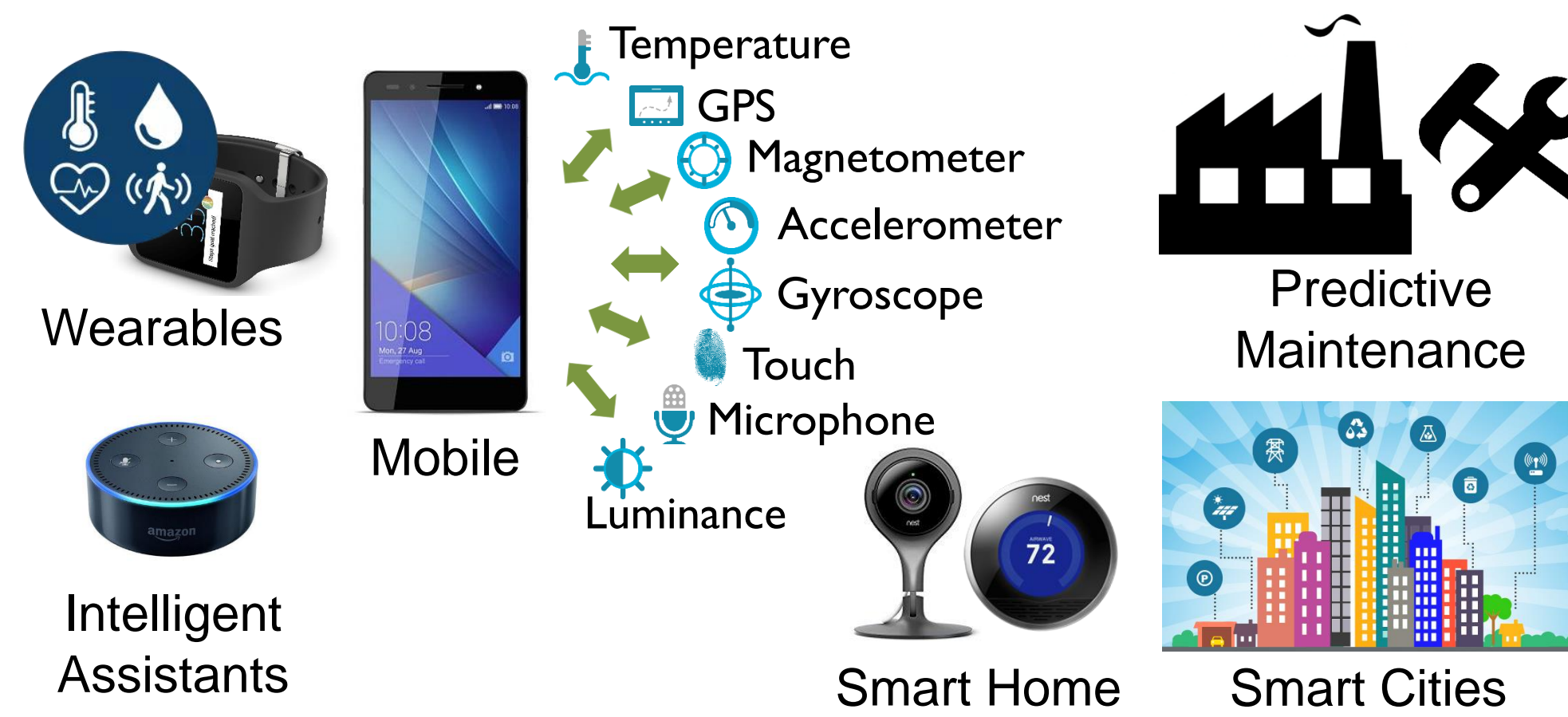


Ternary Hybrid Neural-Tree Networks for Highly Constrained IoT Applications

Dibakar Gope, Ganesh Dasika, Matthew Mattina
Arm ML Research Lab

Challenge

- ML algorithms are increasingly deployed at the edge in IoT devices



- These devices are highly constrained in both memory and compute budget
- Aggressive model compression is required to
 - Target severely constrained microcontrollers
 - Deploy more IoT applications on them



BBC Micro:Bit
Arm Cortex M0 (16KB SRAM)



LPCXpresso 1125
Arm Cortex M0 (8KB SRAM)

