Putting Al on a Diet: TinyML and Efficient Deep Learning

Song Han **Assistant Professor** Massachusetts Institute of Technology



https://songhan.mit.edu





Challenges for Deep Learning:

computation-hungry, data-hungry

AlphaGo: 1920 CPUs and 280 GPUs, \$3000 electric bill per game





1417



Cloud AI: increases TCO & carbon Mobile AI: drains battery Tiny AI: can't fit memory











Lots of efforts to collect/label







TinyML and Efficient Deep Learning



A lot of computation A lot of carbon









Many engineers

A lot of data





TinyML and Efficient Deep Learning Three aspects



less computation less carbon









fewer engineers











Make AI run Fast and Efficiently with Limited Hardware Resource

Large Neural Networks





Model Compression & TinyML



Han, Mao, Dally, Deep Compression, ICLR'16, best paper award

Low-Power Hardware





Deep Compression





Original ResNet-50

with Deep Compression



Han, Mao, Dally, Deep Compression, ICLR'16, best paper award

Make AI run Fast and Efficiently with Limited Hardware Resource



weights (32 bit float)				cluster index (2 bit uint)				centre		
2.09	-0.98	1.48	0.09		3	0	2	1	3:	2.0
0.05	-0.14	-1.08	2.12	cluster	1	1	0	3	2:	1.5
-0.91	1.92	0	-1.03		0	3	1	0	1:	0.0
1.87	0	1.53	1.49		3	1	2	2	0:	-1.0

100MB

6MB 17x compression





Pruning & Sparsity

Increased attention Publications per Year since 2015



Optimal Brain Damage

Yann Le Cun, John S. Denker and Sara A. Solla AT&T Bell Laboratories, Holmdel, N. J. 07733





the curve credits to NVIDIA



Pruning & Sparsity applied by industry:





EIE, Han et al, ISCA'16



ESE, Han et al, FPGA'17



SpArch, Zhang et al, HPCA'20 SpAtten, Wang et al, HPCA'21

Шï



Dense Matrix



(FP32)



2:4 sparsity in A100 GPU 2X peak performance, 1.5X measured BERT speedup



Reduce model complexity by 5x to 50x with minimal accuracy impact. Deep Compression takes the performance of your AI inference to the next level.









TinyML and Efficient Deep Learning

- AutoML and NAS
 - Once-for-all Network [ICLR'19]
- Less Computation
 - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight]
 - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Less Training Data





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]





TinyML and Efficient Deep Learning

- AutoML and NAS - Once-for-all Network [ICLR'19]
- Less Computation - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight] - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Less Training Data





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]







Manual Design: black magic, nutoriously hard to tune



Automatic Design: Synthesize NNs to fit latency/accuracy/memory constraints

EDA tool : Circuit = NAS : Neural net









AutoML and Neural Architecture Search



integrated by Intel OpenVINO

Combine them together: APQ: Joint Search for Network Architecture, Pruning and Quantization Policy [CVPR'20]

1st Place of Visual Wakeup Words (VWW) Challenge, host by Google @CVPR2019

push-button solution for efficient NN design





Proxyless Neural Architecture Search

[ICLR 2019]

integrated by Facebook PyTorch and Amazon AutoGluon





How to design efficient NN models for diverse hardware platforms?



GPU, 16GB







Lots of hand tuning for different devices! # layers, # channels, resolution, kernel size...









2019

200

For training iterations: forward-backward();



The design cost is calculated under the assumption of using MobileNet-v2.

Design Cost (GPU hours)







2019





The design cost is calculated under the assumption of using MnasNet. [1] Tan, Mingxing, et al. "Mnasnet: Platform-aware neural architecture search for mobile." *CVPR*. 2019.

Design Cost (GPU hours)





Diverse Hardware Platforms







The design cost is calculated under the assumption of using MnasNet. [1] Tan, Mingxing, et al. "Mnasnet: Platform-aware neural architecture search for mobile." CVPR. 2019.

