# Comprehensive and Efficient Data Labeling via Adaptive Model Scheduling

**Mu Yuan[1], Lan Zhang[1], Xiang-Yang Li[1], Hui Xiong[2]**

**[1]University of Science and Technology of China**

**[2]Rutgers University**

# Outline

- **Resource-wasting multi-model inference workloads**
- Rule-based scheduler
- Learning-based scheduler
- Evaluation
- Conclusion

walking dogs

**Image Retrieval**
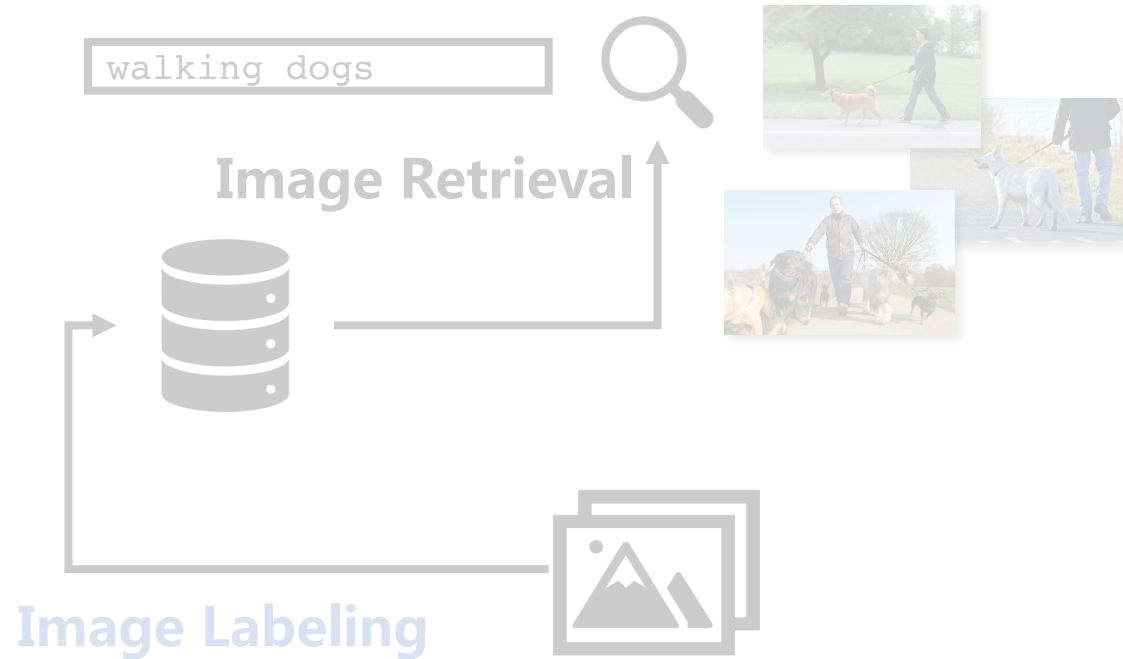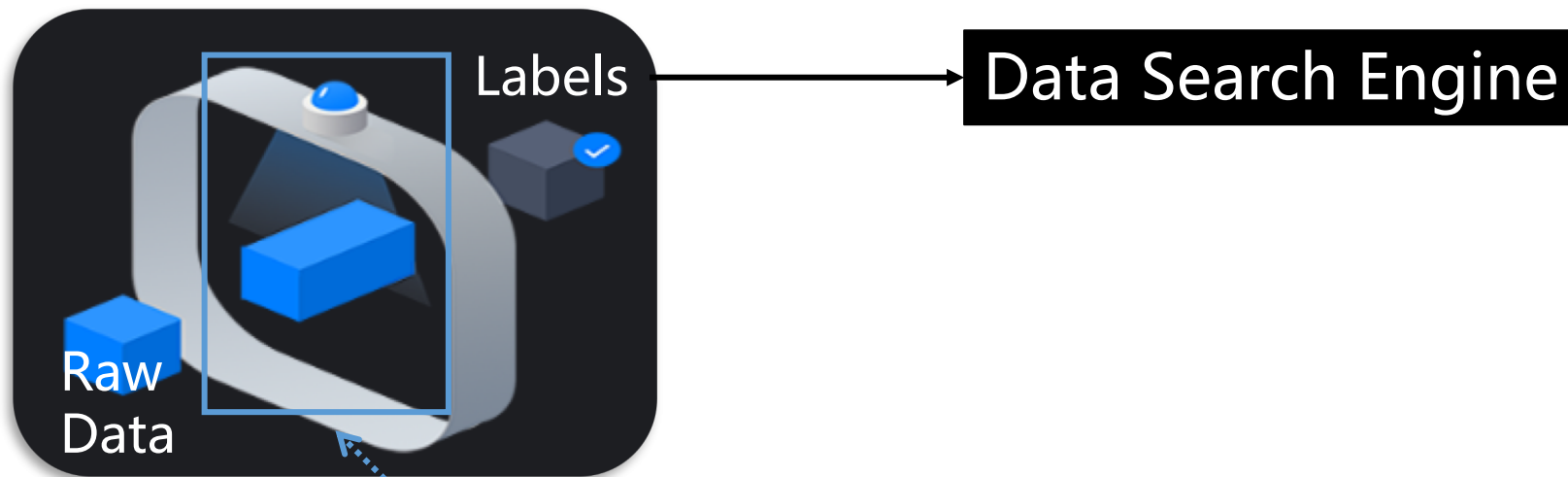
**Image Labeling**

walking dogs

**Image Retrieval**

**Image Labeling**

ILLEGAL

20 mph

Degree of Congestion:
**HIGH**

# Data Trading Platform



Labels → Data Search Engine

Raw Data

**Multi-Model Data Labeling Workloads**

| | | |
|---|---|---|
| **Object Detector** | → | person, 0.994<br>chair, 0.565<br>tv monitor, 0.996 |
| **Scene Classifier** | → | pub, 0.727<br>beer hall, 0.198 |
| **...** | → | ... |

**CV Models**

**Raw Images**

| | | | | | |
|---|---|---|---|---|---|
| **Pose Estimator** | | Body Keypoints | | | Body Keypoints |
| **Face Detector** | | | Face Location | | |
| **Object Detector** | Dog (0.96) | Person (0.43) | Person (0.96) | | Bike (0.97) |
| **Action Classifier** | | Fall Down (0.87) | Make Up (0.9) | | Ride Bike (0.92) |
| **Scene Classifier** | Lawn (0.85) | Lobby (0.91) | Bathroom (0.14) | Mall (0.89) | Mountain (0.75) |
| **Dog Classifier** | Akita (0.91) | | | | |

☐ *Useful Output*     ☐ *No Output*     ☐ *Low-Confidence Output*

| | | | | | |
|---|---|---|---|---|---|
| |  |  |  |  |  |
| **Pose Estimator** | | Body Keypoints | | | Body Keypoints |
| **Face Detector** | | | Face Location | | |
| **Object Detector** | Dog (0.96) | Person (0.43) | Person (0.96) | | Bike (0.97) |
| **Action Classifier** | | Fall Down (0.87) | Make Up (0.9) | | Ride Bike (0.92) |
| **Scene Classifier** | Lawn (0.85) | Lobby (0.91) | Bathroom (0.14) | Mall (0.89) | Mountain (0.75) |
| **Dog Classifier** | Akita (0.91) | | | | |

**Wasted computing resources!**

| | | |
