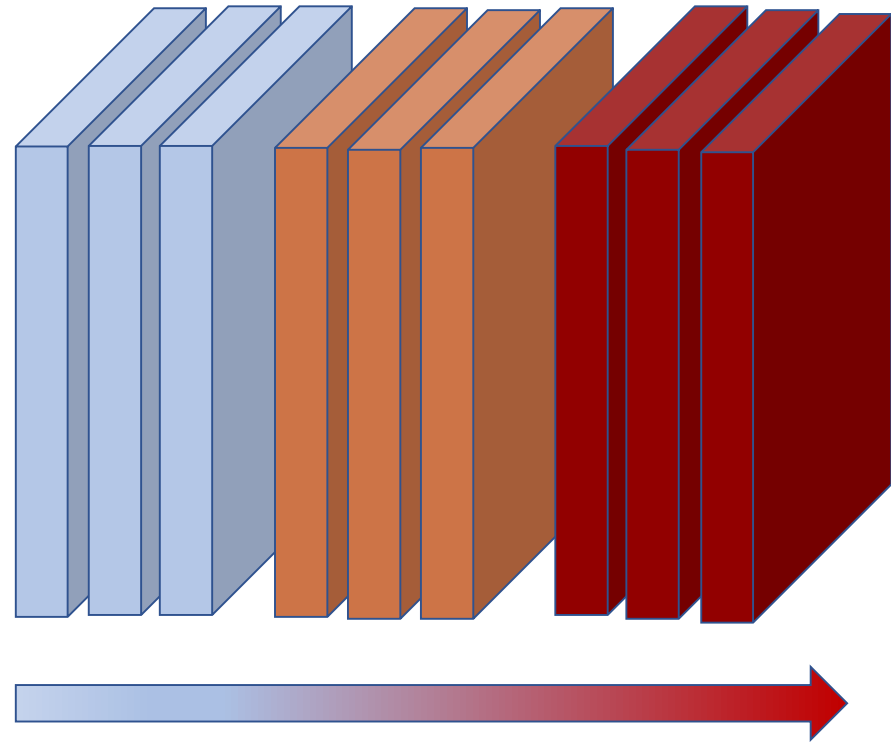
Courtesy:

SGD with Warm restarts, Loschilov [2017]

Differential learning



Gear-box need not spin all gears equally!

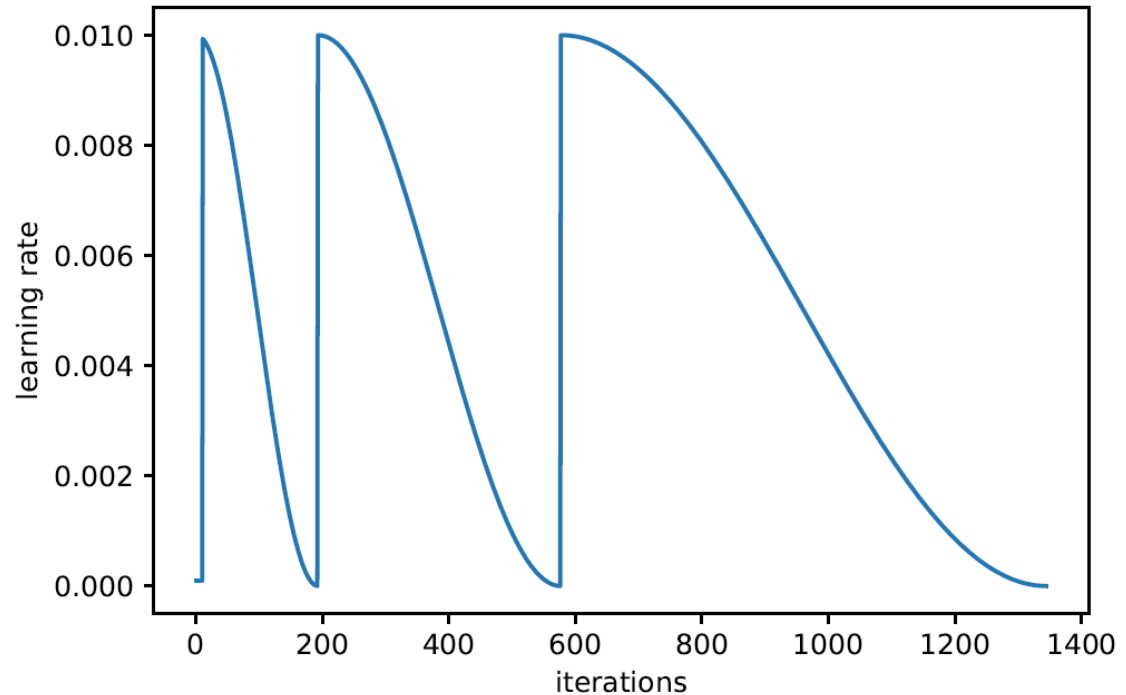


Reduce computational overhead by assigning different learning rates.

