# tinyML® Summit

*Miniature dreams can come true...*

## March 28-30, 2022 | San Francisco Bay Area

### TINY
### ML

## www.tinyML.org

# TinyML for All: Full-stack Optimization for Diverse Edge AI Platforms

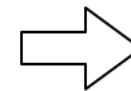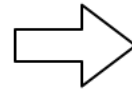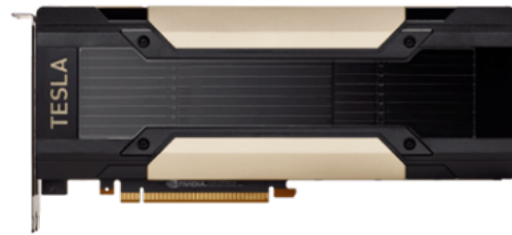**Di Wu**

*Co-founder and CEO, OmniML*

**Song Han**

*Assistant Professor, MIT EECS*

*Co-founder and Chief Scientist, OmniML*

# TinyML is about Constraints

Mismatch: AI has been evolving unconstrained for many years



| | Cloud AI | Mobile AI | Tiny AI |
|---|---|---|---|
| Computation | 10 TFLOPS | GFLOPS | MFLOPS |
| Memory | 32GB | 4GB | 256KB |

**100,000x smaller**

# Everything Together: Real-world AI on Tiny MCUs

Two Generations of Innovations: MCUNet-v1 (2020), MCUNet-v2 (2021)

Facemask Detection

Person Detection

Works on Cortex M7 MCU

# Brief History of MCUNets

Reducing the model sizes with increasing accuracy

# Opportunity in Fundamental ML Algorithms

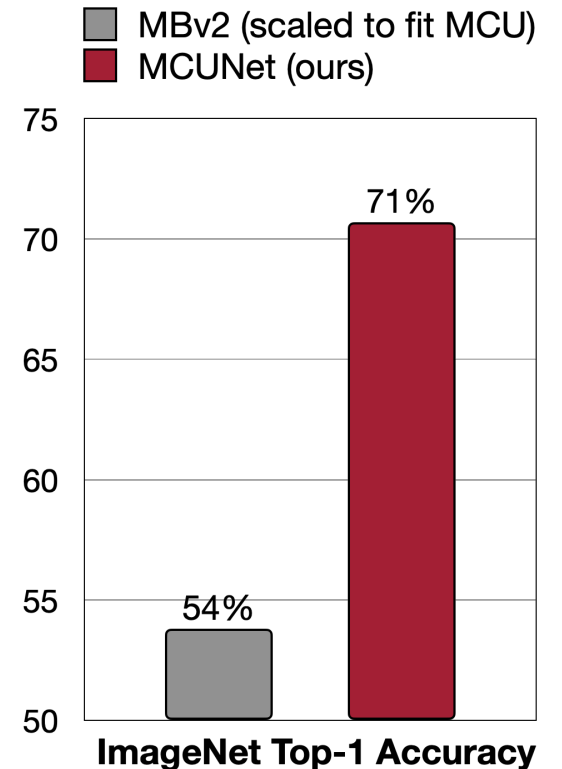Making algorithm more efficient under existing constraints

Faster than Moore's Law:
3.5x model size reduction every 12 months
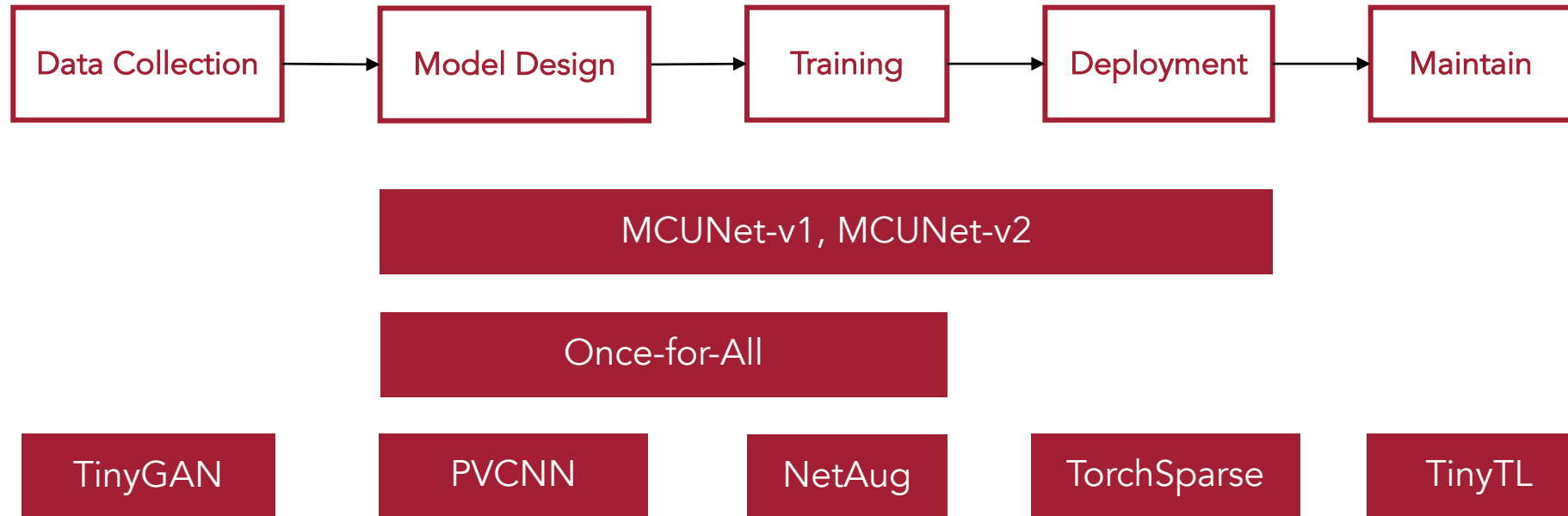
Improving efficiency means _more accurate_ models, too

