

Lecture 12: Model Serving

CSE599W: Spring 2018

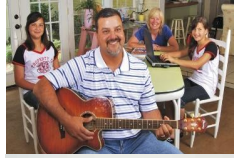
Deep Learning Applications



"That drink will get you to 2800 calories for today"



"I last saw your keys in the store room"

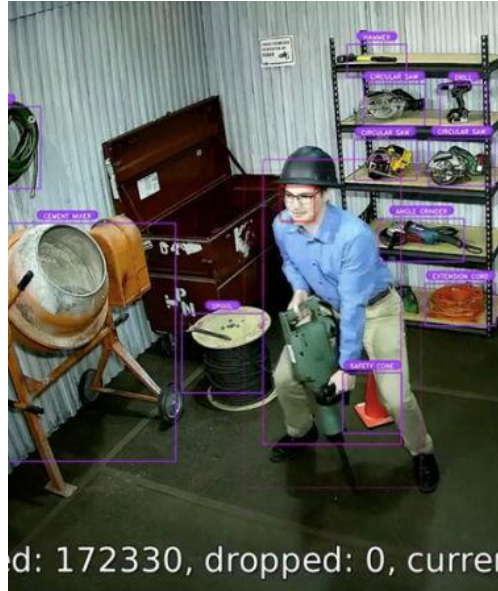


"Remind Tom of the party"

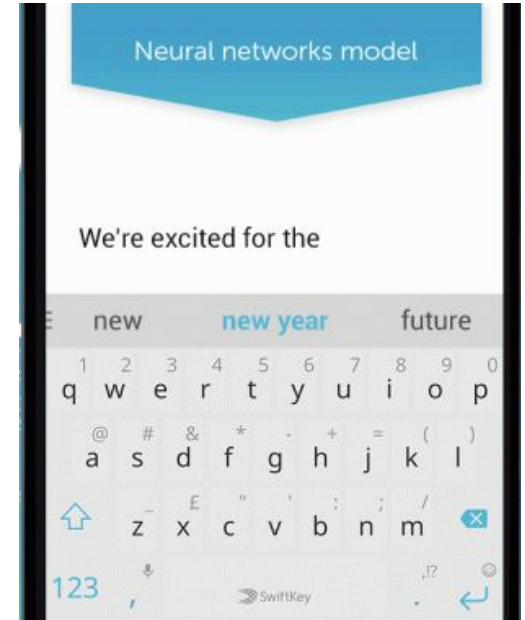


"You're on page 263 of this book"

Intelligent assistant



Surveillance / Remote assistance

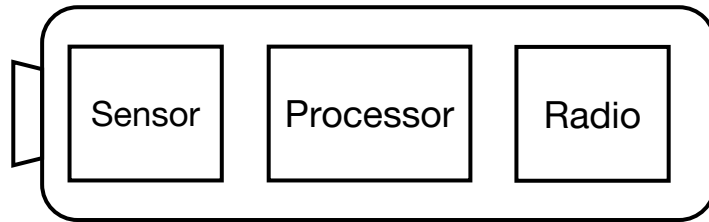


Input keyboard

Model Serving Constraints

- Latency constraint
 - Batch size cannot be as large as possible when executing in the cloud
 - Can only run lightweight model in the device
- Resource constraint
 - Battery limit for the device
 - Memory limit for the device
 - Cost limit for using cloud
- Accuracy constraint
 - Some loss is acceptable by using approximate models
 - Multi-level QoS

Runtime Environment



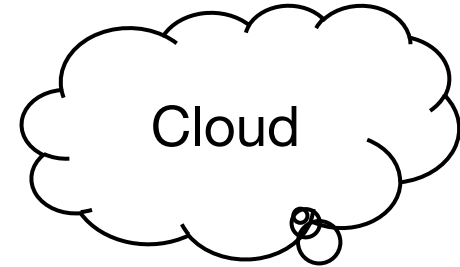
Camera
Touch screen
Microphone

CPU
GPU
FPGA
ASIC

Wifi
4G / LTE
Bluetooth



streaming

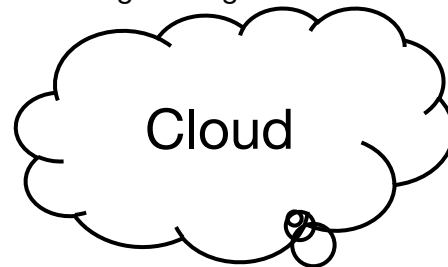
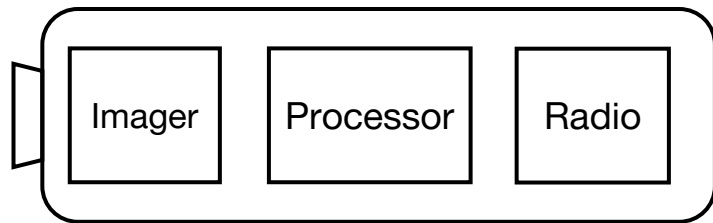


CPU server
GPU server
TPU / FPGA server

Resource usage for a continuous vision app

Omnivision OV2740 90mW
Tegra K1 GPU 290GOPS@10W = 34pJ/OP
Qualcomm SD810 LTE >800mW
Atheros 802.11 a/g 15Mbps@700mW = 47nJ/b

Amazon EC2
CPU c4.large 2x400GFLOPS \$0.1/h
GPU g2.2xlarge 2.3TFLOPS \$0.65/h



Workload

Deep learning 300GFLOPS @ 30GFLOPs/frame, 10fps

Budget

Device power
30% of 10Wh for 10h = 300mW

Cloud cost
\$10 person/year

Compute power

9GFLOPS

3.5GFLOPS (GPU) / 8GFLOPS (CPU)

Huge gap between workload and budget

Outline

- Model compression
- Serving System

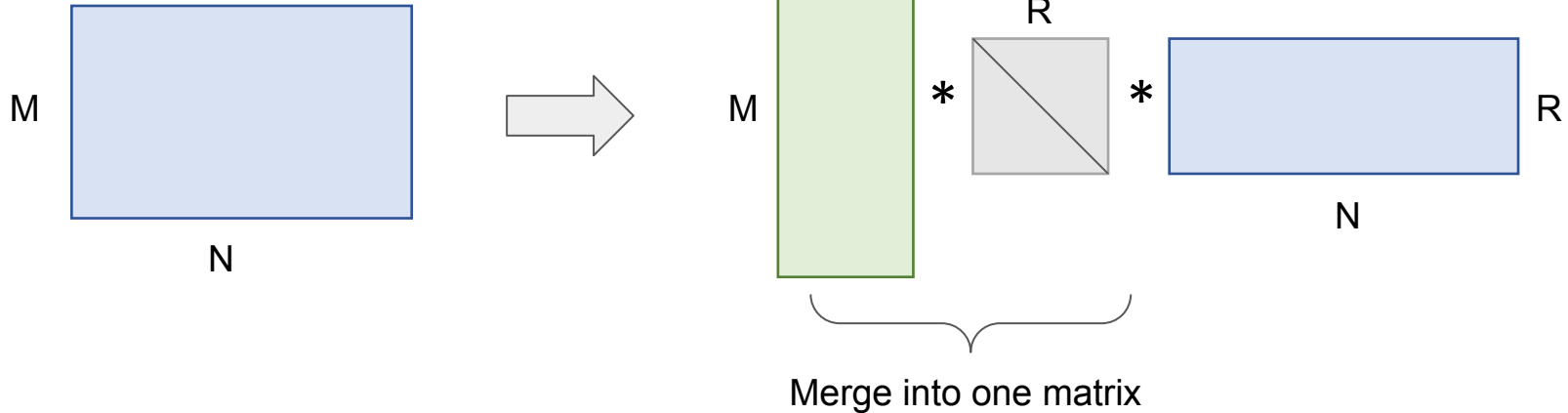
Model Compression

- Tensor decomposition
- Network pruning
- Quantization
- Smaller model

Matrix decomposition

Fully-connected layer

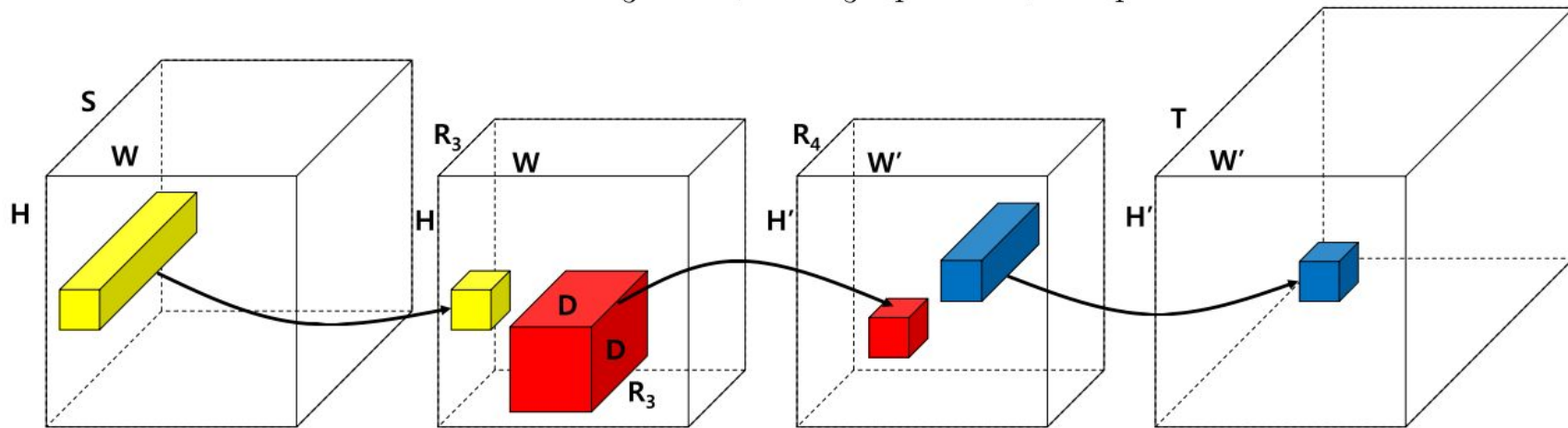
- Memory reduction: $\frac{MN}{(M+N)R}$
- Computation reduction: $\frac{MN}{(M+N)R}$