|---|---|---|
| *Useful Output* | *No Output* | *Low-Confidence Output* |

**7**

# Observation



| | | | | | |
|---|---|---|---|---|---|
| **Pose Estimator** | | Body Keypoints | | | Body Keypoints |
| **Face Detector** | | | Face Location | | |
| **Object Detector** | Dog (0.96) | Person (0.43) | Person (0.96) | | Bike (0.97) |
| **Action Classifier** | | Fall Down (0.87) | Make Up (0.9) | | Ride Bike (0.92) |
| **Scene Classifier** | Lawn (0.85) | Lobby (0.91) | Bathroom (0.14) | Mall (0.89) | Mountain (0.75) |
| **Dog Classifier** | Akita (0.91) | | | | |

**Wasted computing resources!**

■ *Useful Output*   □ *No Output*   ■ *Low-Confidence Output*

394,170 images
30 CV models

394,170 images
30 CV models

Executing **ALL** models.

# Data-Driven Analysis

```
394,170 images
30 CV models
```

Executing models that output **VALUABLE** labels.

394,170 images
30 CV models



```
Valuable(model, img) = (Output(model, img).conf > 0.5).any() ? True : False;
```

394,170 images
30 CV models



```
Valuable(model, img) = (Output(model, img).conf > 0.5).any() ? True : False;
```

394,170 images
30 CV models



Valuable(model, img) = (Output(model, img).conf > **0.5**).**any()** ? True : False;

Executing models **RANDOMLY**.

**Without assumption of input data distribution,
how to predict the value of models before execution?**

**Without assumption of input data distribution,
how to predict the value of models before execution?**

**IMPOSSIBLE** for single-model workloads …
**But not for the multi-model tasks!**

**Without assumption of input data distribution,
how to predict the value of models before execution?**

IMPOSSIBLE for single-model workloads ...
But not for the multi-model tasks!
**We can get hints from executed models.**

# Outline

- Resource-wasting multi-model inference workloads

- **Rule-based scheduler**

- Learning-based scheduler

- Evaluation

- Conclusion

- Hard to express the non-pairwise rules.