TinyML is about improving the entire stack: from design to deployment, from computation to data
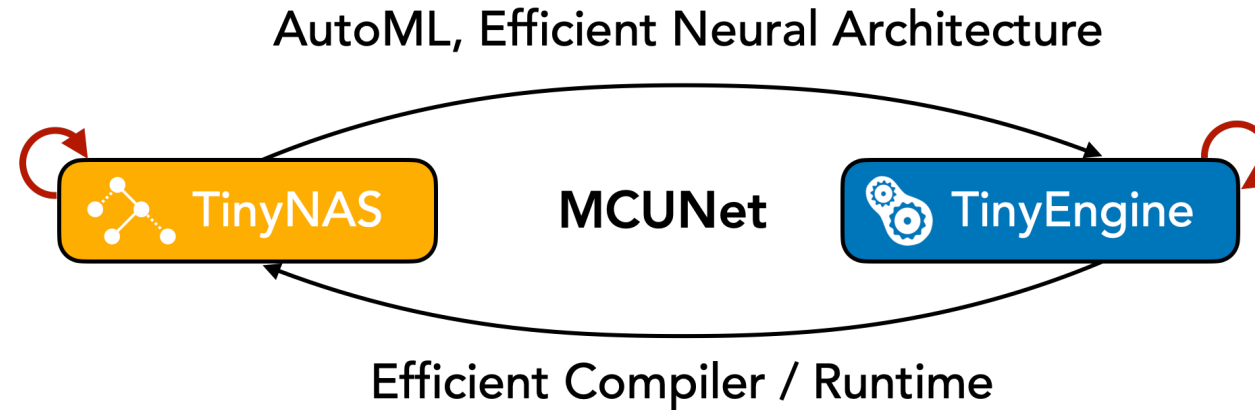


MBv2 (scaled to fit MCU)
MCUNet (ours)

71%

54%

**ImageNet Top-1 Accuracy**

# Agenda

Focus on Constraints on the Entire Stack

| Data Collection | → | Model Design | → | Training | → | Deployment | → | Maintain |
|---|---|---|---|---|---|---|---|---|

MCUNet-v1, MCUNet-v2

Once-for-All

| TinyGAN | PVCNN | NetAug | TorchSparse | TinyTL |
|---|---|---|---|---|

# MCUNet-v1: TinyNAS+TinyEngine Co-design

AutoML, Efficient Neural Architecture

TinyNAS — MCUNet — TinyEngine

Efficient Compiler / Runtime

TinyNAS:
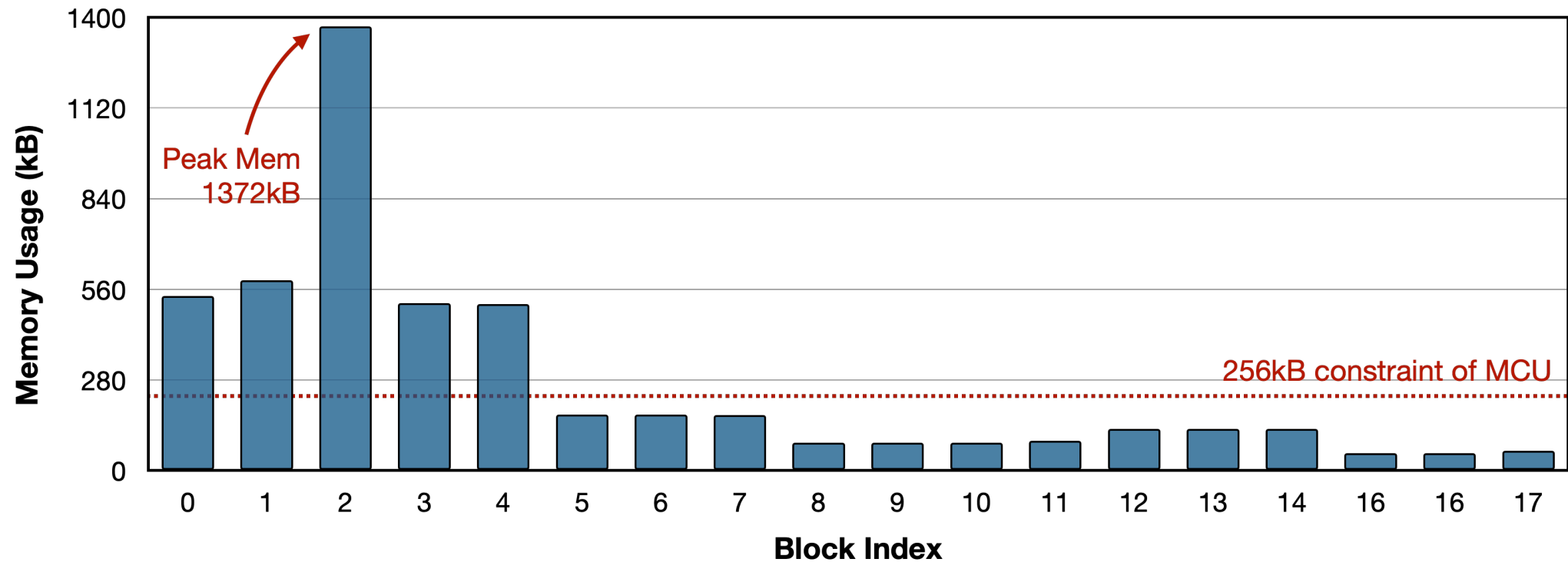- Re-design the design space
- Latency-aware
- Energy-aware
- Once-for-all Network

TinyEngine:
- Co-design, specialization
- Offload run-time to compile-time
- Graph optimizations
- Memory-aware scheduling
- Low-precision
- Assembly-level optimizations

