



# CISPA

HELMHOLTZ CENTER FOR  
INFORMATION SECURITY

# ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models

Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang,  
Mario Fritz, Michael Backes

CISPA Helmholtz Center for Information Security, Swiss Data Science Center





*The more a model learns about  
you, the better it gets*

## Cars



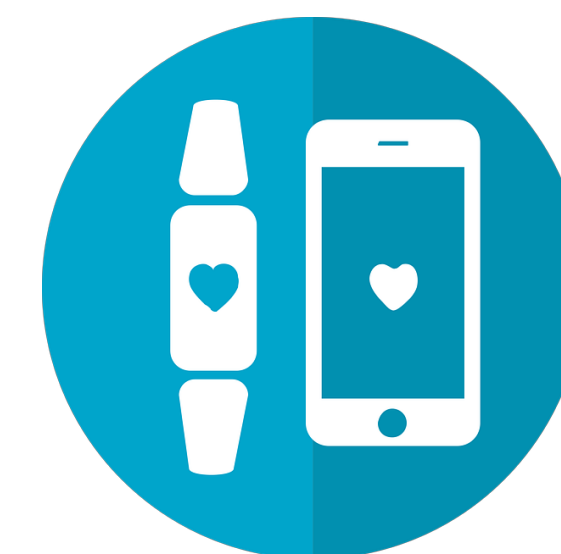
*The more a model learns about  
you, the better it gets*



## Cars



## Activity Tracker



*The more a model learns about  
you, the better it gets*

## Cars



## Activity Tracker



*The more a model learns about you, the better it gets*

## Social Media





### Cars



### Activity Tracker



*The more a model learns about you, the better it gets*

### Social Media



### Personal Assistant



# How Sensitive The Data Can Be?

---



# How Sensitive The Data Can Be?

---

- Financial data

# How Sensitive The Data Can Be?

---

- Financial data
- Location and activity data



# How Sensitive The Data Can Be?

---

- Financial data
- Location and activity data
- Biomedical data

# How Sensitive The Data Can Be?

---

- Financial data
- Location and activity data
- Biomedical data
- etc.

# Privacy in Machine Learning

---

# Privacy in Machine Learning

---

- ML models are trained on sensitive data



# Privacy in Machine Learning

---

- ML models are trained on sensitive data
- Main focus is performance

# Privacy in Machine Learning

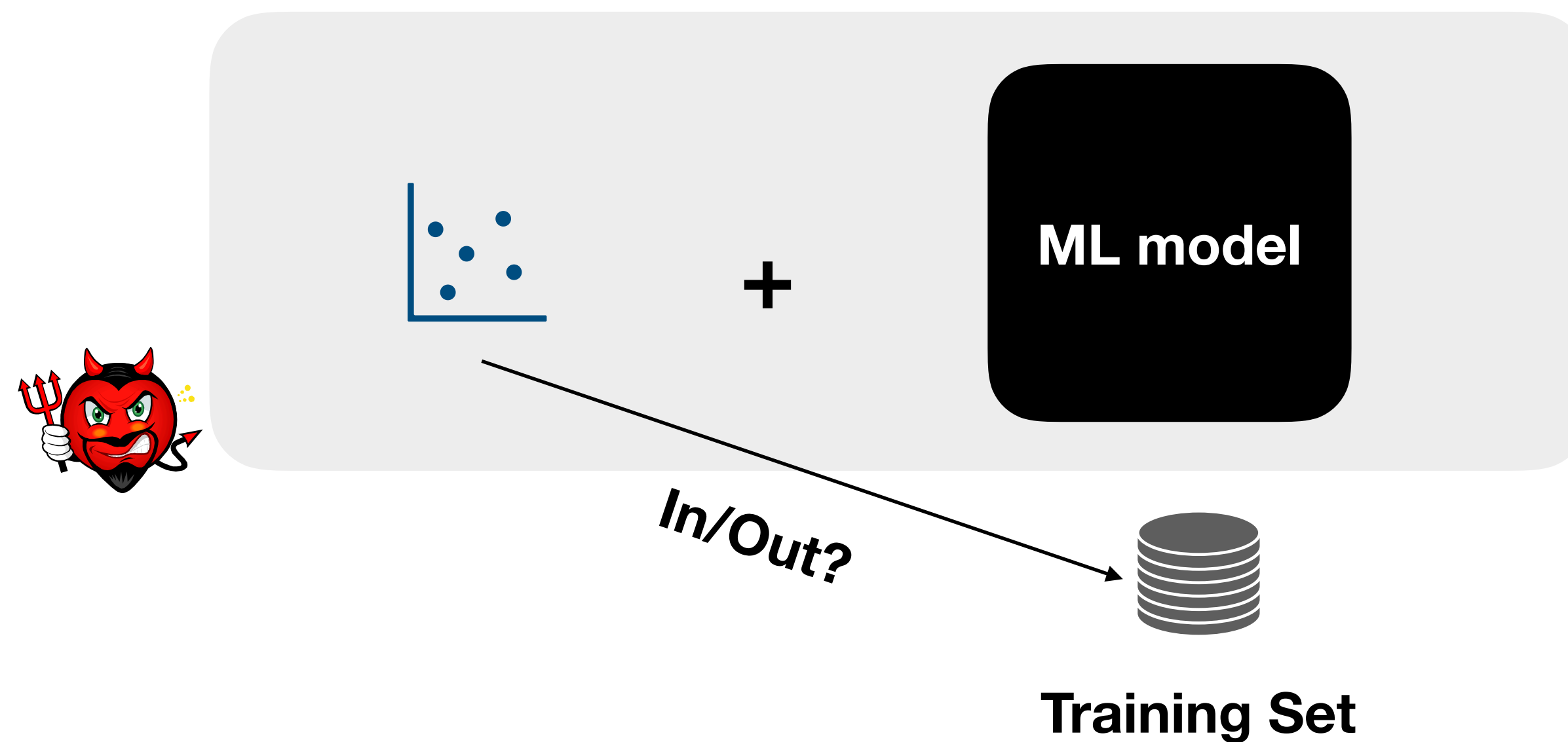
---

- ML models are trained on sensitive data
- Main focus is performance
  - ▶ Largely overlooked

# Membership Inference

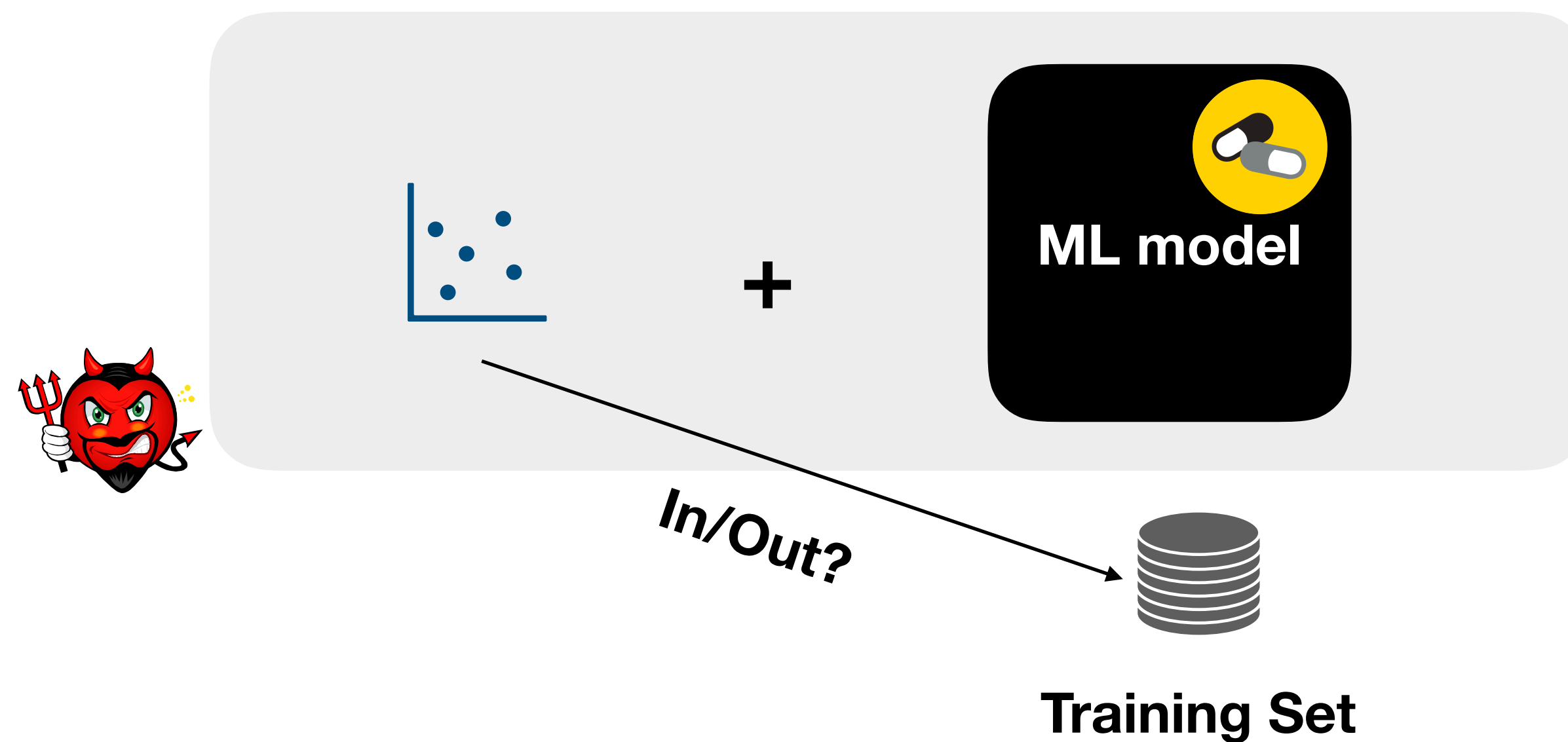


# Membership Inference

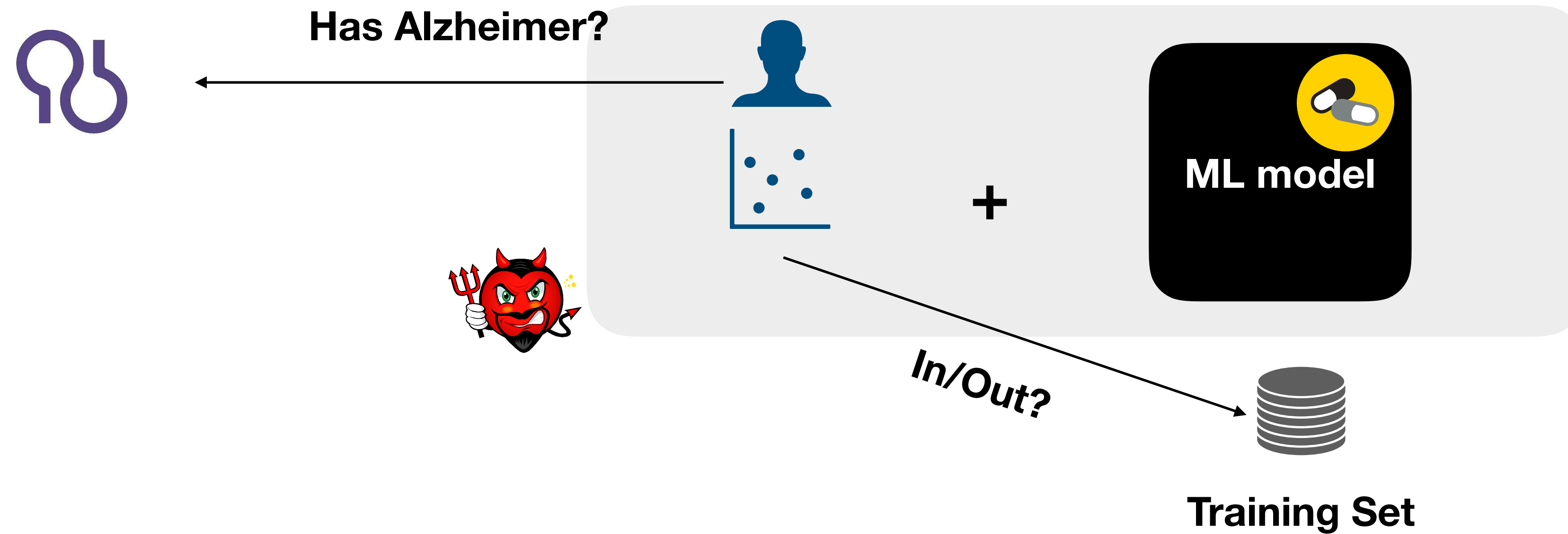




# Membership Inference

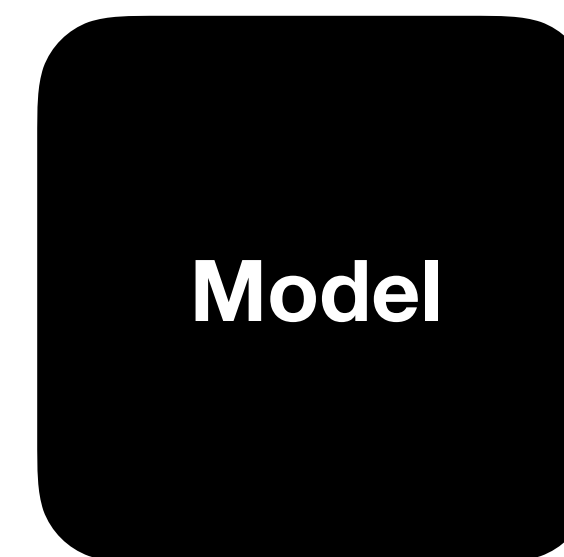
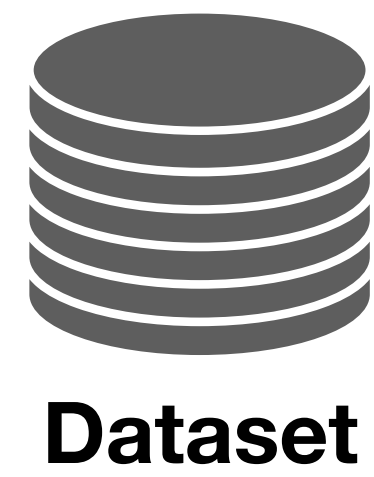


