

# Learning Without Forgetting

Zhizhong Li, Derek Hoiem {zli115,dhoiem}@illinois.edu  
Department of Computer Science, University of Illinois at Urbana-Champaign

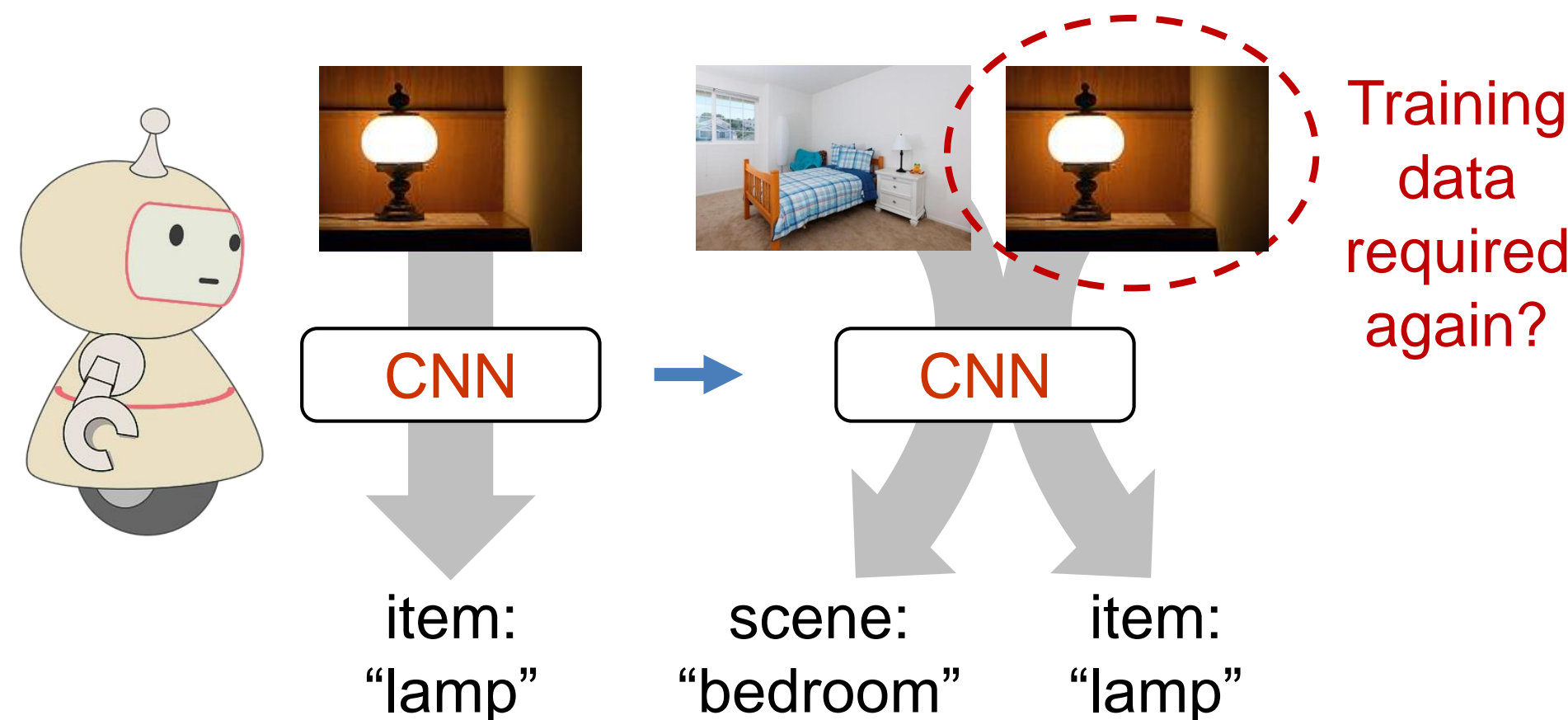
<http://zli115.web.engr.illinois.edu/learning-without-forgetting/>



## Motivation

When expanding the capability of a vision system...

- Fine-tuning? (old task suffers)
- Feature extraction? (new task suffers)
- Joint training:



What if the original dataset...