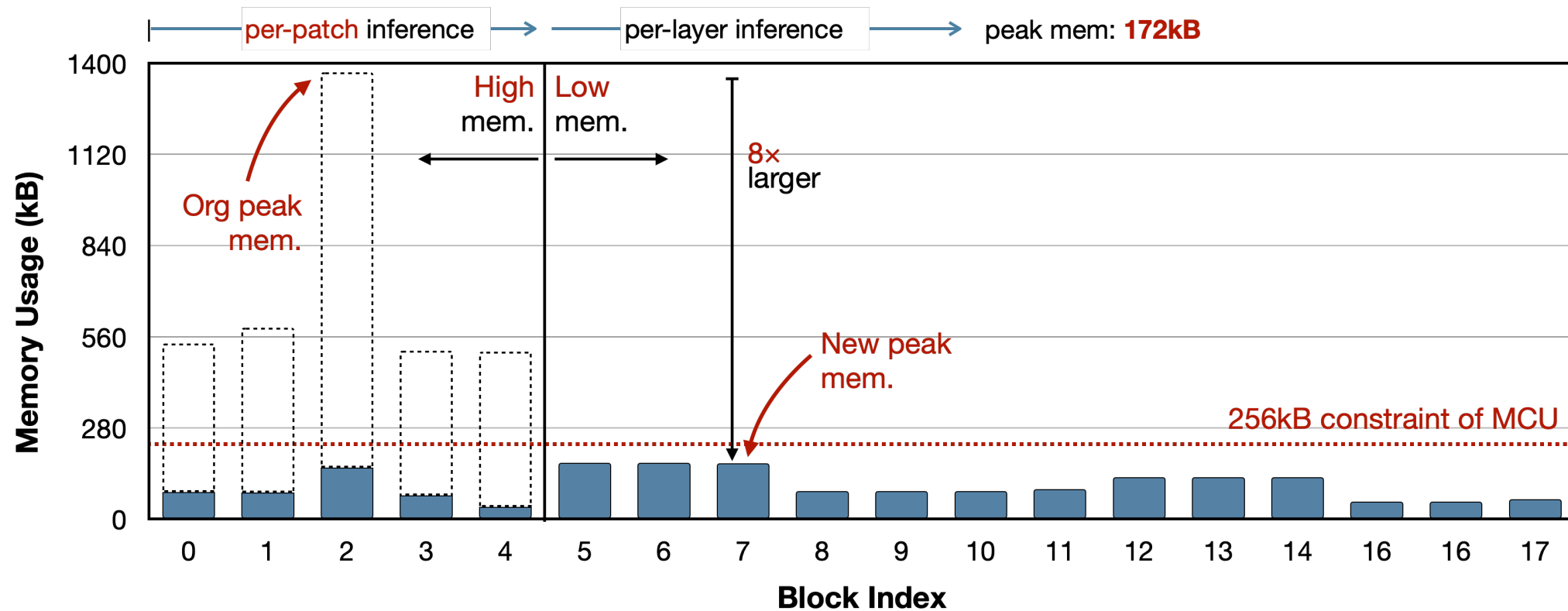# New Problem: Imbalanced Memory Distribution of CNNs [MCUNet-v2, NeurIPS'21]

Per-block memory usage of MobileNetV2

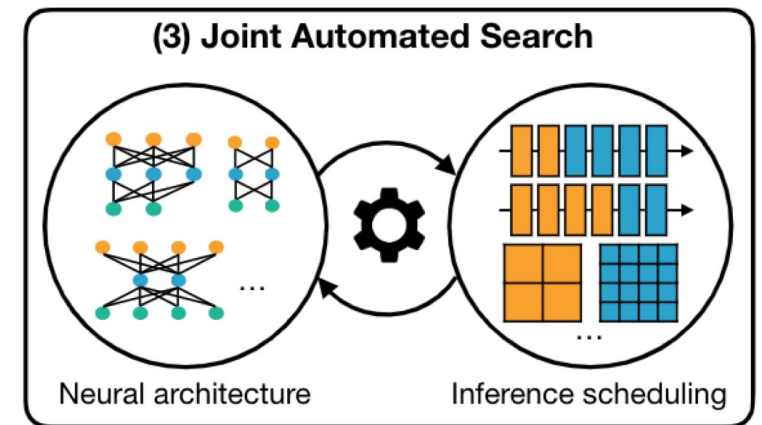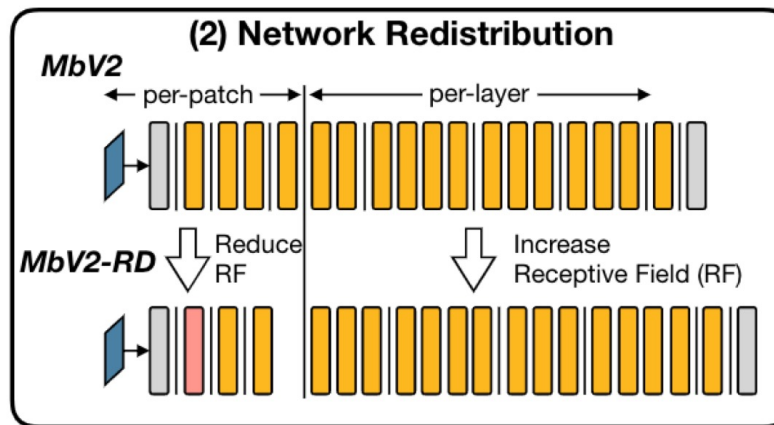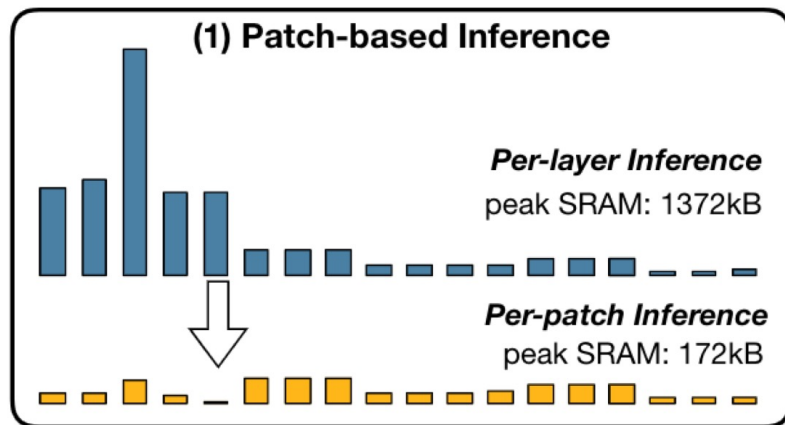# Solving the Imbalance with Patch-based Inference

