Virtual:Tech Talk

Faster time-to-production for Computer Vision AI on Arm-powered Edge Devices

Deeplite Inc.

Davis Sawyer, CPO
Charles Marsh, CCO
Feb 8th, 2022





Welcome!

* Tweet us: <u>@ArmSoftwareDev</u> -> #AIVTT

Check out our Arm Software Developers YouTube <u>channel</u>
*

*Signup now for our next Al Virtual Tech*Talk: developer.arm.com/techtalks

Our upcoming Arm AI Tech Talks

Date	Title	Host
February 22 nd , 2022	* New Arm ML Quarterly Research Special * Federated Learning Based on Dynamic Regularization to Debias Model Updates	Arm ML Research



Presenters



Davis Sawyer

Davis is the Co-founder and CPO at Deeplite, where he leads product development and strategy. Recognized as a C2 Montreal Emerging Entrepreneur, he also serves as chairperson for tinyML Montreal and Toronto.



Charles Marsh

Charles is the CCO at Deeplite, leading sales and business development. He specializes in business growth for high-tech start-up companies, with two successful exits todate.



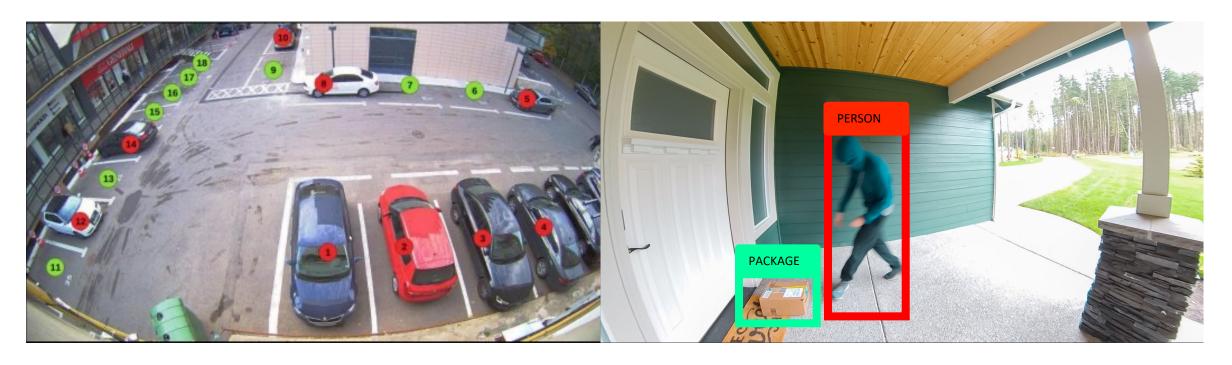
What we've learned as an industry:

Productizing computer vision AI at the edge is hard



Latest wave of accurate DNNs for Computer Vision

Making real world vision-powered products possible



Intelligent Video Analytics

Smart Doorbell Cameras



Latest wave of accurate DNNs for Computer Vision

How do we go from research

NVIDIA Tesla V100 GPU 250 Watt \$2250-8000 USD



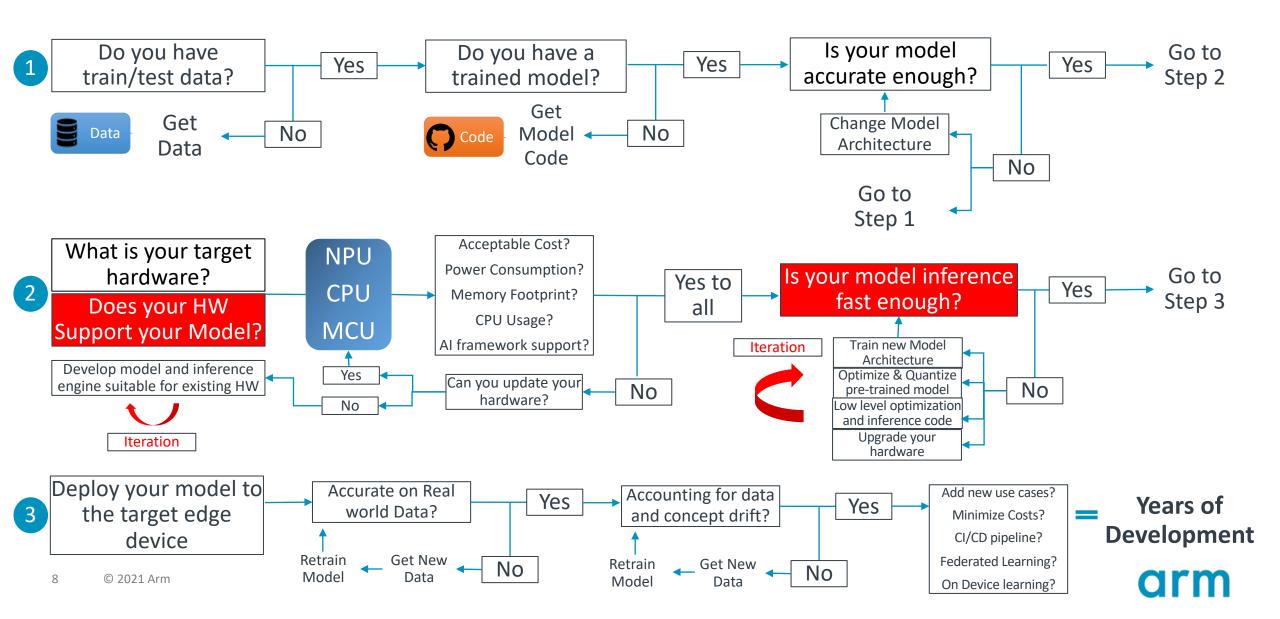
To the real-world:

Quadcore Arm Cortex-A53 CPU 200mW \$30-35 USD (RPi3B)