- Too expensive for large-scale semantic labels.

- Challenging to tune the effects of multiple rules.

- Hard to express the non-pairwise rules.

- Large-scale semantic labels (>1000 labels in our workloads).

- Effects of multiple rules are difficult to tune.

### *Deep learning could help!*

- Resource-wasting multi-model inference workloads

- Rule-based scheduler

- **Learning-based scheduler**

- Evaluation

- Conclusion

One dimension per label.

One **action** per model.

For **versatility** of DRL agent, resource constraints are left to scheduling algorithms.

**Reward:**

$$r(w,d) = log(\theta_m \sum_{l_i \in O'(\{m\},d)} l_i.conf + 1), O'(\{m\},d) \neq \emptyset$$

$$r(w,d) = -1, O'(\{m\},d) = \emptyset$$

**Reward:**

$$r(w, d) = log(\theta_m \sum_{l_i \in O'(\{m\}, d)} l_i.conf + 1), O'(\{m\}, d) \neq \emptyset$$

$$r(w, d) = -1, O'(\{m\}, d) = \emptyset$$

**newly updated labels**

**Reward:**

$$r(w, d) = log(\theta_m \sum_{l_i \in O'(\{m\}, d)} l_i.conf + 1), O'(\{m\}, d) \neq \emptyset$$

$$r(w, d) = -1, O'(\{m\}, d) = \emptyset$$

**Priority parameter**    **Sum of label confidence**

**Reward:**

$$r(w, d) = log(\theta_m \sum_{l_i \in O'(\{m\}, d)} l_i.conf + 1), O'(\{m\}, d) \neq \emptyset$$

$$r(w, d) = -1, O'(\{m\}, d) = \emptyset$$

**Logarithmic smoothing**

**Reward:**

$$r(w, d) = log(\theta_m \sum_{l_i \in O'(\{m\}, d)} l_i.conf + 1), O'(\{m\}, d) \neq \emptyset$$

$$r(w, d) = -1, O'(\{m\}, d) = \emptyset$$

**Punishment**

- Deep Q-Network
  - input -> dense-layer -> ReLU -> (N+1)-output
  - The **+1** action is an **END** action

- Deep Q-Network

  - 1104-input -> 256-dense -> ReLU -> (30+1)-output

  - The *+1* action is an ***END*** action

  - The reward of selecting ***END*** is *0*

- Deep Q-Network
  - 1104-input -> 256-dense -> ReLU -> (30+1)-output
  - The *+1* action is an *END* action
  - The reward of selecting *END* is *0*
- Training mechanisms:
  - Original DQN
  - Double DQN
  - Dueling DQN
  - Deep SARSA

- An **adaptive submodular function maximization** problem.
- Existing approximate algorithms with performance guarantee is infeasible for our problem, since they require partial permutation on items.

- Two common constraints of computing resources are studied:

### 1-D Deadline Constraint

**Algorithm 1** Scheduling under deadline constraint.

**Input:** model set $M$, time budget $B_{time}$, DRL agent $Q$
**Output:** model subset $S$
1: $S \leftarrow \emptyset$
2: **while** $B_{time} > 0$ **do**
3:     Filter out models that $m.time > B_{time}$
4:     $m^* \leftarrow \arg\max\limits_{m \in M \setminus S} \frac{Q(m,d)}{m.time}$
5:     $S \leftarrow S \cup \{m^*\}$, $B_{time} \leftarrow B_{time} - m^*.time$
6: **end while**
7: **return** $S$

### 2-D Deadline-Memory Constraints

**Algorithm 2** Scheduling under deadline-memory constraints.

**Input:** model set $M$, time budget $B_{time}$, memory budget $B_{mem}$, DRL agent $Q$
**Output:** model scheduling policy $S$
1: $S \leftarrow [\ ]$, $TimeCost \leftarrow 0$, $S_t \leftarrow \emptyset$
2: **while** $TimeCost < B_{time}$ **do**
3:     Filter out models that $m.mem > B_{mem}$
4:     $m_1^* \leftarrow \arg\max\limits_{m \in M \setminus S} \frac{Q(m,d)}{m.time \times m.mem}$
5:     $S_t \leftarrow S_t \cup \{m_1^*\}$, $B_{time}^t \leftarrow TimeCost + m_1^*.time$
6:     Filter out models by temporary deadline $B_{time}^t$
7:     **while** $B_{mem} > 0$ **do**
8:         $m_2^* \leftarrow \arg\max\limits_{m \in M \setminus S} \frac{Q(m,d)}{m.mem}$
9:         $S_t \leftarrow S_t \cup \{m_2^*\}$, $B_{mem} \leftarrow B_{mem} - m_2^*.mem$
10:    **end while**
11:    $S.append(S_t)$, Wait until model $m_3^* \in S_t$ finishes
12:    $B_{mem} \leftarrow B_{mem} + m_3^*.mem$, $S_t \leftarrow S_t \setminus \{m_3^*\}$
13: **end while**
14: **return** $S$

- An adaptive submodular function maximization problem.

