



On-Device Deep Learning for IoT-based Wireless Sensing Applications

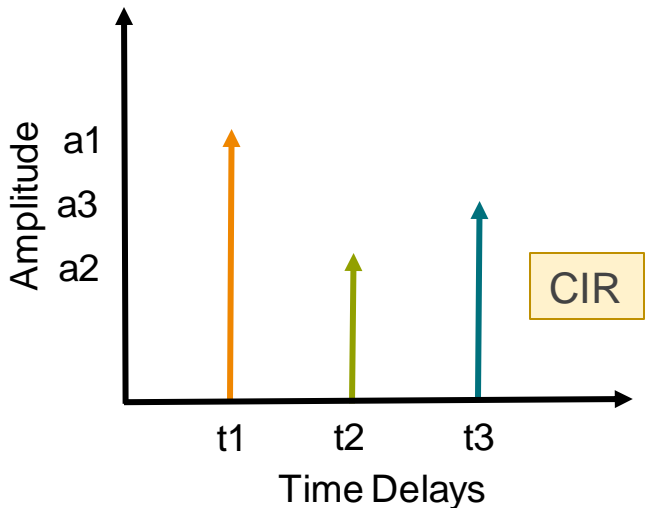
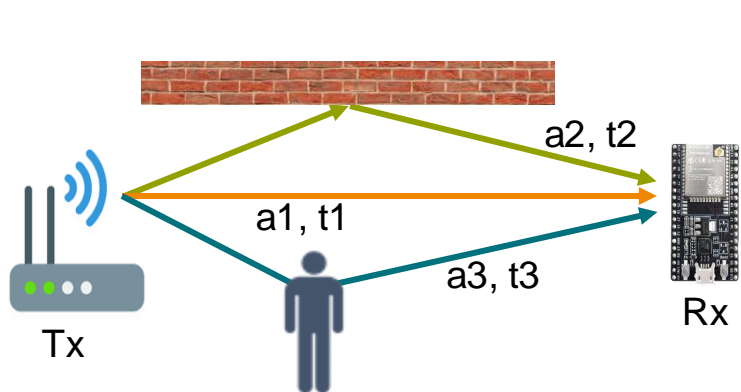
Manoj Lenka and Ayon Chakraborty

SENSE Lab
IIT Madras

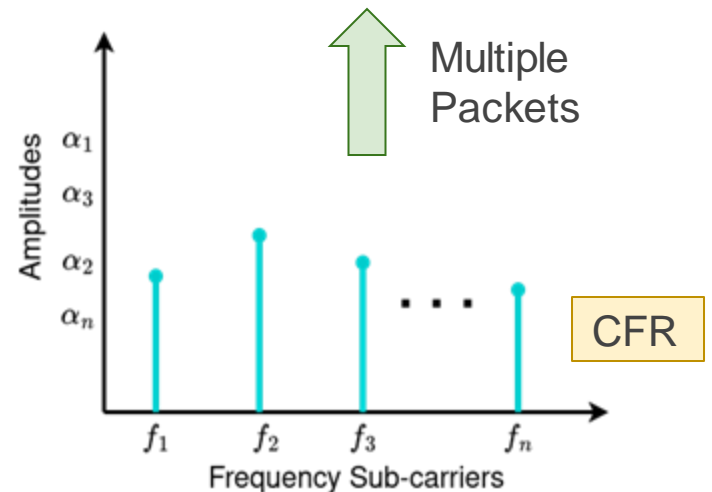
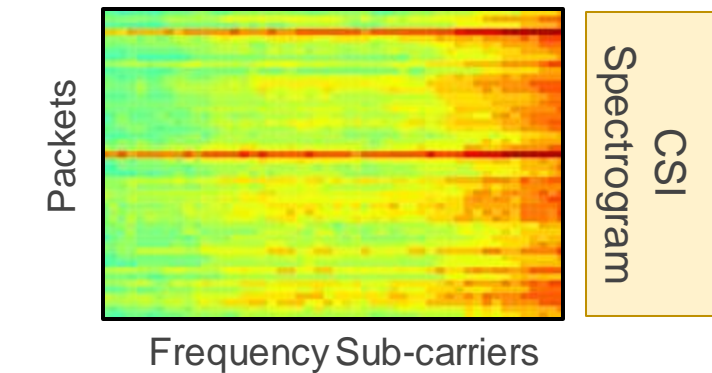