After applying Patch-based Inference
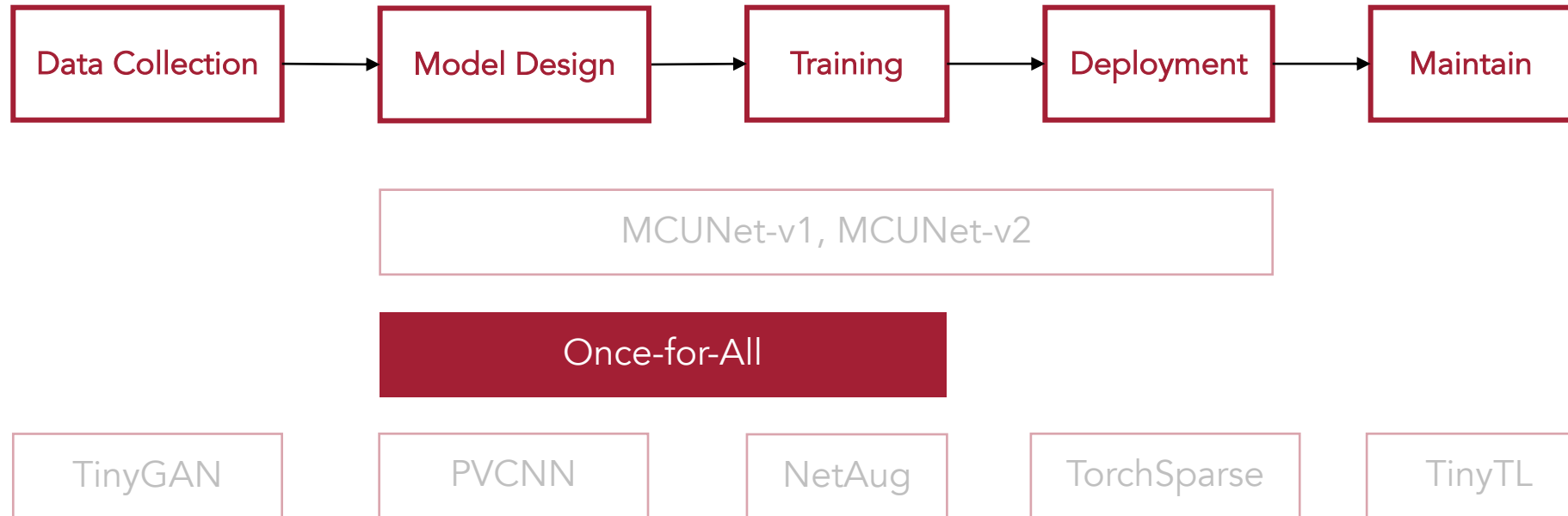
# MCUNet-v2 Takeaways

Solving inference bottleneck (peak memory) results in smaller and better models

# Agenda

Focus on Constraints on the Entire Stack

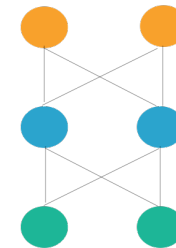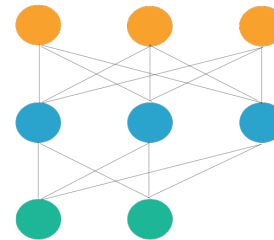# Once-for-All Network
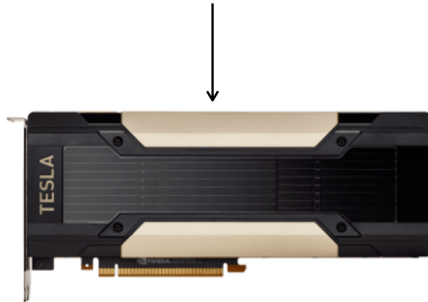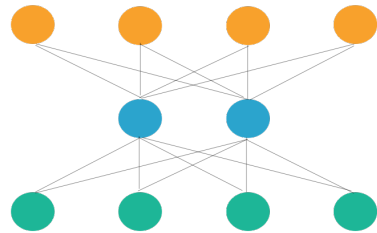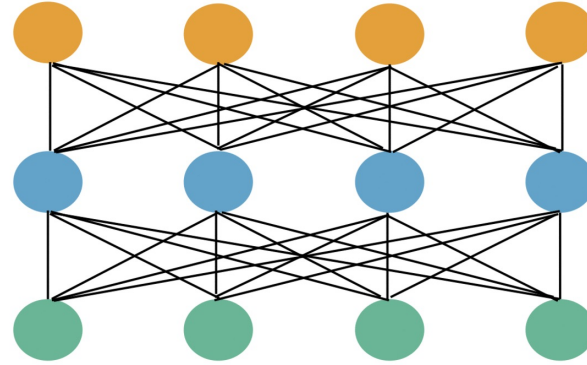
Train once, get many; Fit diverse hardware constraints

# Better Results with Much Smaller Training Cost

Reduce the search cost from 42,000 GPU hours (Google) to 200 GPU hours



Existing:        Lots of hand tuning for different devices and latency.

OFA:             Auto design the NN architecture at low cost

# Agenda

Focus on Constraints on the Entire Stack

# Problem in Training for Tiny Models

## Existing Training Techniques don't Apply to TinyML



Mixup



DropBlock



AutoAugment



ResNet50 (4.1G MACs)



MobileNetV2-Tiny (23.5M MACs)

