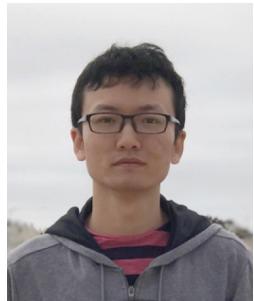
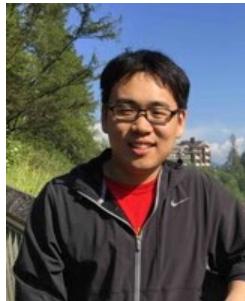


Automatic Neural Network Compression by Sparsity-Quantization Joint Learning: A Constrained Optimization-based Approach



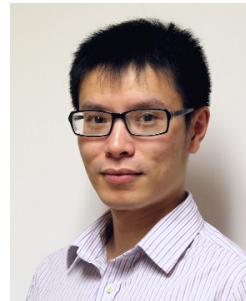
Haichuan Yang¹



Shupeng Gui¹



Yuhao Zhu¹



Ji Liu²



¹University of Rochester



²Kwai Inc.



DNN on resource-constrained devices



Face verification



Speech recognition



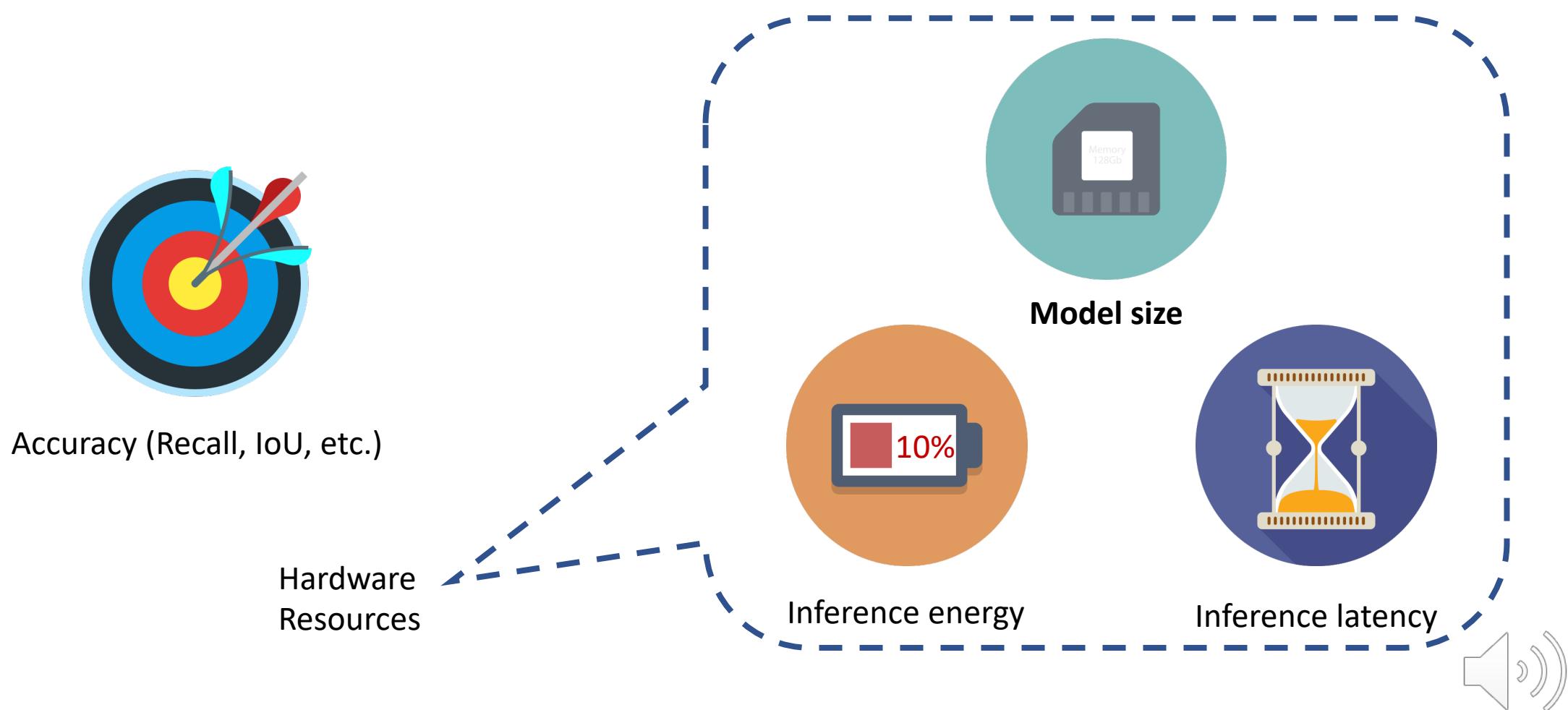
Self driving*



VR/AR



DNN Compression: maximize accuracy within resource constraint



The overall objective

$$\begin{aligned} \min_W \ell(W) & \xrightarrow{\text{Bitwidth of } W^{(i)}} \\ \text{s. t. } \sum_{i=1}^k b(W^{(i)}) \|W^{(i)}\|_0 & \leq B \xrightarrow{\text{Model size budget}} \\ & \text{Jointly consider pruning and quantization} \\ & \text{Each layer can have different bitwidth / sparsity} \end{aligned}$$

