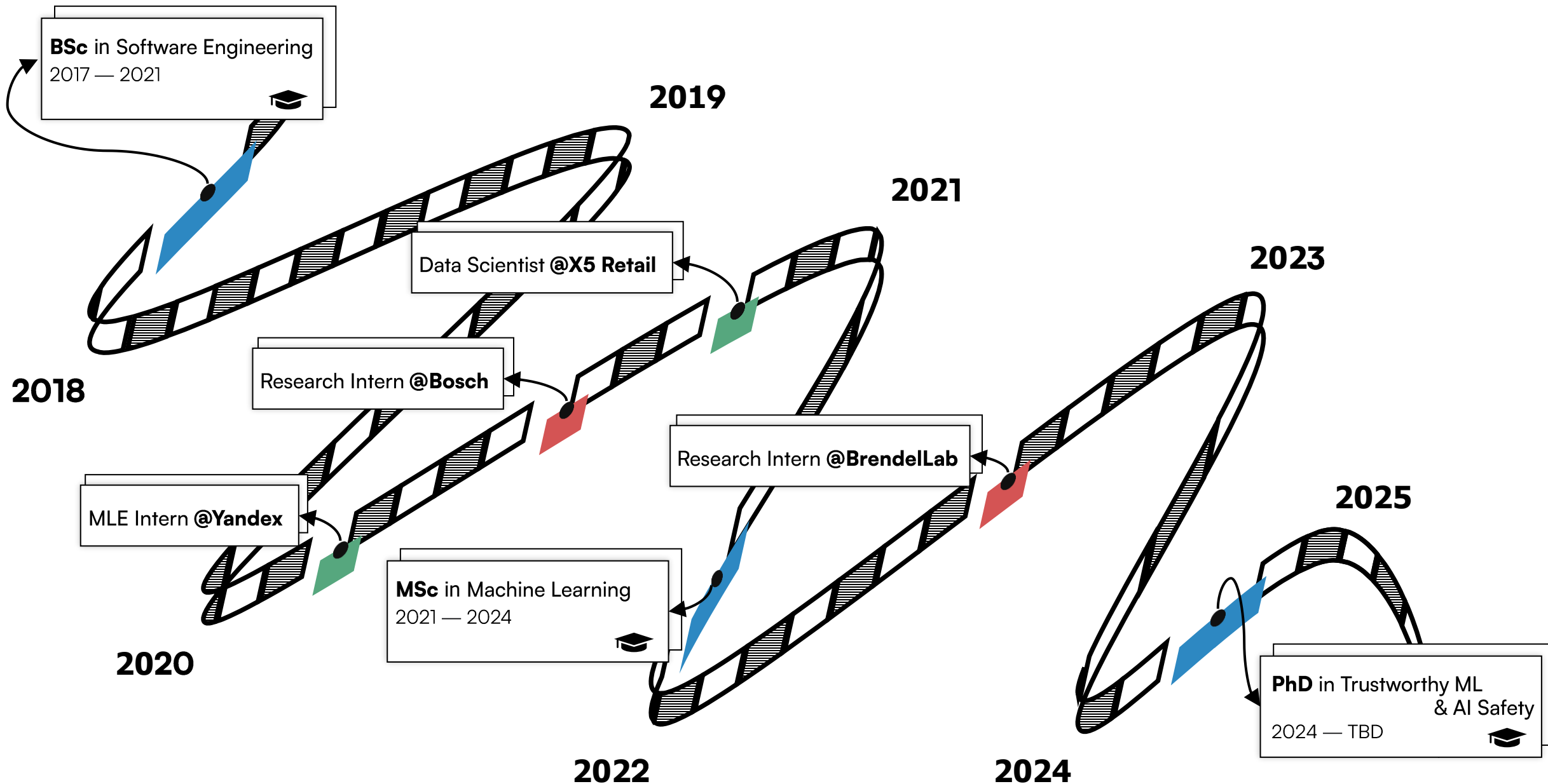


# Out-of-the- $(l_p)$ -Box:

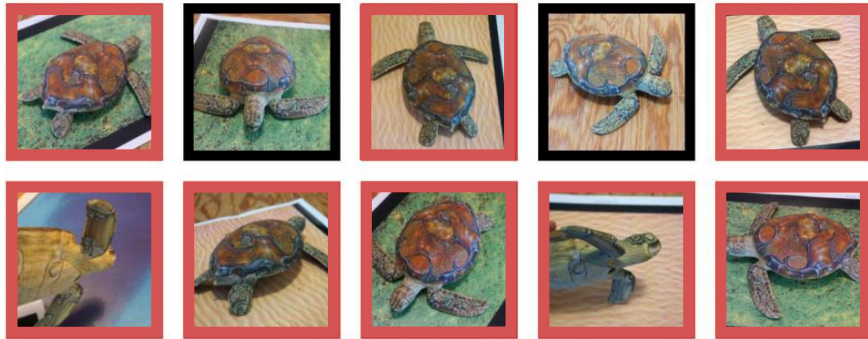
**Exploiting Adversarials,  
Exploring Compositionality,  
and Exposing New AI Threats**

