

# Adversarial Continual Learning



Sayna Ebrahimi  
UC Berkeley



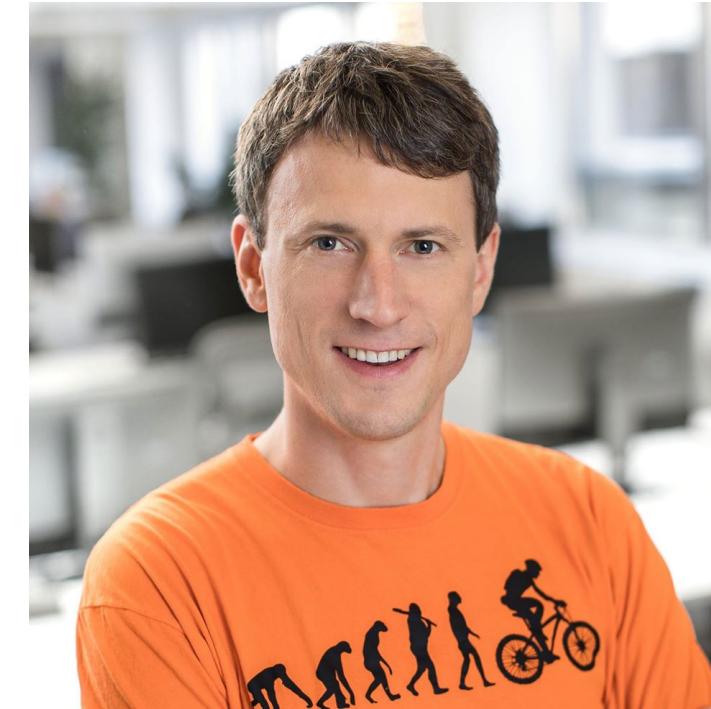
Franziska Meier  
Facebook AI Research



Roberto Calandra  
Facebook AI Research



Trevor Darrell  
UC Berkeley



Marcus Rohrbach  
Facebook AI Research

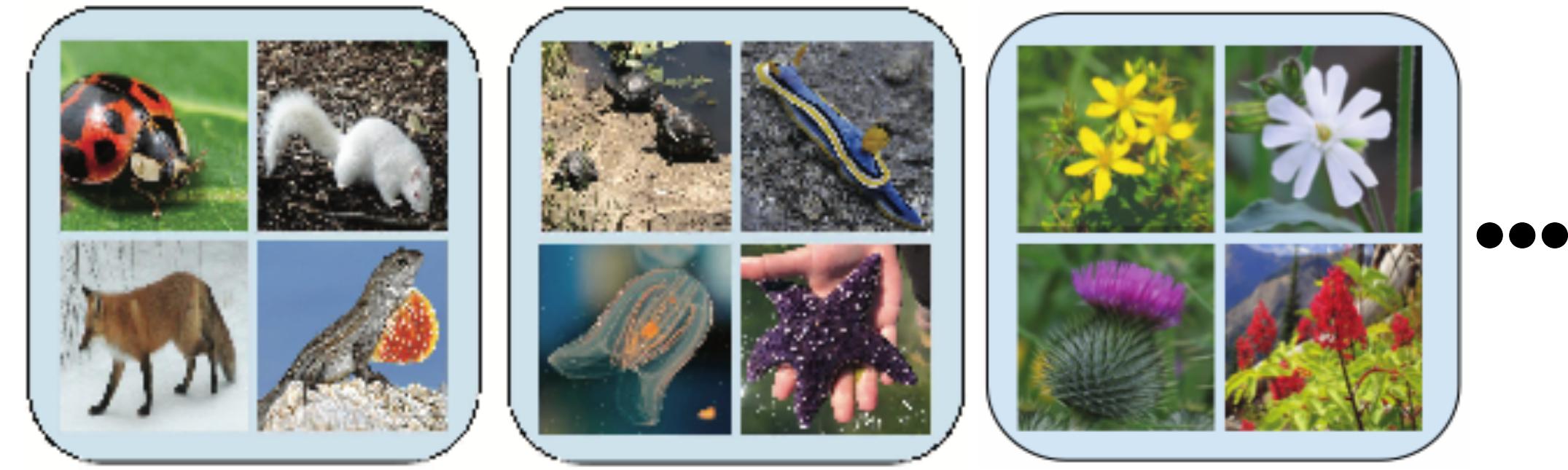
# What is Continual Learning?

## Definition:

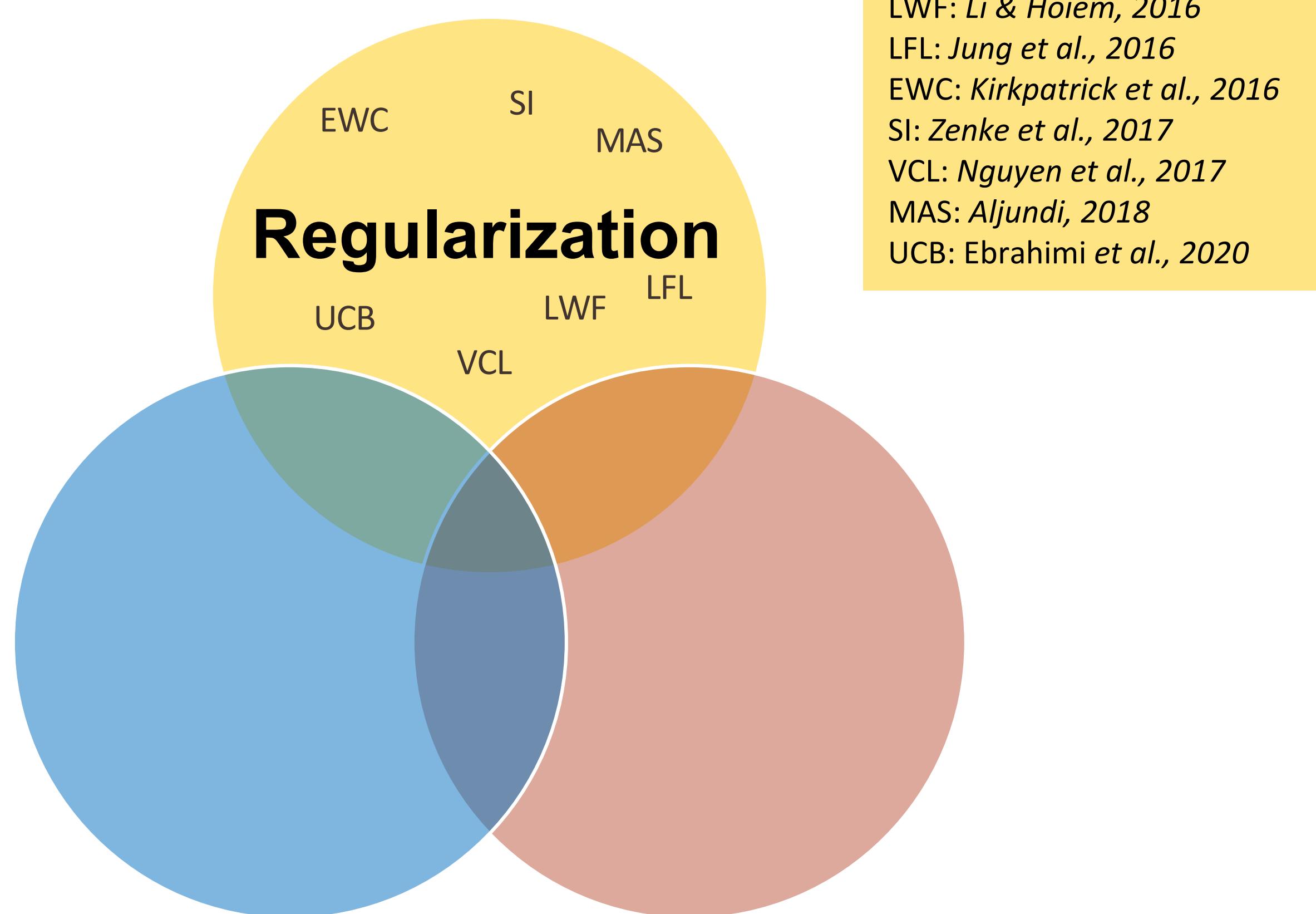
learning a sequence of tasks and performing well on all of them

## Objectives:

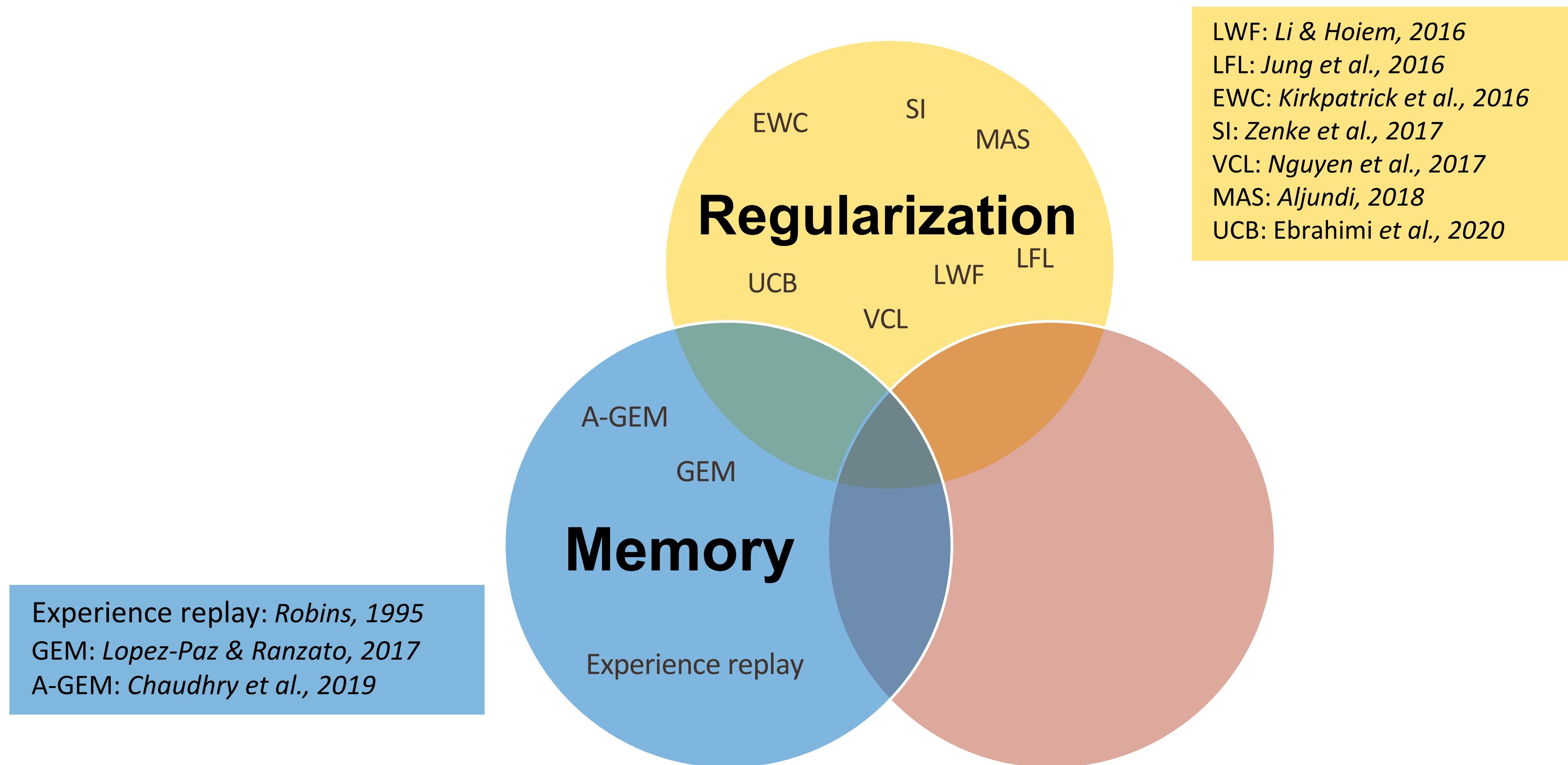
- No forgetting
- Data streams and revisiting is not allowed/limited
- High knowledge transferability
- Efficiency and scalability
- No/limited task information at test time
- etc.



