# MCDNN: An Approximation-Based Execution Framework for Deep Stream Processing Under Resource Constraints
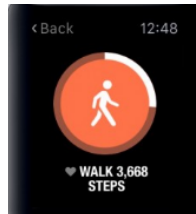
Haichen Shen (haichen@cs.washington.edu)

with Seungyeop Han, Matthai Philipose,

Sharad Agarwal, Alec Wolman, Arvind Krishnamurthy

University of Washington                    Microsoft Research

# Wearable computing → more data



???

# When computer vision meets wearable

"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"



**Life Sciences Industry Report**
Commentary on current industry topics

**Reducing human error in pharmaceutical manufacturing**

Few experts would dispute that human error is the cause of most pharmaceutical manufacturing failures. **Some estimate it to be as high as 80 percent!**

While there are many types of workflow software available, those which allow embedded workflow design are the most effective. These enable seamless integration of document, industrial process, human activity, and operations management workflow via a single user interface. Having to leave one application and open one or more others to resolve any situation, jeopardizes productivity and timeliness, it also makes the process more susceptible to error.

**GOVERNING**
THE STATES AND LOCALITIES
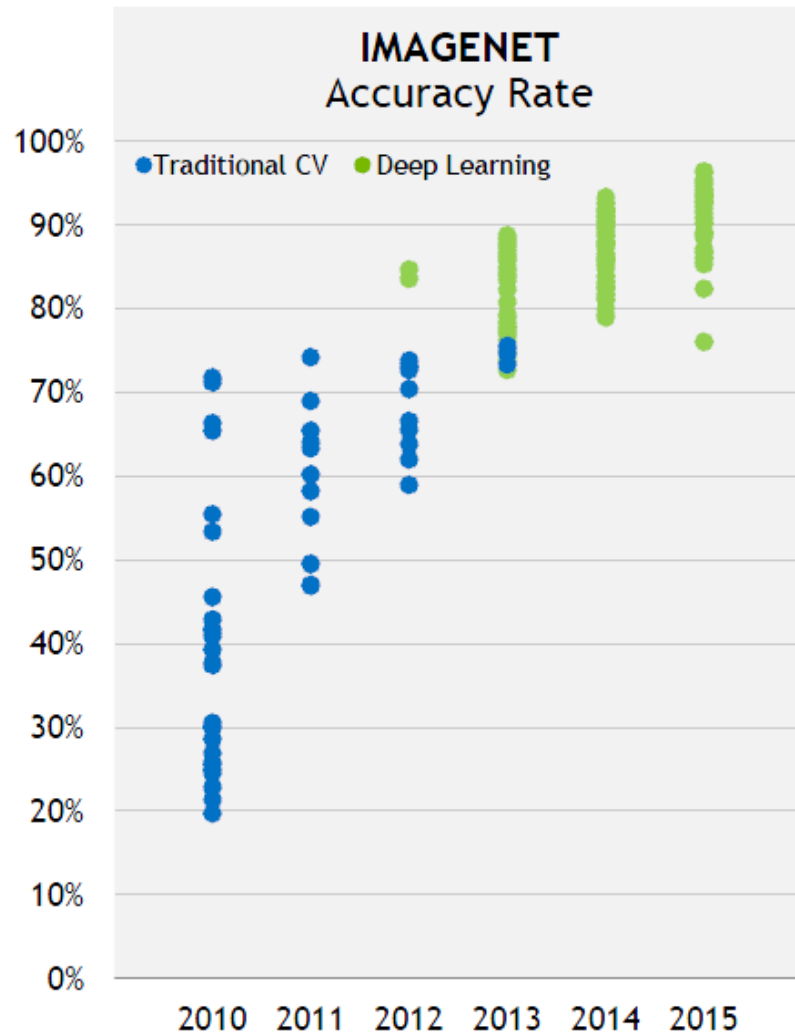
**PUBLIC SAFETY & JUSTICE**

**Can Body Cameras Really Reduce Ferguson Police's Use of Force?**

Ferguson police are the latest of more than 1,000 departments to wear body cameras, which are proven to reduce officers' use of force and citizens' complaints against cops.

BY TOD NEWCOMBE | SEPTEMBER 4, 2014

Consumer

Manufacturing

Public Safety

# Deep learning makes vision work



IMAGENET
Accuracy Rate

But...

| Recognition Task | face | scene* | object* |
|---|---|---|---|
| Accuracy | 97% | 88% | 92% |
| Compute/frame (FLOPs) | 1.00G | 30.9G | 39.3G |
| Compute@1-30fps (FLOPS) | 1-30G | 30-900G | 40G-1.2T |

Do we have enough resources
to run deep learning?

