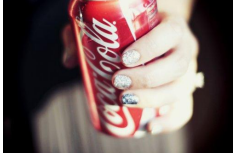


Lecture 14: Model Serving System

CSE599G1: Spring 2017

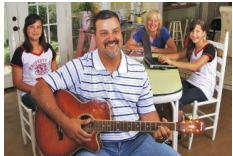
Deep Learning Applications



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



"I last saw your keys in the store room"

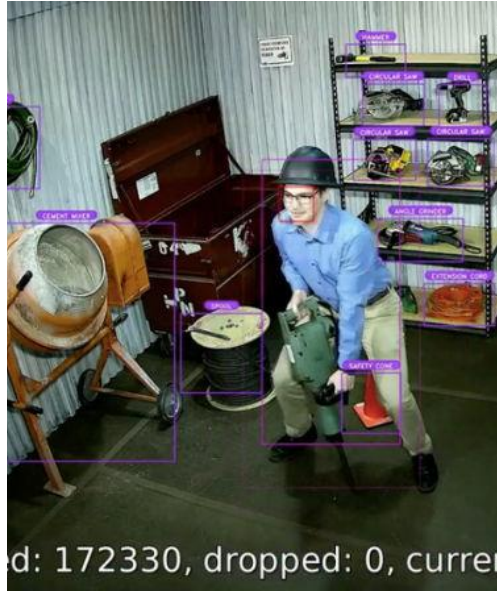


"Remind Tom of the party"

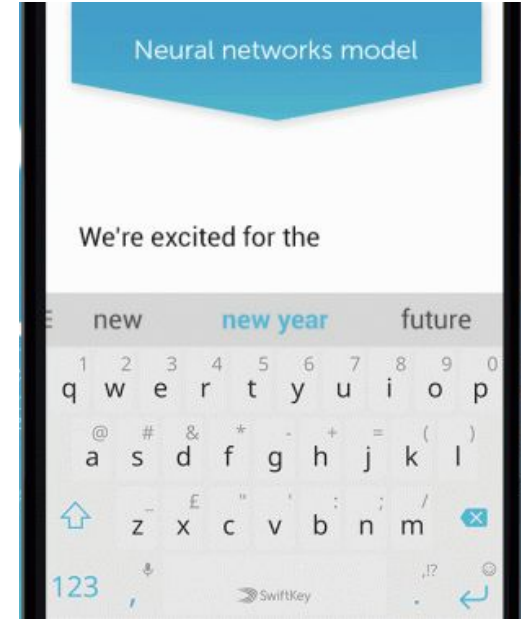


"You're on page 263 of this book"

Intelligent assistant

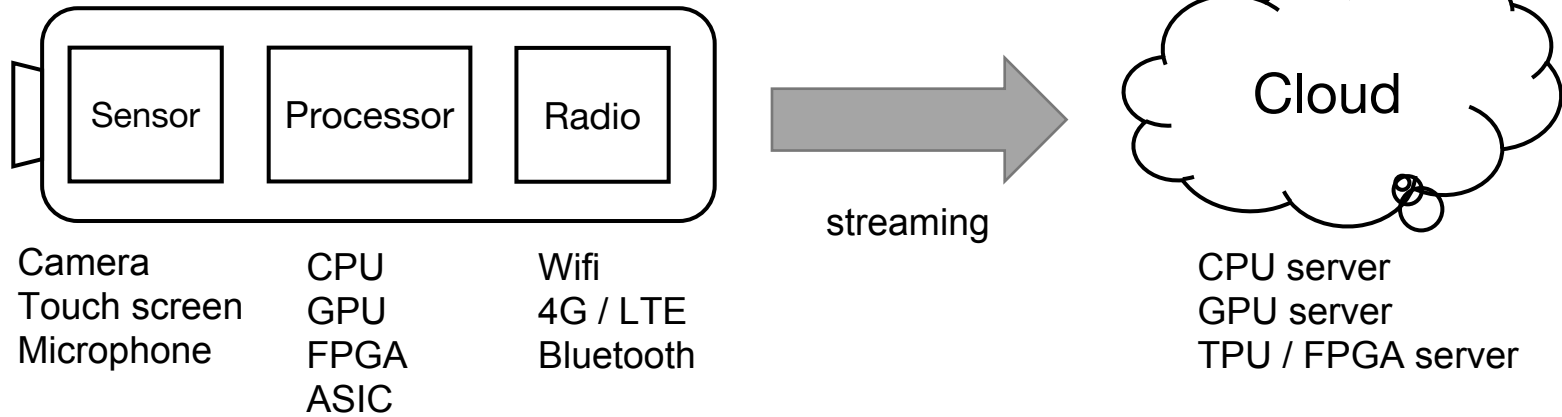


Surveillance / Remote assistance



Input keyboard

Runtime Environment



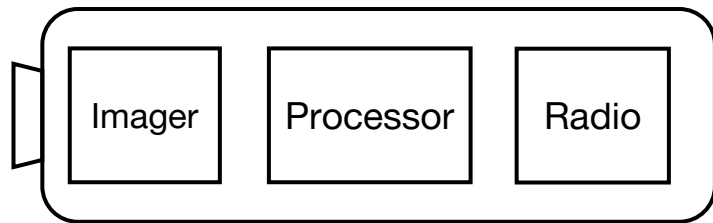
Resource usage for a continuous vision app

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

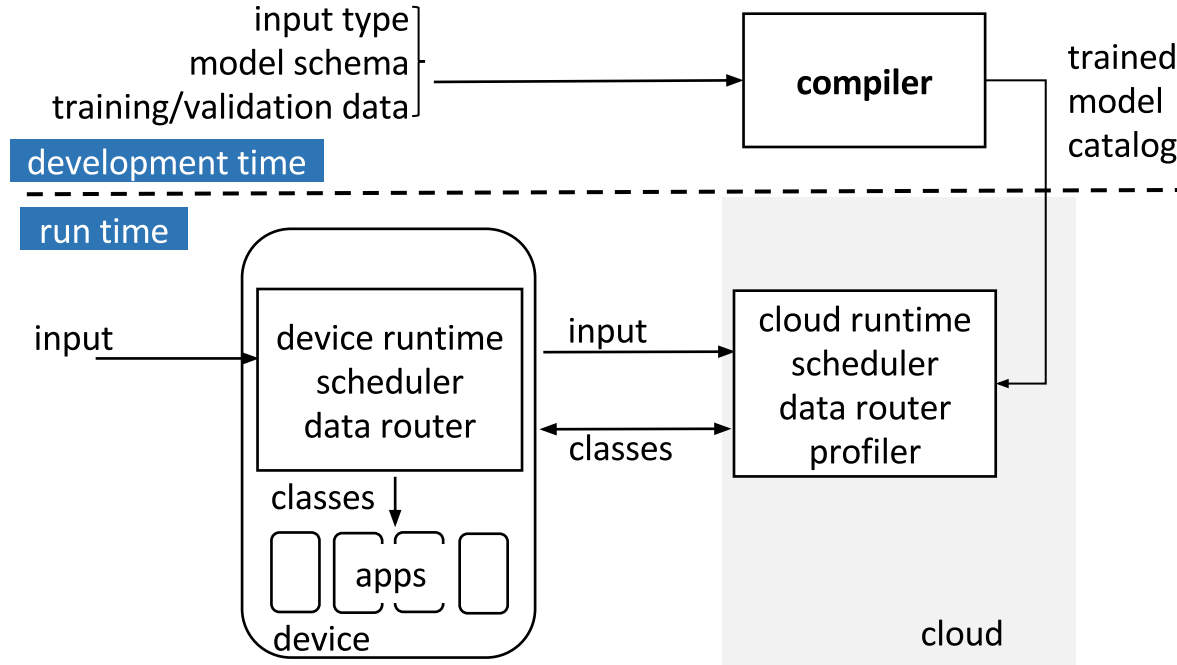
3.5GFLOPS (GPU) / 8GFLOPS (CPU)

