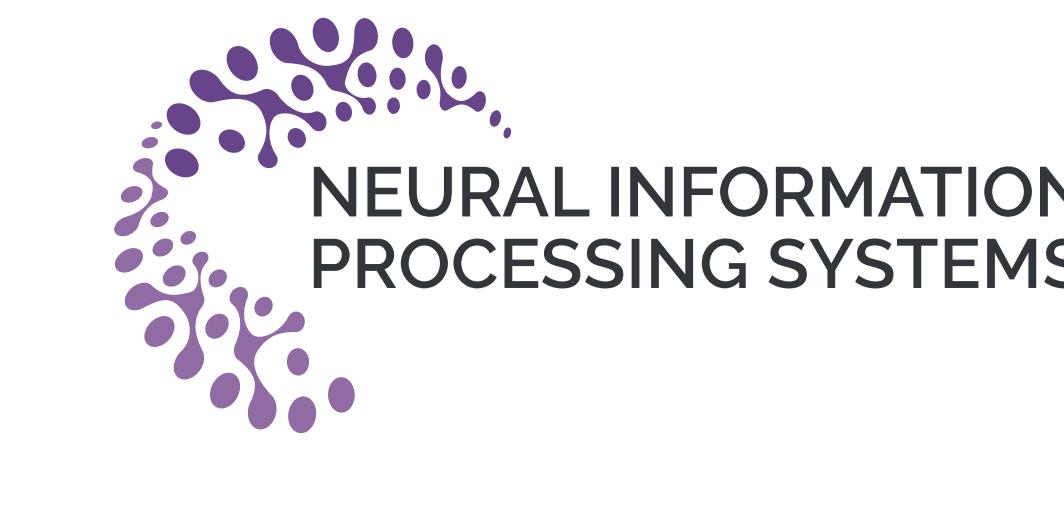


MCUNet: Tiny Deep Learning on IoT Devices

Ji Lin¹, Wei-Ming Chen^{1,2}, Yujun Lin¹, John Cohn³, Chuang Gan³, Song Han¹
¹MIT ²National Taiwan University ³MIT-IBM Watson AI Lab



The Era of TinyML on AIoT with Microcontrollers (MCUs)

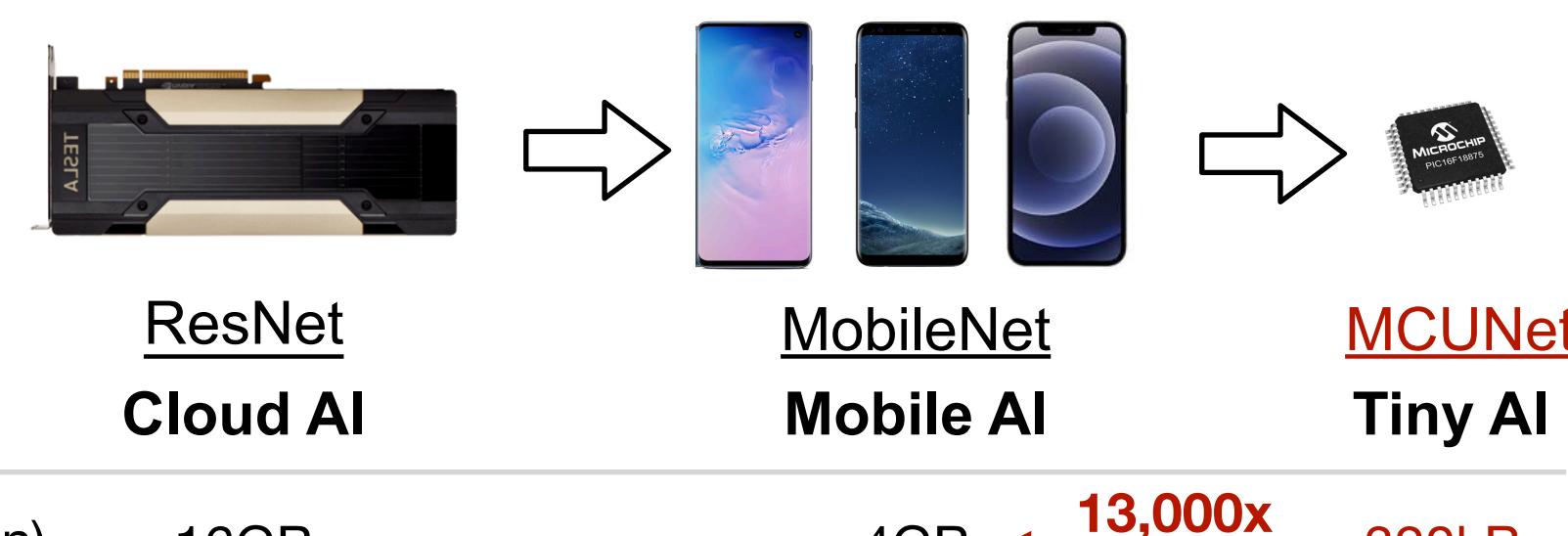
Microcontrollers: low-cost, low-power, widely deployed