Sensing and Networked Systems Engineering
@IIT Madras

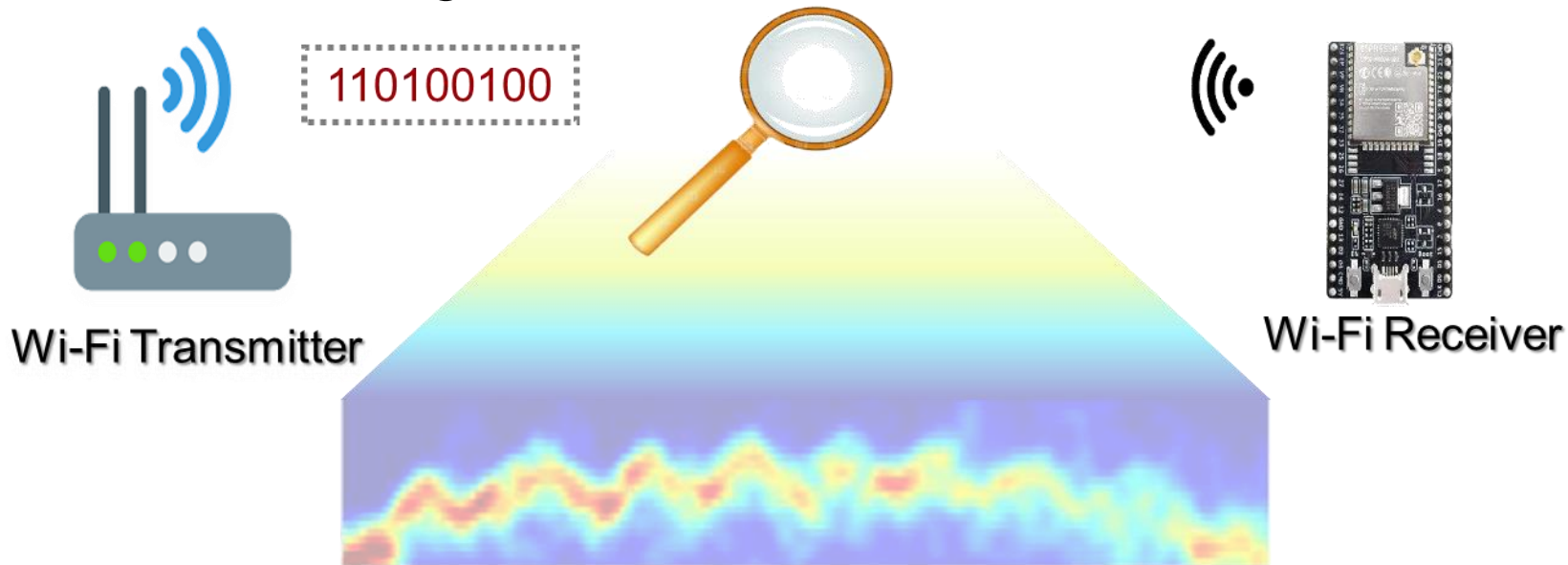
Wireless Sensing 101



FFT



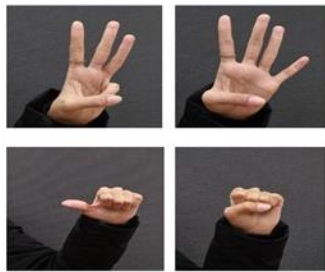
CSI as a Sensing Primitive



Fall Detection



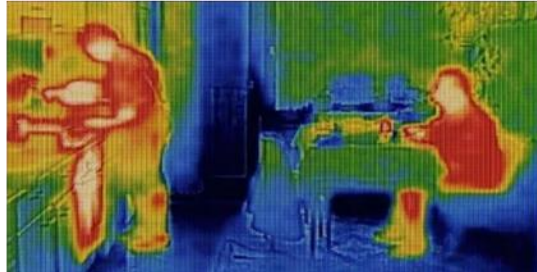
Gestures



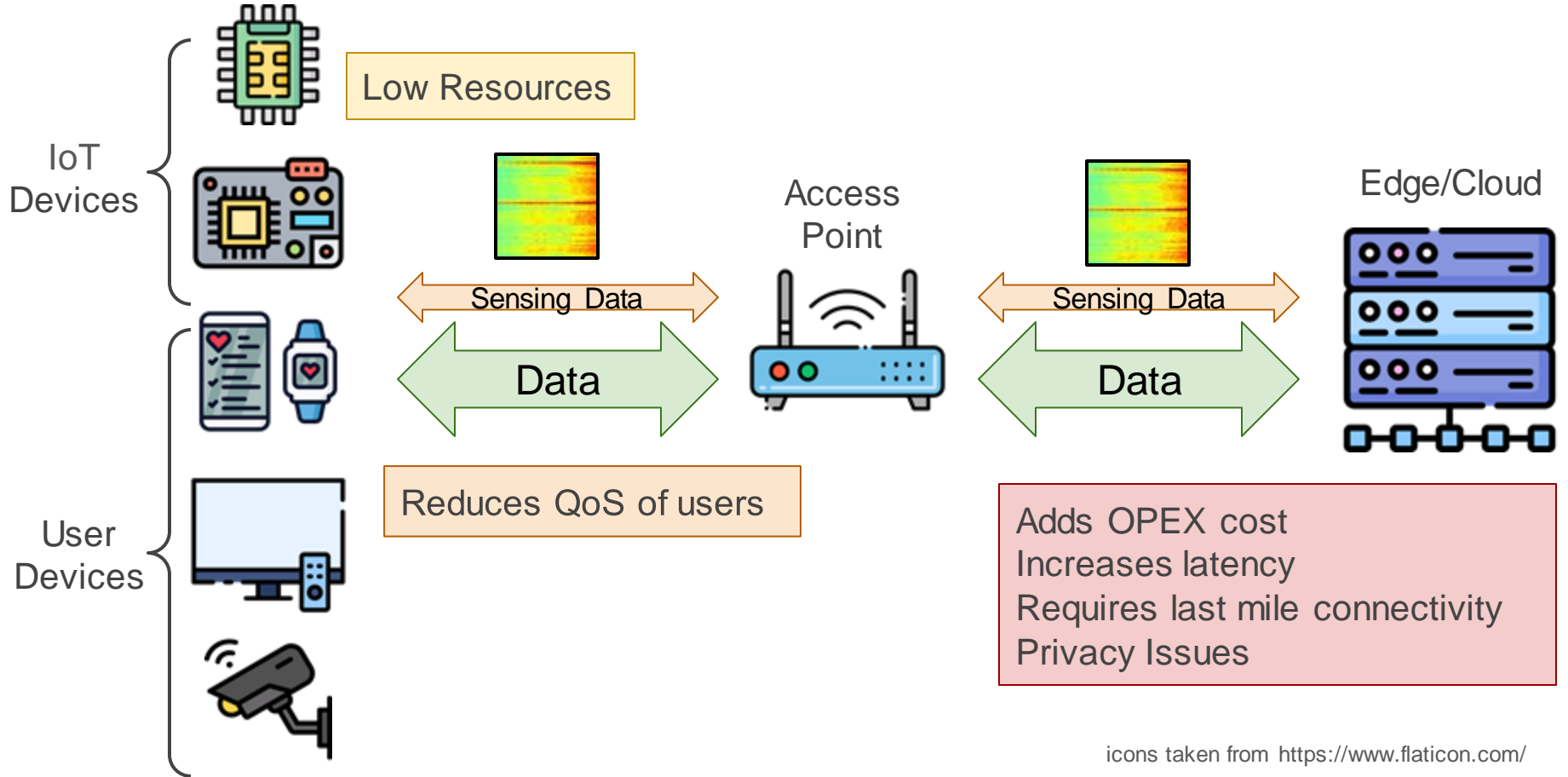
Positioning



Seeing Through Walls



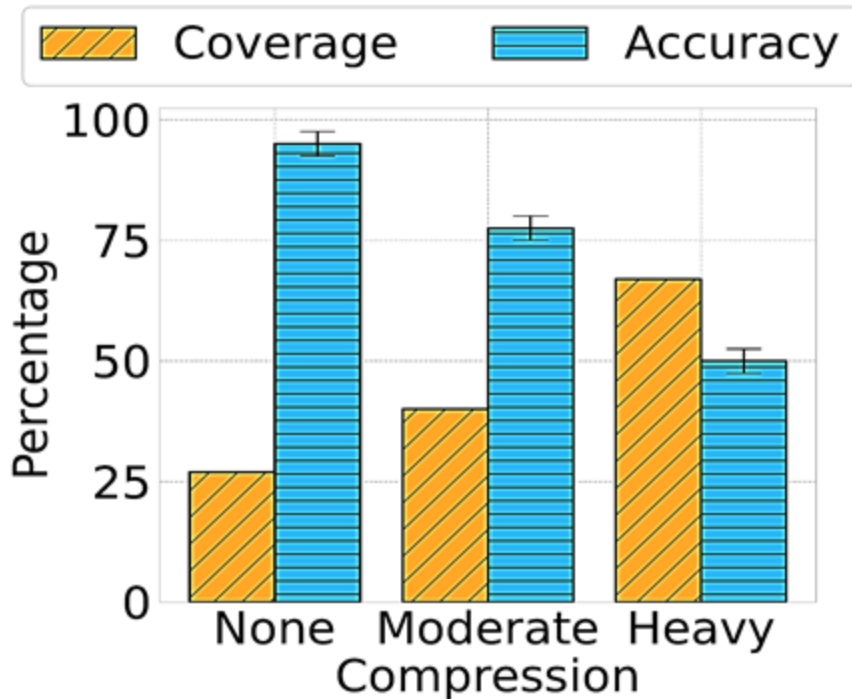
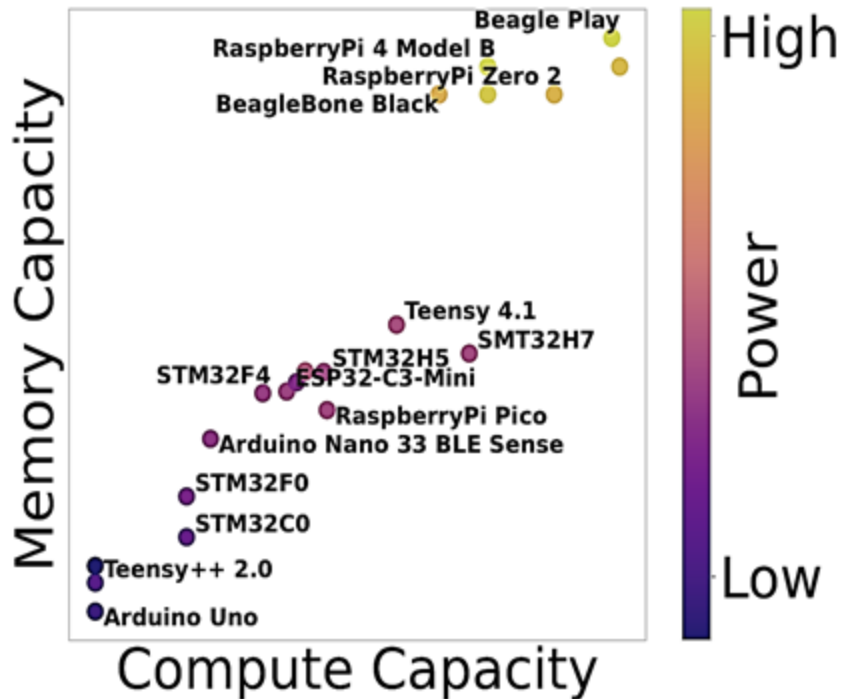
Why Not Edge Based Sensing?