IMPRS Scientific Talk by Alexander Panfilov  
02/02/2024, Tübingen

**Background**

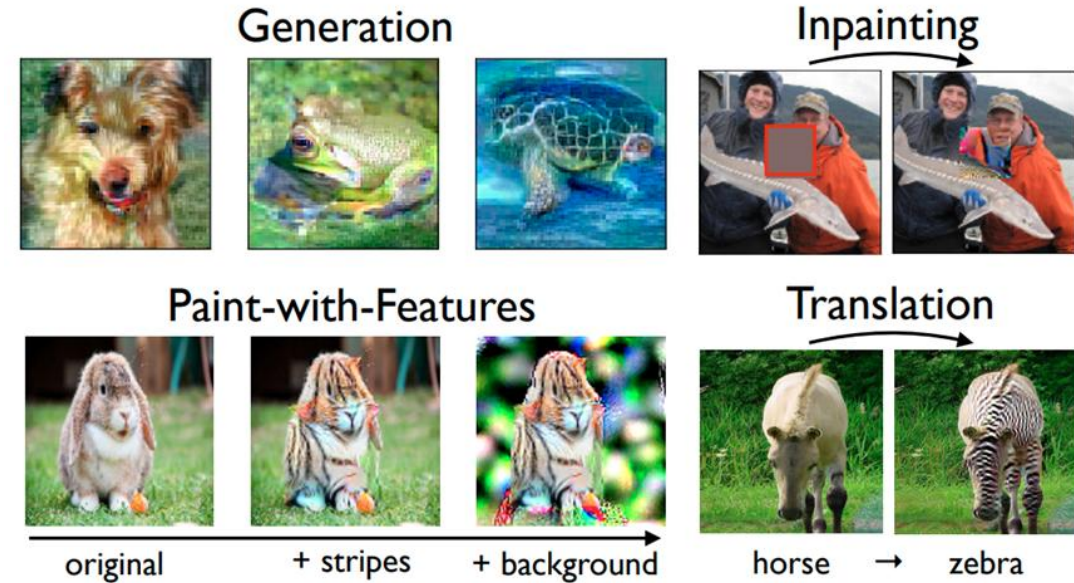


# Adversarial Threats and Robustness



■ classified as turtle    ■ classified as rifle  
■ classified as other

## Safety Problems<sup>[1]</sup>

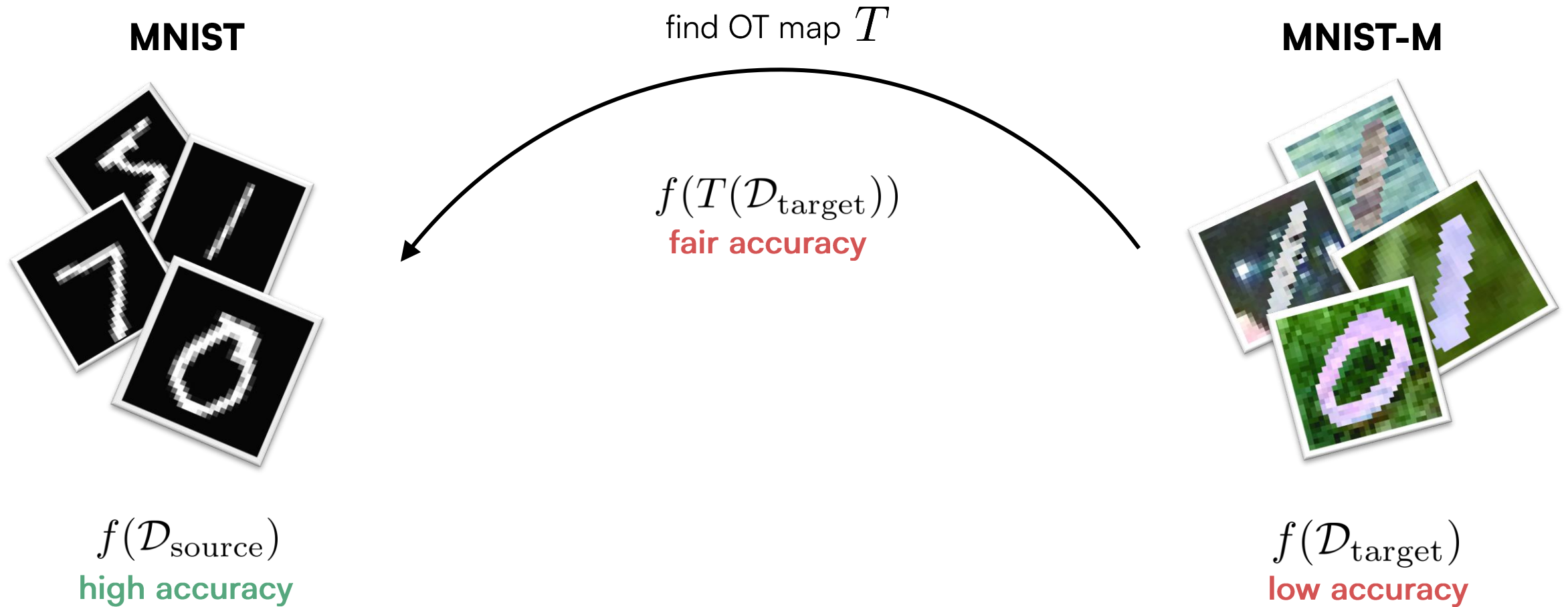


## Intriguing Properties<sup>[2]</sup>

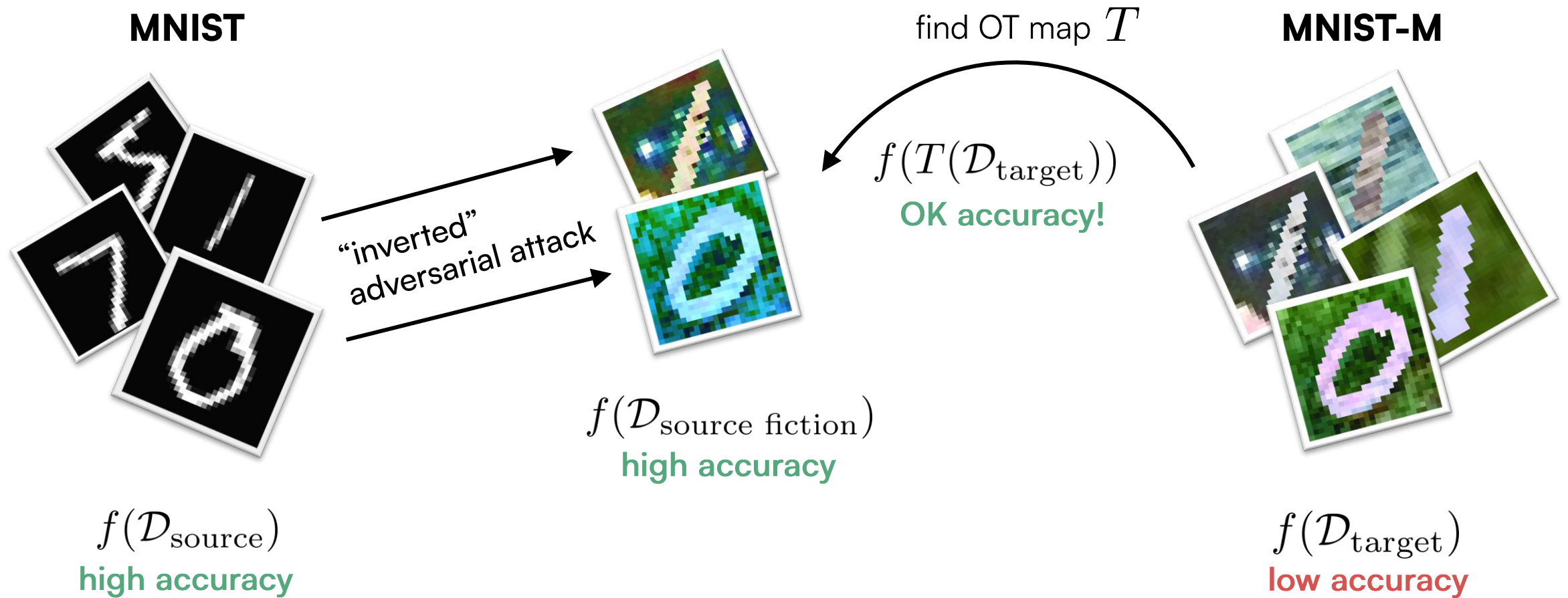
[1] Athalye, A., Engstrom, L., Ilyas, A., & Kwok, K. (2018). Synthesizing robust adversarial examples. In International Conference on Machine Learning (pp. 284-293). PMLR.