# Membership Inference

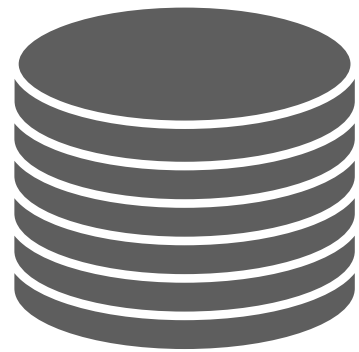


# Machine Learning Pipeline

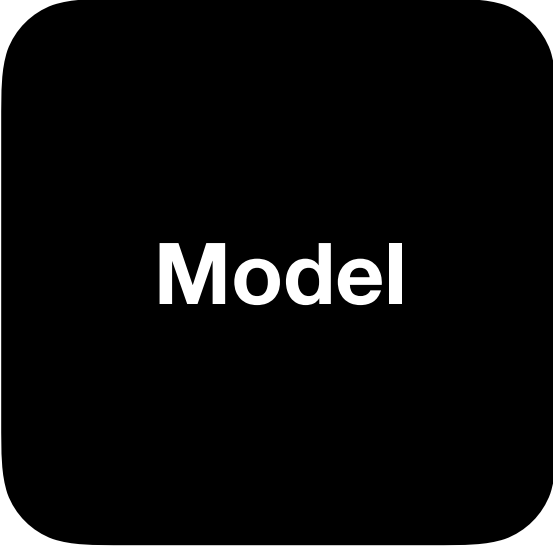
---



# Machine Learning Pipeline

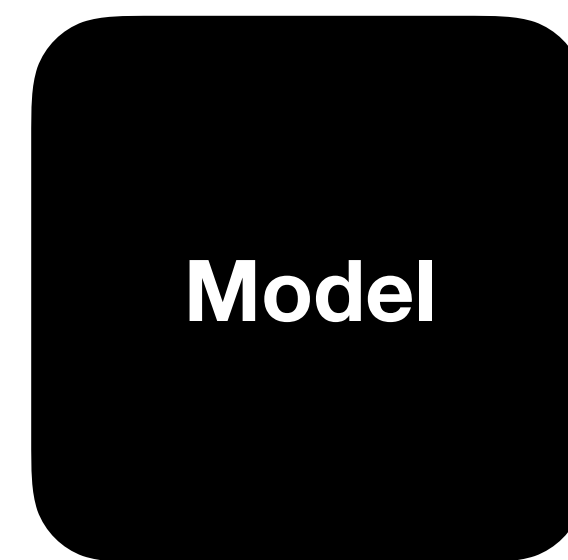
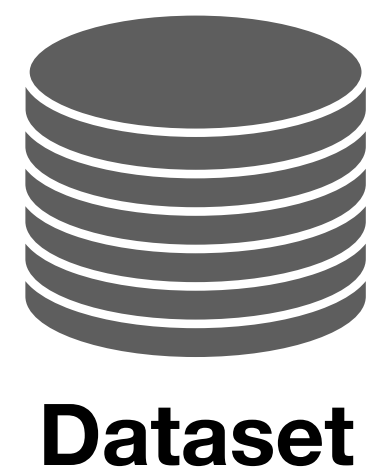
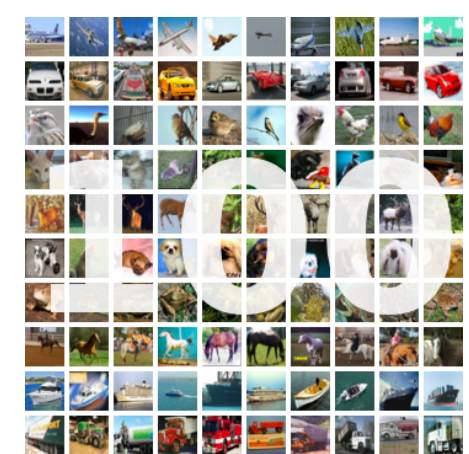


**Dataset**

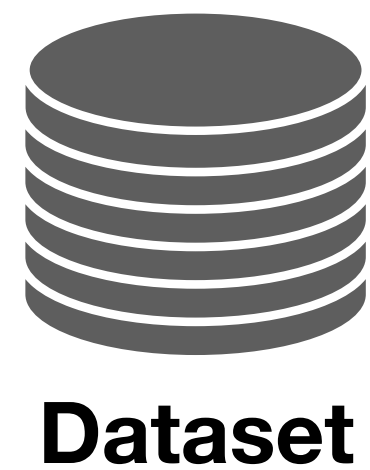




# Machine Learning Pipeline



# Machine Learning Pipeline



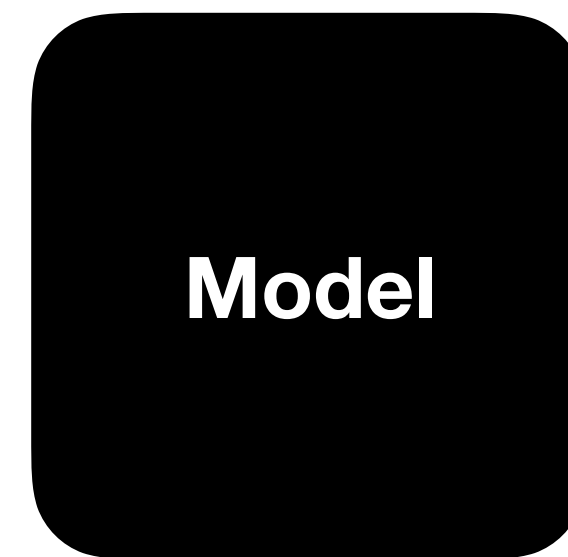
**Dataset**



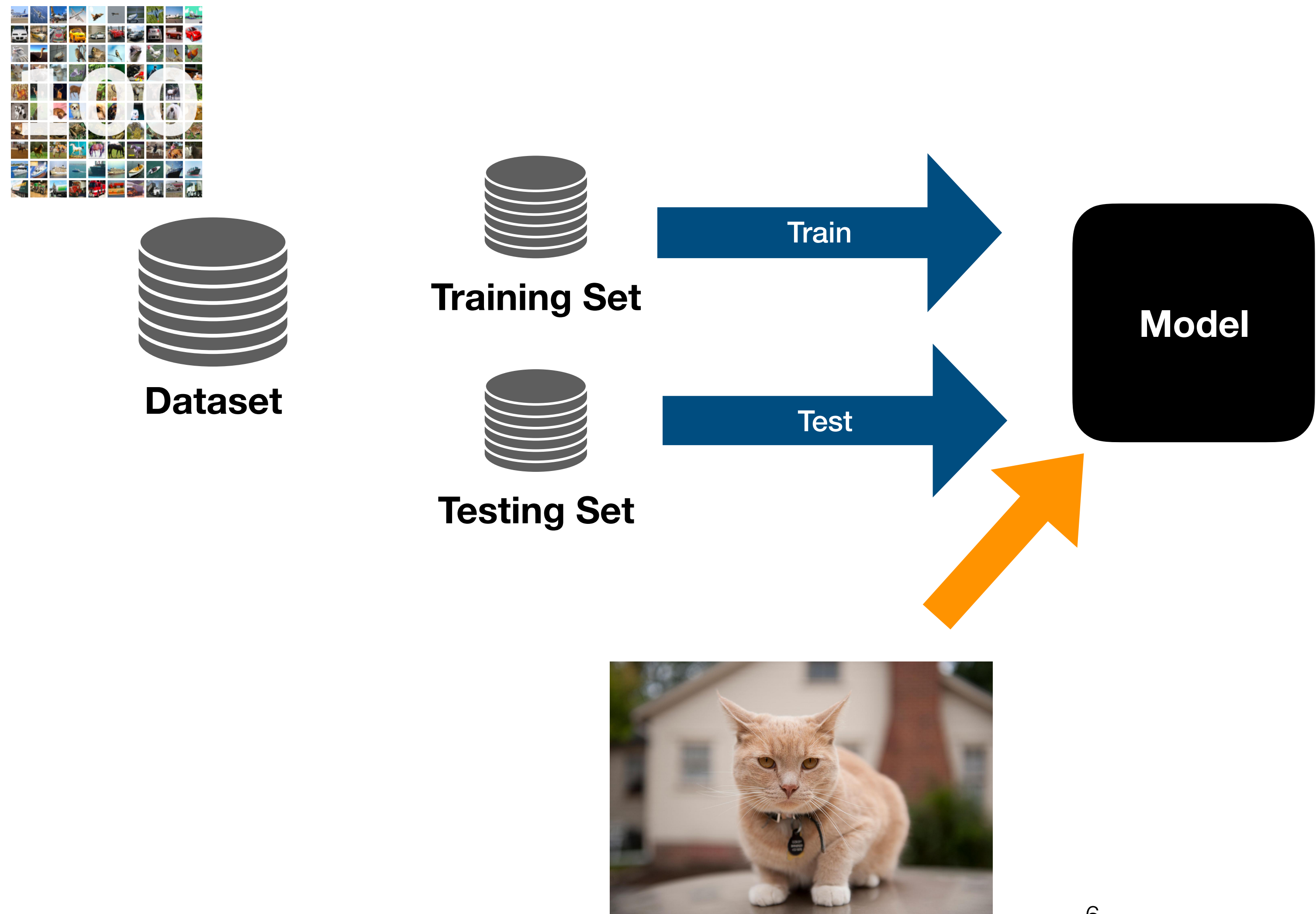
**Training Set**



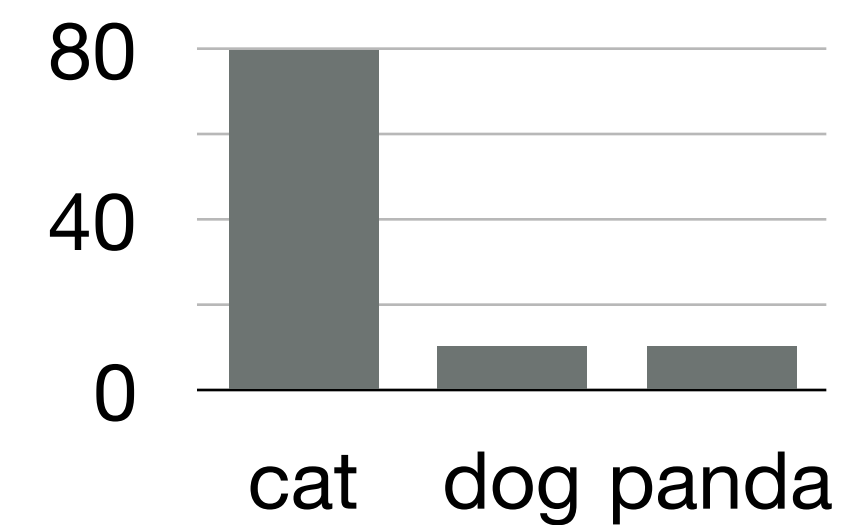
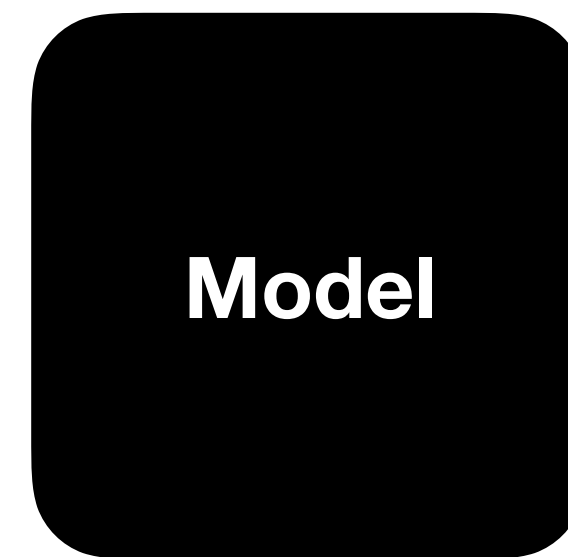
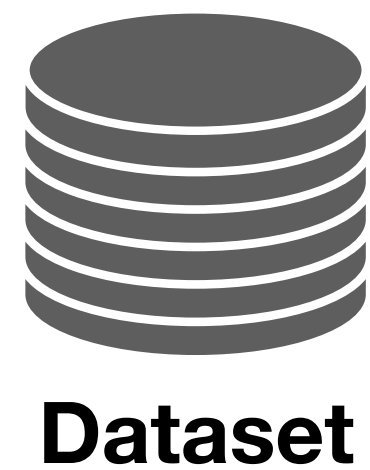
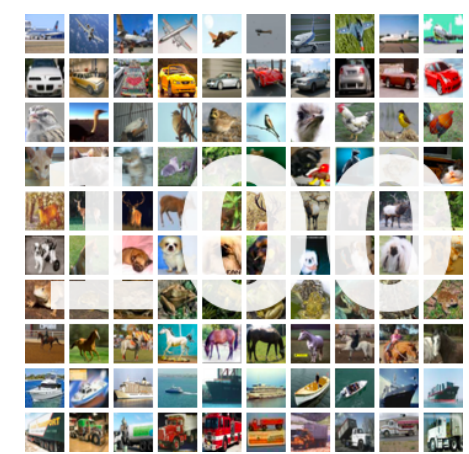
**Testing Set**



# Machine Learning Pipeline



# Machine Learning Pipeline

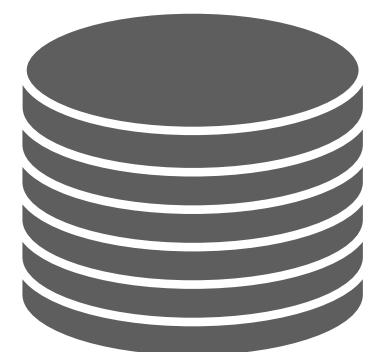




# Threat Model



# State Of The Art (Shokri et al.)

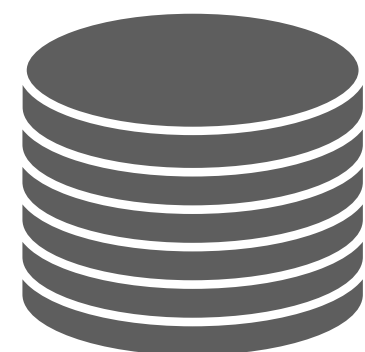


**Target Dataset**

**Target  
Model**



# State Of The Art (Shokri et al.)



**Target Dataset**

**Target  
Model**



**Attack  
Models**



**Attack  
Models**

# State Of The Art (Shokri et al.)



Ground Truth?



