

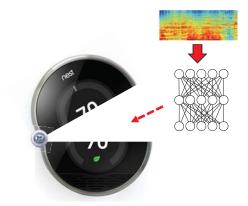
SAMSUNG Samsung Al Cambridge, UK

Machine Learning Systems: On-Device AI and Beyond

Nicholas D. Lane

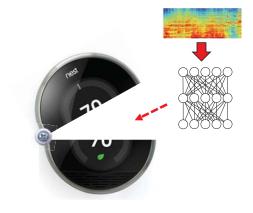
TECoSA Seminar Series Thursday December 3rd 2020 @niclane7 | http://niclane.org
ndl32@cam.ac.uk

Efficient ML Revolution



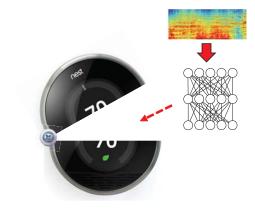


Efficient ML Revolution





Efficient ML Revolution



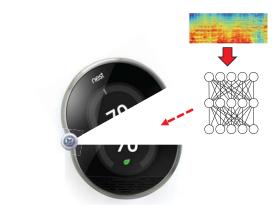


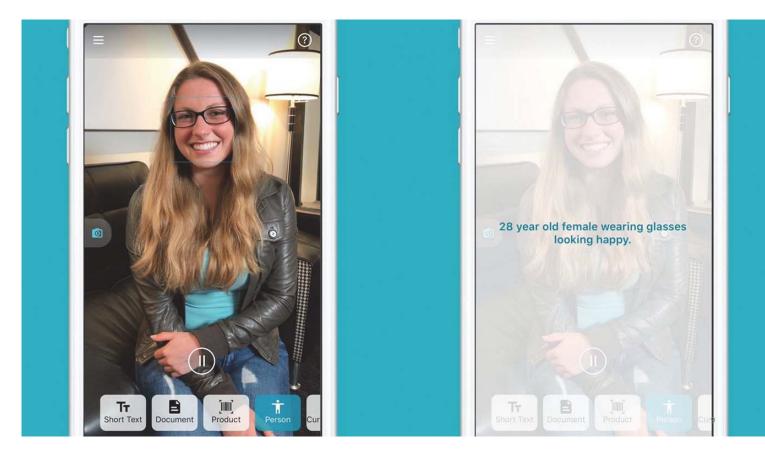
Efficient ML Revolution



- User Privacy
- Strongest algorithms for the hardest problems
- Robust devices without a
 network dependency
- Low energy and latency

Efficient ML Revolution

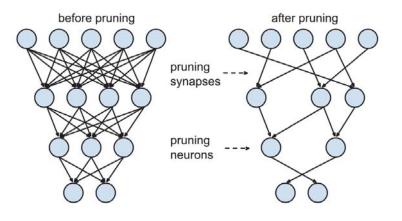




From Late '14: Efficient ML Revolution

2014 T	1 st Proof-of-Concept DL of Audio/IMU on Mobile [HotMobile '15]				
2015 -	DeepEar (1 st DSP-based DL General Audio Sensing) [UbiComp '15]	Community Innovations Algorithmic & Architecture Advances			
2016 -	DL Smartwatch Activity Recognition [WristSense '16] 1 st time: VGG executing directly on a commodity SmartWatch	 Node Pruning, Leverage Sparsity SqueezeNet (50x AlexNet reduction) Low Precision (8/4 bit), Binarization MobileNet, MCDNN, Custom Nets 			
2017	1 st time: Smartphone-scale DL on embedded processors (e.g., M0/M3) [SenSys '16] 1 st time: Multiple DL Vision Models on Wearable [MobiSys '17]	 Hardware Innovations Diannao and Cnvlutin2 Front-ends e.g., SNPE - Qualcomm TPU, FPGAs / Hybrids Analog from Digital Approaches Spiking H/W & Approx. Compute 			

Model Compression Example: Node Pruning

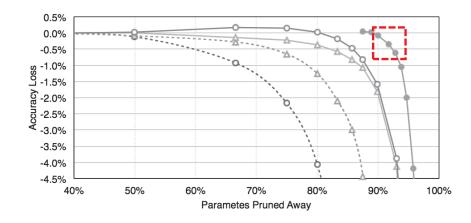


Related Methodologies

- Leveraging Sparsity
- Low Precision Results (8-bit etc.)
- Binarization of Networks
- MobileNet, MCDNN, Custom Nets
- .. and even hardware approaches

Song Han, Jeff Pool, John Tran, William J. Dally, "Learning both Weights and Connections for Efficient Neural Networks", NIPS 2015