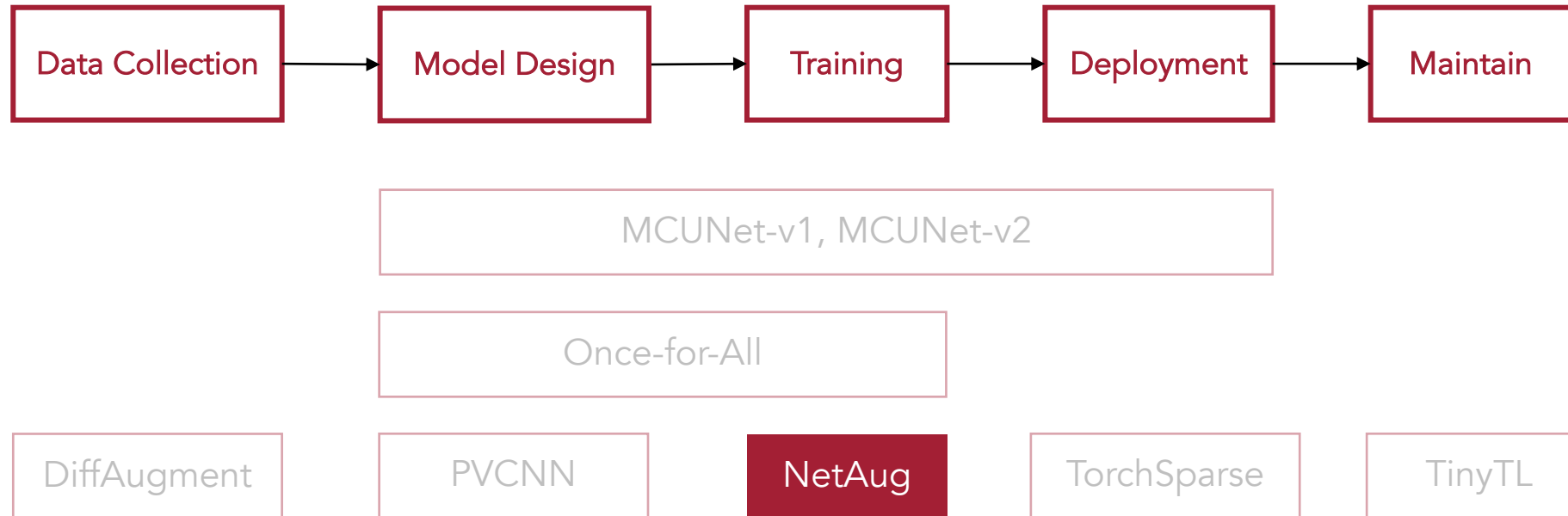# NetAug for TinyML

Augment Model Rather than Data

# Agenda
## Focus on Constraints on the Entire Stack

Data Collection → Model Design → Training → Deployment → Maintain

MCUNet-v1, MCUNet-v2

Once-for-All

TinyGan   PVCNN   NetAug   TorchSparse   TinyTL

# Problem: Training Memory is much Larger

Bottleneck is Activation rather than Parameters

# TinyTL: Up to 6.5x Memory Saving without Accuracy Loss

Use Fine-Tune Bias Only and Lite Residual Learning                    [TinyTL, NeurIPS'20]

# Agenda

Focus on Constraints on the Entire Stack



Data Collection → Model Design → Training → Deployment → Maintain

MCUNet-v1, MCUNet-v2

Once-for-All

TinyGAN · PVCNN · NetAug · TorchSparse · TinyTL

# Data is Also Constrained

Many TinyML Applications Have Limited Access to Data



Rare Defects



Specific Tasks



Privacy Concerns

# Differentiable Augmentation

Photo-realistic and Smooth Generation with 100 Training Images

# Agenda

Focus on Constraints on the Entire Stack

# TinyML for LIDAR & Point Cloud

Challenge: High Algorithm Complexity vs. Limited Computational Resource

[PCVNN, NeurIPS'19]
[SPVNAS, ECCV'20]
[PointAcc, Micro'21]
[TorchSparse, MLSys'22]

# Point Cloud Processing
## rence Library

**ball / kNN query**

$f_4 = \max(f_2w, f_3w, f_4w)$

$p_2$ $p_3$ $p_4$ $p_4$

**Point-Based Feature Transformation** (*Fine-Grained*)

**Multi-Layer Perceptron**

**Normalize**

**Add**

**Voxelize** → **Voxel Conv** → **Devoxelize**

**Voxel-Based Feature Aggregation** (*Coarse-Grained*)

[Point-Voxel CNN, NeurIPS'19]
New design space for Point Cloud

Fine-Grained Channel + Elastic Depth | Weight Sharing

Stage I (Depth: 3)
Stage II (Depth: 2,3)
Stage III (Depth: 1,2,3)