- Two common constraints are studied:

**1-D Deadline Constraint**

**Algorithm 1** Scheduling under deadline constraint.

**Input:** model set $M$, time budget $B_{time}$, DRL agent $Q$

**Output:** model subset $S$

1: $S \leftarrow \emptyset$
2: **while** $B_{time} > 0$ **do**
3:      Filter out models that $m.time > B_{time}$
4:      $m^* \leftarrow \arg \max\limits_{m \in M \setminus S} \dfrac{Q(m,d)}{m.time}$     ----- **Based on model profiling on sampled data.**
5:      $S \leftarrow S \cup \{m^*\}, B_{time} \leftarrow B_{time} - m^*.time$
6: **end while**
7: **return** $S$

- An adaptive submodular function maximization problem.

- Two common constraints are studied:

**1-D Deadline Constraint**

**Algorithm 1** Scheduling under deadline constraint.

**Input:** model set $M$, time budget $B_{time}$, DRL agent $Q$
**Output:** model subset $S$
1: $S \leftarrow \emptyset$
2: **while** $B_{time} > 0$ **do**
3:     Filter out models that $m.time > B_{time}$
4:     $m^* \leftarrow \arg \max_{m \in M \setminus S} \frac{Q(m,d)}{m.time}$
5:     $S \leftarrow S \cup \{m^*\}$, $B_{time} \leftarrow B_{time} - m^*.time$
6: **end while**
7: **return** $S$

TIME

Img1

| M1 | M5 | M8 |

Img4 ...

Img2

| M3 | M4 |

Img5 ...

Img3

| M6 | M1 | M2 | M7 |

Img6 ...

- An adaptive submodular function maximization problem.

- Two common constraints are studied:

**2-D Deadline-Memory Constraints**

---
**Algorithm 2** Scheduling under deadline-memory constraints.

---
**Input:** model set $M$, time budget $B_{time}$, memory budget $B_{mem}$, DRL agent $Q$

**Output:** model scheduling policy $S$

1: $S \leftarrow [\ ]$, $TimeCost \leftarrow 0$, $S_t \leftarrow \emptyset$
2: **while** $TimeCost < B_{time}$ **do**
3:      Filter out models that $m.mem > B_{mem}$
4:      $m_1^* \leftarrow \arg \max_{m \in M \setminus S} \dfrac{Q(m,d)}{m.time \times m.mem}$
5:      $S_t \leftarrow S_t \cup \{m_1^*\}$, $B_{time}^t \leftarrow TimeCost + m_1^*.time$
6:      Filter out models by temporary deadline $B_{time}^t$

**Step#1:**
**Determine the temporary *deadline*.**

- An adaptive submodular function maximization problem.

- Two common constraints are studied:

**2-D Deadline-Memory Constraints**

5:  $S_t \leftarrow S_t \cup \{m_1^*\}$, $B_{time}^t \leftarrow TimeCost + m_1^*.time$
6:  Filter out models by temporary deadline $B_{time}^t$
7:  **while** $B_{mem} > 0$ **do**
8:  $\quad m_2^* \leftarrow \arg \max_{m \in M \setminus S} \frac{Q(m,d)}{m.mem}$
9:  $\quad S_t \leftarrow S_t \cup \{m_2^*\}$, $B_{mem} \leftarrow B_{mem} - m_2^*.mem$
10:  **end while**
11:  $S.append(S_t)$, Wait until model $m_3^* \in S_t$ finishes
12:  $B_{mem} \leftarrow B_{mem} + m_3^*.mem$, $S_t \leftarrow S_t \setminus \{m_3^*\}$
13:  **end while**
14:  **return** $S$