[2] Santurkar, S., Ilyas, A., Tsipras, D., Engstrom, L., Tran, B., & Madry, A. (2019). Image synthesis with a single (robust) classifier. Advances in Neural Information Processing Systems, 32.

# “Easy Features” for Domain Adaptation with OT



# “Easy Features” for Domain Adaptation with OT





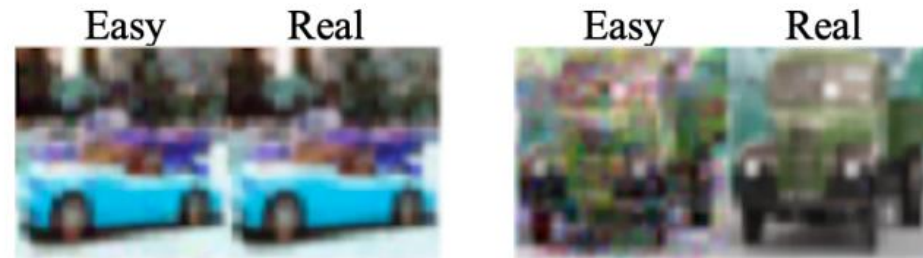
# “Easy Features” for Domain Adaptation with OT

Method	M / S	S / M	M / U	M / MM
EMD	21.2 ±3	68.7±3	79.2±2	56.1±3
EMD( <i>sf</i> )	<b>23.0±3</b>	<b>86.3±3</b>	<b>83.1±2</b>	<b>62.7±2</b>
OTLin	21.8±4	69.9±4	84.1±7	62.3±1
OTLin( <i>sf</i> )	<b>25.5±4</b>	<b>88.4±4</b>	<b>89.3±6</b>	<b>64.5±3</b>
Sinkh	21.8±4	68.8±2	82.1±7	55.7±12
Sinkh( <i>sf</i> )	<b>25.5±4</b>	<b>86.2±4</b>	<b>83.8±6</b>	<b>62.9±4</b>
SinkhLp	21.8±4	68.8±6	84.8±16	55.7±19
SinkhLp( <i>sf</i> )	<b>25.5±4</b>	<b>86.3±7</b>	<b>88.3±19</b>	<b>63.0±27</b>
SinkhL2	21.8±4	68.8±4	84.8±2	55.7±4
SinkhL2( <i>sf</i> )	<b>25.5±4</b>	<b>86.3±2</b>	<b>88.3±2</b>	<b>63.0±2</b>