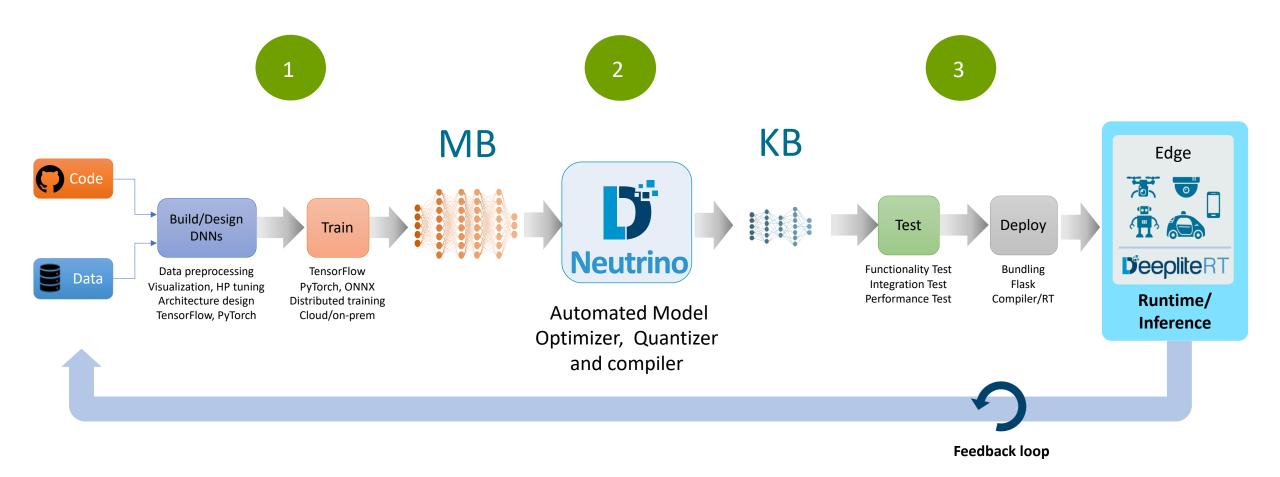




Simplified process for an edge computer vision AI use case



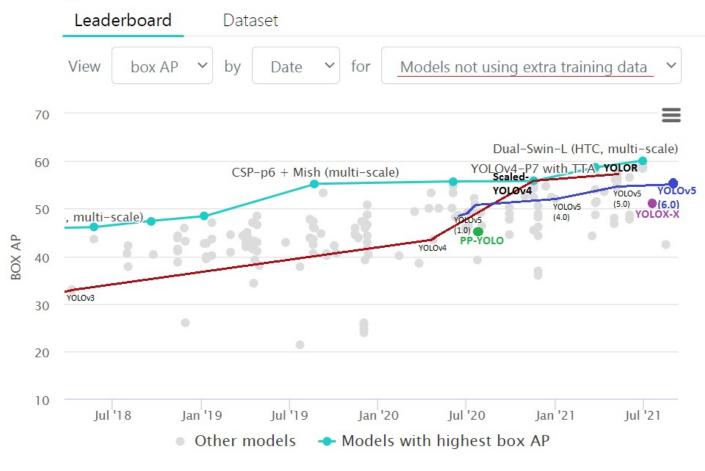
Really simplified process for an edge computer vision AI use case





Best Practices for Model Training / Development (Step 1)

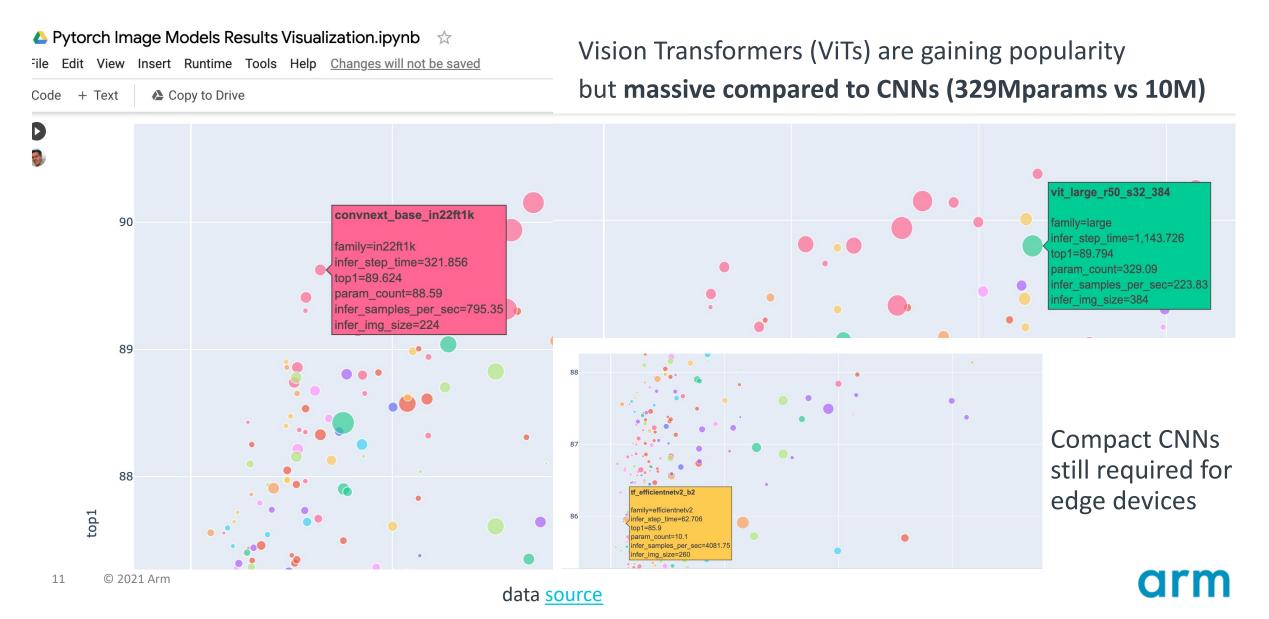
Object Detection on COCO test-dev



- Surge in last few years in accurate, off-the-shelf deep neural networks for vision related tasks
- Model Size/Inference Speed tradeoffs to be made when moving from training to inference hardware