# State Of The Art (Shokri et al.)



⋮



⋮





# State Of The Art (Shokri et al.)



⋮

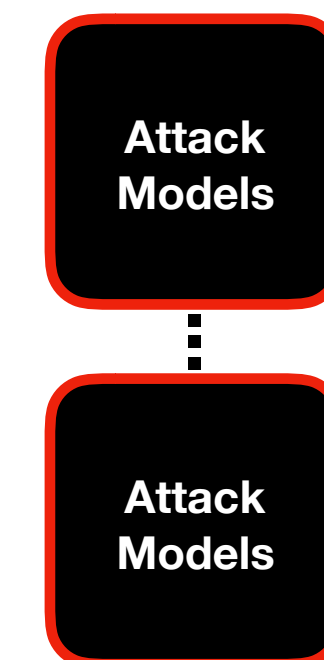
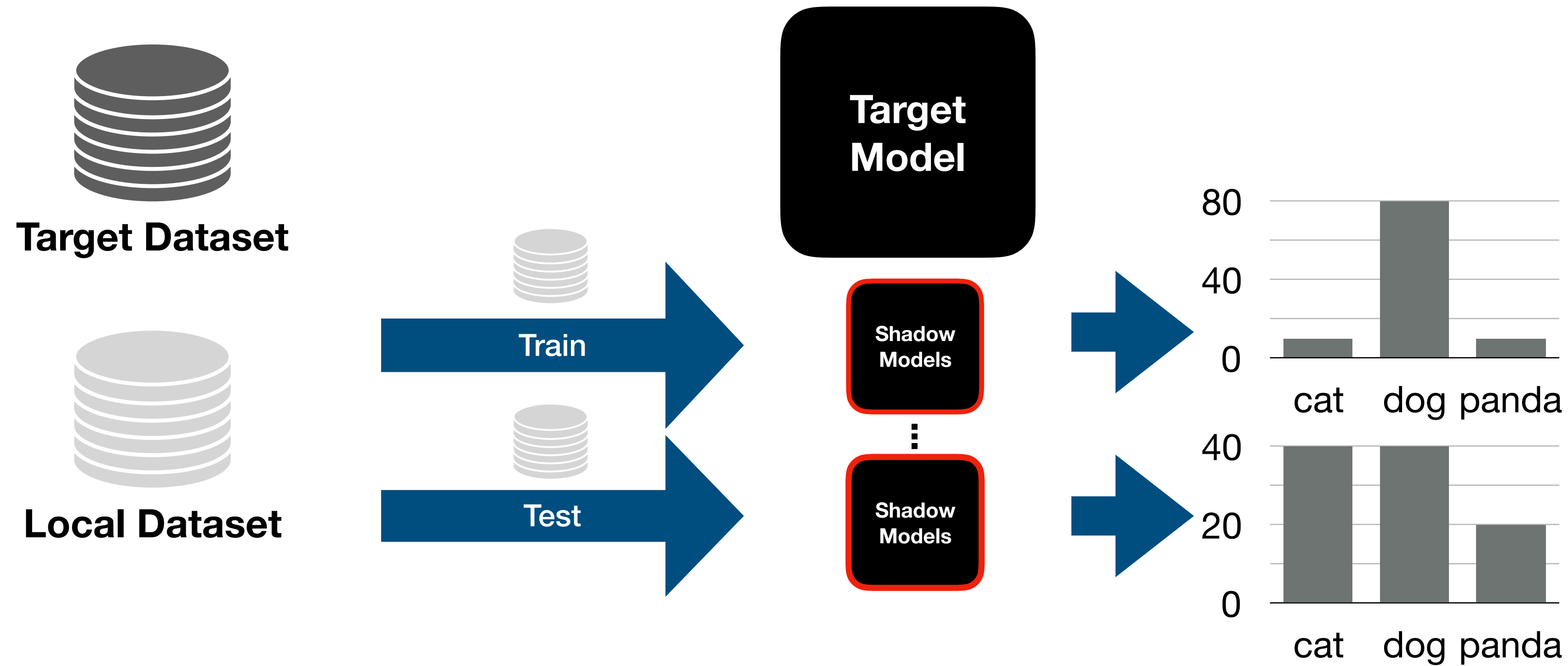


⋮

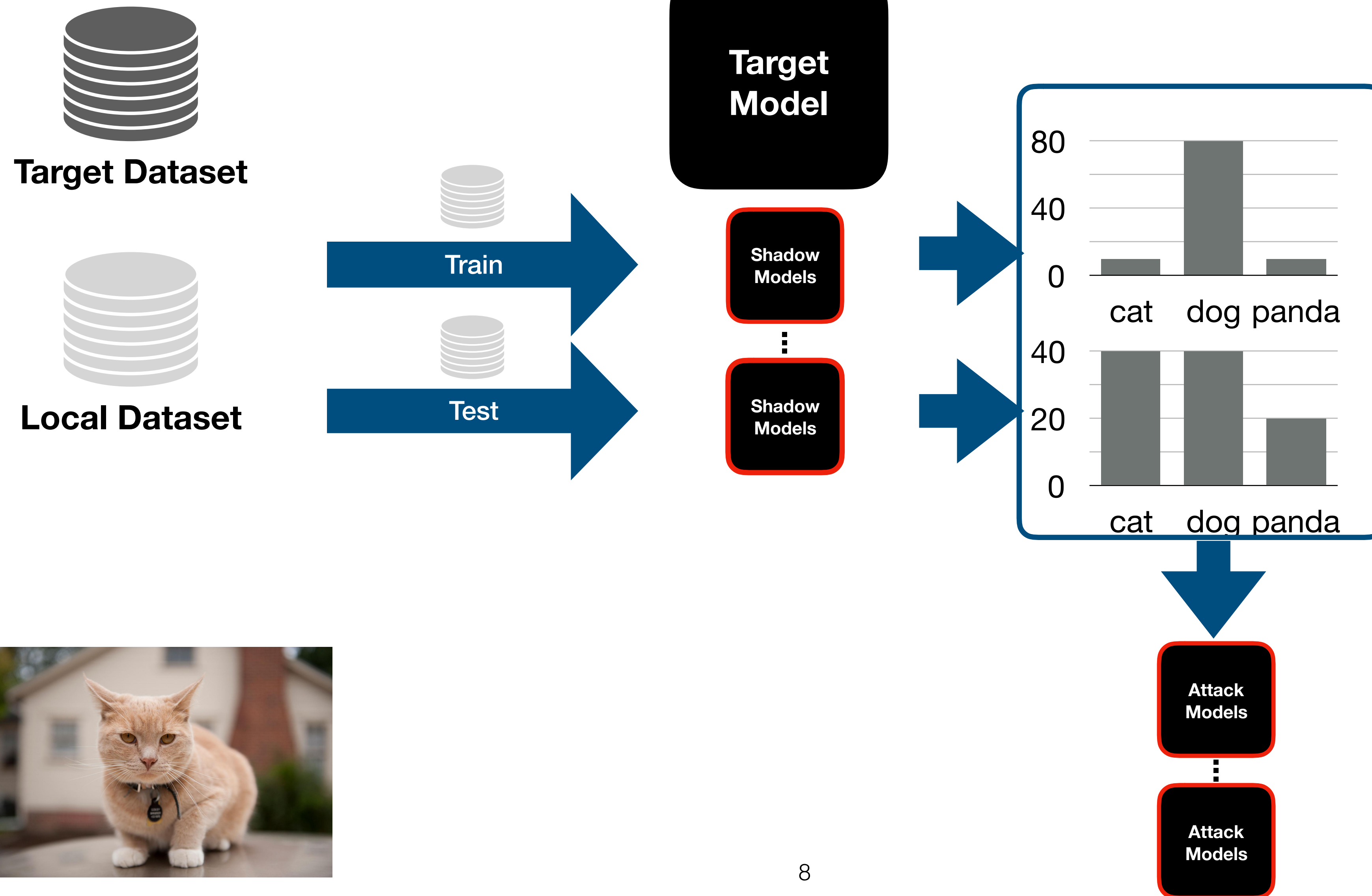




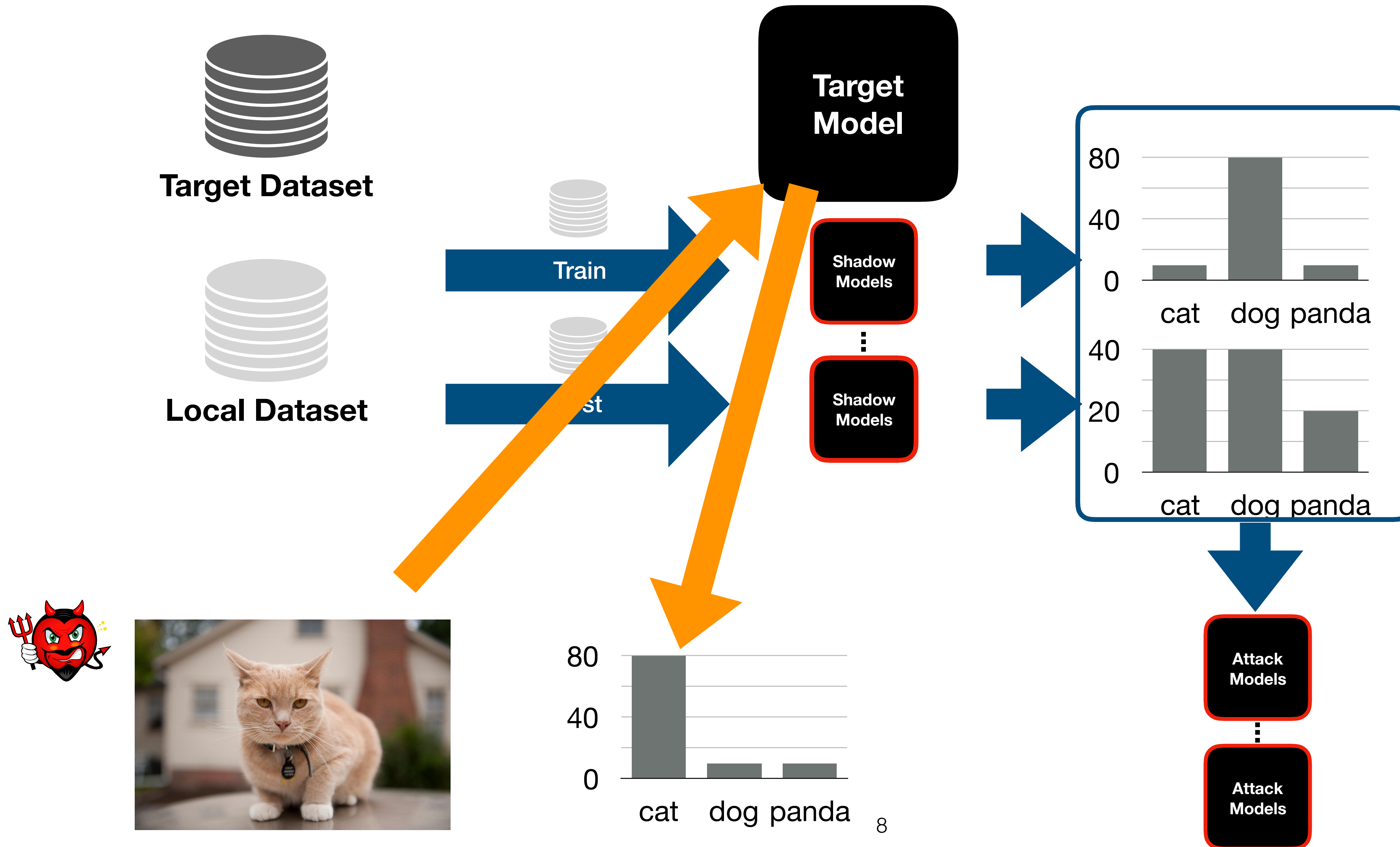
# State Of The Art (Shokri et al.)



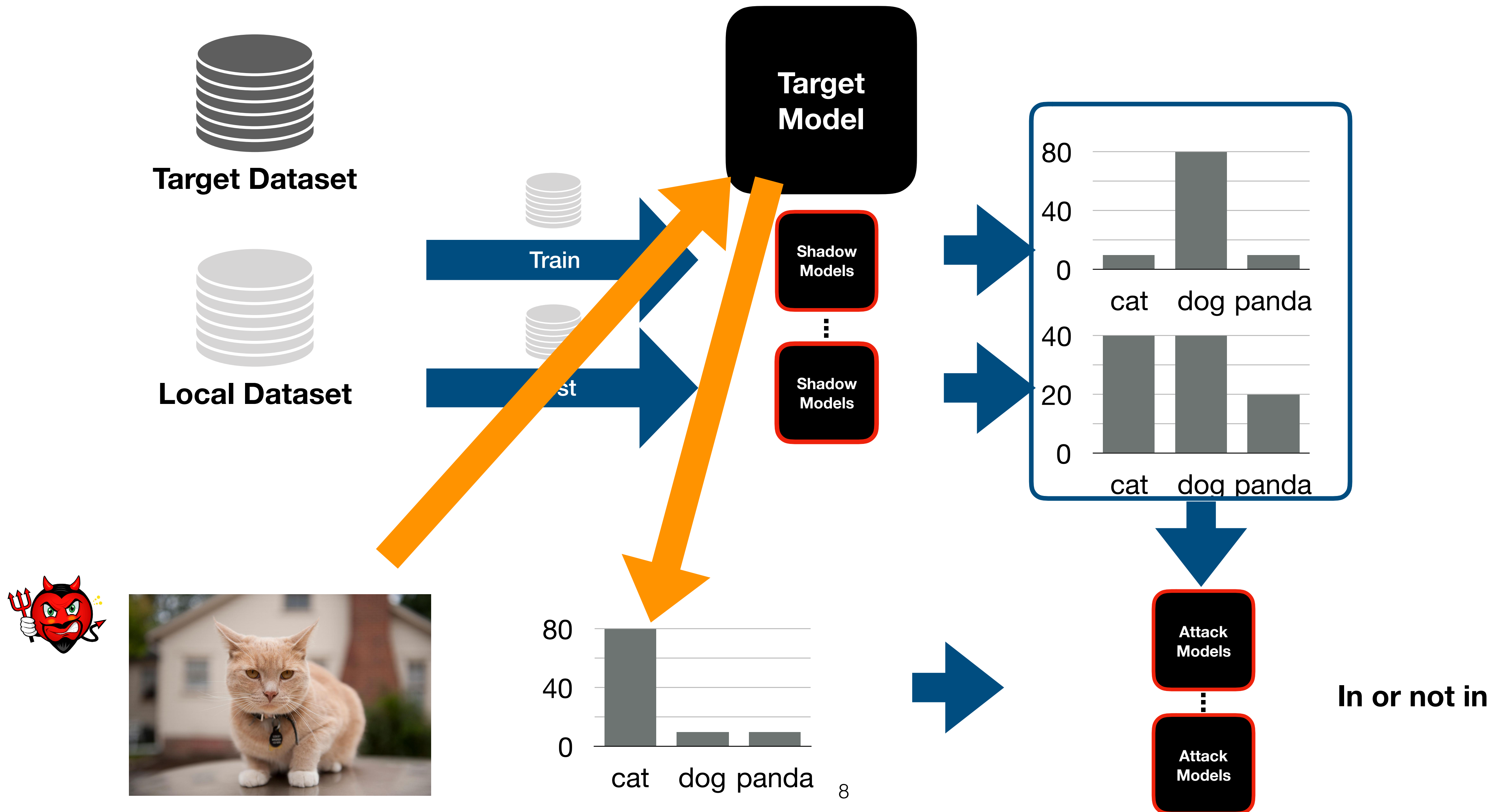
# State Of The Art (Shokri et al.)