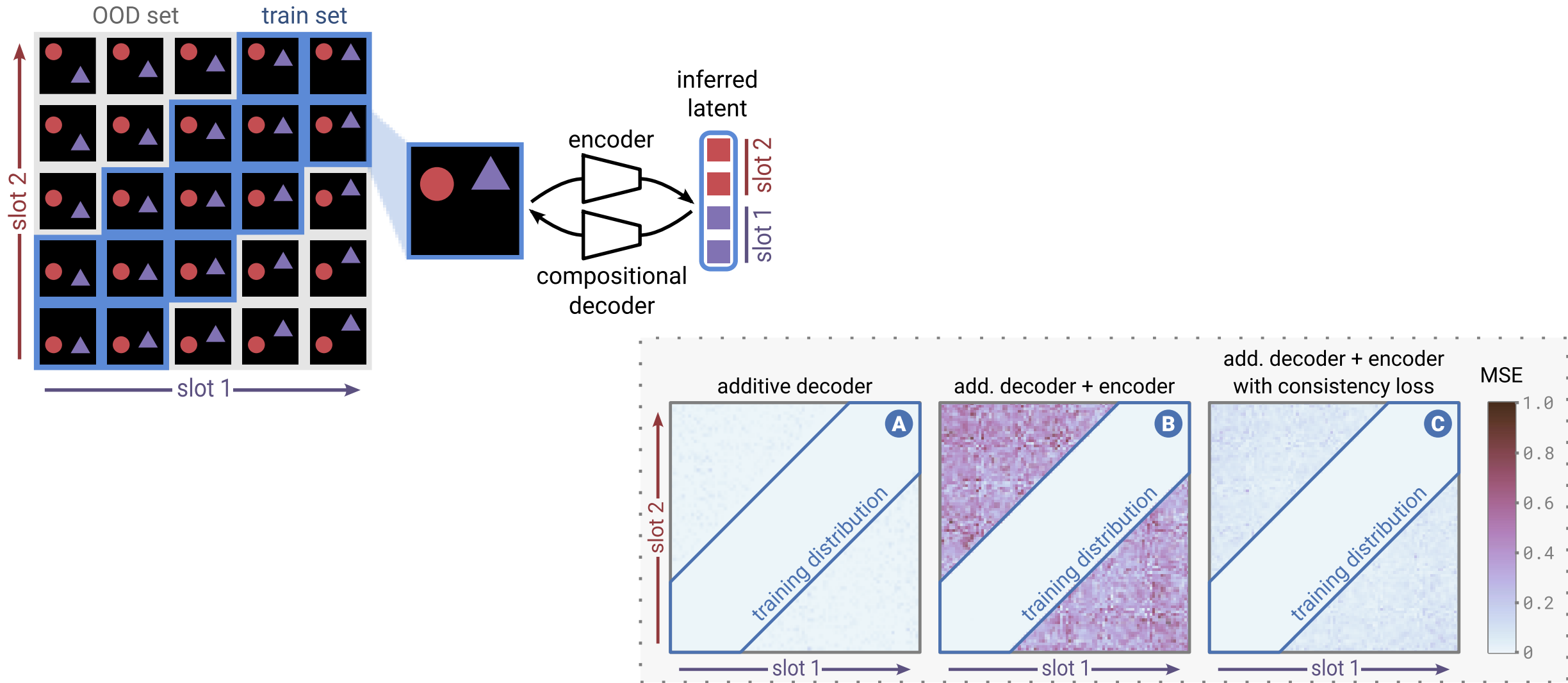
(a) Digits dataset domains.