* top-5 accuracy is shown in the table

# Resource usage for continuous vision

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

| Imager | Processor | Radio |

Cloud

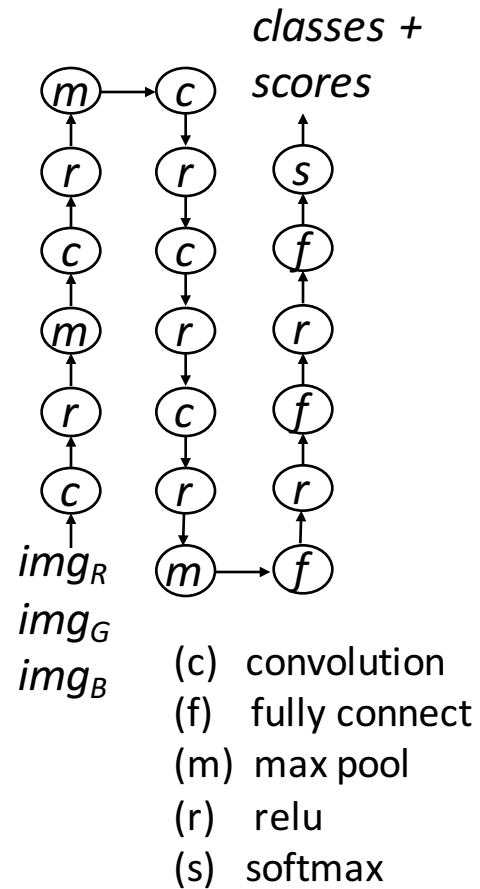| | Device | Cloud |
|---|---|---|
| **Workload** | Deep learning **300GFLOPS** @ 30GFLOPs/frame, 10fps | |
| **Budget** | Device power<br>30% of 10Wh for 10h = 300mW | Cloud cost<br>$10 person/year |
| **Compute power** | **9GFLOPS** | **3.5GFLOPS** (GPU) / **8GFLOPS** (CPU) |

Huge gap between workload and budget

# Neural network

*classes +*
*scores*

```
m ──→ c

r     r     s

c     c     f

m     r     r

r     c     f

c     r     r

imgR  m ──→ f
imgG
imgB
```

(c) convolution
(f) fully connect
(m) max pool
(r) relu
(s) softmax

# Neural network ≈ matrix multiplications



classes + scores

$img_R$
$img_G$
$img_B$

(c)  convolution
(f)  fully connect
(m)  max pool
(r)  relu
(s)  softmax

classes + scores

$img_R$
$img_G$
$img_B$

Architectural changes
(J. Ba, et al. 2014)

X

Low rank approximation
(Y. Kim, et al. 2016)

X

Matrix sparsification
(S. Han, et al. 2015)

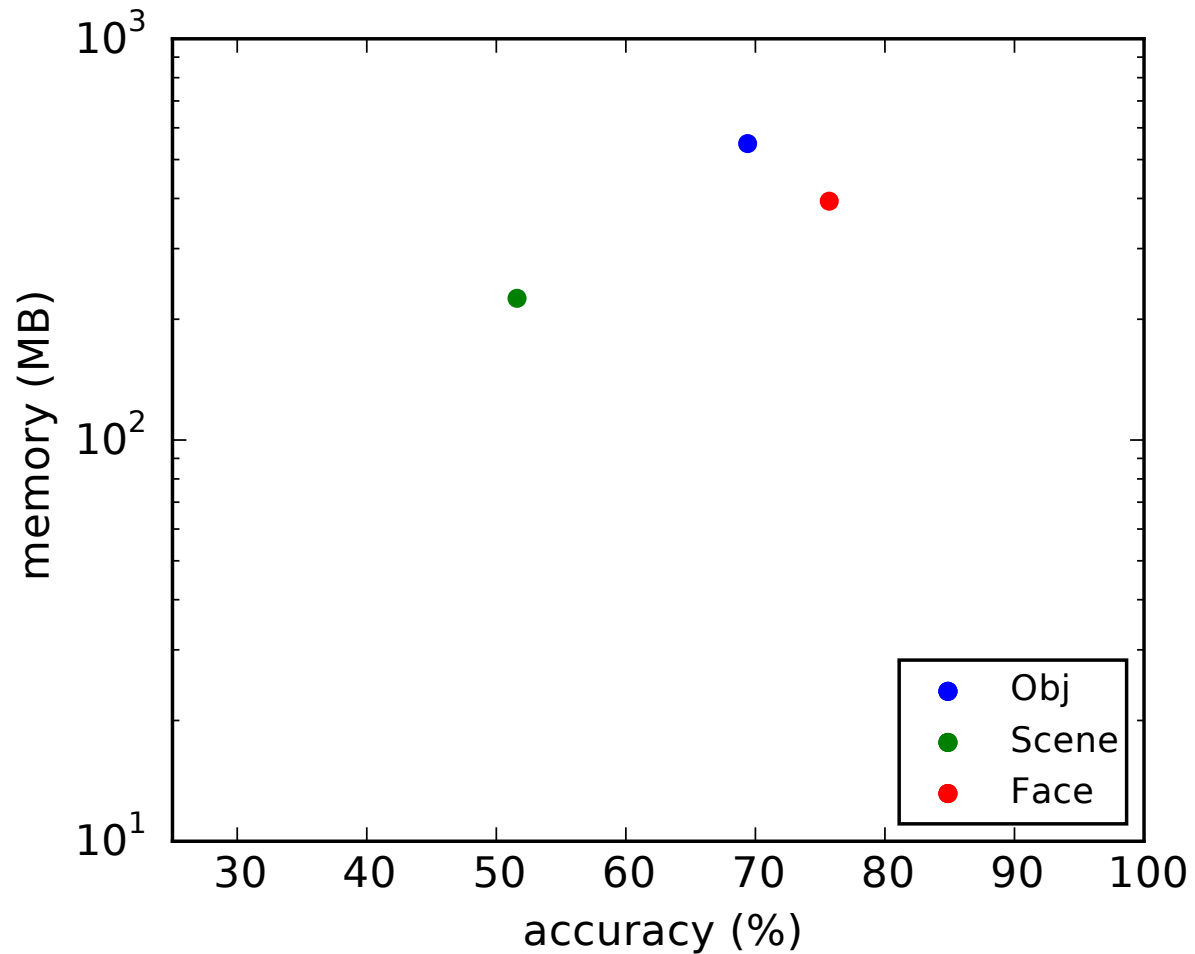# Managing the approx. / resource trade-off