**Step#2:**
**Greedily fill in the memory pool.**

- An adaptive submodular function maximization problem.
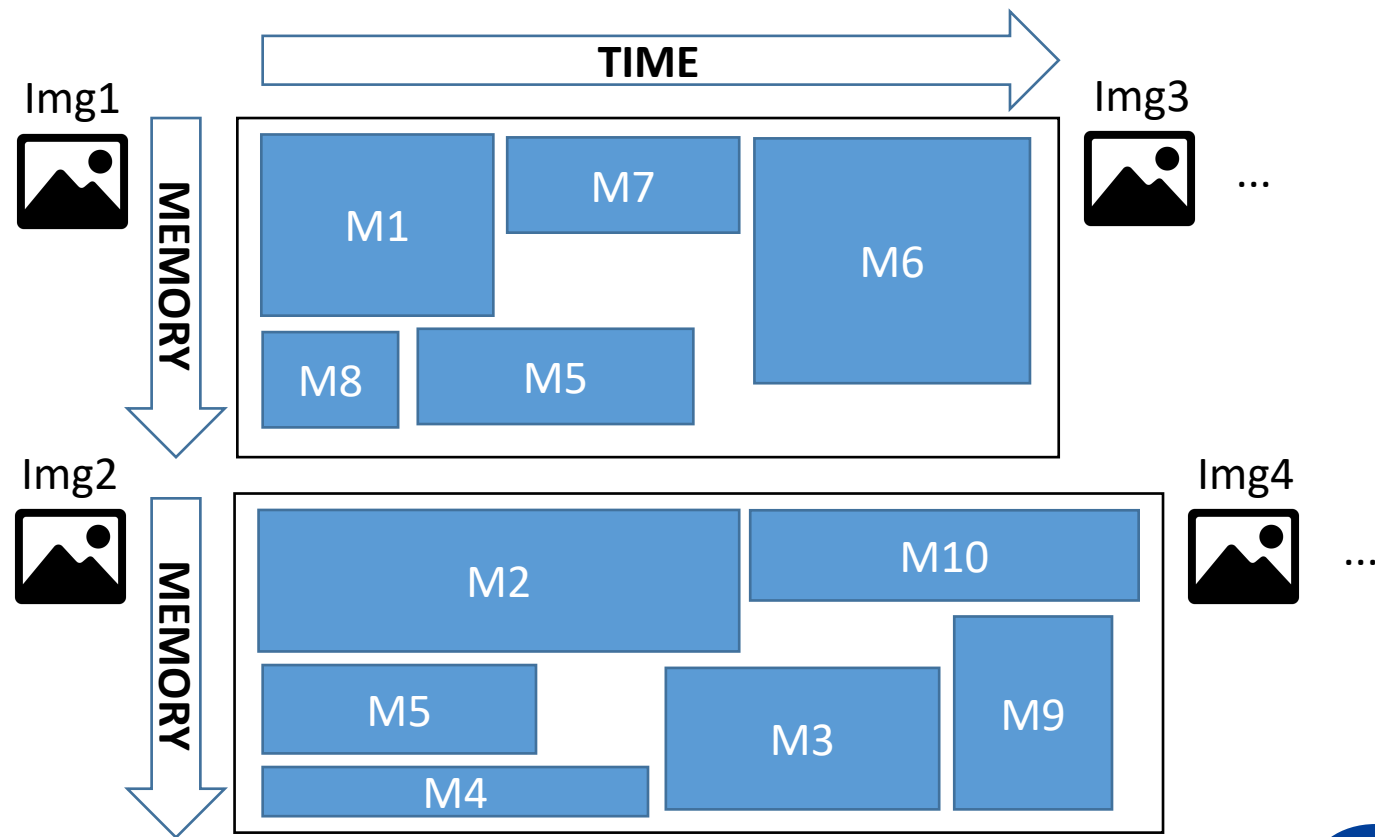
- Two common constraints are studied:

**2-D Deadline-Memory Constraints**



**Algorithm 2** Scheduling under deadline-memory constraints.

**Input:** model set $M$, time budget $B_{time}$, memory budget $B_{mem}$, DRL agent $Q$

**Output:** model scheduling policy $S$

1: $S \leftarrow [\,], TimeCost \leftarrow 0, S_t \leftarrow \emptyset$
2: **while** $TimeCost < B_{time}$ **do**
3:    Filter out models that $m.mem > B_{mem}$
4:    $m_1^* \leftarrow \arg \max_{m \in M \setminus S} \frac{Q(m,d)}{m.time \times m.mem}$
5:    $S_t \leftarrow S_t \cup \{m_1^*\}, B_{time}^t \leftarrow TimeCost + m_1^*.time$
6:    Filter out models by temporary deadline $B_{time}^t$
7:    **while** $B_{mem} > 0$ **do**
8:       $m_2^* \leftarrow \arg \max_{m \in M \setminus S} \frac{Q(m,d)}{m.mem}$
9:       $S_t \leftarrow S_t \cup \{m_2^*\}, B_{mem} \leftarrow B_{mem} - m_2^*.mem$
10:    **end while**
11:    $S.append(S_t)$, Wait until model $m_3^* \in S_t$ finishes
12:    $B_{mem} \leftarrow B_{mem} + m_3^*.mem, S_t \leftarrow S_t \setminus \{m_3^*\}$
13: **end while**
14: **return** $S$