Huge gap between workload and budget

Model Serving System 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

System overview



Outline

- Model compression
- Serving Backend
- Runtime scheduling between device and cloud

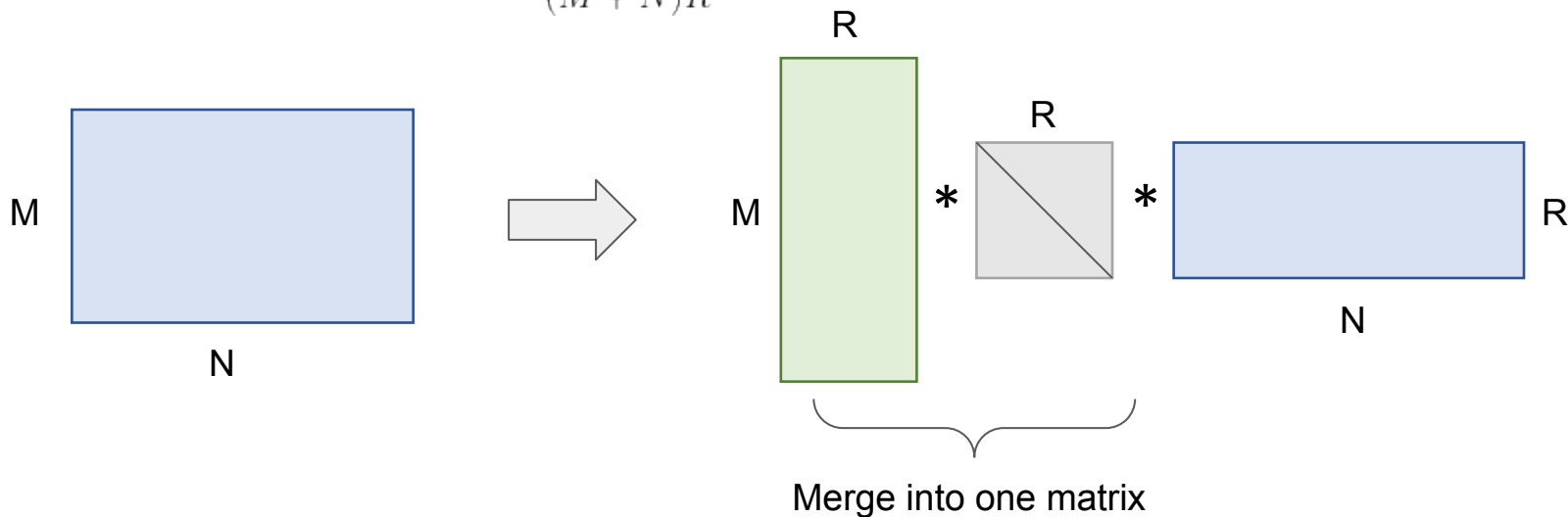
Model Compression

- Tensor decomposition
- Quantization
- Smaller model

Matrix Decomposition

Fully-connected layer

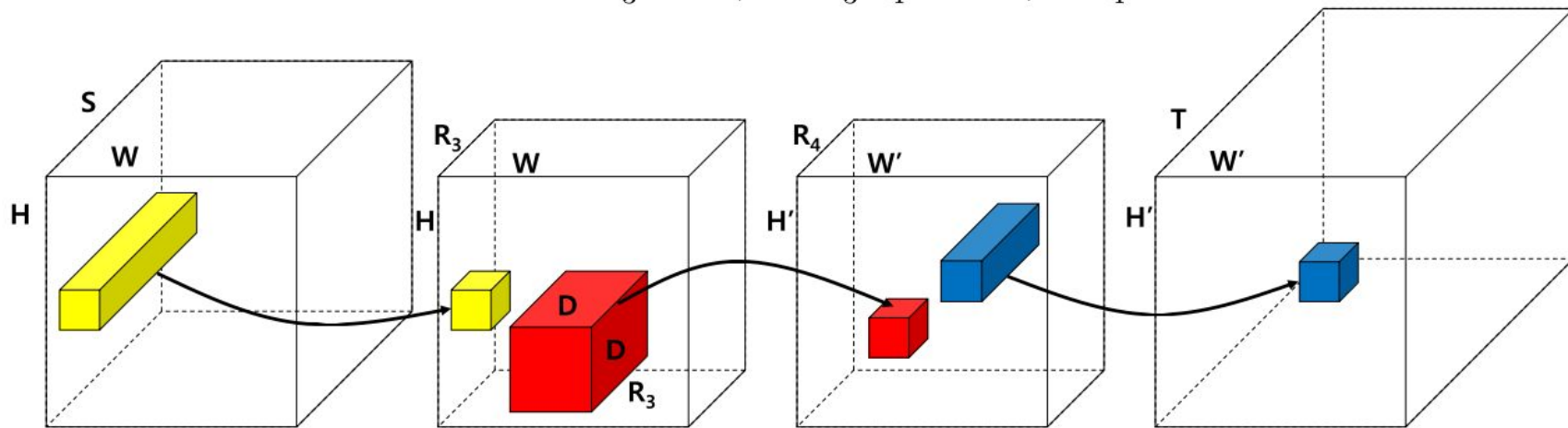
- Memory reduction: $\frac{MN}{(M+N)R}$
- Computation reduction: $\frac{MN}{(M+N)R}$