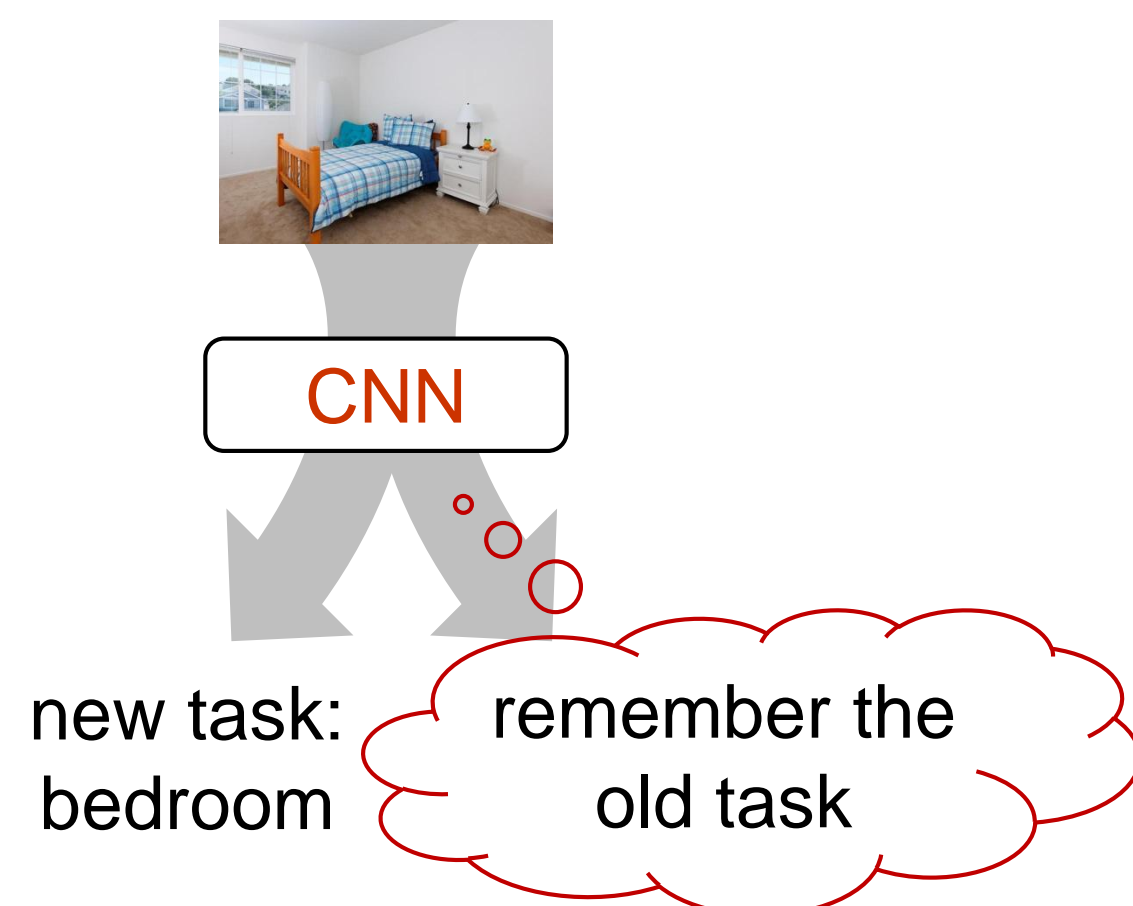
- Is not recorded?
- Is proprietary?
- Is too cumbersome?

But we want...

- Benefit of shared representation
- No or little degradation of the original capability
- Without the need to access original task dataset?

**Goal:**

Add new capabilities to a CNN-based vision system using *only data from the new task*.