# Outline

- Resource-wasting multi-model inference workloads

- Rule-based scheduler

- Learning-based scheduler

- **Evaluation**

- Conclusion

## Multi-Model Image Labeling Workloads

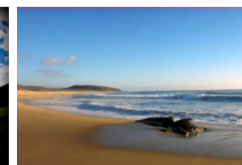| Task | Label# |
|---|---|
| Object Detection | 80 |
| Scene Classification | 365 |
| Face Detection | 1 |
| Face Landmark Localization | 70 |
| Pose Estimation | 17 |
| Emotion Classification | 7 |
| Gender Classification | 2 |
| Action Classification | 40 |
| Hand Landmark Localization | 42 |
| Dog Classification | 120 |
| **10 Tasks** | **1104 Labels** |

## MSCOCO-2017



## MIRFLICKR-25k



by Silke Gerstenkorn    by Dave Wild    by Hugo A.B. Olivas
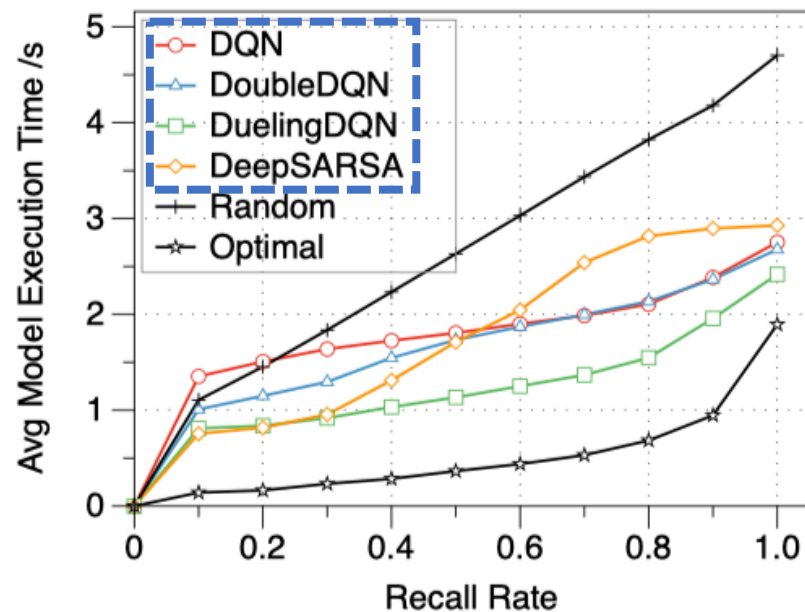
by Martin P. Szymczak    by Mani Babbar    by Lee Otis
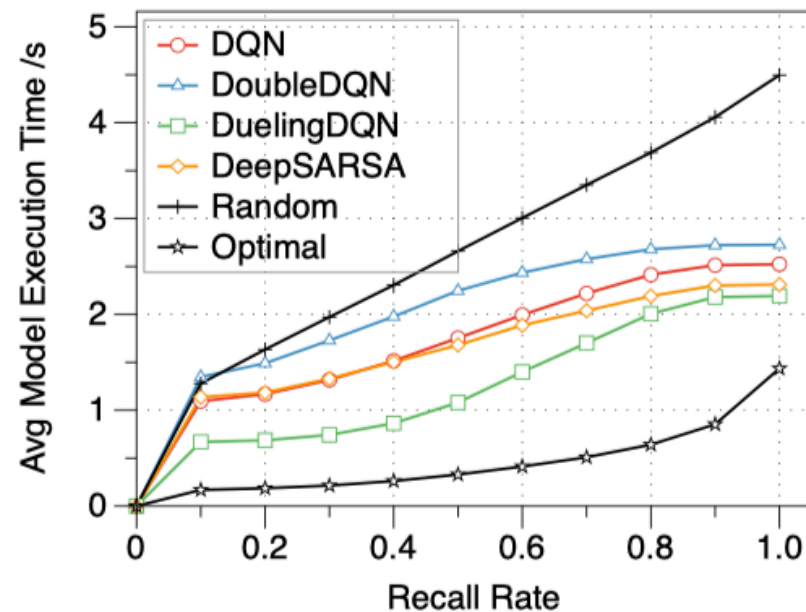
**Different Training Mechanisms**