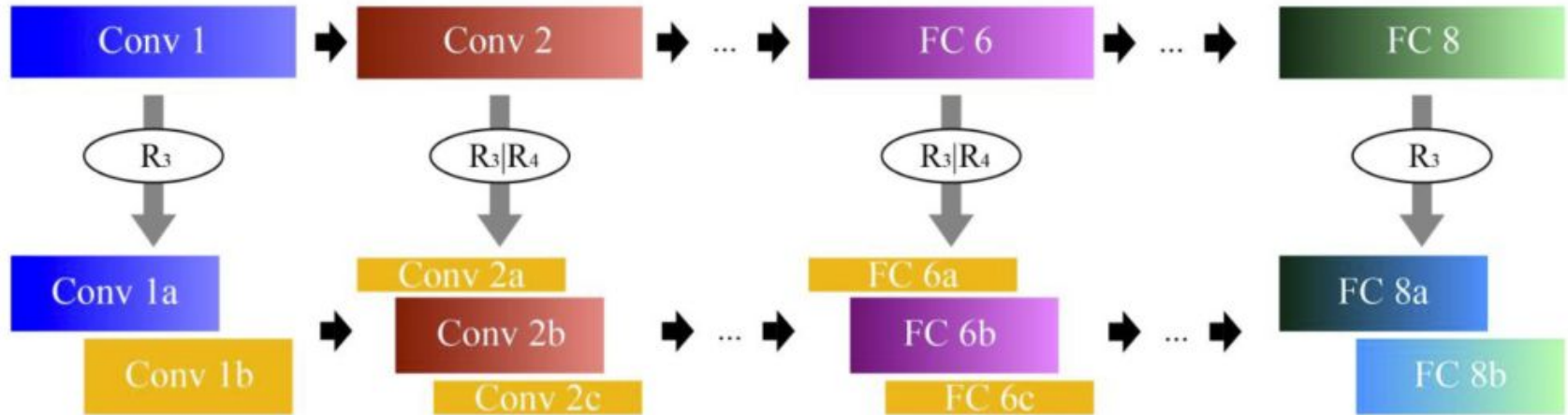
Tensor Decomposition

Convolutional layer

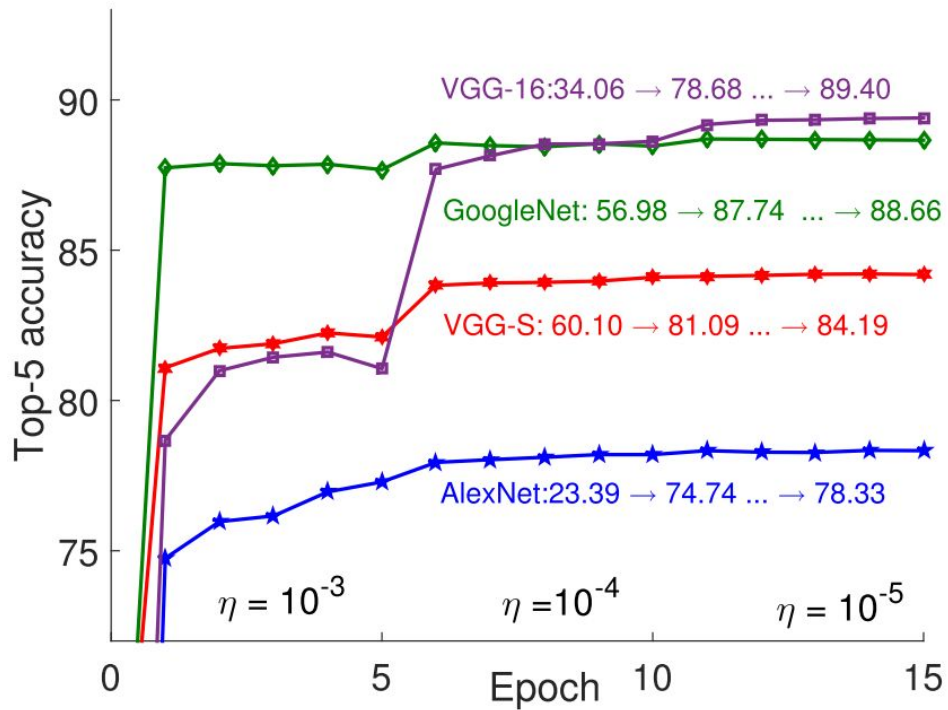
- Memory reduction: $\frac{D^2 ST}{SR_3 + D^2 R_3 R_4 + TR_4}$
- Computation reduction: $\frac{D^2 ST H' W'}{SR_3 H W + D^2 R_3 R_4 H' W' + TR_4 H' W'}$