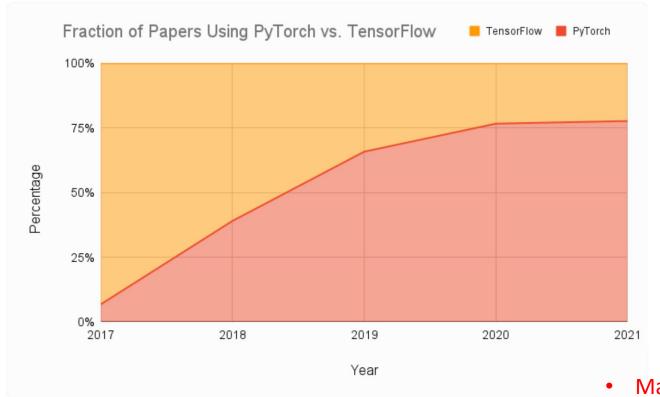


Best Practices for Model Training / Development (Step 1)



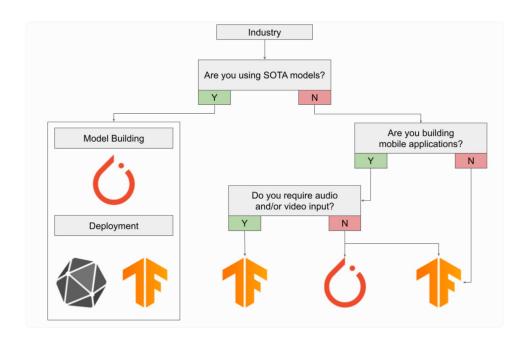
Best Practices for Model Training / Development (Step 1)

Al Frameworks



data source

What if I'm in Industry?



Many CPUs and MCUs don't support PyTorch/ONNX for inference

SOTA Vision models are increasingly available in PyTorch



Best Practices for Inference Optimization (Step 2)

CV Model Architecture Optimization



Built-in Capabilities

- Weight Pruning
- Pruning+XNN Pack
- NAS



3rd Party Tools

- Open-Source Distillation & Compression
- SW Accelerators

Application	Initial Architecture	Deeplite Neutrino Compression (FP32)			Accuracy Change (%)	Dataset
		Original Size	Optimized Size	Improvement		
	VGG19	80MB	2.16MB	x37	<1.00% (Top1)	CIFAR100
Image	Mobilenet- v1.0	12.8MB	530KB	x22	~1.50% (Top1)	Visual Wake Words
classification	ResNet18	45MB	1.8MB	x27	<1.00% (Top1)	Subset of ImageNet
	ResNet50	97MB	26MB	х3.72	~1.50% (Top1)	ImageNet
	ResNet50- SSD300	54MB	18MB	х3	~ 0.00 (mAP)	Subset of COCO2017
Object Detection	ResNet34- SSD300	32MB	11MB	x3.75	~0.00 (mAP)	VOC 2012
	Yolov3	235MB	28MB	х8	<0.02 (mAP)	
	Unet- ResNet18	83.34MB	34.7MB	x2.4	~ 0.00 (mloU)	
Segmentation	Deeplab- Mobilenetv2	29.1MB	13.3MB	x2.2	< 0.02 (mIoU)	
	U-Net	65.9MB	1.04MB	x63.5	< 0.01 (mloU)	Carvana



Best Practices for Inference Optimization (Step 2)

