



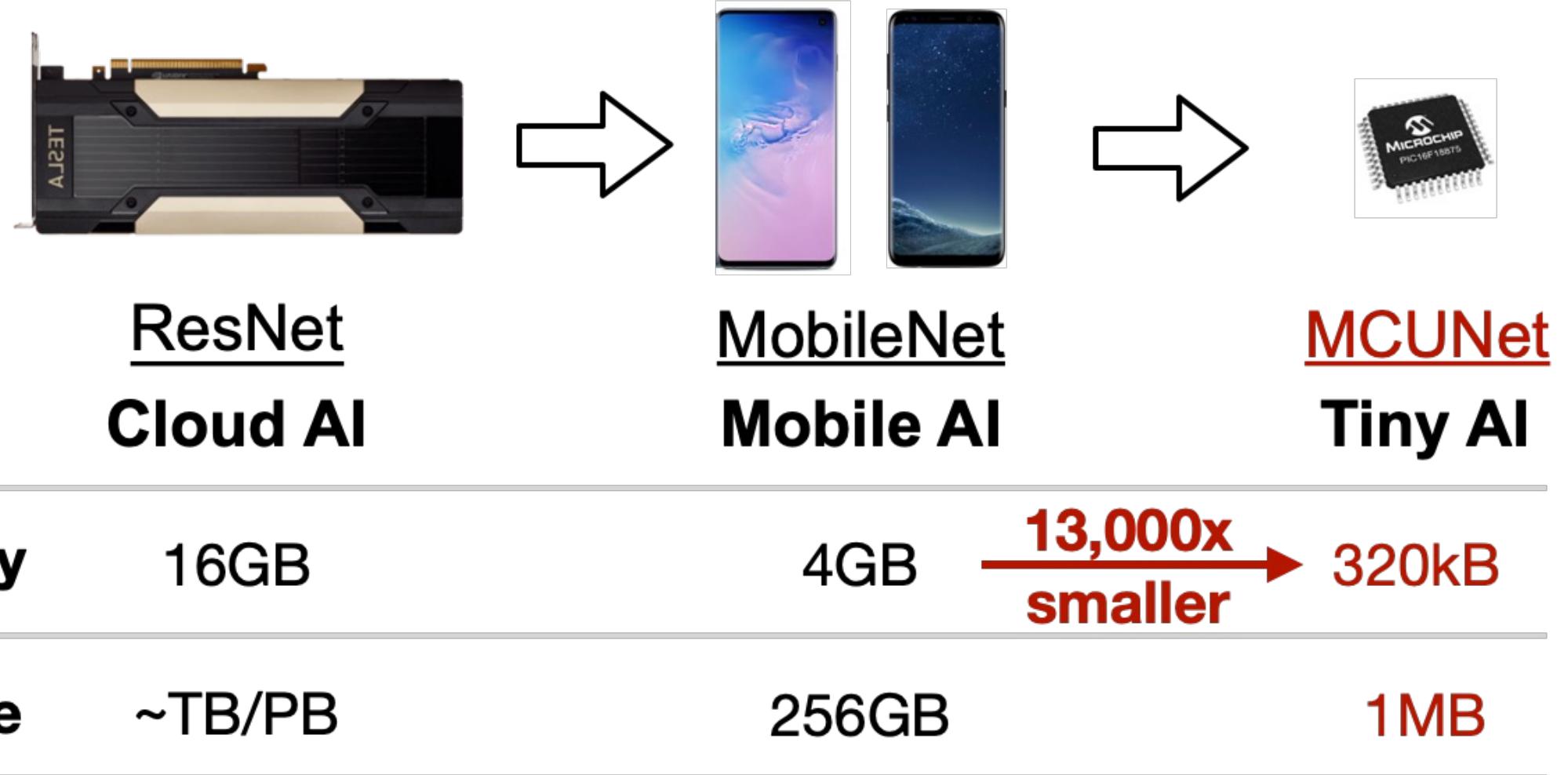
Wei-Chen Wang  
 Song Han  
 Department of EECS  
 MIT  
 wweichen@mit.edu