Decompose the entire model



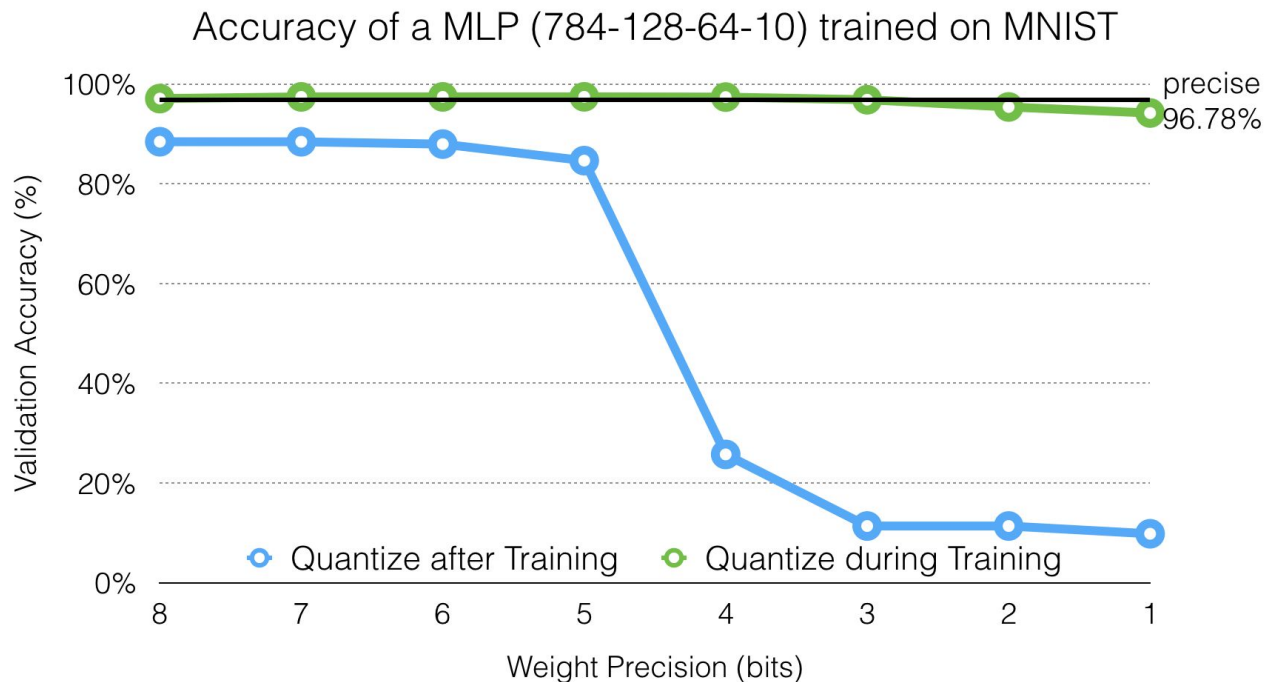
Fine-tuning after decomposition



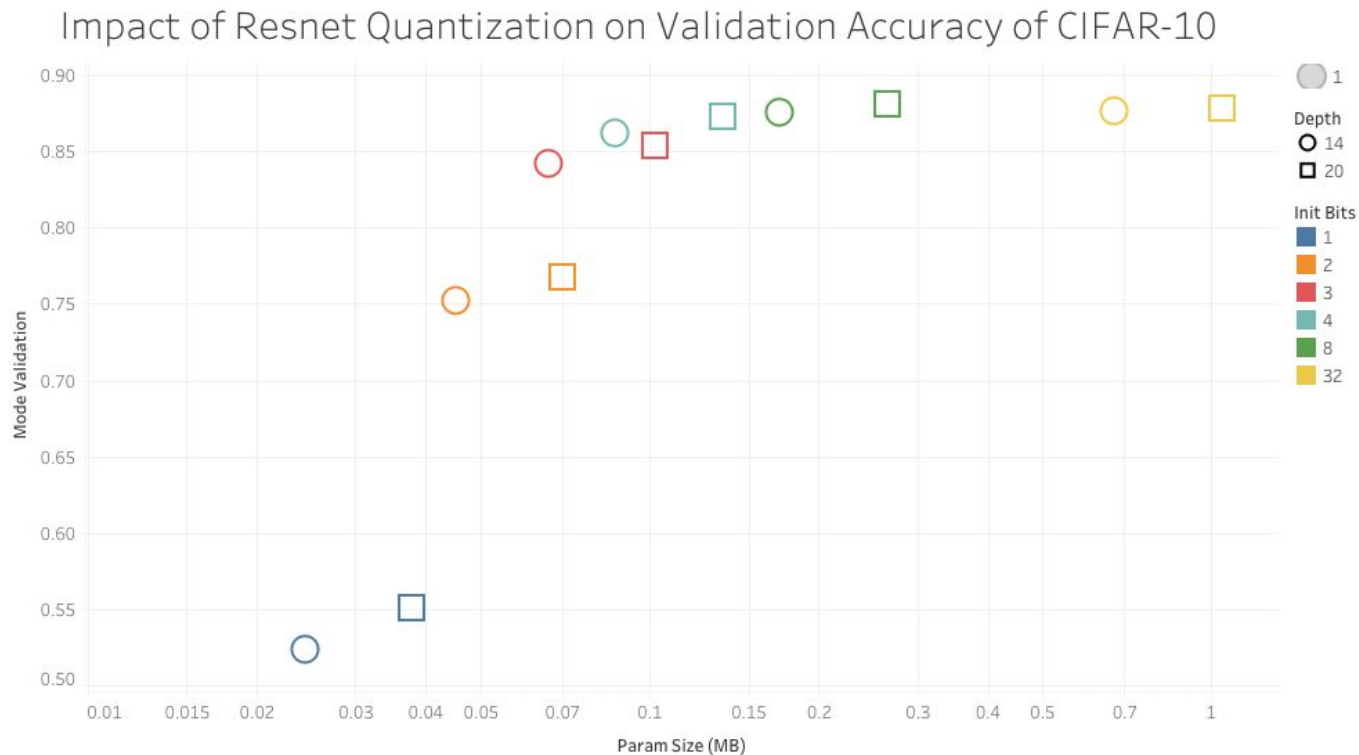
Accuracy & Latency after Decomposition

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

Quantization



Quantization



Train smaller model

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

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

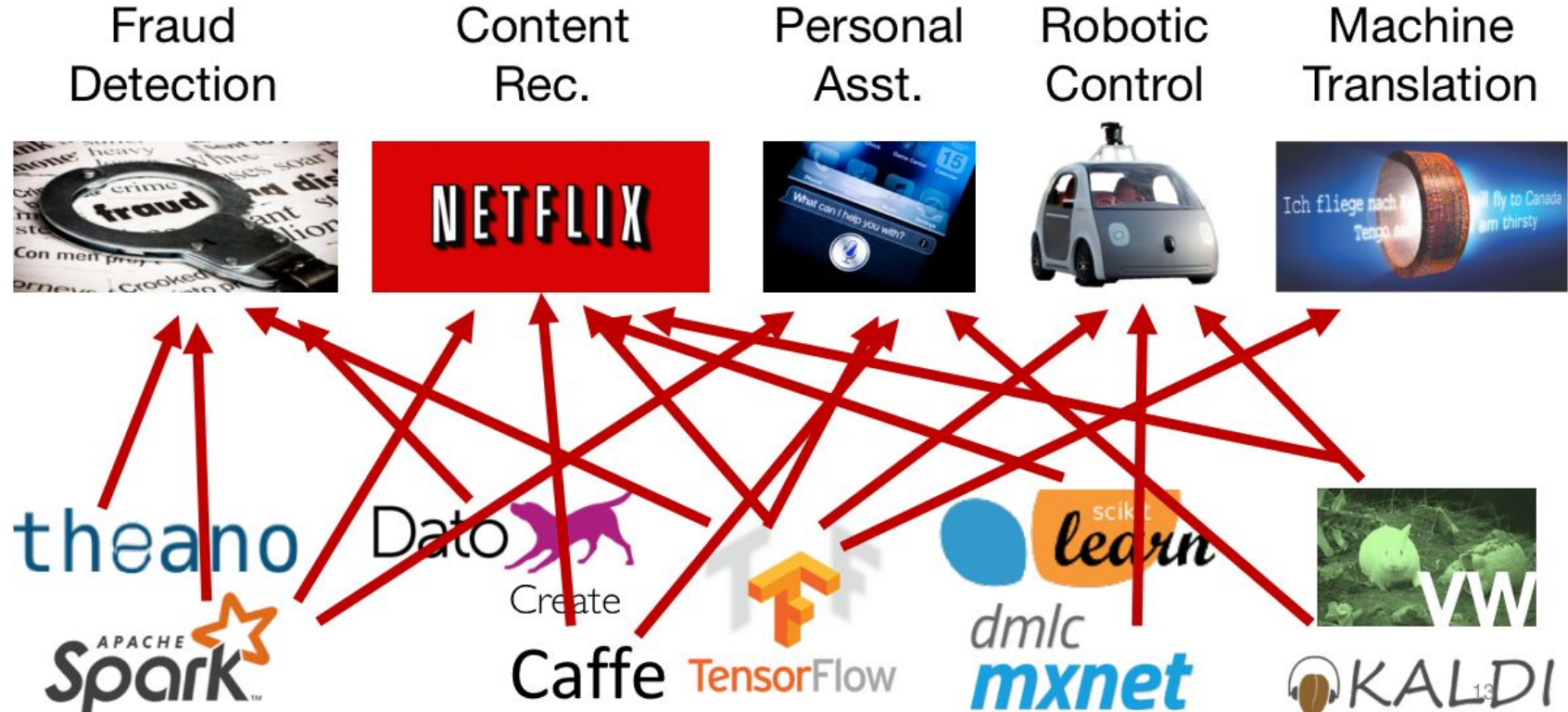
Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	64.96%
Teacher	~9M	63.54%
<i>State-of-the-art methods</i>		
Maxout		61.43%
Network in Network		64.32%
Deeply-Supervised Networks		65.43%

Table 2: Accuracy on CIFAR-100

Serving backend

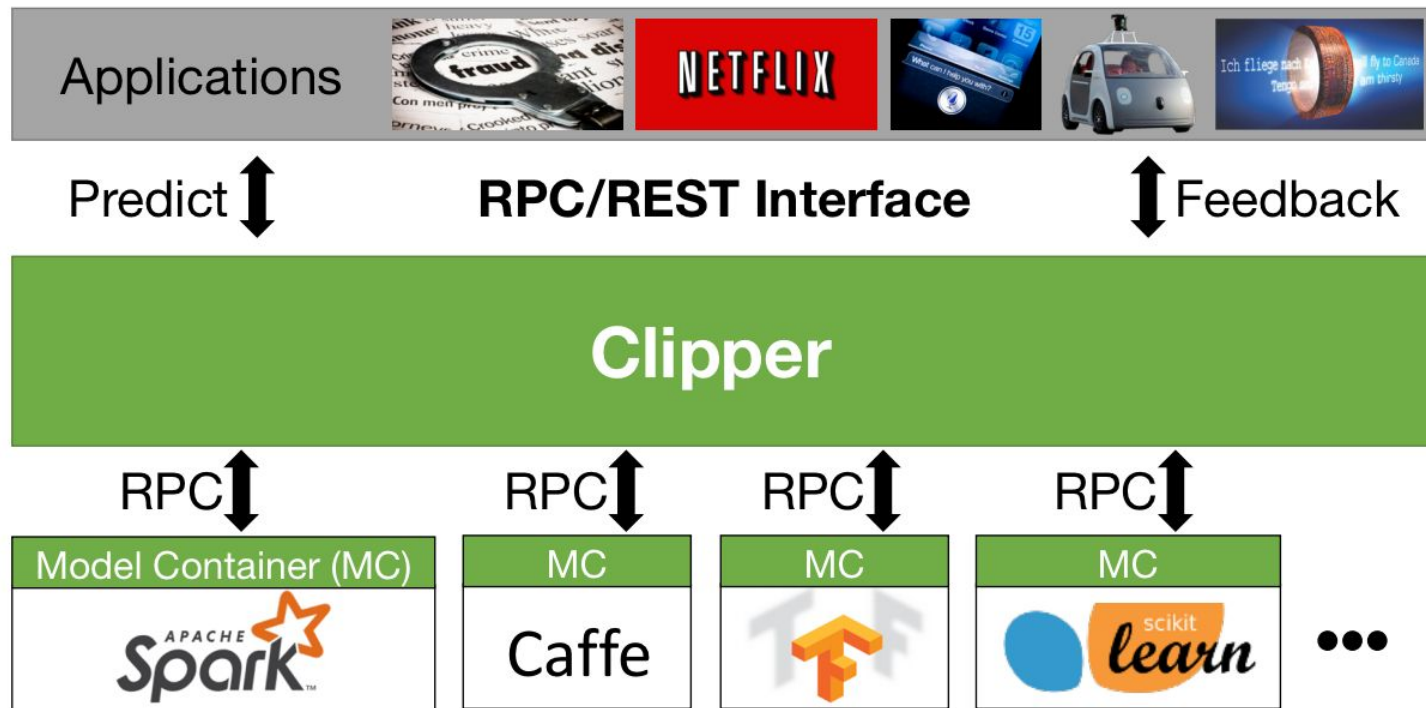
- Provide common abstraction for different frameworks
- Decide the batch size
- Load balancing and scheduling (on-going research)