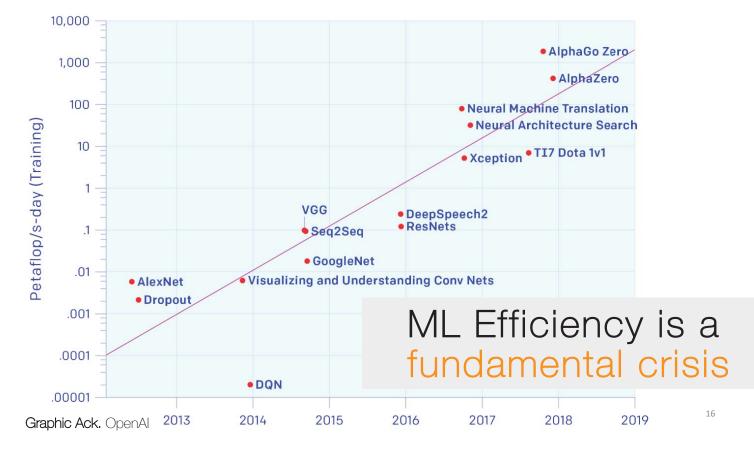
Model Compression Example: Node Pruning

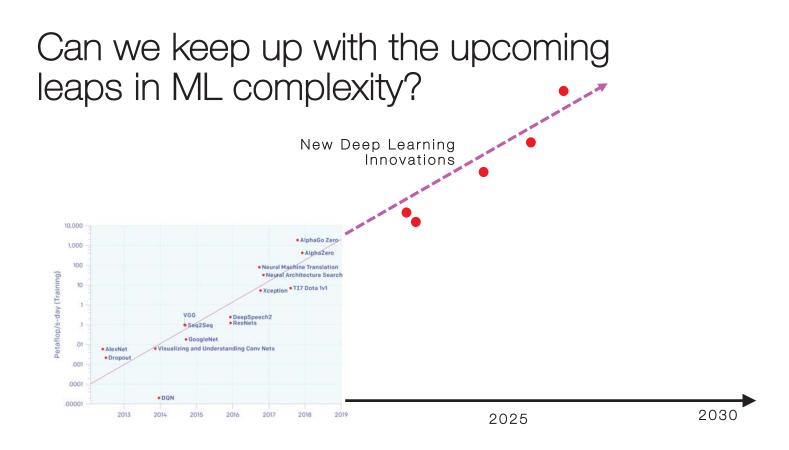


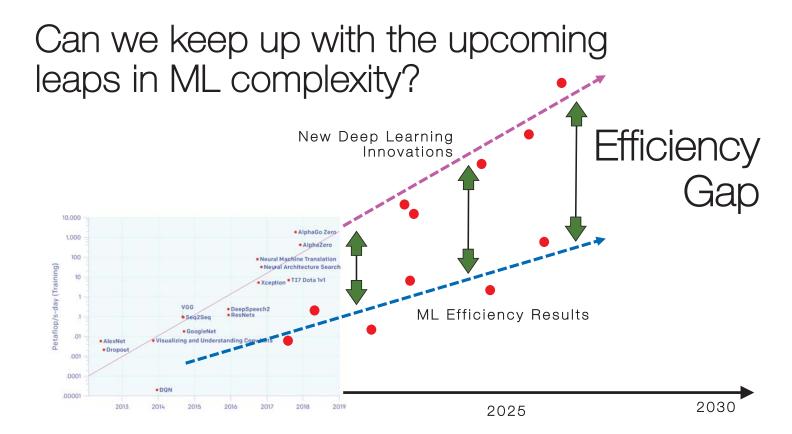
~50x gains tiny accuracy loss

Song Han, Jeff Pool, John Tran, William J. Dally, "Learning both Weights and Connections for Efficient Neural Networks", NIPS 2015

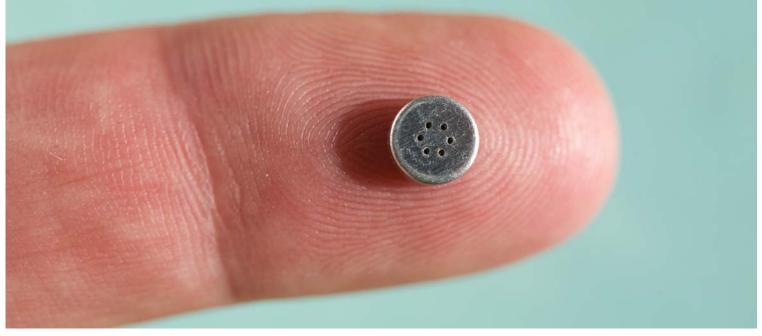








What will an Efficient ML device need look like in 2030?



What will an Efficient ML device need look like in 2030?