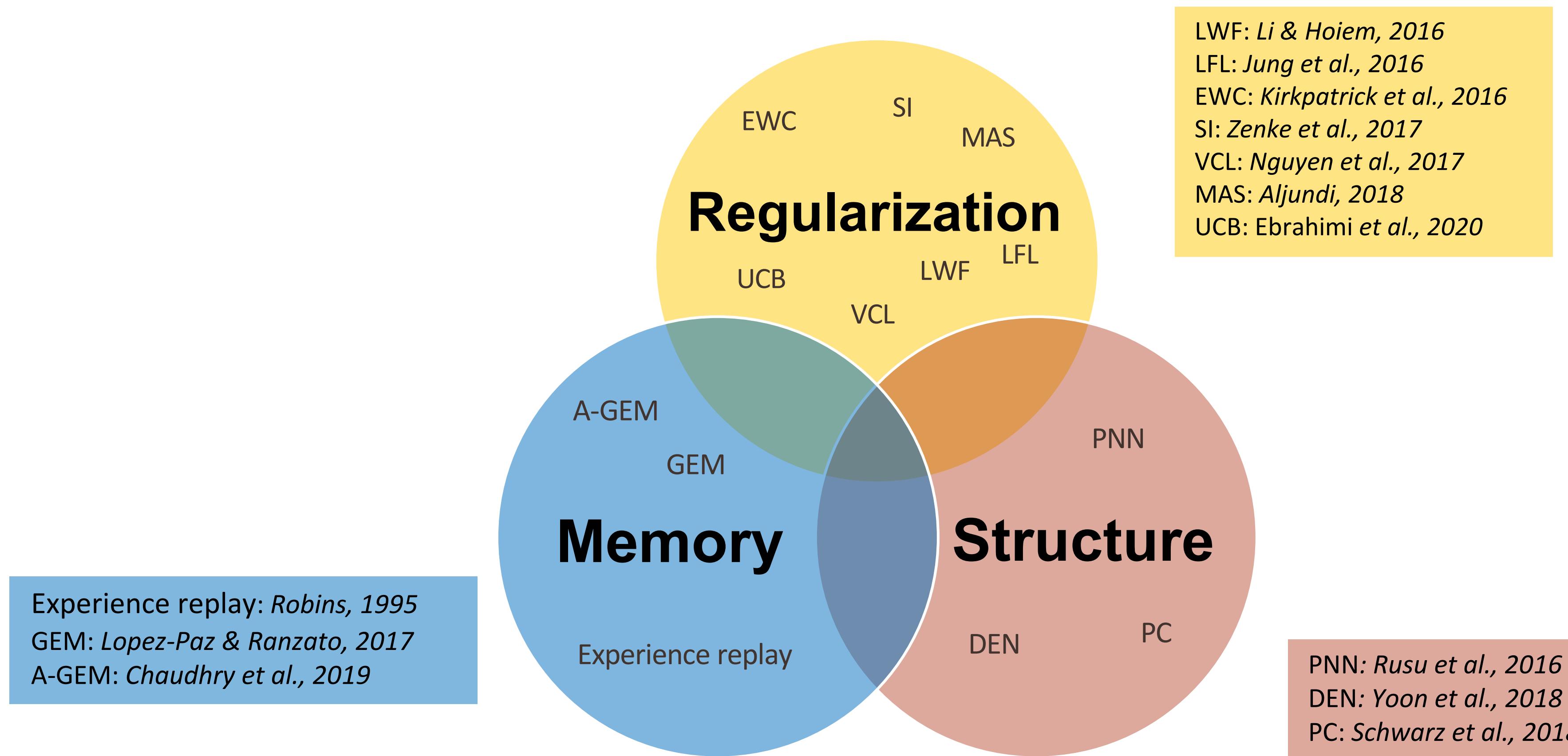
# Approaches in Continual Learning



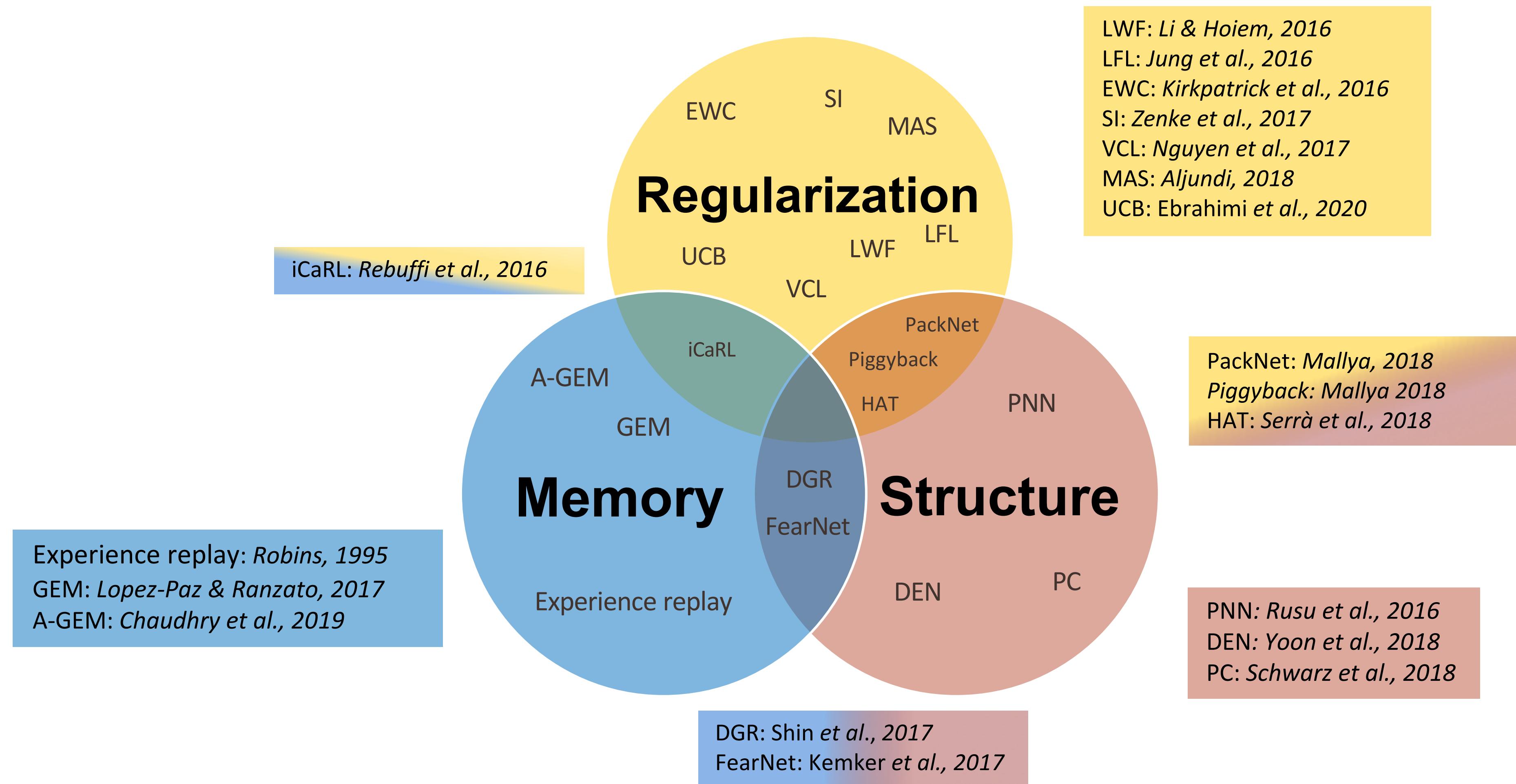
# Approaches in Continual Learning



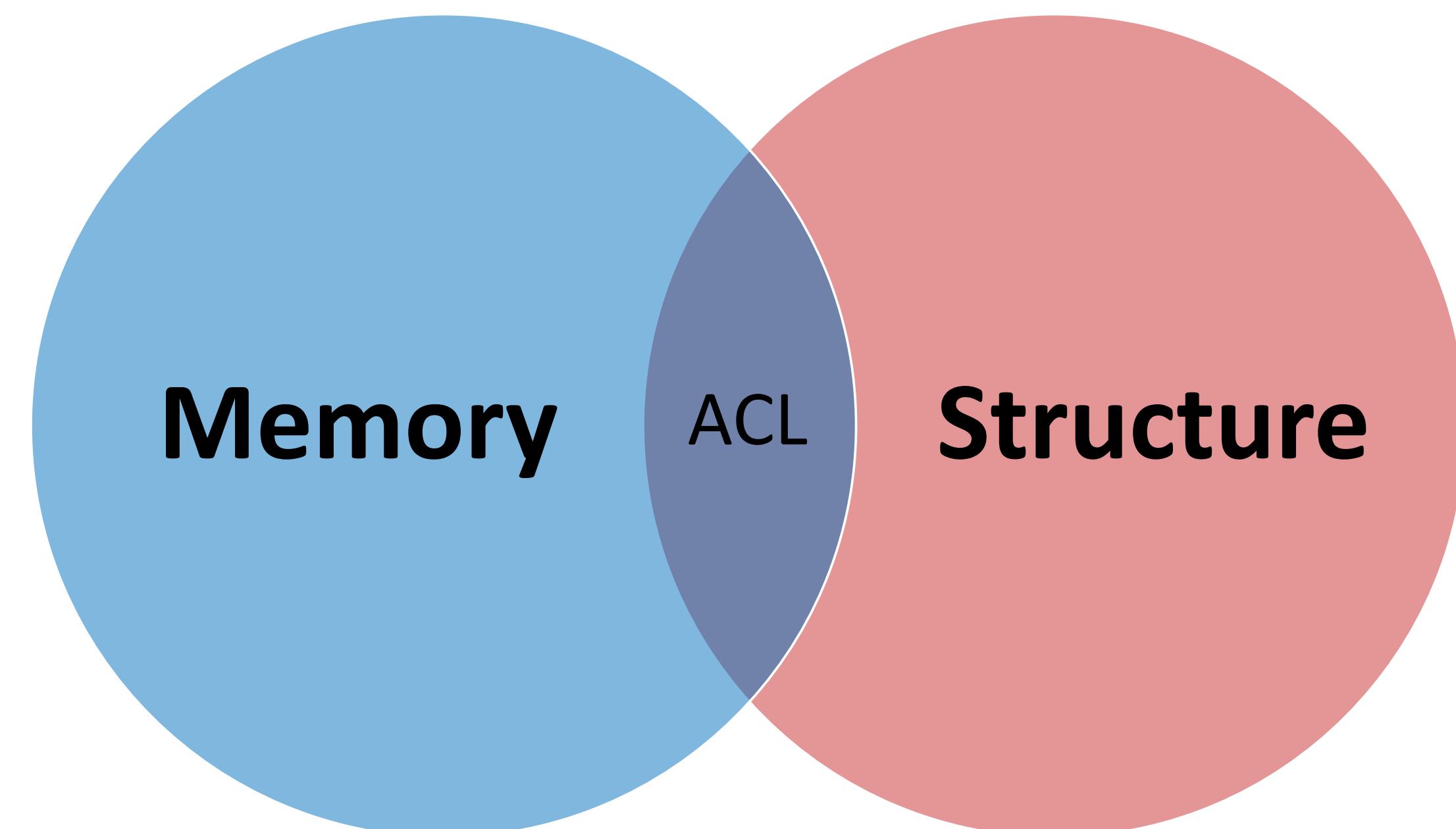
# Approaches in Continual Learning



# Approaches in Continual Learning

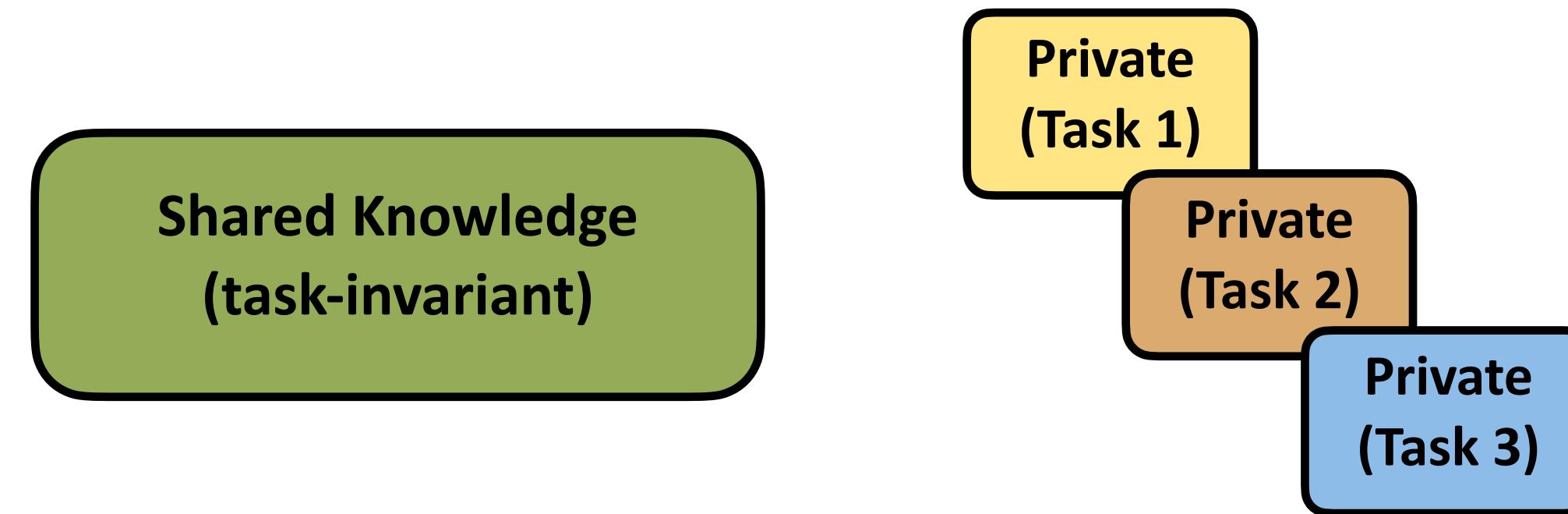


# Where does our approach (ACL) stand?



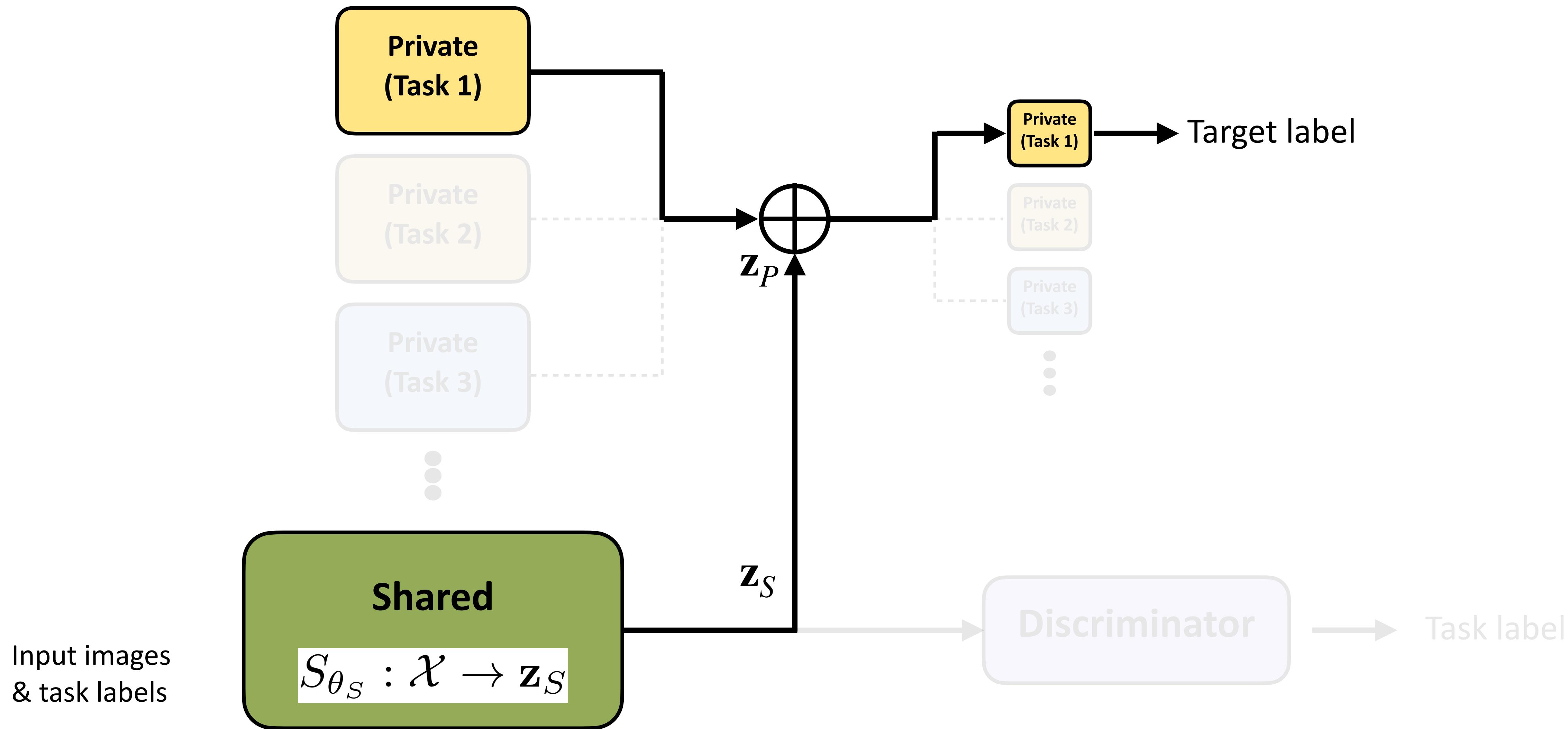