Miscellaneous DL frameworks and models



* Crankshaw, Daniel, et al. "Clipper: A Low-Latency Online Prediction Serving System." presentation for NSDI (2017). https://www.usenix.org/sites/default/files/conference/protected-files/nsdi17_slides_crankshaw.pdf

Clipper Architecture



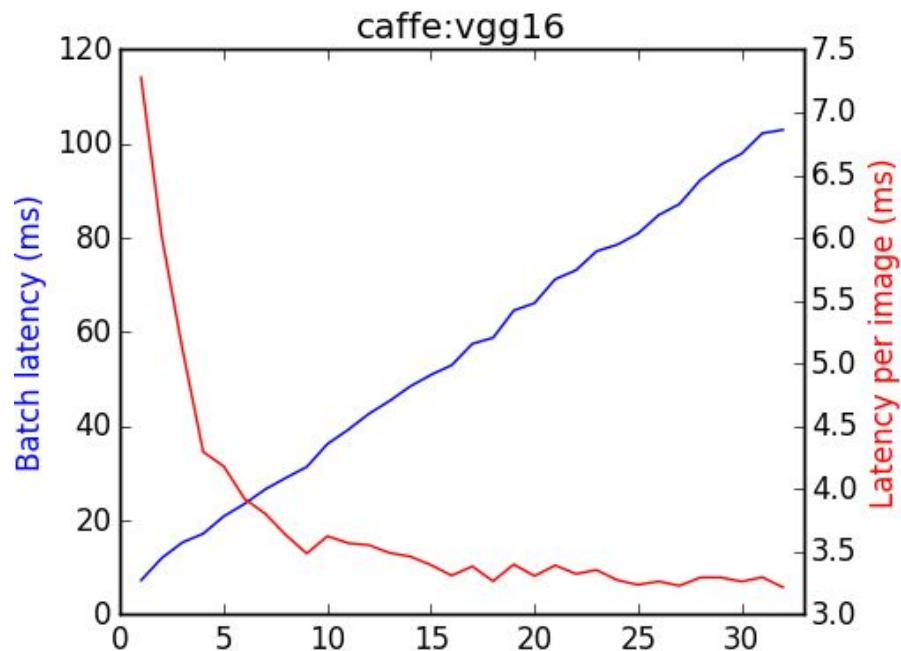
Model Abstraction

- Model container

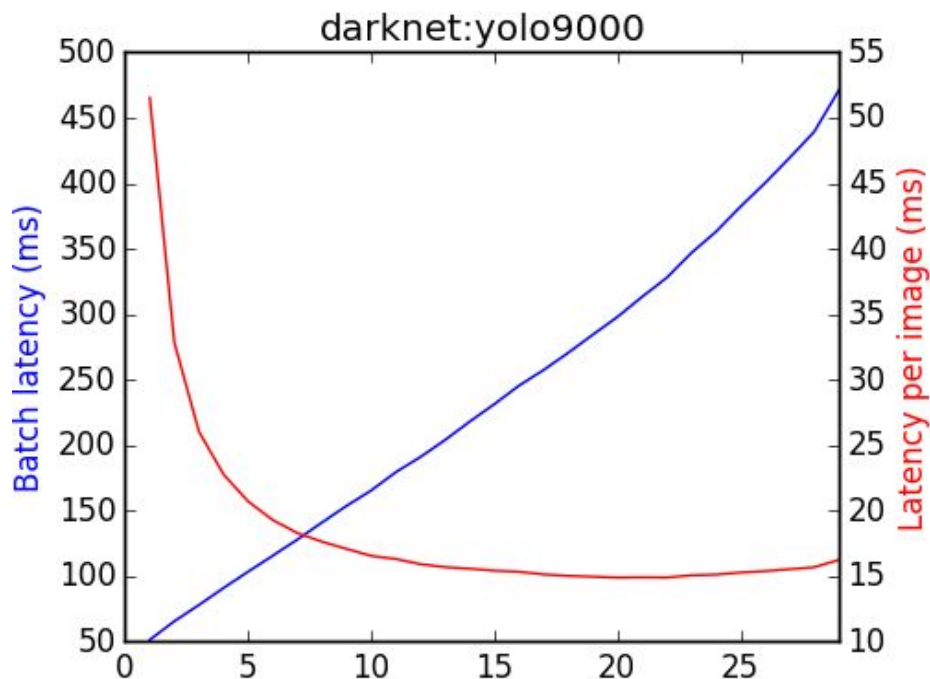
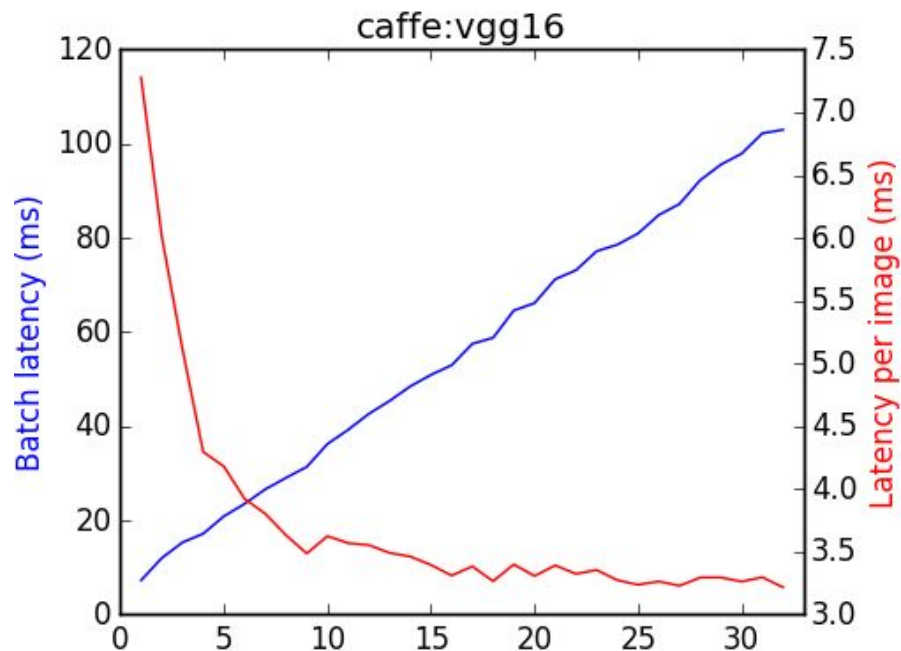
```
class ModelContainer:  
    def __init__(model_data)  
    def predict_batch(inputs)
```

- Evaluate models using original framework
- Model run in separate process as Docker containers