Wide Applications

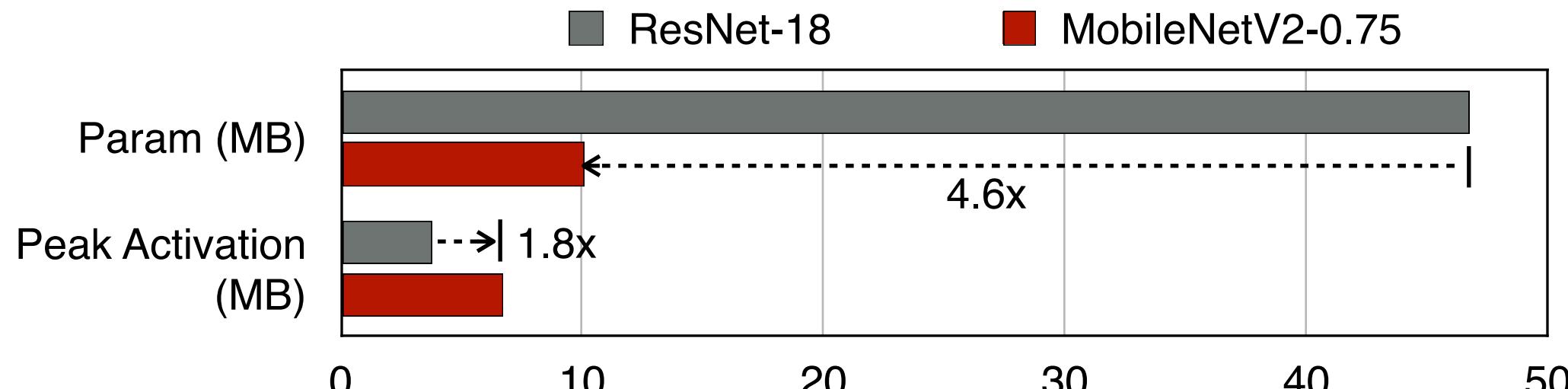


Challenge: Memory Too Small to Hold DNNs

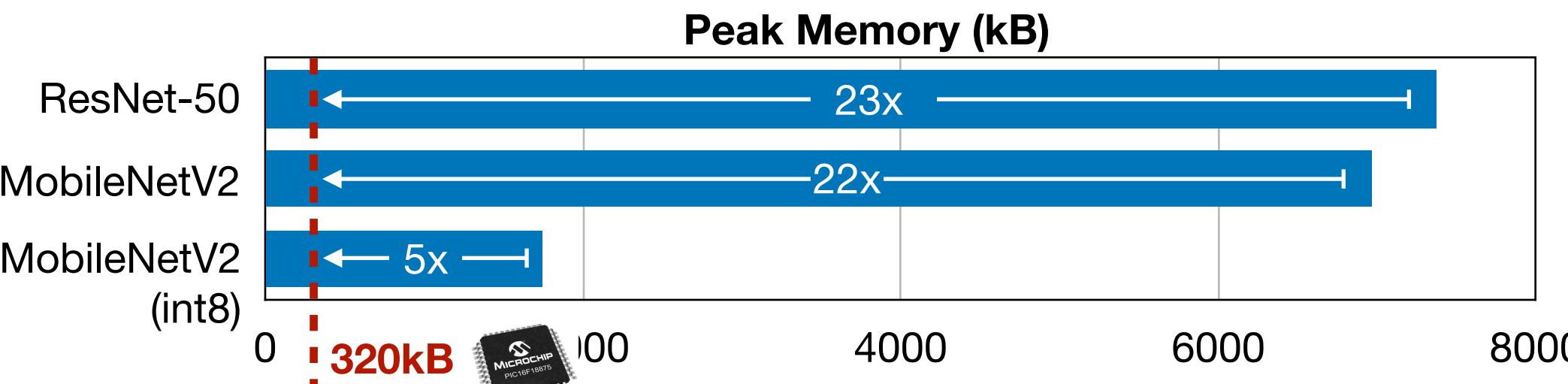


	ResNet Cloud AI	MobileNet Mobile AI	MCUNet Tiny AI
Memory (Activation)	16GB	4GB	320kB
Storage (Weights)	~TB/PB	256GB	1MB