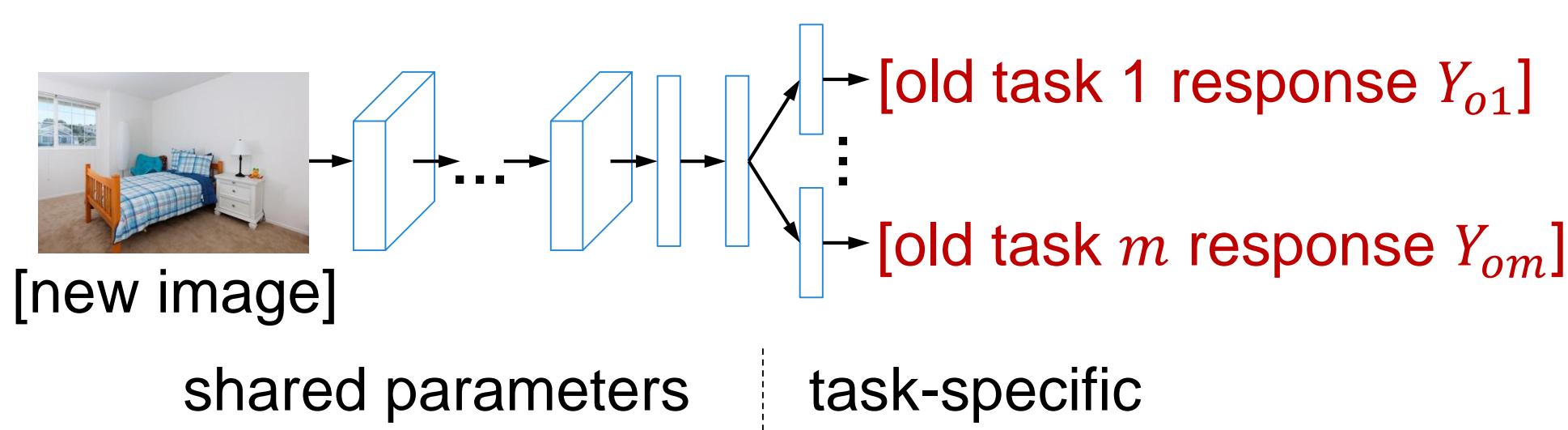
**Our strengths:**

- Outperforms the widely-used fine-tuning on *both original and new task*.
- Outperforms feature extraction on the new task.
- Simple to implement and deploy
- Training efficiency compared to joint training