(a) MSCOCO 2017
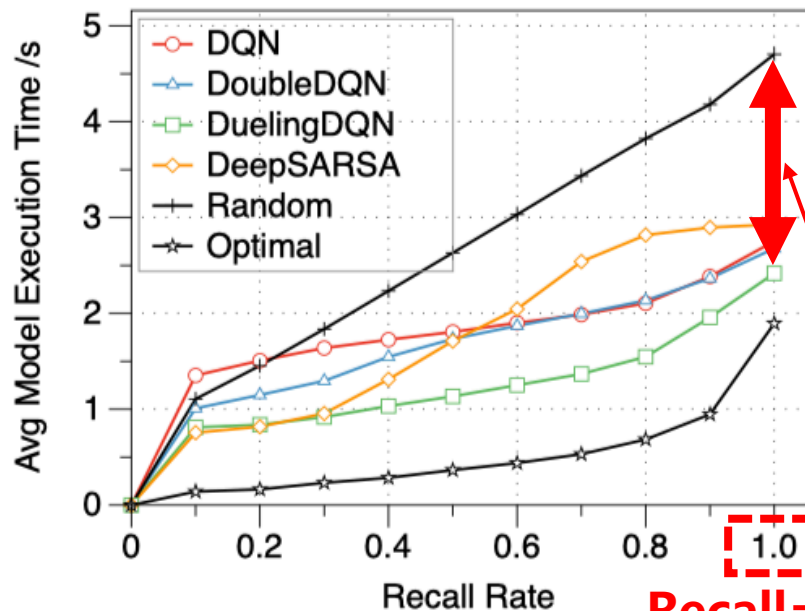
(b) MirFlickr25

(a) MSCOCO 2017

(b) MirFlickr25

Recall=1.0

Better

Saving **48.6-51.2%** execution time **without loss** of valuable labels.

45

**10 Manual Rules**

| Current Model Task | Output Label | Rule |
|---|---|---|
| Object Detection | person | $2 \times \mathcal{P}(\text{Pose Estimation})$ |
| Object Detection | person | $2 \times \mathcal{P}(\text{Gender Classification})$ |
| Object Detection | dog | $2 \times \mathcal{P}(\text{Dog Classification})$ |
| Face Detection | face | $2 \times \mathcal{P}(\text{Face Landmark Localization})$ |
| Face Detection | face | $2 \times \mathcal{P}(\text{Emotion Classification})$ |
| Pose Estimation | body keypoints | $2 \times \mathcal{P}(\text{Action Classification})$ |
| Pose Estimation | wrist keypoints | $2 \times \mathcal{P}(\text{Hand Landmark Localization})$ |
| Place Classification | indoor places | $0.5 \times \mathcal{P}(\text{Animal-Object Detection})$ |
| Place Classification | indoor places | $0.5 \times \mathcal{P}(\text{Sport-Action Classification})$ |

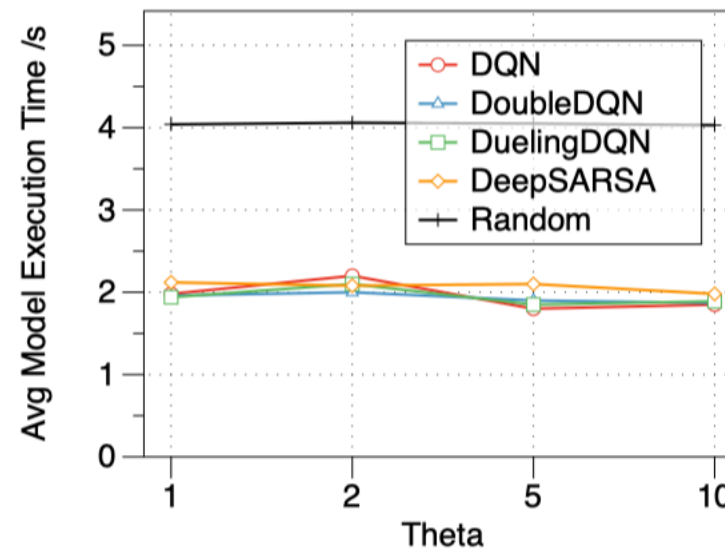Saving only **1.4%** execution time
when required recall is 1.0.

**Adjusting priority of the face-detection model.**

$$r(w, d) = log(\theta_m \sum_{l_i \in O'(\{m\}, d)} l_i.conf + 1), O'(\{m\}, d) \neq \emptyset$$
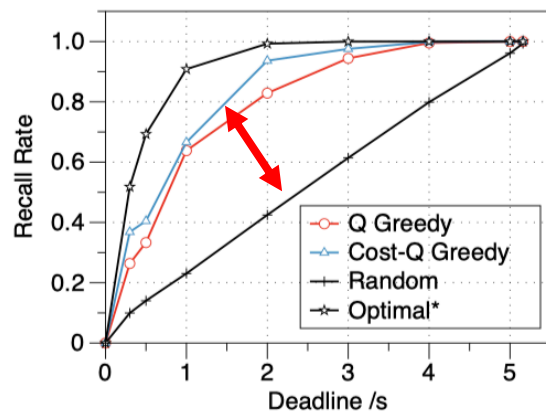
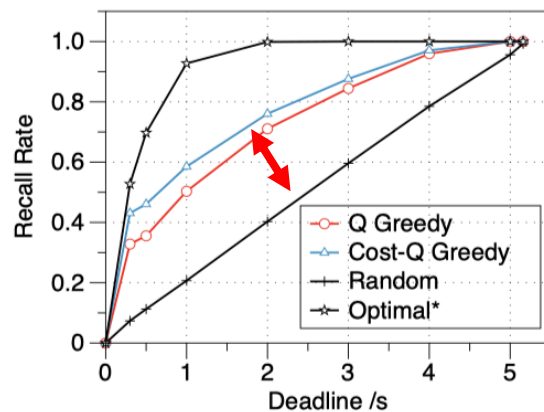**Priority parameter**