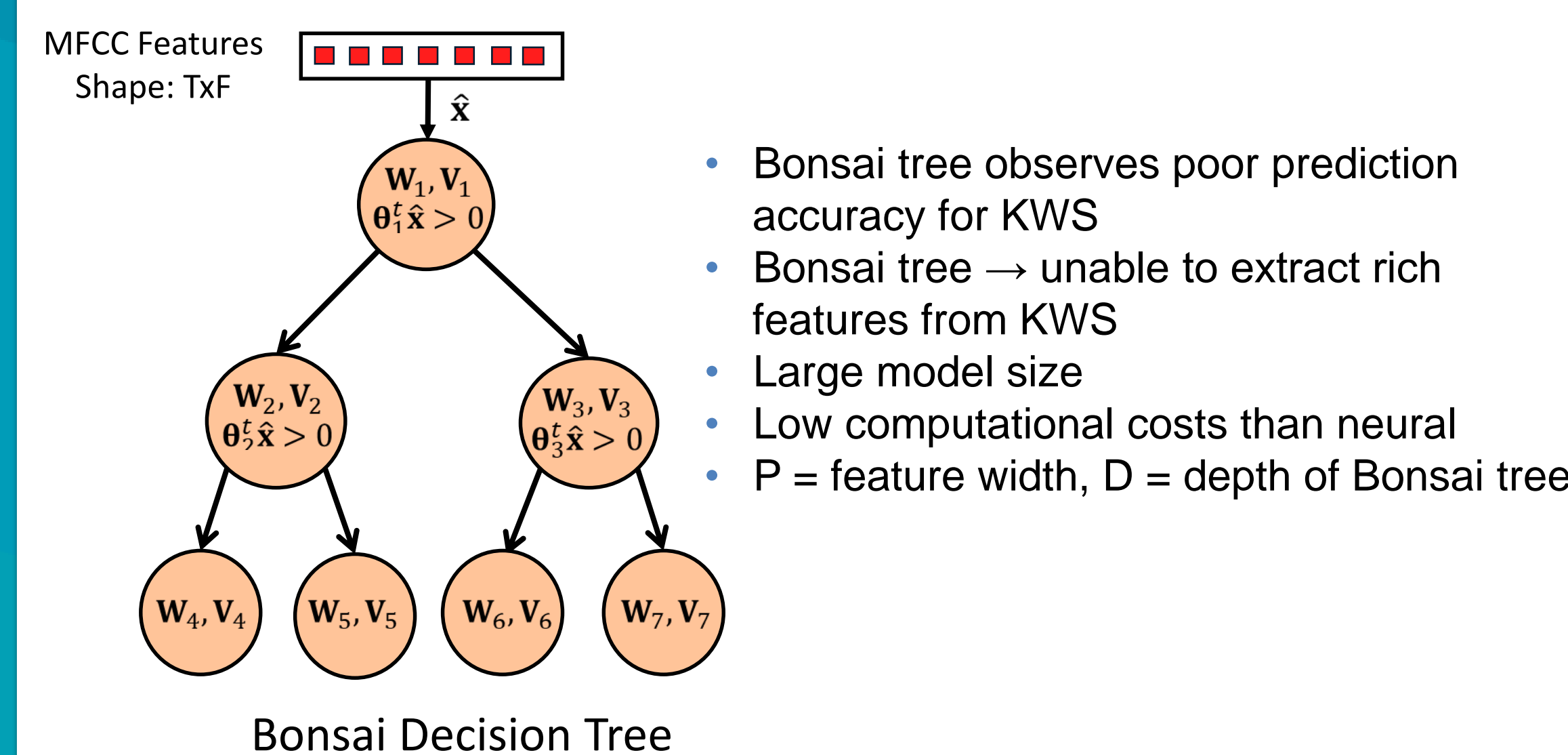
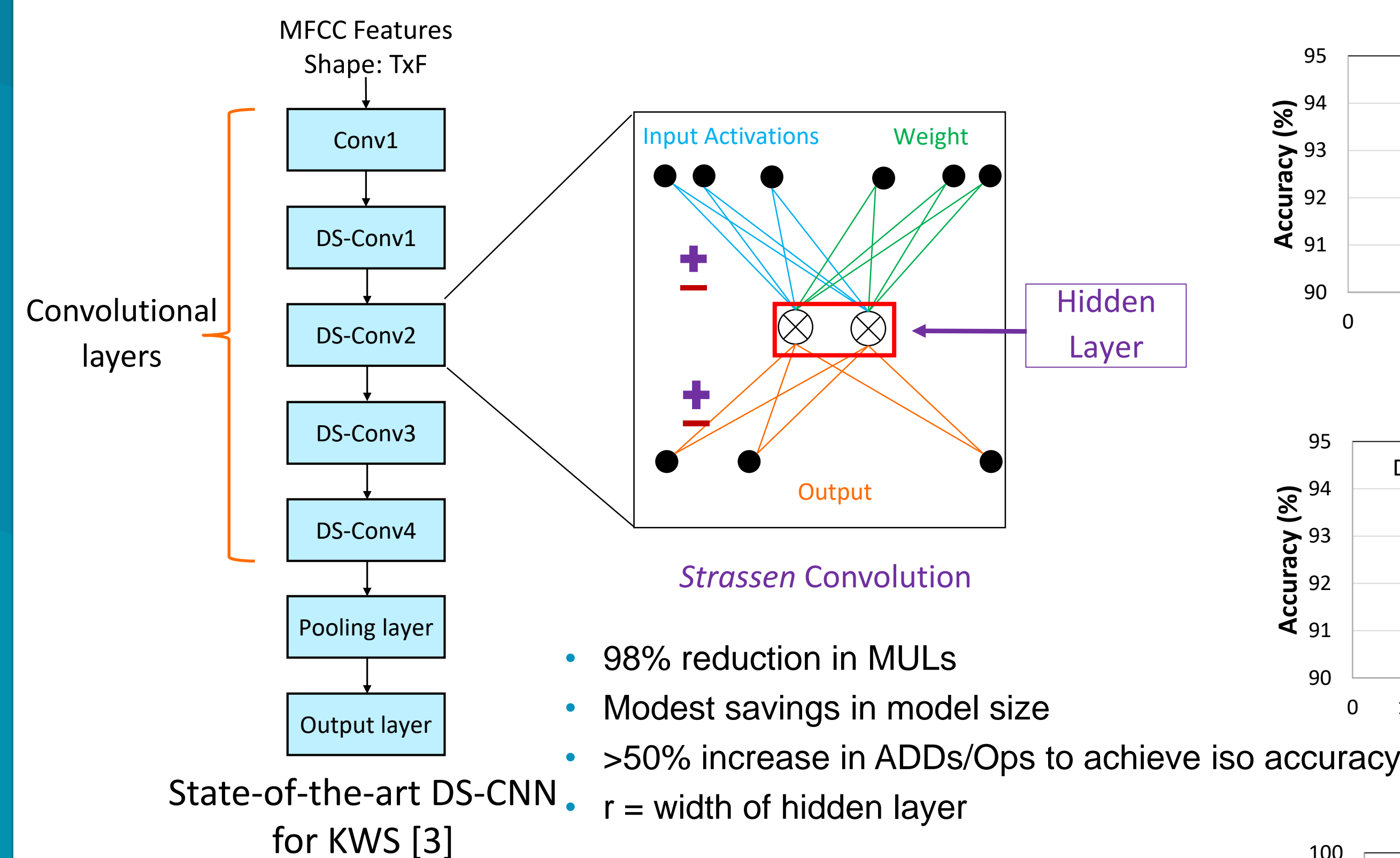
Prior Solutions