Batch / latency trade-off



Batch / latency trade-off



Runtime scheduling between device and cloud

- Schedule execution between device and cloud, and select approximate models
- Manage the battery and memory constraints for the device
- Manage the cost constraint for the cloud
- Goal: Maximize the overall accuracy

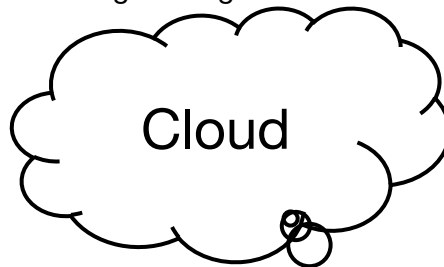
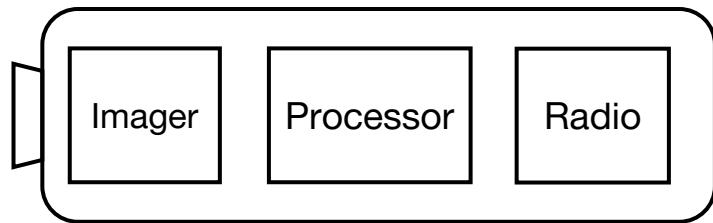
Resource usage for a continuous vision app

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

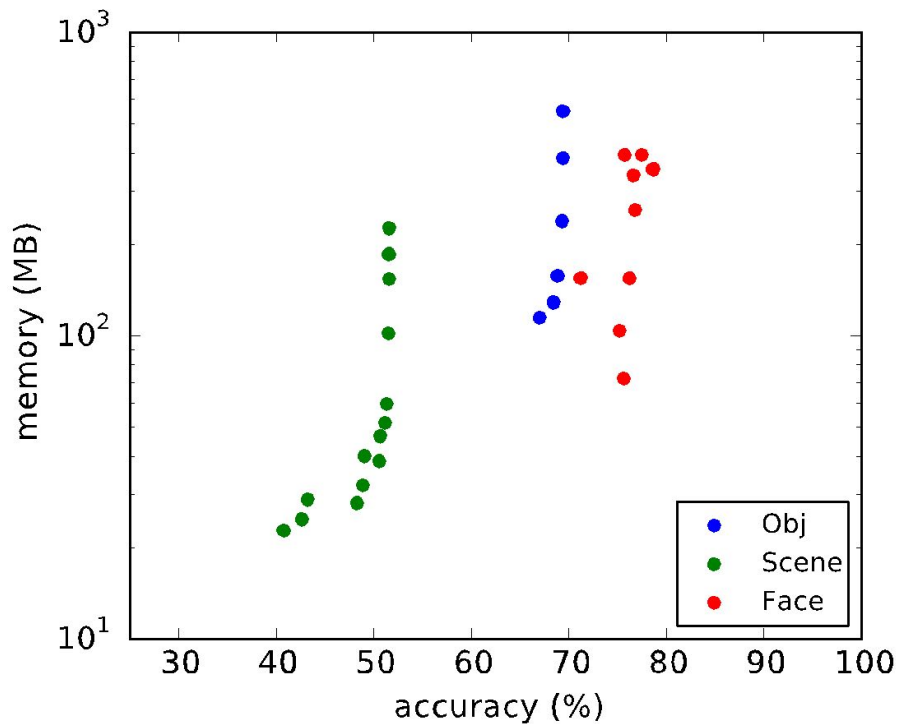
\$10 person/year

Compute power

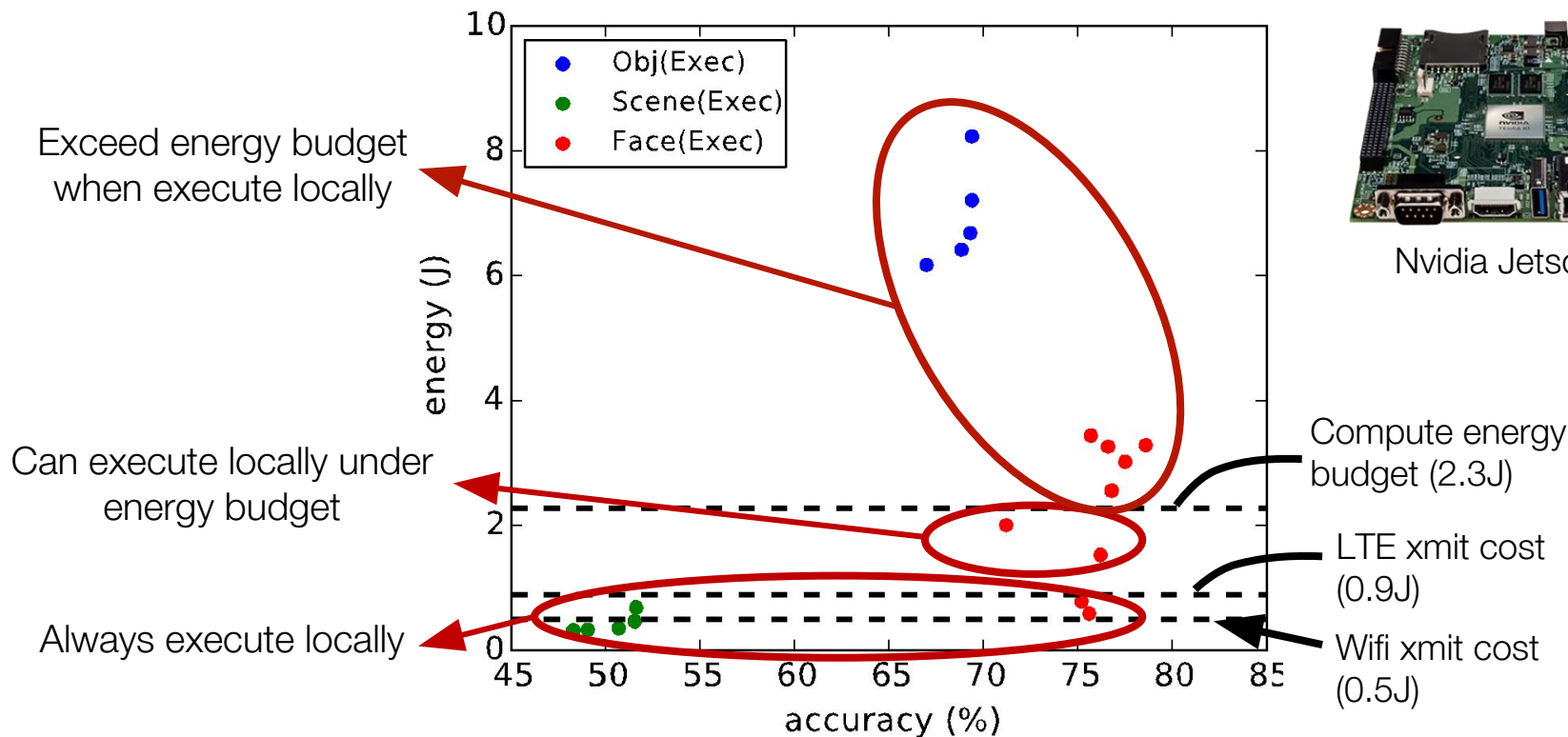
9GFLOPS

3.5GFLOPS (GPU) / 8GFLOPS (CPU)