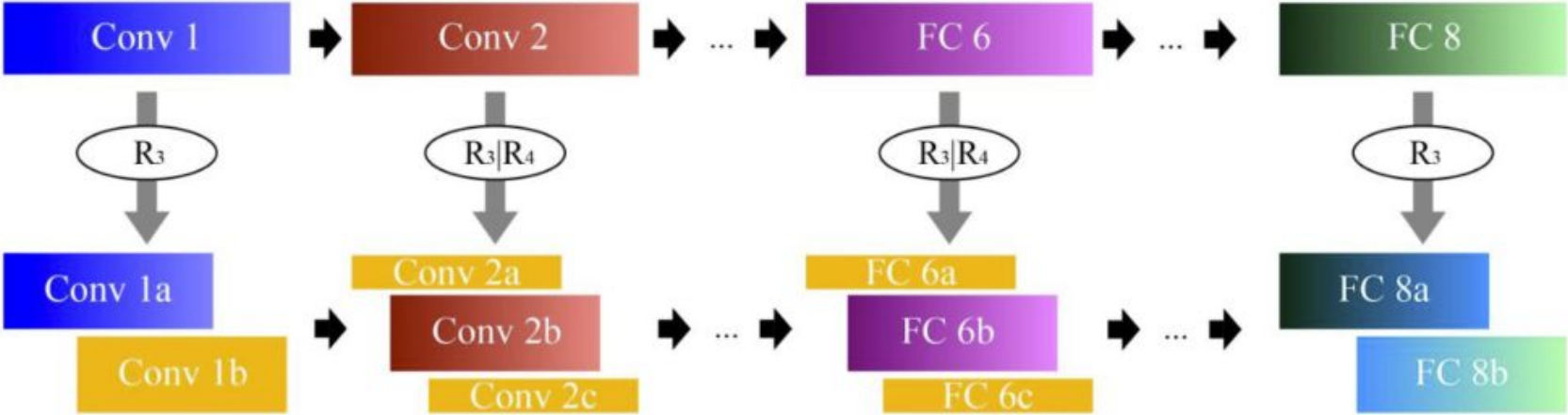
Tensor decomposition

Convolutional layer

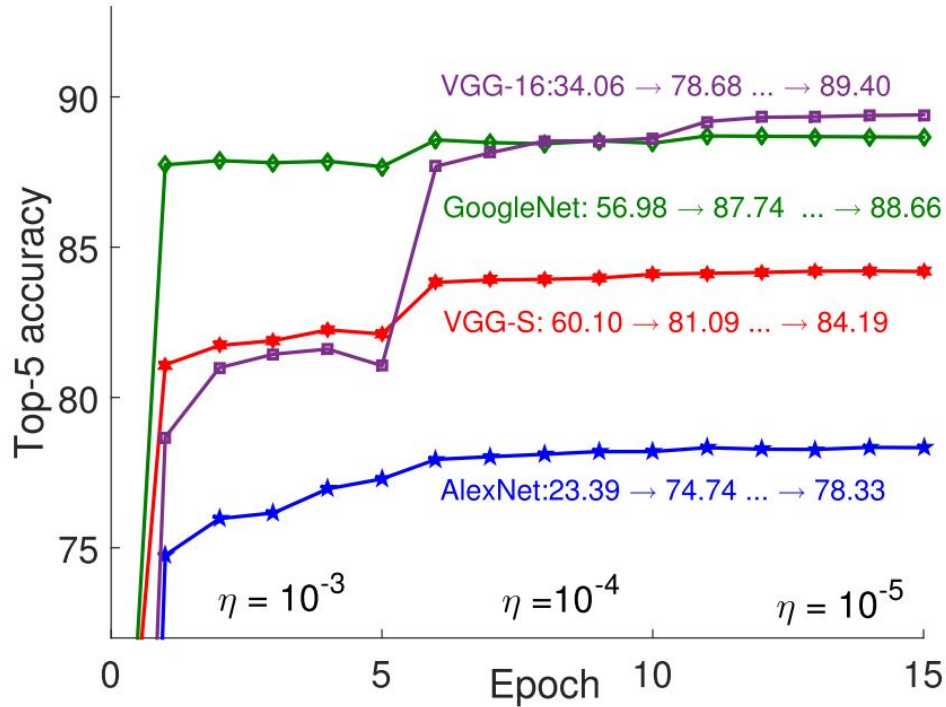
- Memory reduction: $\frac{D^2 ST}{SR_3 + D^2 R_3 R_4 + TR_4}$
 - Computation reduction: $\frac{D^2 ST H' W'}{SR_3 HW + D^2 R_3 R_4 H' W' + TR_4 H' W'}$
- $\xrightarrow{\text{bounded by}} ST/R_3 R_4$



Decompose the entire model



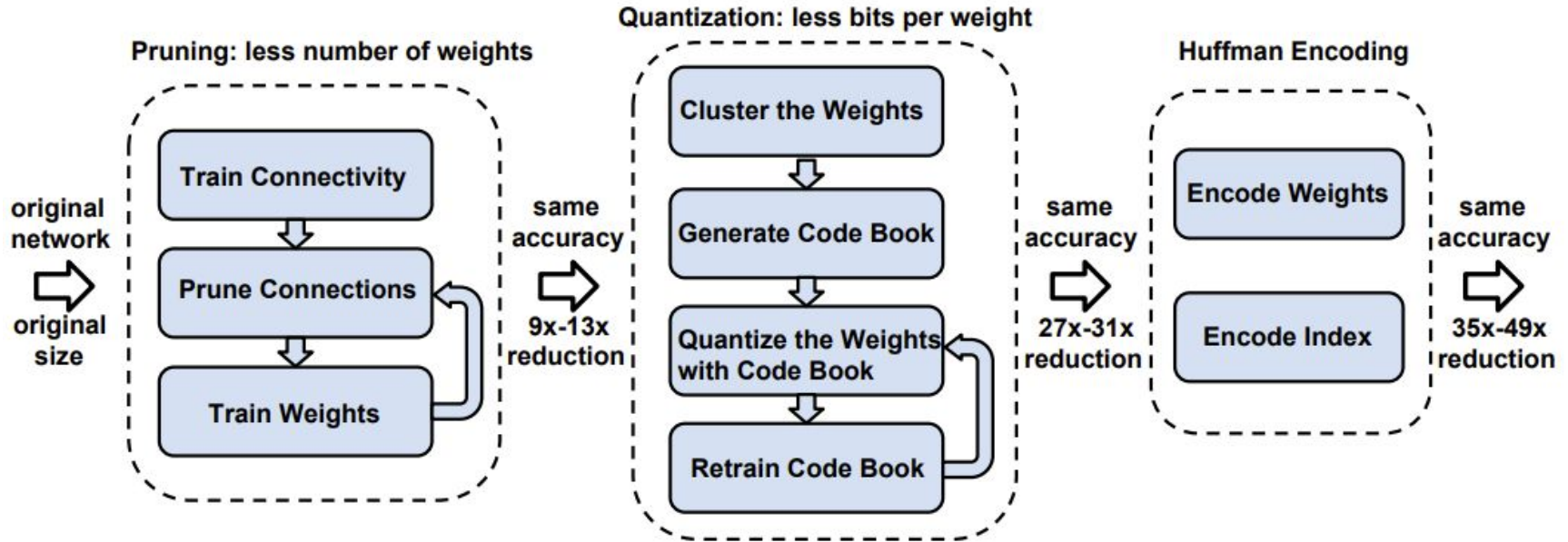
Fine-tuning after decomposition



Accuracy & Latency after Decomposition

Model	Top-5	Weights	FLOPs	S6		Titan X
<i>AlexNet</i>	80.03	61M	725M	117ms	245mJ	0.54ms
<i>AlexNet*</i> (imp.)	78.33 (-1.70)	11M ($\times 5.46$)	272M ($\times 2.67$)	43ms ($\times 2.72$)	72mJ ($\times 3.41$)	0.30ms ($\times 1.81$)
<i>VGG-S</i>	84.60	103M	2640M	357ms	825mJ	1.86ms
<i>VGG-S*</i> (imp.)	84.05 (-0.55)	14M ($\times 7.40$)	549M ($\times 4.80$)	97ms ($\times 3.68$)	193mJ ($\times 4.26$)	0.92ms ($\times 2.01$)
<i>GoogLeNet</i>	88.90	6.9M	1566M	273ms	473mJ	1.83ms
<i>GoogLeNet*</i> (imp.)	88.66 (-0.24)	4.7M ($\times 1.28$)	760M ($\times 2.06$)	192ms ($\times 1.42$)	296mJ ($\times 1.60$)	1.48ms ($\times 1.23$)
<i>VGG-16</i>	89.90	138M	15484M	1926ms	4757mJ	10.67ms
<i>VGG-16*</i> (imp.)	89.40 (-0.50)	127M ($\times 1.09$)	3139M ($\times 4.93$)	576ms ($\times 3.34$)	1346mJ ($\times 3.53$)	4.58ms ($\times 2.33$)

Network pruning: Deep compression



* Song Han, et al. "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding." ICLR (2016).