▸ Detailed characterization of the approximation / resource trade-off for many optimizations

▸ Two new optimizations for streaming, multi-application settings

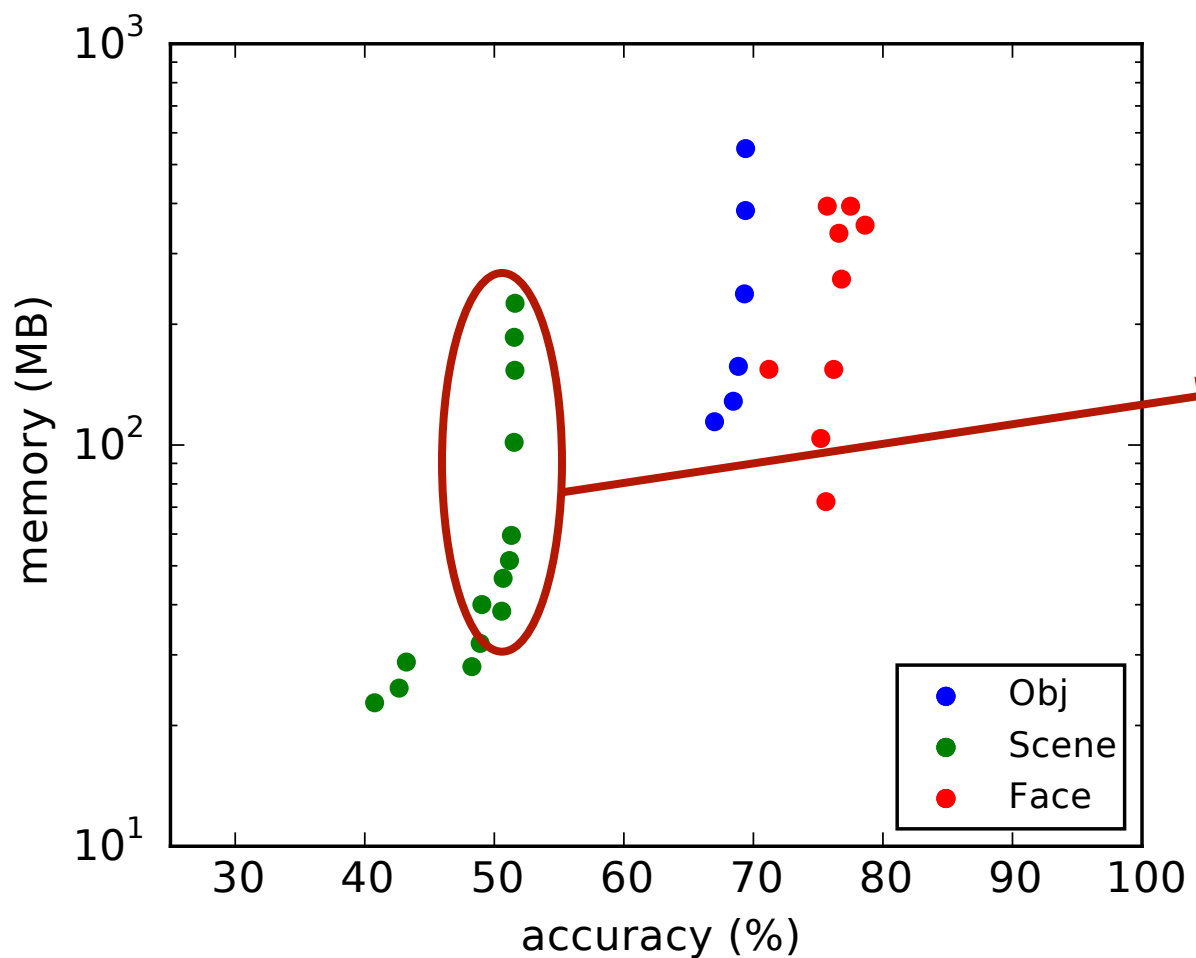▸ New scheduling problem, Approximate Model Scheduling, with a heuristic solution