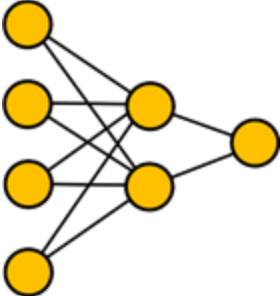
Challenges of Inferencing on IoT Devices

More Resources = More Energy

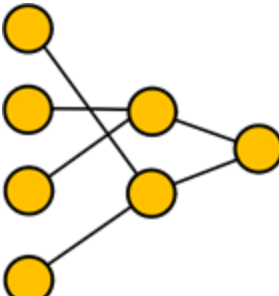
No one-size-fits-all solution



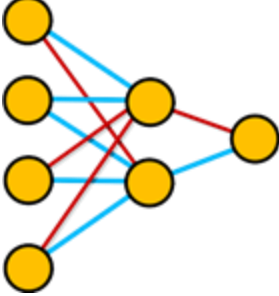
Compressing a Neural Network



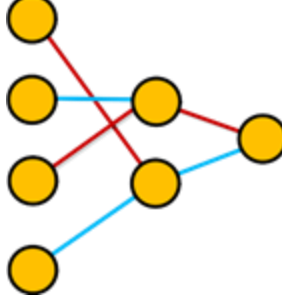
Uncompressed



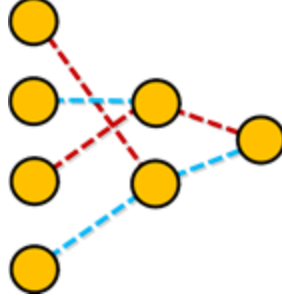
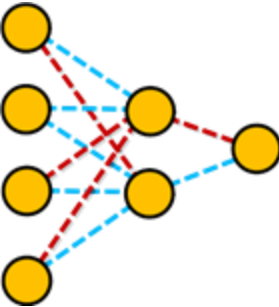
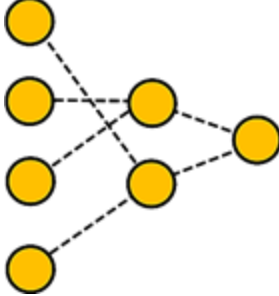
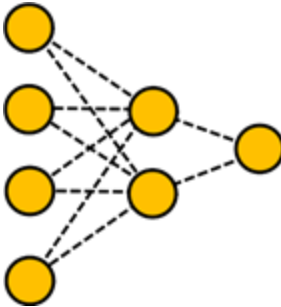
Prune (P)



Cluster (C)



P & C



Quantize (Q)

Related Work and Research Gaps

Traditional Wi-Fi Sensing

Mainly focuses on improving performance and finding new and innovative applications. Less interest in actual system implementation

System Consideration for IoT


Some recent work do look into on-device Wi-Fi sensing on microcontrollers (like ESP-32) from a quantization perspective. Works like EfficientFi look into edge-based deployment.


TinyML Related Work

Has developed techniques like quantization, pruning, etc. Tools like TensorFlow Lite and Micro. Does not specifically focus on wireless sensing

Design a framework that provides a best-effort compressed neural network for a Wi-Fi sensing application such that the user can tune the trade off between performance and cost


WISDOM: Inputs

w_{acc}  A_{min}
Accuracy

w_{inf}  I_{min}
Inference Rate

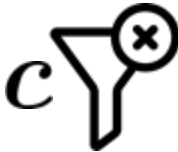
w_{eng}  E_{max}
Energy per Inference

w_{flh}  F_{max}
Flash Consumed

w_{ram}  R_{max}
RAM Required


Weights



c 
Filters

WISDOM: Outputs



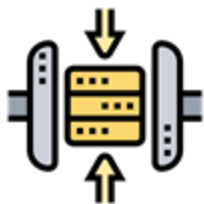
Architecture
Type

$$\mathbf{T} = \{FCN, CNN, LSTM\}$$



Number of
Parameters

$$\mathbf{N} = \{250, 1.5K, 3K, 6K, \dots, 180K\}$$



Compression
Techniques

$$\mathbf{O} = \{none, prune, cluster, qat, \dots, pcptq\}$$



Neural
Network
Configuration

$$\mathbf{I} = \{[t, n, o] | t \in \mathbf{T}, n \in \mathbf{N}, o \in \mathbf{O}\}$$

Total of 324 models

WISDOM: Utility Function

Utility



$$U = P - C$$

Performance



Cost



$$U_{w,c} : \mathbf{I} \rightarrow \mathbb{R}$$

Accuracy



Inference Rate



Energy per Inference RAM Required Flash Consumed



$$C = w_{eng}\mathcal{E} + w_{ram}\mathcal{R} + w_{flh}\mathcal{F} \text{ where}$$