Network pruning: prune the connections

Initialization: $W^{(0)}$ with $W^{(0)} \sim N(0, \Sigma)$, $iter = 0$.

Hyper-parameter: *threshold*, δ .

Output: $W^{(t)}$.

————— *Train Connectivity* —————

```
while not converged do
  |  $W^{(t)} = W^{(t-1)} - \eta^{(t)} \nabla f(W^{(t-1)}; x^{(t-1)});$ 
  |  $t = t + 1;$ 
end
```

————— *Prune Connections* —————

// initialize the mask by thresholding the weights.

$Mask = \mathbf{1}(|W| > threshold);$

$W = W \cdot Mask;$

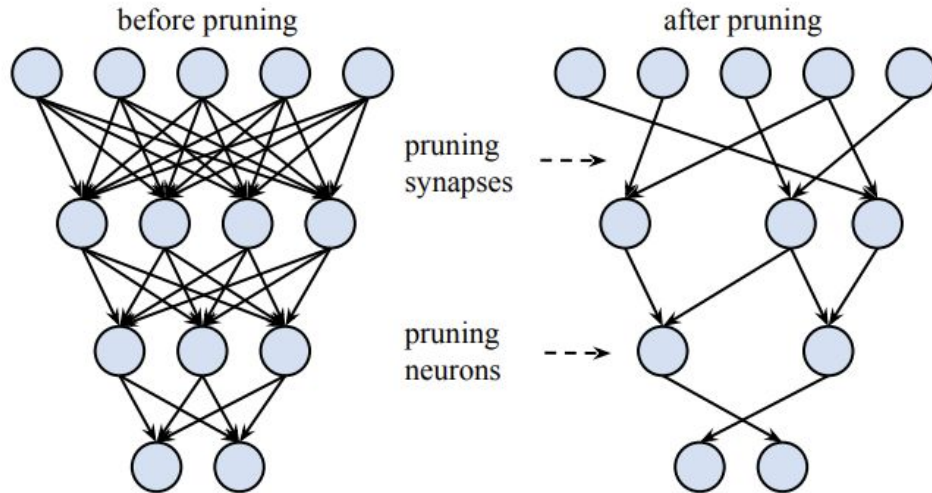
————— *Retrain Weights* —————

```
while not converged do
  |  $W^{(t)} = W^{(t-1)} - \eta^{(t)} \nabla f(W^{(t-1)}; x^{(t-1)});$ 
  |  $W^{(t)} = W^{(t)} \cdot Mask;$ 
  |  $t = t + 1;$ 
end
```

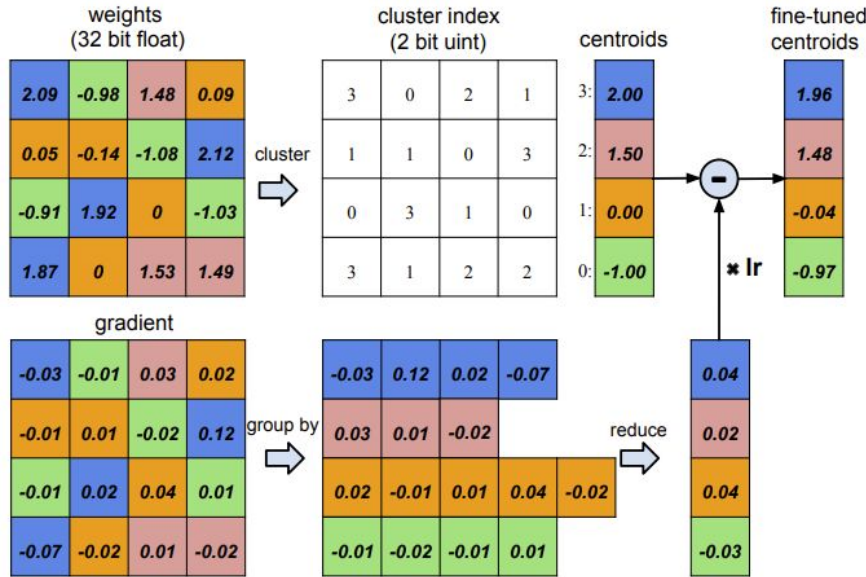
————— *Iterative Pruning* —————

$threshold = threshold + \delta[iter ++];$

goto *Pruning Connections*;



Network pruning: weight sharing



Weight sharing by scalar quantization (top) and centroids fine-tuning (bottom).

1. Use k-means clustering to identify the shared weights for each layer of a trained network. Minimize

$$\arg \min_C \sum_{i=1}^k \sum_{w \in c_i} |w - c_i|^2$$

2. Finetune the neural network using shared weights.

Network pruning: accuracy

Table 1: The compression pipeline can save $35\times$ to $49\times$ parameter storage with no loss of accuracy.

Network	Top-1 Error	Top-5 Error	Parameters	Compress Rate
LeNet-300-100 Ref	1.64%	-	1070 KB	
LeNet-300-100 Compressed	1.58%	-	27 KB	40\times
LeNet-5 Ref	0.80%	-	1720 KB	
LeNet-5 Compressed	0.74%	-	44 KB	39\times
AlexNet Ref	42.78%	19.73%	240 MB	
AlexNet Compressed	42.78%	19.70%	6.9 MB	35\times
VGG-16 Ref	31.50%	11.32%	552 MB	
VGG-16 Compressed	31.17%	10.91%	11.3 MB	49\times

Network pruning: accuracy vs compression

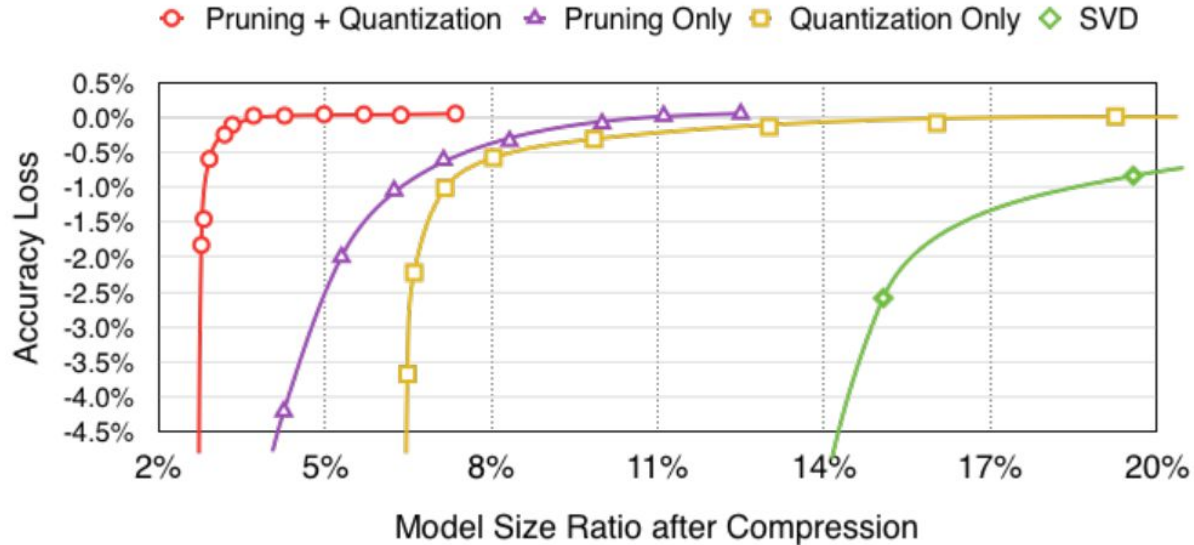
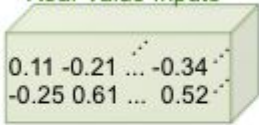
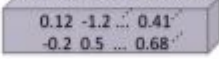
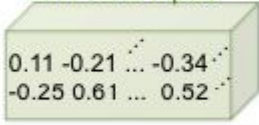
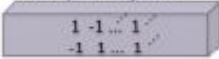

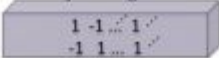


Figure 6: Accuracy v.s. compression rate under different compression methods. Pruning and quantization works best when combined.

Quantization: XNOR-Net

	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)
Standard Convolution	<p>Real-Value Inputs</p>  <p>Real-Value Weights</p> 	$+, -, \times$	1x	1x	%56.7
Binary Weight	<p>Real-Value Inputs</p>  <p>Binary Weights</p> 	$+, -$	$\sim 32x$	$\sim 2x$	%56.8
BinaryWeight Binary Input (XNOR-Net)	<p>Binary Inputs</p>  <p>Binary Weights</p> 	XNOR, bitcount	$\sim 32x$	$\sim 58x$	%44.2

* Mohammad Rastegari, et al. "XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks." ECCV (2016).

Quantization: binary weights

Algorithm 1 Training an L -layers CNN with binary weights:

Input: A minibatch of inputs and targets (\mathbf{I}, \mathbf{Y}) , cost function $C(\mathbf{Y}, \hat{\mathbf{Y}})$, current weight \mathcal{W}^t and current learning rate η^t .

Output: updated weight \mathcal{W}^{t+1} and updated learning rate η^{t+1} .