# Intuition

- Tasks in a sequence
  - Have *task-invariant (shared)* knowledge in common
  - Require *task-specific (private)* features to master them

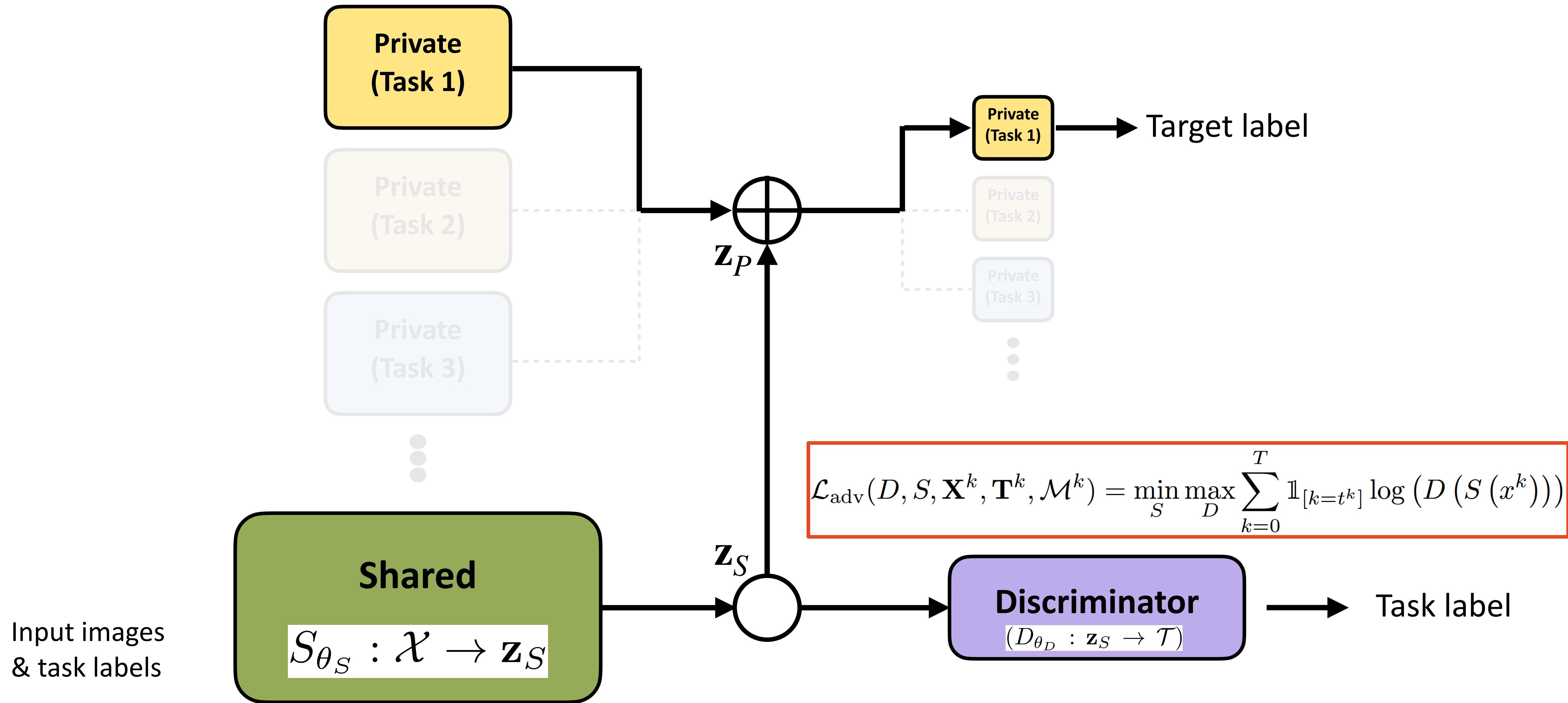


How can we factorize *task-invariant* from *task-specific* features?

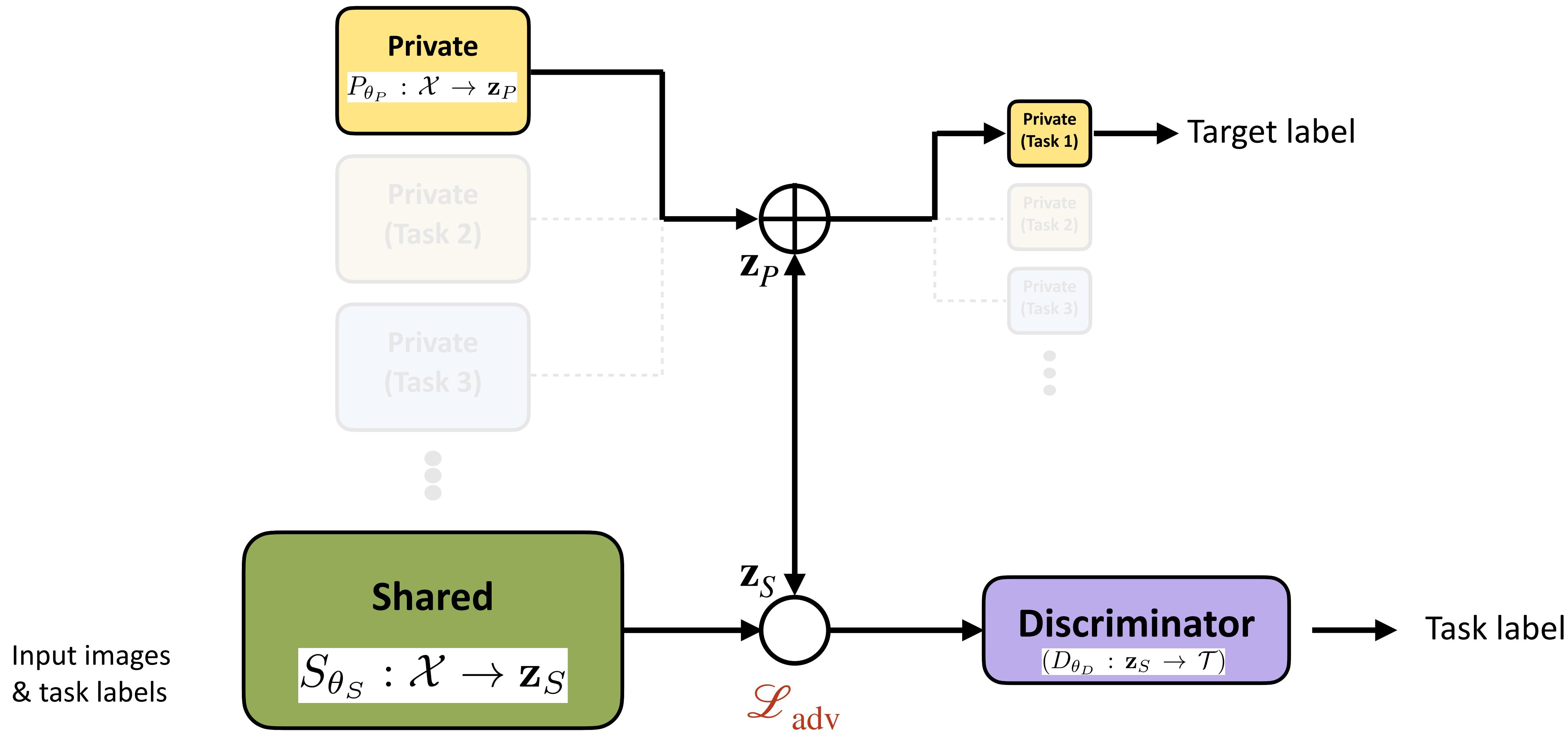
# Our Approach: Adversarial Continual Learning