Method	A / S	S / A	A / W	W / A	S / W	W / S
EMD	38.4±3	9.3±5	45.2±3	45.6±5	13.6±3	36.7±3
EMD( <i>sf</i> )	<b>56.8±3</b>	<b>29.7±4</b>	<b>64.9±3</b>	<b>73.9±4</b>	<b>40.1±3</b>	<b>60.1±2</b>
OTLin	37.1±3	11.0±3	38.7±3	47.5±3	6.2±3	39.6±4
OTLin( <i>sf</i> )	<b>58.5±3</b>	<b>29.8±3</b>	<b>65.2±3</b>	<b>74.4±3</b>	<b>40.1±3</b>	<b>63.1±5</b>
Sinkh	38.0±3	10.1±4	44.7±6	45.5±3	13.1±7	37.2±3
Sinkh( <i>sf</i> )	<b>57.0±3</b>	<b>31.0±4</b>	<b>65.2±7</b>	<b>73.9±3</b>	<b>39.9±4</b>	<b>60.0±2</b>
SinkhLp	38.1±6	10.4±8	45.2±7	45.3±5	13.1±7	37.2±3
SinkhLp( <i>sf</i> )	<b>57.2±6</b>	<b>31.0±11</b>	<b>65.2±8</b>	<b>74.0±5</b>	<b>40.1±5</b>	<b>60.1±4</b>
SinkhL2	38.1±4	10.4±7	45.0±4	45.3±6	13.1±6	37.2±3
SinkhL2( <i>sf</i> )	<b>57.2±4</b>	<b>31.0±7</b>	<b>65.2±4</b>	<b>74.0±6</b>	<b>40.1±5</b>	<b>60.1±4</b>

(b) Modern Office-31 dataset domains.