# State Of The Art (Shokri et al.)



# State Of The Art (Shokri et al.)





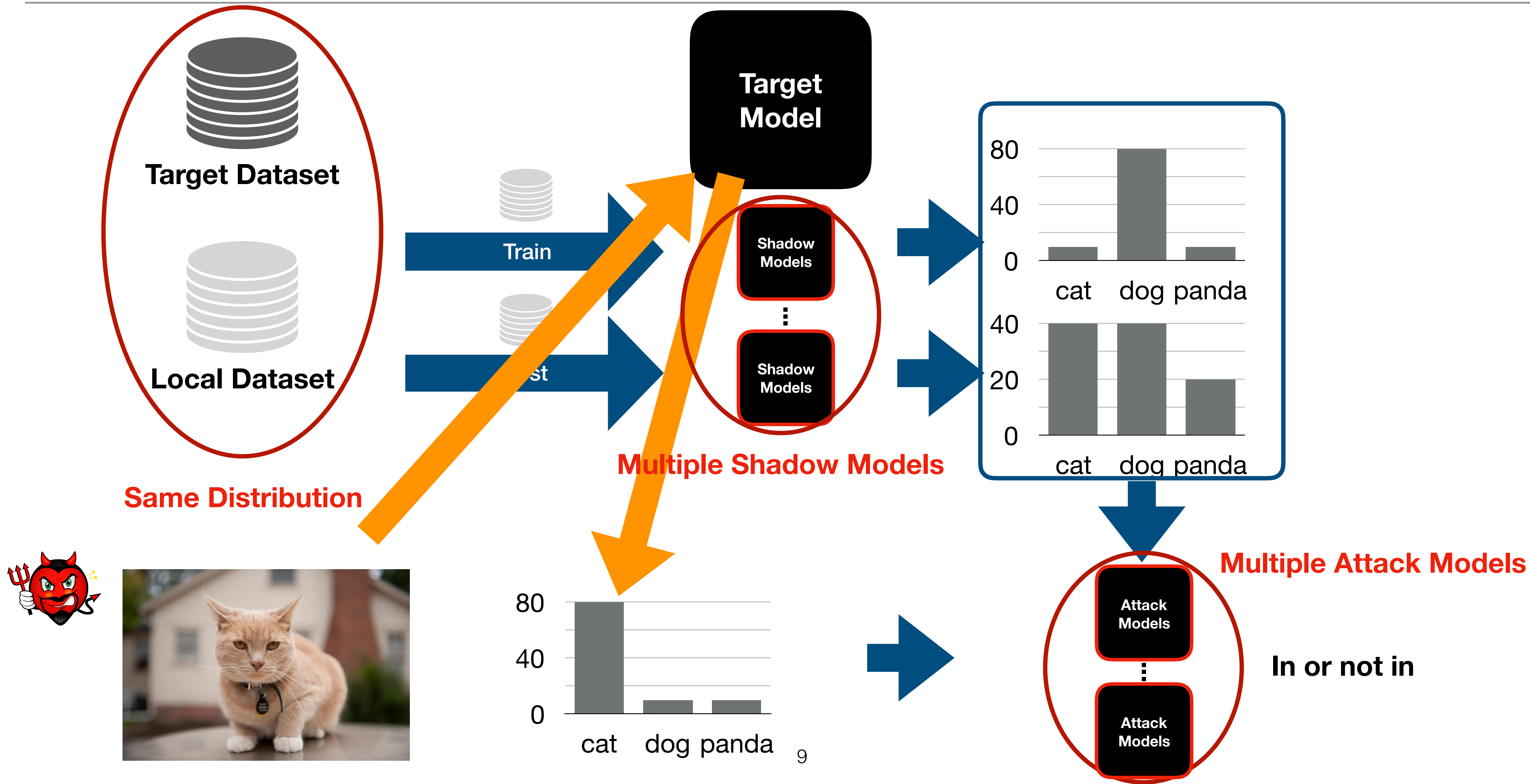




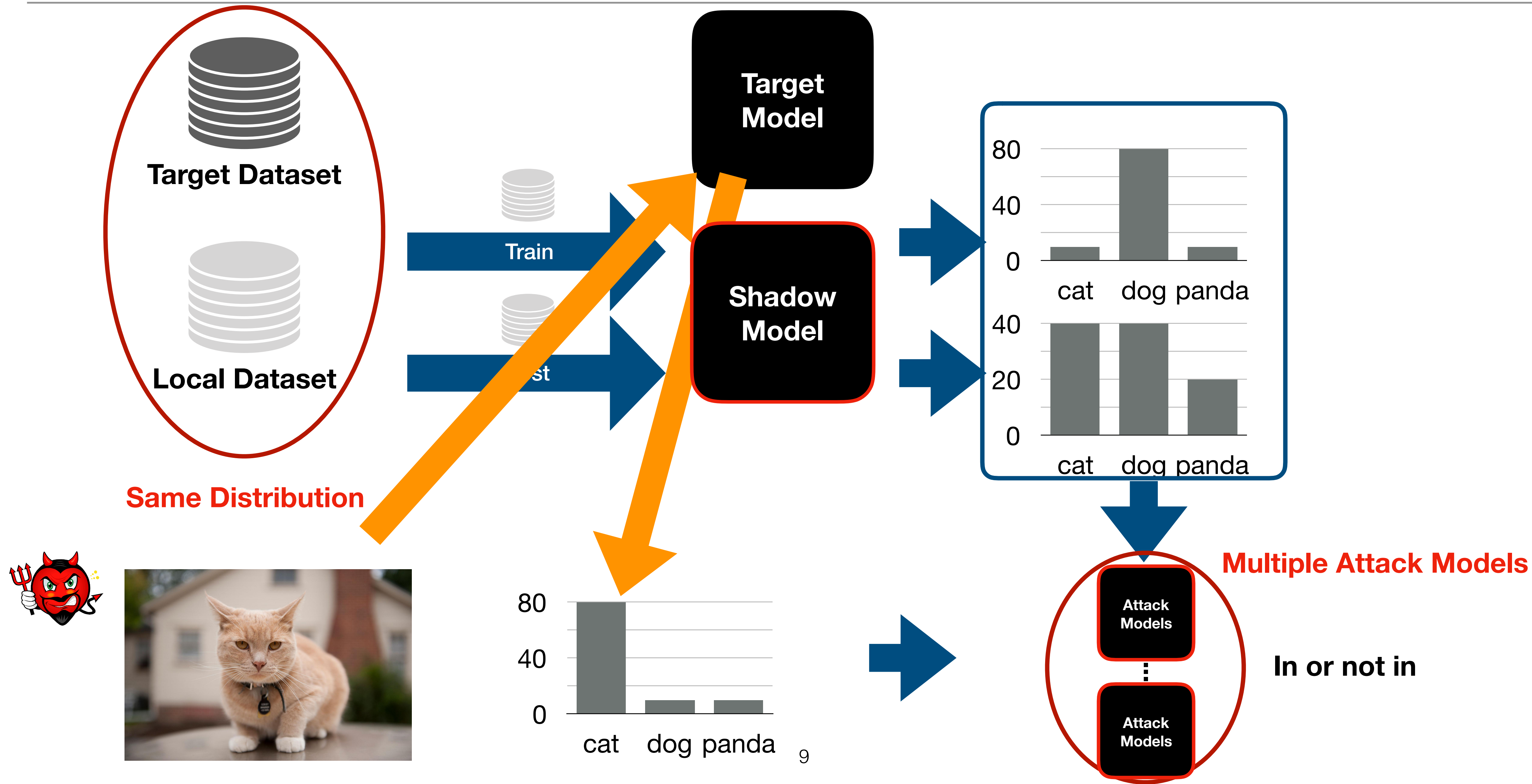




# Our First Attack (Adversary 1)

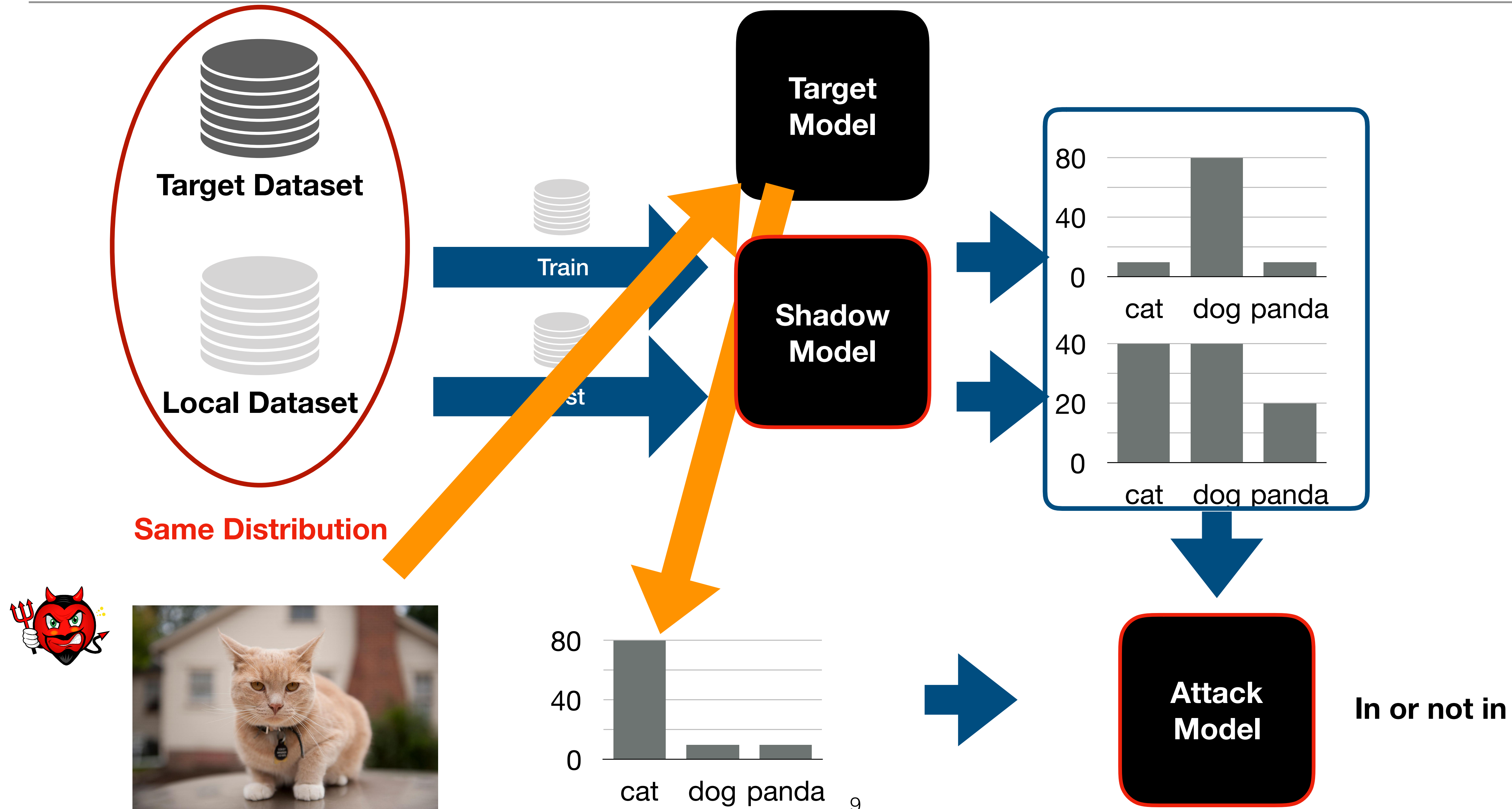


# Our First Attack (Adversary 1)

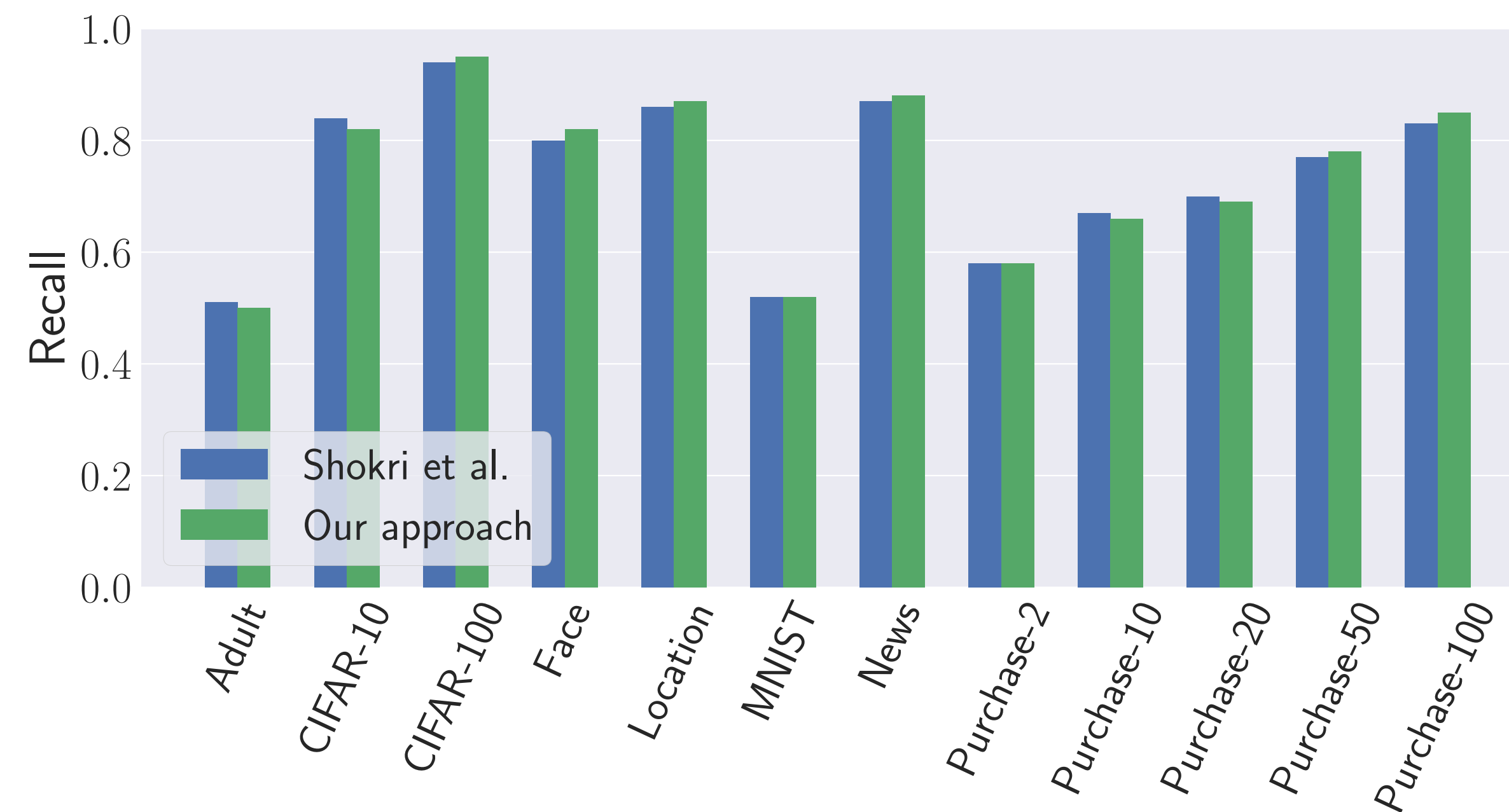
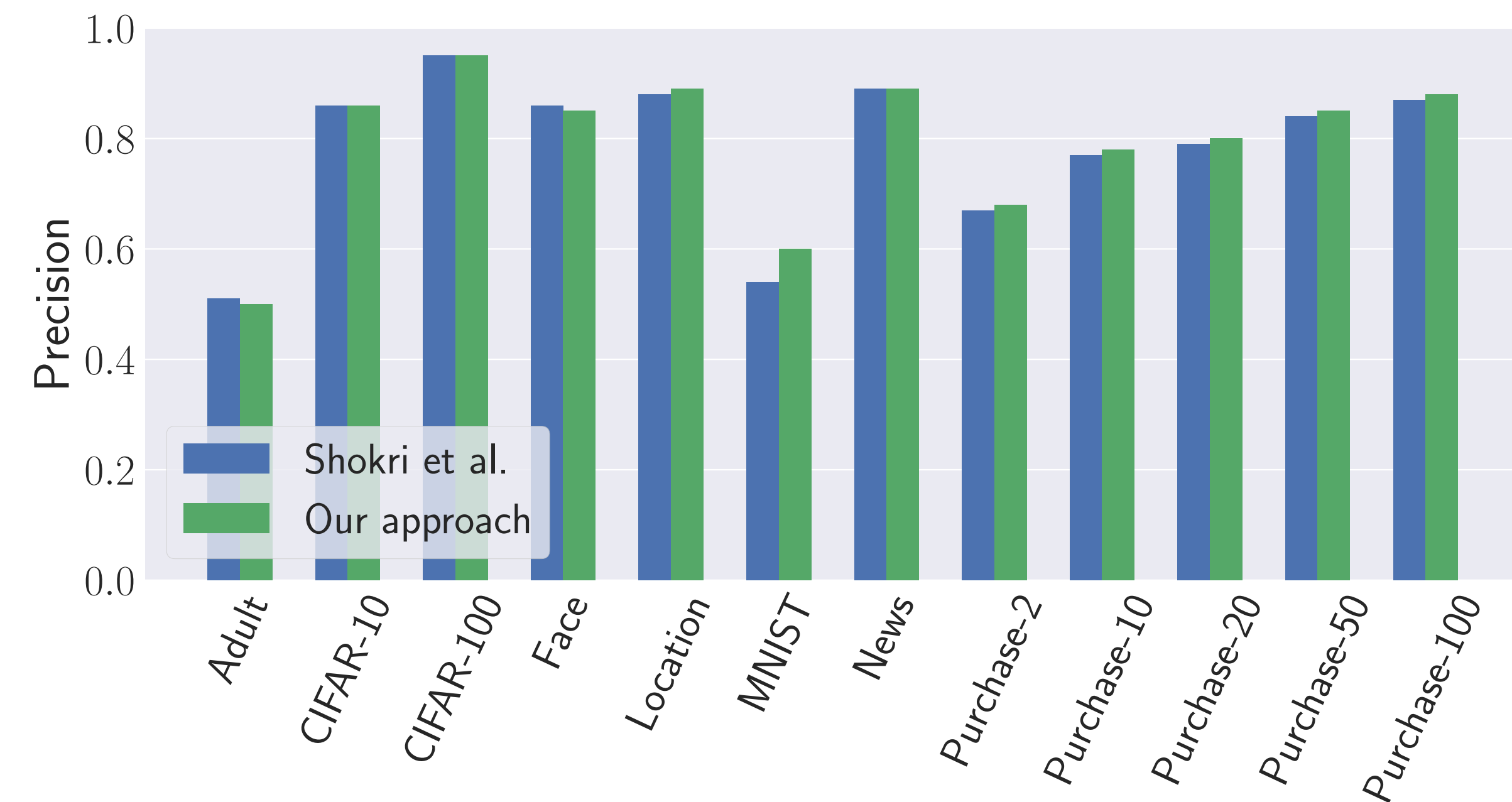