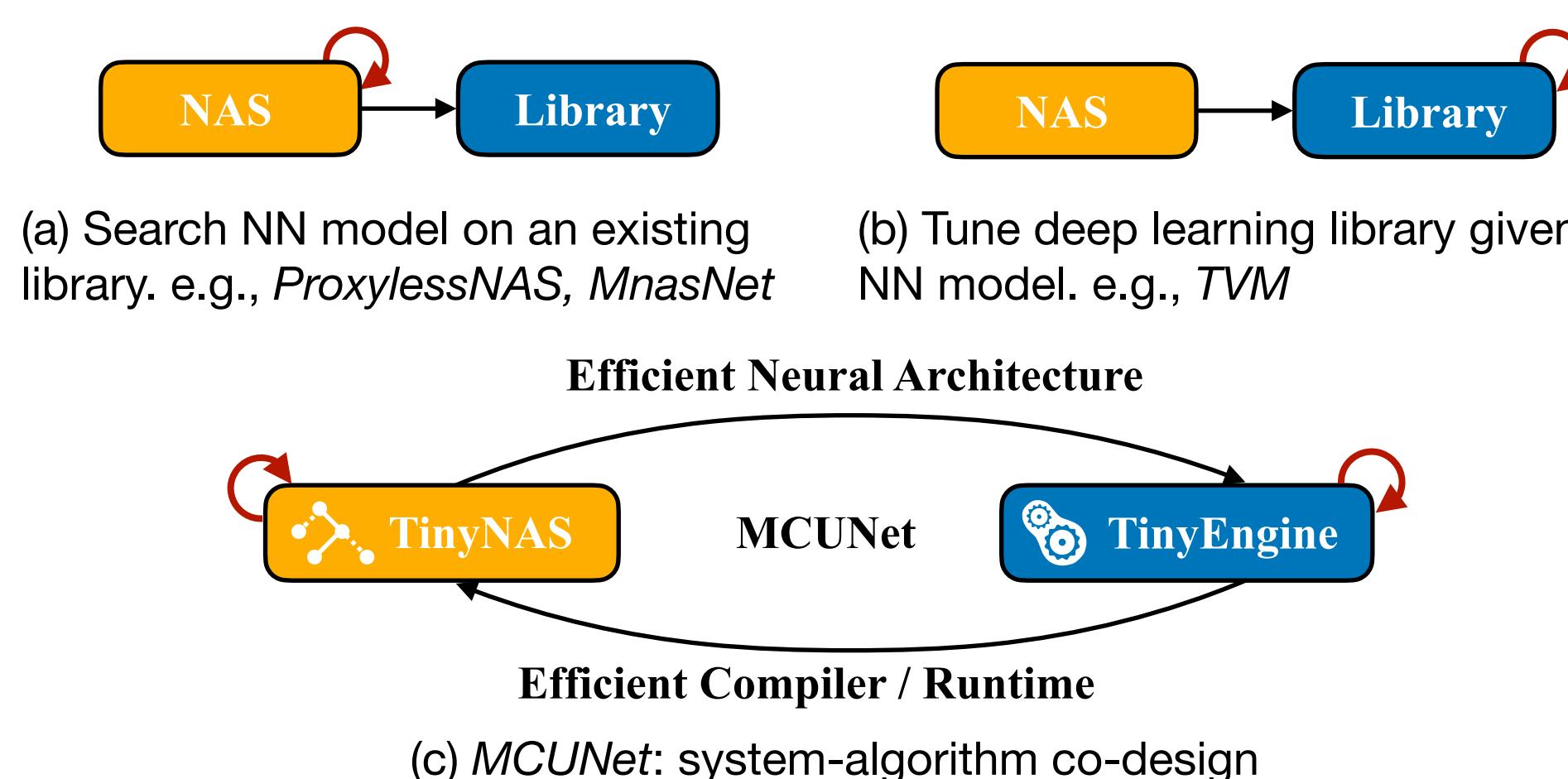
Existing Methods Reduce Model Size, but not the Activation Size