- Architectural optimization
 - Depthwise separable (DS) convolution
 - Bonsai decision trees [1]
- Model quantization
 - Binary/ternary quantization
 - StrassenNets [2]
- Model pruning
- Low-rank matrix factorization

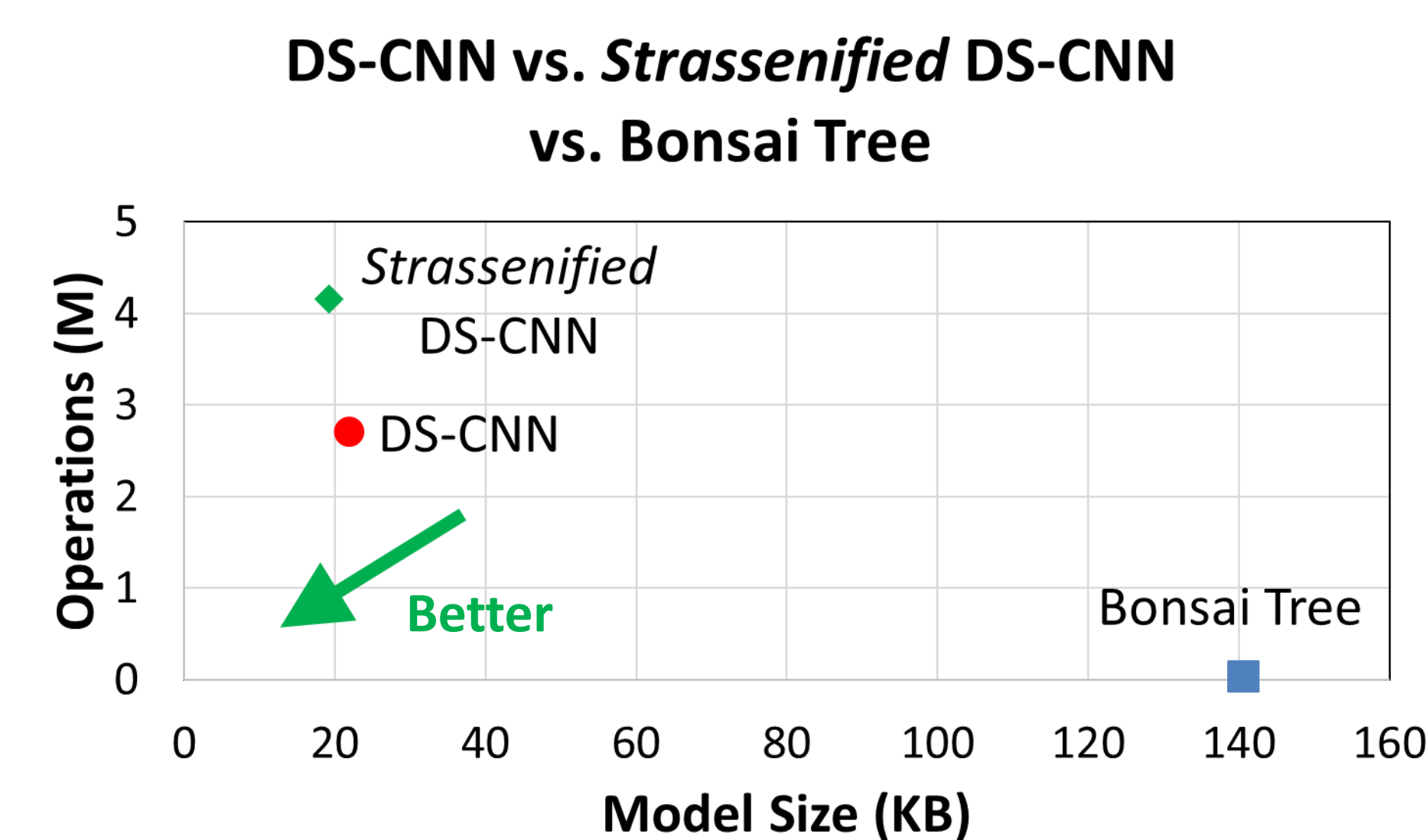
Problem

- Prior solutions cannot offer better compression than the state-of-the-art for popular IoT applications like keyword spotting (KWS)
- They come with their own advantages and limitations