$$\mathcal{E} \leq E_{max}$$

$$\mathcal{R} \leq R_{max}$$

$$\mathcal{F} \leq F_{max}$$

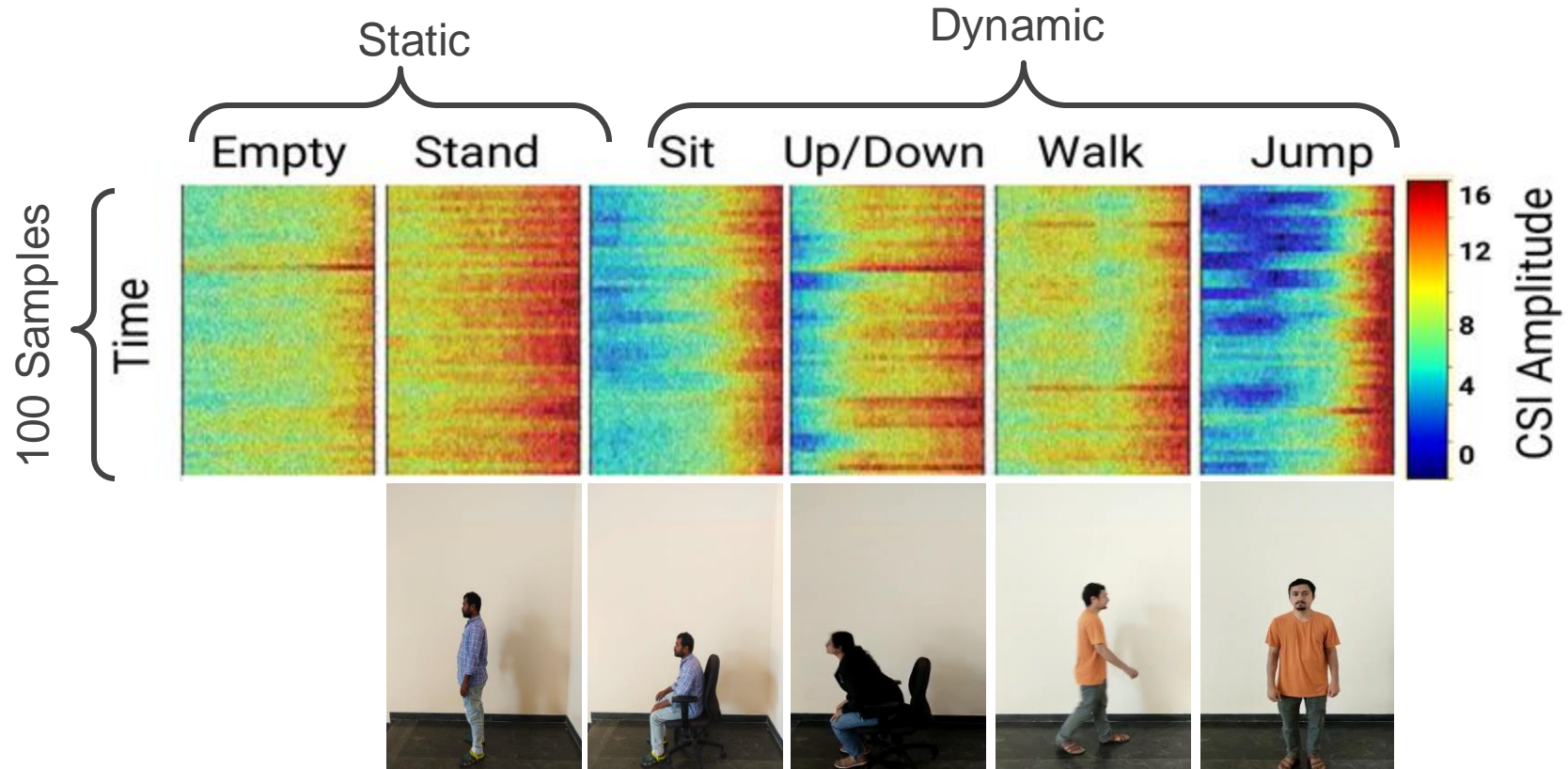
$$P = w_{acc}A + w_{inf}I \text{ where}$$

$$A \geq A_{min}$$

$$I \geq I_{min}$$

All metrics are normalized between 0 and 1 for a fair comparison

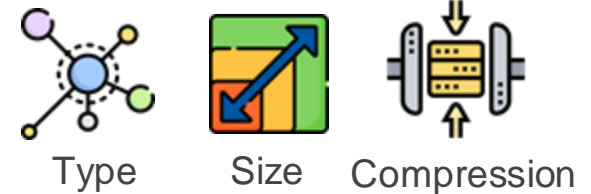
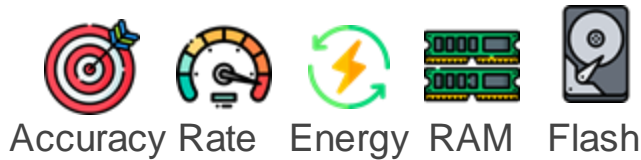
Application Use Case: Human Activity Recognition



WISDOM: Optimization Problem

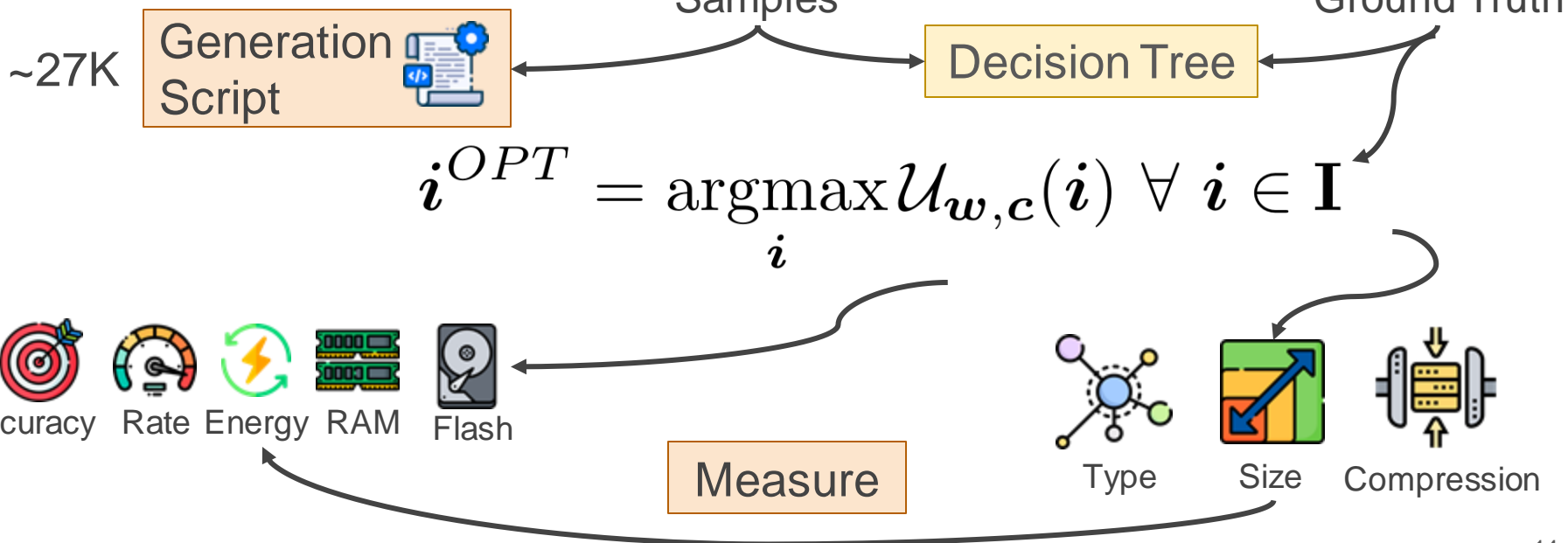
$$\mathcal{U}_{w,c}(\text{Wisdom}(w, c)) \approx \mathcal{U}_{w,c}(i^{OPT})$$

$$i^{OPT} = \underset{i}{\operatorname{argmax}} \mathcal{U}_{w,c}(i) \quad \forall i \in \mathbf{I}$$