- 1: Binarizing weight filters:
 - 2: **for** $l = 1$ to L **do**
 - 3: **for** k^{th} filter in l^{th} layer **do**
 - 4: $\mathcal{A}_{lk} = \frac{1}{n} \|\mathcal{W}_{lk}^t\|_{\ell_1}$
 - 5: $\mathcal{B}_{lk} = \text{sign}(\mathcal{W}_{lk}^t)$
 - 6: $\widetilde{\mathcal{W}}_{lk} = \mathcal{A}_{lk} \mathcal{B}_{lk}$
 - 7: $\hat{\mathbf{Y}} = \mathbf{BinaryForward}(\mathbf{I}, \mathcal{B}, \mathcal{A})$ // standard forward propagation except that convolutions are computed using equation 1 or 11
 - 8: $\frac{\partial C}{\partial \widetilde{\mathcal{W}}} = \mathbf{BinaryBackward}(\frac{\partial C}{\partial \hat{\mathbf{Y}}}, \widetilde{\mathcal{W}})$ // standard backward propagation except that gradients are computed using $\widetilde{\mathcal{W}}$ instead of \mathcal{W}^t
 - 9: $\mathcal{W}^{t+1} = \mathbf{UpdateParameters}(\mathcal{W}^t, \frac{\partial C}{\partial \widetilde{\mathcal{W}}}, \eta^t)$ // Any update rules (e.g.,SGD or ADAM)
 - 10: $\eta^{t+1} = \mathbf{UpdateLearningrate}(\eta^t, t)$ // Any learning rate scheduling function
-

Quantization: binary input and weights

Binary Dot Product: To approximate the dot product between $\mathbf{X}, \mathbf{W} \in \mathbb{R}^n$ such that $\mathbf{X}^\top \mathbf{W} \approx \beta \mathbf{H}^\top \alpha \mathbf{B}$, where $\mathbf{H}, \mathbf{B} \in \{+1, -1\}^n$ and $\beta, \alpha \in \mathbb{R}^+$, we solve the following optimization:

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \underset{\alpha, \mathbf{B}, \beta, \mathbf{H}}{\operatorname{argmin}} \|\mathbf{X} \odot \mathbf{W} - \beta \alpha \mathbf{H} \odot \mathbf{B}\| \quad (7)$$

$$\mathbf{C}^* = \operatorname{sign}(\mathbf{Y}) = \operatorname{sign}(\mathbf{X}) \odot \operatorname{sign}(\mathbf{W}) = \mathbf{H}^* \odot \mathbf{B}^*$$

$$\gamma^* = \frac{\sum |\mathbf{Y}_i|}{n} = \frac{\sum |\mathbf{X}_i| |\mathbf{W}_i|}{n} \approx \left(\frac{1}{n} \|\mathbf{X}\|_{\ell_1} \right) \left(\frac{1}{n} \|\mathbf{W}\|_{\ell_1} \right) = \beta^* \alpha^*$$

Quantization: accuracy

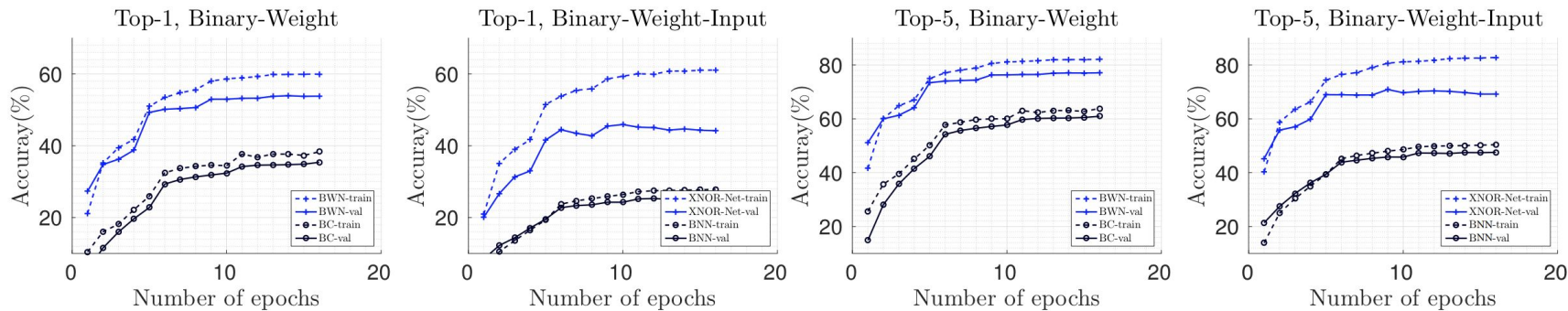
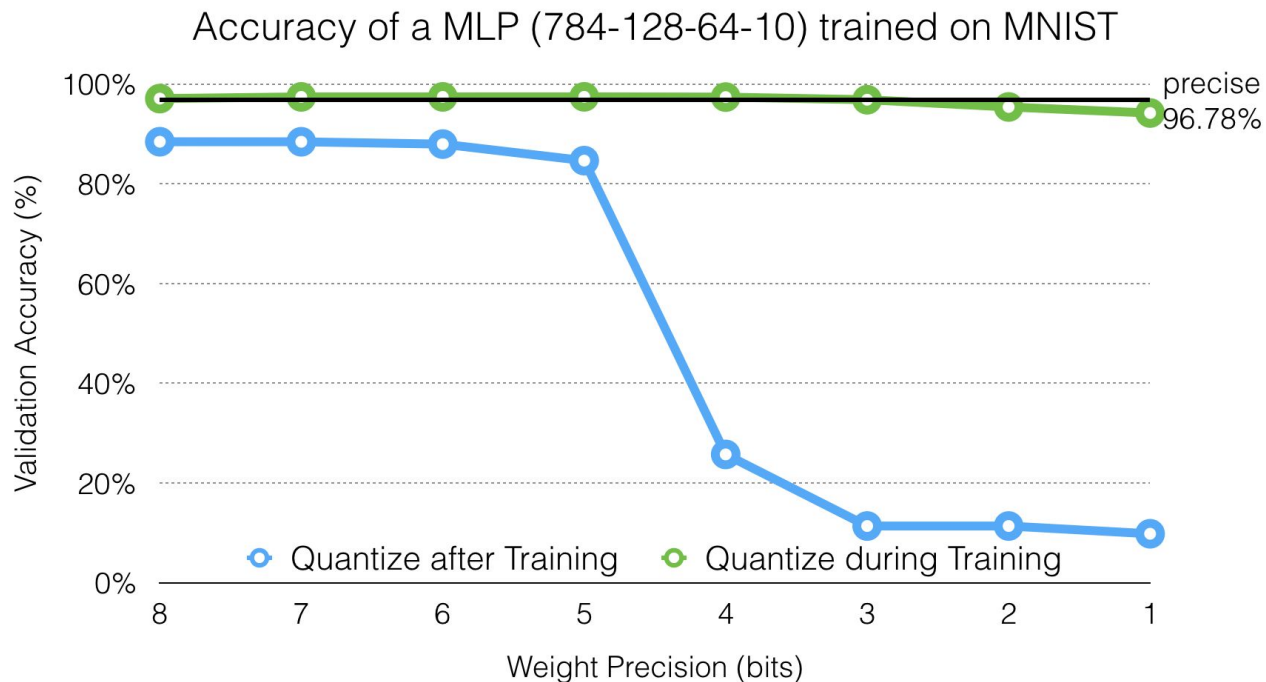
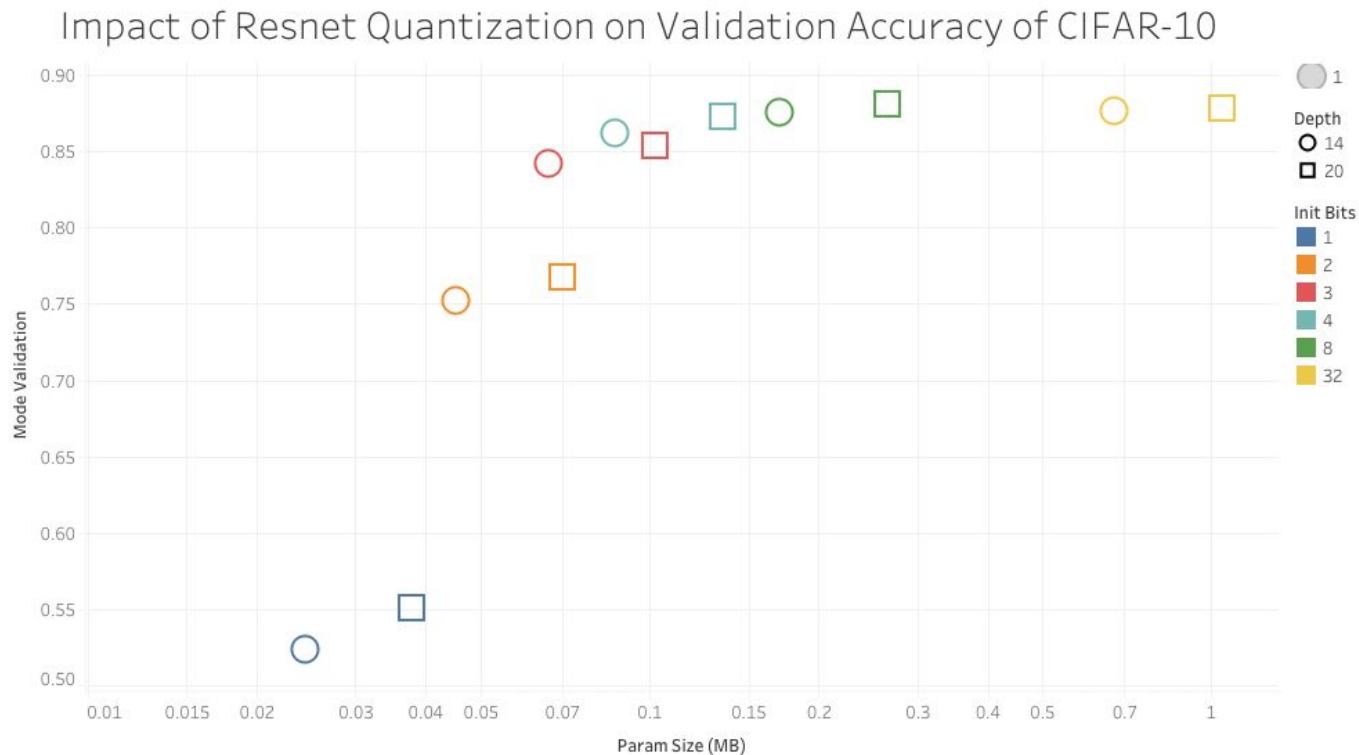