# MCUNetV1 & V2: On-Device Inference of Tiny Deep Learning on IoT Devices

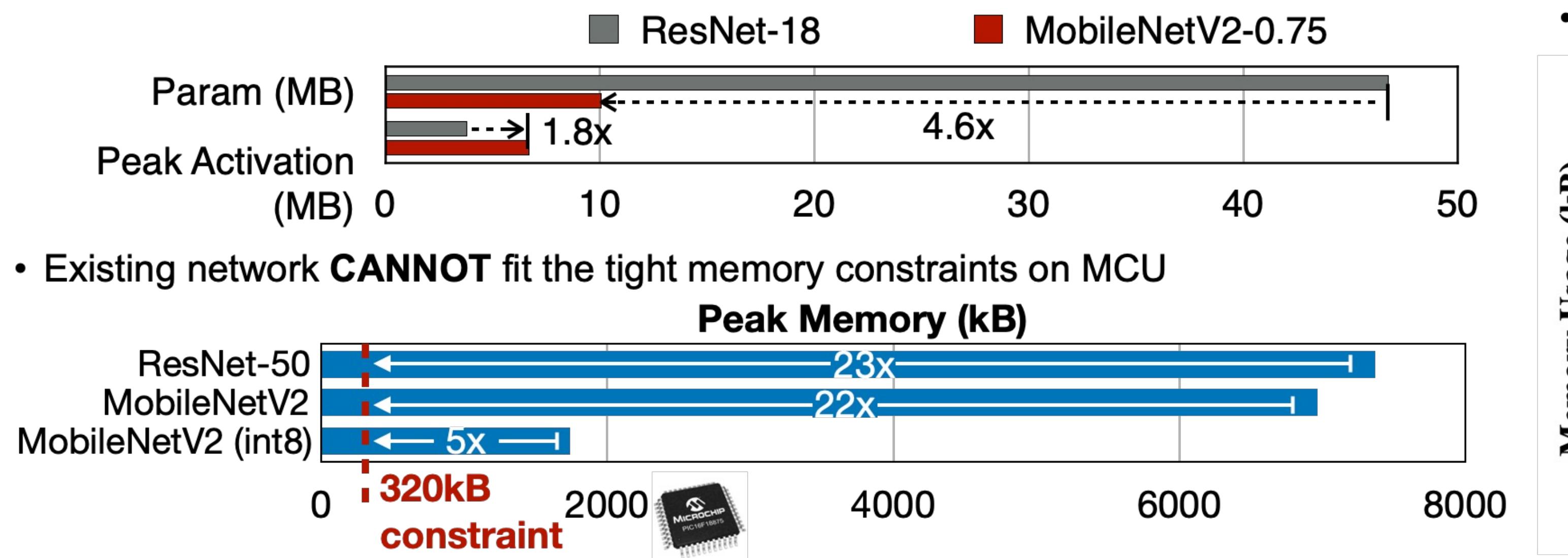
Ji Lin, Wei-Ming Chen, Yujun Lin, Han Cai, John Cohn, Chuang Gan, Song Han  
 MIT, MIT-IBM Watson AI Lab