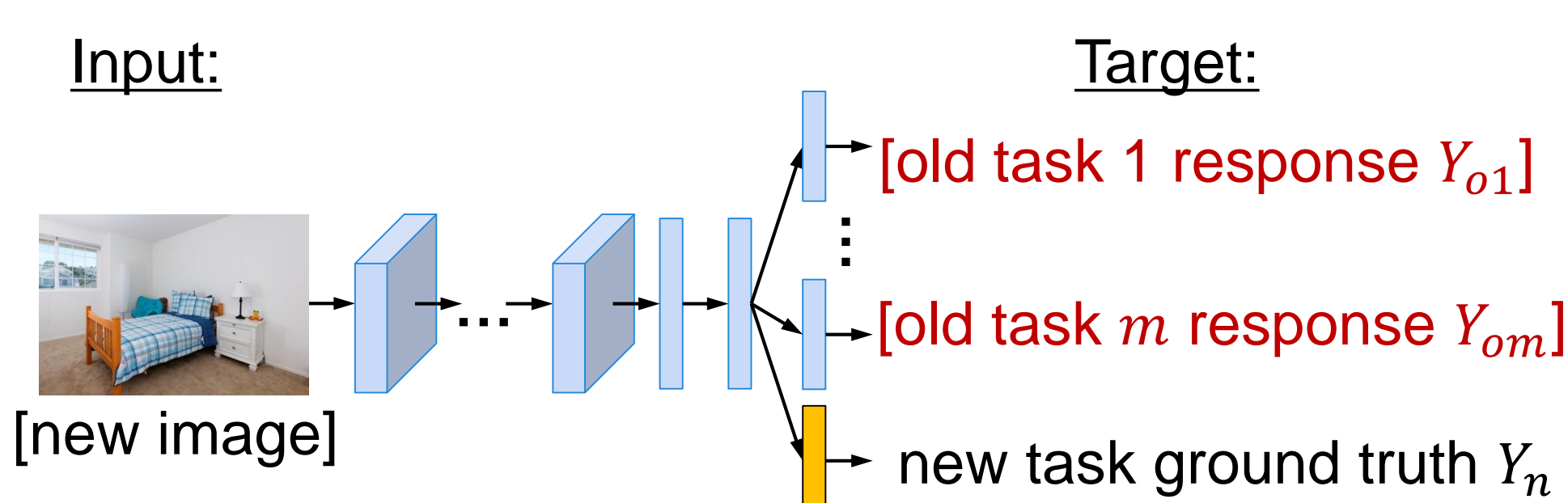
## Method

**Outline**

1. Obtain old task responses



2. Train on new images



**Training: loss**

$$\theta_s^*, \theta_o^*, \theta_n^* \leftarrow \underset{\theta_s, \theta_o, \theta_n}{\operatorname{argmin}} \left( \sum_i^m \mathcal{L}_{old}(Y_o, \hat{Y}_o; \hat{\theta}_s, \hat{\theta}_o) + \mathcal{L}_{new}(Y_n, \hat{Y}_n; \hat{\theta}_s, \hat{\theta}_n) + \mathcal{R}(\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n) \right)$$