Memory / accuracy trade-off



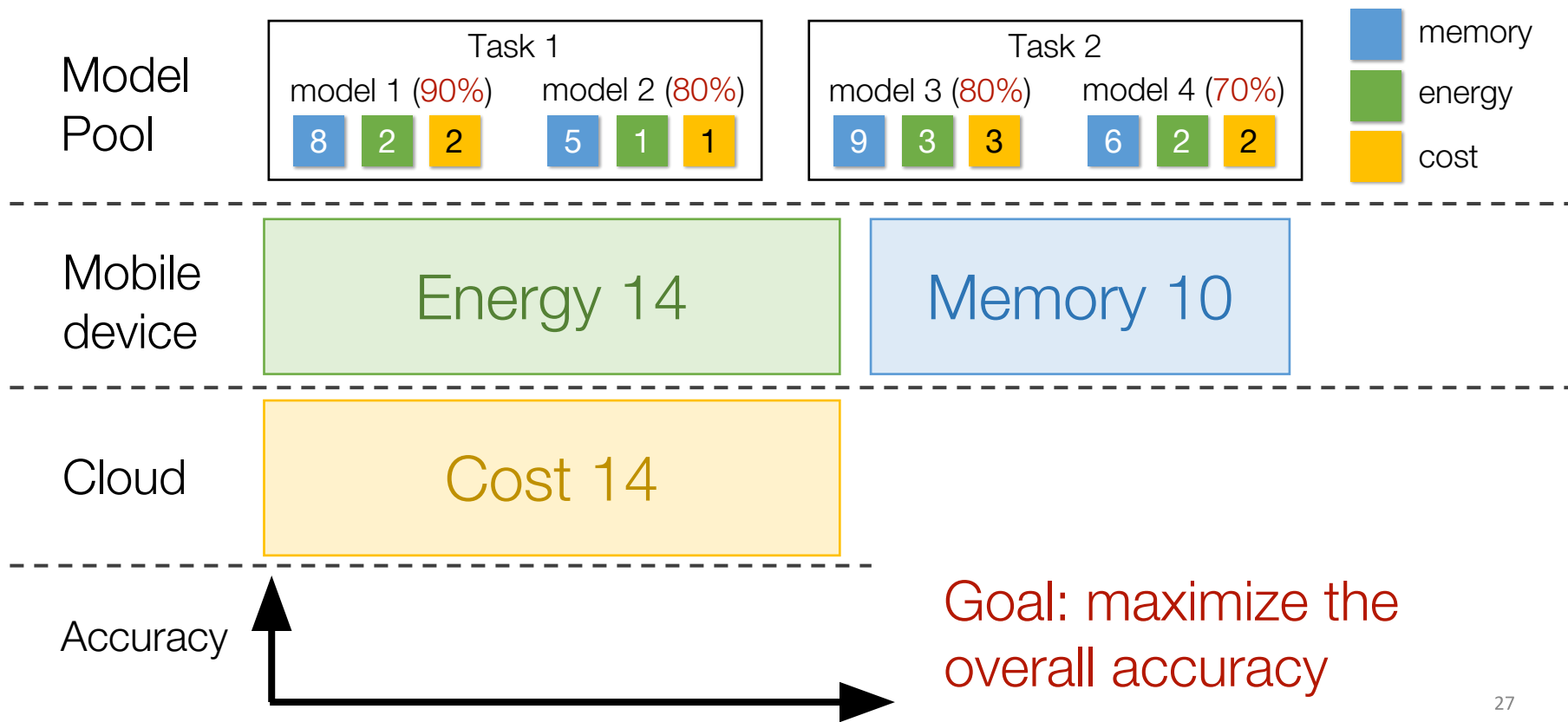
Energy / accuracy trade-off



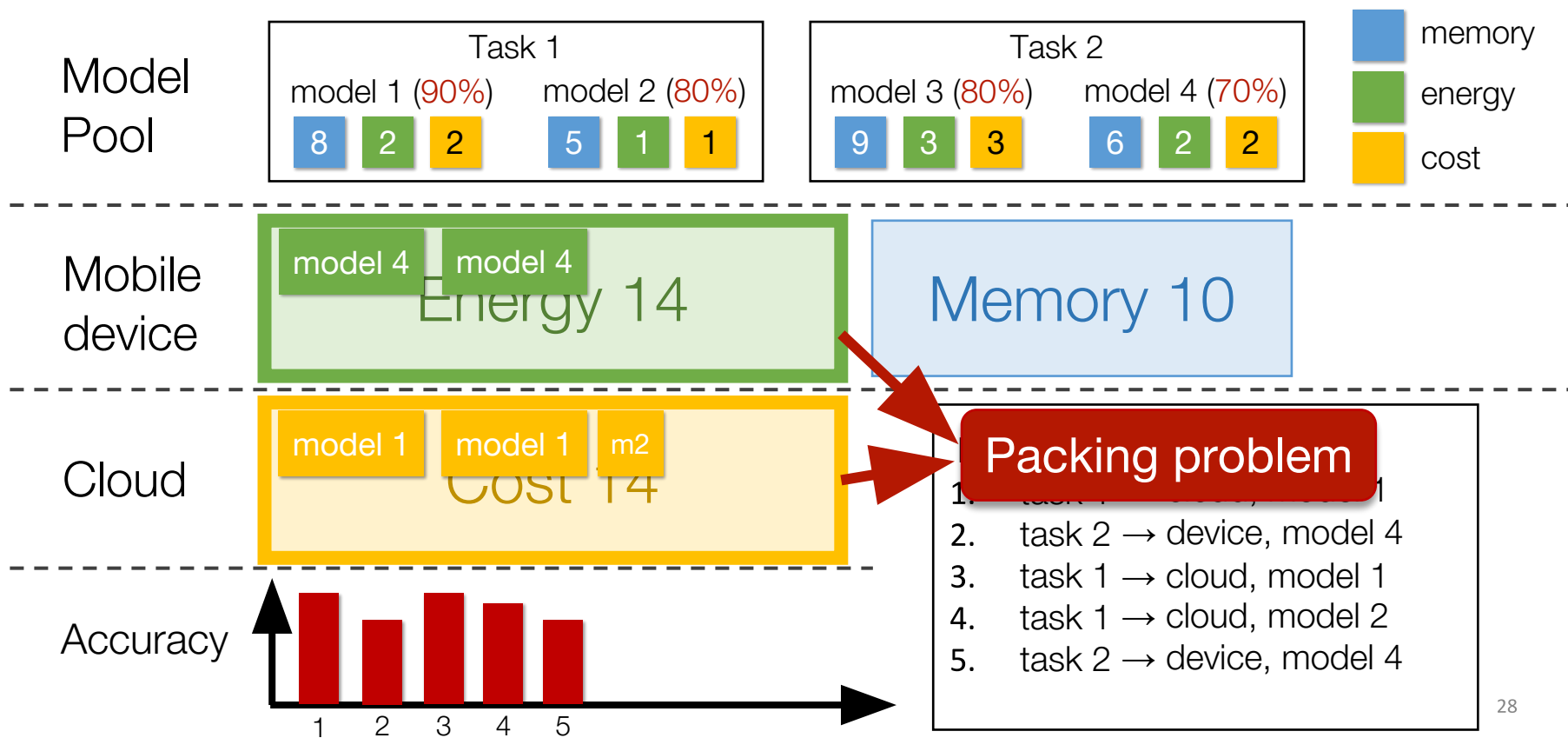
Nvidia Jetson TK1

$$\text{energy budget} = \text{total energy} / \text{total time}(10\text{h}) / \text{requests per second}(1 \text{ req/sec})$$

Approximate model scheduling

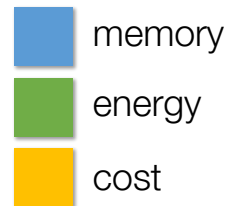
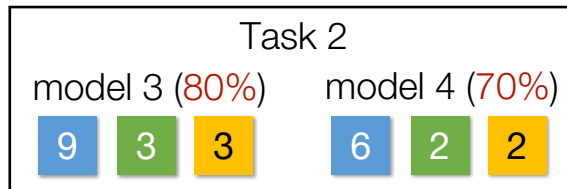
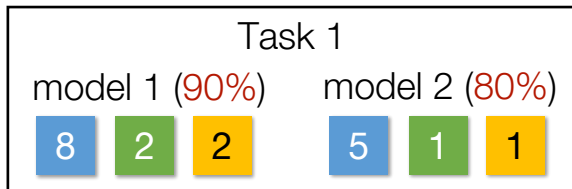