Design Cost (GPU hours)

160K







Cloud AI (10^{12} FLOPS)

Plii



Mobile AI (10^9 FLOPS)





Diverse Hardware Platforms



Tiny AI (10^6 FLOPS)

Design Cost (GPU hours)

160K

1600K

[1] Tan, Mingxing, et al. "Mnasnet: Platform-aware neural architecture search for mobile." CVPR. 2019.







Cloud AI (10^{12} FLOPS)



Mobile AI (10^9 FLOPS)

For many devices:	Des			
For search episodes: // meta controller	40K → 11			
For training iterations:				
forward-backward(); Expensive!!				
If good_model: break;				
For post-search training iterations:				
forward-backward(); Expensive!!				



Diverse Hardware Platforms



Tiny AI (10^6 FLOPS)

sign Cost (GPU hours)

1.4k lbs CO₂ emission

 $160K \rightarrow 45.4k$ lbs CO₂ emission

1600K \rightarrow **454.4k** lbs CO₂ emission

Strubell, Emma, et al. "Energy and policy considerations for deep learning in NLP." ACL. 2019.









Today's NAS is too expensive We need Green Al

Common carbon footprint benchmarks

in lbs of CO2 equivalent



Chart: MIT Technology Review · Source: Strubell et al. · Created with Datawrapper









Artificial intelligence / Machine learning

Training a single Al model can emit as much carbon as five cars in their lifetimes

Deep learning has a terrible carbon footprint.

by Karen Hao

Evolved Transformer

ICML'19, ACL'19









Train once, get many Reduce the design cost Fit diverse hardware constraints





Users may have the high-end phones and also low-end phones; We want to be inclusive for users who have low-end phones.

But, it's expensive to design NN of difference sizes



iPhone 11

A13 Bionic, 2019





A12 Bionic, 2018









Train once, get many Reduce the design cost Fit diverse hardware constraints







Train once, get many Reduce the design cost Fit diverse hardware constraints







Weight-sharing, Decouple Training and Search

Conventional NAS





=>



Once-for-All:





Once for All Network: Train 10¹⁹ networks at the same time



Reducing the carbon footprint of artificial intelligence MIT system cuts the energy required for training and running neural networks.

Rob Matheson | MIT News Office April 23, 2020



or Search Browse MIT researchers have developed a new automated AI system with improved computational efficiency and a much smaller carbon footprint. The researchers' system trains one large neural network comprising many pretrained subnetworks of different sizes MCU that can be tailored to diverse hardware platforms without retraining. Image: MIT News, based on figures courtesy of the researchers

Press Inquiries

RELATED

Human brain activates sparsely

OFA network containts many child networks that are sparsely activated

Child networks share the weights with the Once-for-All network, trained joinly





Make the OFA network elastic: How to prevent large & small sub networks from interfering with each other?











progressive shrinking: kernel size











progressive shrinking: kernel size



progressive shrinking: # layers











progressive shrinking: kernel size







progressive shrinking: # channels







Subnets sampled from OFA SuperNet outperforms training from scratch!



















Sub-networks under various architecture configurations D: depth, W: width, K: kernel size

Progressive shrinking consistently improves accuracy of sub-networks on ImageNet. \bullet





Once-for-All, ICLR'20



Once-for-All Network Evolutionary Architecture Search

Latency measurement

ШiГ



IANIA



Once-for-All Network Evolutionary Architecture Search Q











Once-for-All Network Train only once, handle diverse hardware constraints

PHIT



Google Pixel1 CPU Latency (ms)









1417

Train only once, generate the entire Pareto curve

OFA



24 Once-for-All, ICLR'20

MobileNetV3





Automatically Synthesize Neural Nets given different latency/accuracy constraints

Drag the bar to target different latency. ← Slide left for faster and less accurate models \rightarrow Slide right for slower but more accurate models

Specialize for 35(ms) on Note10 Device, top1 78.47(%)



https://hanlab.mit.edu/projects/ofa/demo/





Automatically Synthesize Neural Nets

given different latency/accuracy constraints



Platform: Intel(R) Xeon(R) Platinum 8280 CPU @ 2.70GHz (28 cores), HT ON, turbo ON, Total Memory: 192GB (12 slots / 16 GB / 2933 MHz) Framework: PyTorch v.1.5.0 with Intel®MKL-DNN 0.14 enabled, data type FLOAT 32 System: BIOS SE5C620.86B.02.01.0008.031920191559 / Ubuntu 18.04.4 LTS, Kernel 4.15.0-88-generic










