

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

MCUNetV3: On-Device Training Under 256KB Memory

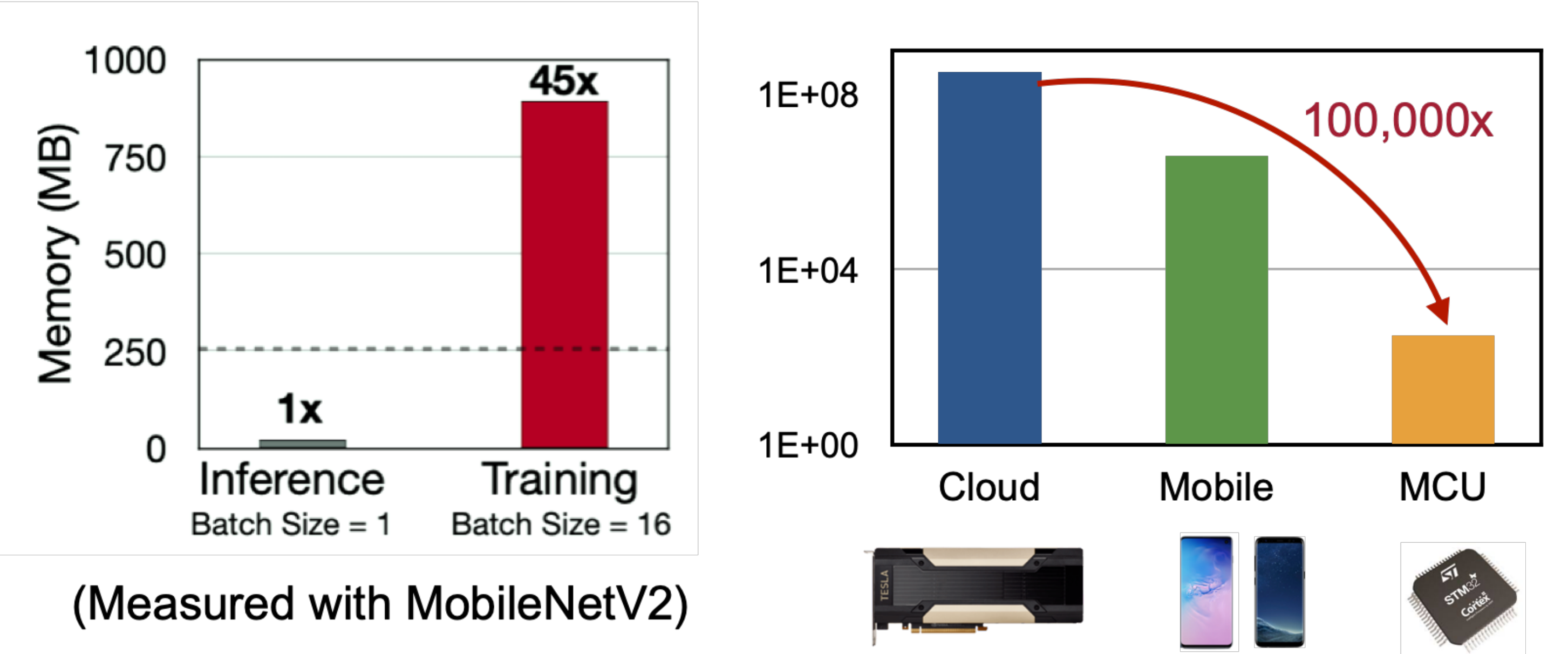
Ji Lin*, Ligeng Zhu*, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, Song Han
MIT, MIT-IBM Watson AI Lab

On-Device Training Enables Customization and Continual Learning



- Security:** Data **never leaves devices**, thus promises security and regularization.
- Customization:** Models **continually adapt** to new data from the sensors.