Low precision quantization for efficient inference



Post-Training Quantization (PTQ)
FP16, INT8

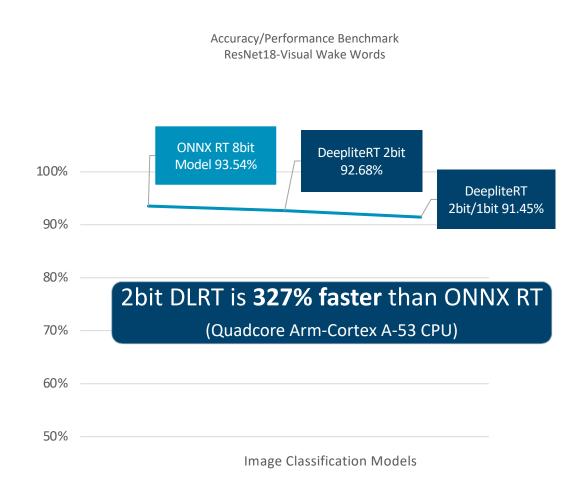
Quantization-Aware Training (QAT)
INT8



Post-Training Quantization (PTQ)
FP16, INT8

Quantization-Aware Training (QAT)
INT8, 2bit, 1bit

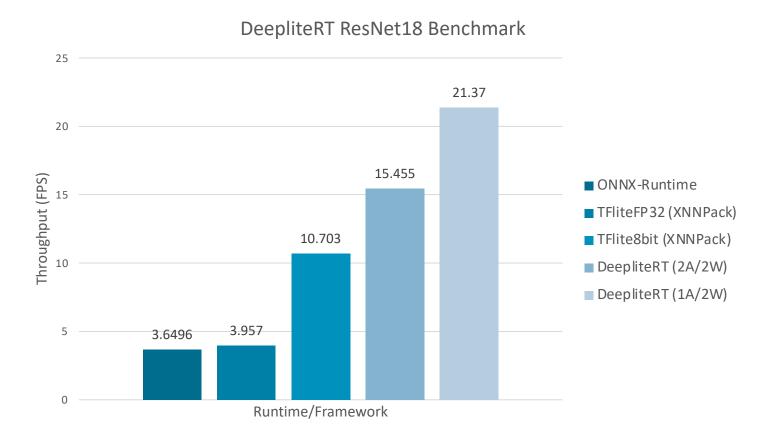
DeepliteRT





Best Practices for Streamlined Edge Deployment (Step 3)

Approaching GPU inference speed on Cortex-A CPU





Visual Wake Up + Detection on Arm Cortex-A CPU

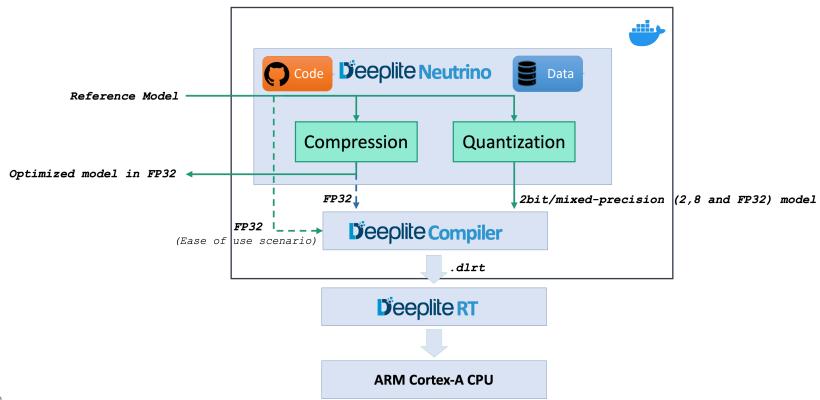
2bit YOLOv5 Running on Raspberry Pi4 (Arm Cortex-A72 @ 1.5GHz Quad-core CPU)