$\#C_{out}$
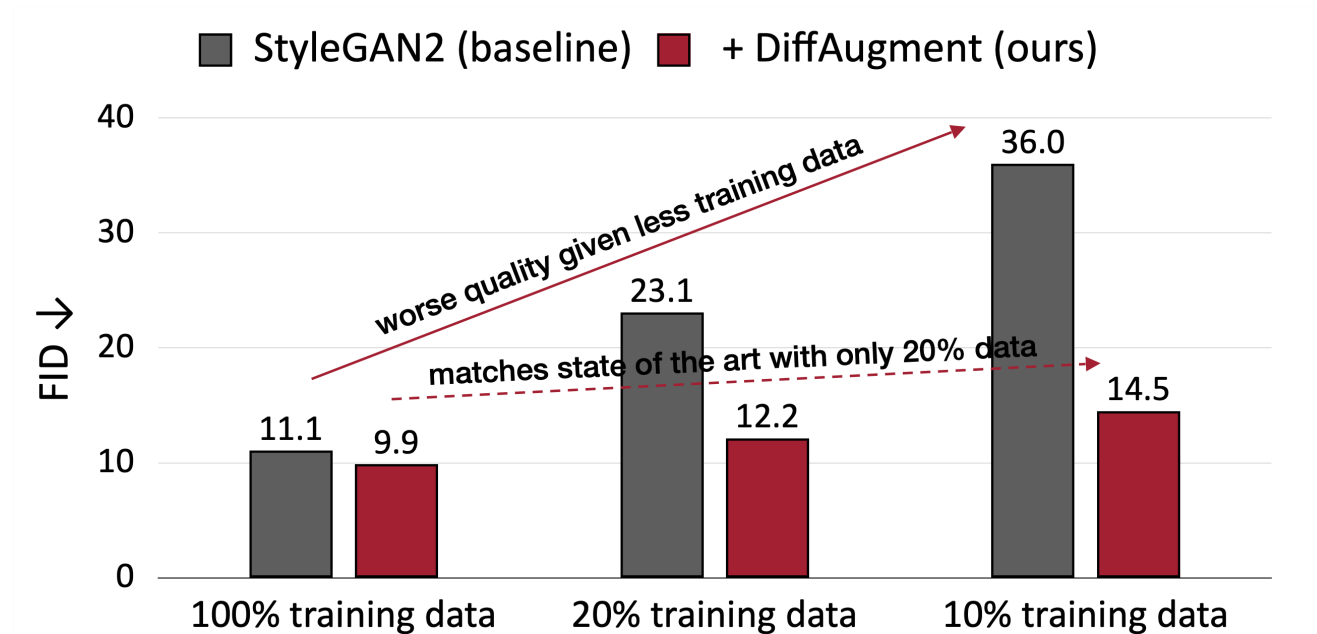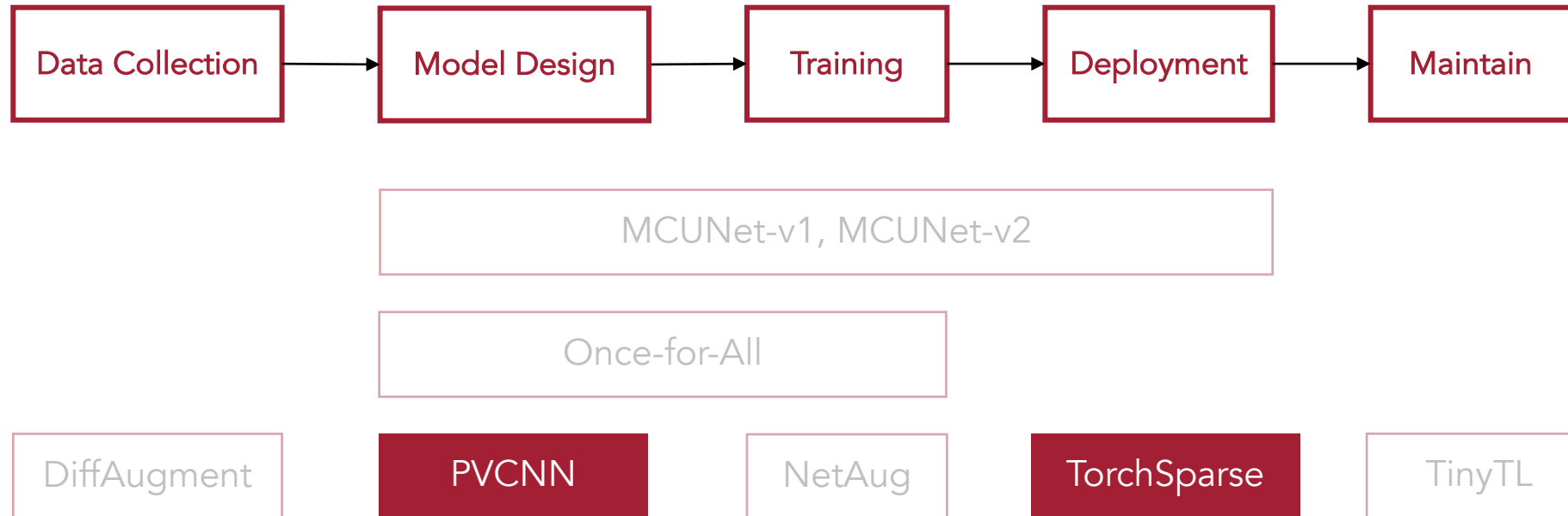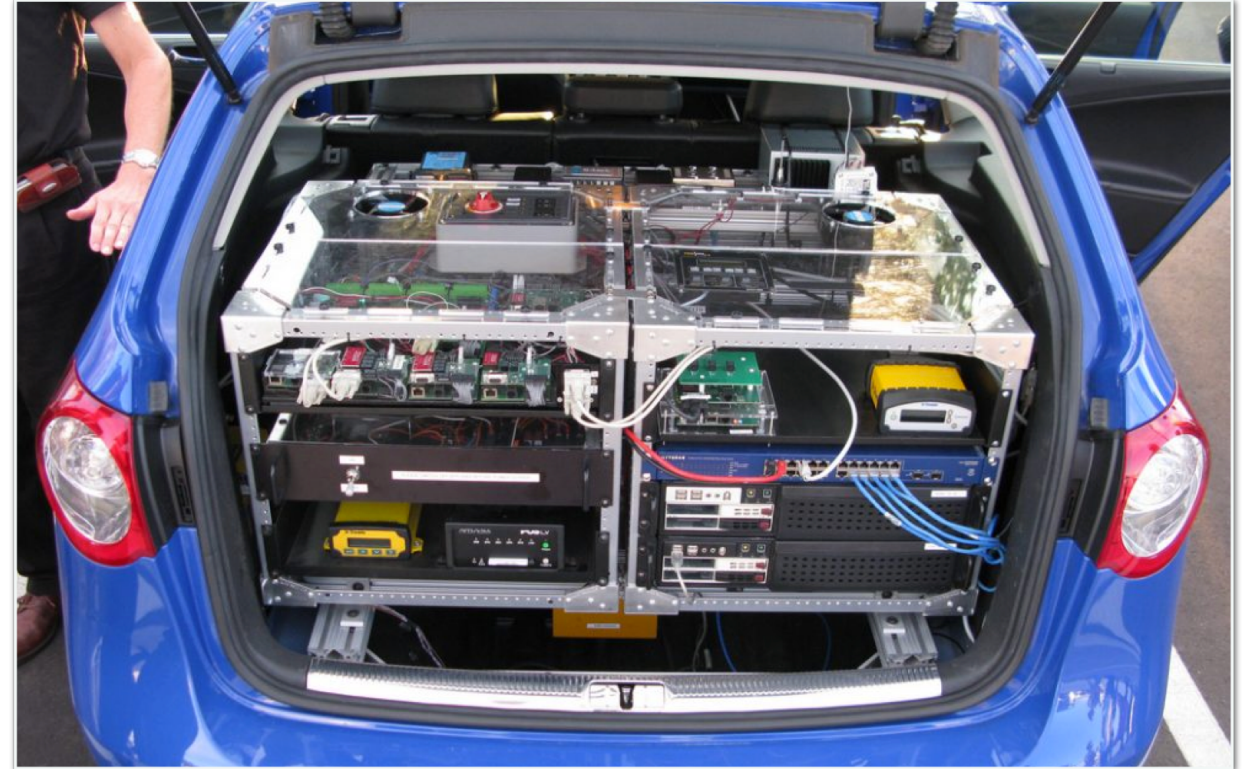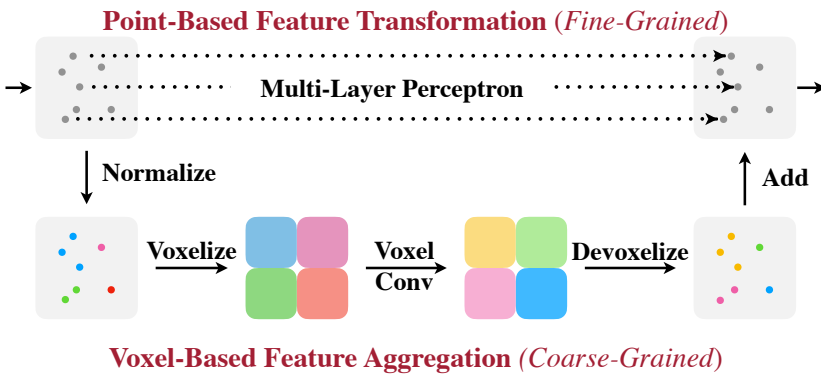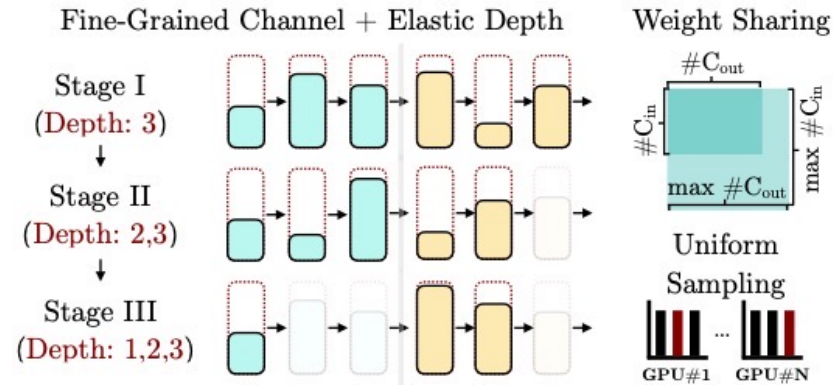$\#C_{in}$
max $\#C_{out}$
max $\#C_{in}$