(1) Rich Powerful ML Tasks From Classification to Open-world Weakly-supervised Reasoning

(2) 100s of ML Models per device

(3) On-Device *Learning* is routine

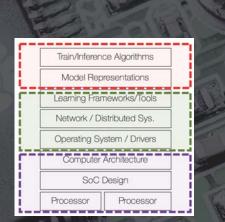


Fundamental Efficient ML Challenges

#1: Think (i.e., learn) Different

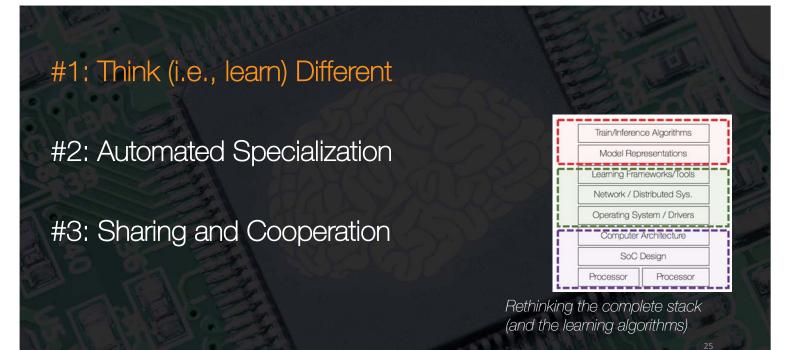
#2: Automated Specialization

#3: Sharing and Cooperation

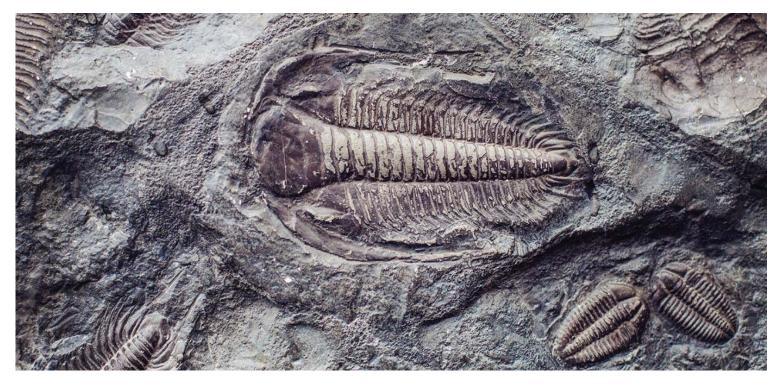


Rethinking the complete stack (and the learning algorithms)