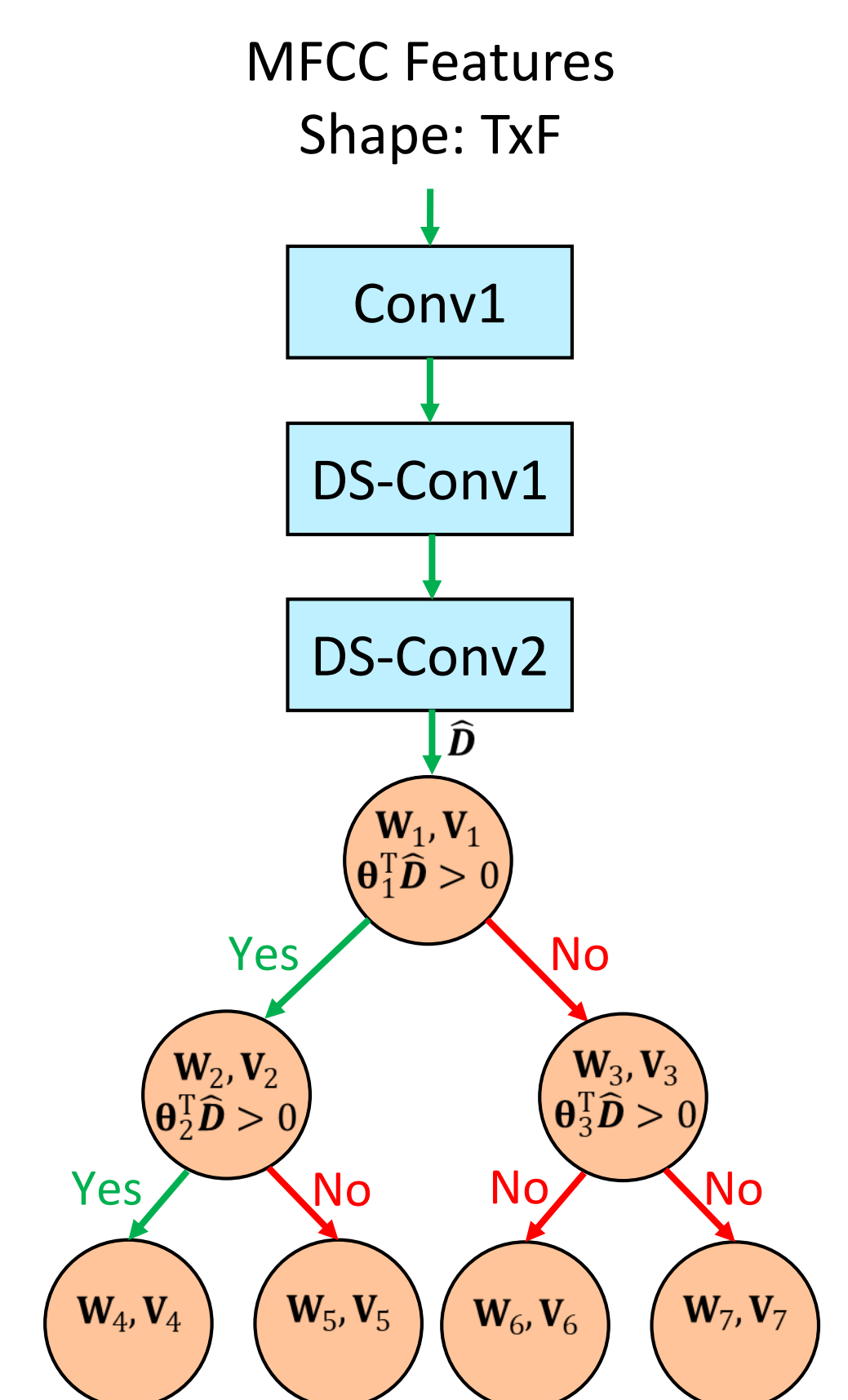
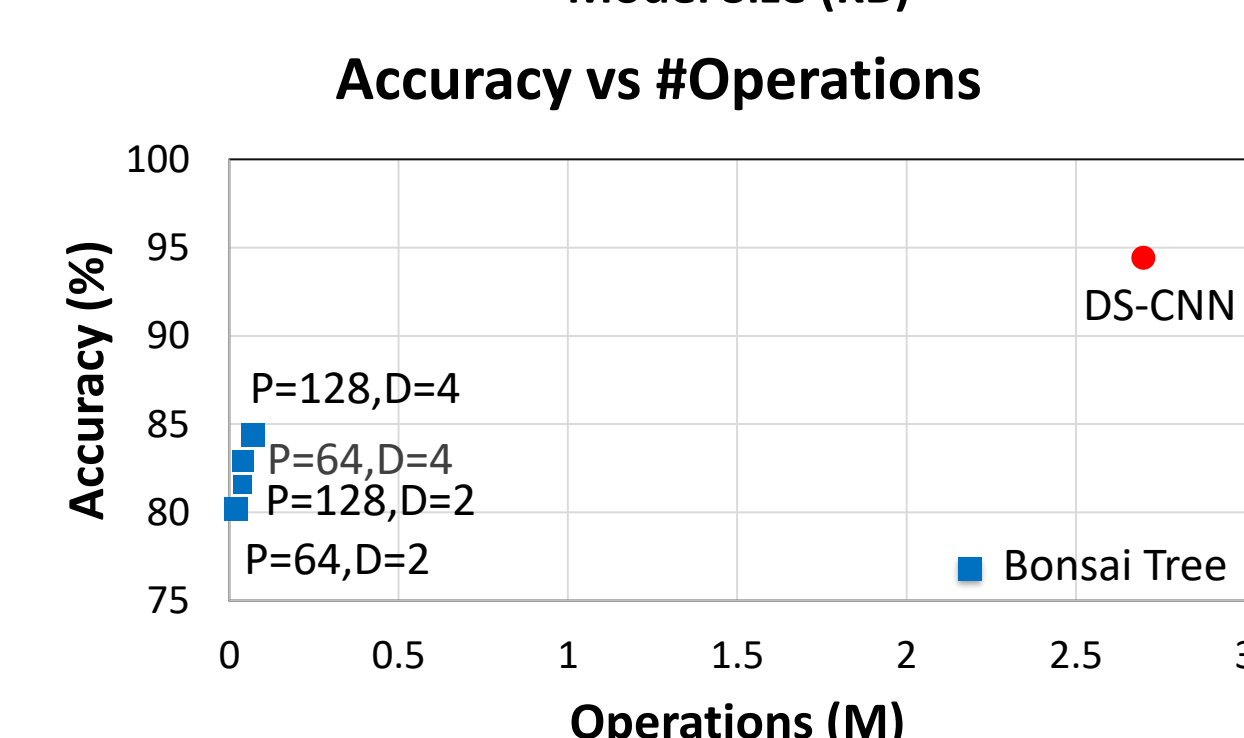