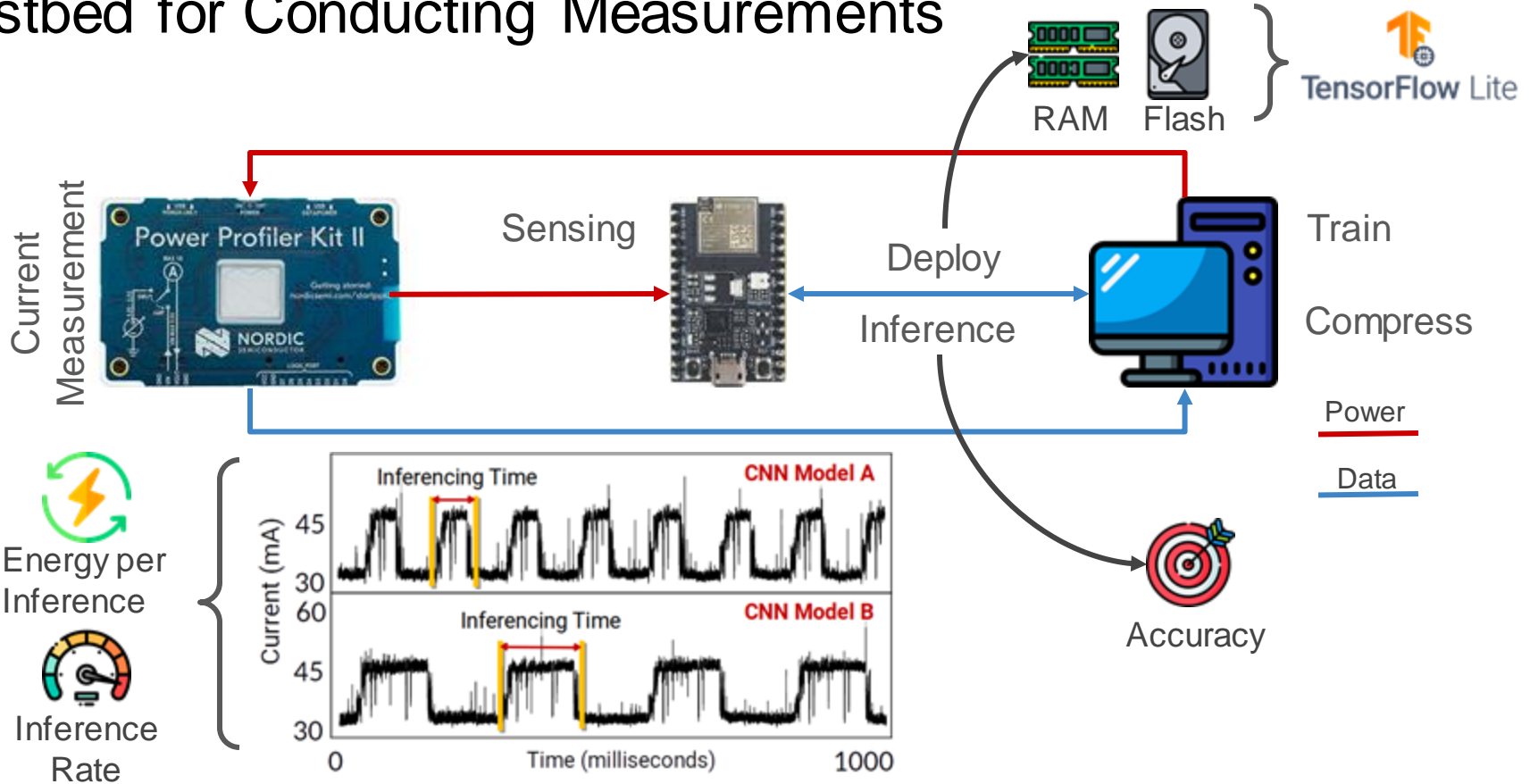


WISDOM: Training

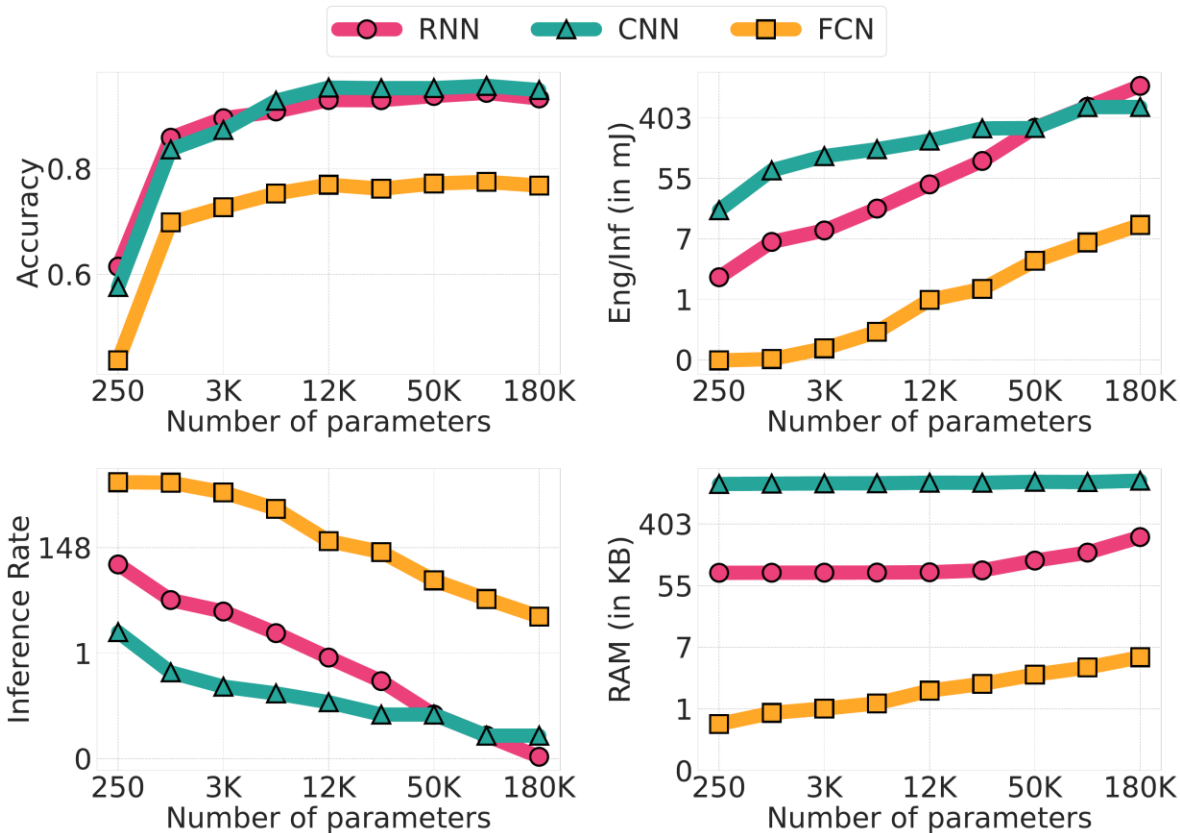
$$\mathcal{U}_{w,c}(\underbrace{\text{Wisdom}(w, c)}_{\text{Samples}}) \approx \mathcal{U}_{w,c}(\underbrace{i^{OPT}}_{\text{Ground Truth}})$$