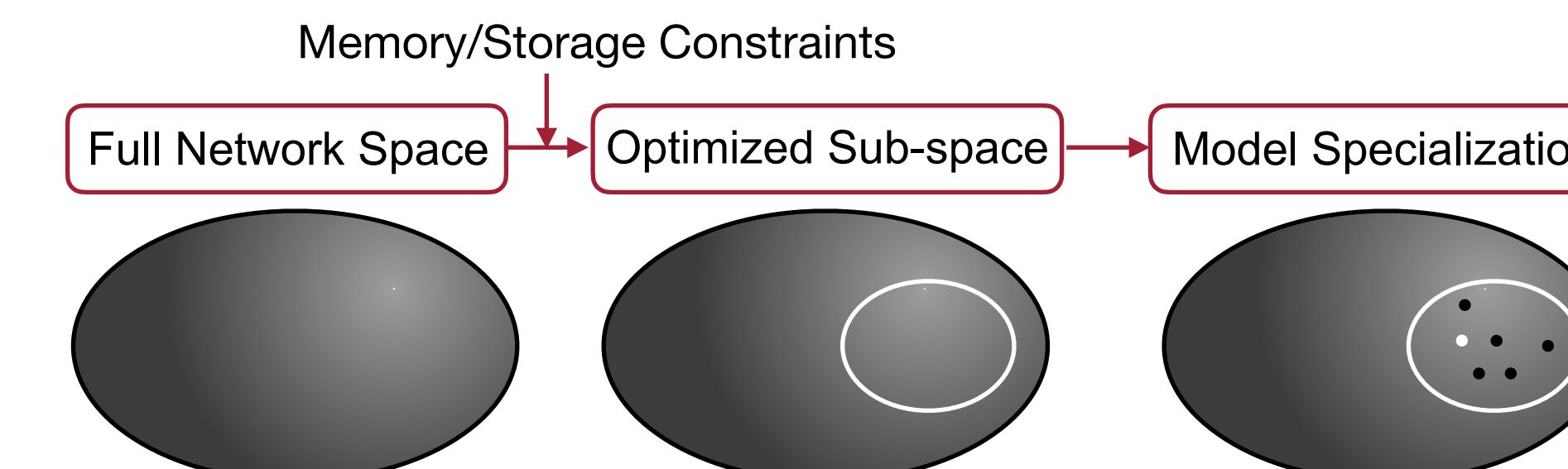
- Existing network **CANNOT** fit the tight memory constraints on MCU