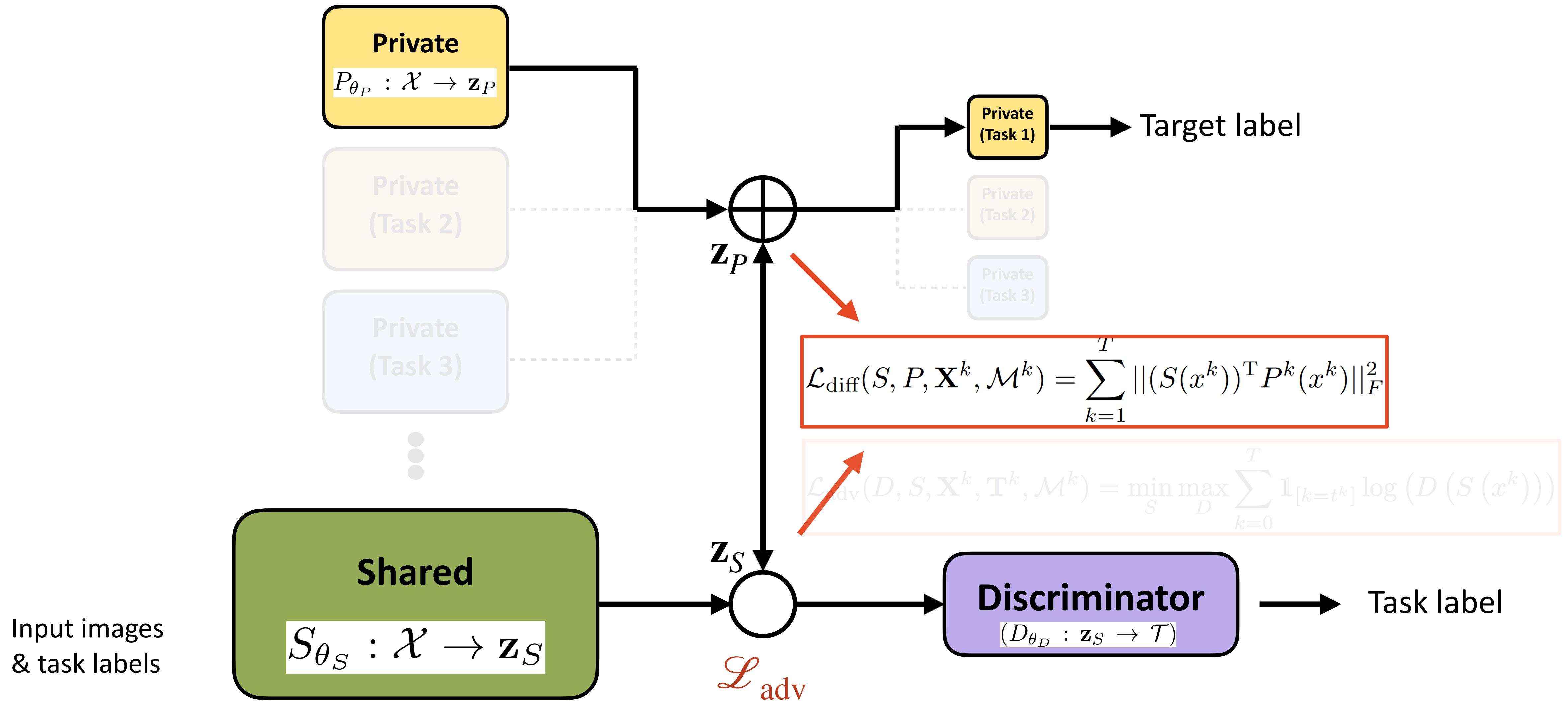
# Our Approach: Adversarial Continual Learning



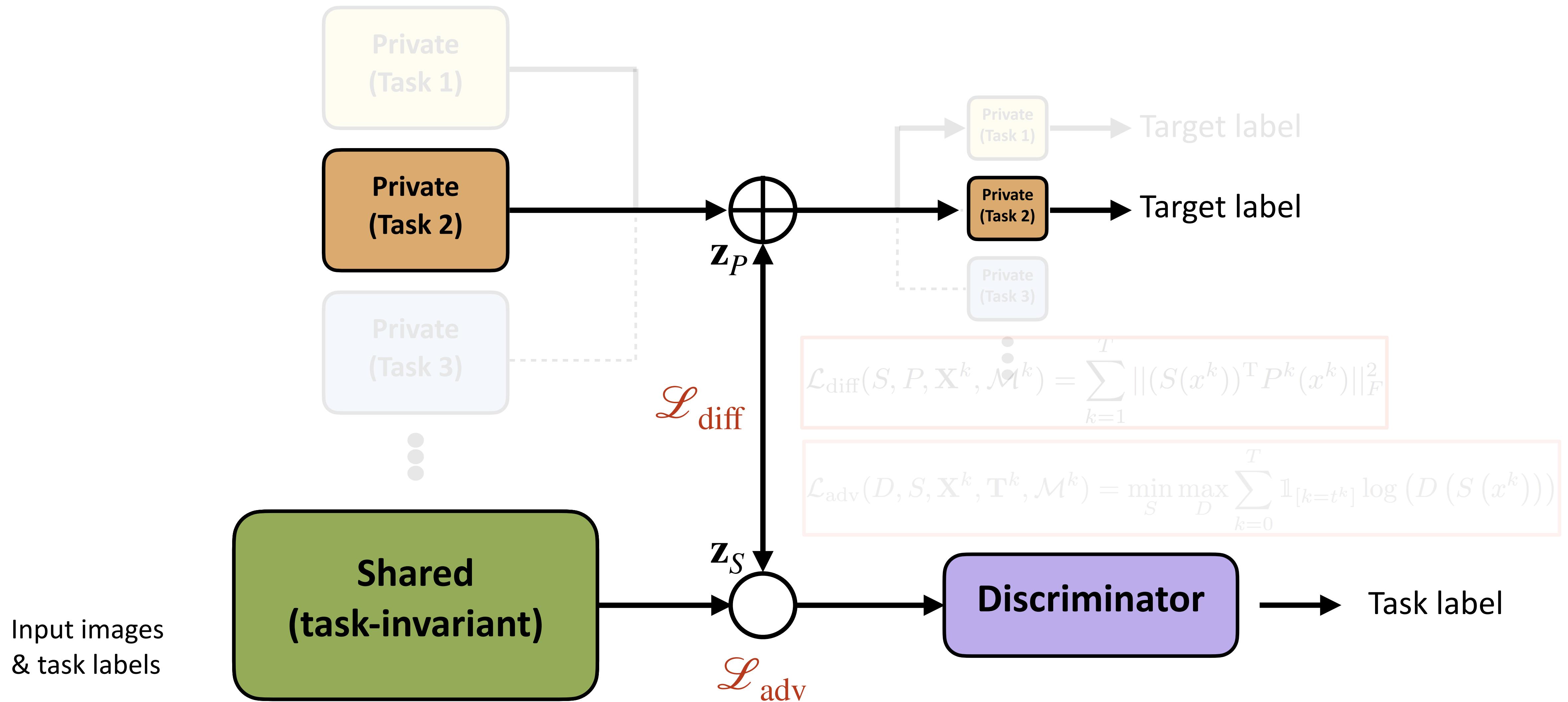
# Our Approach: Adversarial Continual Learning



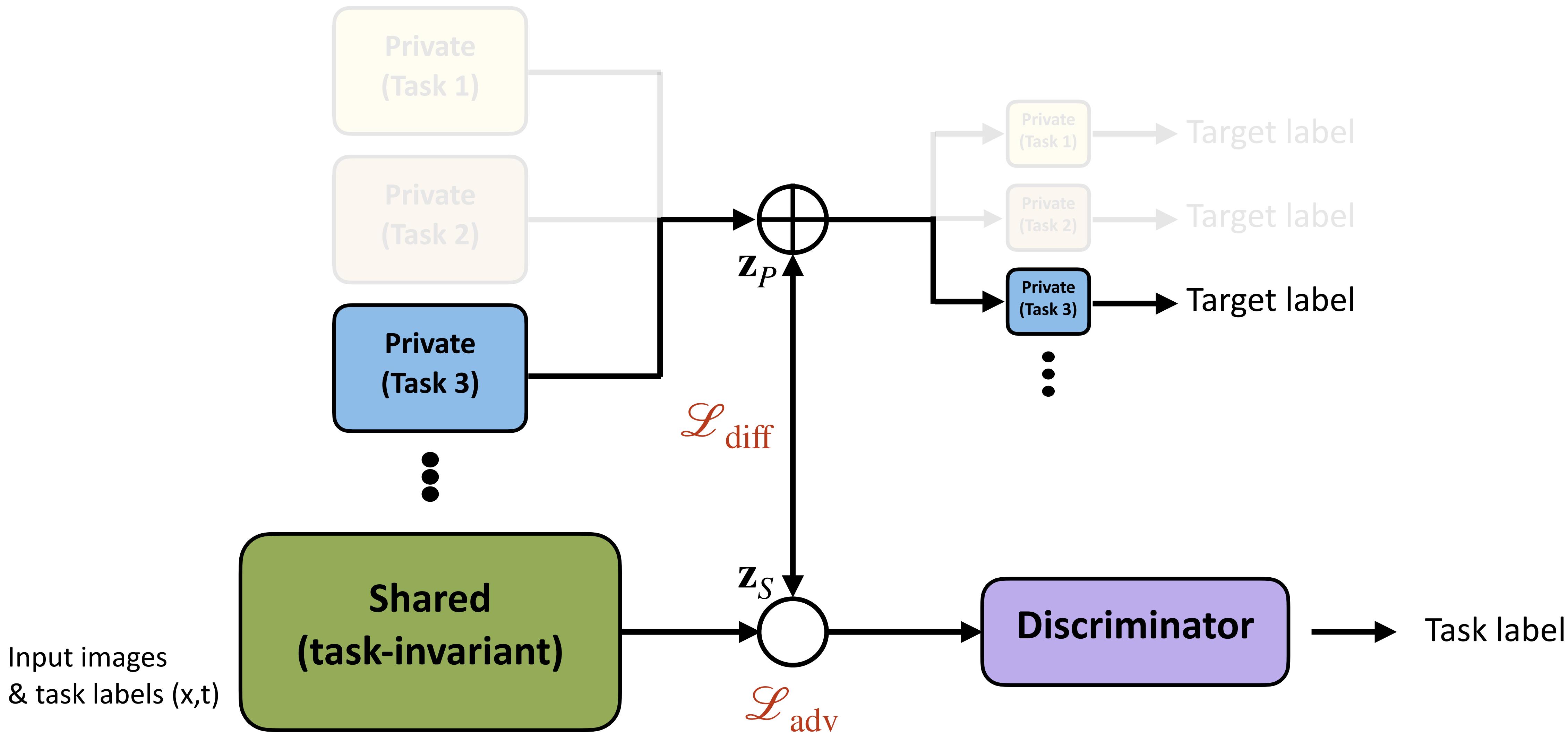
# Our Approach: Adversarial Continual Learning