Testbed for Conducting Measurements

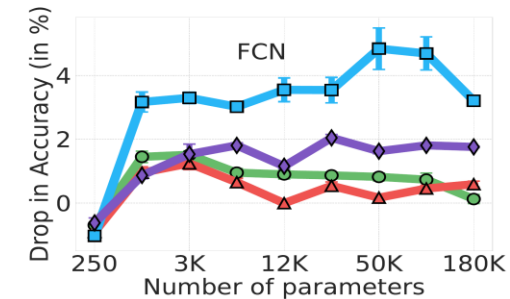
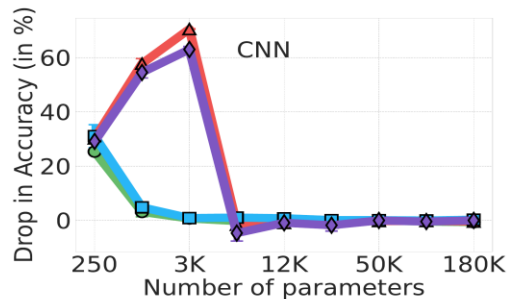
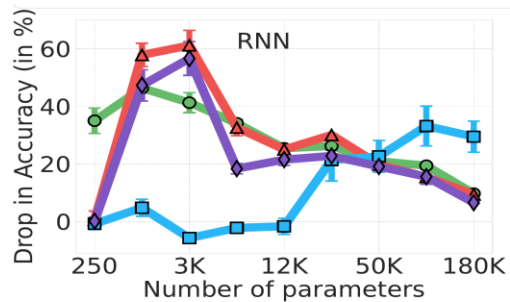
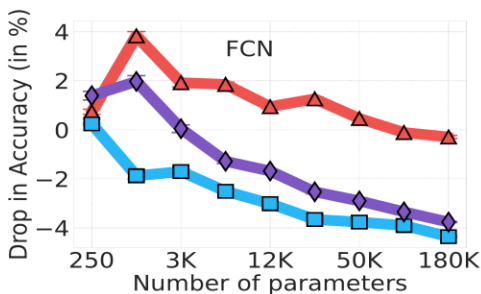
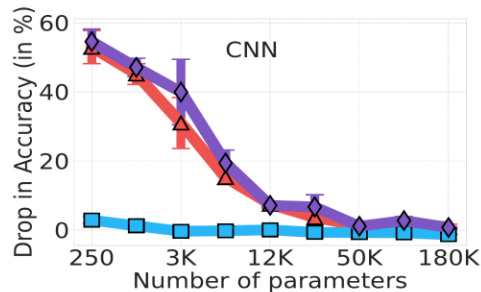
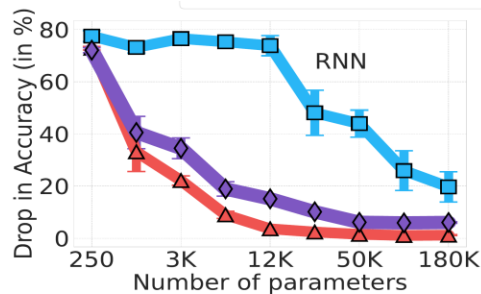


Key Insights (1)

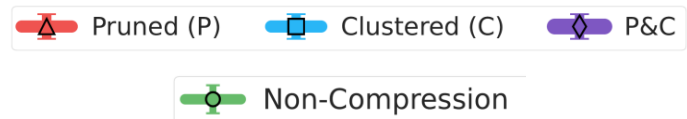


CNNs/RNNs are more accurate but also consume more resources compared to FCNs

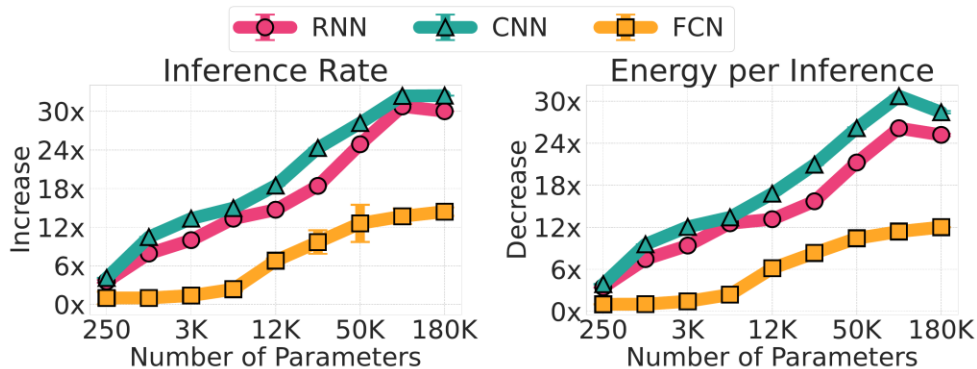
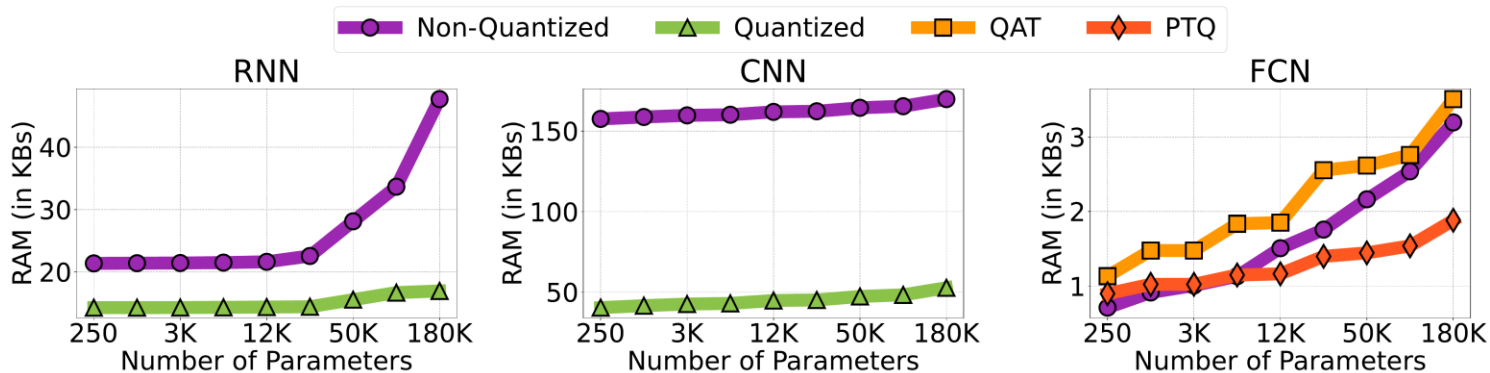
Key Insights (2)



RNNs are more adversely affected by compression compared to CNNs/FCNs

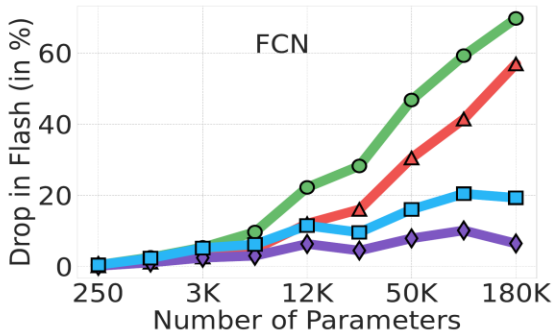
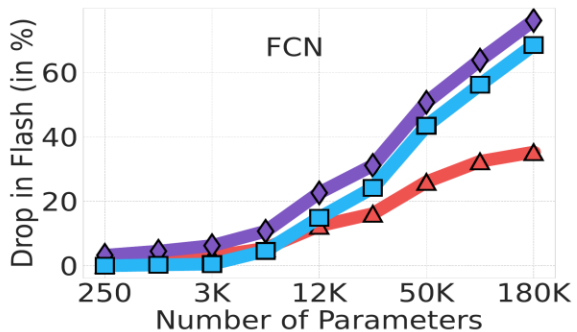
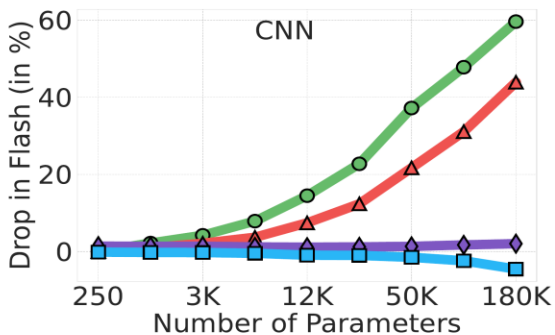
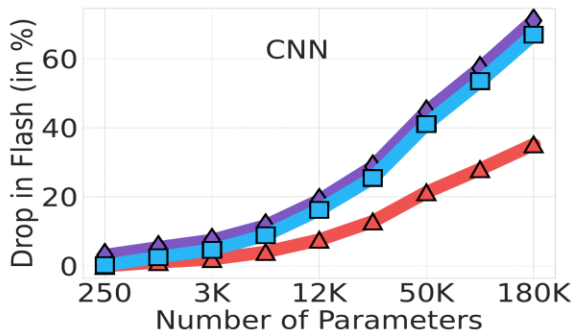
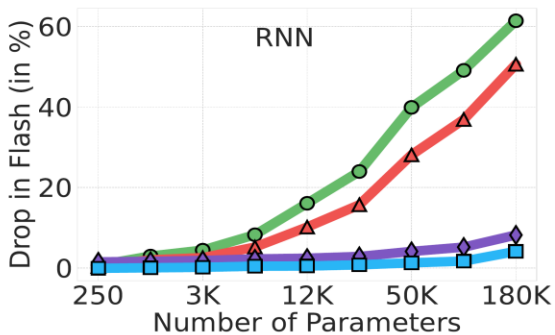
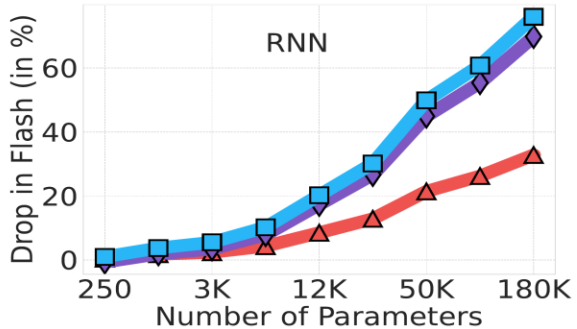