Once-for-All Network





MobileNetV3 MobileNetV2 \diamond



Once-for-All Network

Consistently outperforms human baselines, world-record on MLPerf Turn-key solution for co-design



Plii

1.078M images per second on eight A100 GPUs



Award Winning Technology



CPU detection **FPGA** detection



5th Low-Power Computer Vision Challenge

Challenge



Visual Wake Words on TF-lite



Visual Wake Words Challenge @CVPR 2019 **AI Driving Olympics @ICRA 2021**



CPU classification CPU detection



DSP Recognition

4th Low-Power Computer Vision

3th Low-Power Computer Vision Challenge

3D Semantic Segmentation

NLP track Language Model

MicroNet Challenge @NeurIPS 2019





Industry Adoption

MLPerf

Once-for-All (OFA) Network adopbed by Alibaba received a world-record in the

open division of <u>MLPerf Inference Benchmark</u>, achieving 1.078M images per Cond on eight A100 GPUs



maxim integrated

Once-for-All (OFA) Network adopted by Maxim Integrated provides 6% accuracy increase in image recognition and 2% accuracy increase in speech command recognition, with >100x energy efficiency compared to Cortex-M4.



Proxyless Neural Architecture Search, an efficient neural architecture search algorithm with light-weight model for mobile AI is integrated by AWS AutoGluon and Facebook PyTorch.



NVIDIA

HAQ: Hardware-Aware Automated Quantization with Mixed Precision is integrated by Intel OpenVINO Toolkit. Efficiently search over the bitwidth space for mixedprecision machine learning inference (2, 4, 8 bits)

integrated by **NVIDIA** for video classification.



TSM: Temporal Shift Module for Efficient Video Understanding is





Anycost GAN

original







* Status: ready

Demo:



AnyCost GAN, CVPR'21











Anycost GAN



MACs:

Smaller, faster child networks are nested in larger ones



100% 1.0x reduction





GAN Compression Accelerating Horse2zebra by GAN Compression



Demo:



Large Neural Networks



1

Original CycleGAN; FLOPs: 56.8G; FPS: 12.1; FID: 61.5

GAN Compression; FLOPs: 3.50G (16.2x); FPS: 40.0 (3.3x); FID: 53.6



Small Neural Networks



Efficient NLP

Once-for-all Transformer









IIANI_A

Efficient NLP

On WMT'14 En-Fr Task





HAT: Hardware-Aware Transformers, ACL 2020

HAT, ACL'20 SpAtten, HPCA'21















OFA for NLP

Once-for-all Transformer



1

2

3

Feb 10, 2021 🛛 <u>Full story</u> Share: 🎔 🕇









I-IANI_AI=

NAAS: Neural Accelerator Architecture Search





Neural Accelerator Architecture Search, DAC'21













TinyML and Efficient Deep Learning

- AutoML and NAS - Once-for-all Network [ICLR'19]
- TinyML
 - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight] - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Data-Efficiency





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]





MCUNet: Tiny Deep Learning on IoT Devices

Chuang Gan³ Wei-Ming Chen^{1,2} Yujun Lin¹ John Cohn³ Song Han¹ Ji Lin¹

²National Taiwan University ³MIT-IBM Watson AI Lab ¹MIT



NeurIPS 2020 (spotlight)





Background: The Era of AloT on Microcontrollers (MCUs)

• Low-cost, low-power







Background: The Era of AloT on Microcontrollers (MCUs)

Low-cost, low-power



#Units (Billion)



Rapid growth







Background: The Era of AloT on Microcontrollers (MCUs)

Low-cost, low-power MCU



Wide applications

Smart Retail



Personalized Healthcare Smart Manufacturing





Rapid growth



Autonomous Driving







Deep Learning Going "Tiny"



Cloud Al => ResNet

Data centers, connection required, privacy issue

Smartphones process locally

- The future belongs to Tiny AI.
- Much cheaper, much smaller, almost everywhere in our lives.
- democratize AI and extend the applications of deep learning.