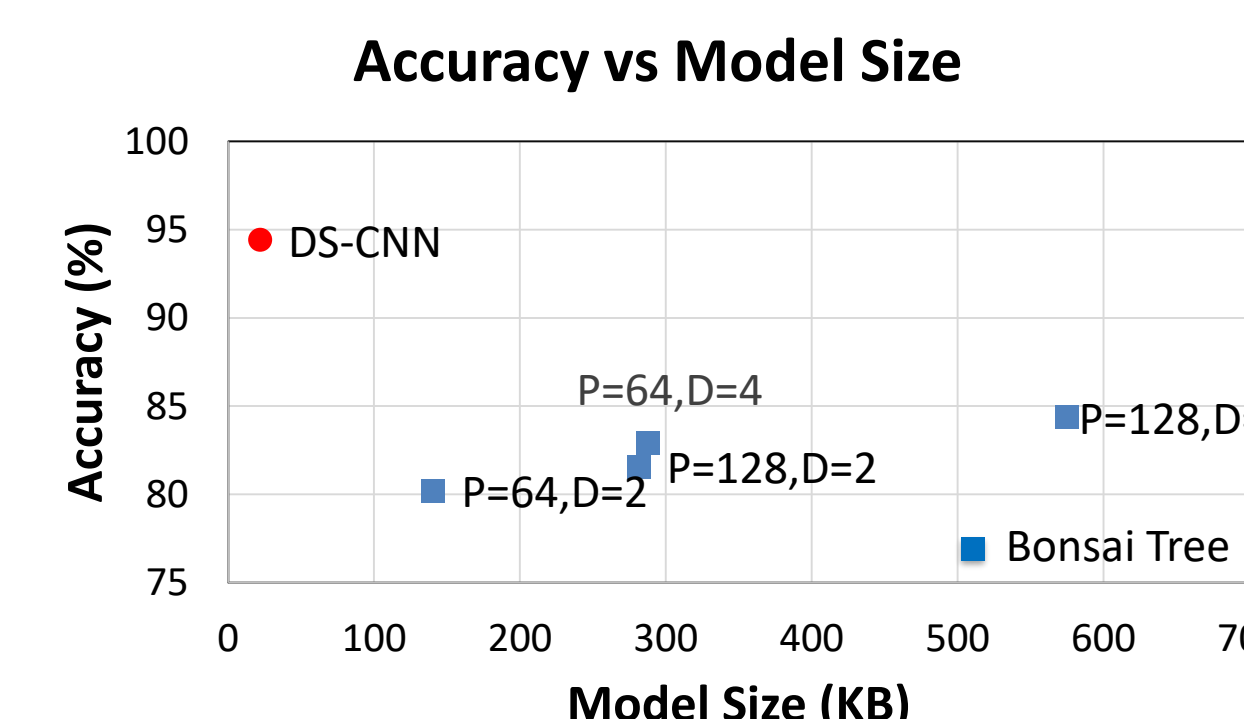
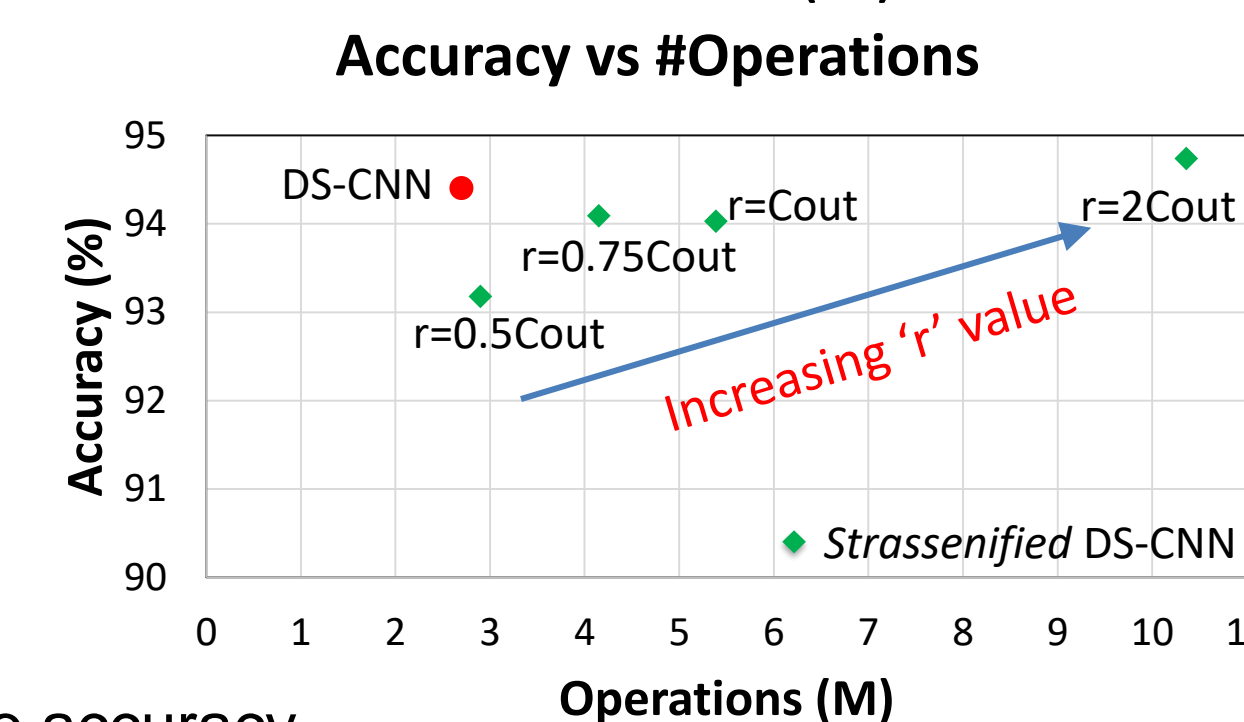