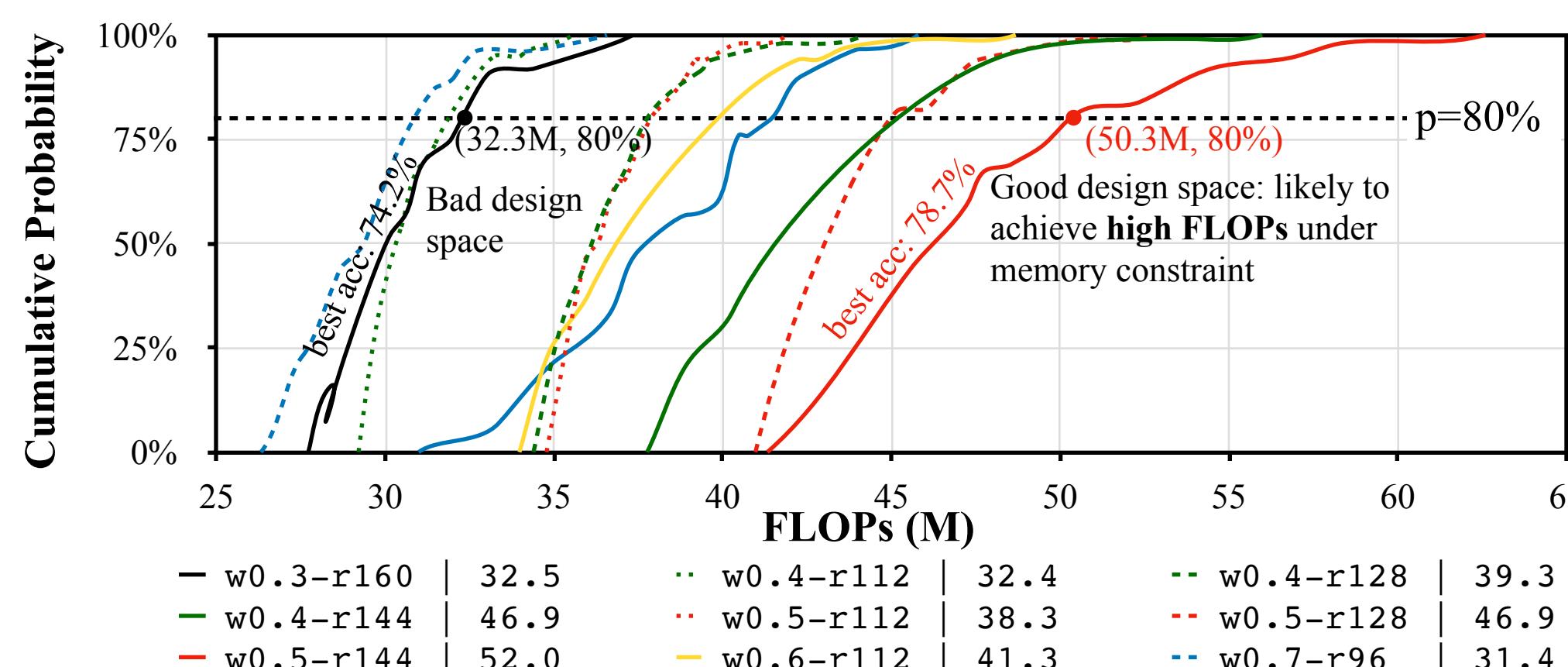
MCUNet: System-Algorithm Co-design



1. TinyNAS: Two-Stage NAS for Tiny Memory



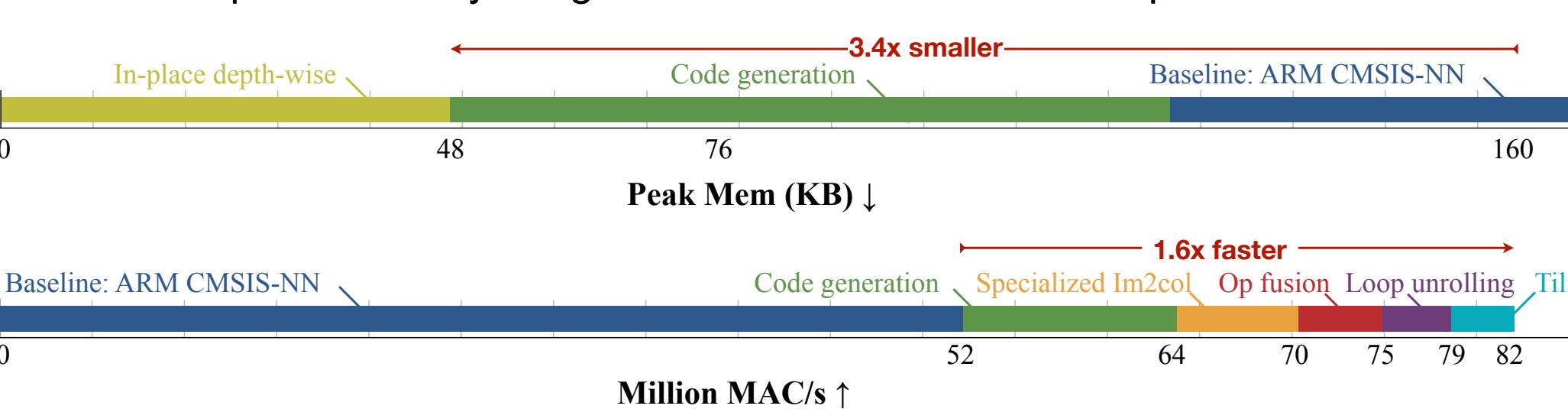
- Analyzing **FLOPs distribution** of satisfying models in each search space:
Larger FLOPs → Larger model capacity → More likely to give higher accuracy