Fig. 5: This figure compares the imagenet classification accuracy on Top-1 and Top-5 across training epochs. Our approaches BWN and XNOR-Net outperform BinaryConnect(BC) and BinaryNet(BNN) in all the epochs by large margin($\sim 17\%$).

Quantization

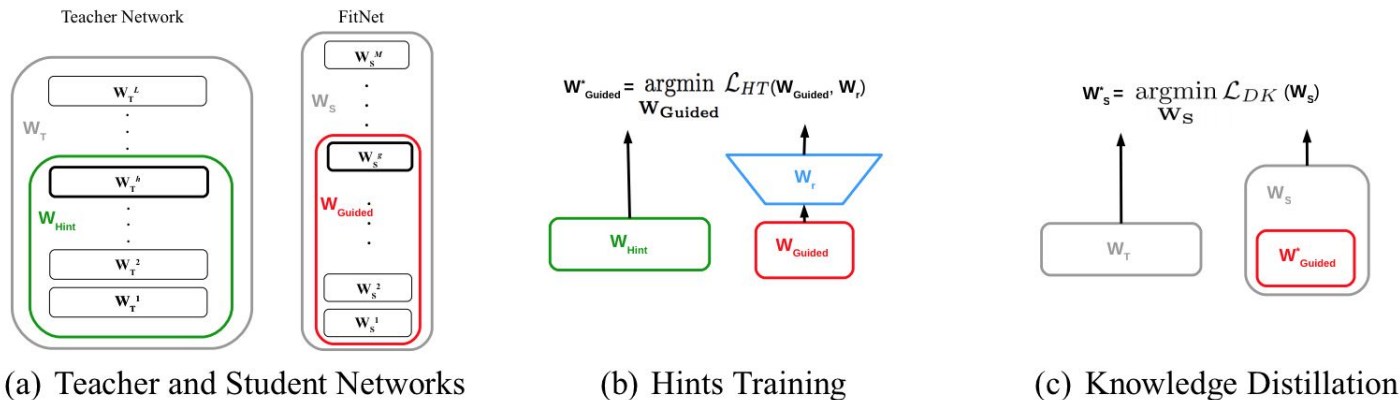


Quantization



Smaller model: Knowledge distillation

- Knowledge distillation: use a teacher model (large model) to train a student model (small model)



$$\mathcal{L}_{HT}(\mathbf{W}_{\text{Guided}}, \mathbf{W}_{\mathbf{r}}) = \frac{1}{2} \| |u_h(\mathbf{x}; \mathbf{W}_{\text{Hint}}) - r(v_g(\mathbf{x}; \mathbf{W}_{\text{Guided}}); \mathbf{W}_{\mathbf{r}}) | \|^2,$$

* Romero, Adriana, et al. "Fitnets: Hints for thin deep nets." ICLR (2015).

Smaller model: accuracy

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	91.61%
Teacher	~9M	90.18%
Mimic single	~54M	84.6%
Mimic single	~70M	84.9%
Mimic ensemble	~70M	85.8%
<i>State-of-the-art methods</i>		
Maxout		90.65%
Network in Network		91.2%
Deeply-Supervised Networks		91.78%
Deeply-Supervised Networks (19)		88.2%

Table 1: Accuracy on CIFAR-10

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	64.96%
Teacher	~9M	63.54%
<i>State-of-the-art methods</i>		
Maxout		61.43%
Network in Network		64.32%
Deeply-Supervised Networks		65.43%

Table 2: Accuracy on CIFAR-100

Discussion

What are the implications of these model compression techniques for serving?

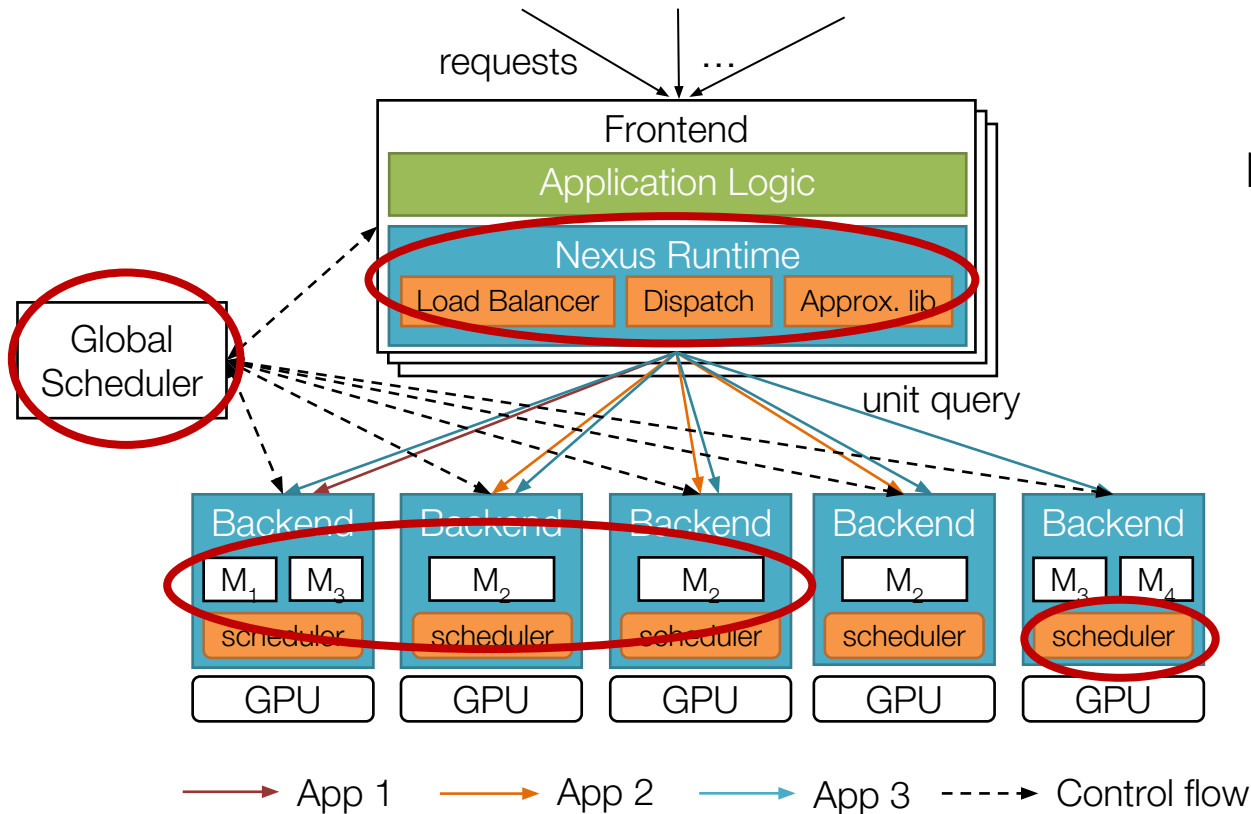
- Specialized hardware for sparse models
 - Song Han, et al. "EIE: Efficient Inference Engine on Compressed Deep Neural Network." ISCA 2016
- Accuracy and resource trade-off
 - Han, Seungyeop, et al. "MCDNN: An Approximation-Based Execution Framework for Deep Stream Processing Under Resource Constraints." MobiSys (2016).