SGD-R + CLM + DLR

Cycle Length Multiply
by integral powers of
2.

Cover larger # epochs
per anneal cycle

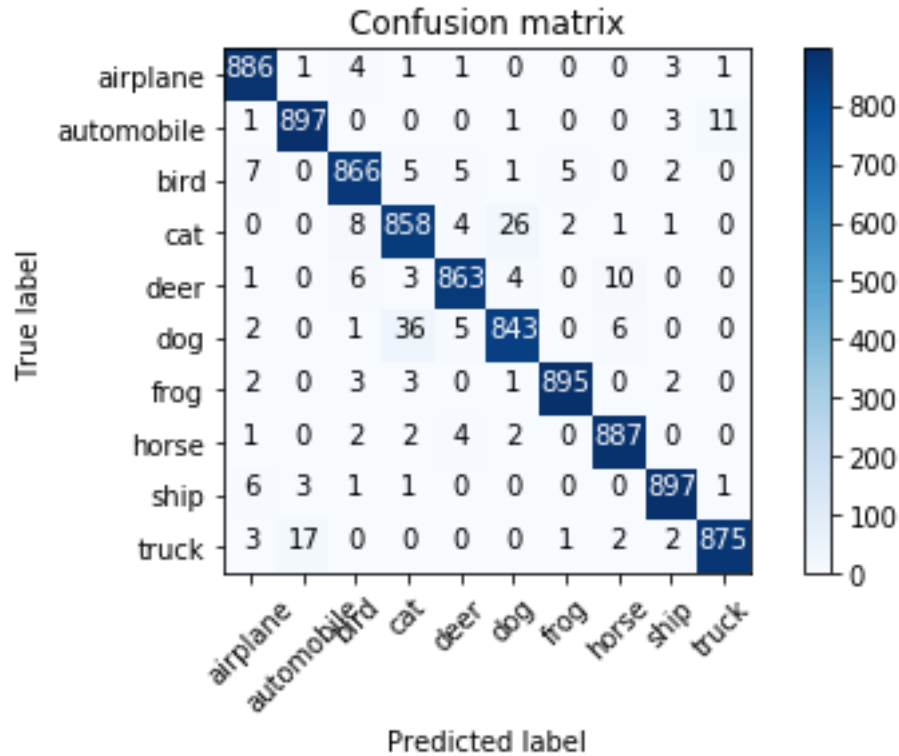


Perform tighter fit over all layers

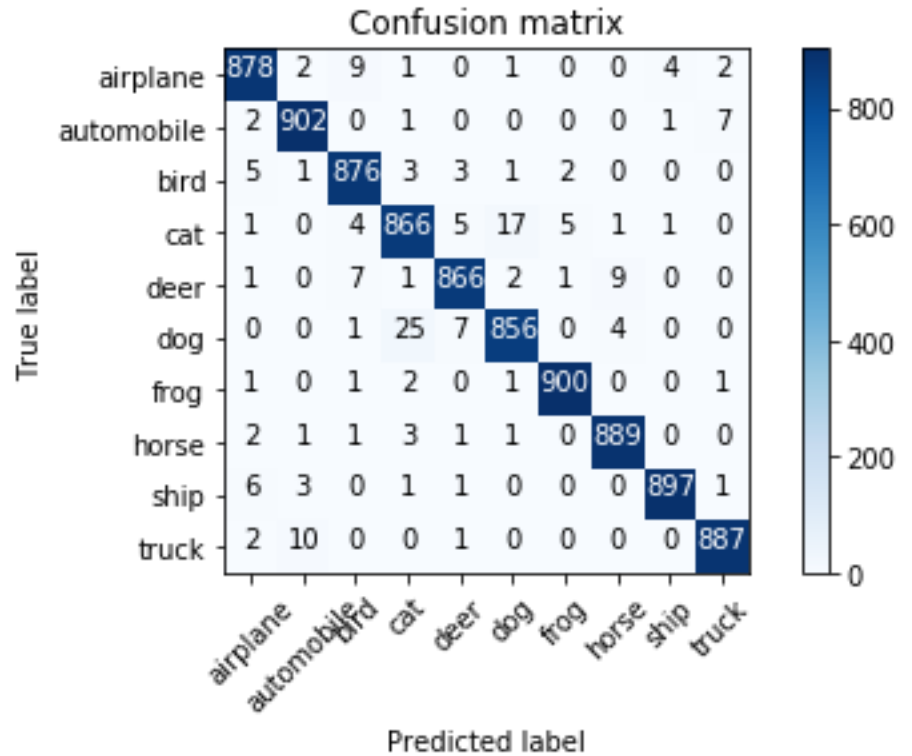
Results

Architecture	Accuracy (Top-1)	Time (s)	η
ResNet 34	96.84%	9,565	1.84
ResNet-50	96.82%	11,817	2.88
ResNet-101	97.61%	6,673	9.09
ResNet-152	97.78%	9,012	10.2
DenseNet-161	97.15%	7,195	7.59