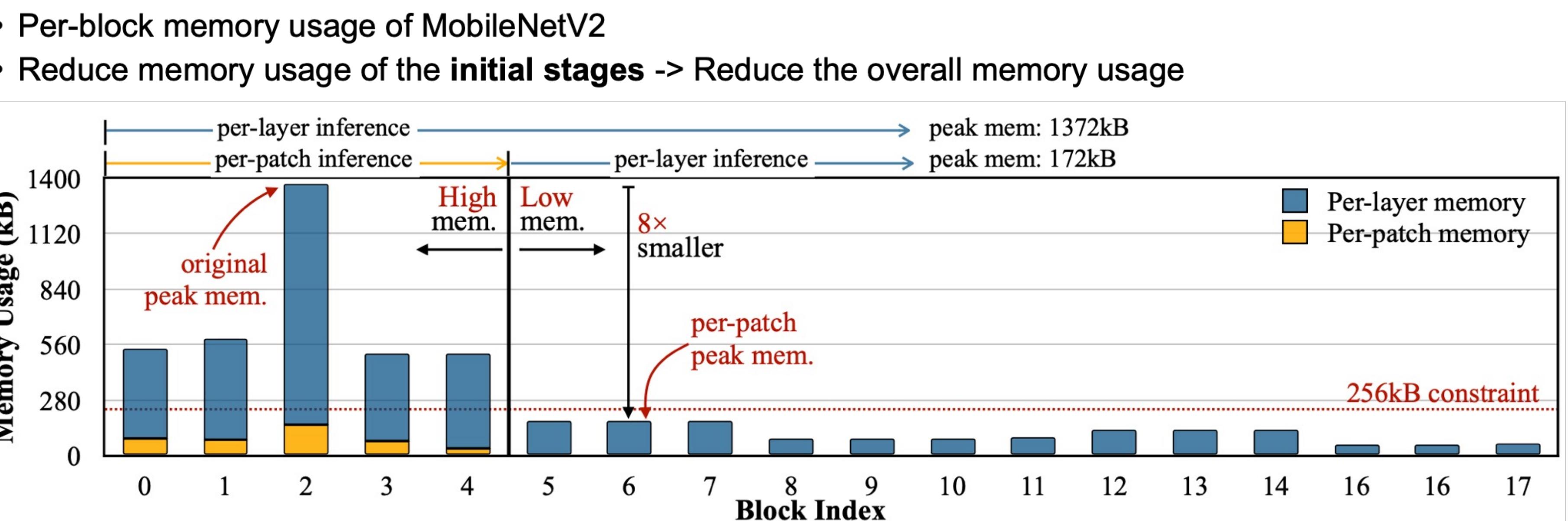
## Tiny Machine Learning (TinyML) Faces Challenge of Limited Memory