# Outline

▸ Detailed characterization of the approximation / resource trade-off for many optimizations

▸ Two new optimizations for streaming, multi-application settings

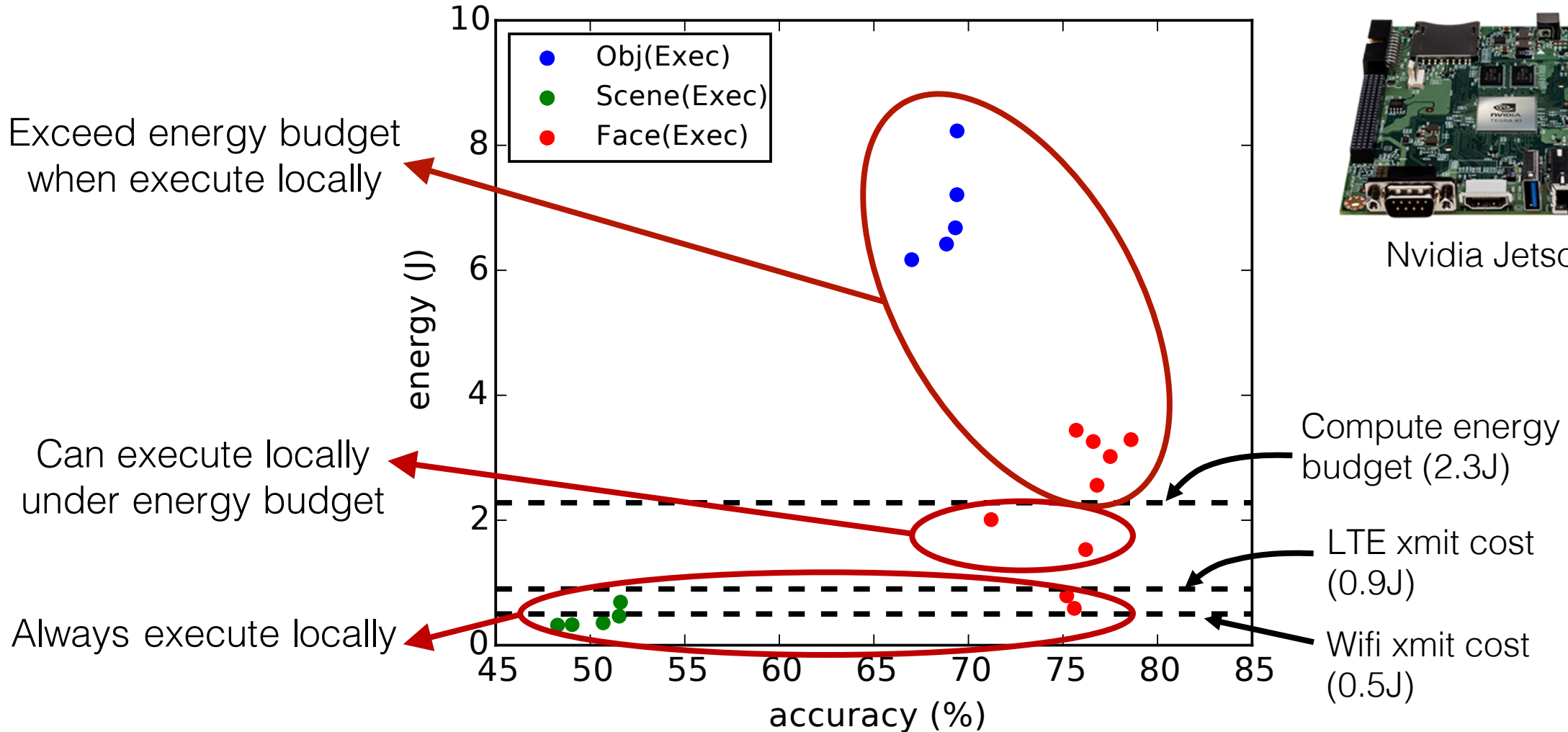▸ New scheduling problem, Approximate Model Scheduling, with a heuristic solution

# Memory / accuracy trade-off

# Memory / accuracy trade-off



Substantially reduce memory use with gradual accuracy loss

# Energy / accuracy trade-off



Nvidia Jetson TK1

Exceed energy budget
when execute locally

Can execute locally
under energy budget

Always execute locally

Compute energy
budget (2.3J)

LTE xmit cost
(0.9J)

Wifi xmit cost
(0.5J)

Legend: Obj(Exec), Scene(Exec), Face(Exec)

energy budget = total energy / total time(10h) / requests per second(1 req/sec)

12

# Outline

‣ Detailed characterization of the approximation / resource trade-off for many optimizations

‣ Two new optimizations for streaming, multi-application settings
  ‣ Specialization
  ‣ Model sharing

‣ New scheduling problem, Approximate Model Scheduling, with a heuristic solution