Uniform Sampling

GPU#1 ... GPU#N

[SPVNAS, ECCV'20]
3D neural architecture search

**Algorithm**

**Hardware** ← → **System**

Global Buffer

Sorter Buffer | Input Feature Buffers | Merger Buffer | Output Feature Buffers | Weight Buffers

**Coordinates**

**Features & Weights**

Fetch Coords | Fetch Coords
Calculate Distance | MergeSort
Sort | Detect Intersection

**Mapping Unit (MPU)**

Memory Meta Container
Tile 2 | Tile 1 | Tile 0

Map FIFO
p index
q index
w index

Addr. Gen

Tag Array

**Address**

**Maps**

**Memory Management Unit (MMU)**

Systolic Array
⊗ ⊗ ⊗
⊗ ⊗ ⊗
⊗ ⊗ ⊗

**Matrix Computing Unit (MCU)**

**Maps (In, Out, Wgt)**

$(P_0, Q_1, W_{-1,-1})$
$(P_3, Q_4, W_{-1,-1})$
$(P_1, Q_3, W_{-1,0})$
$(P_0, Q_0, W_{0,0})$
$(P_1, Q_1, W_{0,0})$
$(P_2, Q_2, W_{0,0})$
$(P_3, Q_3, W_{0,0})$
$(P_4, Q_4, W_{0,0})$
$(P_3, Q_1, W_{1,0})$
$(P_1, Q_0, W_{1,1})$
$(P_4, Q_3, W_{1,1})$

**Gather By Weight** ✕

F0 F3 ✕ W_{-1,-1} = PSUM 1 / PSUM 4
F1 ✕ W_{-1,0} = PSUM 3
F0 F1 F2 F3 F4 ✕ W_{0,0} = PSUM 0 / PSUM 1 / PSUM 2 / PSUM 3 / PSUM 4
F3 ✕ W_{1,0} = PSUM 1
F1 F4 ✕ W_{1,1} = PSUM 0 / PSUM 3
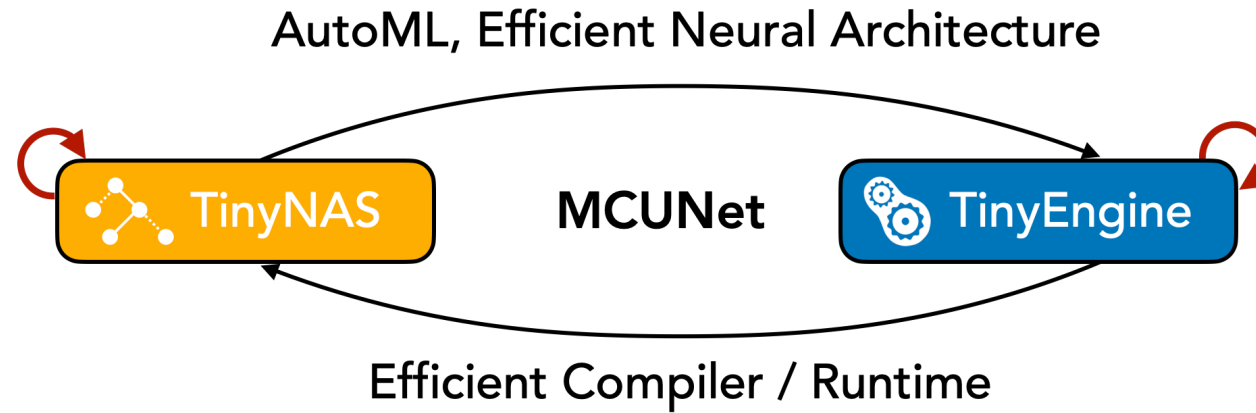
✕ n

**Scatter By Output** ✕

**Reduction**

[PointAcc, MICRO'21]
Hardware accelerator for point cloud

[TorchSparse, MLSys'22]
GPU library for 3D sparse convolution

MIT | HAN LAB

# Takeaways: Coming Back to MCUNets

AutoML, Efficient Neural Architecture



**MCUNet**

**TinyNAS** ⟷ **TinyEngine**

Efficient Compiler / Runtime

Co-optimization on the entire stack is the key to unlock the most potential for TinyML
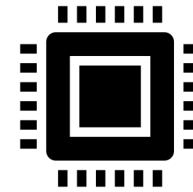
# Fundamental Problems in TinyML