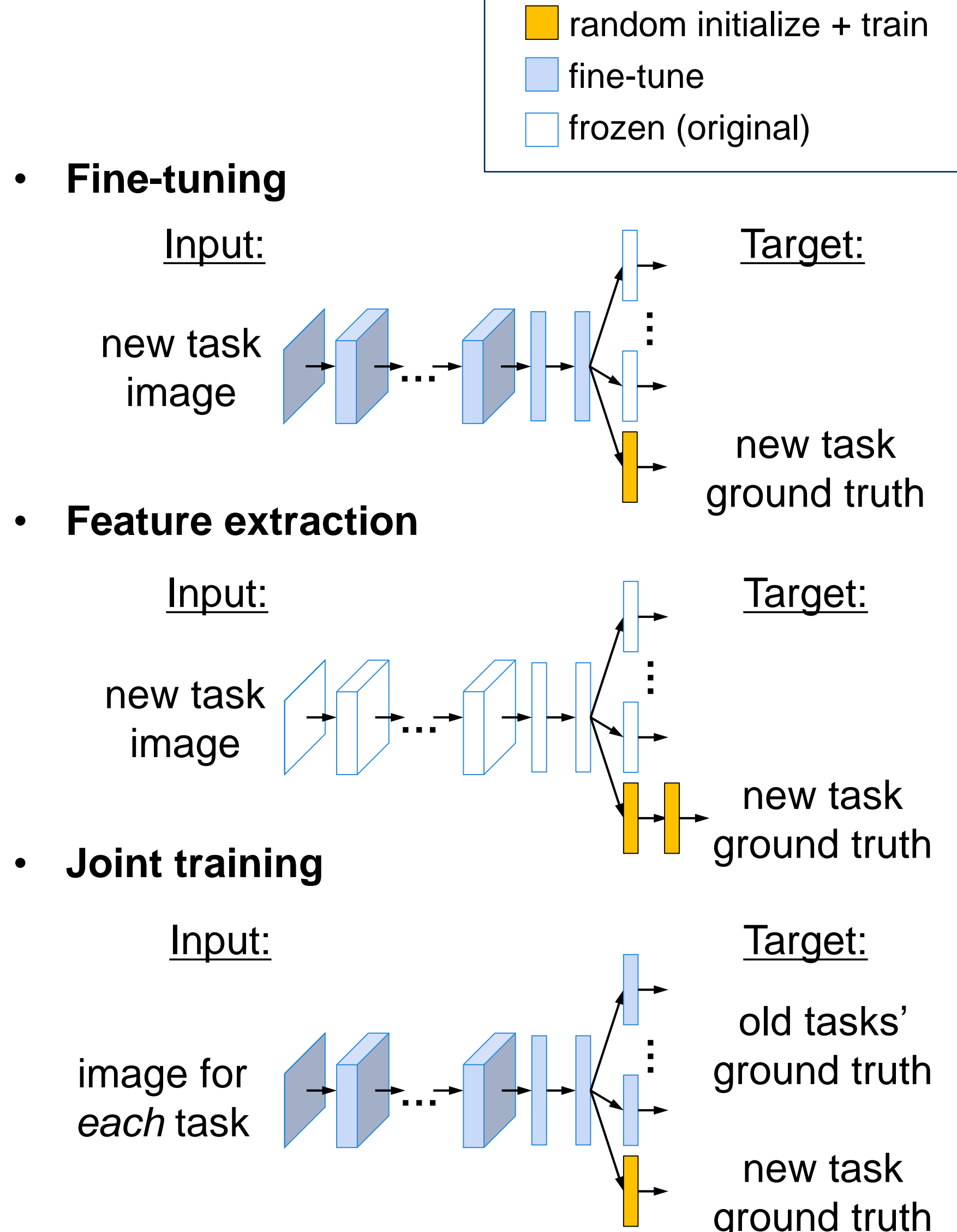
shared/old/new parameters

old task response preservation loss

new task classification loss

regularization

## Compared methods

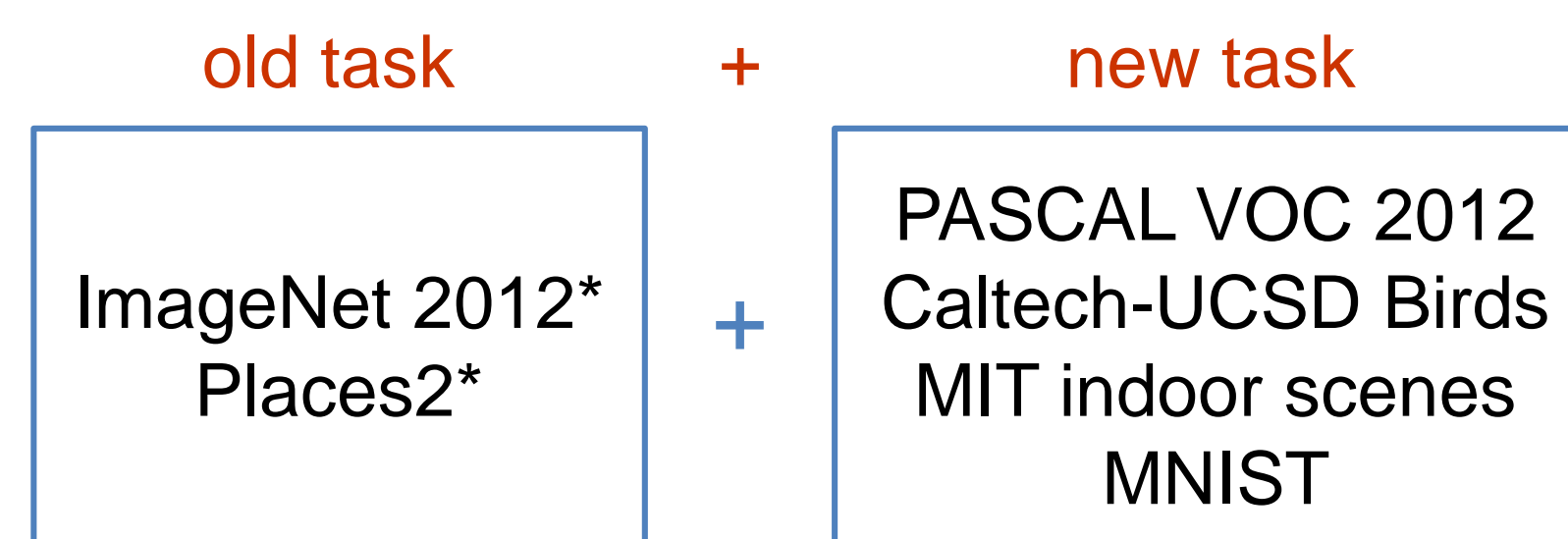


## Limitations of existing methods

	Fine Tuning	Duplicating and Fine Tuning	Feature Extraction	Joint Training	Learning without Forgetting
new task performance	good	good	X medium	best	✓ best
original task performance	X bad	good	good	good	✓ good
training efficiency	fast	fast	fast	X slow	✓ fast
testing efficiency	fast	X slow	fast	fast	✓ fast
storage requirement	medium	X large	medium	X large	✓ medium
requires old task data	no	no	no	X yes	✓ no

## Experiment Settings

**Datasets**



\* Pre-trained AlexNet obtained from authors

**Efficiency:**

- Training: forward-pass shared parameters once. Faster than joint training, similar to fine-tuning
- Test: same as compared methods; more efficient than keeping different networks for each task

**Design choices and alternatives**

We experimented with some variations:

- Possibly: more layers as task-specific parameters.
  - Possibly: add nodes to earlier layers
  - Possibly: use alternative loss for  $\mathcal{L}_{old}(Y_o, \hat{Y}_o)$
  - Possibly: just reduce fine-tuning learning rate
- These variations provided insignificant or inconsistent improvements, if any.