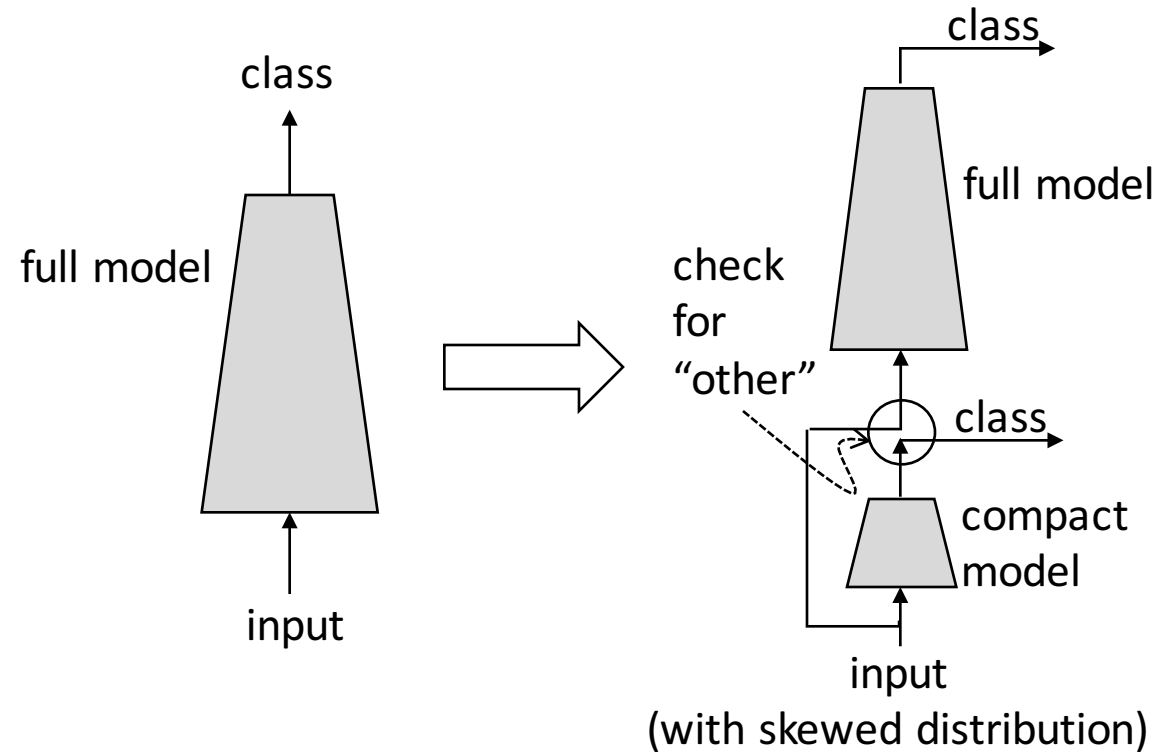# Exploiting stream locality by specialization

