



RETHINKING OF DEEP MODELS PARAMETERS WITH RESPECT TO DATA DISTRIBUTION

#### **Shitala Prasad**

Scientist, Visual Intelligence

Co-authors: Yiqun Li, Dongyun Lin, Dong Sheng, Oo Zaw Min I2R / A\*Star Singapore



# Agenda

Introduction & Motivations	3-5 Page
Methodology	6-12 Page
Experimental Results	13-15 Page
Conclusion	16 Page



## **Introduction and Motivation**

The revolutionization of deep learning in field of computer vision has changed the way the visual learning takes place.

The performance of deep learning models are driven by various parameters but to tune all of them every time, for every given dataset, is a heuristic practice.

In this paper, unlike the common practice of decaying the learning rate, we propose a step-wise training strategy where the learning rate and the batch size are tuned based on the dataset size.



## **Introduction and Motivation**

Deep learning has always such questions:

- Can the learning performance be improved without additional data?
- Can accuracy be increased for the same architecture with the same dataset?
- Can data distribution strategy boosts the accuracy and reduces the training cost?
- Can it reduce network over-fitting?



## **Introduction and Motivation**

The proposed training strategy randomly splits the data into small subsets and incrementally train the classifier by updating the dataset with other subsets in a progressive step-wise fashion

✓ a simple step-wise hyper-parameter tuning strategy to boost the network classification performance without using any additional data, consistently valued for several state-of-the-art image classification network architectures and reduces the overall training cost by ~40%





# **Standard Baseline Approach**

#### In fine-tuning deep models, mainly:





# **Standard Baseline Approach**

#### In fine-tuning deep models, mainly:



Epoch



## **Standard Baseline Approach**

In fine-tuning deep models, other than batch size and learning rate, dataset size is also a key factor:





# **Proposed Training Approach**





# **Proposed Training Approach**





# **Proposed Training Approach**

### Batch size + learning rate + dataset size: Train from **scratch**



Epoch

12



# **Experimental Results**





### Datasets

Pindonan Datacat	Res2Net-50							
Birdshap Dataset	top-1	top-5	Epochs					
with set O (baseline)	61.7901	83.3790	100	_				
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	62.1391	85.150	20 each	]				
Res2NeXt-50		0	-					
Birdsnap Dataset	top-1	top-5	Epochs					
with set O (baseline)	61.9900	84.0025	100	_				
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	62.4611	84.9875	20 each					
					Res2Net-50			
	Food-101 Dataset		Pre-train	ed=False	alse Pre-trained=True		Enochs	Weight Shared
_			top-1	top-5	top-1	top-5	Lpoths	Weight Shared
_	with set O	(baseline)	53.5584	81.0297	85.7347	96.8832	100	$\checkmark$
_	with se	et A	12.0832	32.2733	75.5881	93.2198	20	
	with se	et B	15.7109	38.2733	80.2851	95.1723	20	
	with se	et C	17.0257	40.8119	82.0371	95.8020	20	×
	with se	et D	18.2574	42.7287	83.9247	96.6059	20	
_	with se	et O	37.8614	67.5485	86.2574	97.1406	20	
	with se	et A	13.0218	33.4020	75.5881	93.2198		
	21	$\stackrel{b}{\rightarrow}$ B	21.8970	48.4634	78.7247	94.3842	20 each	
	31	<sup>b</sup> → C	30.2851	59.0416	80.6337	95.0416	20 Cach	×
	41	b D	38.2772	66.9782	81.8614	95.5762		
	b	→ O	54.9901	81.5010	86.7089	97.2554		



### Datasets

Clubbed together: 1. covid mask images: https://www.kaggle.com/danielferrazcampos/face-mask-images