(a) Average execution order

(b) Average time cost
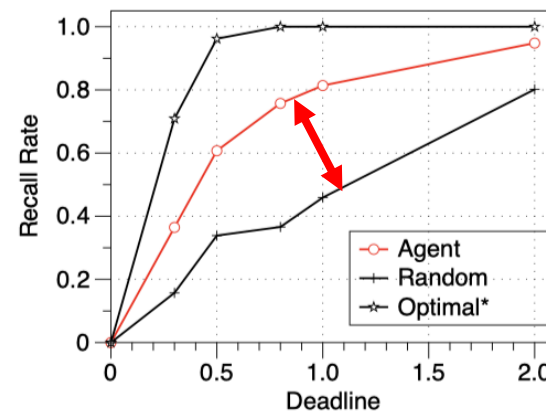
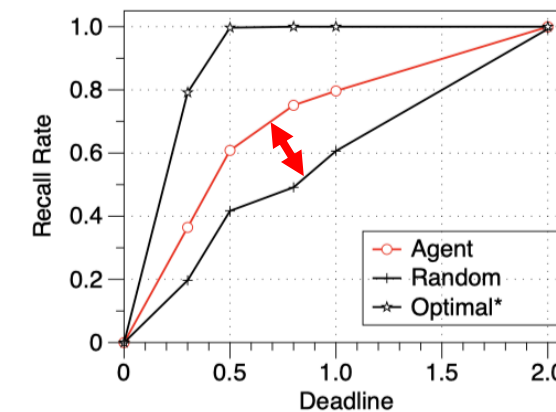## 1-D Deadline Constraint



(a) MSCOCO 2017

(b) MirFlickr25

Boosts **188.7-309.5%** recall of valuable labels under 0.5s deadline constraint.

## 2-D Deadline-Memory Constraints



(a) 8GB Memory

(b) 12GB Memory

Boosts **106.9%/52.8%** recall of valuable labels under 0.8s deadline and 8/12GB memory constraints.

- Resource-wasting multi-model inference workloads
- Rule-based scheduler
- Learning-based scheduler
- Evaluation
- **Conclusion**

**Adaptive model scheduling can improve the efficiency of multi-model inference workloads by avoiding valueless execution.**

## 12 Faculty members, 2 Post-Doc, 3 Secretary; 7 with PhD from abroad

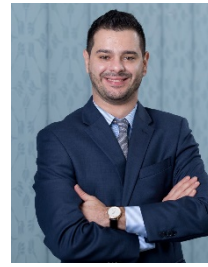**XiangYang Li**
IEEE Fellow,
ACM Fellow,
ACM China Co-Chair

**Yanyong Zhang**
IEEE Fellow,
Prof. in Rutgers,
NSF Career

**Panlong Yang**
CCF Dist Speaker
Wireless network、
Mobile computing

**Nikolaos M.Freris**
USA NYU A.P.
CPS, Algorithms,
Distributed optimization
Machine learning

**Lan Zhang**
CCF, ACM China Doctor
Thesis Award, Youqing,
Qingcheng Award
Data understanding/trading,
privacy protection

**Bei Hua**
High performance
computing、
Edge computing

**Yu Zhang**
system software ,Software
optimization/security,
Quantum software

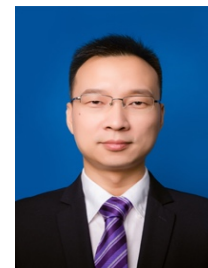**Hao Zhou**
Japan NTII
Wireless Network Resource
Management

**Haisheng Tan**
HK, Tsinghua Post-
Doc
Cloud computing、
Algorithms Analysis

**YuBo Yan**
Wireless/Passive
network, IntelliSense,
IoT, SDR

**Xin He**
Doc. University of Oulu
Passive network,
Theories of Information
and Coding

**Xin Guo**
Edge computing,
Security of IoT

**Xuerong Huang**
Master in HKBU
Research assistant

**Ludi Xue**
Research assistant

## Mu Yuan

University of Science and Technology of China
School of Computer Science and Technology
ym0813@mail.ustc.edu.cn

## Lan Zhang

University of Science and Technology of China
School of Computer Science and Technology
zhanglan@ustc.edu.cn