- Standard deep neural network recognizes 4000 people
- Most of videos are dominated by less than 10 faces over minutes
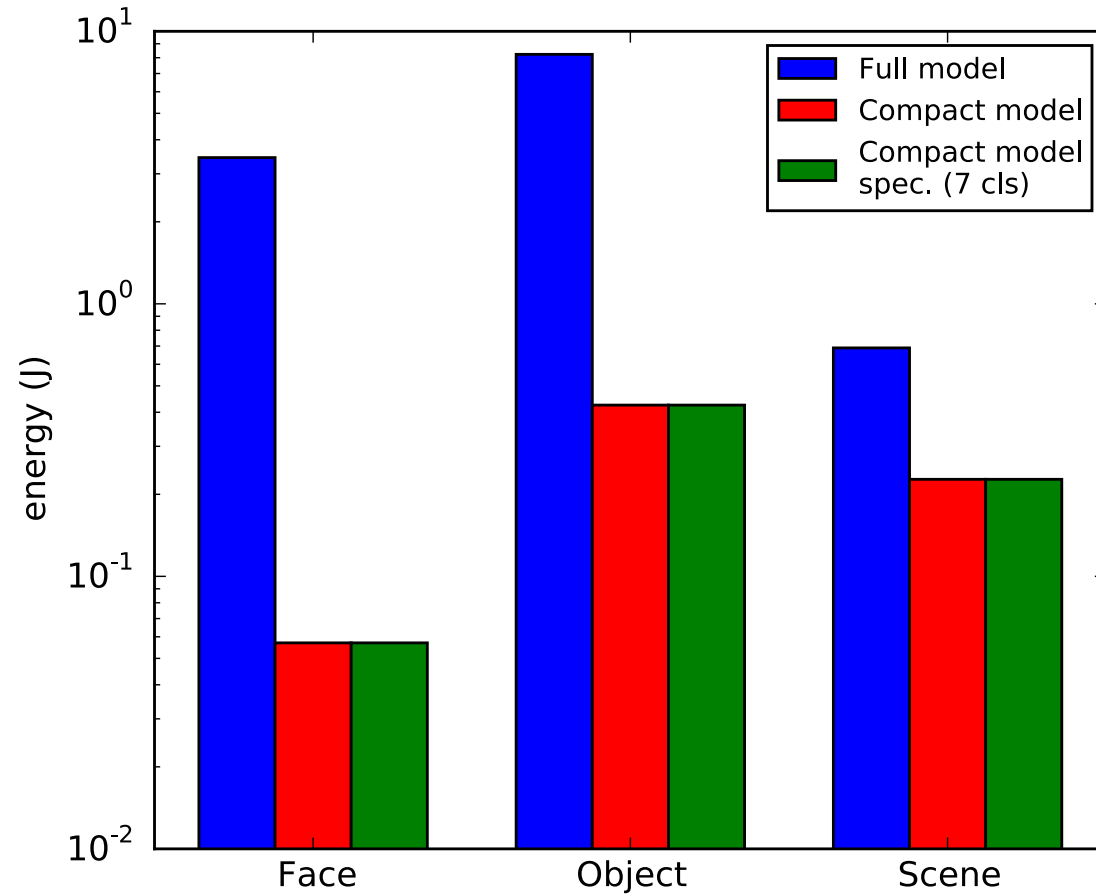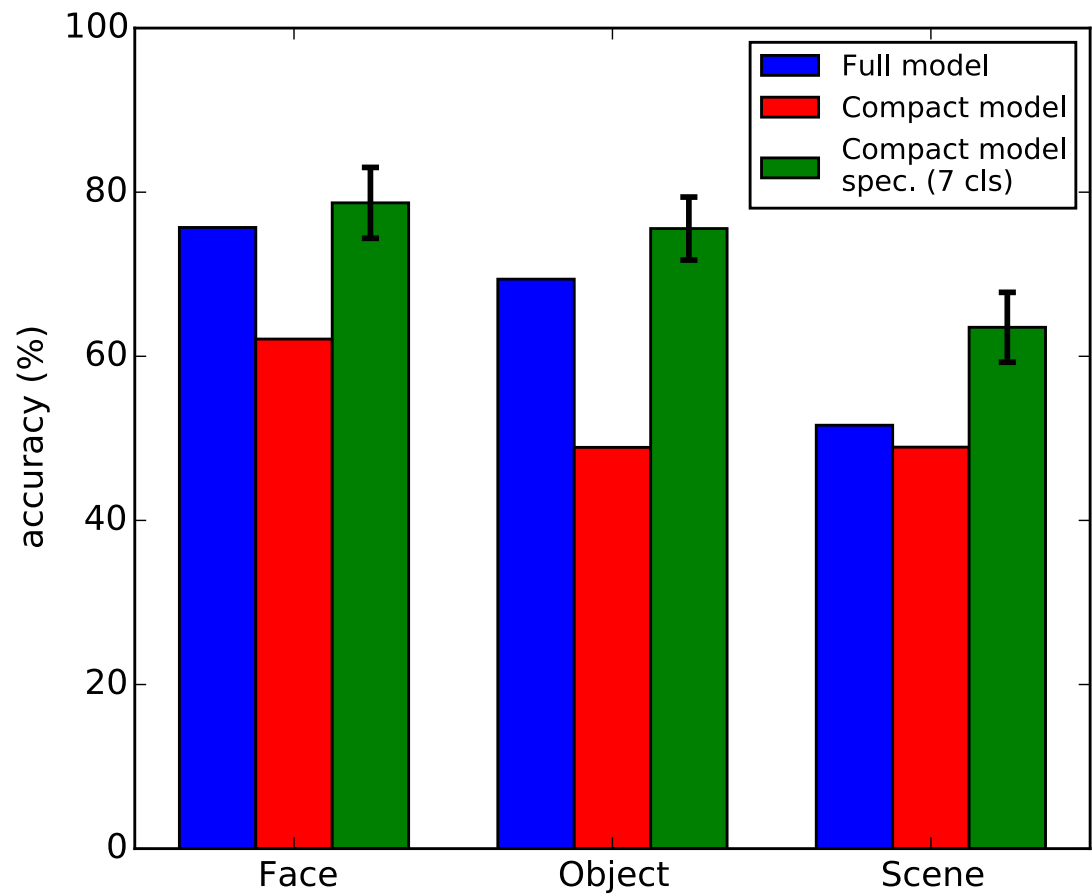
Timeline



**Produce more compact models for skewed classes**

# Specialization runtime

# Better resource/accuracy trade-off

# Outline

‣ Detailed characterization of the approximation / resource trade-off for many optimizations

‣ Two new optimizations for streaming, multi-application settings

‣ **New scheduling problem, Approximate Model Scheduling, with a heuristic solution**