Serving system

Serving system

- **Goals:**
 - High flexibility for writing applications
 - High efficiency on GPUs
 - Satisfy latency SLA
- **Challenges**
 - Provide common abstraction for different frameworks
 - Achieve high efficiency
 - Sub-second latency SLA that limits the batch size
 - Model optimization and multi-tenancy causes long tail

Nexus: efficient neural network serving system



Remarks

- Frontend runtime library allows arbitrary app logic
- Packing models to achieve higher utilization
- A GPU scheduler allows new batching primitives
- A batch-aware global scheduler allocates GPU cycles for each model

Flexibility: Application runtime

```
class ModelHandler:
```

```
    # return output future
```

```
    def Execute(input)
```

Async RPC, execute the model remotely

```
class AppBase:
```

```
    # return ModelHandler
```

Send load model request to global scheduler

```
    def GetModelHandler(framework, model, version, latency_sla)
```

```
    # Load models during setup time, implemented by developer
```

```
    def Setup()
```

```
    # Process requests, implemented by developer
```

```
    def Process(request)
```

Application example: Face recognition

```
class FaceRecApp(AppBase):
```

```
    def Setup(self):
```

```
        self.m1 = self.GetModelHandler("caffe", "vgg_face", 1, 100)
```

```
        self.m2 = self.GetModelHandler("mxnet", "age_net", 1, 100)
```

Load model from different framework

```
    def Process(self, request):
```

```
        ret1 = self.m1.Execute(request.image)
```

```
        ret2 = self.m2.Execute(request.image)
```

Execute concurrently on remote GPUs

```
        return Reply(request.user_id, ret1["name"], ret2["age"])
```

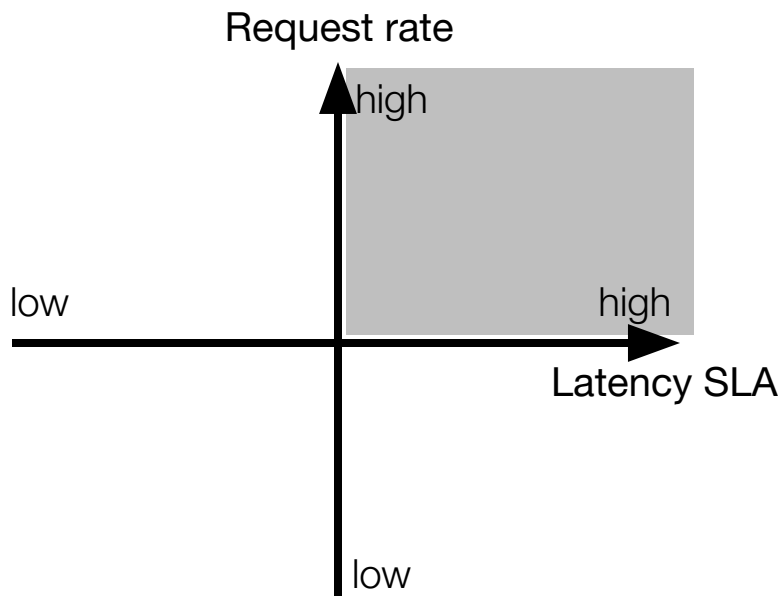
Force to synchronize when accessing future data

Application example: Traffic Analysis

```
class TrafficApp(AppBase):
    def Setup(self):
        self.det = self.GetModelHandler("darknet", "yolo9000", 1, 300)
        self.r1 = self.GetModelHandler("caffe", "vgg_face", 1, 150)
        self.r2 = self.GetModelHandler("caffe", "googlenet_car", 1, 150)

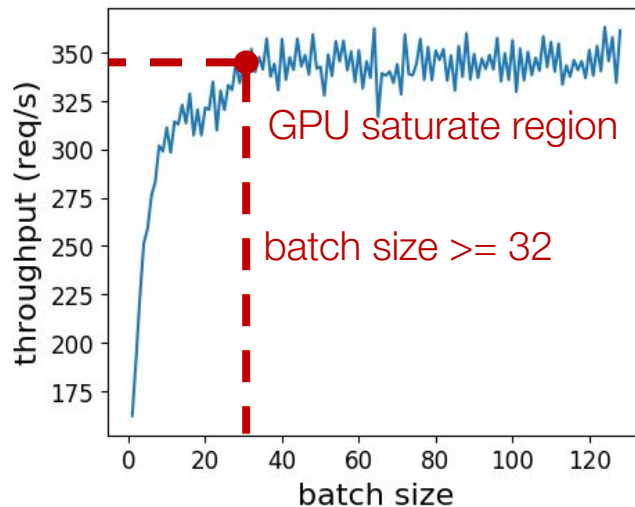
    def Process(self, request):
        persons, cars = [], []
        for obj in self.det.Execute(request.image):
            if obj["class"] == "person":
                persons.append(self.r1.Execute(request.image[obj["rect"]]))
            elif obj["class"] == "car":
                cars.append(self.r2.Execute(request.image[obj["rect"]]))
        return Reply(request.user_id, persons, cars)
```


High Efficiency



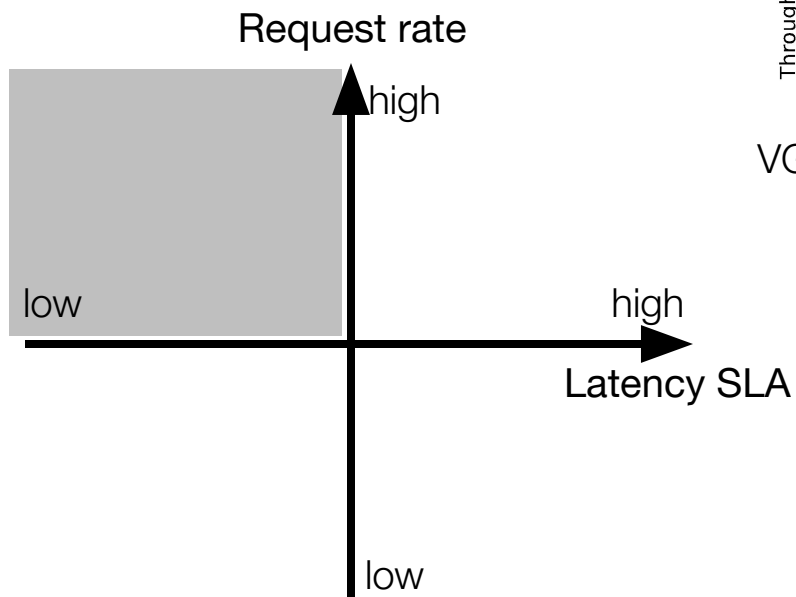
Workload characteristic

- For high request rate, high latency SLA workload, saturate GPU efficiency by using large batch size

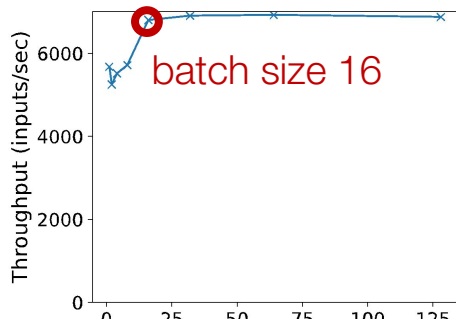


VGG16 throughput vs batch size

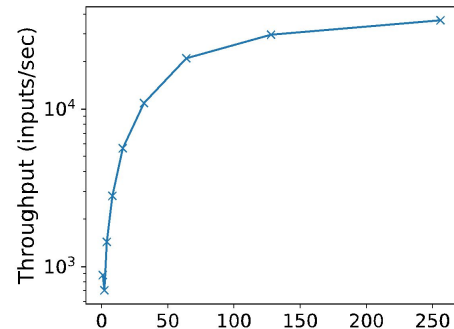
High Efficiency



Workload characteristic

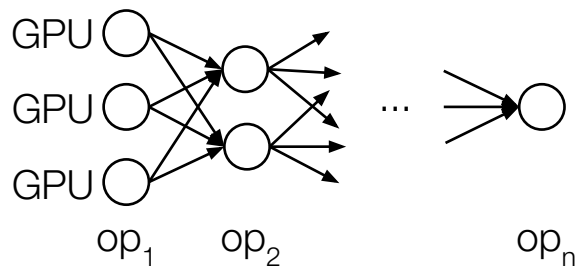


VGG16 conv2_1 (winograd)



VGG16 fc6

- Suppose we can choose a different batch size for each op (layer), and allocate dedicated GPUs for each op.



Split Batching

- Optimization problem

$$\max_{b_1 \dots b_n} \frac{\max_j tp_j(b_j)}{\sum_i \max_j tp_j(b_j) / tp_i(b_i)}$$

equivalent to

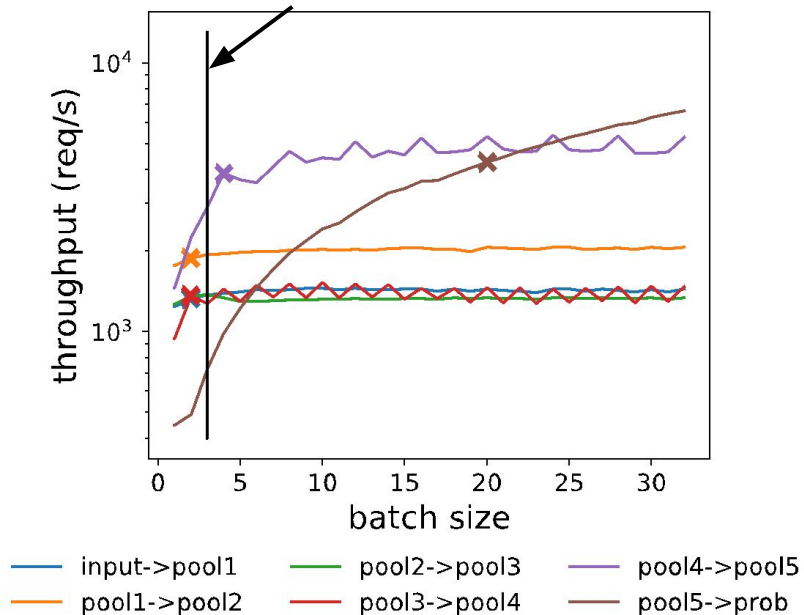
$$\min_{b_1 \dots b_n} \sum_i 1/tp_i(b_i)$$

such that

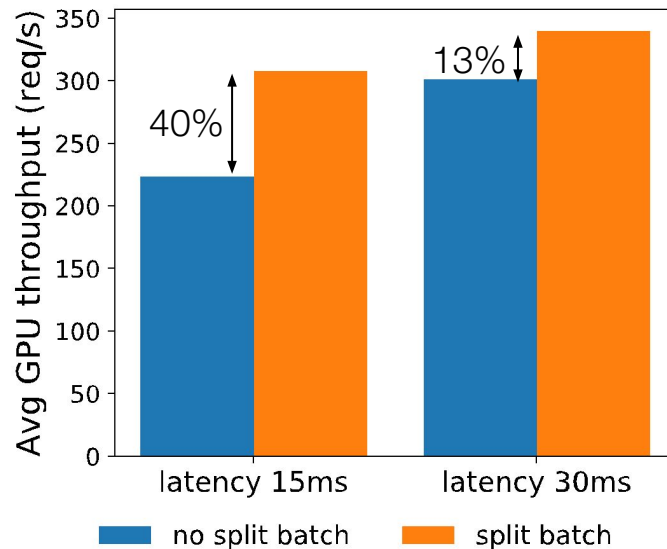
$$\sum_i lat_i(b_i) + \sum_{b_i \neq b_{i+1}} overhead(out_i * b_i) \leq latency_sla$$