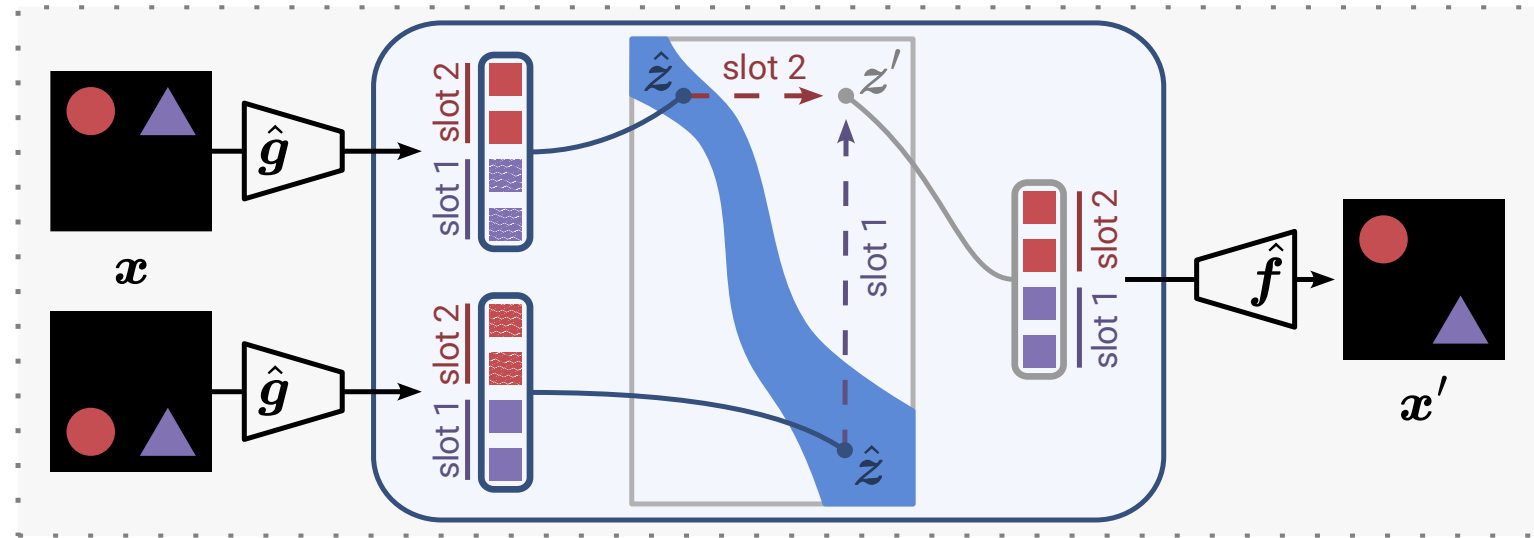
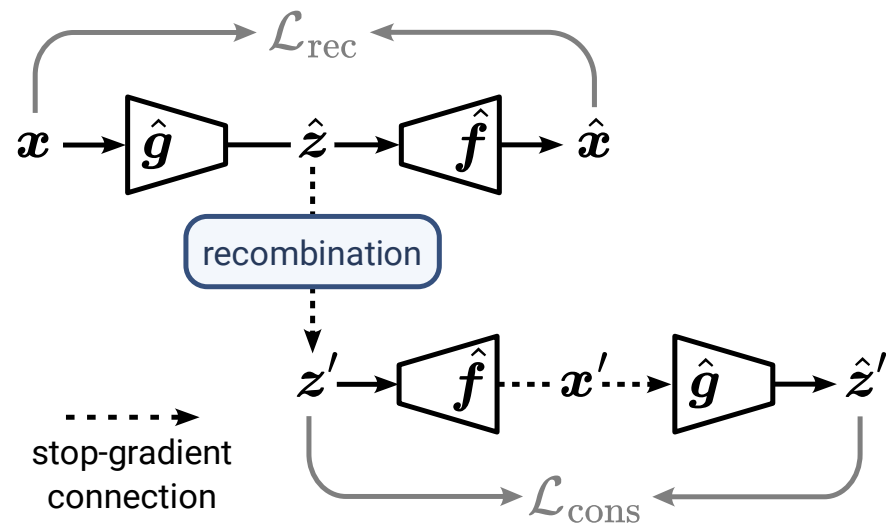