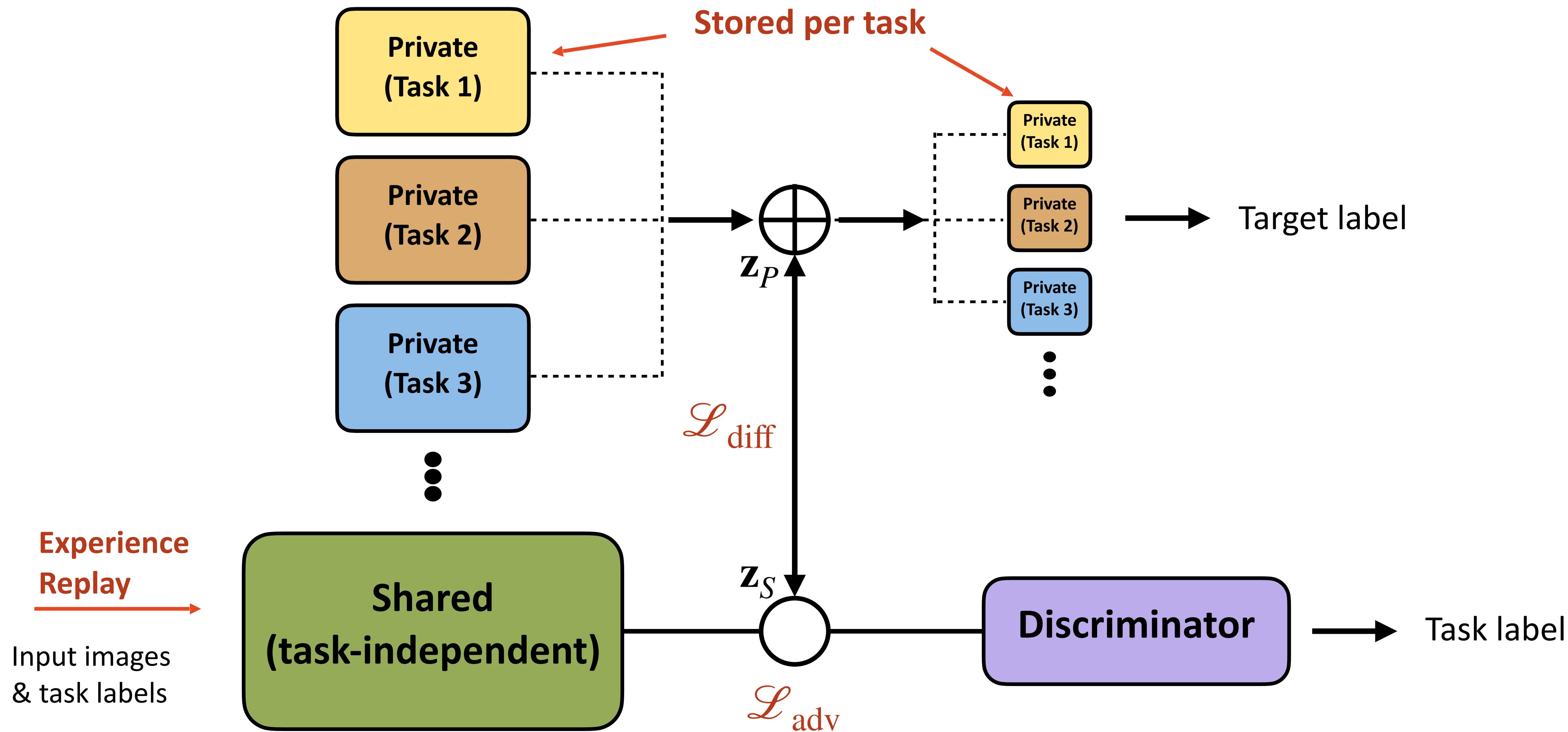
# Our Approach: Adversarial Continual Learning



# Our Approach: Adversarial Continual Learning



# Avoiding Forgetting in ACL



# Experiments

minilmageNet (20 Tasks)

CIFAR100 (20 Tasks)

Split MNIST (5 Tasks)

Permuted MNIST (10/20/30/40 Tasks)

Sequence of 5 datasets:

(SVHN, CIFAR10, MNIST, FashionMNIST, NotMNIST)

Datasets

Average Accuracy

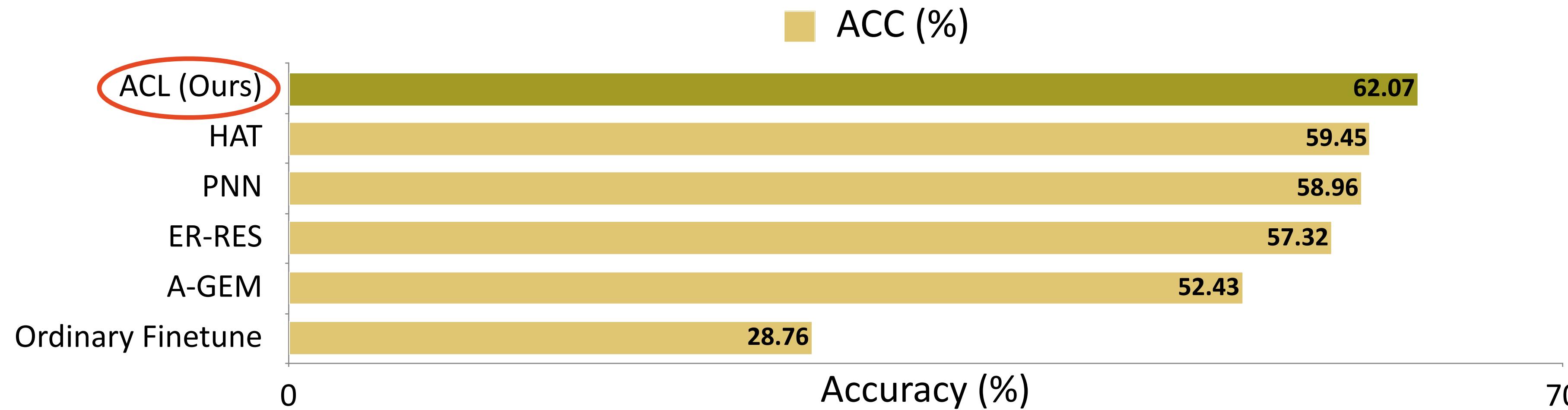
$$\text{ACC} = \frac{1}{n} \sum_{i=1}^n R_{i,n}$$

Backward Transfer:

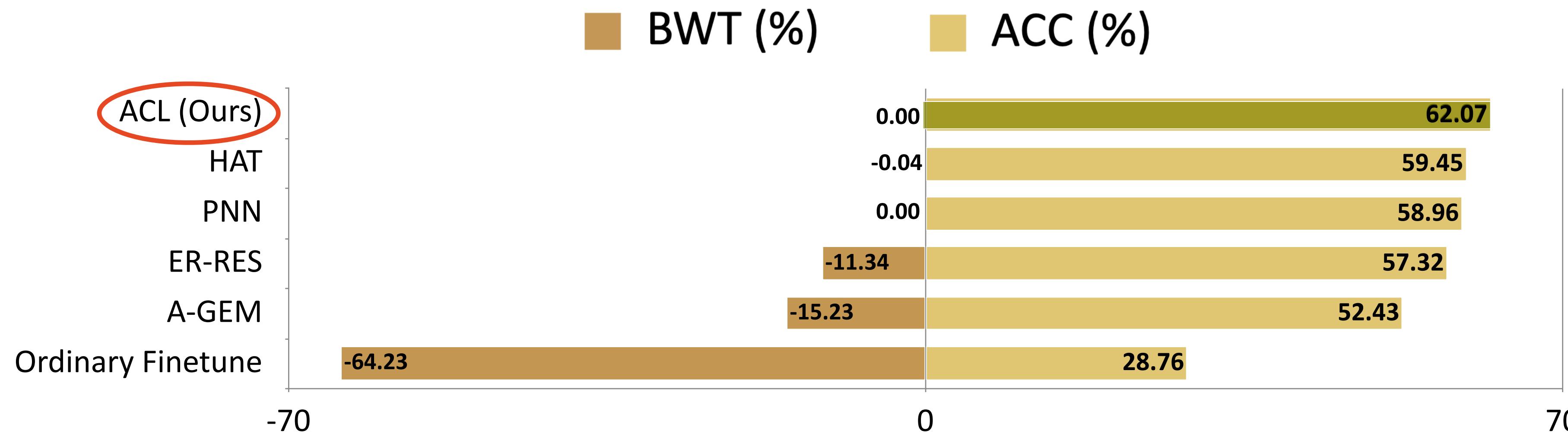
$$\text{BWT} = \frac{1}{n-1} \sum_{i=1}^n R_{i,n} - R_{i,i}$$

Evaluation metrics

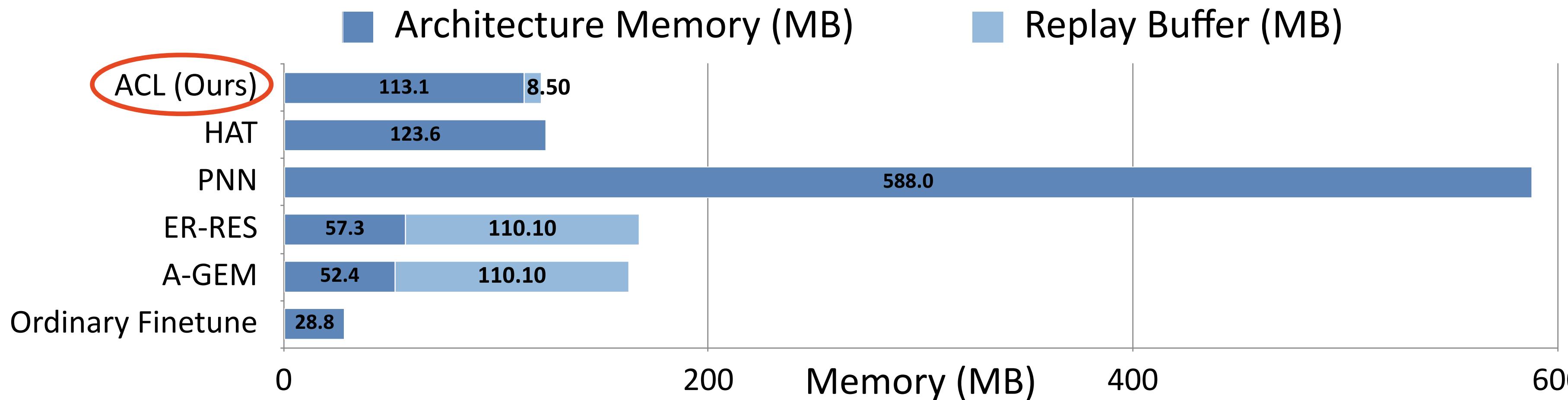
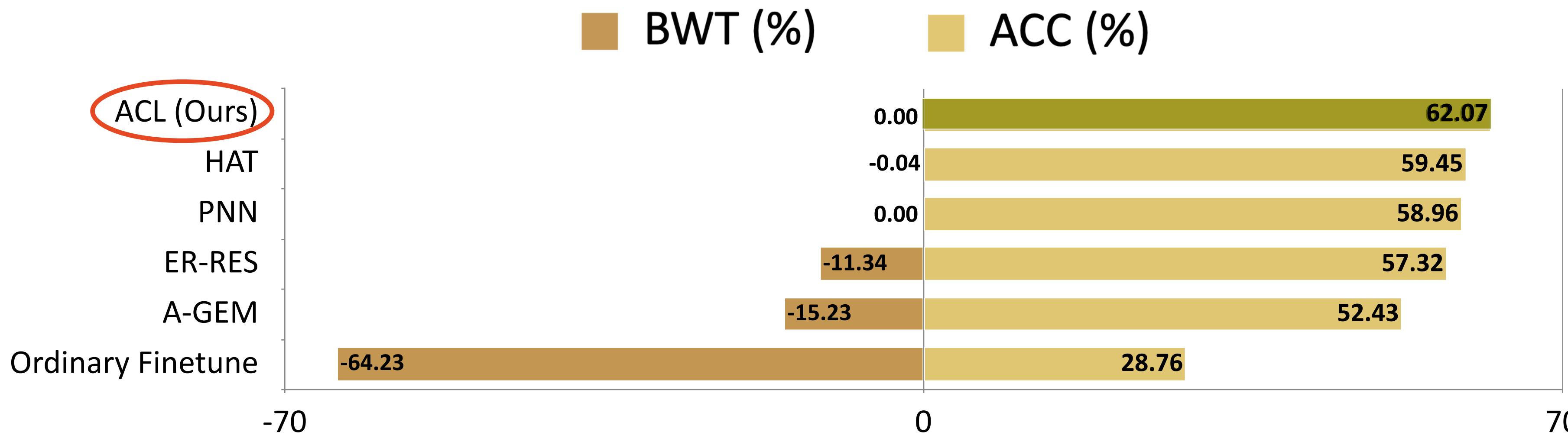
# Results on 20-Split MinilmageNet



# Results on 20-Split MinilmageNet



# Results on 20-Split MinilmageNet

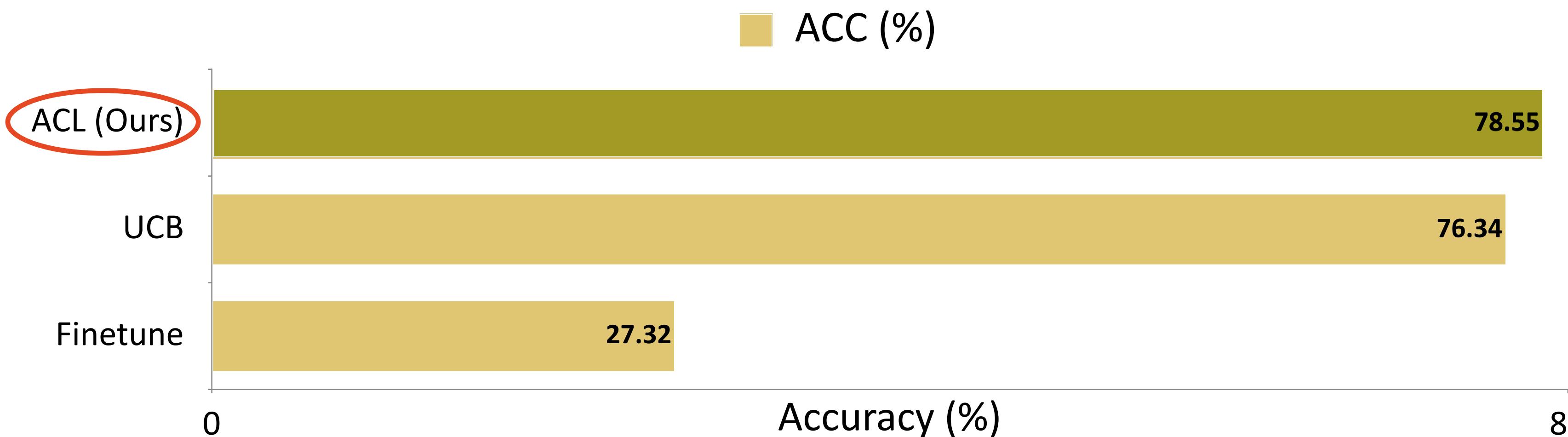


# Results on Sequence of 5-Datasets

(SVHN, CIFAR10, MNIST, FashionMNIST, NotMNIST)