# Approximate model scheduling

# Approximate model scheduling

**Model Pool**

Task 1
model 1 (90%)    model 2 (80%)
8  2  2          5  1  1

Task 2
model 3 (80%)    model 4 (70%)
9  3  3          6  2  2

- memory
- energy
- cost

**Mobile device**

model 4    model 4    Energy 14          Memory 10

**Cloud**

model 1    model 1    m2    Cost 14

**Packing problem**

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

**Accuracy**

1  2  3  4  5

# Approximate model scheduling



**Model Pool**

Task 1
model 1 (90%)　　model 2 (80%)
8　2　2　　　5　1　1

Task 2
model 3 (80%)　　model 4 (70%)
9　3　3　　　6　2　2

memory
energy
cost

**Mobile device**

model 4　model 4
Energy 14
model 4
model 2

**Cloud**

model 1　model 1　m2

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**

**Accuracy**

1　2　3　4　5

# Approximate model scheduling

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

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

- 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

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

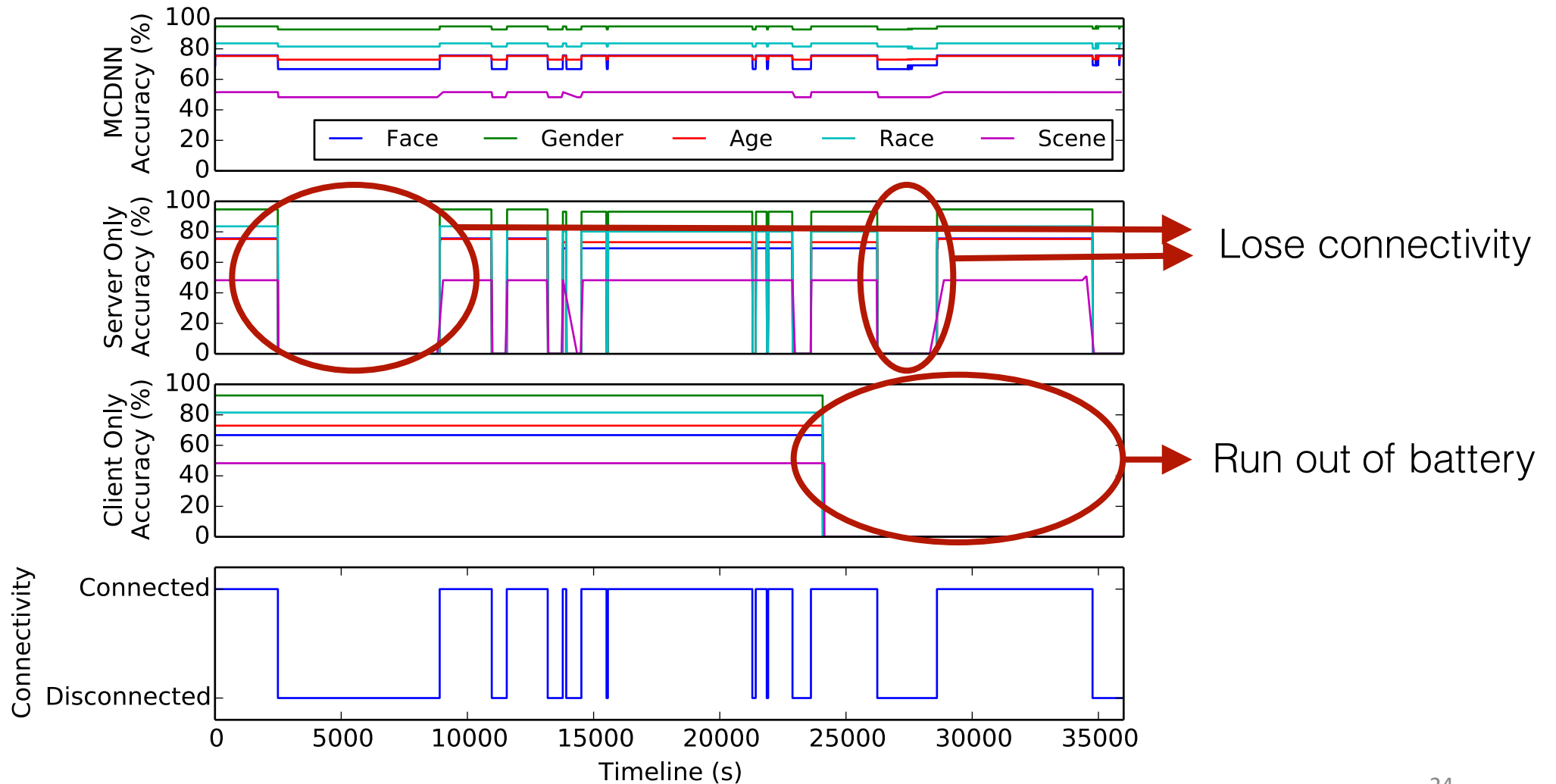No known optimal online algorithms

# 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

# MCDNN framework

input type
model schema
training/validation data

**compiler**

trained
model
catalog

development time

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

run time

input → **device runtime**
scheduler
data router

input →

**cloud runtime**
scheduler
data router
profiler

← classes →

classes ↓