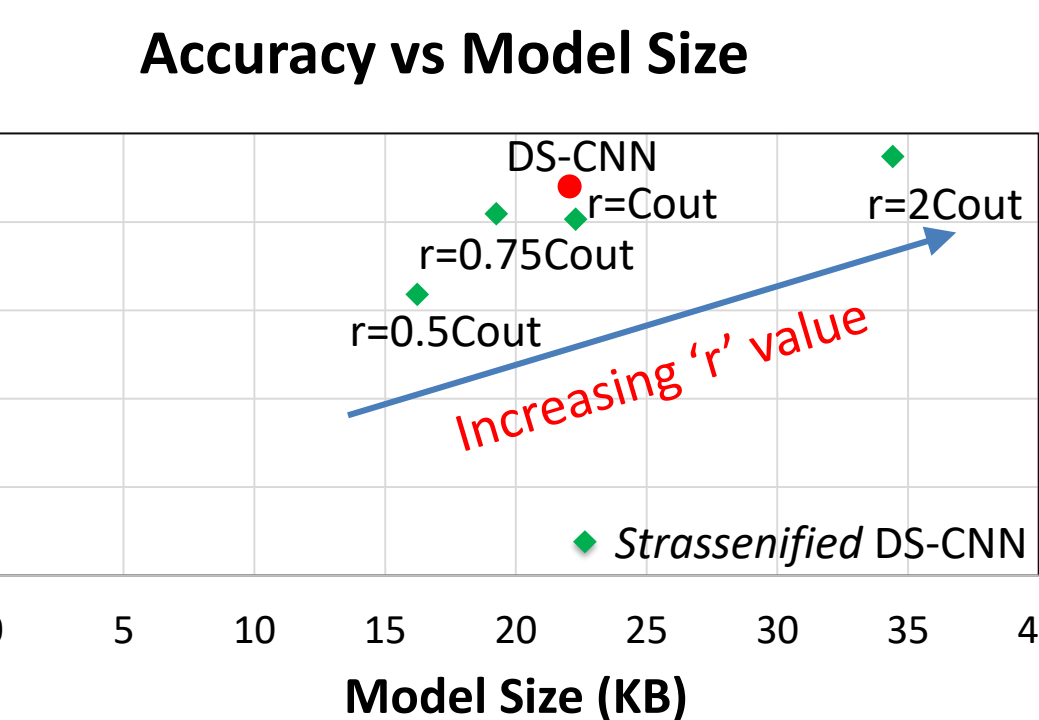
Observations with Prior Solutions