Key Insights (3)



Quantization reduces the RAM and energy requirement, while increasing inferencing rate

Clustering provides significant reduction in flash, but quantization along with clustering is not reasonably effective



Key Insights (4)

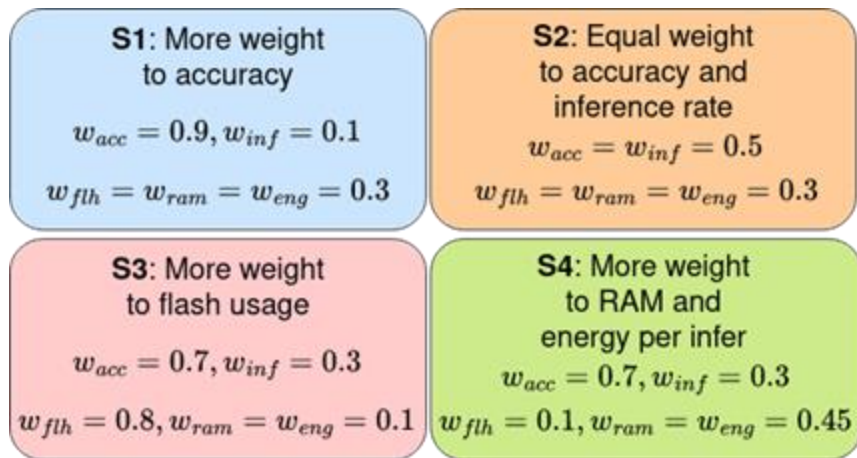
Key Insights (5)

Compressing a model with higher number of parameters yields a more accurate model than an uncompressed model with lesser number of parameters

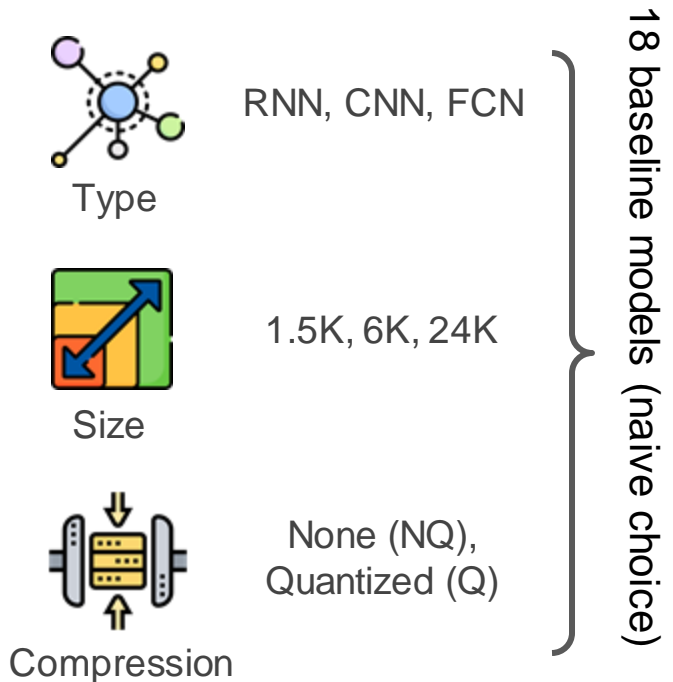
(while having a similar footprint in terms of energy and memory)

Baseline Models and Scenarios Used for Testing

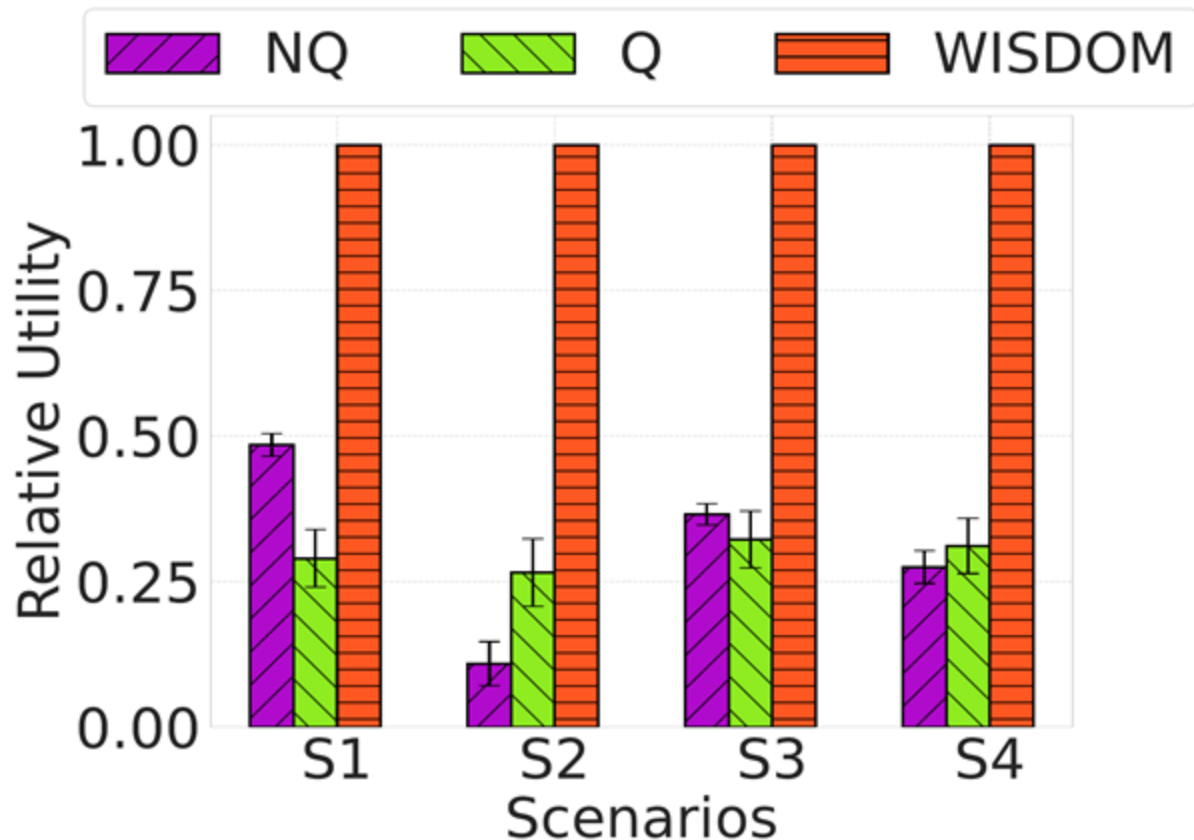
There are additional 126 different test cases with different combination of weights