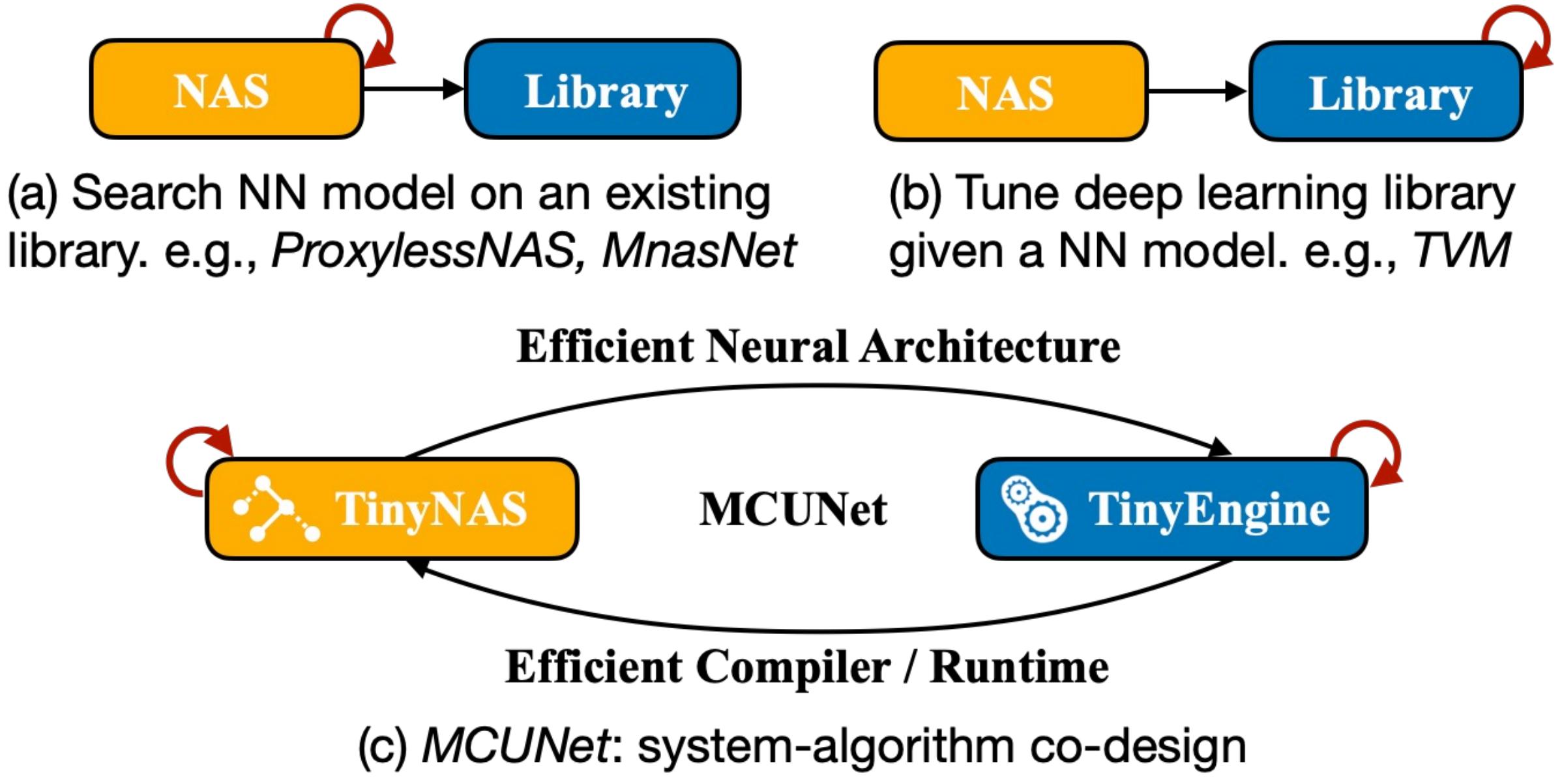
## Challenge 1: Existing Methods Reduce Model Size, but not the Activation Size



## Challenge 2: Efficient CNNs Have Imbalanced Memory Distribution



## MCUNet: System-Algorithm Co-design



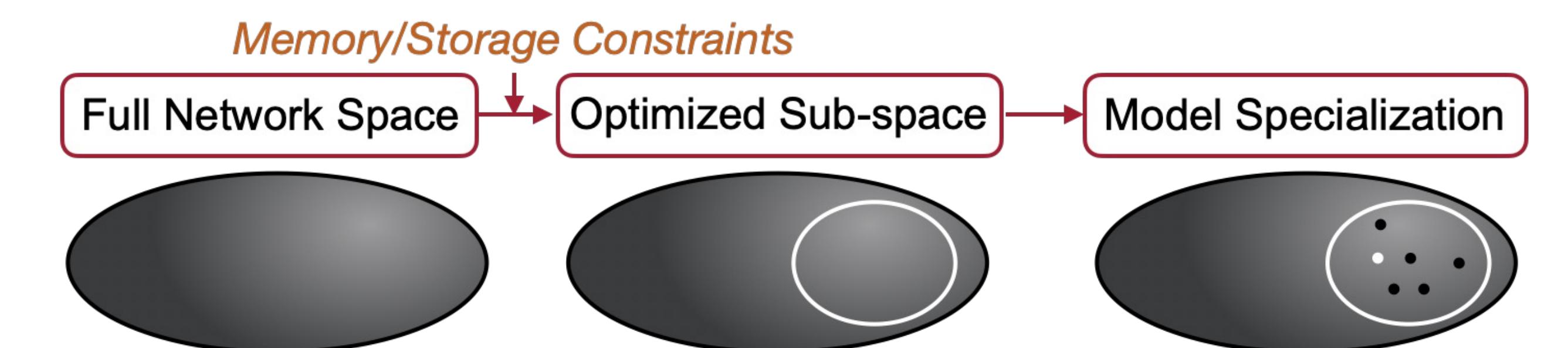
**TinyNAS** analyzes **FLOPs distribution** of satisfying models in each search space:

- Larger FLOPs -> Larger model capacity -> More likely higher accuracy

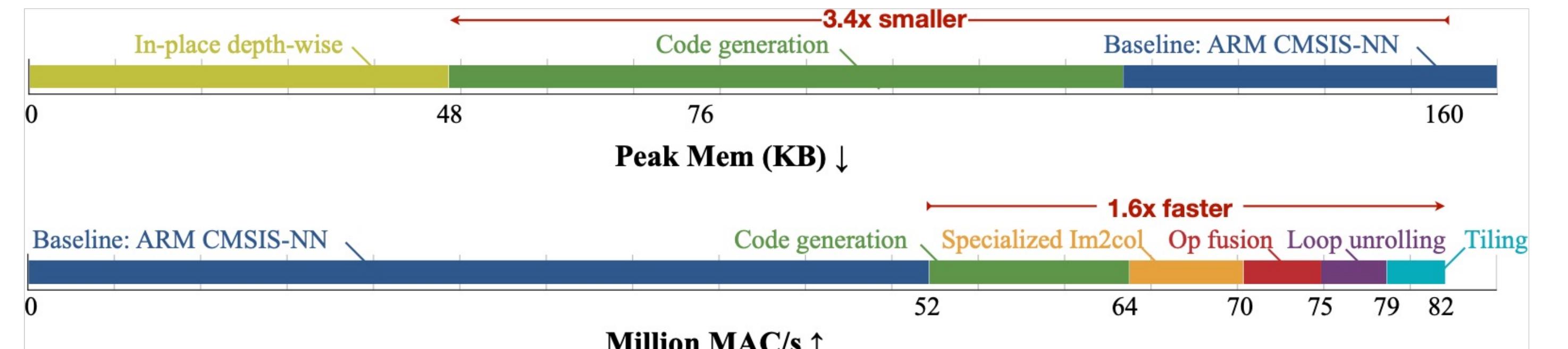
**TinyEngine** allows fitting a larger model at the same hardware resource by:

- Reduced peak memory usage • Accelerated inference speed

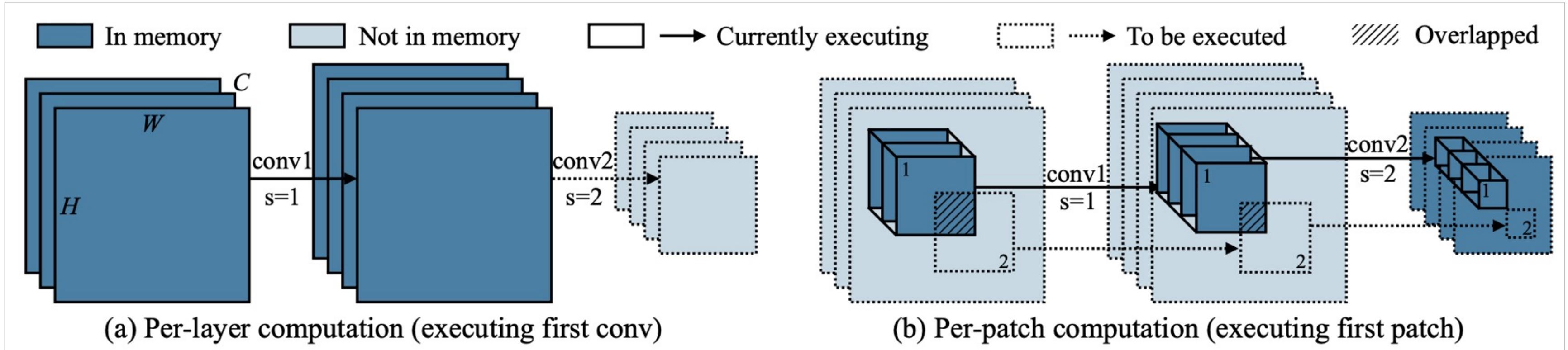
## 1. TinyNAS: Two-Stage NAS for Tiny Memory