apps

device

cloud

# MCDNN framework



input type
model schema
training/validation data

**compiler**

trained
model
catalog

development time

specialization time

**specializer**

run time

specialized models

stats

input

**device runtime**
**scheduler**
data router

input

**cloud runtime**
**scheduler**
data router
profiler

classes

classes

apps
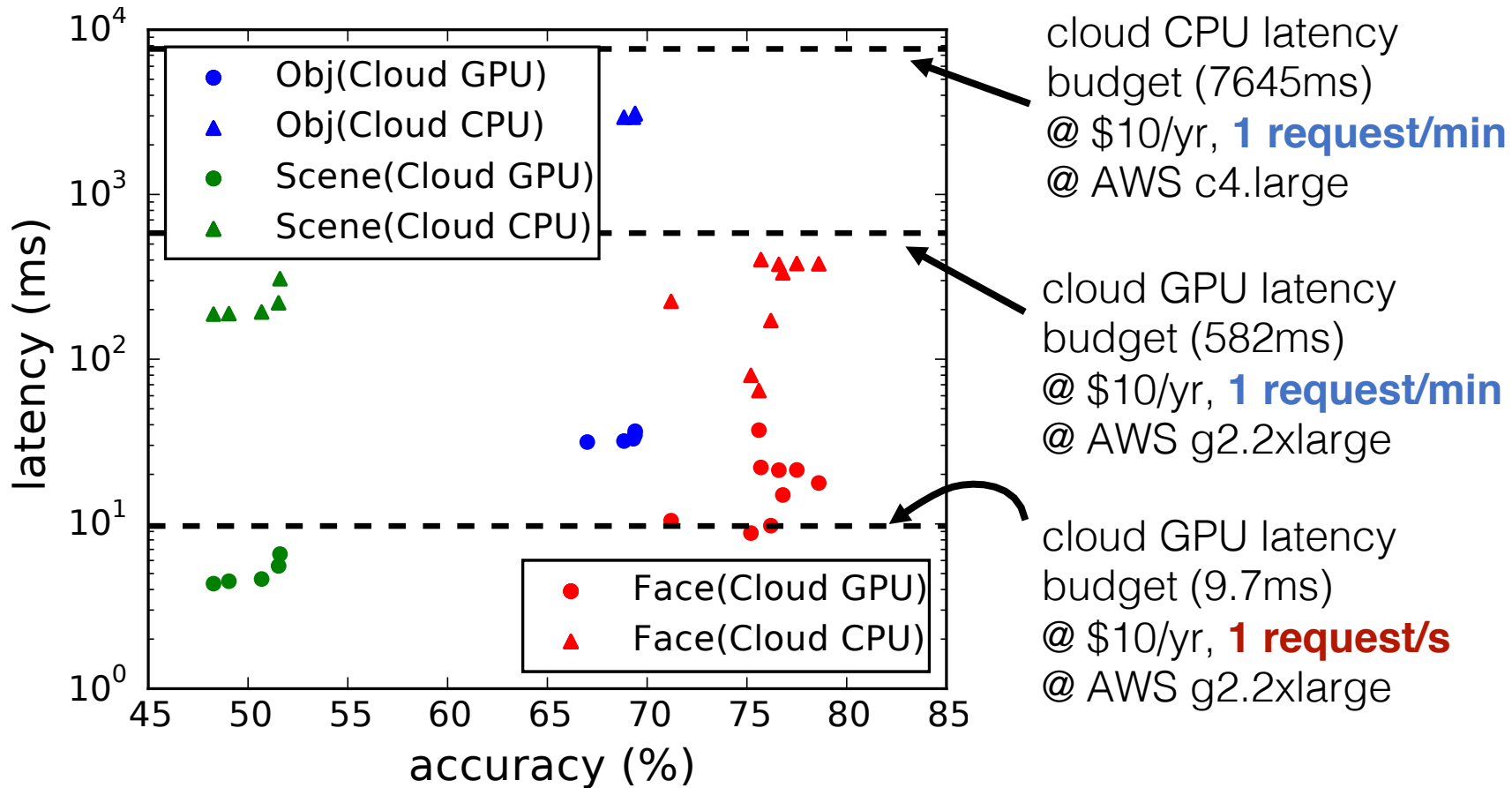
device

cloud

# Conclusion

- MCDNN makes efficient trade-offs between resource use and accuracy

- Formulate the approximate model scheduling problem and devise a heuristic algorithm

- Design a generic approximation-based execution framework for continuous mobile vision

## Thank you! Questions?

# Backup Slides

# Cloud cost / accuracy trade-off



cloud CPU latency budget (7645ms) @ $10/yr, **1 request/min** @ AWS c4.large

cloud GPU latency budget (582ms) @ $10/yr, **1 request/min** @ AWS g2.2xlarge

cloud GPU latency budget (9.7ms) @ $10/yr, **1 request/s** @ AWS g2.2xlarge

Legend:
- Obj(Cloud GPU)
- Obj(Cloud CPU)
- Scene(Cloud GPU)
- Scene(Cloud CPU)
- Face(Cloud GPU)
- Face(Cloud CPU)

latency budget = cost budget / cost per hour / #requests

# Model sharing



face
ID
race
age
gender
input
intermediate
values

model-fragment cache
...
router
input

# Dynamically-sized caching scheme



% Requests Serviced

Energy Consumed/Request

Average Accuracy

fraction of requests served

energy/request (J)

accuracy/request (%)

cache load rate (MB/request)

Original Model
Best Model
All Models