Approximate model scheduling

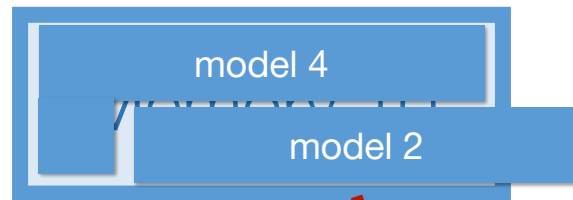
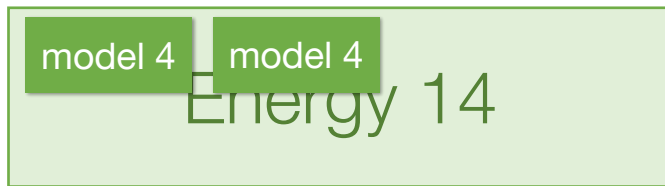


Approximate model scheduling

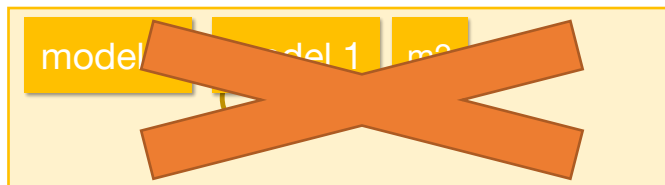
Model
Pool



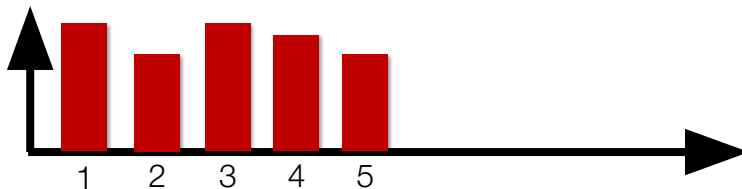
Mobile
device



Cloud



Accuracy



Requests:

1. task 1 → cloud, model 1
- 2.
- 3.
4. task 1 → cloud, model 2
5. task 2 → device, model 4
6. task 1

Paging problem

Approximate model scheduling

- Packing problem: pick versions that satisfy energy/cost budgets

$$\sum_t e_i x_{it} \leq E, \sum_t c_i x'_{it} \leq C \quad (x_{it}, x'_{it} \in [0,1], x_{it} \cdot x'_{it} = 0)$$

- Paging problem: pick versions that fit in memory

$$\forall 1 \leq t \leq T, \sum_{i=1}^n s_i x_{it} \leq S$$

- Goal: maximize the accuracy

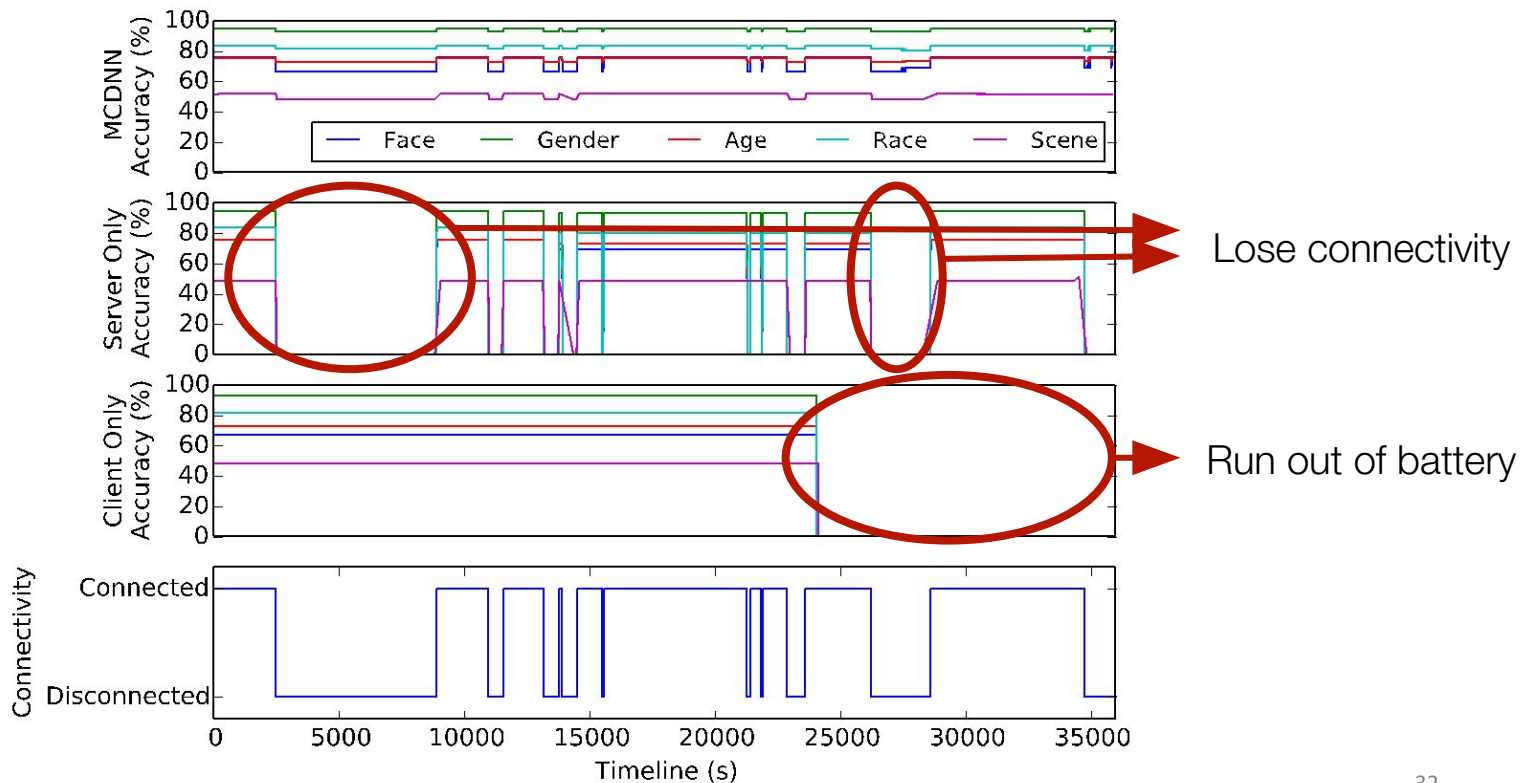
$$\max_x \sum_t \sum_i a_i (x_{it} + x'_{it})$$

No known optimal online algorithms

Heuristic scheduler

- Estimate future resource use and compute the budget for each request
- Account for paging cost to reduce oscillations
- Use increasingly more accurate versions of more heavily used models

Trace-driven evaluation



Reference

- Kim, Yong-Deok, et al. "Compression of deep convolutional neural networks for fast and low power mobile applications." ICLR (2016).
- Han, Seungyeop, et al. "MCDNN: An Approximation-Based Execution Framework for Deep Stream Processing Under Resource Constraints." MobiSys (2016).
- Romero, Adriana, et al. "Fitnets: Hints for thin deep nets." ICLR (2015).
- Crankshaw, Daniel, et al. "Clipper: A Low-Latency Online Prediction Serving System." NSDI (2017).