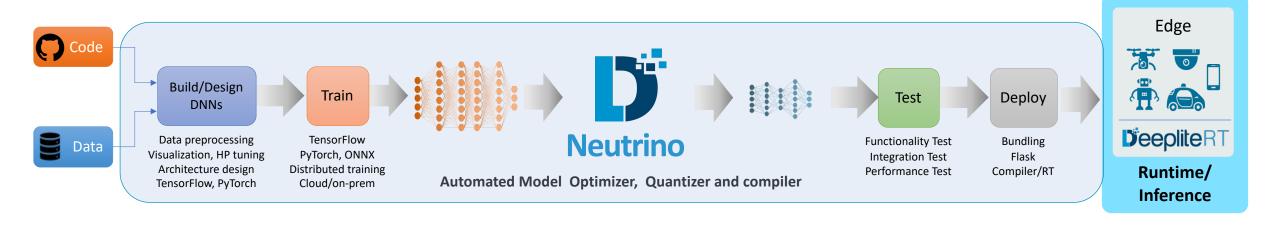
Best Practices for Streamlined Edge Deployment (Step 3)

- Save on memory, cost and power for inference
- Save on development effort with docker & pip install for deployment





Putting it all Together in the Edge AI Workflow





Person Detection on low-power Arm CPU

Problem

- Increasing cloud costs to support application
- 1000s of cameras in field (Cortex-A53)
- Train a multi class object detection model
- Small, fast & accurate enough to run on device
- Product differentiation



Solution & Result

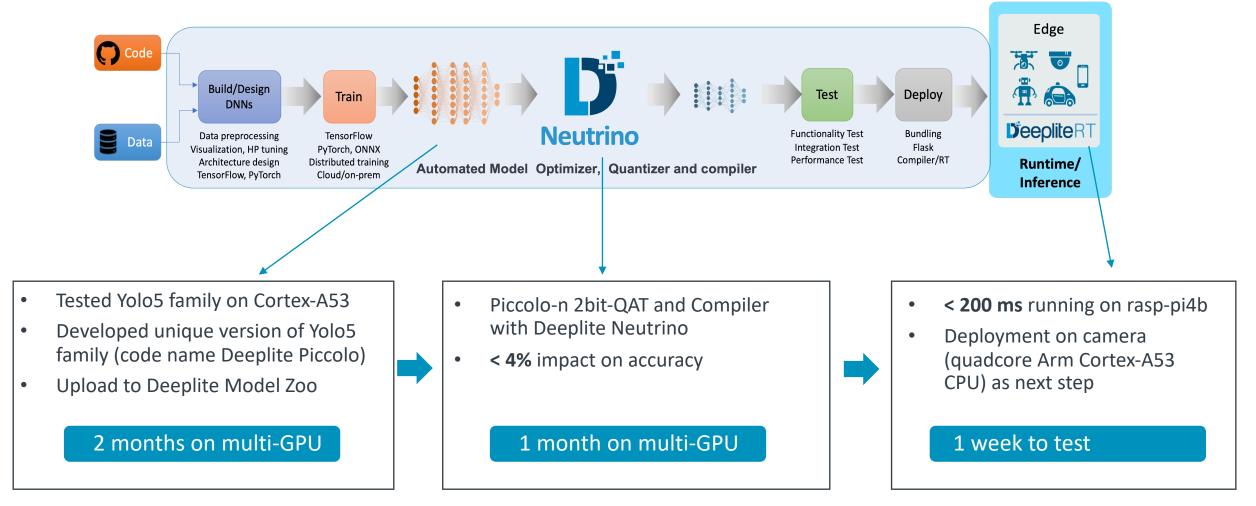
- Developed unique version of Yolov5 family (code name **Deeplite Piccolo**)
- 2-bit quantization via Neutrino engine
- Running via DeepliteRT on camera
- Ground-breaking performance
 - < 200 ms
 - < 4% impact on accuracy





Person Detection on low-power Arm CPU

Approach





Arm CPU-Powered Person Detection in Action!