# Classes	Training	Test
50	212,785	48,365

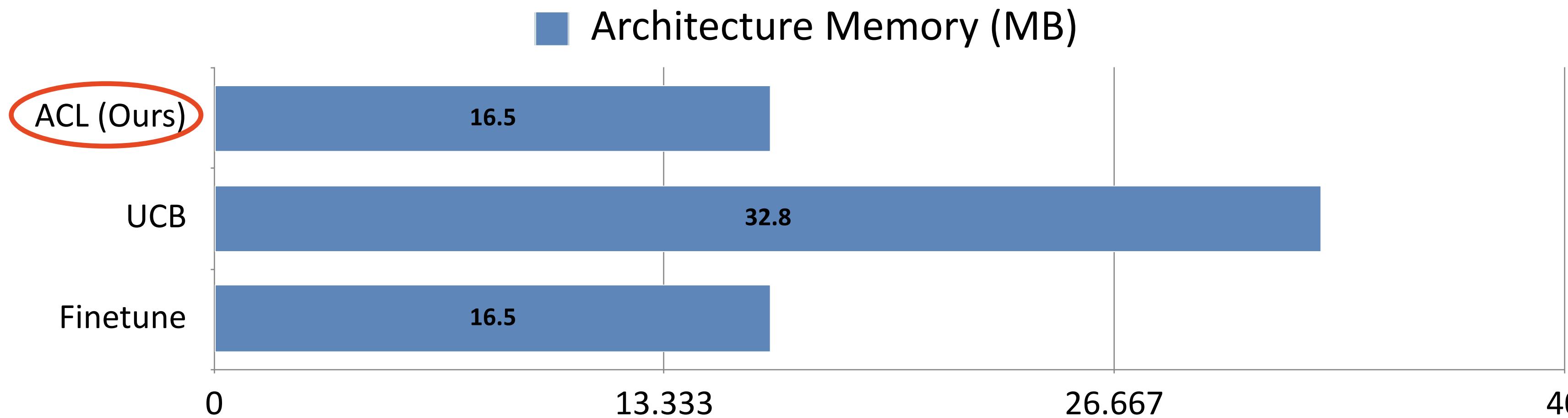
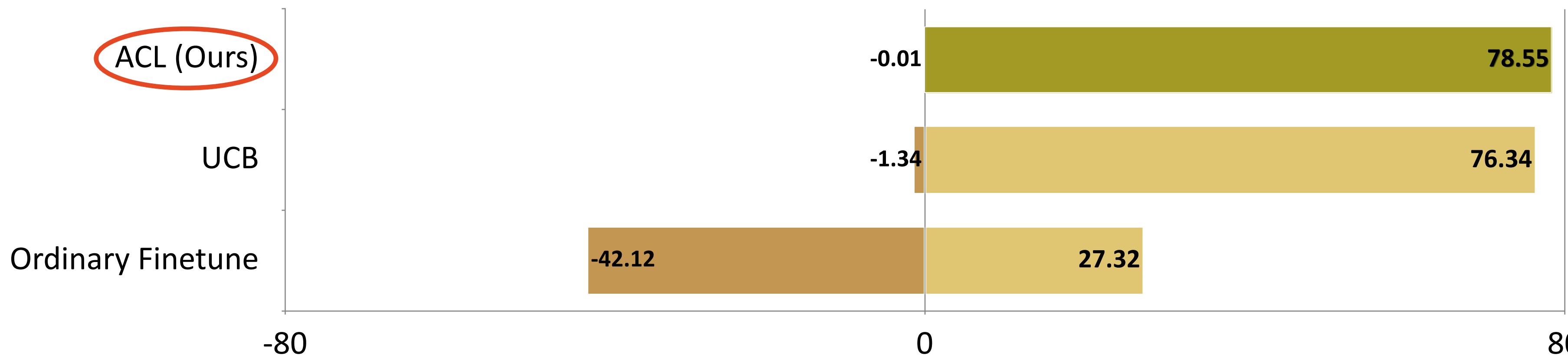
# Results on Sequence of 5-Datasets



(SVHN, CIFAR10, MNIST, FashionMNIST, NotMNIST)

# Classes	Training	Test
50	212,785	48,365

# Results on Sequence of 5-Datasets



# Ablation Study on 20-Split minilmageNet

Discriminator	Replay buffer	$\mathcal{L}_{\text{diff}}$	ACC (%)	BWT (%)
x	x	x	62.07	0.00

# Ablation Study on 20-Split minilmageNet

Discriminator	Replay buffer	$\mathcal{L}_{\text{diff}}$	ACC (%)	BWT (%)
-	X	X	52.07	-0.01
X	X	X	<b>62.07</b>	<b>0.00</b>

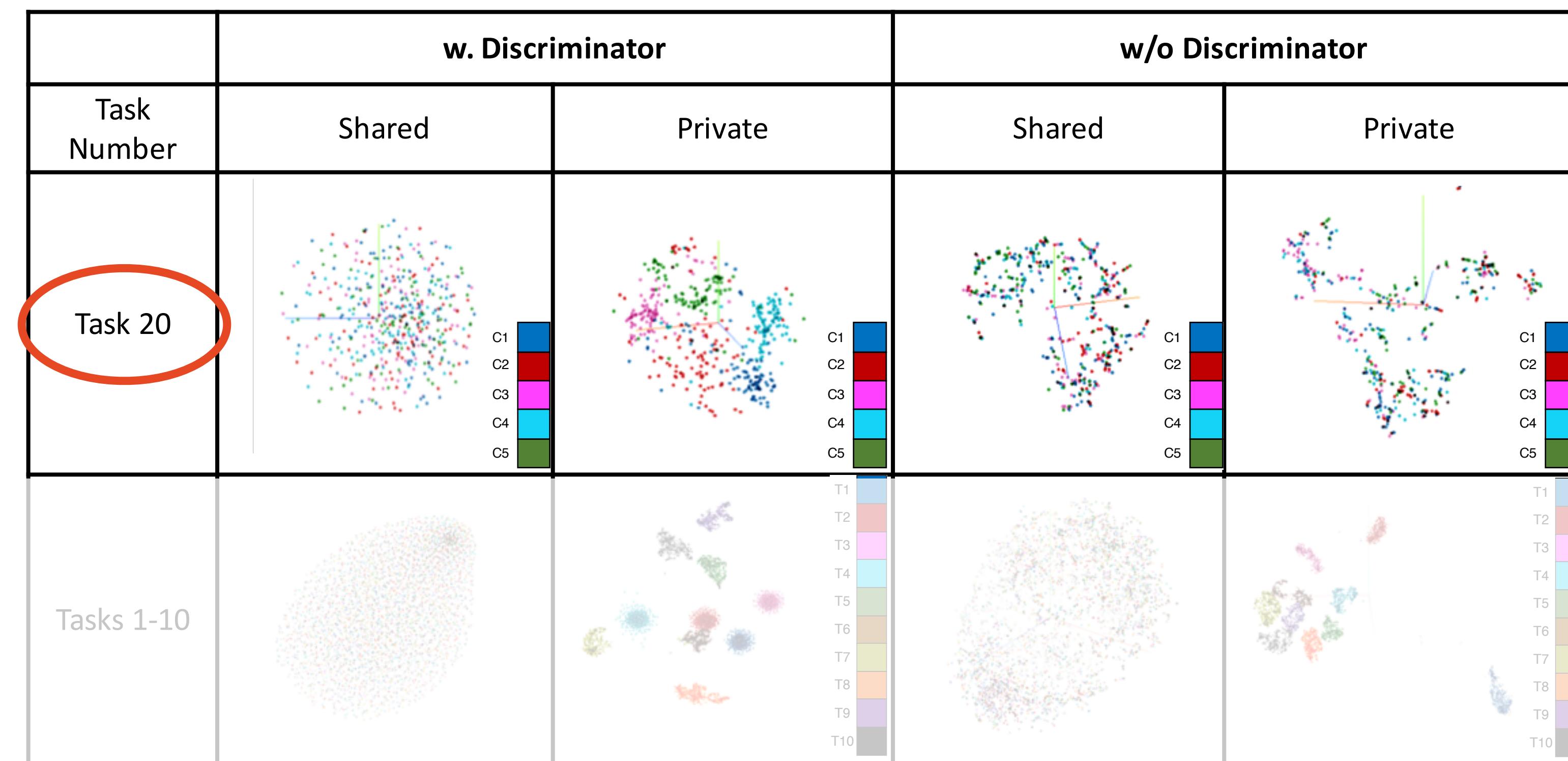
# Ablation Study on 20-Split minilmageNet

Discriminator	Replay buffer	$\mathcal{L}_{\text{diff}}$	ACC (%)	BWT (%)
-	X	X	52.07	-0.01
X	-	X	57.66	-3.71
X	X	X	<b>62.07</b>	<b>0.00</b>

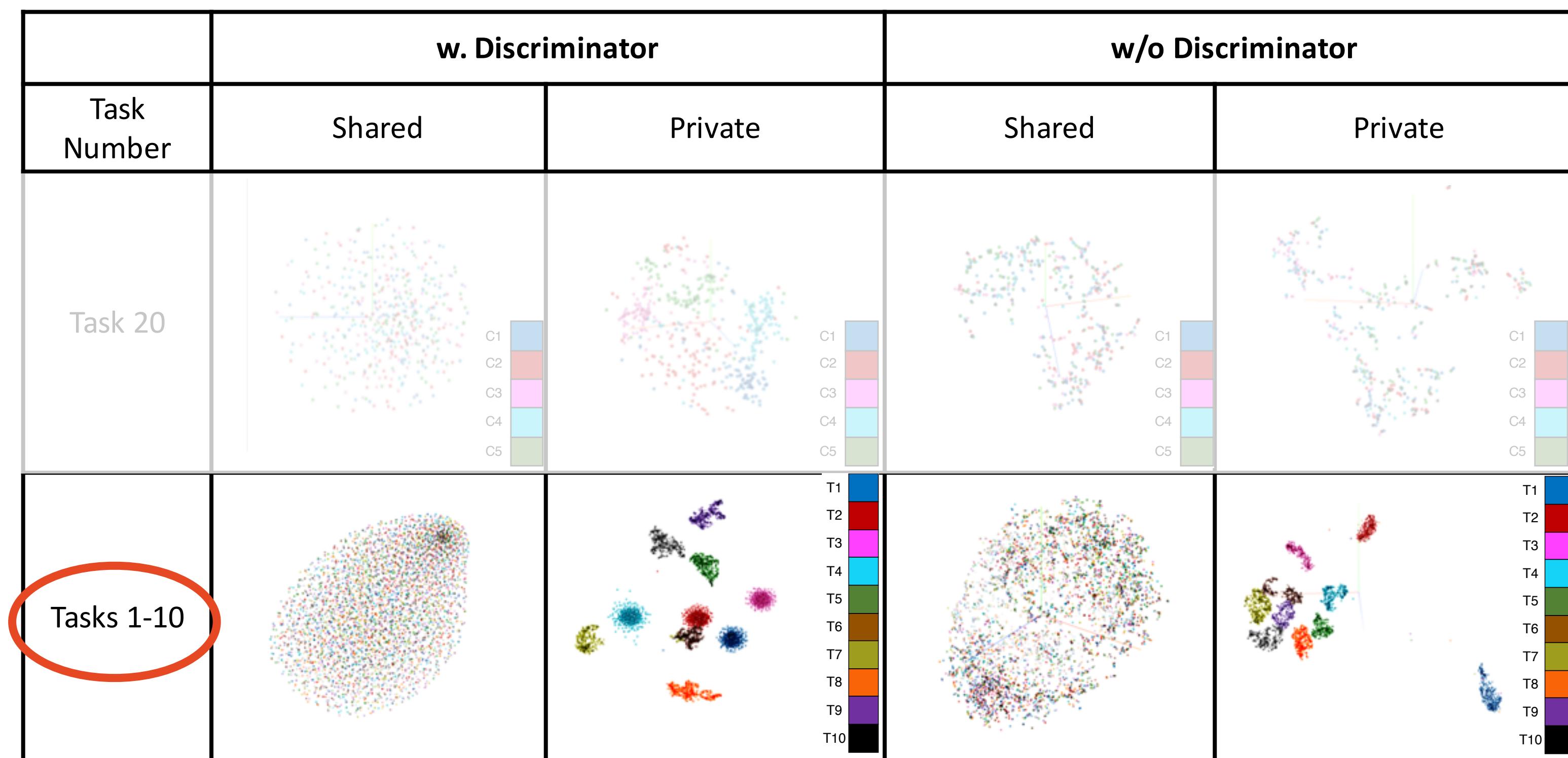
# Ablation Study on 20-Split minilmageNet