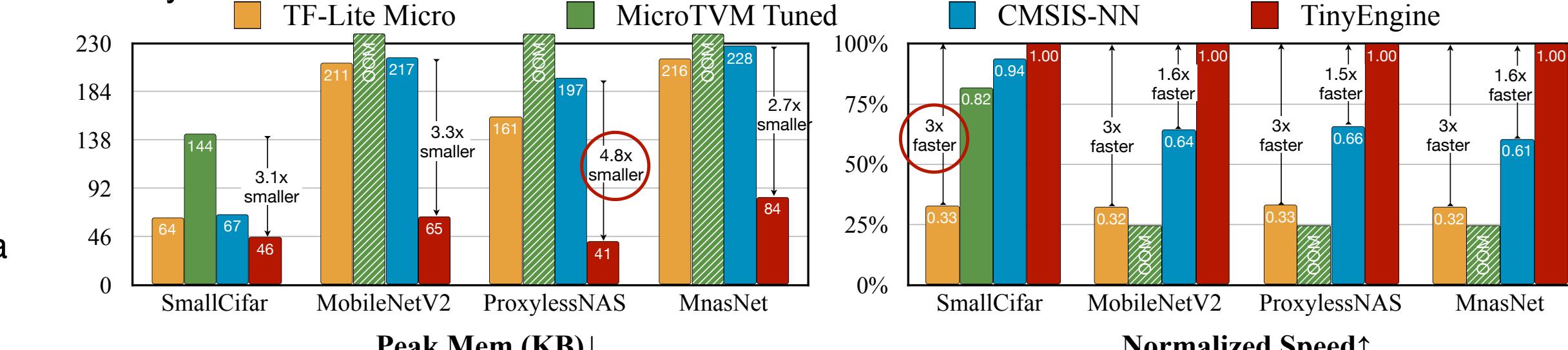
2. TinyEngine: Memory-Efficient Inference Library

TinyEngine allows to fit a larger model at the same hardware resource by:

- Reduced peak memory usage
- Accelerated inference speed

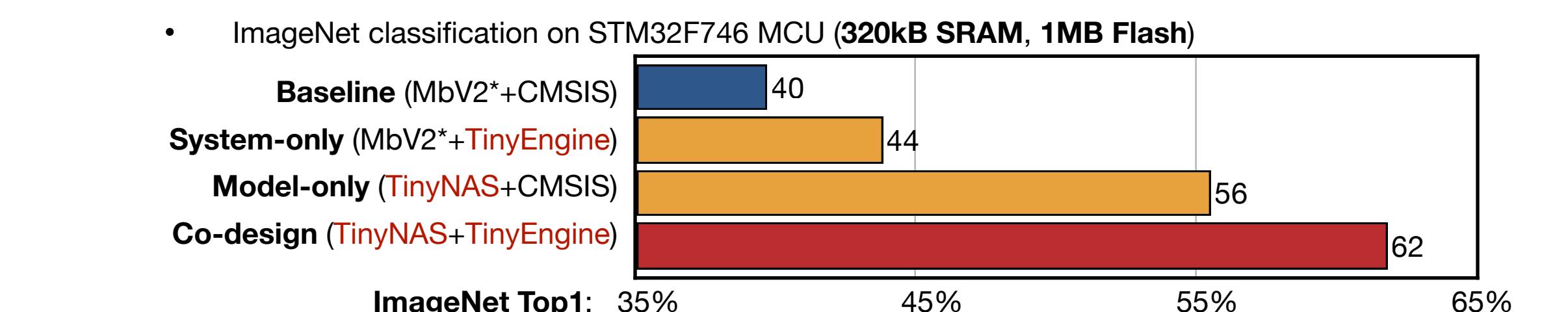


- TinyEngine consistently improves on different networks

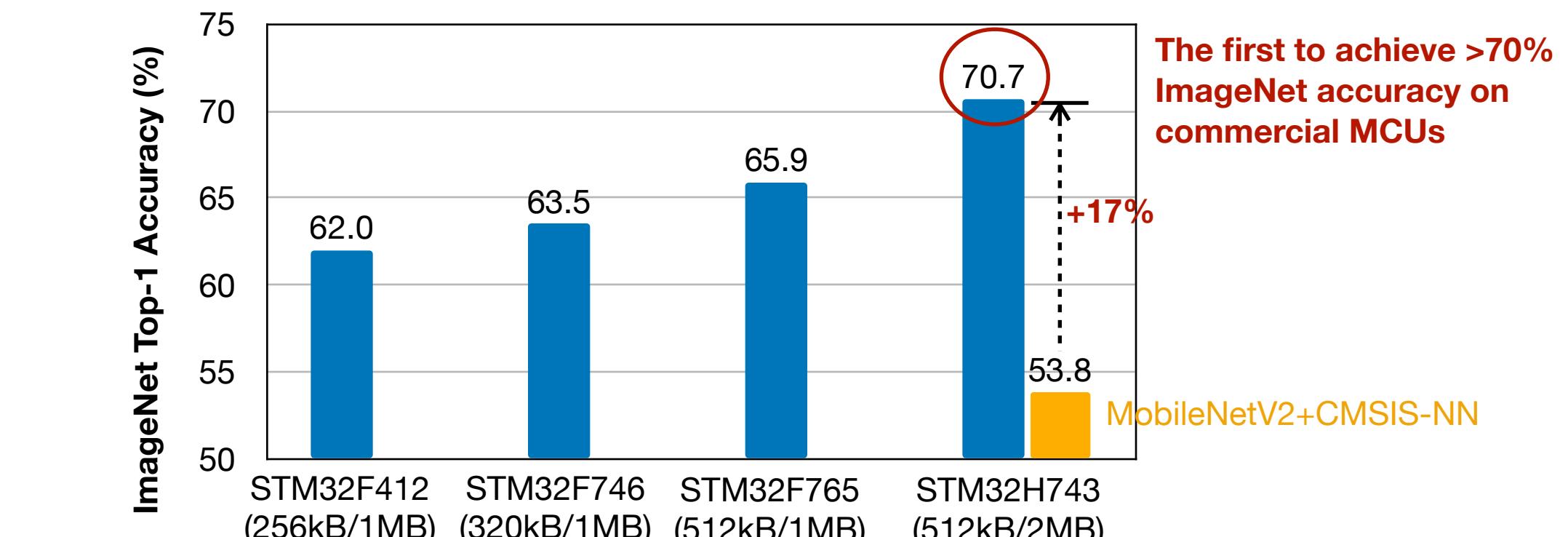


Experimental Results

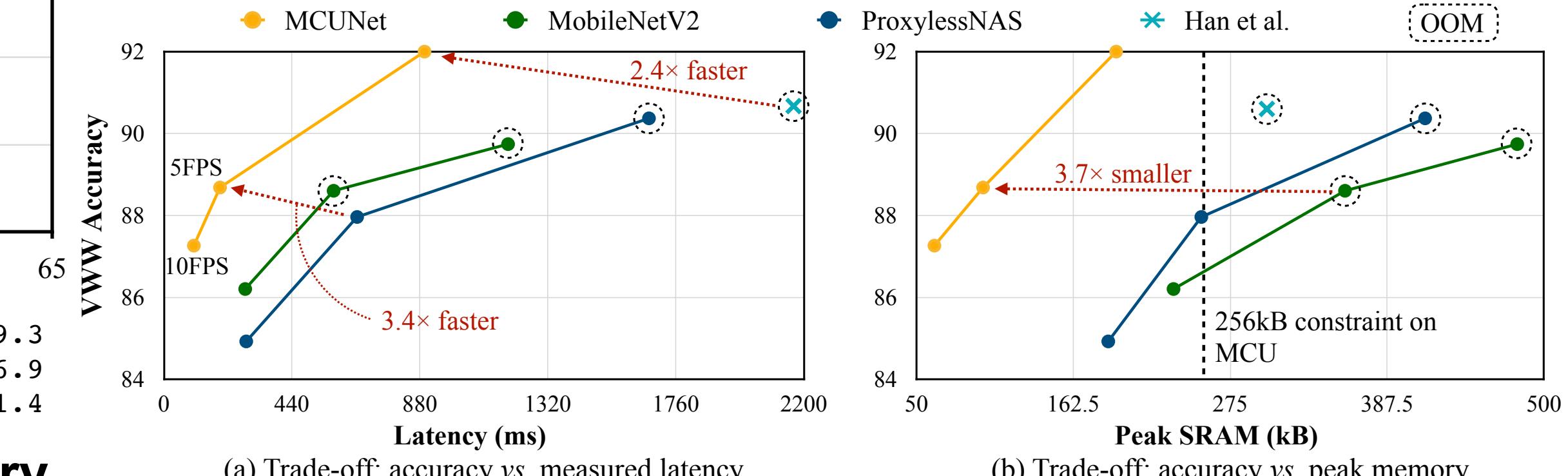
- System-algorithm co-design** gives the best performance



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



- MCUNet significantly outperforms existing solutions on **visual/audio keyword spotting**



- On-device deployment** demo or visual wake word (person present or not)