Mobile Al => MobileNet

Tiny Al => MCUNet

IoT devices, cheap, small, low-power, rapid growth

- There are billions of IoT devices around the world based on microcontrollers - If we can enable powerful AI algorithms on those IoT devices, we can greatly





Challenge: Memory Too Small to Hold DNN









Activation is the bottleneck, not parameters





(all with ~70% ImageNet Top-1)





MCUNet: TinyNAS + TinyEngine

Search space design is crucial for NAS performance There is no prior expertise on MCU model design













Reduce Both Model Size and Activation Size







MCUNet: TinyNAS+TinyEngine Co-design



- TinyNAS:
 - Re-design the design space
 - Latency-aware
 - Energy-aware
 - Once-for-all Network: train once, get many



• TinyEngine:

- Co-design, specialization
- Run time => Compile time
- Graph optimizations
- Memory-aware scheduling
- Low-precision
- Assembly-level optimizations





MCUNet: TinyNAS+TinyEngine

ImageNet classification on STM32F746 MCU (320kB SRAM, 1MB Flash) lacksquare



* scaled down version: width multiplier 0.3, input resolution 80









MCUNet: TinyNAS+TinyEngine

ImageNet classification on STM32F746 MCU (**320kB SRAM**, **1MB Flash**) \bullet

Baseline (MbV2*+CMSIS) **System-only** (MbV2**+TinyEngine) **Model-only** (TinyNAS+CMSIS)

ImageNet Top1: 35%

* scaled down version: width multiplier 0.3, input resolution 80 ** scaled down version: width multiplier 0.35, input resolution 144











MCUNet: TinyNAS+TinyEngine

ImageNet classification on STM32F746 MCU (**320kB SRAM**, **1MB Flash**) \bullet



* scaled down version: width multiplier 0.3, input resolution 80 ** scaled down version: width multiplier 0.35, input resolution 144









Bring Al to IoT Devices



face/mask detection, person detection, VWW on STM32F746 (1MB Flash)





tinyml.mit.edu



TinyEngine: Memory Saving









OctoML





TinyEngine: Speedup







Once-for-All Network + MCUNet

Design specialized models for different MCUs

75 70 65 63.5 62.0 60 55 50 STM32F412 STM32F746 (<u>256kB</u>/1MB) (<u>320kB</u>/1MB) (<u>512kB</u>/1MB) (SRAM/Flash)









ImageNet Top-1 Accuracy (%)









Visual Wake Words (VWW)









360







Audio Wake Words (Speech Commands)





yes













Demo: Visual Wake Words on MCU

MBv1+TFLite-Micro PERSON fps:2.873

75% accuracy, fps: 2.9







87% accuracy, fps: 7.3

- Detecting if there is person
- STM32F746
- 320KB SRAM
- 1MB Flash
- ARM Cortex-M7 @216MHz



Demo: Face Mask Detection on MCU





- Detecting faces & masks
- STM32F746
- 320KB SRAM
- 1MB Flash
- ARM Cortex-M7 @216MHz





Demo: Person Detection on MCU





- Detecting persons
- STM32F746
- 320KB SRAM
- 1MB Flash
- ARM Cortex-M7 @216MHz






NuScenes LiDAR Segmentation Challenge



nuScenes Dataset LiDAR Segmentation Challenge @ICRA

1 Spinning LiDAR (20 Hz, 32 channels) 16 40,000 annotated frames 35,0 1,000 driving sequences 100m x 1



16 semantic classes 35,000 points per frame 100m x 100m x 20m spatial range

Point-Voxel Convolution (PVConv)



Low-Cost Point-Based Branch Provides Fine Details of the Scene





Comparisons with **A** Point-Based Models

Voxel-based branch conducts convolution over a regular grid representation, which resolves the issue of random memory access.

Comparisons with :::: Voxel-Based Models

Point-based branch captures high-resolution information efficiently, which resolves the issue of large memory footprint.