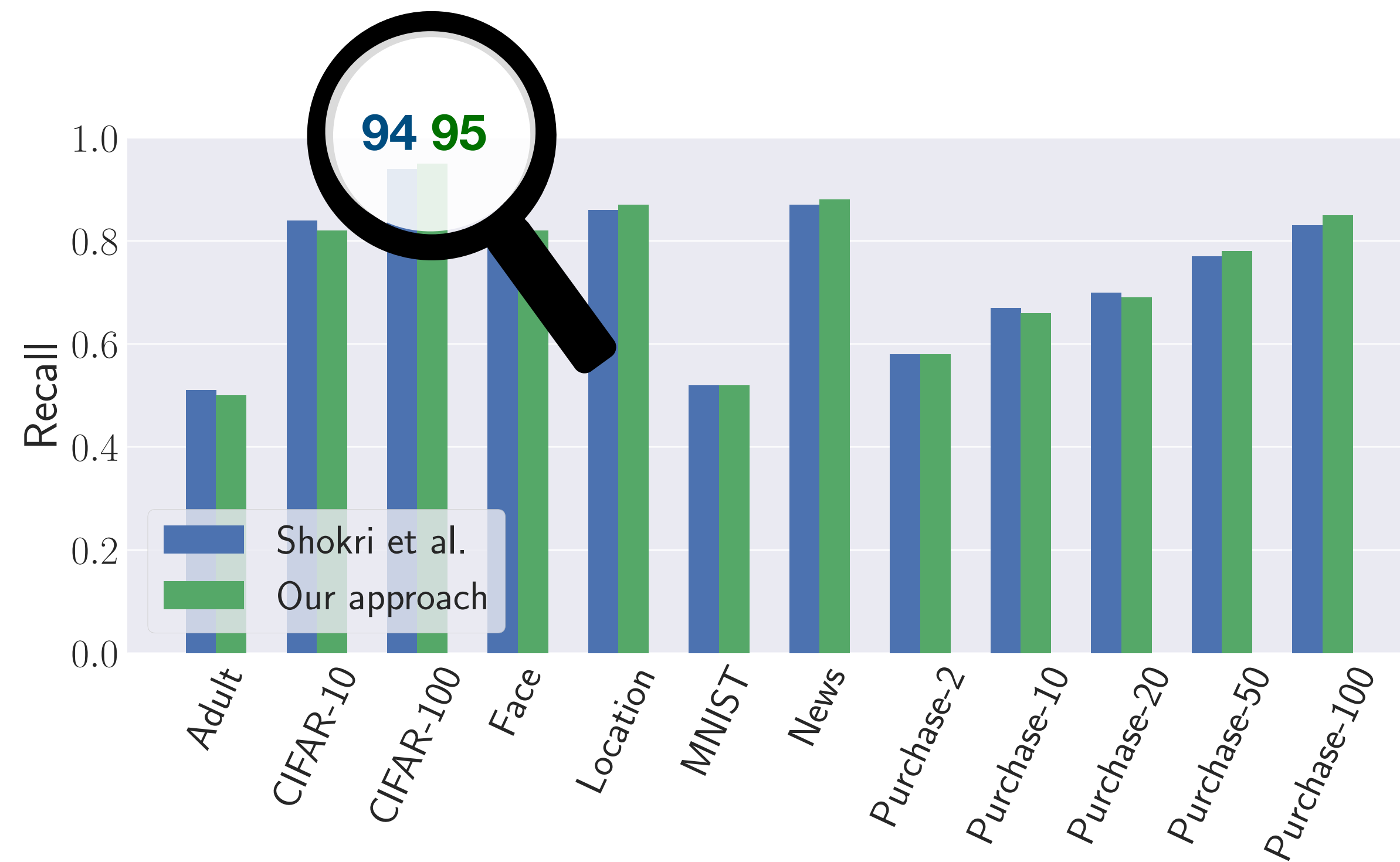
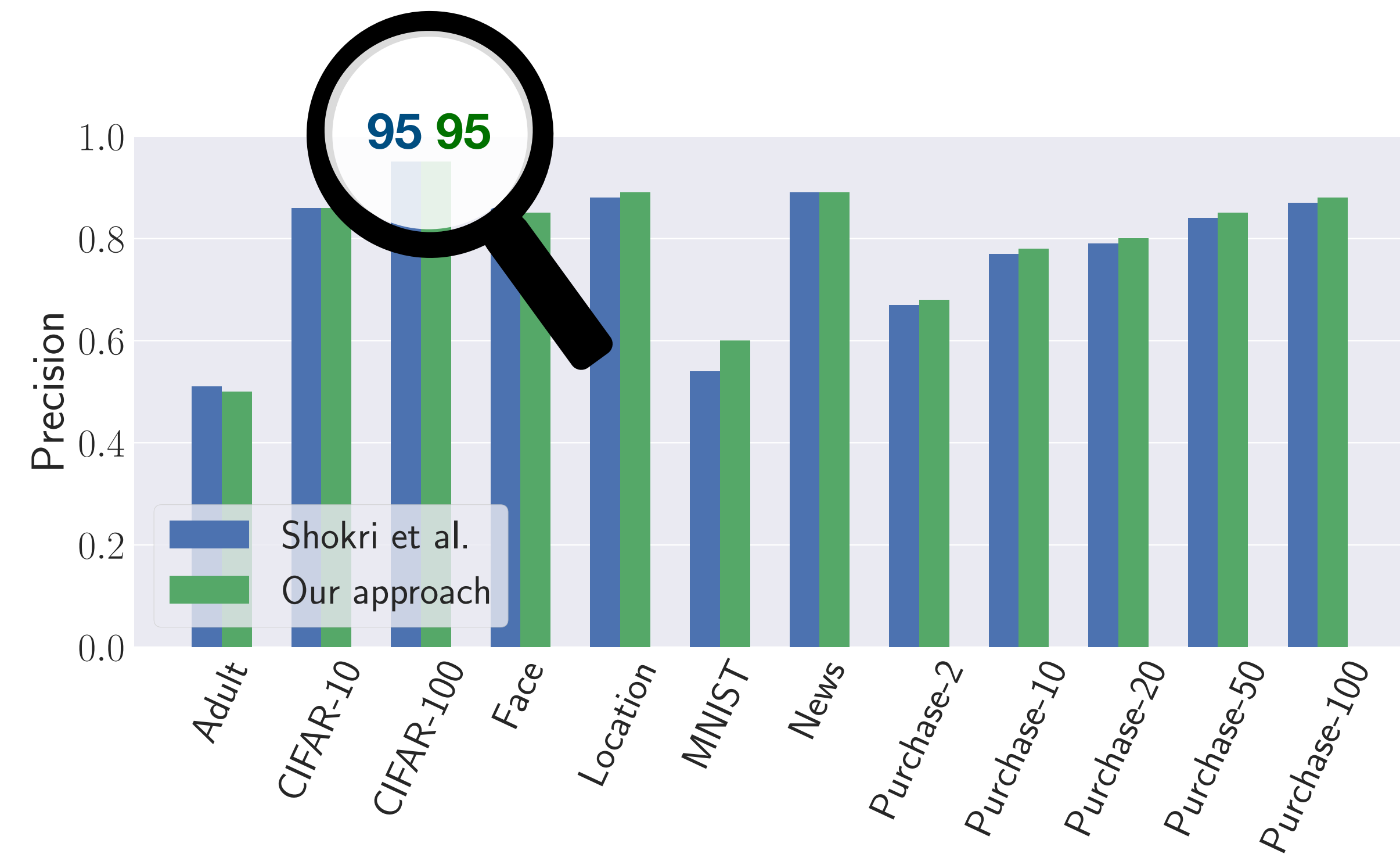
# Our First Attack (Adversary 1)



# Performance Comparison

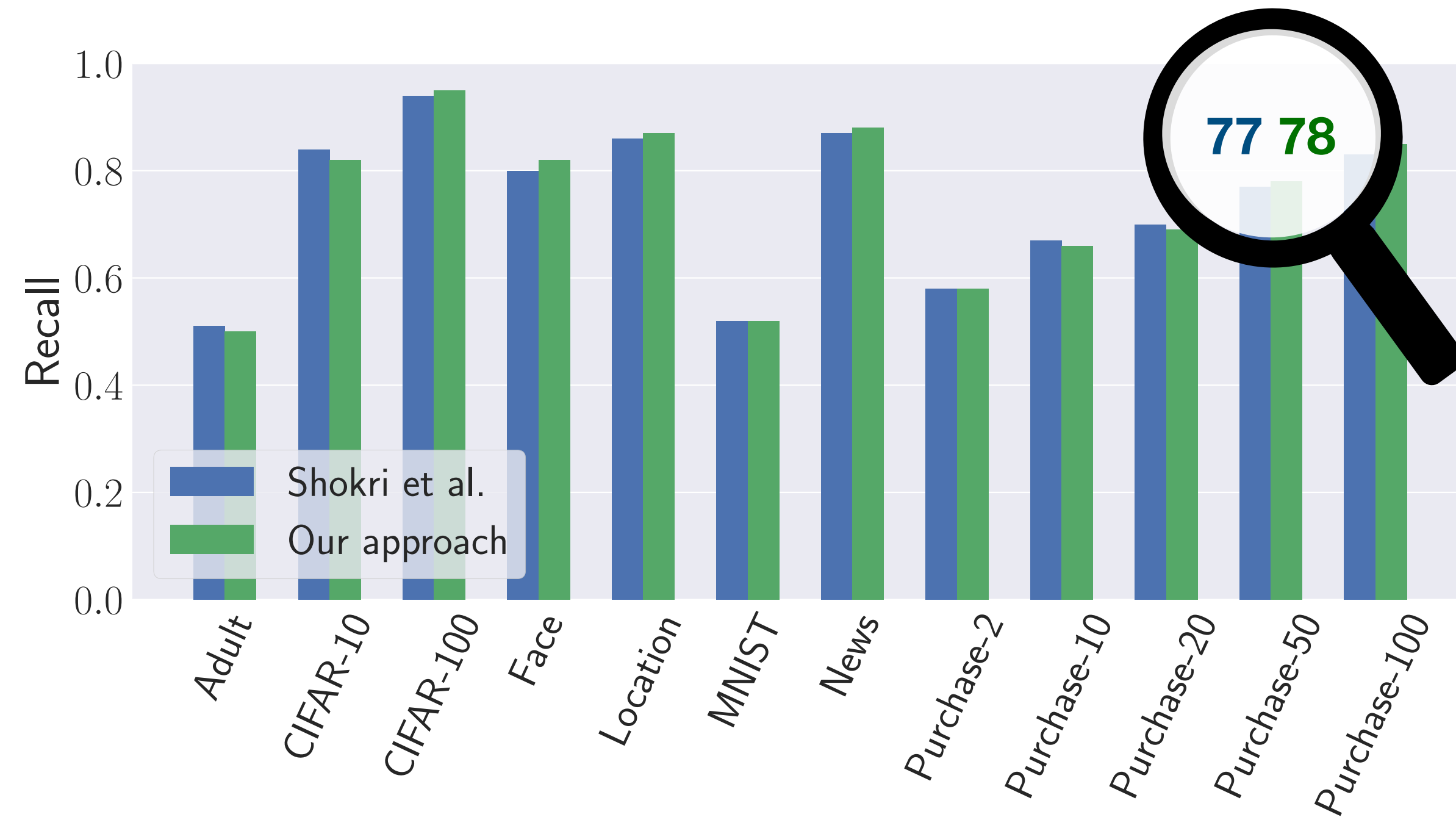
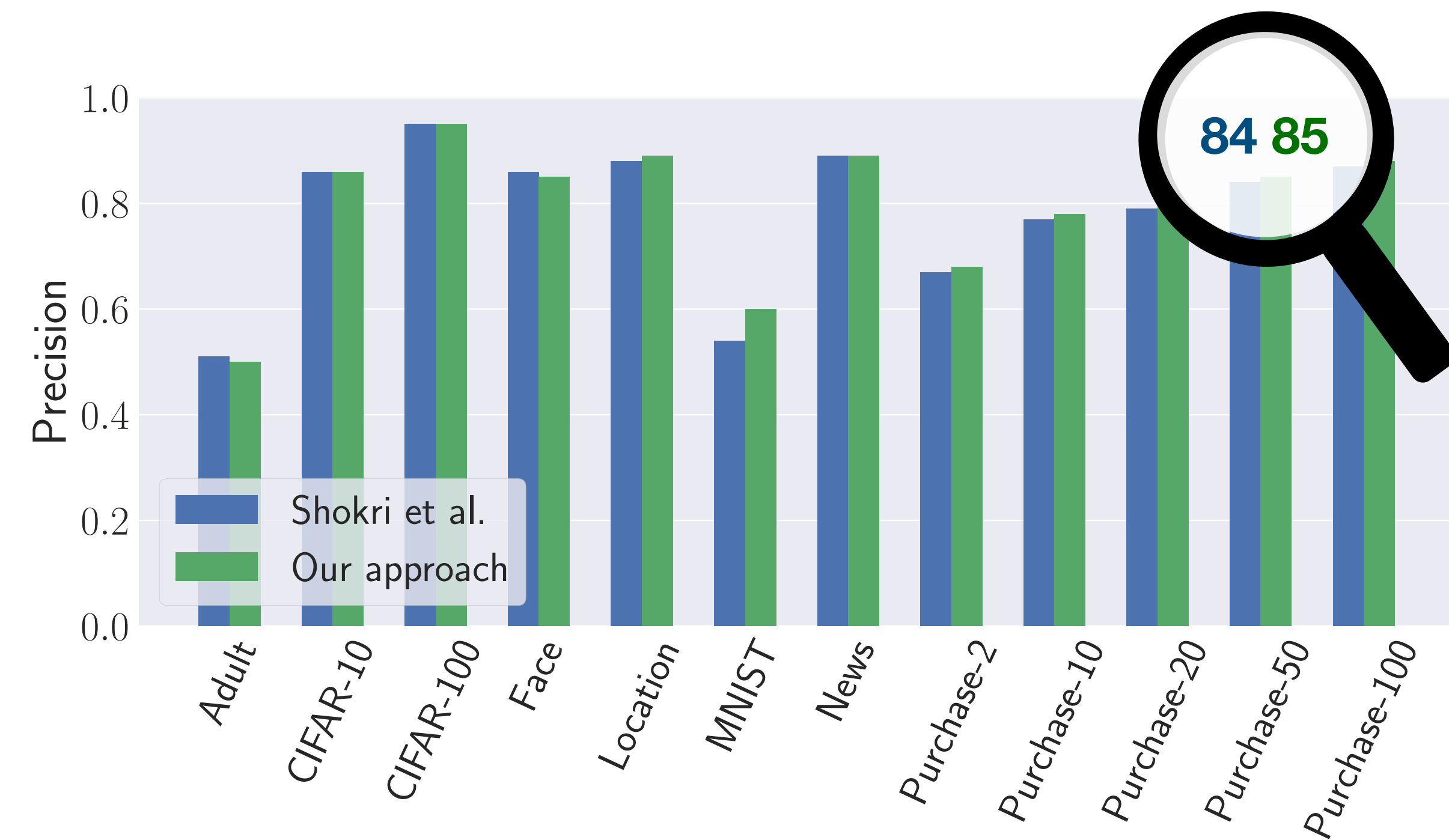