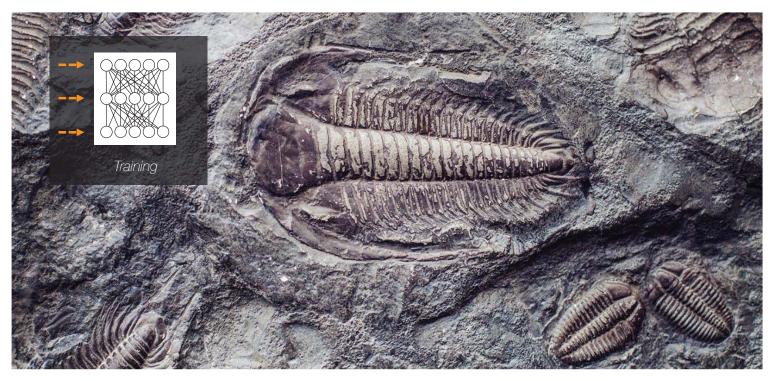
Fundamental Efficient ML Challenges



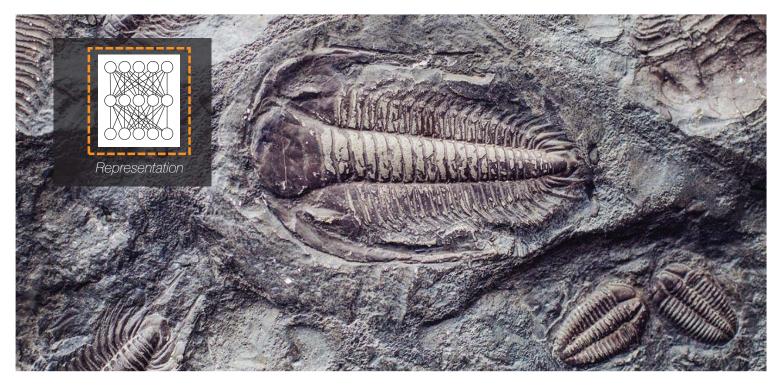
#1 Think (i.e., learn) Different



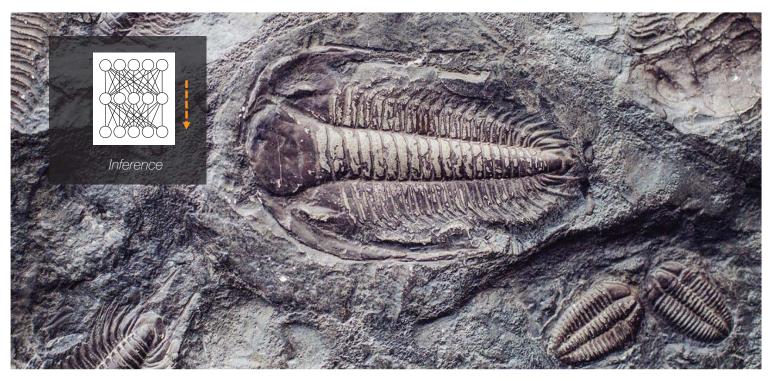
#1 Think (i.e., learn) Different



#1 Think (i.e., learn) Different

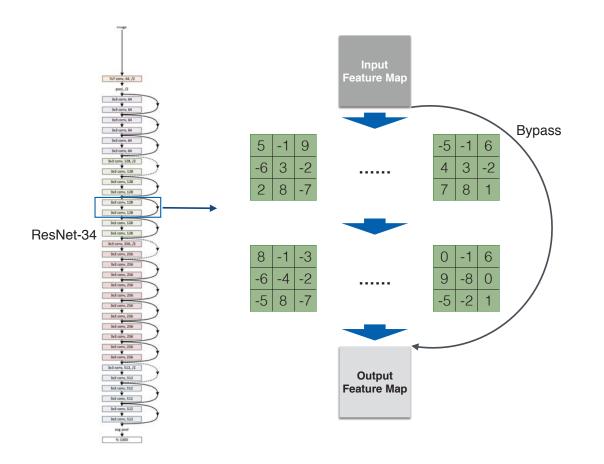


#1 Think (i.e., learn) Different

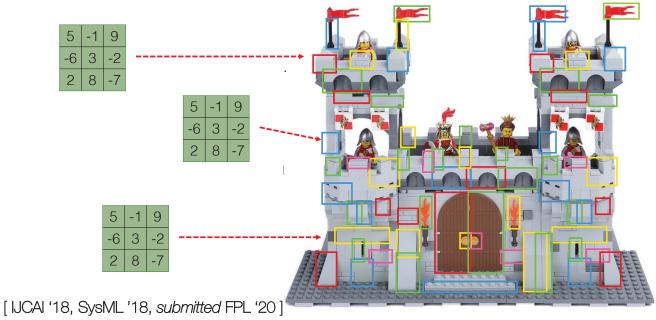


#1 Think (i.e., learn) Different

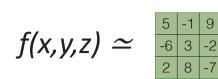




Inflating deep models from functions at inference: a new form of trading memory for compute

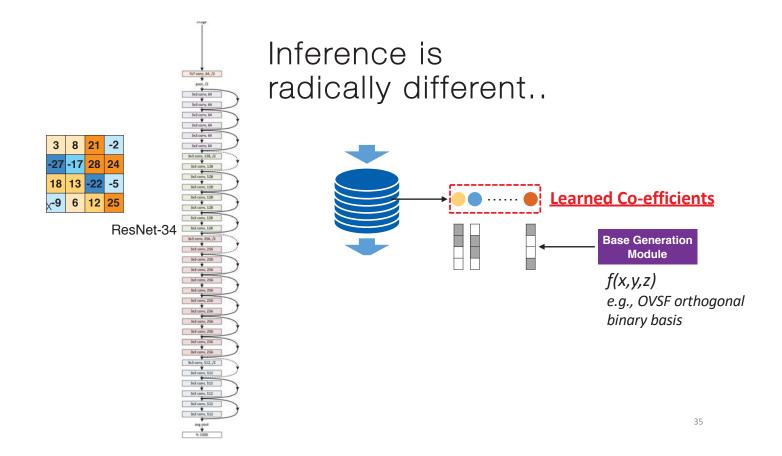


Inflating deep models from functions at inference: a new form of trading memory for compute





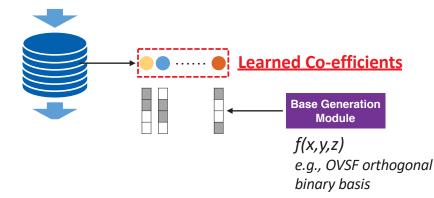
[IJCAI '18, SysML '18, submitted FPL '20]



Inference is radically different..