Results

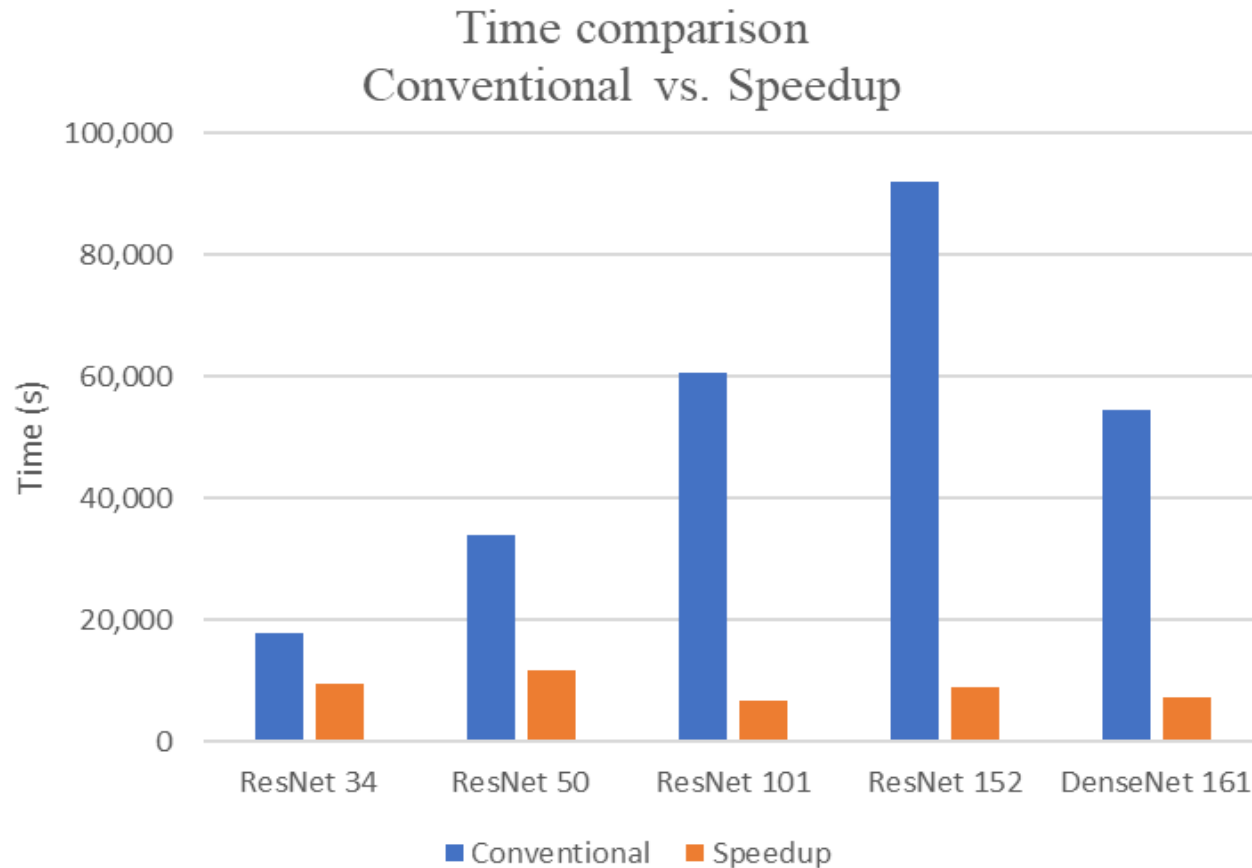


DenseNet 161



ResNet 152

Results



Higher dividends when architecture size grows larger.
Possible by offsetting the computation overhead by DLR

Application

Dermoscopic detection as diagnostic aid:

- 10 common classes & homogenous
- Requires rapid retrain
- Severe class imbalance

- >10% acc. improvement
- 3.1x - 5.7x in time
- \$ 6.9 (old) vs \$ 0.7 (new)

Confusion matrix

acne	186	2	1	2	1	2	4	1	0	1
alopecia	1	143	0	1	0	2	0	2	0	0
blister	2	0	117	6	13	4	1	0	7	0
crust	3	0	8	128	1	0	0	1	9	0
erythema	6	0	4	4	108	2	11	5	2	8
leuko	4	3	2	0	3	127	7	0	1	2
macula	3	0	1	4	14	3	115	7	1	2
tumor	0	0	2	4	2	2	6	173	11	0
ulcer	1	0	0	5	5	0	1	18	170	0
wheal	0	0	0	0	7	0	0	0	0	143

Actual

Predicted

Application - Samples



Acne



Alopecia



Crust



Tumor