# Performance Comparison

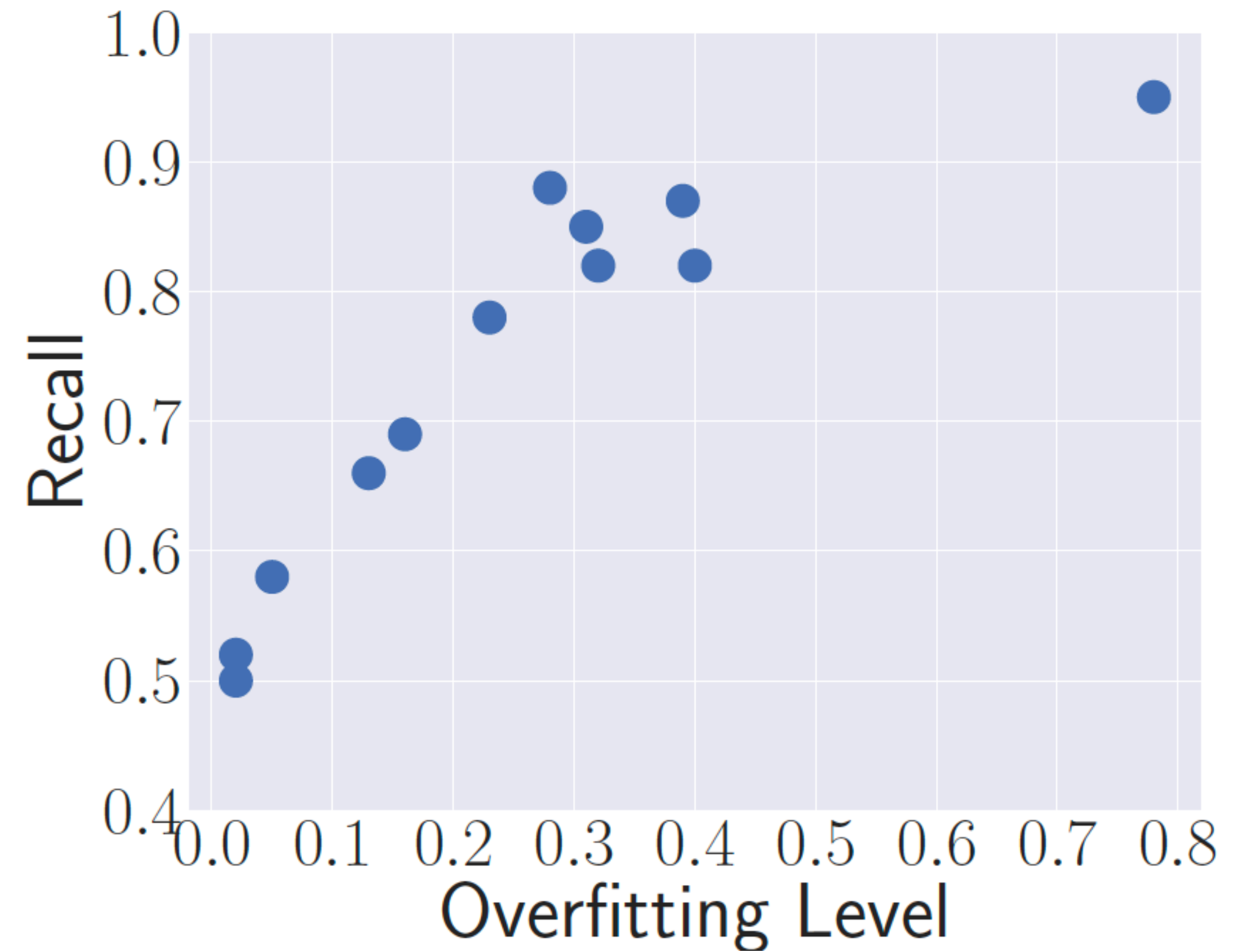
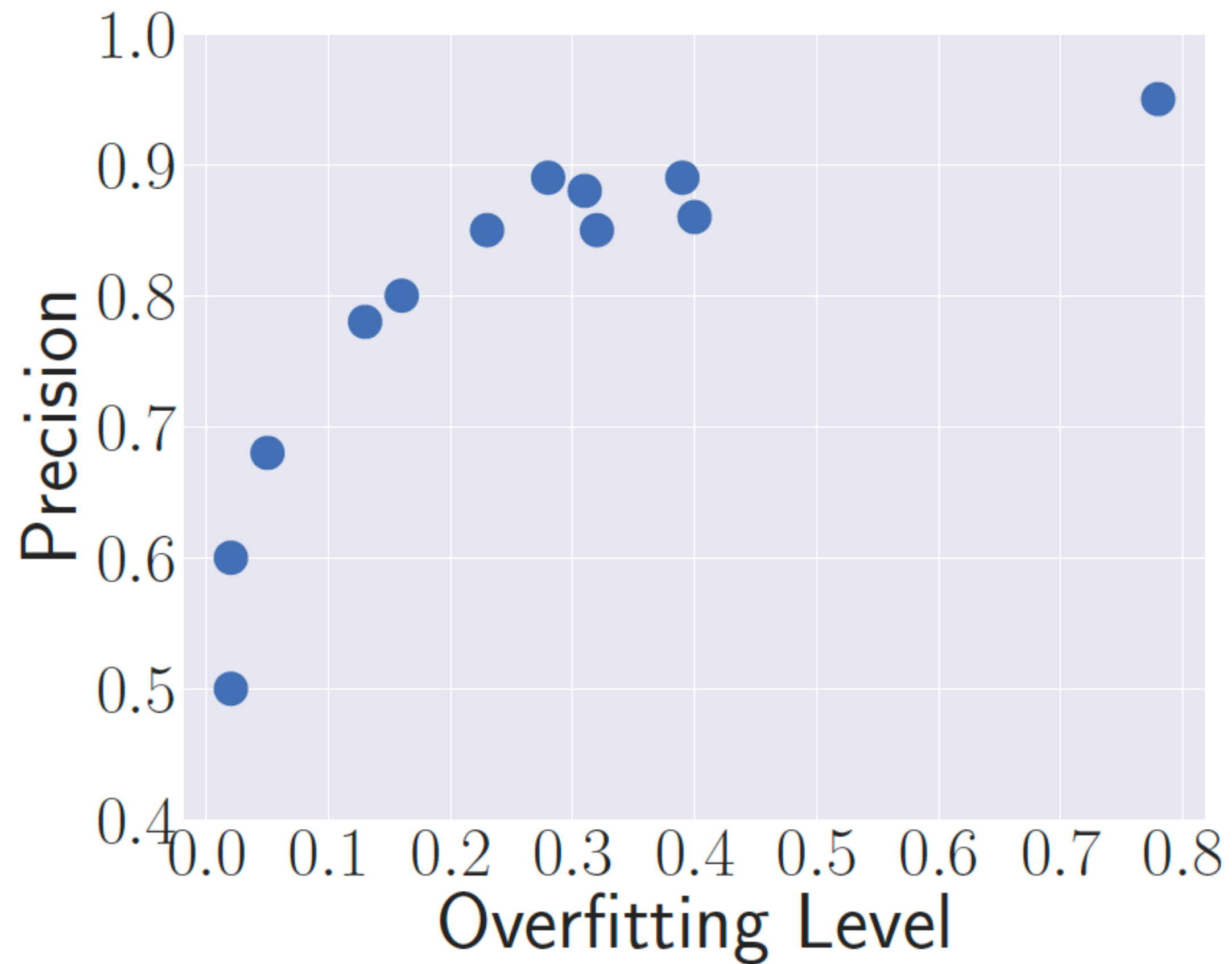




# Performance Comparison



# Overfitting vs Attack's Performance



# Towards a Data Independent Attack

---

# Towards a Data Independent Attack

---

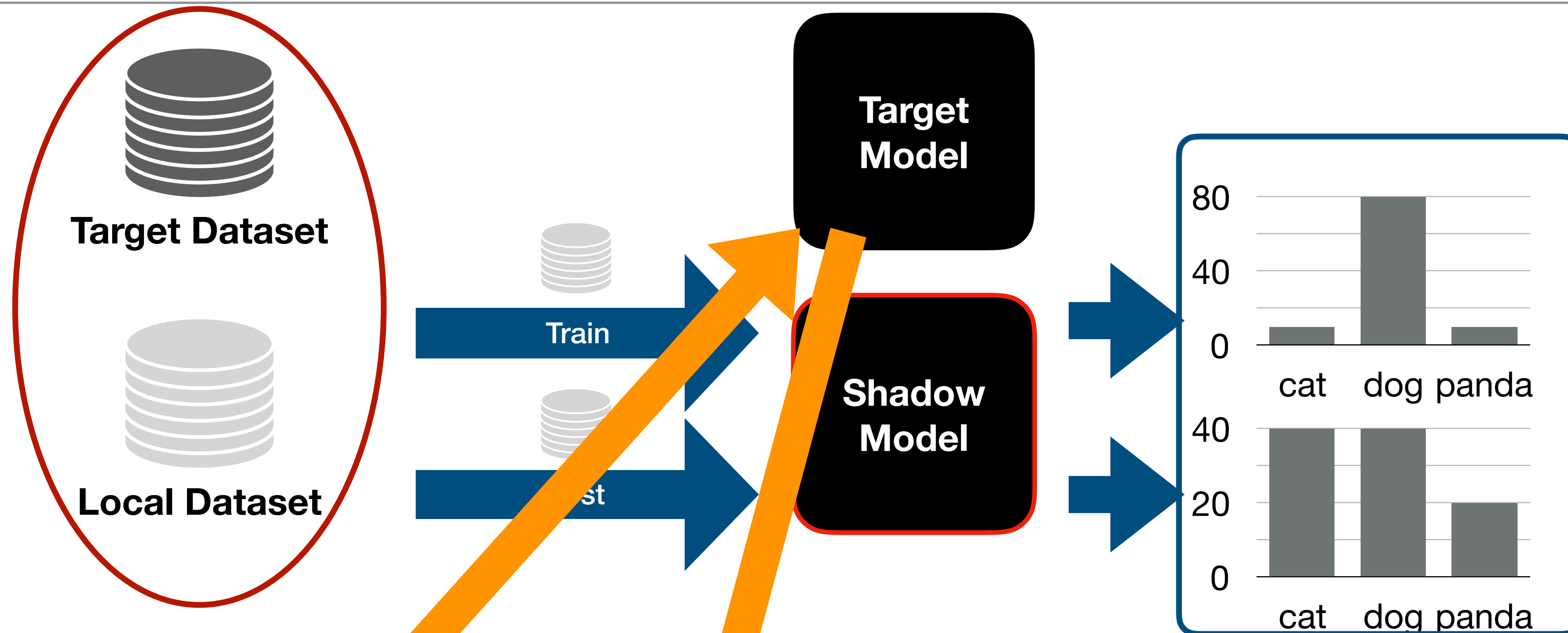
- Assumption on the datasets' distribution