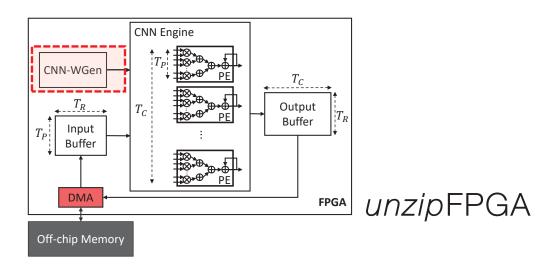
	Top1 Accuracy (%)
ResNet-18	91.15
ResNet-18 (OVSF)	91.02
ResNet-34	92.46
ResNet-34 (OVSF)	92.32
SqueezeNet	91.16
SqueezeNet (OVSF)	91.33

Dataset: CIFAR-10

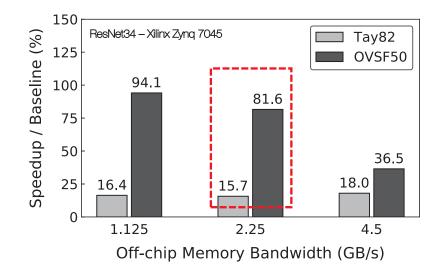


[IJCAI '18, SysML '18, submitted FPL '20]

Maximizing efficiency potential through hardware



Maximizing efficiency potential through hardware

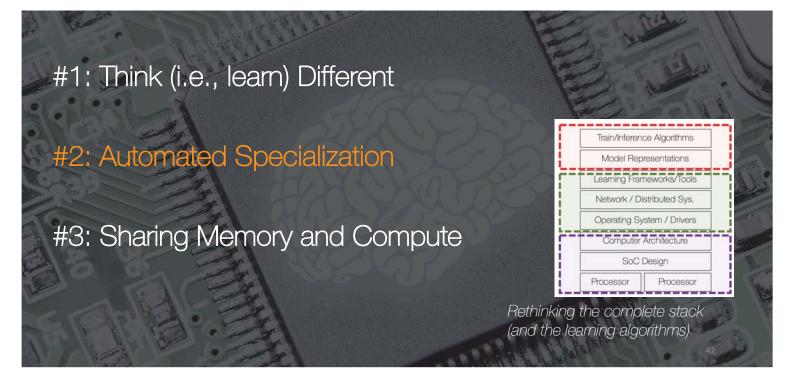


	Dataset: ImageNet		
Compression	Params	Accuracy	
(ResNet34)	(millions)	(%)	
None	21.8	73.3	
Tay82	17.4	72.7	
OVSF50	17.2	72.1	

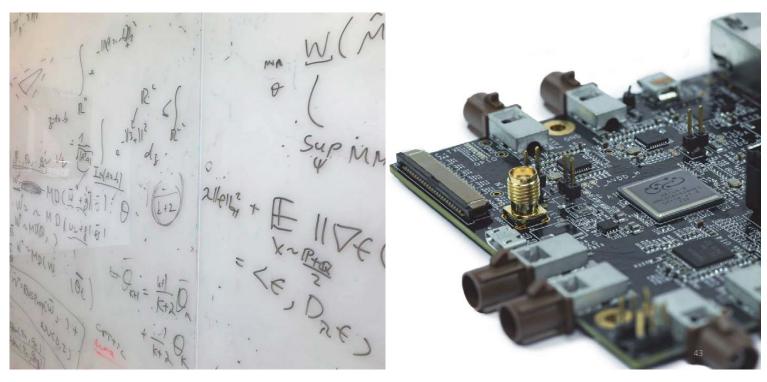
[IJCAI '18, SysML '18, submitted FPL '20]