Challenge: Training is Expensive for Edge



Quantization-Aware Scaling (QAS)

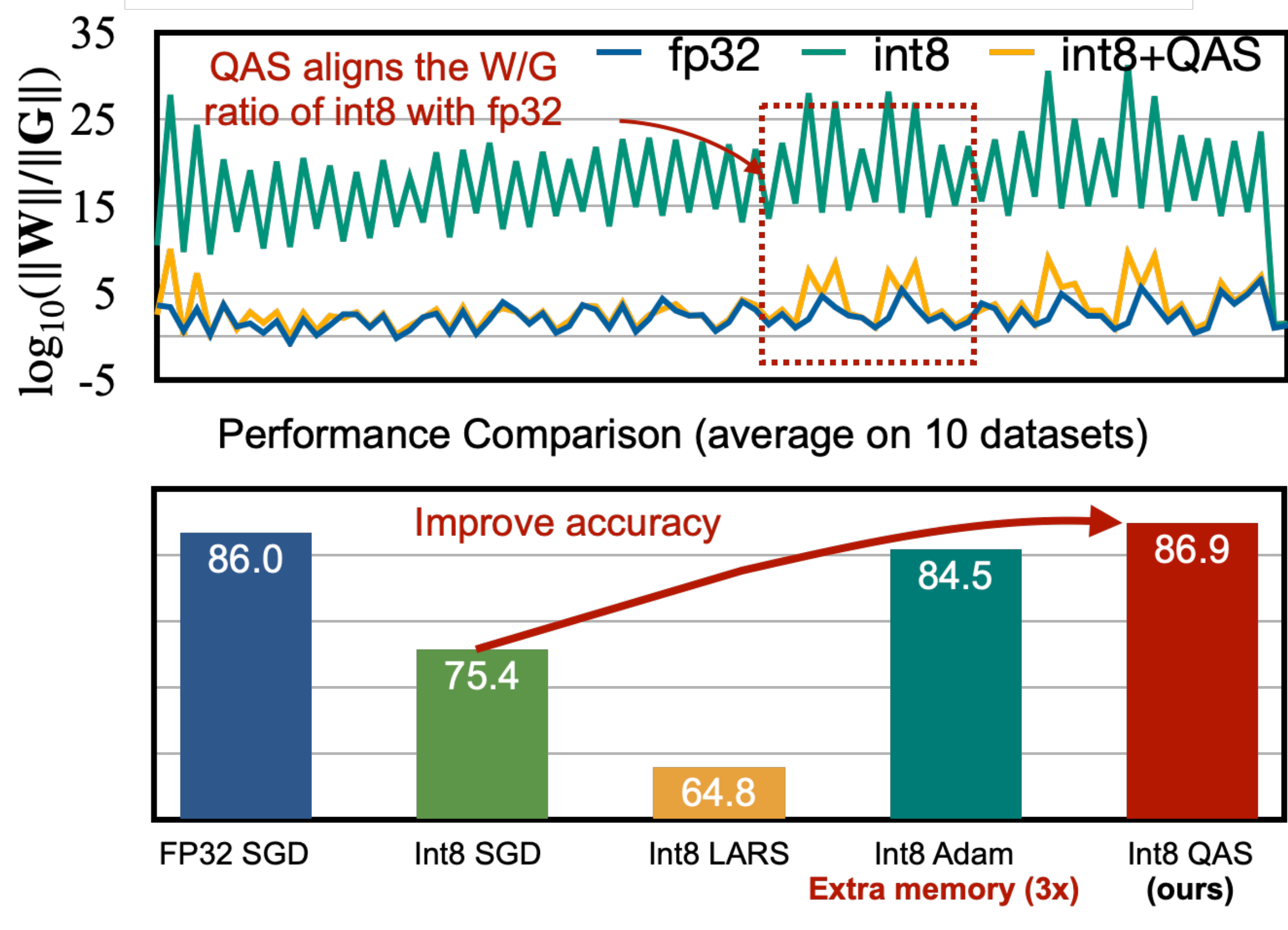
- Re-scale gradients** to help convergence with **quantized training**