2. mask dataset: https://www.kaggle.com/ahmetfurkandemr/mask-datasets-v1

3. COVID19 mask image dataset: https://github.com/UniversalDataTool/coronavirus-mask-image-dataset

COVID 10 mask nomask Dataset	Res2Net-50				
COVID-19_mask-nomask Dataset	top-1	Method			
with set O (baseline) $(b=16)$	98.90				
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	$97.50 \xrightarrow{2b} 98.75 \xrightarrow{3b} 99.06 \xrightarrow{4b} 98.90 \xrightarrow{5b} 99.67$	$\mathscr{G} \approx \frac{\lambda_i M_i}{h_i}$			
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	$97.50 \xrightarrow{2b} 98.75 \xrightarrow{3b} 99.06 \xrightarrow{4b} 98.90 \xrightarrow{b} 99.84$	01			
with set O (baseline) $(b=24)$	99.53				
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	97.65 $\xrightarrow{2b}$ 97.18 $\xrightarrow{3b}$ 97.34 $\xrightarrow{4b}$ 97.50 $\xrightarrow{b}$ 100				
with set O (baseline) $(b=24)$	99.53	$(R_{a}, \lambda_{i}M_{i})$			
with set $A \xrightarrow{2b} B \xrightarrow{3b} C \xrightarrow{4b} D \xrightarrow{b} O$	$97.65 \xrightarrow{2b} 97.65 \xrightarrow{3b} 97.81 \xrightarrow{4b} 98.28 \xrightarrow{b} 97.97$	$\mathcal{G} \approx \frac{1}{b}$			
with set O (baseline) $(h-24)$	00.52				
with set $O(baseline)(b=24)$	99.35	$\mathcal{A} \sim \lambda M_i$			
with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O	99.55 97.65 $\xrightarrow{2b}$ 97.50 $\xrightarrow{3b}$ 98.12 $\xrightarrow{4b}$ 98.75 $\xrightarrow{b}$ 99.69	$\mathscr{G} \approx \frac{\lambda M_i}{b_i}$			
with set O (baseline) $(b=24)$ with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set O (baseline) $(b=24)$	$99.53$ $97.65 \xrightarrow{2b} 97.50 \xrightarrow{3b} 98.12 \xrightarrow{4b} 98.75 \xrightarrow{b} 99.69$ $99.53$	$\mathscr{G} \approx \frac{\lambda M_i}{b_i}$			
	COVID-19_mask-nomask Dataset with set O (baseline) (b=16) with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set O (baseline) (b=24) with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set O (baseline) (b=24) with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set O (baseline) (b=24) with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O with set A $\xrightarrow{2b}$ B $\xrightarrow{3b}$ C $\xrightarrow{4b}$ D $\xrightarrow{b}$ O	Res2ret-30top-1with set O (baseline) (b=16)98.90with set A $\stackrel{2b}{\rightarrow}$ B $\stackrel{3b}{\rightarrow}$ C $\stackrel{4b}{\rightarrow}$ D $\stackrel{b}{\rightarrow}$ O97.50 $\stackrel{2b}{\rightarrow}$ 98.75 $\stackrel{3b}{\rightarrow}$ 99.06 $\stackrel{4b}{\rightarrow}$ 98.90 $\stackrel{5b}{\rightarrow}$ 99.67with set A $\stackrel{2b}{\rightarrow}$ B $\stackrel{3b}{\rightarrow}$ C $\stackrel{4b}{\rightarrow}$ D $\stackrel{b}{\rightarrow}$ O97.50 $\stackrel{2b}{\rightarrow}$ 98.75 $\stackrel{3b}{\rightarrow}$ 99.06 $\stackrel{4b}{\rightarrow}$ 98.90 $\stackrel{b}{\rightarrow}$ 99.84with set O (baseline) (b=24)with set O (baseline) (b=24)with set O (baseline) (b=24)97.65 $\stackrel{2b}{\rightarrow}$ 97.18 $\stackrel{3b}{\rightarrow}$ 97.34 $\stackrel{4b}{\rightarrow}$ 97.50 $\stackrel{b}{\rightarrow}$ 100with set O (baseline) (b=24)99.53with set A $\stackrel{2b}{\rightarrow}$ B $\stackrel{3b}{\rightarrow}$ C $\stackrel{4b}{\rightarrow}$ D $\stackrel{b}{\rightarrow}$ O97.65 $\stackrel{2b}{\rightarrow}$ 97.65 $\stackrel{3b}{\rightarrow}$ 97.81 $\stackrel{4b}{\rightarrow}$ 98.28 $\stackrel{b}{\rightarrow}$ 97.97with set O (baseline) (b=24)99.53with set A $\stackrel{2b}{\rightarrow}$ B $\stackrel{3b}{\rightarrow}$ C $\stackrel{4b}{\rightarrow}$ D $\stackrel{b}{\rightarrow}$ O97.65 $\stackrel{2b}{\rightarrow}$ 97.65 $\stackrel{3b}{\rightarrow}$ 97.81 $\stackrel{4b}{\rightarrow}$ 98.28 $\stackrel{b}{\rightarrow}$ 97.97			



COVID-19 mask-nomask dataset. Details: train (1050 mask, 1459 no-mask) and test (264 mask, 375 no-mask) images.



### Take back

The  $M_i$  trained network weights are uses as the new initializer for Mi+1 subset training that boost the learning curve without saturating Analyzes a close interrelation between M, b and  $\lambda$  and propose a step-wise training to up rise the performance instead of traditional baseline training without performing any change in the network architecture The proposed stepwise training reduces the risk of over-fitting by adopting different b and also reduces the training cost by 40%

In future, we would like to explore other aspects of CV such as object detection and segmentation where annotation is the biggest challenge





# **THANK YOU**

www.a-star.edu.sg