40

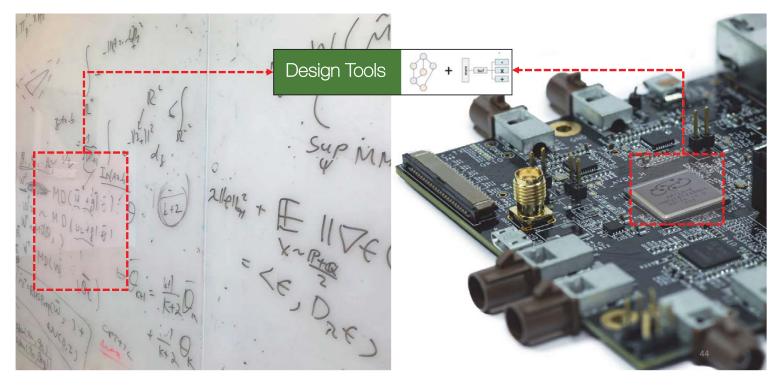
Fundamental Efficient ML Challenges



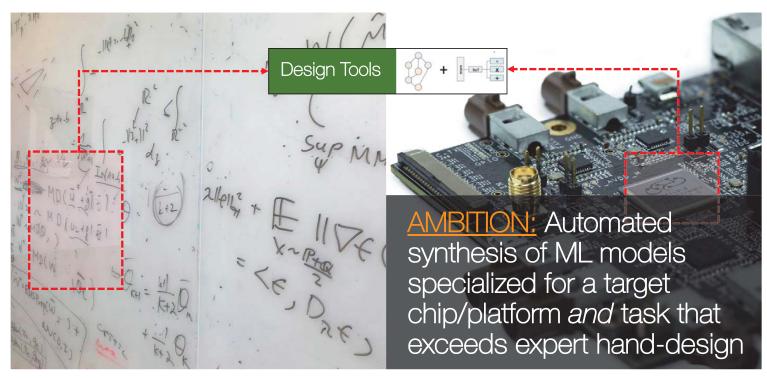
#2 Automated Specialization



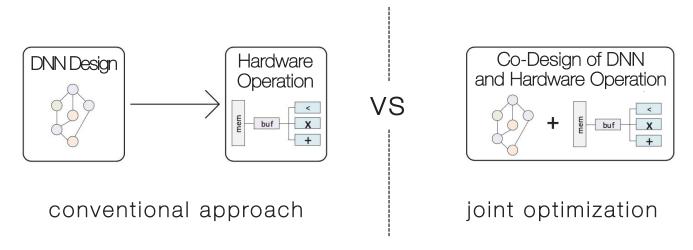
#2 Automated Specialization



#2 Automated Specialization



Joint Optimization of Hardware Operation and Neural Architecture



Joint Optimization of Hardware Operation and Neural Architecture

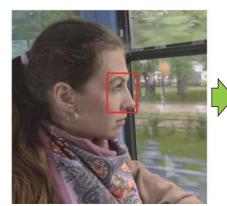
Dataset: Librispeech



[INTERSPEECH '19, INTERSPEECH '20, EECV '20, DAC '20, NeurIPS '20]

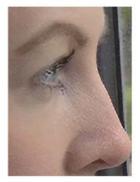
Joint Optimization of Hardware Operation and Neural Architecture

Dataset: Set5			
	LPIPS	Mult-Adds (G)	Params (K)
ESRGAN	0.074	1034.1	16,697
FEQE	0.091	5.6	96
AutoCAML	<u>0.076</u>	<u>3.6</u>	<u>61</u>



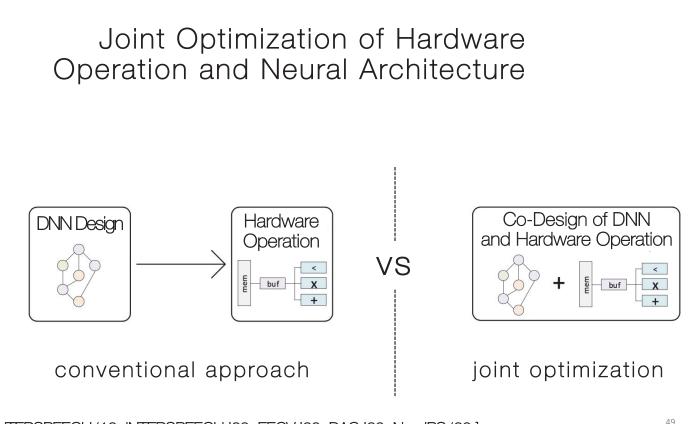


AutoCAML (61k)



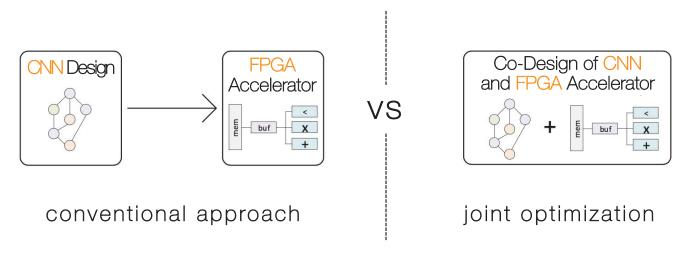
47

ESRGAN (16,697k) 48

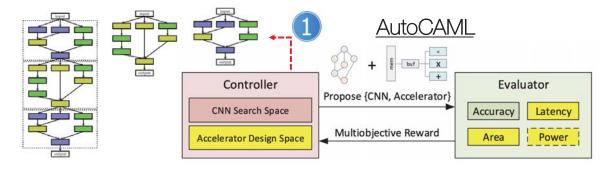


[INTERSPEECH '19, INTERSPEECH '20, EECV '20, DAC '20, NeurIPS '20]