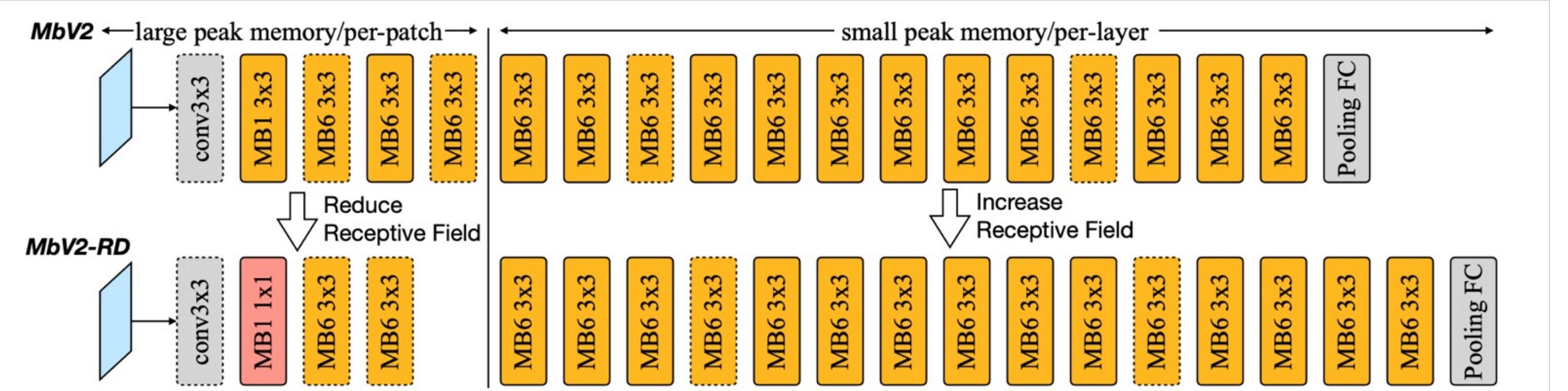
## 2. TinyEngine: Memory-Efficient Inference Library



## 3. Patch-based Inference to Break Memory Bottleneck



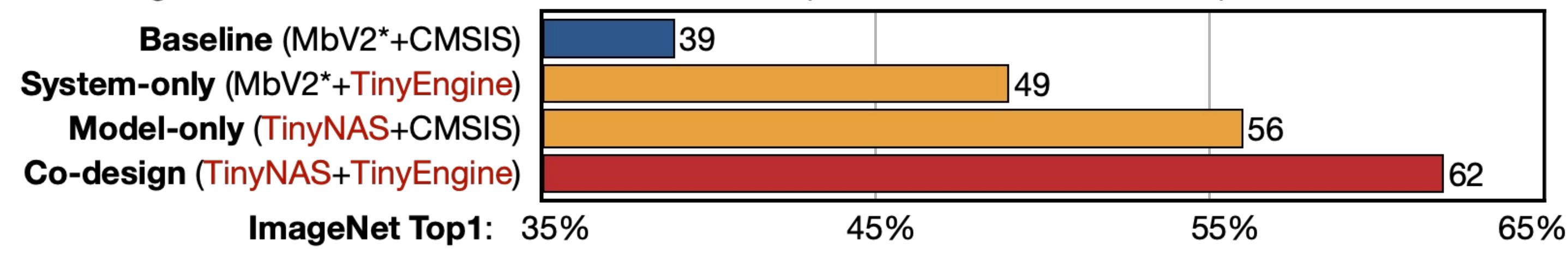
## 4. Network Redistribution to Reduce Computation Overhead