# Provable Compositional Generalization for Object-Centric Learning

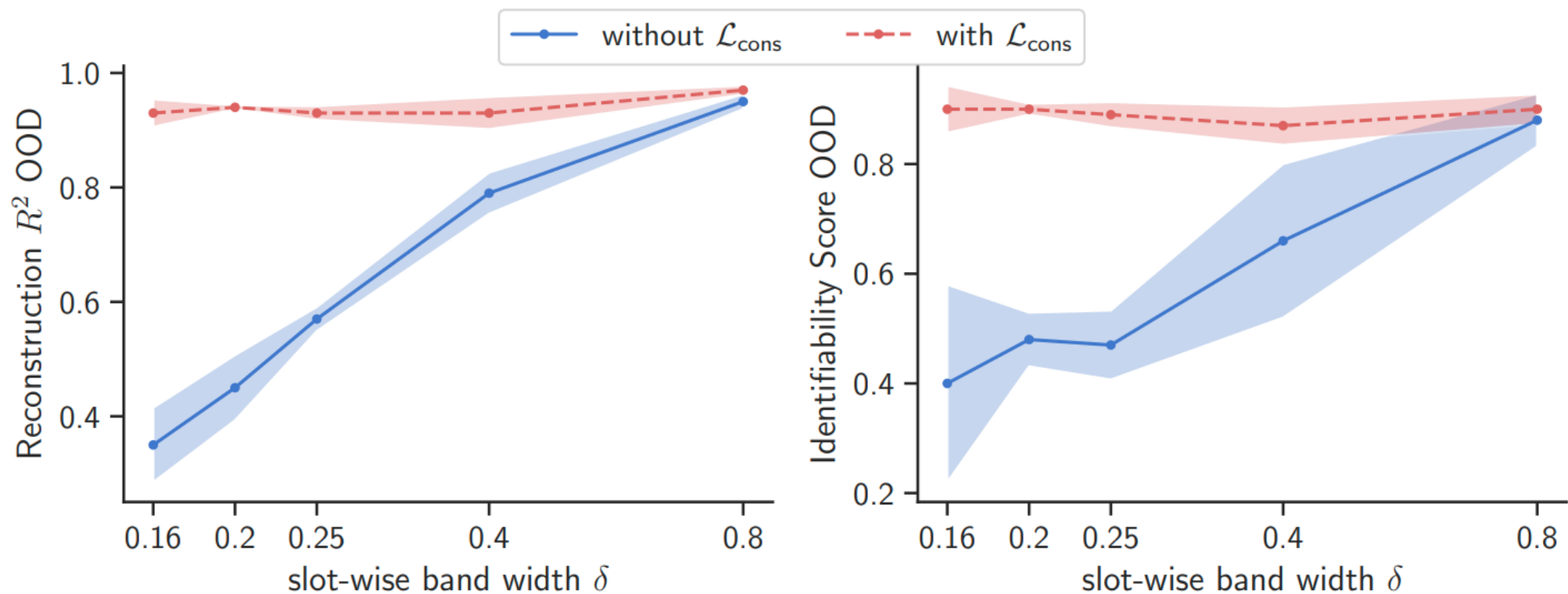




# Provable Compositional Generalization for Object-Centric Learning



# Provable Compositional Generalization for Object-Centric Learning



# **Topics of Interest**

# Topics of Interest



**Different Aspects of Generalization**



**Safety & Security Problems**

# Topics of Interest

## Different Aspects of Generalization

- Diffusion Models
  - Compositional (Concepts) Generalization
  - Shape Bias
  - Generative Classifiers

## Safety & Security Problems

# Topics of Interest

## Different Aspects of Generalization

- Diffusion Models
  - Compositional (Concepts) Generalization
  - Shape Bias
  - Generative Classifiers

## Safety & Security Problems

- Adversarial / Backdoor Attacks
- Data Memorization / Extraction
- Data Poisoning
- Watermarking

# Topics of Interest

## Different Aspects of Generalization

- Diffusion Models
  - Compositional (Concepts) Generalization
  - Shape Bias
  - Generative Classifiers
- Capabilities vs. Alignment

## Safety & Security Problems

- Adversarial / Backdoor Attacks
- Data Memorization / Extraction
- Data Poisoning
- Watermarking
- New Threat Models
  - Author profiling
  - Deceptive Behavior
  - Malicious ChatBots



# Topics of Interest

## Different Aspects of Generalization

- Diffusion Models
  - Compositional (Concepts) Generalization
  - Shape Bias
  - Generative Classifiers
- Capabilities vs. Alignment

## Safety & Security Problems

- Adversarial / Backdoor Attacks
- Data Memorization / Extraction
- Data Poisoning
- Watermarking
- New Threat Models
  - Author profiling
  - Deceptive Behavior
  - Malicious ChatBots

# Topics of Interest

## Different Aspects

- Diffusion Models
- Compositionality
- Shape Bias
- Generative Models
- Capabilities vs. Alignment

## Security Problems

- Backdoor Attacks
- Jailbreaks / Prompt Injection / Extraction
- Data Poisoning
- Model Inversion
- Model Extraction
- Model Stealing
- Model Profiling
- Model Behavior
- ChatBots

Some alignment researchers [...] believe that sufficiently advanced language models should be aligned to prevent an existential risk [...] to humanity: if this were true, **an attack** that causes such a model to become misaligned **would be devastating**. [1]

[1] Carlini, N., Nasr, M., Choquette-Choo, C. A., Jagielski, M., Gao, I., Awadalla, A., ... & Schmidt, L. (2023). Are aligned neural networks adversarially aligned?. arXiv preprint arXiv:2306.15447.

**Thank you!**