## Conclusions

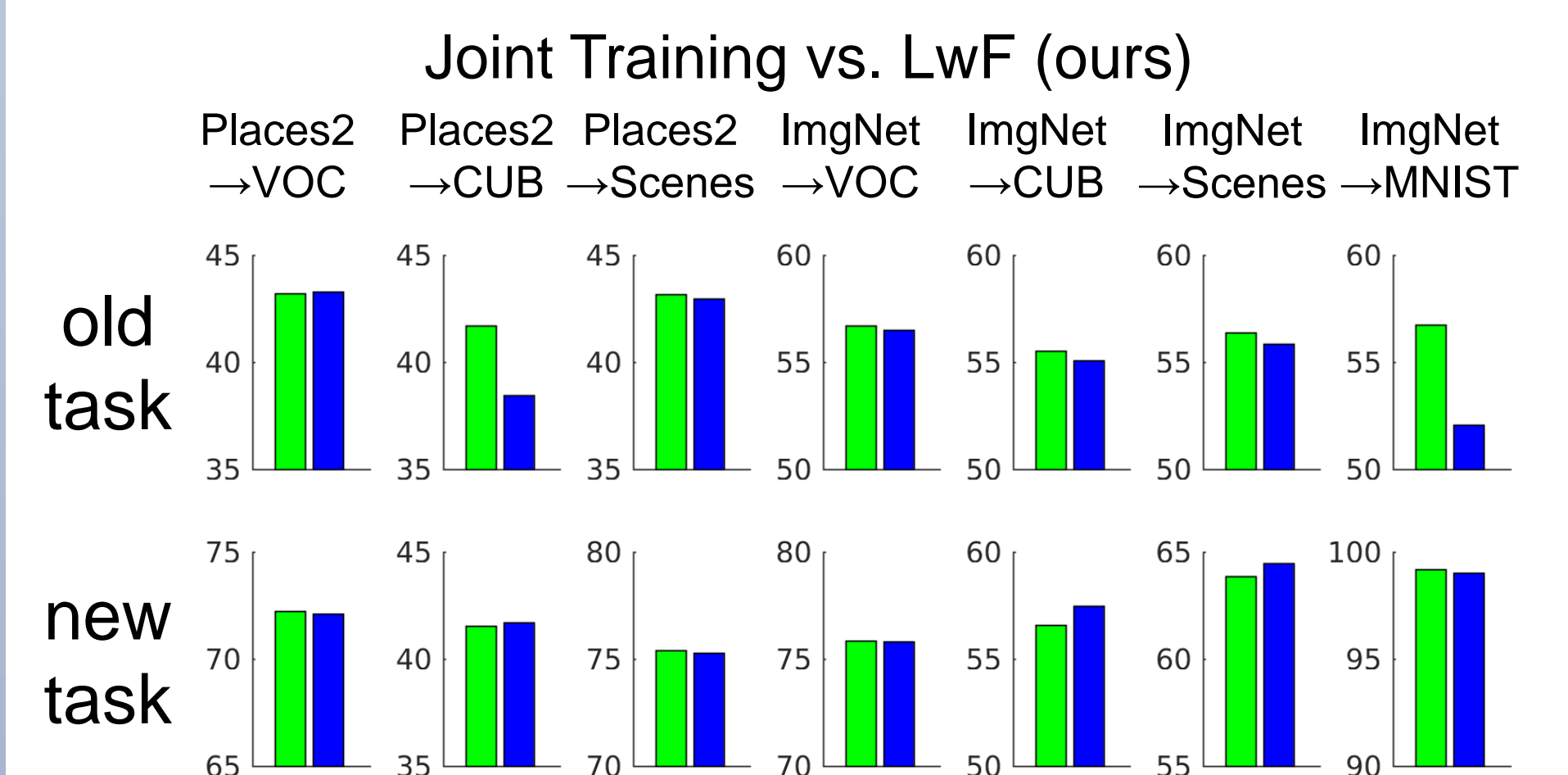
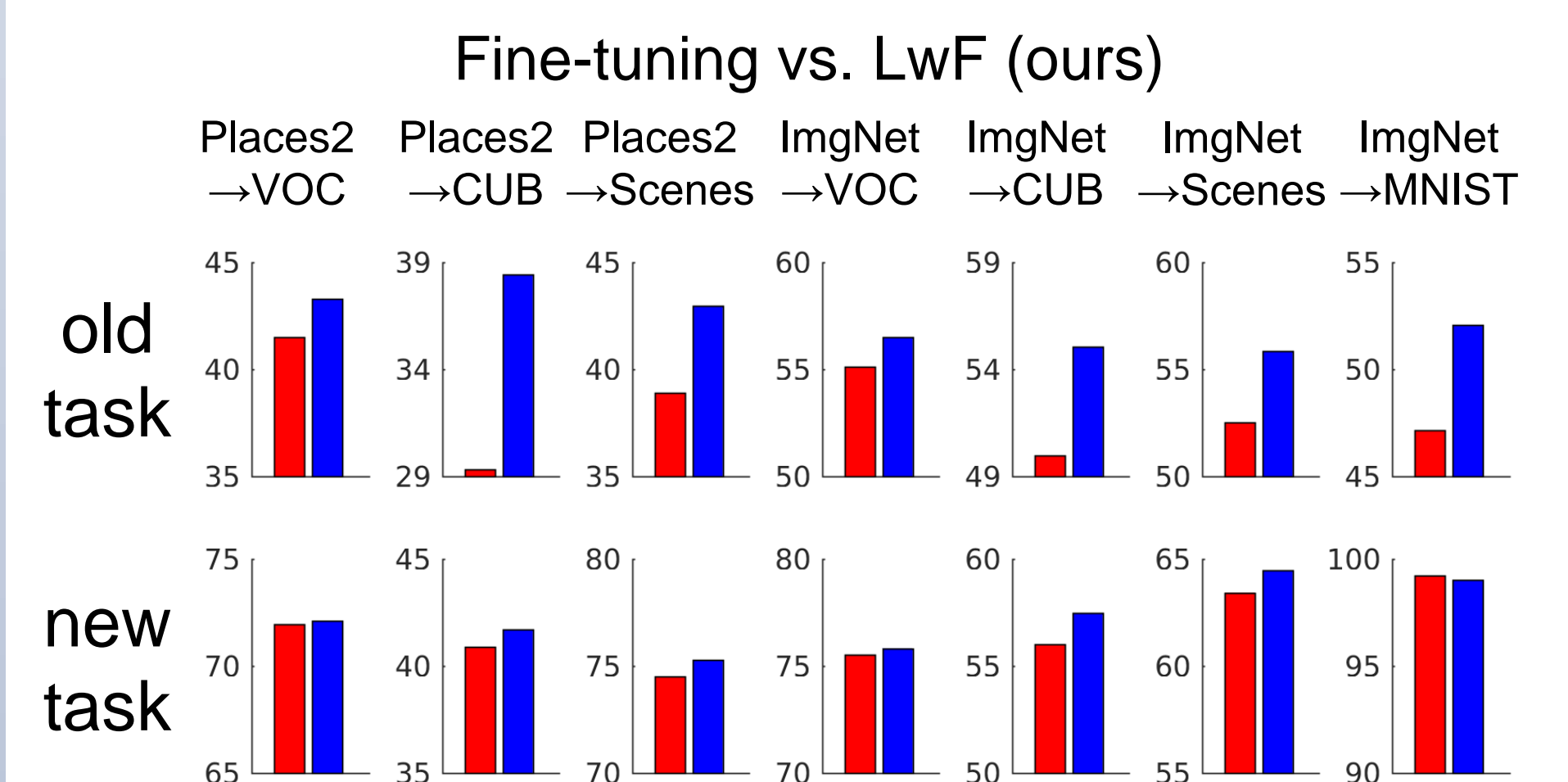
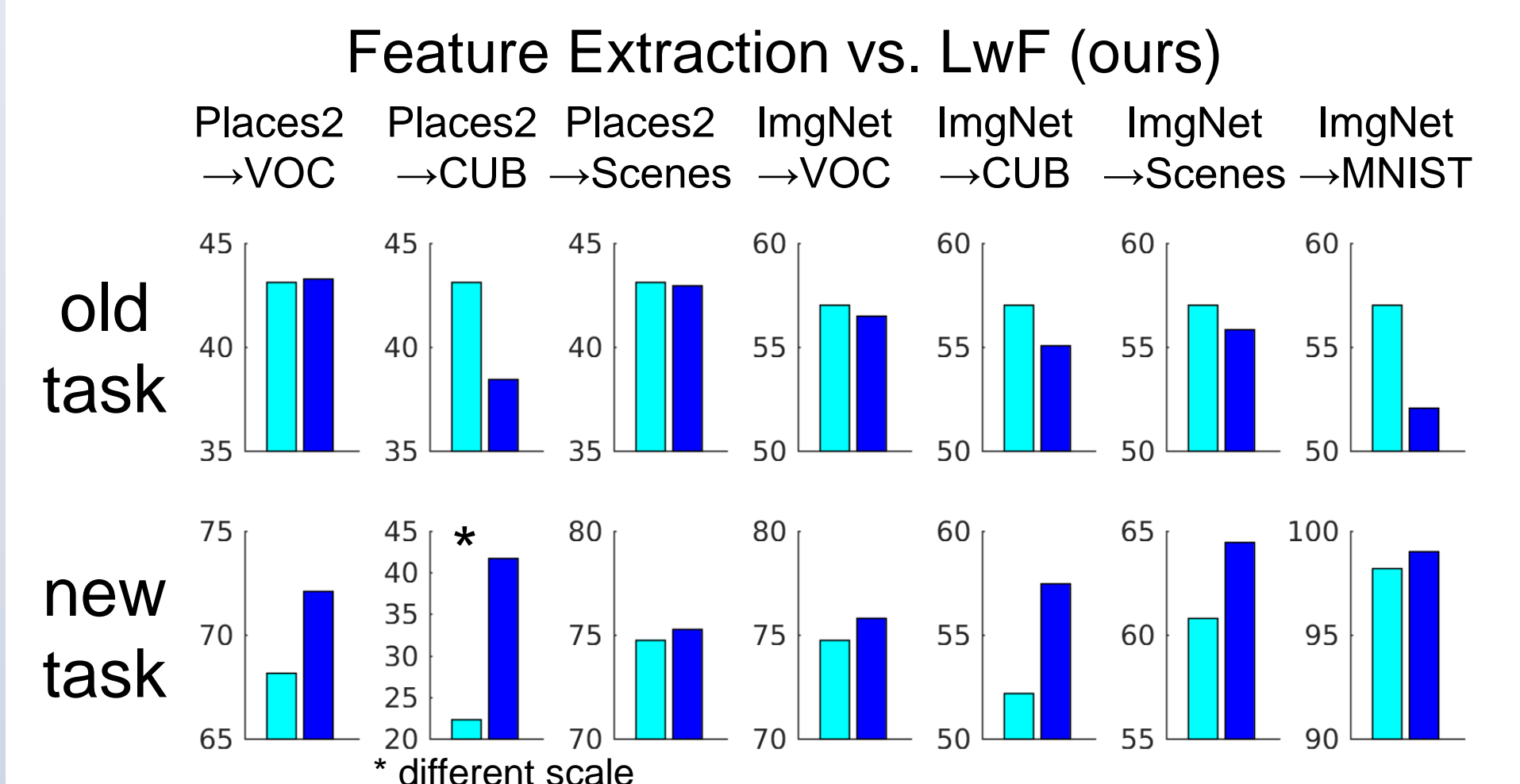
- Vs. Feature Extraction: LwF outperforms on new task; underperforms on old task
- Vs. Fine-tuning: LwF outperforms on both tasks, as keeping old responses regularizes model
- Vs. Joint Training: LwF performs nearly as well as joint training
- Dissimilar new tasks degrade old task performance
- Similar results and same observations for adding multiple new tasks

## Results

**Single new task scenario**



\* Accuracy (average precision for VOC)  
\* Using AlexNet



\* Validation set results shown. Test set results similar  
\* VGG-16 network results are mostly similar, however Joint Training outperforms our method more on both tasks (0.8%~2.5%)

**Multiple new task scenario**