Scenarios



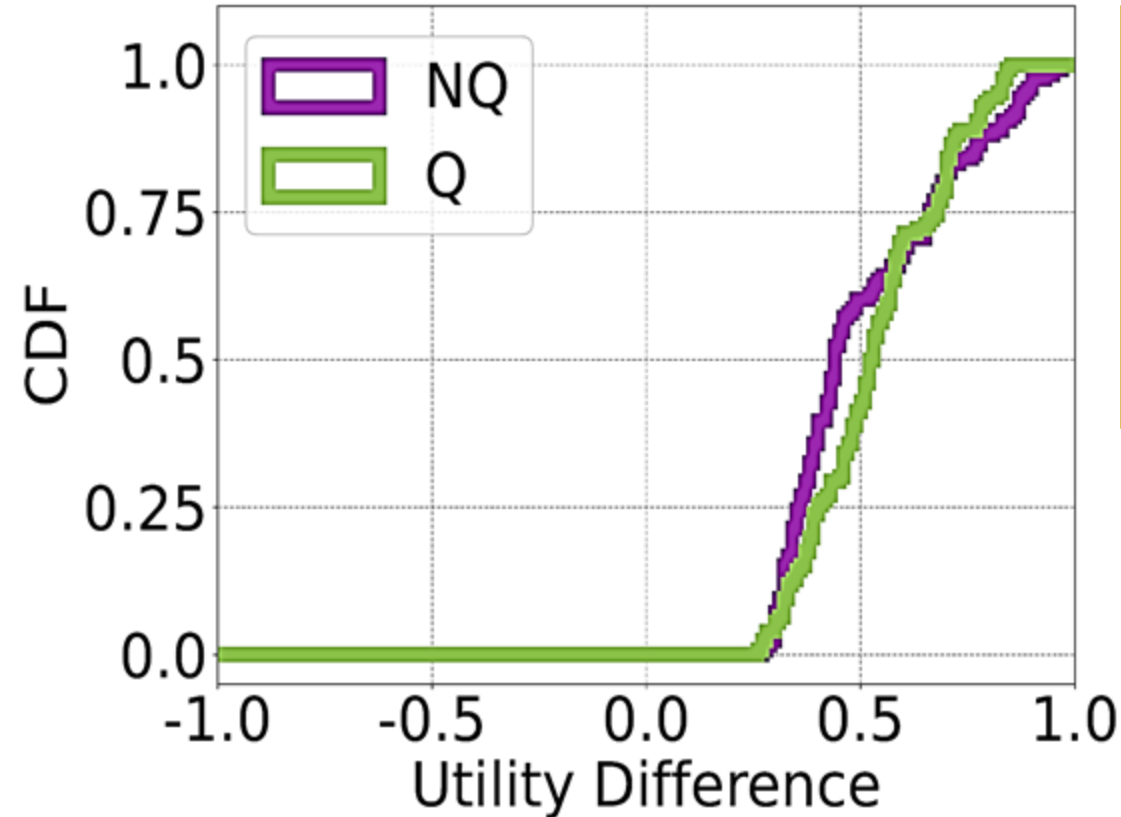
Results: Models Chosen by WISDOM have Higher Utility



Relative utility of NQ and Q models are lower than WISDOM chosen models.

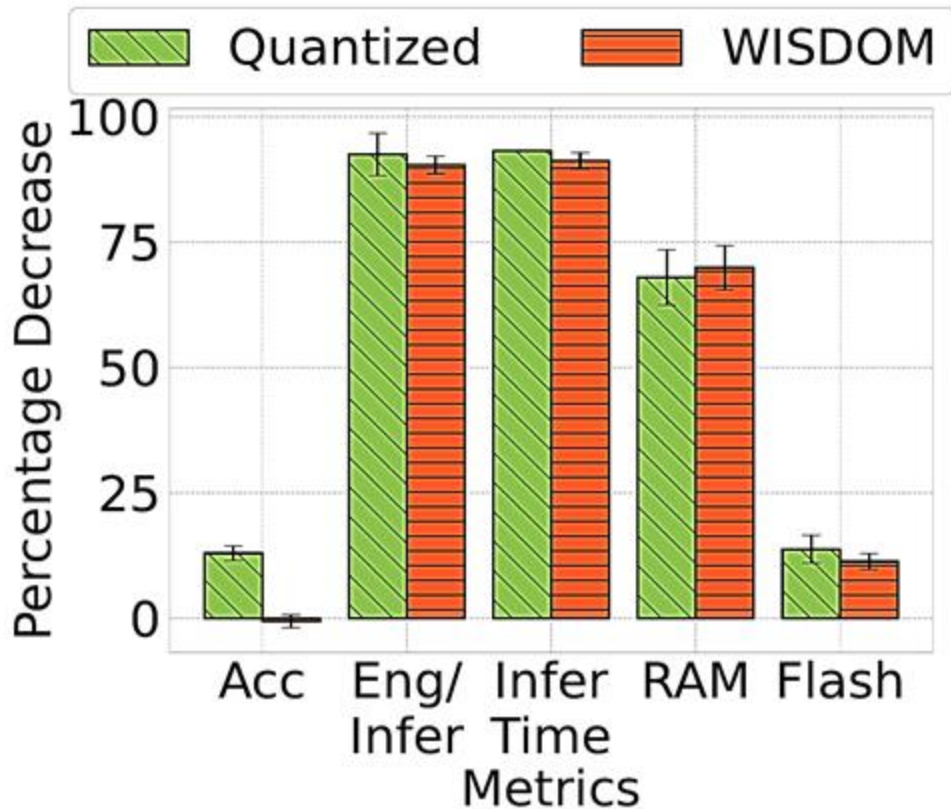
Relative utility is w.r.t to the optimal model i.e., $\frac{U(i)}{U(i^{OPT})}$

Results: Models Chosen by WISDOM have Higher Utility



The CDF of utility difference between Q or NQ model and WISDOM chosen model for all 126 test cases is always positive and starts increasing after 0.5

Results: Models Chosen by WISDOM Uses Less Resources While Maintaining High Accuracy



WISDOM chosen models show a percentage decrease similar to Q models for resource consumption, but still maintains higher accuracy of ~15% compared to Q models. The percentage decrease is w.r.t NQ models

WISDOM chosen model outperforms the best quantized model 83% of time, and the best non-compressed model 99% of time

Thank You, Questions?



Manoj Lenka, lenka98.github.io and **Ayon Chakraborty**, cse.iitm.ac.in/~ayon

Contact: cs22s008@cse.iitm.ac.in, ayon@cse.iitm.ac.in

Artifacts available at: <https://cse.iitm.ac.in/~sense/wisdom/>