Huge Improvements on Safety-Critical Small Objects





Rank 1: SPVCNN++

Authors: Zhijian Liu, Haotian Tang, Kevin Shao, Song Han Affiliation: MIT HAN Lab

Method: Use point cloud segmentation network composed of multiple Sparse Point-Voxel Convolutions. More computationally heavy variant of SPVCNN, longer training, Lovasz-Softmax-loss with standard cross entropy loss. Lidar-only. 13 FPS.

Link: https://spvnas.mit.edu, "Searching Efficient 3D Architectures with Sparse **Point-Voxel Convolution**"

Mean IOU: 81.1% +0.8%

Award: Best lidarseg submission

Previous awards: 4th ranked lidarseg submission in NeurIPS 2020





NUSC		Motional	Home r	nuScenes	nulmages	Re	sources	- Task	s ~	About ~	Sign	l
			Method							Metrics		
Date	Name		Modalities				Map data	External data	mIOU	fwlOU	FPS (Hz)	
			Any			•	All 👻	All 👻				
2021-05-27	SPVCNN++	Ours (2021)	Lidar				no	no	0.811	0.910	n/a	
2021-05-26	GU-Net	UISEE	Lidar				no	no	0.803	0.910	n/a	
2021-05-26	2D3DNet	Google S	dustry ^{Camera, Lidar} ubmissions				no	yes	0.800	0.901	n/a	
2020-12-09	AF2S3Net	Huawei	Lidar				no	no	0.783 (+2.8)	0.885 (+2.5)	3.7 (x3.5)	
2020-12-09	Cylinder3D++	CUHK	Lidar				no	no	0.779	0.899	10.6	
2020-12-08	CPFusion	Horizon Roboti	CS Camera, Lidar, Ra	dar			no	no	0.777	0.892	n/a	
2020-12-09	SPVNAS	Ours (2020)	Lidar				no	no	0.774	0.897	15.9	

Ranks 1st on the NuScenes Leaderboard Best LiDARSeg Submission @ ICRA 2021







Efficient Point Cloud



AR/VR: a whole backpack of computer



Self-driving: a whole trunk of GPU



Mobile phone: limited battery







MinkowskiNet: 3.4 FPS



SPVNAS (Ours): 9.1 FPS



Efficient Point Cloud

State-of-the-Art Accuracy

Approach	Paper	Code	mloU	Classes (IoU)
SPVNAS	Å		67.0	
TORNADONet	×		63.1	
KPRNet	2		63.1	
Cylinder3D	لح	0	61.8	

First Place on SemanticKITTI Leaderboard (as of Fall 2020)





PVCNN, NeurIPS'19 SPVNAS, ECCV'20 FastLidarNet, ICRA'21

Real-Time Inference on Edge Devices



Project of the Month" by NVIDIA Jetson Community

5.7x memory reduction

Outdoor Scene Semantic Segmentation



2.7x measured speedup **7.6x** computation reduction



Efficient Point Cloud



3D LiDAR Sensor



3D Point Cloud: 2M points/s

Demo:







SemAlign: Annotation-Free Camera-LiDAR Calibration with Semantic Alignment Loss

After **10** seconds of optimization by SemAlign:















SemAlign: Annotation-Free Camera-LiDAR Calibration with Semantic Alignment Loss











SemAlign: Annotation-Free Camera-LiDAR Calibration with Semantic Alignment Loss [IROS'21]











SemAlign: Annotation-Free Camera-LiDAR Calibration with Semantic Alignment Loss





Specialized Accelerator Design: PointAcc [Micro'21]

Significant Speedups and Energy Savings over CPU, GPU and TPU



Lower Latency and Higher Accuracy with PointAcc (Edge)

Point Cloud Networks

MIT Driverless Accuracy: 95.0% Range: 8 meters Latency: 2 ms/object

PVCNN (Ours) Accuracy: 99.9% Range: 12 meters Latency: 1.25 ms/object















Rank 1: SPVCNN++

Authors: Zhijian Liu, Haotian Tang, Kevin Shao, Song Han Affiliation: MIT HAN Lab

Method: Use point cloud segmentation network composed of multiple Sparse Point-Voxel Convolutions. More computationally heavy variant of SPVCNN, longer training, Lovasz-Softmax-loss with standard cross entropy loss. Lidar-only. 13 FPS.