Split Batching

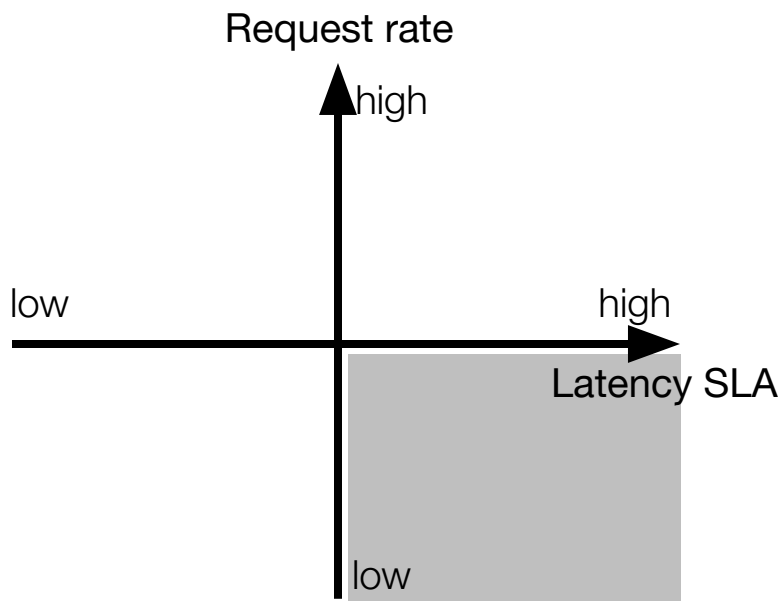
Batch size 3 for entire model



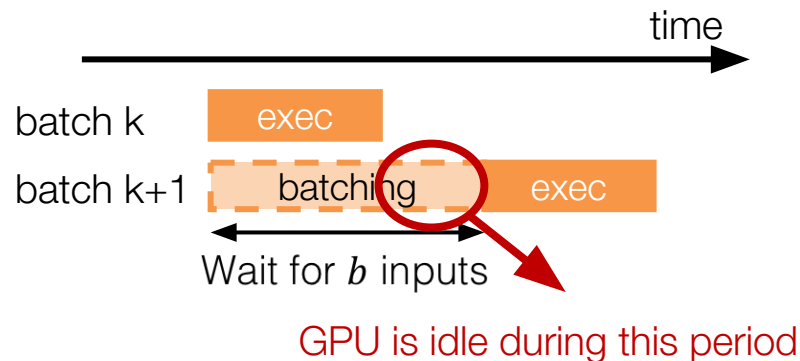
VGG16 optimal batch size for each segment
for latency SLA 15ms



High Efficiency



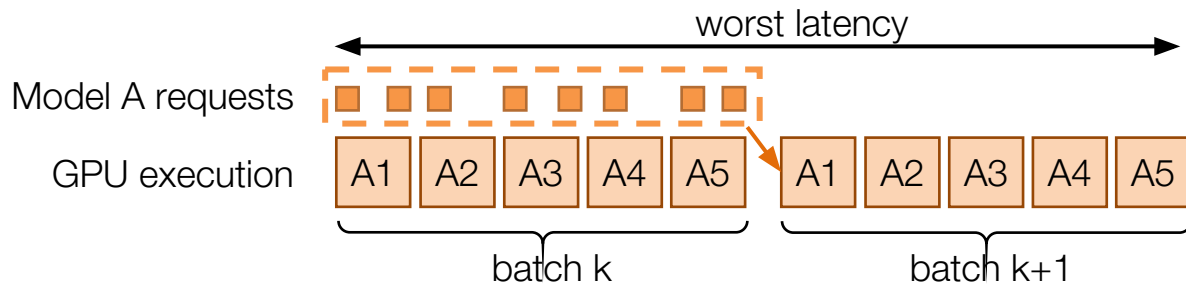
- This type of workload cannot saturate GPU in temporal domain
- Suppose the optimal batch size is b under latency SLA



Workload characteristics

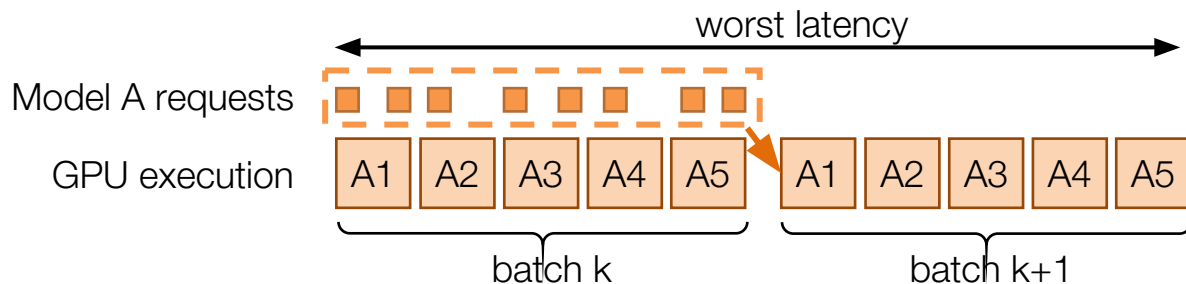
Execute multiple models on one GPU

Execute multiple models on a single GPU

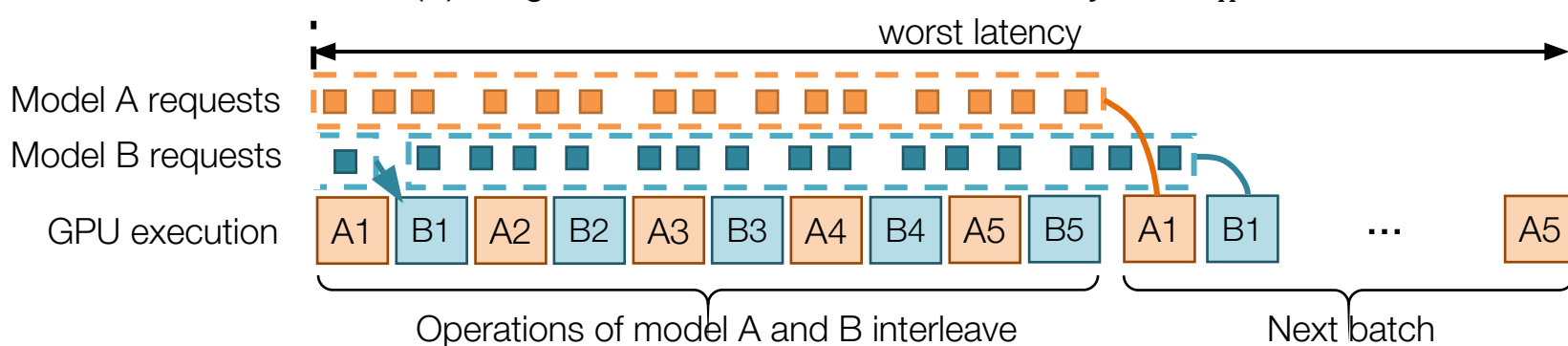


(a) Single model execution: worst latency is $2lat_A$

Execute multiple models on a single GPU

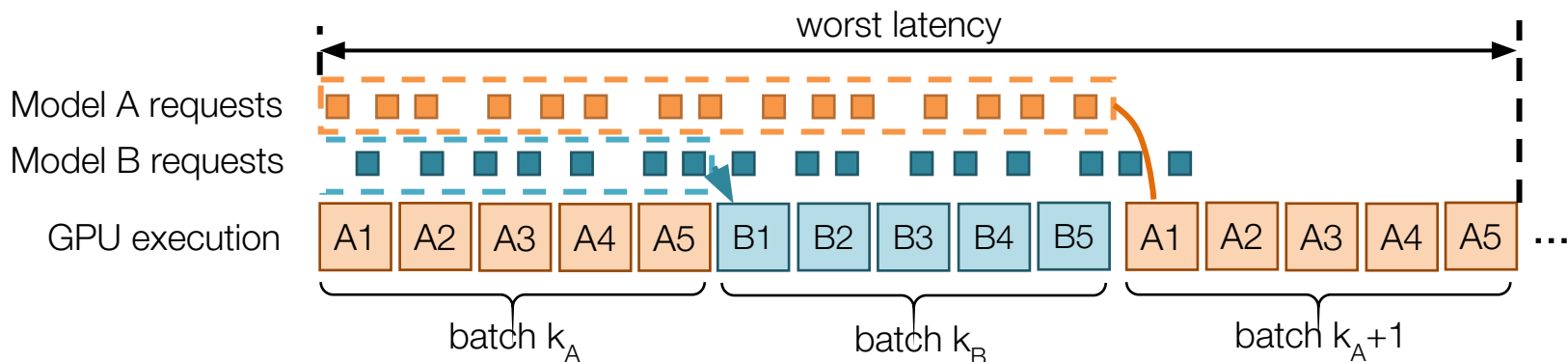


(a) Single model execution: worst latency is $2lat_A$



(b) Execute multiple models **concurrently**:
worst latency for model A and B is $2(lat_A + lat_B)$