Weight quantization

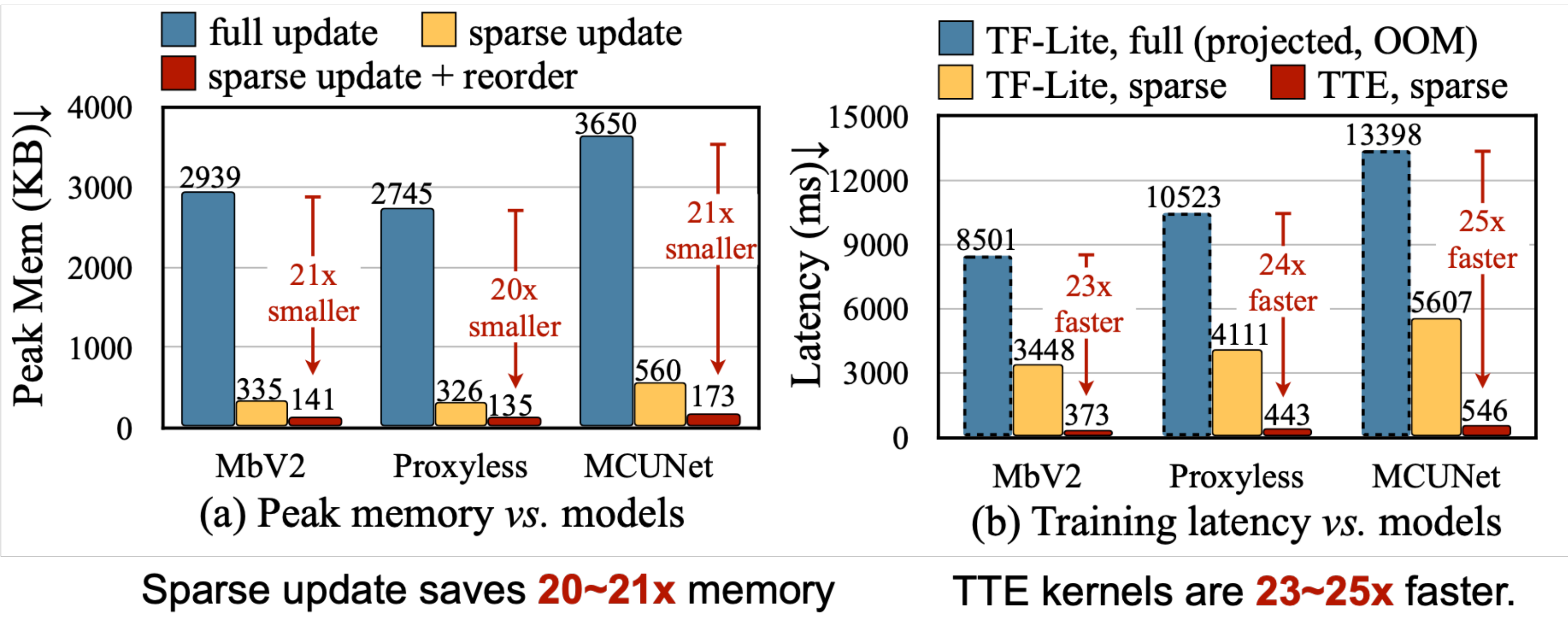
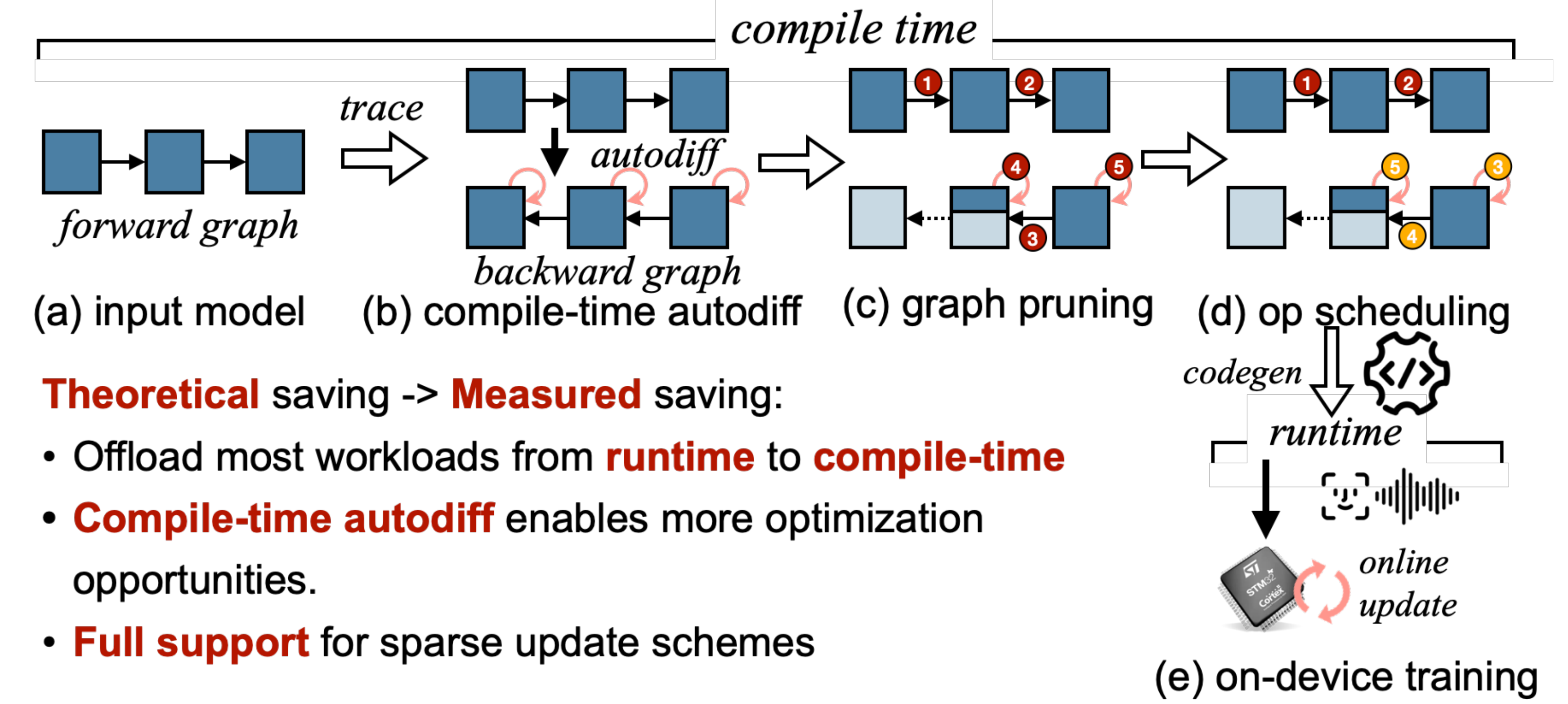
$$W = s_W \cdot (W/s_W) \approx s_W \cdot \bar{W}, \quad G_{\bar{W}} \approx s_W \cdot G_W,$$
 Weight and gradient ratios are off by s_W^{-2}