Link: https://spvnas.mit.edu, "Searching Efficient 3D Architectures with Sparse **Point-Voxel Convolution**"

Mean IOU: 81.1% +0.8%

Award: Best lidarseg submission

Previous awards: 4th ranked lidarseg submission in NeurIPS 2020





NUSC		Motional	Home	nuScenes	nulmages	Re	esources 、	 Task 	s v 🛛	About ~	Sign	1
			Method			-				Metrics		l
Date	Name		Modalities	ŝ			Map data	External data	mIOU	fwlOU	FPS (Hz)	
			Any			*	All -	All 👻				
2021-05-27	SPVCNN++	Ours (2021)	Lidar				no	no	0.811	0.910	n/a	
2021-05-26	GU-Net	Industry	Lidar				no	no	0.803	0.910	n/a	
2021-05-26	2D3DNet	Submissions	Camera, Lid	ar			no	yes	0.800	0.901	n/a	
2020-12-09	AF2S3Net	Winner (2020)	Lidar				no	no	0.783 (+2.8)	0.885 (+2.5)	3.7 (x3.5)	
2020-12-09	Cylinder3D++		Lidar				no	no	0.779	0.899	10.6	
2020-12-08	CPFusion		Camera, Lidar, I	Radar			no	no	0.777	0.892	n/a	
2020-12-09	SPVNAS	Ours (2020)	Lidar				no	no	0.774	0.897	15.9	

Ranks 1st on the NuScenes Leaderboard Best LiDARSeg Submission @ ICRA 2021







Efficient Video Recognition



Prediction: Moving something closer to something

Scaling Down Inference: **Inference on Low-power Edge Devices**



Devices	Jetson Nano		Jetson	n TX2	Rasp.	Note8	Pixel1	
	CPU	GPU	CPU	GPU	T			
FPS	20.9	74.6	27.5	117.6	14.4	29.0	21.1	
Power (watt)	4.8	4.5	5.6	5.8	3.8	- LE	D Bulb	



<u>TSM</u>, ICCV 2019



51

48

46

41

38

0

(%)

Compress the model by 6x, higher accuracy

Spotlight by IBM director Dario Gil @MIT-IBM Watson AI Lab's AI Research Week



Scaling Up Training: Large-Scale Distributed Training with 1536 GPUs

		Training Time	Accuracy	Peak GPU Performance
	1 Nodes (6 GPUs)	49h 50min	74.1%	46.5 TFLOP/s
	256 Nodes (1536 GPUs)	14min	74.0%	11,978 TFLOP/ s
Level!				













TinyML and Efficient Deep Learning

- AutoML and NAS - Once-for-all Network [ICLR'19]
- TinyML
 - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight]
 - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Data-Efficiency





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]





TinyTL: Reduce Memory, not Parameters for Efficient On-Device Learning

Han Cai¹ Chuang Gan² Ligeng Zhu¹ Song Han¹

¹MIT ²MIT-IBM Watson AI Lab



NeurIPS 2020





The Rise of AloT







IOT + AI = AIOT





Tiny Transfer Learning



Customization: Al systems need to continually adapt to new data collected from the sensors.

Tiny Transfer Learning

- Security: Data cannot leave devices because of security and regularization.
- TinyTL reduces the training memory from 300MB to 16MB

• Customization: Al systems need to continually adapt to new data collected from the sensors.

Training Memory is much Larger than Inference

- neural networks can easily exceed the limit.
- energy-efficient on-chip SRAM will significantly increase the energy cost.

• Edge devices have tight memory constraints. The training memory footprint of

• Edge devices are energy-constrained. Failing to fit the training process into the I-IANI_AI=

Training Memory is much Larger than Inference

- Edge devices have tight memory constraints. The training memory footprint of neural networks can easily exceed the limit.
- Edge devices are energy-constrained. Failing to fit the training process into the energy-efficient on-chip SRAM will significantly increase the energy cost. I-IANI_AI=

Efficient On-device Learning Requires Small Training Memory

neural networks can easily exceed the limit.

• Edge devices have tight memory constraints. The training memory footprint of

Activation is the Memory Bottleneck, not Parameters

• Activation is the main bottleneck for on-device learning, not parameters.