Pros:

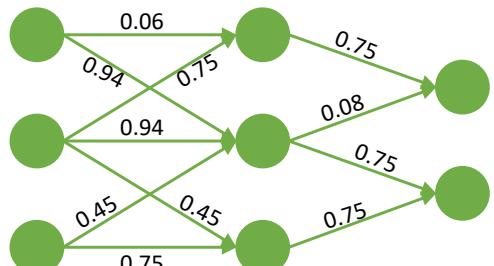
- Directly formulated from the original problem
- No introduced hyper-parameter (*no need to set layer-wise sparsity / bitwidth*)



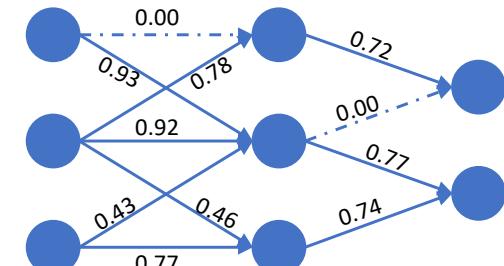
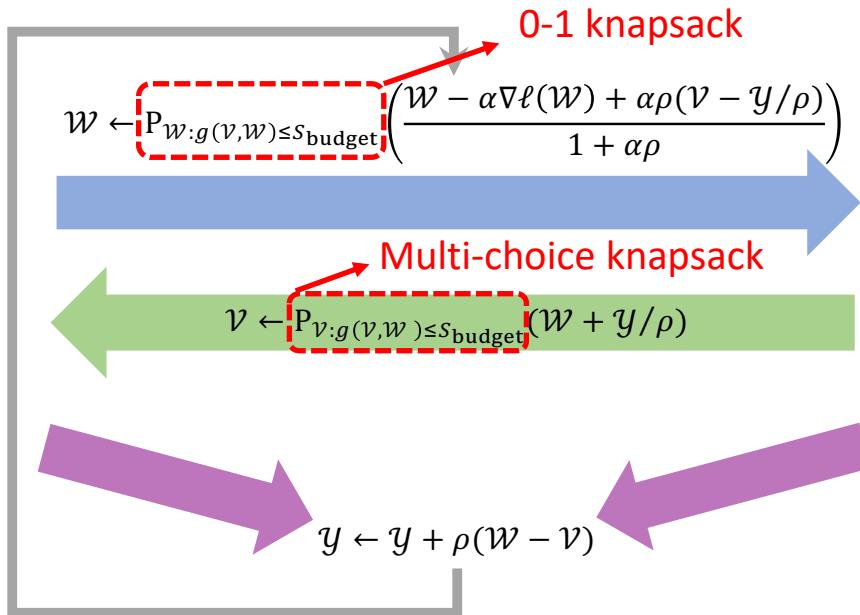
DNN compression framework

$$\min_W \ell(W) \quad \text{s.t.} \quad \sum_{i=1}^k b(W^{(i)}) \|W^{(i)}\|_0 \leq B \quad \xrightarrow{\hspace{1cm}} \quad \begin{aligned} & \min_{W,V} \max_U \ell(W) + \langle U, W - V \rangle + \frac{\rho}{2} \|W - V\|^2 \\ & \text{s.t.} \quad \sum_{i=1}^k b(V^{(i)}) \|W^{(i)}\|_0 \leq B \end{aligned}$$

Alternatively update W , V , and U



$$\mathcal{V}: b(V^{(1)}) = 2, b(V^{(2)}) = 1$$



$$\mathcal{W}: \|W^{(1)}\|_0 = 6, \|W^{(2)}\|_0 = 3$$



Experiment Results

- Table: Comparison across different compression methods on ImageNet.