Discriminator	Replay buffer	$\mathcal{L}_{\text{diff}}$	ACC (%)	BWT (%)
	X	X	52.07	-0.01
X	X		57.66	-3.71
X	X	-	60.28	0.00
X	X	X	<b>62.07</b>	<b>0.00</b>

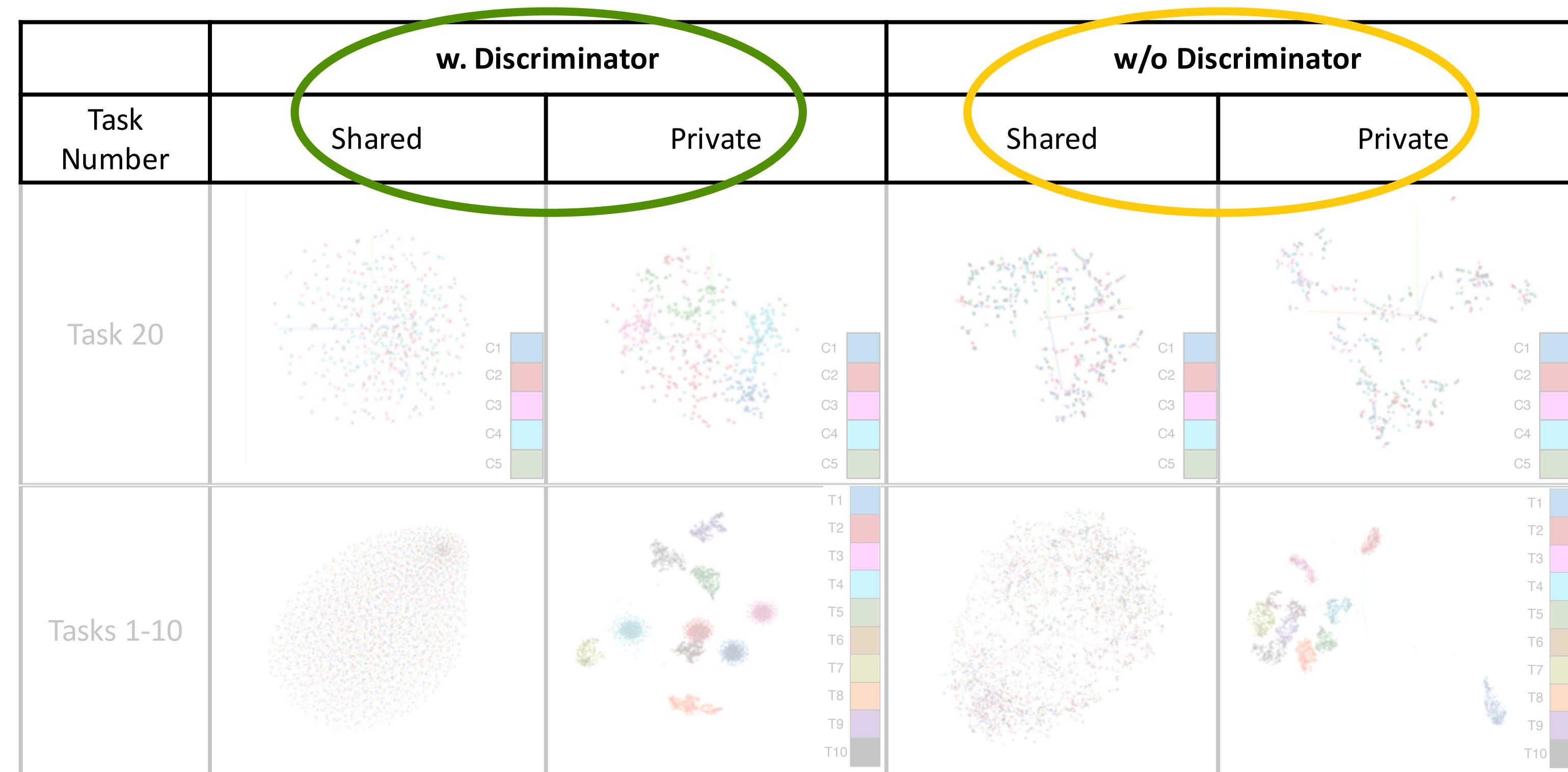
# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



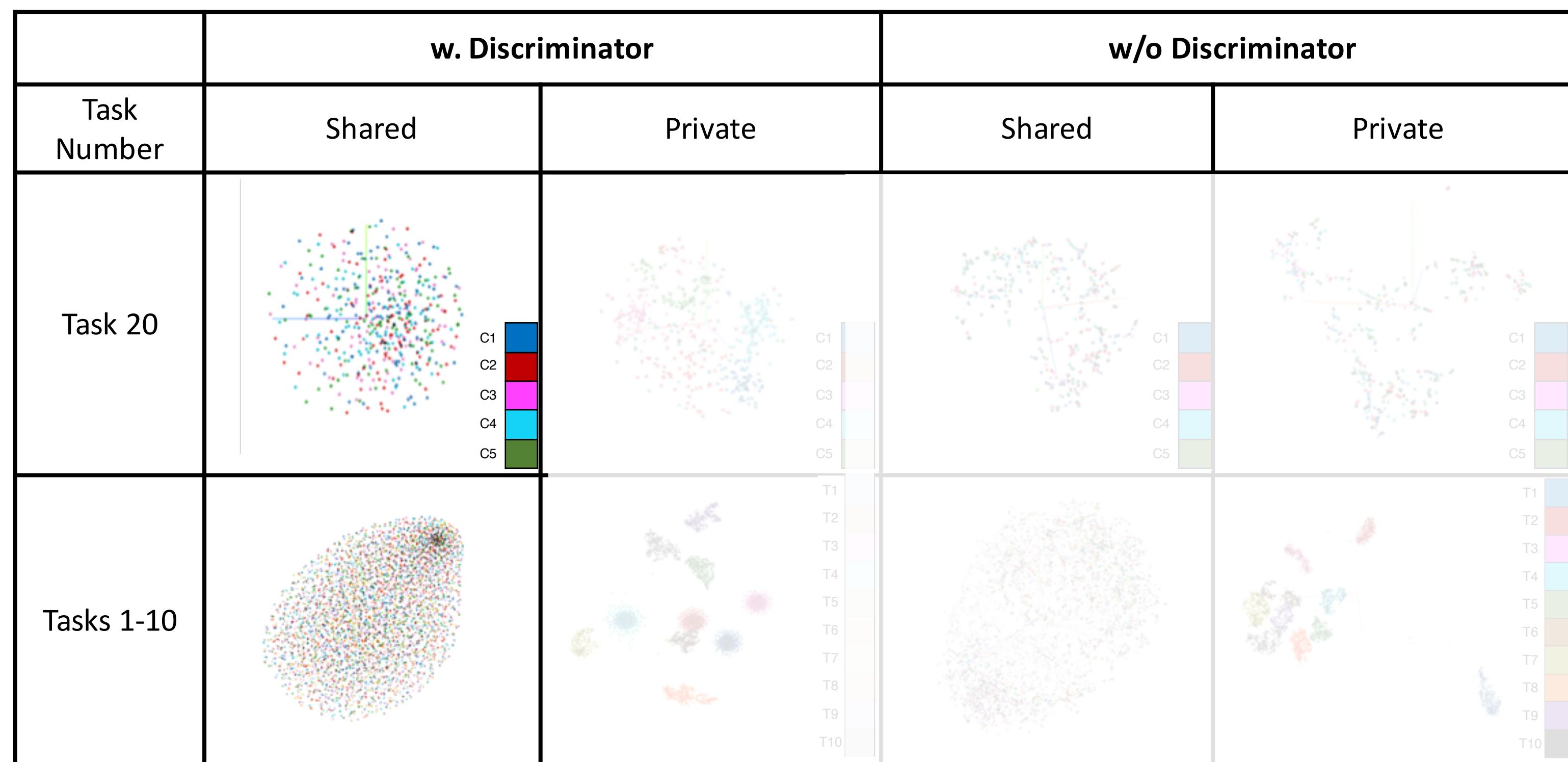
# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



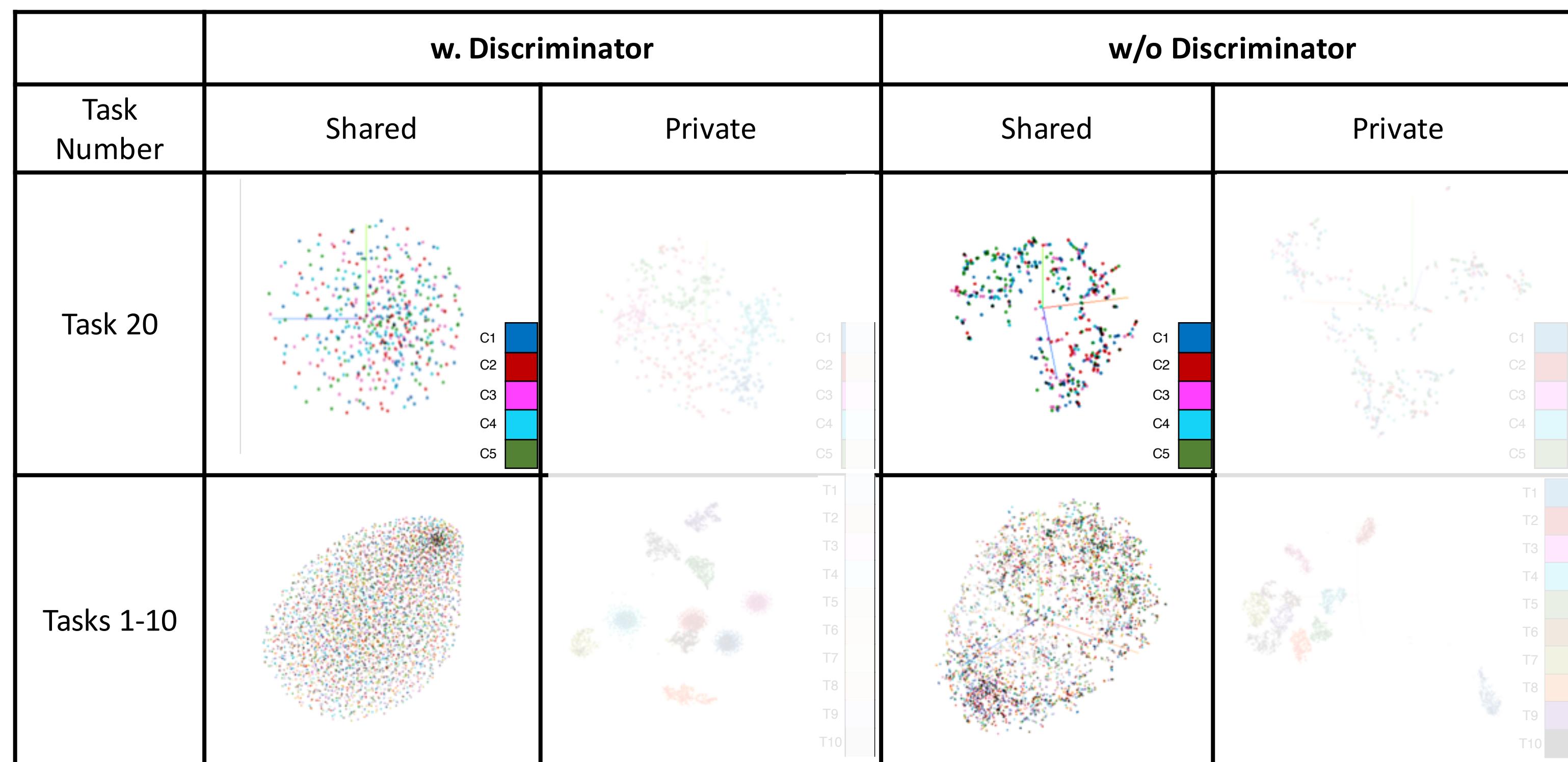
# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