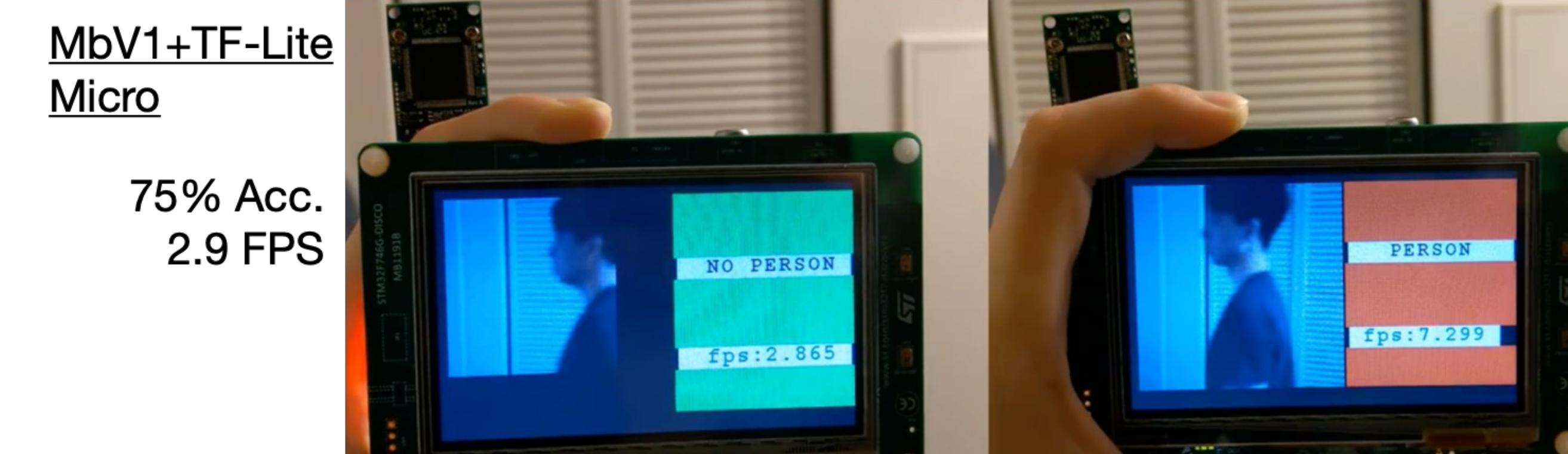
## Experimental Results

### System-algorithm co-design gives the best performance

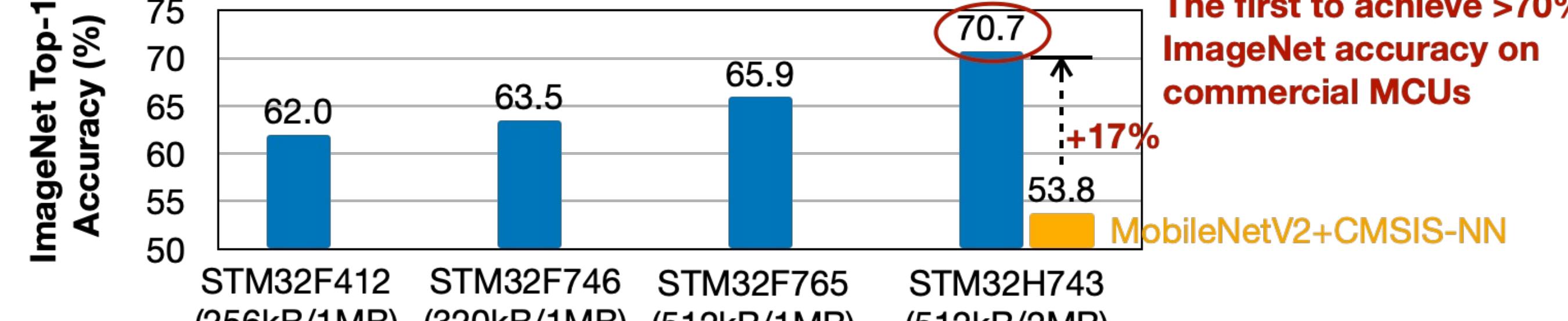
- ImageNet classification on STM32F746 MCU (320kB SRAM, 1MB Flash)