Execute multiple models on a single GPU



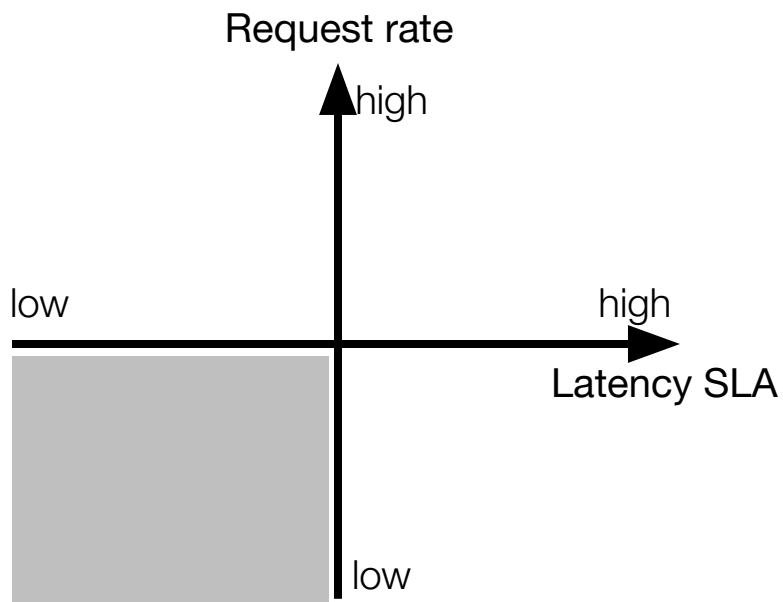
(c) Execute multiple models in **round-robin** fashion:

worst latency for model A is $2lat_A + lat_B$

worst latency for model B is $lat_A + 2lat_B$

Use larger batch size as latency is reduced and predictive

High Efficiency



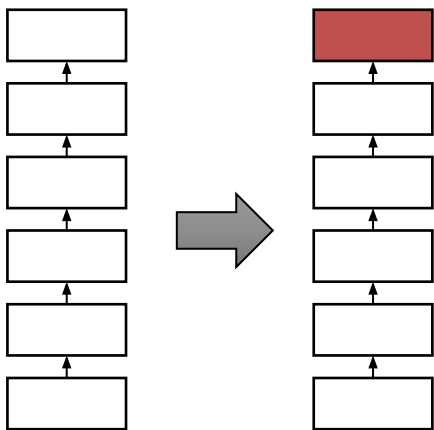
Workload characteristic

Solution depends

- If saturate GPU in temporal domain due to low latency: allocate dedicated GPU(s)
- If not: can use multi-batching to share GPU cycles with other models

Prefix batching for model specialization

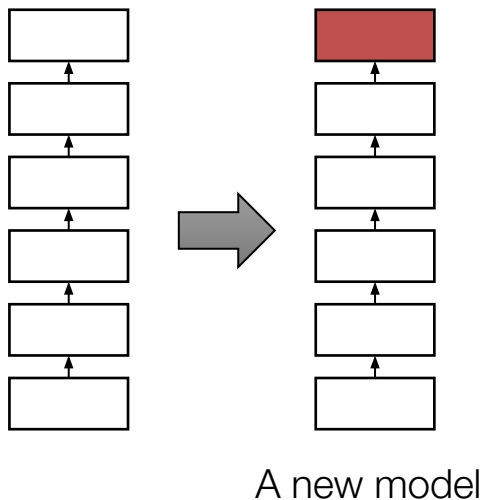
- Specialization (long-term / short-term) re-trains last a few layers



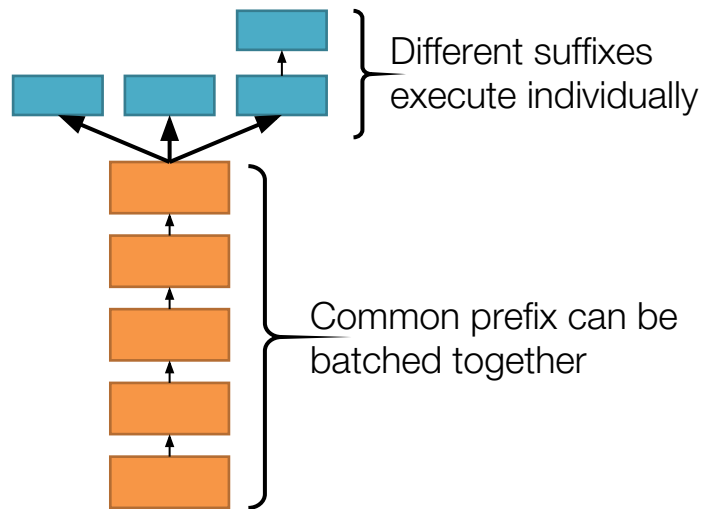
A new model

Prefix batching for model specialization

- Specialization (long-term / short-term) re-trains last a few layers



- Prefix batching allows batch execution for common prefix



Meet Latency SLA: Global scheduler

- First apply **split batching** and **prefix batching** if possible
- **Multi-batching**: bin-packing problem to pack models into GPUs
- Bin-packing optimization goal
 - Minimize the resource usage (number of GPUs)
- Constraint
 - Requests have to be served within latency SLA
- Degrees of freedom
 - Split workloads into smaller tasks
 - Change the batch size

Best-fit decreasing algorithms

1. For each workload, T_i is the max throughput that can be achieved on a GPU within latency SLA, and allocate dedicated GPUs

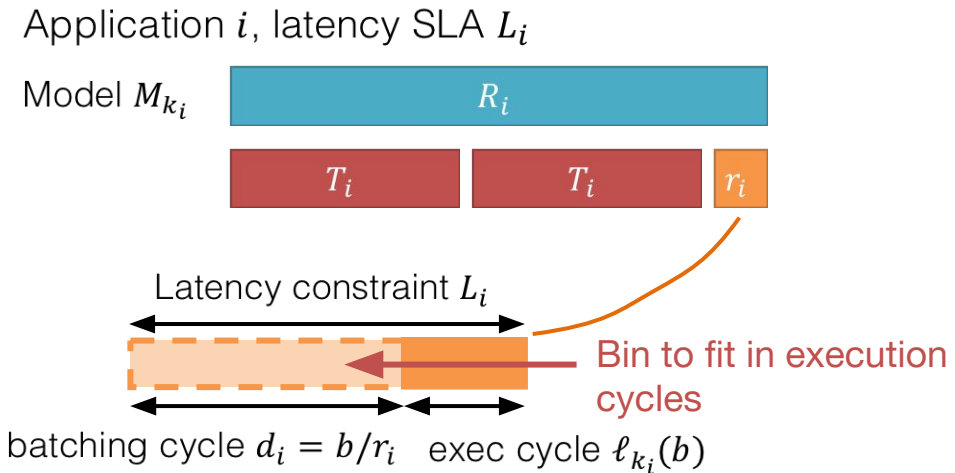
Application i , latency SLA L_i

Model M_{k_i}



Best-fit decreasing algorithms

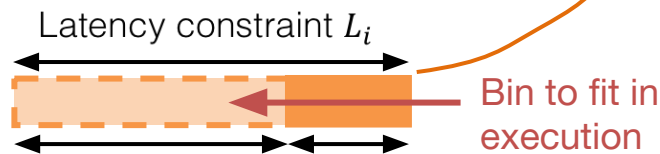
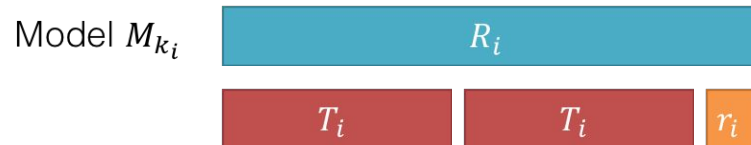
1. For each workload, T_i is the max throughput that can be achieved on a GPU within latency SLA, and allocate dedicated GPUs
2. For each residue workload, split latency SLA into batching cycle and exec cycle



Best-fit decreasing algorithms

1. For each workload, T_i is the max throughput that can be achieved on a GPU within latency SLA, and allocate dedicated GPUs
2. For each residue workload, split latency SLA into batching cycle and exec cycle
3. Sort all residue workloads by occupancy $\ell_{k_i}(B_i)/d_i$, and merge them by best-fit

Application i , latency SLA L_i



batching cycle $d_i = b/r_i$ exec cycle $\ell_{k_i}(b)$



Reference

- Kim, Yong-Deok, et al. "Compression of deep convolutional neural networks for fast and low power mobile applications." ICLR (2016).
- Song Han, et al. "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding." ICLR (2016).
- Romero, Adriana, et al. "Fitnets: Hints for thin deep nets." ICLR (2015).
- Mohammad Rastegari, et al. "XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks." ECCV (2016).
- Crankshaw, Daniel, et al. "Clipper: A Low-Latency Online Prediction Serving System." NSDI (2017).