ML under new HW constraints is very hard

Mismatch

X

**Slow Adoption** ➡ **Less Revenue/Volume**

Typical AI/ML Development

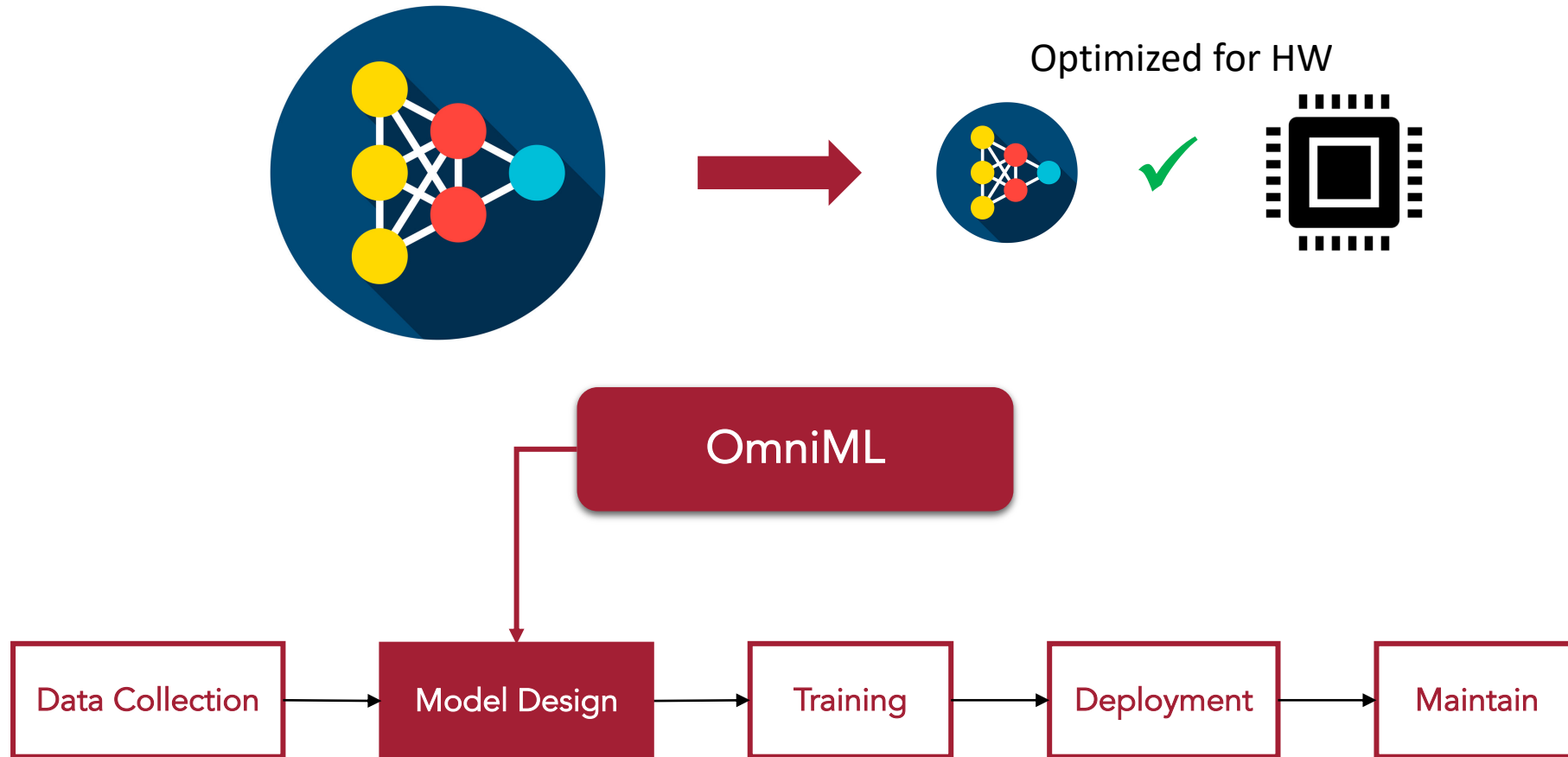| Data Collection | Model Design | Training | Deployment | Maintain |

Designing new models that works on different HW is still a manual and iterative approach

# OmniML "Compress" the Model Before Training

Bring HW deployment constraints into model design and training

Optimized for HW

OmniML

| Data Collection | → | Model Design | → | Training | → | Deployment | → | Maintain |

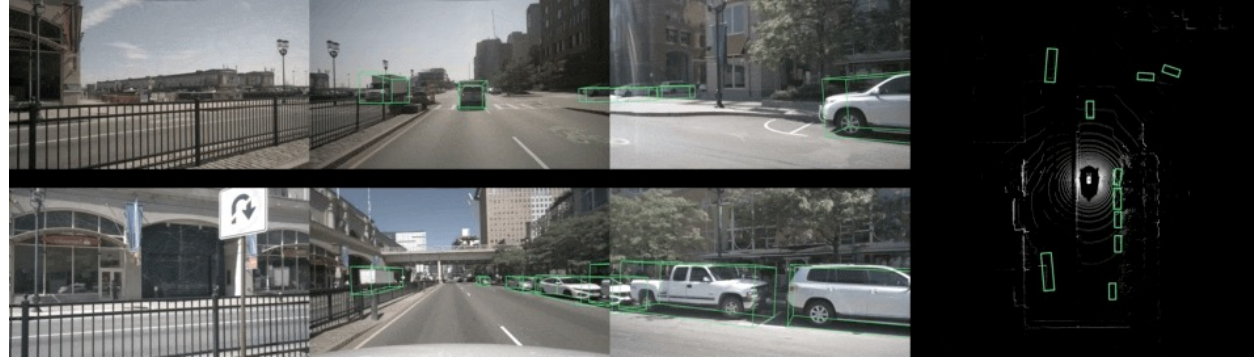# OmniML: Enable TinyML for All Vision Tasks

Create the Best Models on Different Platforms Effortlessly



**CV on Mobile Devices**
- Pose estimation
- Scene Segmentation
- Image denoise, super resolution
- AR/VR

**Sensor Fusion 3D Detection**
Multi-sensor 3D object detection for automotive applications.

**Smarter Cameras**
Turn "dumb" cameras into AI-powered cameras with advanced CV features on low-power, low-cost chip.

**Computer Vision on MCUs**
Not only classification but also object detection on microcontrollers with only 256~512KB of memory.

**40+** Customers Conversations     **10+** POCs     **100K** Installed devices

# Founding Team

## Leading Experts in Efficient Deep Learning

OMNI ᴹᴸ

### Song Han

- Assistant professor at MIT, PhD from Stanford
- Co-founder of DeePhi Tech (acquired by Xilinx)
- "35 Innovators Under 35" by MIT Technology Review
- NSF CAREER Award, IEEE "AIs 10 to Watch"
- Inventor of "Deep Compression"
- 29K Google Scholar citations

### Di Wu

- Previous tech lead at Facebook AI, PyTorch accelerator enablement
- Product and engineering leader at Falcon Computing Solutions (acquired by Xilinx)
- PhD from UCLA, years of experience in customized hardware systems at Intel Lab, MSRA.

### Huizi Mao

- PhD from Stanford. Co-Inventor of "Deep Compression"
- Early member of DeePhi and Megvii.
- Worked at Google Research, Facebook AML and NVIDIA.
- NVIDIA Fellowship Recipient.

![OMNI ML logo]

Thank you

Come talk to us to learn more

Follow us:
https://www.linkedin.com/company/omniml
https://twitter.com/OmniML_AI

We are hiring:
https://omniml.ai/career/
contact@omniml.ai

# tinyML Summit 2022 Sponsors

# Copyright Notice

# www.tinyml.org