Our Solution: Hybrid Neural-Tree Network



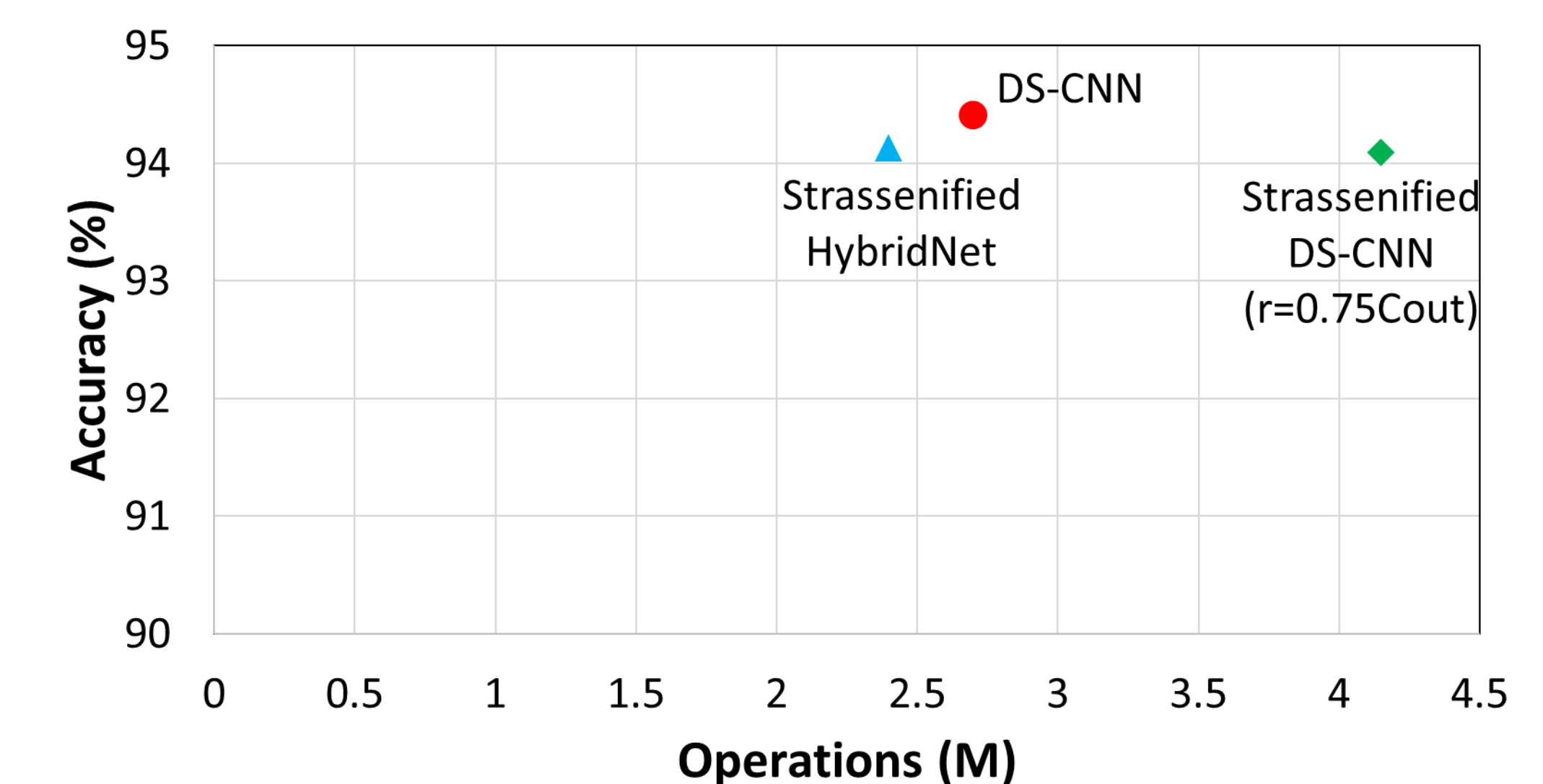
$$\hat{y} = \begin{pmatrix} W_1^T \bar{D} \tanh(V_1^T \bar{D}) \\ + W_2^T \bar{D} \tanh(V_2^T \bar{D}) \\ + W_4^T \bar{D} \tanh(V_4^T \bar{D}) \end{pmatrix}$$