# Towards a Data Independent Attack

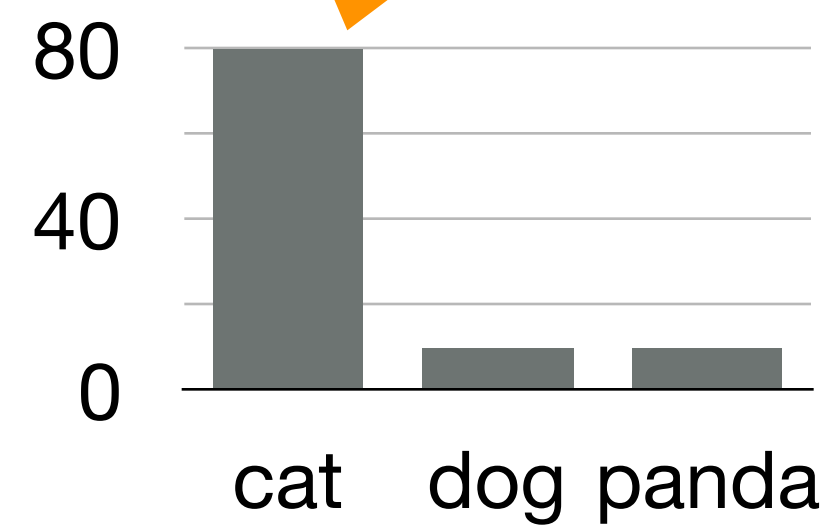
---

- Assumption on the datasets' distribution
  - Data transferring attack

# Data Transferring Attack (Adversary 2)



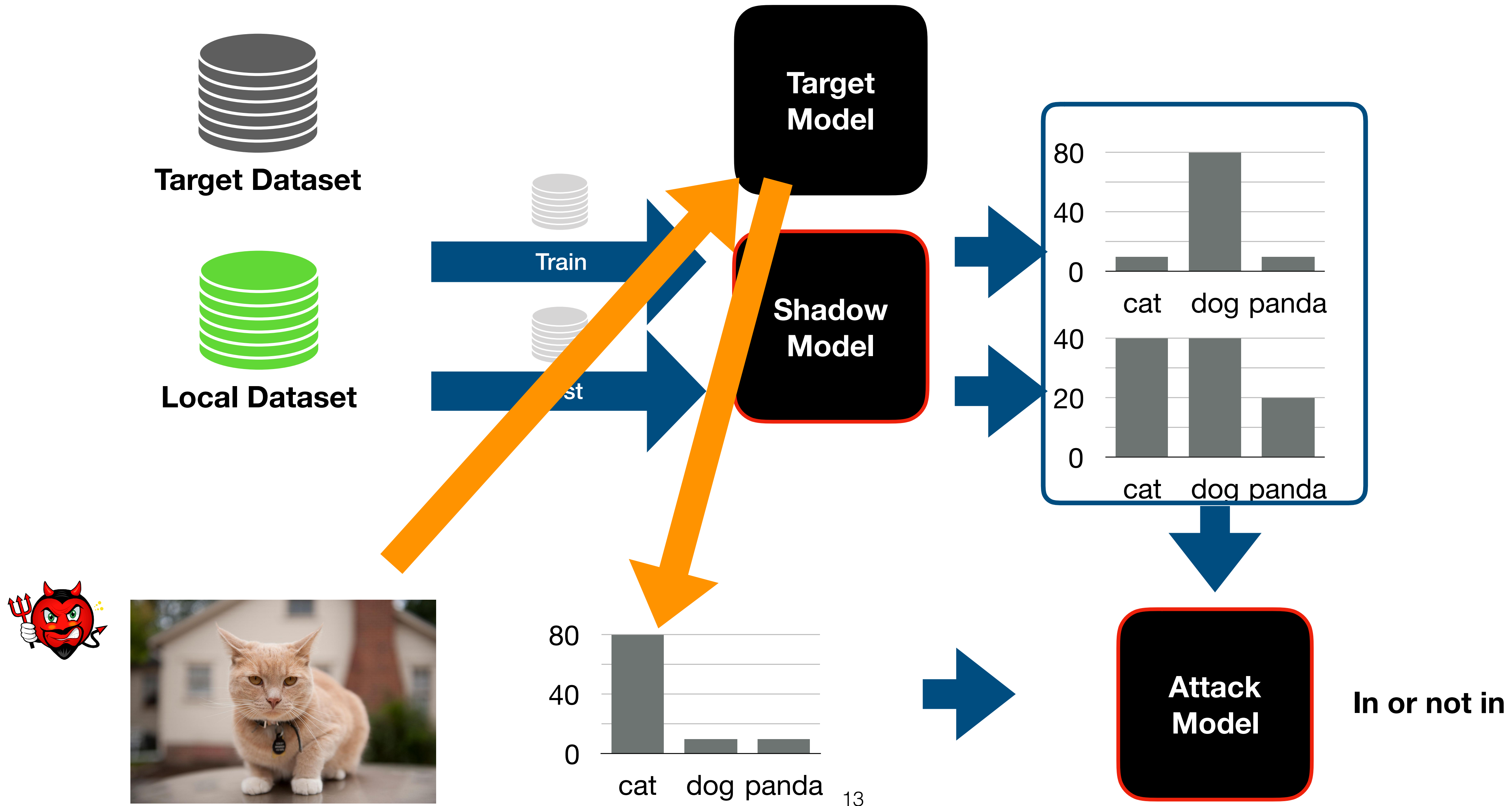
Same Distribution



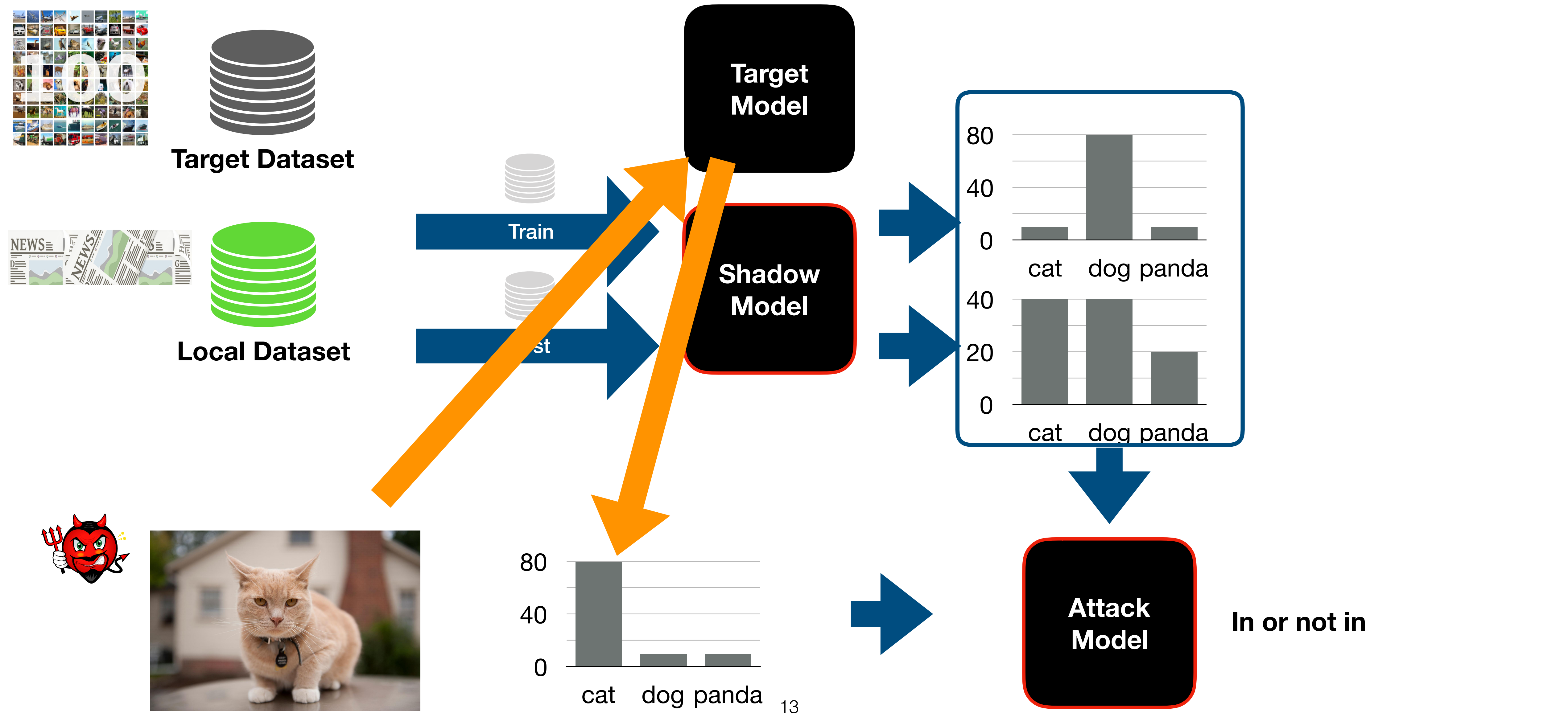
In or not in



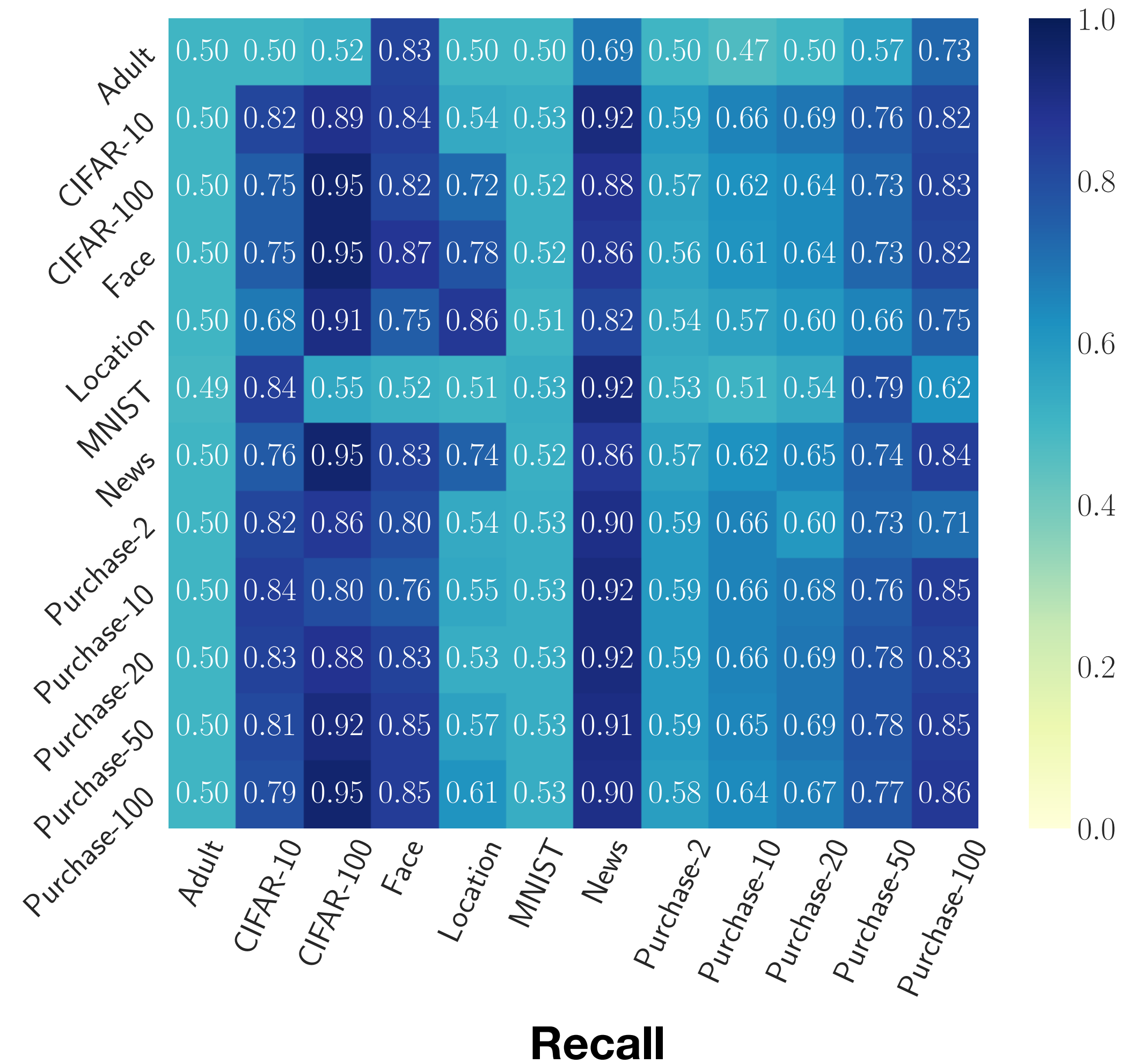
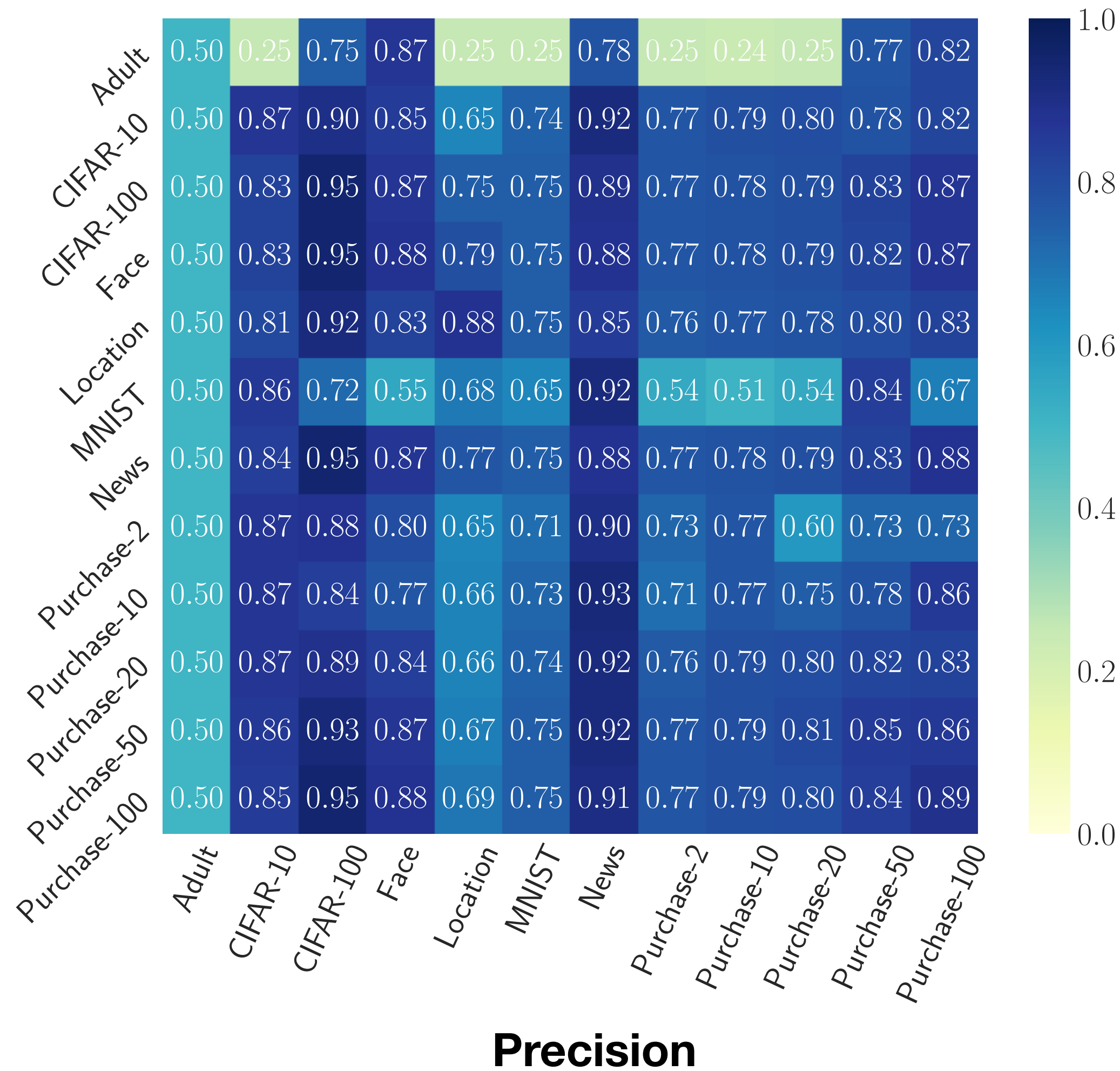
# Data Transferring Attack (Adversary 2)