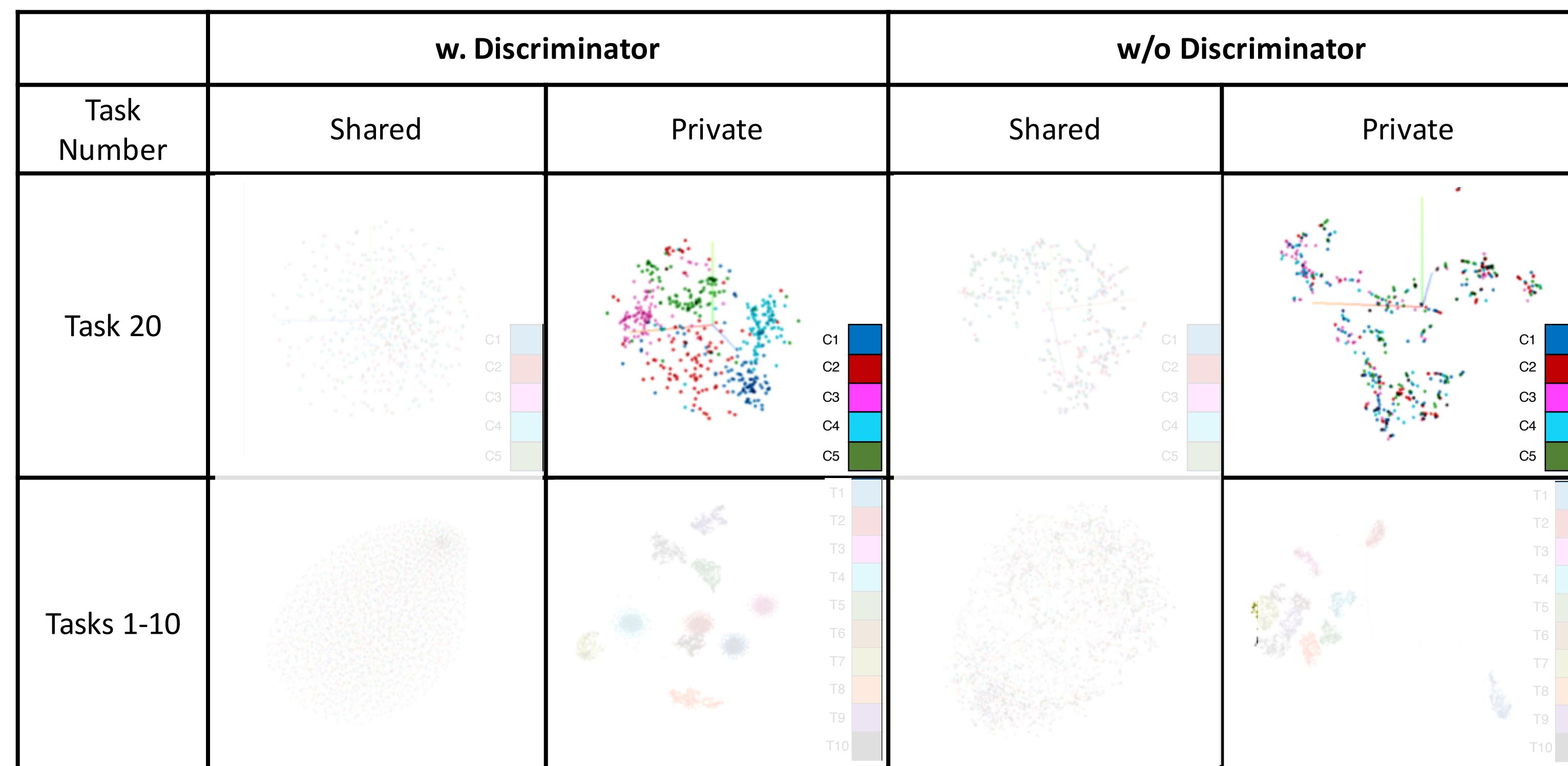
# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



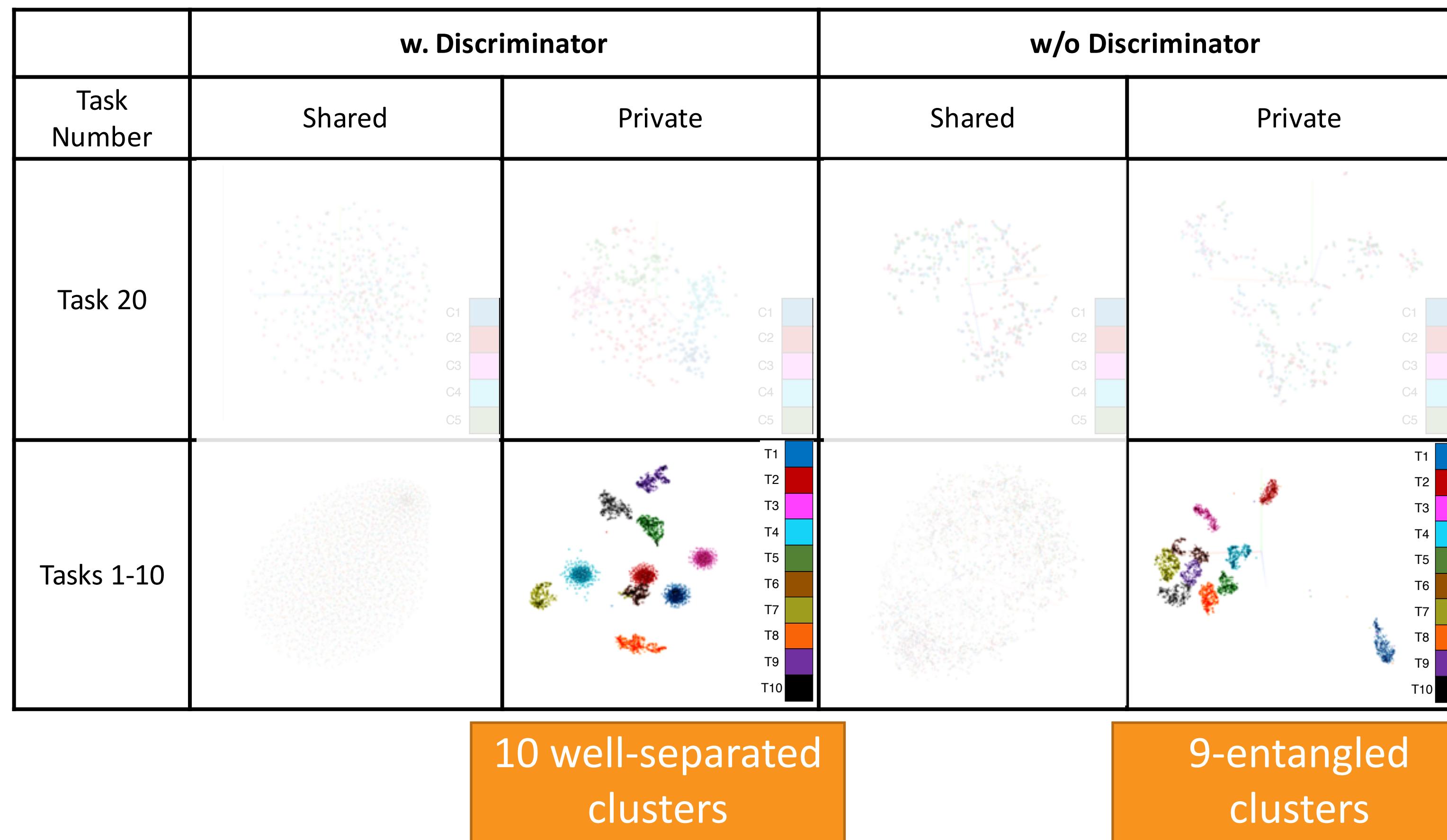
# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



# Visualizing Adversarial Learning Effect (20-Split minilmageNet)



# Conclusion

- ACL is primarily an architecture-based method but can benefit from experience replay if need be
- Uses adversarial learning and an orthogonality constraint to disentangle task-specific and task-invariant features
- Achieves near zero forgetting and state-of-the-art accuracy on image classification benchmarks

Paper: <https://arxiv.org/pdf/2003.09553.pdf>  
Code: <https://github.com/facebookresearch/Adversarial-Continual-Learning>

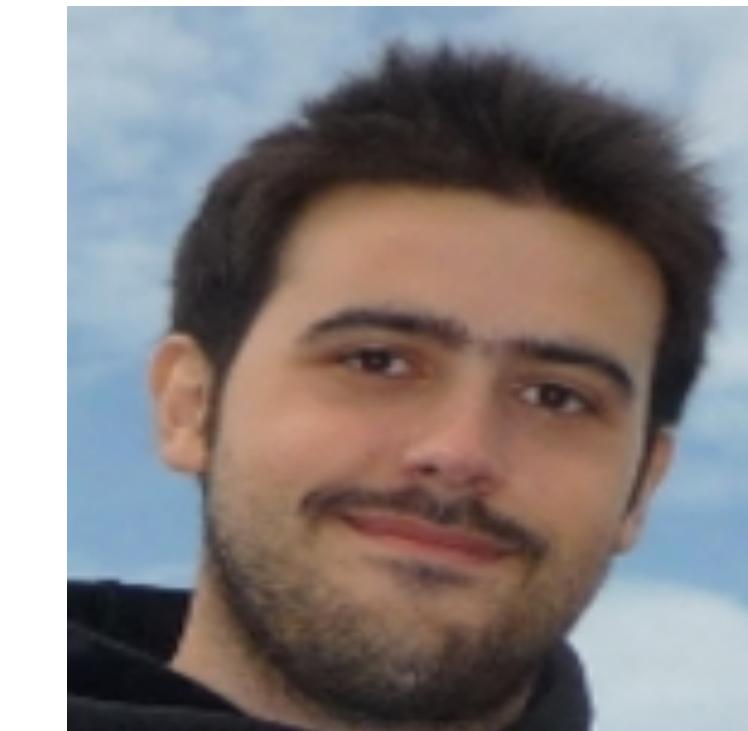
Questions: [sayna@berkeley.edu](mailto:sayna@berkeley.edu)



Sayna Ebrahimi  
UC Berkeley



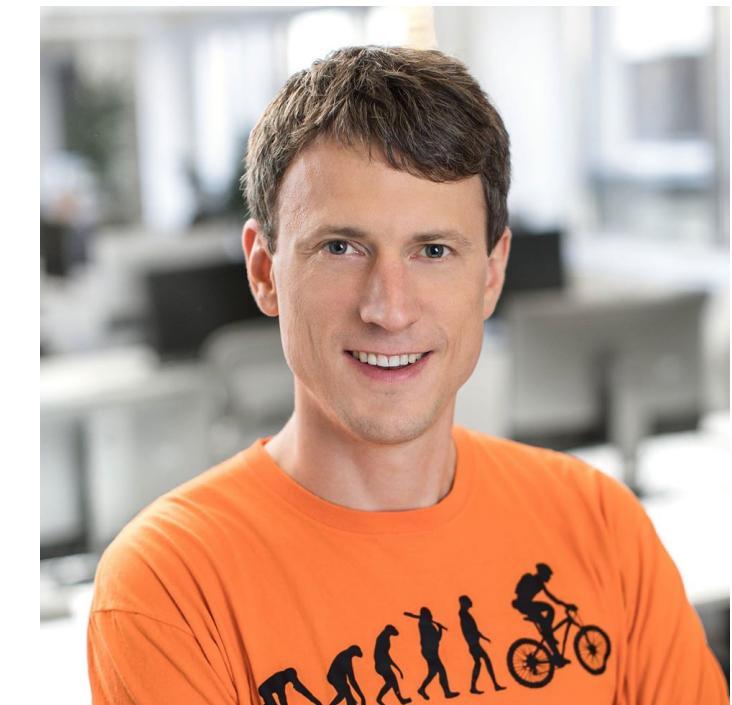
Franziska Meier  
Facebook AI Research



Roberto Calandra  
Facebook AI Research



Trevor Darrell  
UC Berkeley



Marcus Rohrbach  
Facebook AI Research