Joint Optimization of Hardware Operation and Neural Architecture

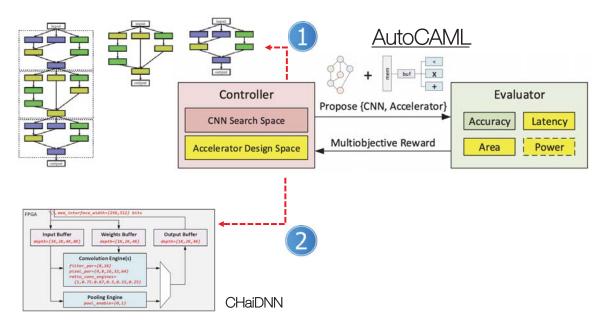


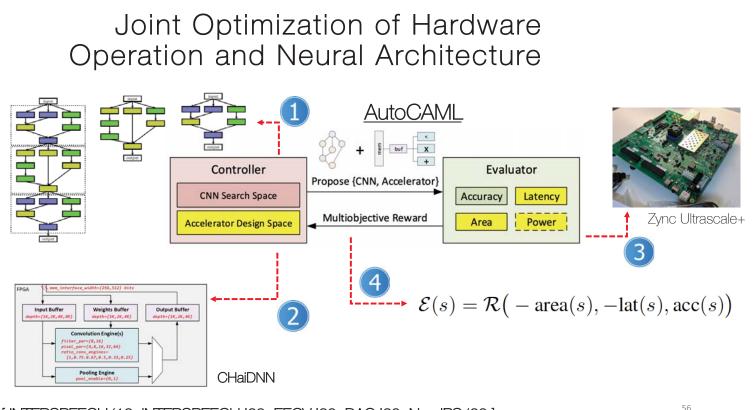
Joint Optimization of Hardware Operation and Neural Architecture



[INTERSPEECH '19, INTERSPEECH '20, EECV '20, DAC '20, NeurIPS '20]

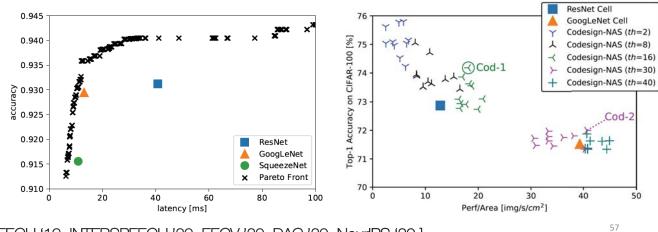
Joint Optimization of Hardware Operation and Neural Architecture



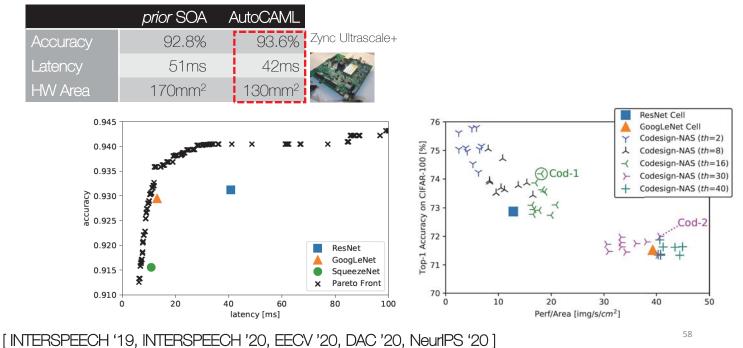


[INTERSPEECH '19, INTERSPEECH '20, EECV '20, DAC '20, NeurlPS '20]

Joint Optimization of Hardware Operation and Neural Architecture



Joint Optimization of Hardware Operation and Neural Architecture



Fundamental Efficient ML Challenges

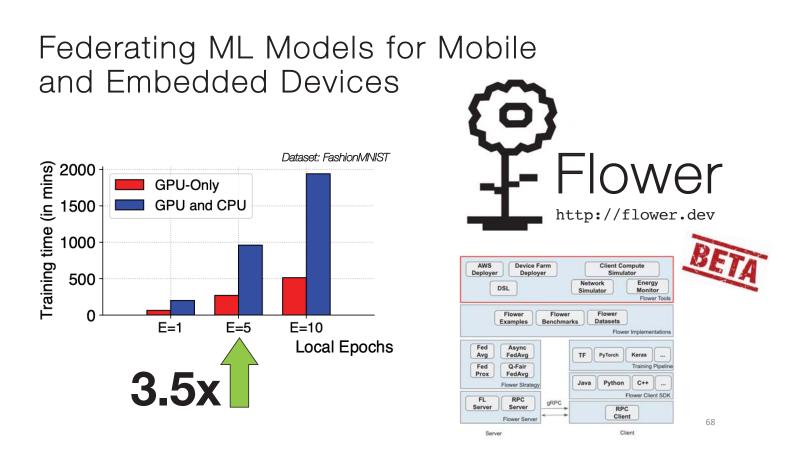
#1: Think (i.e., learn) Different
#2: Automated Specialization
#3: Sharing and Cooperation
Formulation Constraints
Formu