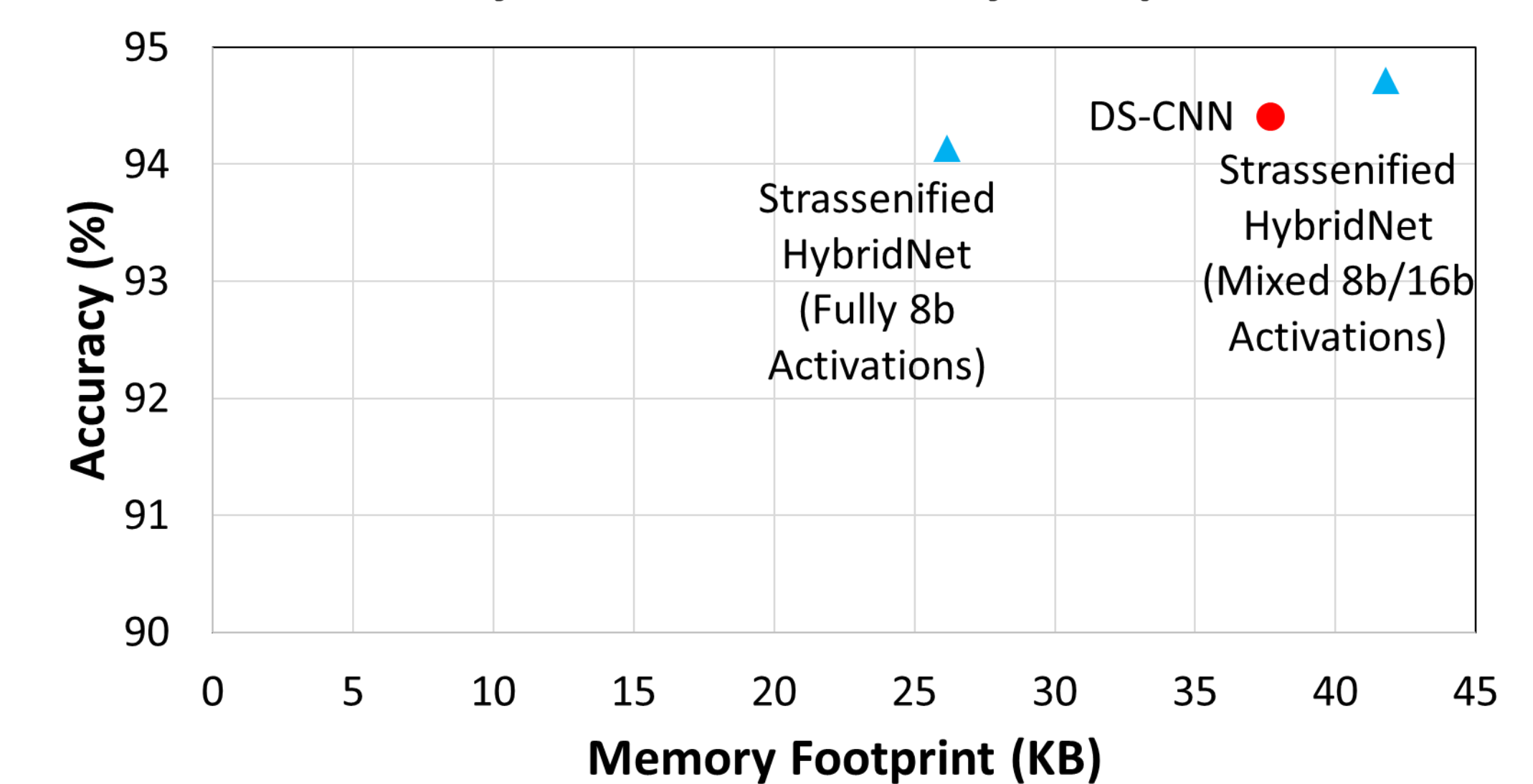
Evaluation Results

- 11.1% reduction in Ops, 98.9% reduction in MULs, 12.2% reduction in ADDs
- 52.2% reduction in model size, 30.6% reduction overall memory footprint (memory footprint = model size + size of intermediate activations)
- 0.27% loss in model accuracy
- Activations quantized to 8b fixed-point format
- 0.27% loss is attributed to high sensitivity of *strassenified* depthwise layers towards 8b quantizations

Accuracy vs #Operations



Accuracy vs Overall Memory Footprint



Accuracy, operations, and model size of hybrid network and improvement over state-of-the-art KWS networks

References

- [1] Kumar et al., "Resource-efficient Machine Learning in 2 KB RAM for the Internet of Things", ICML 2017
- [2] Tschannen et al., "StrassenNets: Deep Learning with a Multiplication Budget", ICML 2018
- [3] Zhang et al., "Hello Edge: Keyword Spotting on Microcontrollers", 2017