Activation is the Memory Bottleneck, not Parameters

- FLOPs, while the main bottleneck does not improve much.

Activation is the main bottleneck for on-device learning, not parameters.

• Previous methods focus on reducing the number of parameters or

- Full: Fine-tune the full network. Better accuracy but highly inefficient.

#Trainable Param (M)

• Last: Only fine-tune the last classifier head. Efficient but the capacity is limited.

- Full: Fine-tune the full network. Better accuracy but highly inefficient.

• Last: Only fine-tune the last classifier head. Efficient but the capacity is limited. • BN+Last: Fine-tune the BN layers and the last layer. Parameter-efficient.

Mudrakarta, Pramod Kaushik, et al. "K for the Price of 1: Parameter-efficient Multi-task and Transfer Learning." ICLR 2019.

IANI_A

Cars Top1 (%)

- Full: Fine-tune the full network. Better accuracy but highly inefficient.
- memory saving is limited.

Шiī

• Last: Only fine-tune the last classifier head. Efficient but the capacity is limited.

• BN+Last: Fine-tune the BN layers and the last layer. Parameter-efficient, but the

Cars Top1 (%)

- Full: Fine-tune the full network. Better accuracy but highly inefficient.
- memory saving is limited. Significant accuracy loss.

IIIiii

Parameter-efficiency does not directly translate to memory-efficiency

Memory Cost (MB)

• Last: Only fine-tune the last classifier head. Efficient but the capacity is limited.

• BN+Last: Fine-tune the BN layers and the last layer. Parameter-efficient, but the

TinyTL: Memory-Efficient Transfer Learning

Cars Top1 (%)

- Full: Fine-tune the full network. Better accuracy but highly inefficient.
- memory saving is limited. Significant accuracy loss.

• Last: Only fine-tune the last classifier head. Efficient but the capacity is limited.

• BN+Last: Fine-tune the BN layers and the last layer. Parameter-efficient, but the

• TinyTL: fine-tune bias only + lite residual learning: high accuracy, large memory saving

Weight update is Memory-expensive; **Bias update is Memory-efficient**

Forward:
$$\mathbf{a}_{i+1} = \mathbf{a}_i \mathbf{W}_i + \mathbf{b}_i$$

- Updating weights requires storing intermediate activations
- Updating biases does not

$$\frac{\partial L}{\partial \mathbf{b}_i} = \frac{\partial L}{\partial \mathbf{a}_{i+1}} = \frac{\partial L}{\partial \mathbf{a}_{i+2}} \mathbf{W}_{i+1}^T$$

TinyTL Idea 1: Fine-tune Bias Only y=Wa+b

- Updating biases does not

$$\frac{\partial L}{\partial \mathbf{b}_i} = \frac{\partial L}{\partial \mathbf{a}_{i+1}} = \frac{\partial L}{\partial \mathbf{a}_{i+2}} \mathbf{W}_{i+1}^T$$

Updating weights requires storing intermediate activations

Freeze weights, only fine-tune biases => save 12x memory

TinyTL Idea 2: Lite Residual Learning

- - (1/6 channel, 1/2 resolution, 2/3 depth)

Add lite residual modules (small memory overhead) to increase model capacity

	Stanford cars	51.6	49.2	47.3	48.8	42.4	51.6	53.4	55.0	
T			2. C r			Indal	o for	Diffa	cont 7	laeke
		MobileNetV2	ResNet-34	MnasNet	ProxylessNAS	MobileNetV3	ResNet-50	ResNet-101	Inception-v3	asns
	Aircraft	44.2	41.3	41.6	43.2	37.4	41.5	41.5	45.9	

The relative accuracy order between different pre-trained models changes significantly among ImageNet and the transfer learning datasets, which motivates personalized and specialized NN architecture for different downstream tasks.

tinyml.mit.edu

TinyTL + Once-for-All Network

difficult dataset

personalized model for different hardware and different tasks

easy dataset

tinyml.mit.edu

Memory Saving

• TinyTL provides 4.6x memory saving without accuracy loss.