$$\|\bar{W}\|/\|G_{\bar{W}}\| \approx \|W/s_W\|/\|s_W \cdot G_W\| = s_W^{-2} \cdot \|W\|/\|G\|.$$
 Thus, re-scale the gradients

$$\tilde{G}_{\bar{W}} = G_{\bar{W}} \cdot s_W^{-2}, \quad \tilde{G}_{\bar{b}} = G_{\bar{b}} \cdot s_W^{-2} \cdot s_x^{-2} = G_{\bar{b}} \cdot s^{-2}$$

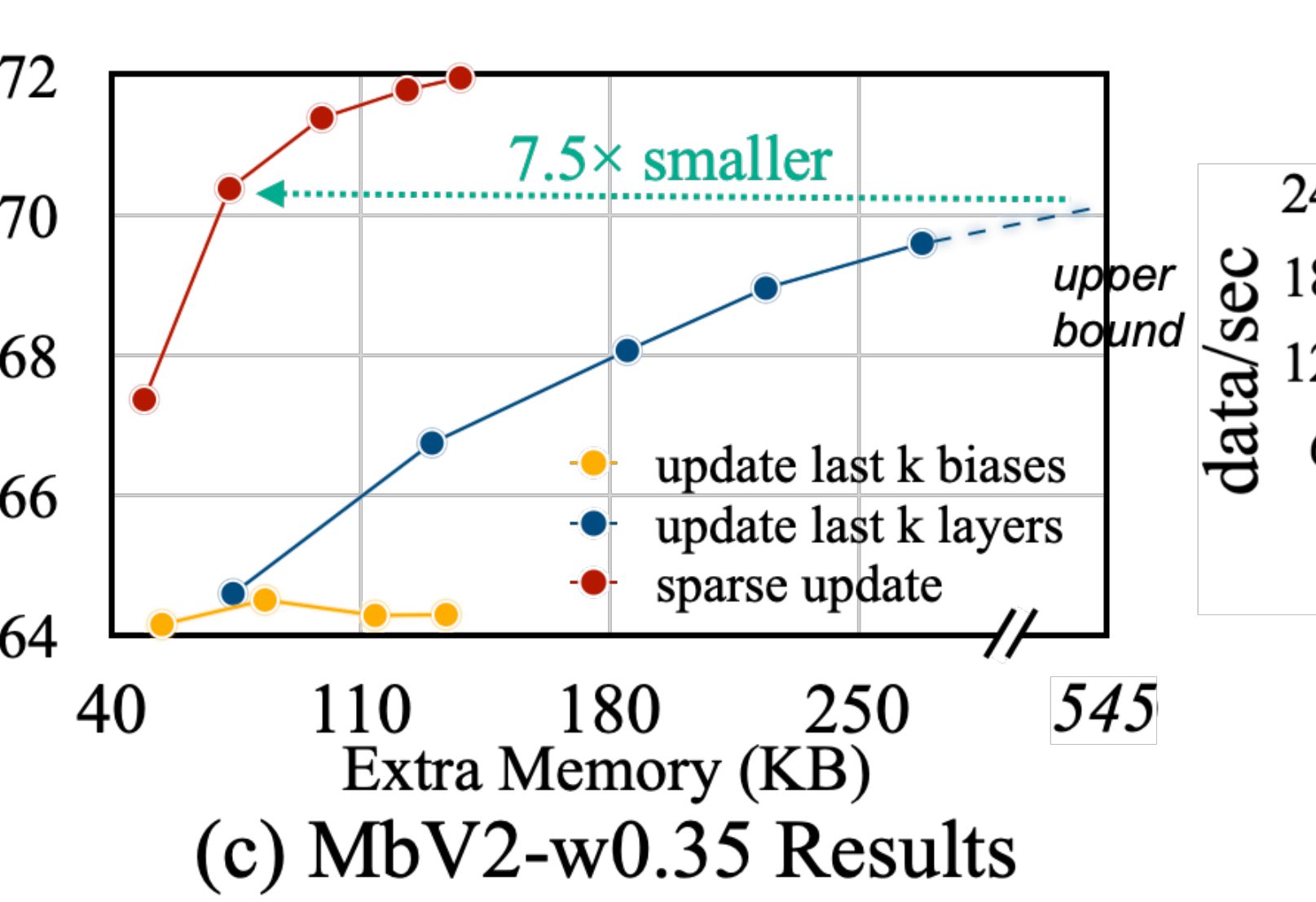
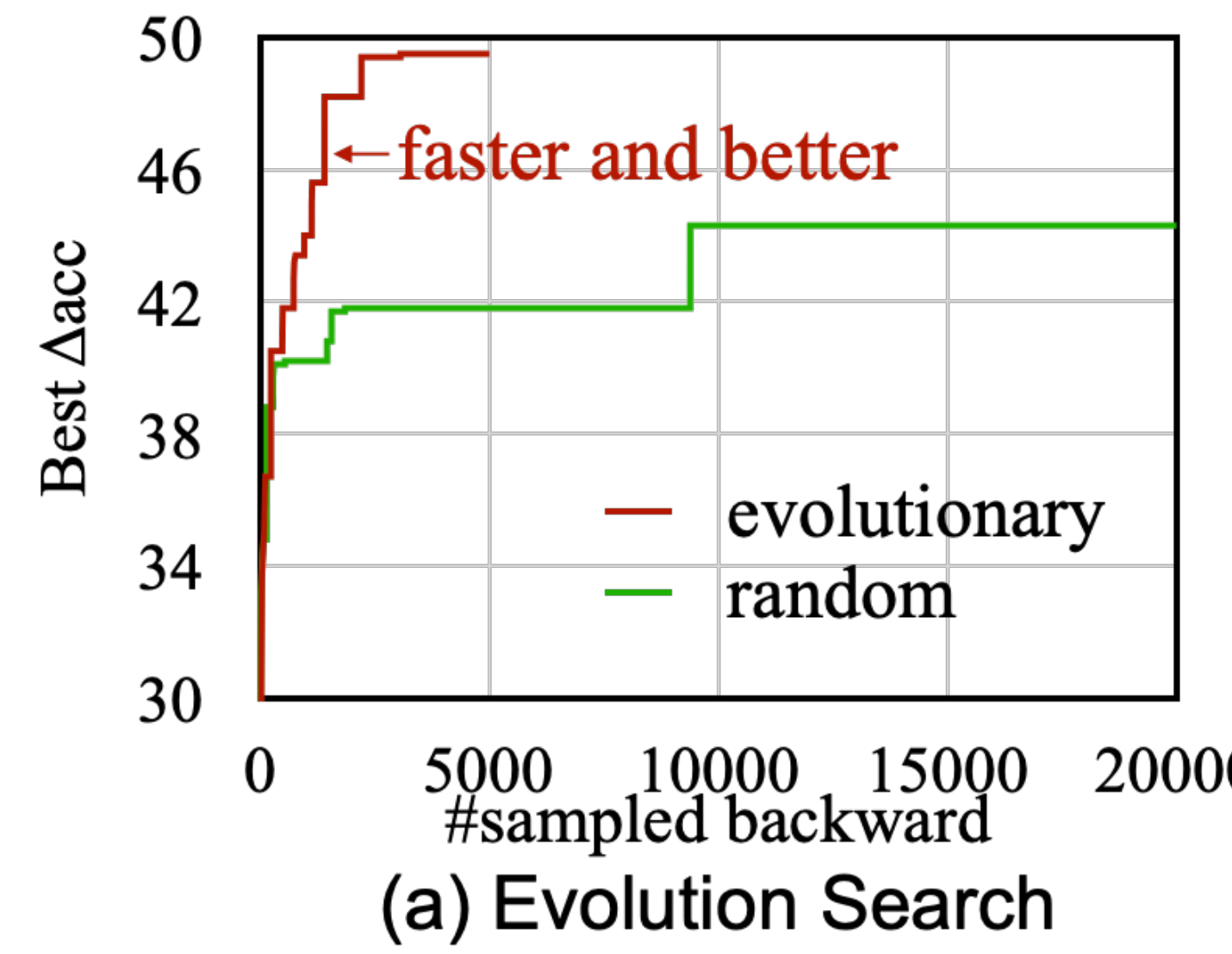
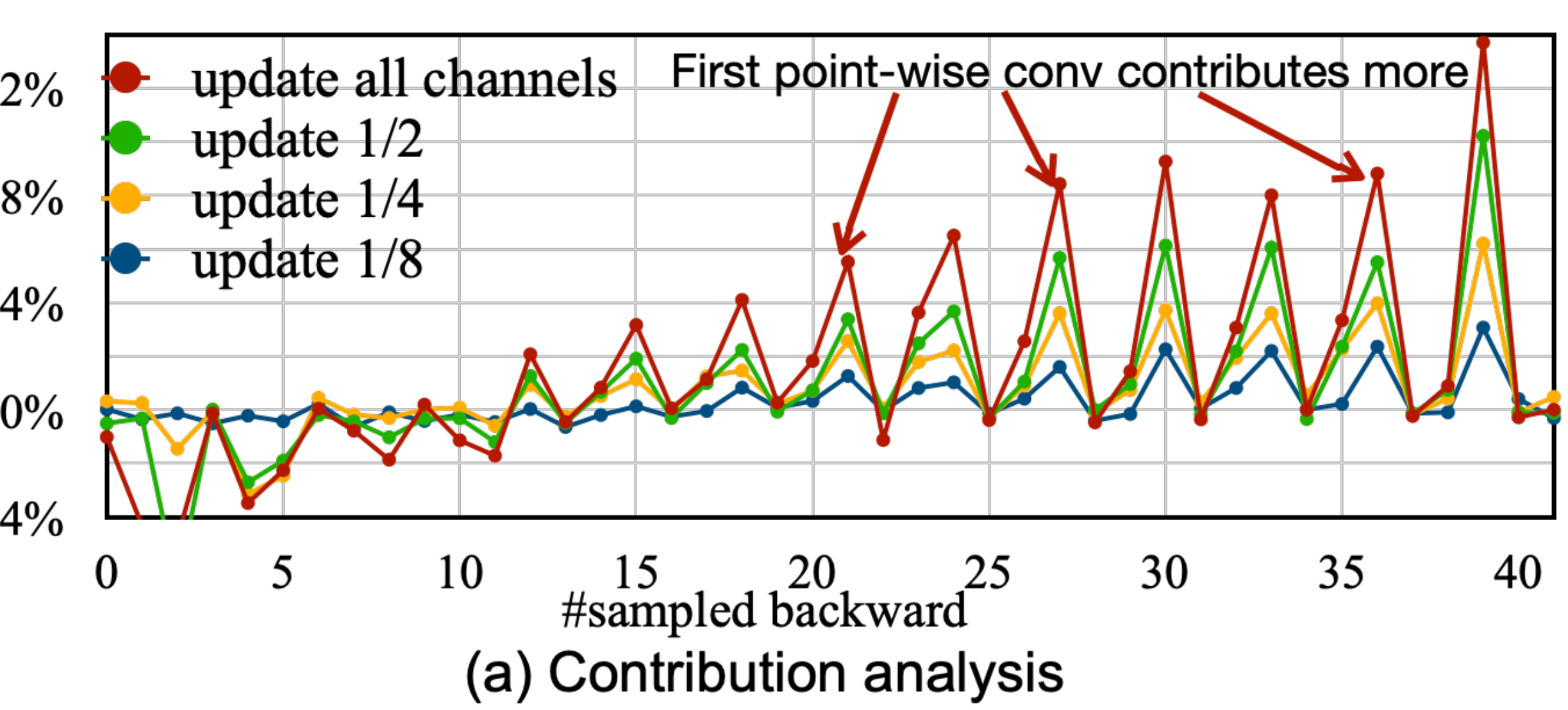
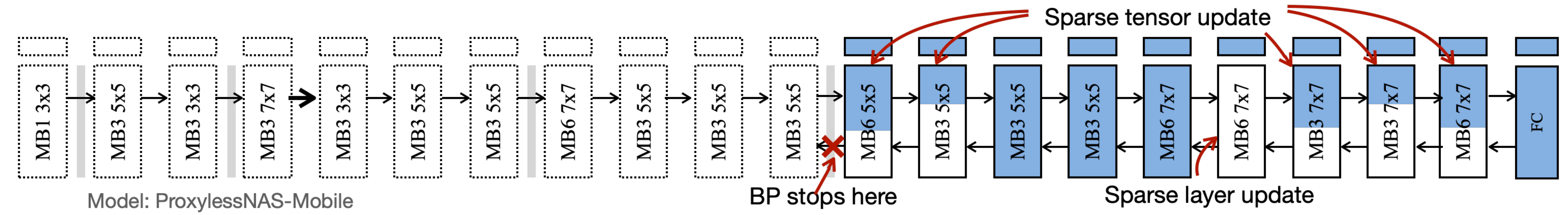


Tiny Training Engine (TTE)



Sparse Layer/Tensor Update

- Conventionally, we update the **full model** for transfer learning
- We find some layers are **more important** than others, so we can **sparingly update**



Algorithm-System Co-Design Results