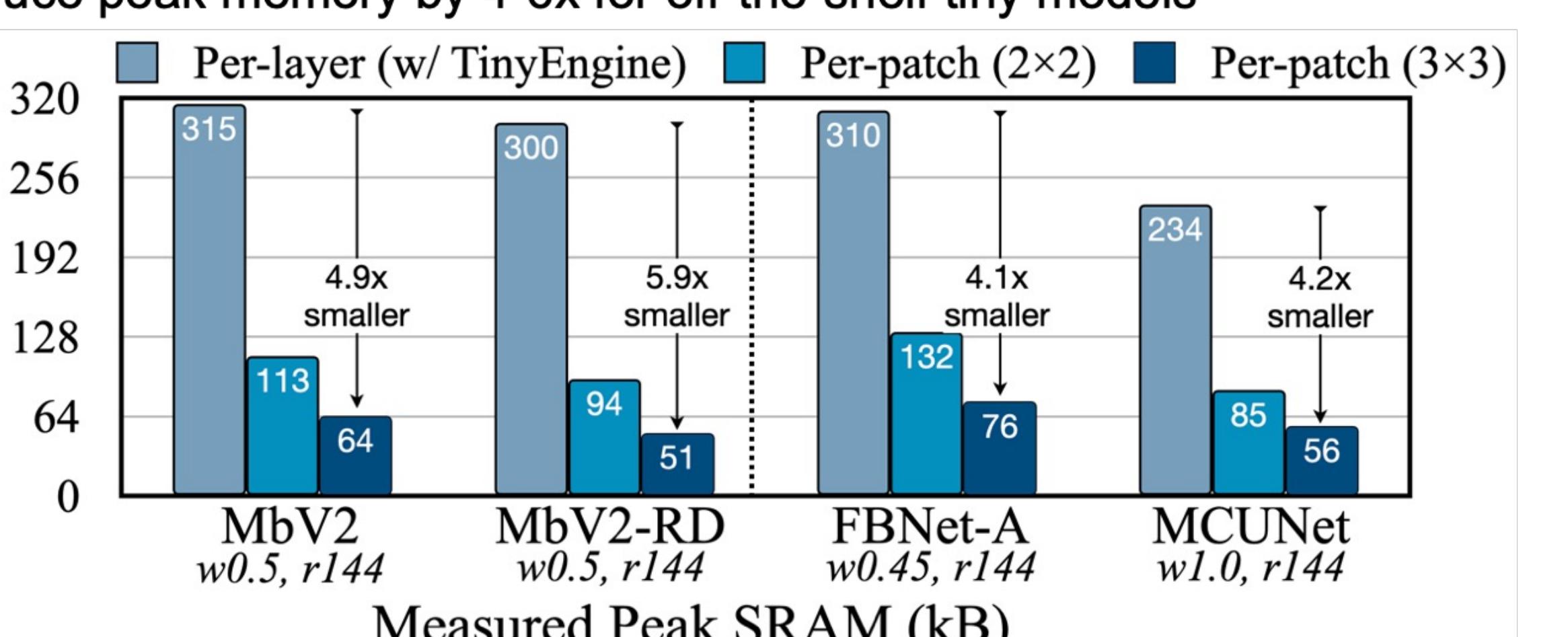
### On-device deployment demo for visual wake word (person or no person)



- MCUNet automatically handles **diverse hardware capacity** by optimizing search spaces



- Reduce peak memory by 4-6x for off-the-shelf tiny models



- Face detection on WIDER Face
- More robust results at a smaller peak memory

- Patch-based inference allows for a **larger resolution**, improving detection performance