- [1] Chatfield, Ken, et al. "Return of the devil in the details: Delving deep into convolutional nets." BMVC 2014.
- [2] Mudrakarta, Pramod Kaushik, et al. "K for the Price of 1: Parameter-efficient Multi-task and Transfer Learning." ICLR 2019.
- [3] Kornblith, Simon, Jonathon Shlens, and Quoc V. Le. "Do better imagenet models transfer better?." CVPR 2019.

TinyTL O Fine-tune BN+Last [1] × Fine-tune Last [2] + Fine-tune Full Network [3]

108

Memory Saving

• Fine-tune BN+Last [1] \times Fine-tune Last [2] + Fine-tune Full Network [3] ▼ TinyTL



- - [1] Chatfield, Ken, et al. "Return of the devil in the details: Delving deep into convolutional nets." BMVC 2014.
 - [2] Mudrakarta, Pramod Kaushik, et al. "K for the Price of 1: Parameter-efficient Multi-task and Transfer Learning." ICLR 2019.
 - [3] Kornblith, Simon, Jonathon Shlens, and Quoc V. Le. "Do better imagenet models transfer better?." CVPR 2019.



On different datasets, TinyTL provides up to 6.5x memory saving without accuracy loss.



109

TinyTL enables in-memory training

✓ TinyTL (batch size 1)

Mir



- TinyTL (tiny transfer learning) supports batch 1 training by group normalization.
- Together with the lite residual model, it further reduces the training memory cost to 16MB (fits L3 cache), enabling fitting the training process into cache, which is much more energy-efficient than training on DRAM.

- ✓ TinyTL + Fine-tune Full Network

300





TinyML and Efficient Deep Learning

- AutoML and NAS - Once-for-all Network [ICLR'19]
- TinyML
 - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight]
 - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Data-Efficiency





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]







FFHQ dataset: **70,000** selective post-processed human faces **ImageNet** dataset: **millions** of images from diverse categories

> Months or even years to collect the data, along with *prohibitive* annotation costs.







GANs Heavily Deteriorate Given Limited Data



100 images

160 images

1417

Generated samples of StyleGAN2 (Karras et al.) using only hundreds of images





GANs Heavily Deteriorate Given Limited Data

StyleGAN2 (baseline)











Discriminator Overfitting









#1 Approach: Augment reals only



Augment reals only: the same artifacts appear on the generated images.



Generated images



Artifacts from Color jittering



Artifacts from Translation



Artifacts from Cutout (DeVries et al.)



#2 Approach: Augment reals & fakes for D only



Augment D only: the unbalanced optimization cripples training.







#3 Approach: Differentiable Augmentation (Ours)



Our approach (DiffAugment): Augment reals + fakes for both D and G











Our Results

StyleGAN2 (baseline) + DiffAugment (ours)





Train GAN with only 100 Images



Generated samples of StyleGAN2 (baseline)

Without our technique:

With our technique:



Generated samples of StyleGAN2 + DiffAugment (ours)







Train GAN with only 100 Images



Without our technique:

With our technique:









Fine-Tuning vs. Ours









Train GAN with only 100 Images



Smooth interpolation, generalize well https://github.com/mit-han-lab/data-efficient-gans







Data-Efficient Deep Learning



Rare incidents





Privacy concerns

Under-represented subpopulations

Various factual and ethical reasons could cause limited data available. This research will help alleviate these limitations.







TinyML and Efficient Deep Learning

Three aspects: computation, engineers, data



A lot of computation A lot of carbon







Many engineers

A lot of data





TinyML and Efficient Deep Learning

- AutoML and NAS
 - Once-for-all Network [ICLR'19]
- TinyML
 - Inference: MCUNet for IoT Devices [NeurIPS'20, spotlight]
 - Training: Tiny On-Device Transfer Learning (TinyTL) [NeurIPS'20]
- Data-Efficient





- Differentiable Augmentation for Data-Efficient GAN Training [NeurIPS'20]





TinyML and Efficient AI



- Initiation <u>songhan.mit.edu</u> <u>tinyml.mit.edu</u> youtube.com/c/MITHANLab

itiative ear Algebra"

n be pruned to very sparse,