Model	Method	Automated	Pruning	Quantization	NZ%	Ave. bits	Comp. rate	Acc.-1↓	Acc.-5↓
MobileNet	Uniform Baseline [17]	✗	✓	✗	61%	-	1.6×	2.50%	1.70%
	Uniform Baseline [17]	✗	✓	✓	61%	8	6.6×	4.10%	2.90%
	Deep Compression [11]	✗	✗	✓	-	2	16×	33.28%	25.59%
	HAQ [42]	✓	✗	✓	-	2	16×	13.76%	8.03%
	Ours	✓	✗	✓	-	2	16×	7.10%	4.40%
	Deep Compression [11]	✗	✗	✓	-	3	10.7×	4.97%	3.05%
	HAQ [42]	✓	✗	✓	-	3	10.7×	3.24%	1.69%
	Ours	✓	✗	✓	-	3	10.7×	1.19%	0.76%
	Ours	✓	✓	✓	42%	2.8	26.7 ×	4.41%	2.61%
AlexNet	Constraint-Aware [3]	✓	✓	✗	4.9%	-	20×	2.57%	-
	Deep Compression [11]	✗	✓	✓	11%	5.4	54×	0.00%	-0.03%
	CLIP-Q [41]	✓	✓	✓	8%	3.3	119×	-0.70%	-
	Ours	✓	✓	✓	7.4%	3.7	118×	-1.00%	-1.15%
	Ye et al. [53]	✗	✓	✓	4%	4.1	210 ×	0.10%	-
	Ours	✓	✓	✓	5%	3.1	205×	-0.08%	-0.56%
MnasNet	Fixed-Bitwidth	✓	✓	✓	50%	4	16×	3.14%	1.86%
	Ours	✓	✓	✓	50%	3.7	17.1×	1.66%	0.92%
	Ours	✓	✓	✓	30%	3.0	35.6 ×	5.82%	3.23%
ProxylessNAS-mobile	Fixed-Bitwidth	✓	✓	✓	50%	4	16×	3.17%	1.73%
	Ours	✓	✓	✓	51%	3.8	16.8×	2.13%	1.16%
	Ours	✓	✓	✓	31%	2.9	35.6 ×	5.21%	2.84%



Table: Comparison results on LeNet-5@MNISTs.

Method	Automated	NZ%	Avg. bits	Comp. rate	Acc. \downarrow
Deep Compression [11]	\times	8.3%	5.3	70 \times	0.1%
BC-GNJ [33]	\times	0.9%	5	573 \times	0.1%
BC-GHS [33]	\times	0.6%	5	771 \times	0.1%
Ye et al. [53]	\times	0.6%	2.8	1,910 \times	0.1%
Ours	\checkmark	1.0%	1.46	2,120\times	0.0%

Table: Comparison results on CIFAR-10.

Model	Method	Pruning	Quantization	NZ%	Ave. bits	Comp. rate	Acc. \downarrow
ResNet-20	ReLeQ [50]	\times	\checkmark	-	2.8	11.4 \times	0.12%
	Ours	\times	\checkmark	-	2	16 \times	0.00%
	Ours	\checkmark	\checkmark	46%	1.9	35.4\times	0.14%
ResNet-50	AMC [14]	\checkmark	\times	60%	-	1.7 \times	-0.11%
	Ours	\checkmark	\times	50%	-	2 \times	-1.51%
	Ours	\checkmark	\checkmark	4.2%	1.7	462 \times	-1.25%
	Ours	\checkmark	\checkmark	3.1%	1.9	565 \times	-0.90%
	Ours	\checkmark	\checkmark	2.2%	1.8	836\times	0.00%



Summary

- We automatically allocate the layer-wise sparsity and bit-width and jointly prune and quantize a DNN with the given model size budget.