Intelligent Speed Adaptation on low-power Arm CPU

Problem

- GPS speed control not efficient & accurate
- Speeding in construction zones = huge expense
- 1000s of units in field (must use existing CPU)
- Train an efficient object detection model
- Small and fast enough to run on Cortex-A53



Solution & Result

- Model training pipeline & testing for Yolo and Resnet-SSD
- Neutrino optimized model (17x size reduction)
- Deeplite platform and run-time produced Yolo5s (2bit quantized model) on rasp-pi4b
 - 256.78 ms
 - < 4% impact on accuracy



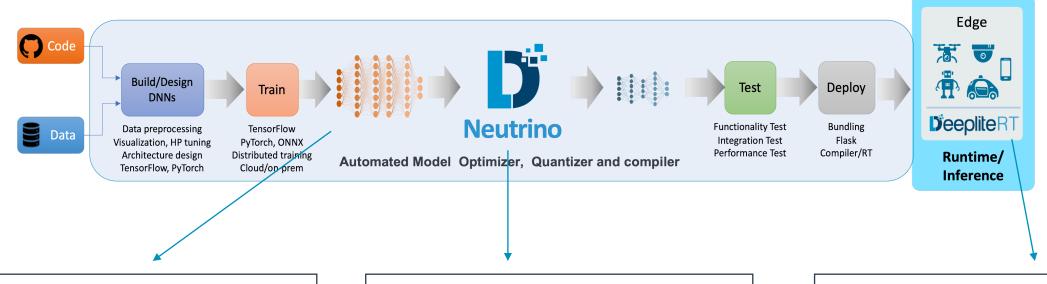






Intelligent Speed Adaptation on low-power Arm CPU

Approach



- Train object detection models for construction zones & speed limits >0.80mAP
- Based on Yolo and Resnet-SSD series of models (PyTorch)

2 months on multi-GPU

- 17x compression over 32bit models with <0.02 mAP drop
- Yolo5s 2bit-QAT and Compiler with Deeplite Neutrino

1 month on multi-GPU



 Successfully deployed & running on camera (quadcore Arm Cortex-A53 CPU)

1 month to test/deploy



Always-on, Optimized Computer Vision with Arm NPU

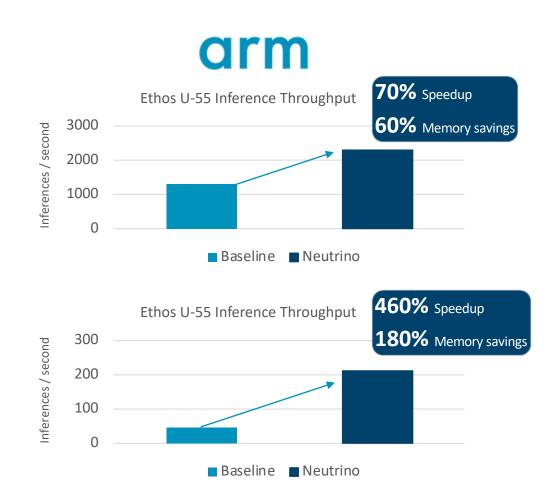
D'eeplite Neutrino

Person Detection, High Efficiency Use Case:

- Baseline MobileNetv1-0.25x model
- Neutrino 2x compression with > 85% final accuracy

Object Classification, High Performance Use Case

- Baseline ResNet-18 model
- Neutrino 22x compression with < 1% drop in accuracy





Go from AI research to the real-world with DeepliteRT

Try it Today: dlrt.deeplite.ai

Username: armdlrt

Password: armdlrt







Thank You Danke Merci 射射 ありがとう Gracias

* Kiitōs 감사합니다 धन्यवाद

شکِرًا

תודה



Thank you!

*Tweet us: <u>@ArmSoftwareDev</u> -> #AIVTT

Check out our Arm Software Developers YouTube <u>channel</u>
*

*Signup now for our next Al Virtual Tech*Talk: developer.arm.com/techtalks