#3 Sharing and Cooperation



#3 Sharing and Cooperation

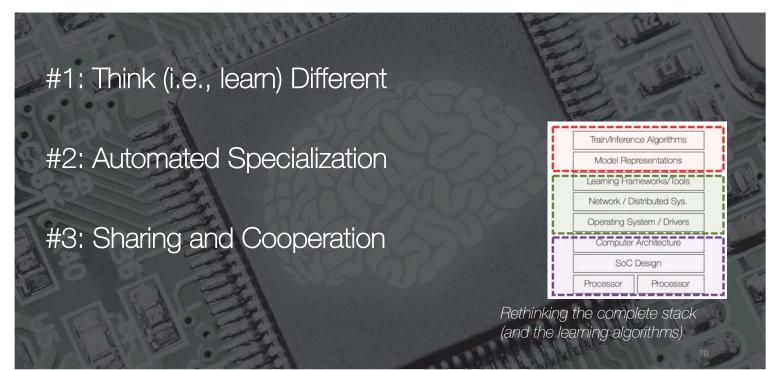




#3 Sharing and Cooperation



Fundamental Efficient ML Challenges



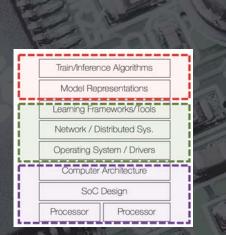
Fundamental Efficient ML Challenges

#1: Think (i.e., learn) Different

#2: Automated Specialization

#3: Sharing and Cooperation

#4: Next Steps in Hardware



Rethinking the complete stack (and the learning algorithms)



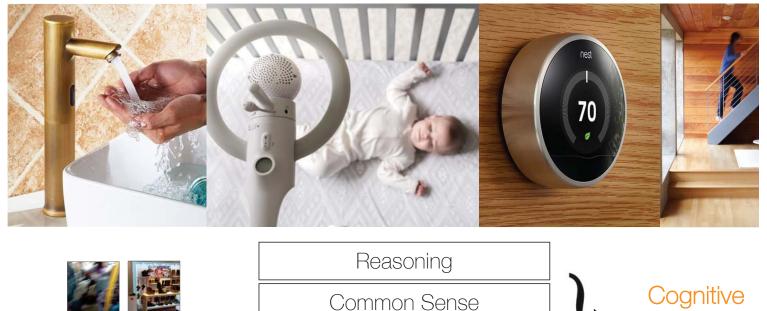
Predictions for the Efficient ML Revolution #1 On-Device ML goes far beyond just classification #2 SOTA Accuracy will come from efficient ML Models #3 Data Center is *replaced* by devices as the heart of ML



#1 Efficient ML Prediction On-Device ML goes far beyond classification

Discriminative Task

person sleeping }



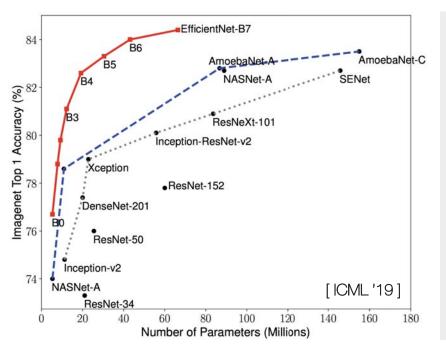
Understanding

Perception (Discriminative)



Stack

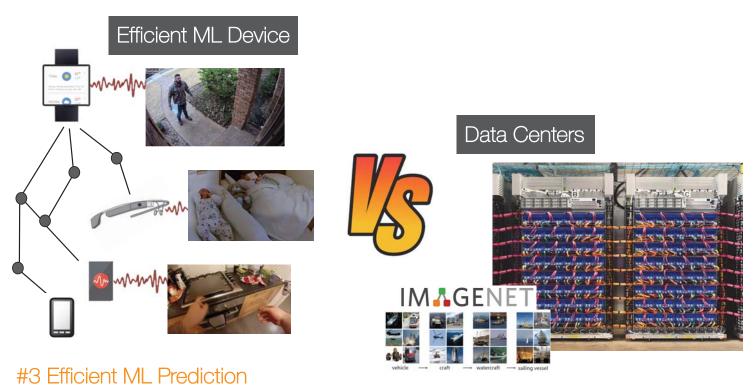
Embedded



Peace Dividends of Efficiency

- Faster exploration & experimentation
- "Intractable approaches" become possible...
- Consuming more data both labeled and unlabeled varieties
- Larger and larger architectures
- Able to heavily exploit **automated** methods like architecture search
- Brand new ML tasks

#2 Efficient ML Prediction SOTA Accuracy will come from Efficient ML



Devices replace Data Centers as core of ML