# Data Transferring Attack (Adversary 2)

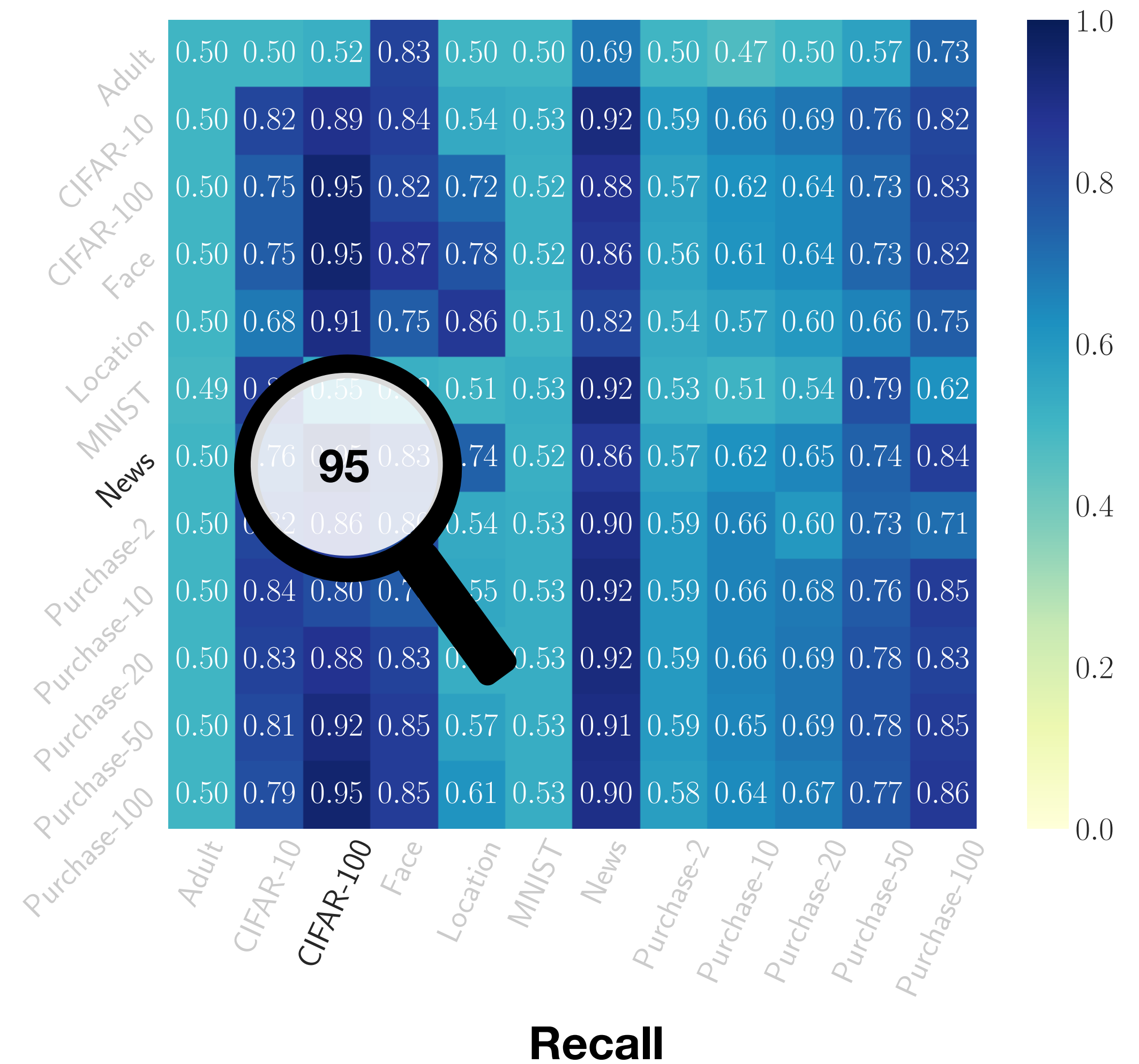
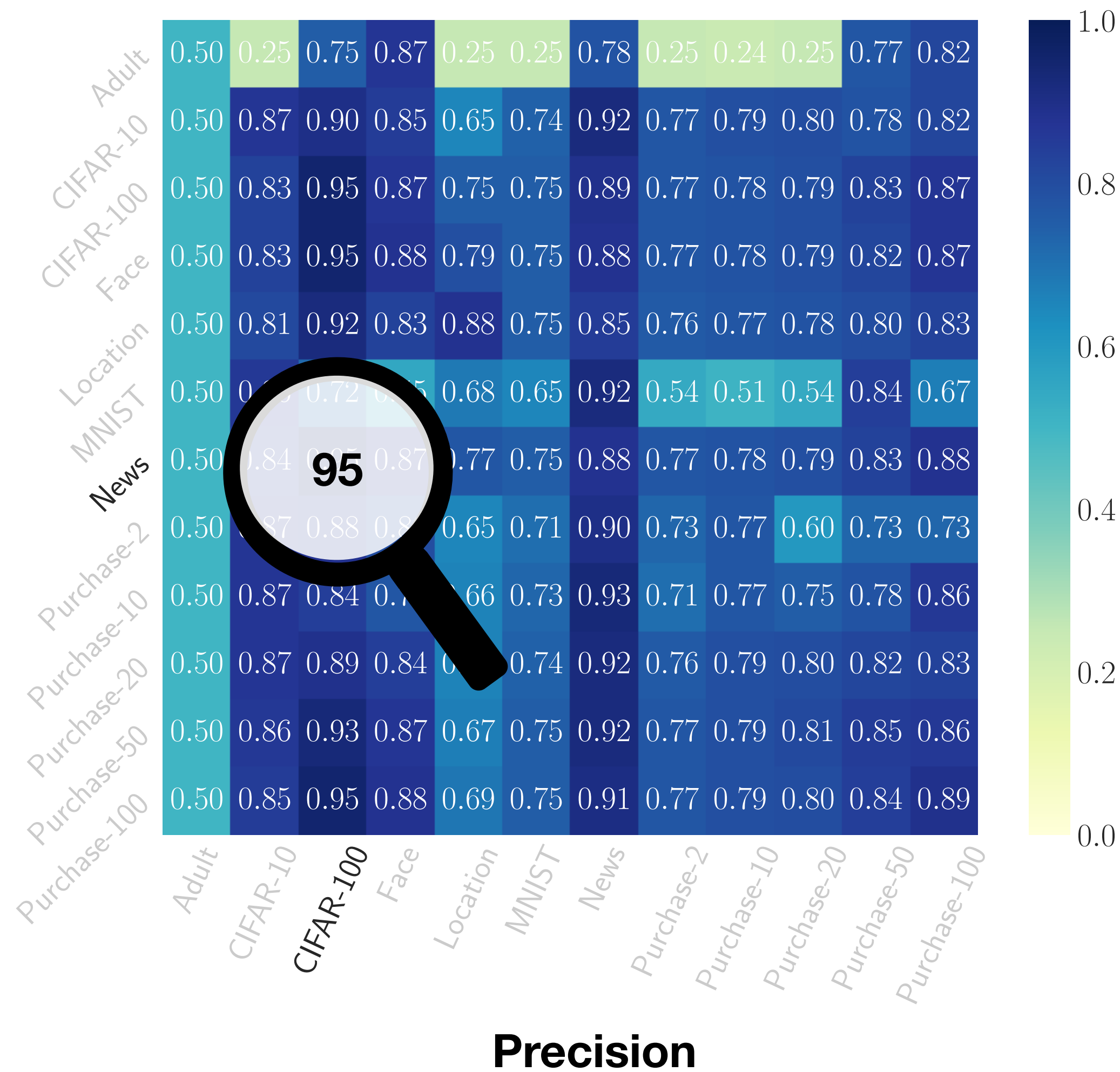


# Performance

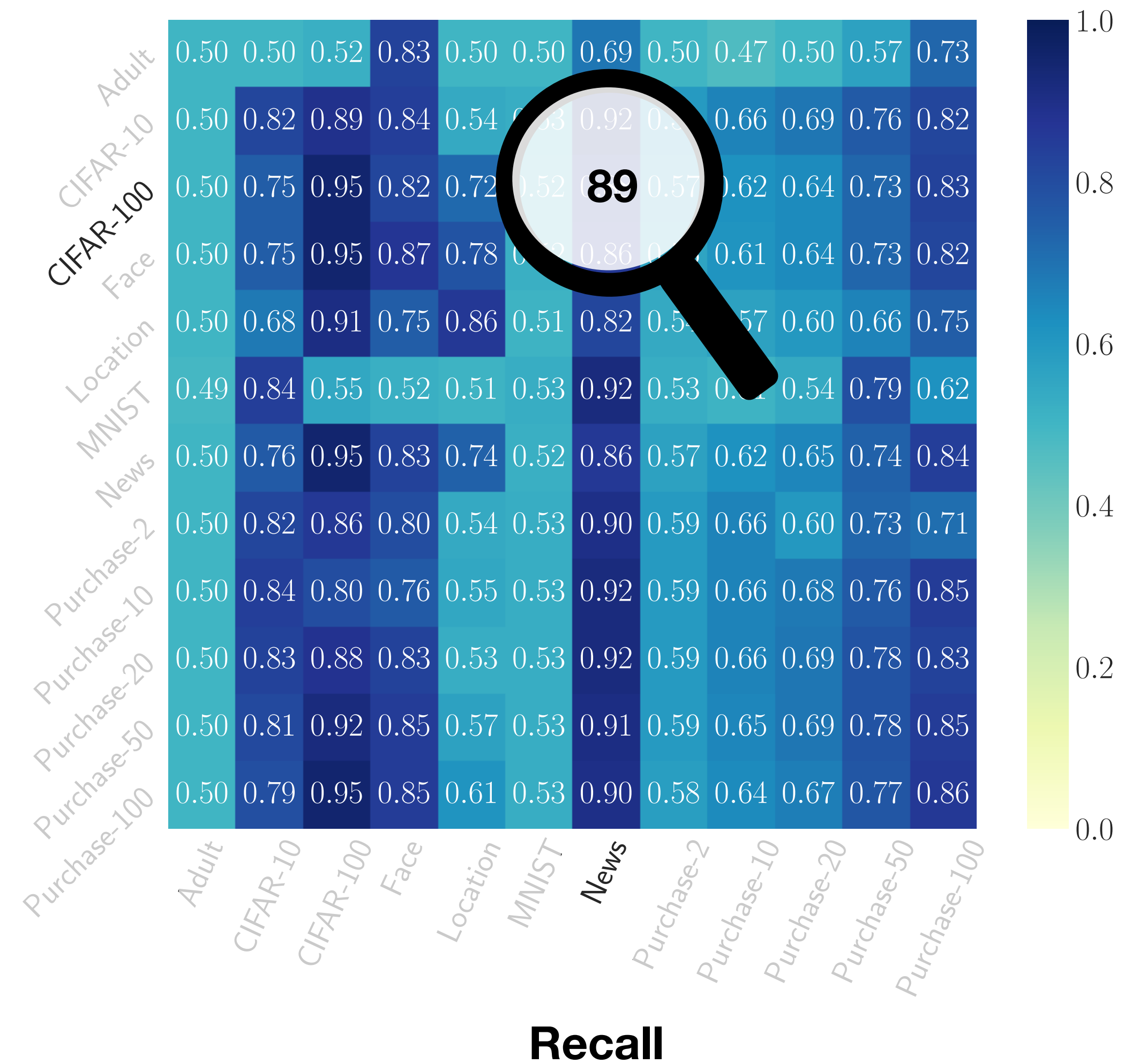
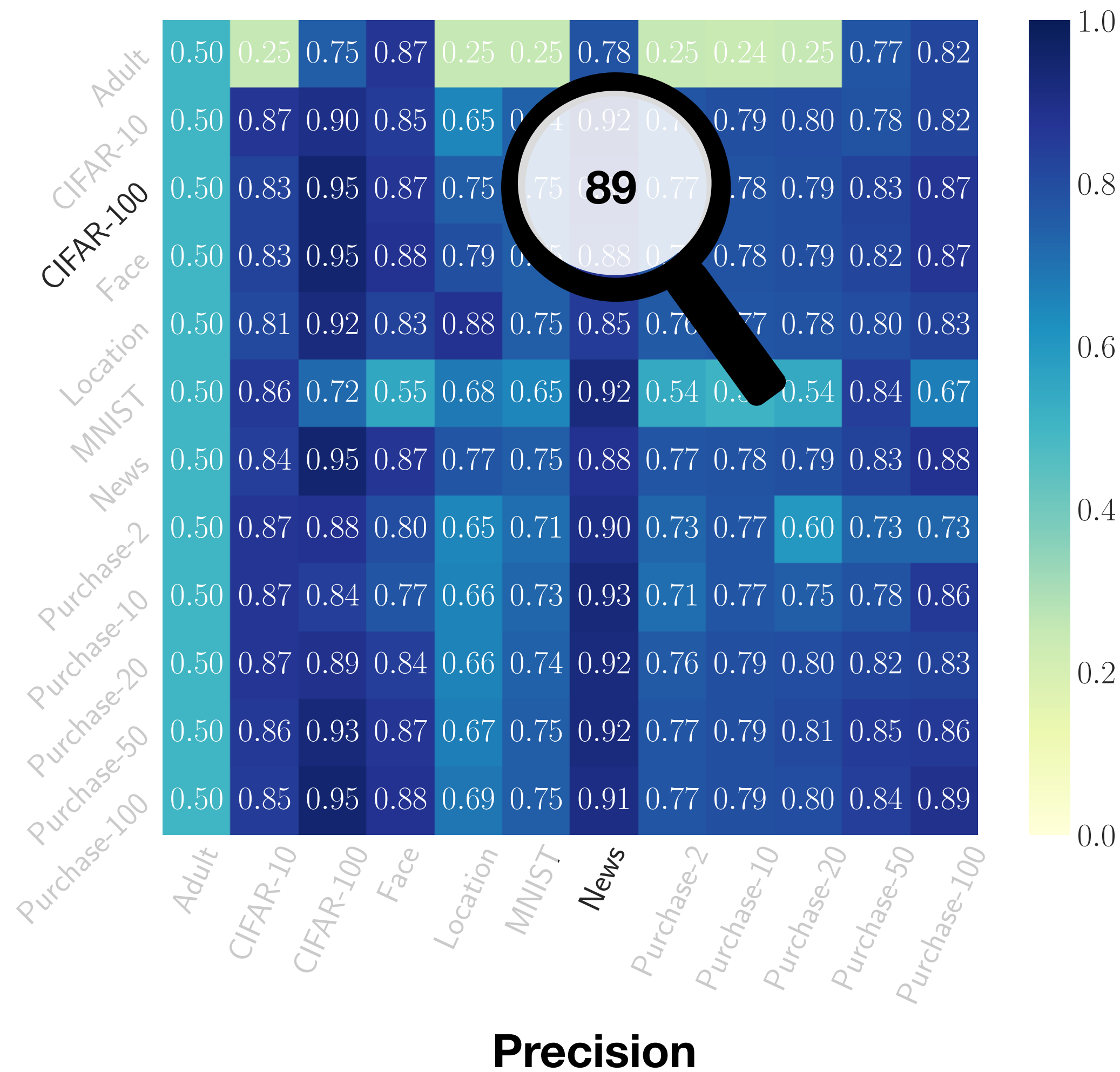




# Performance

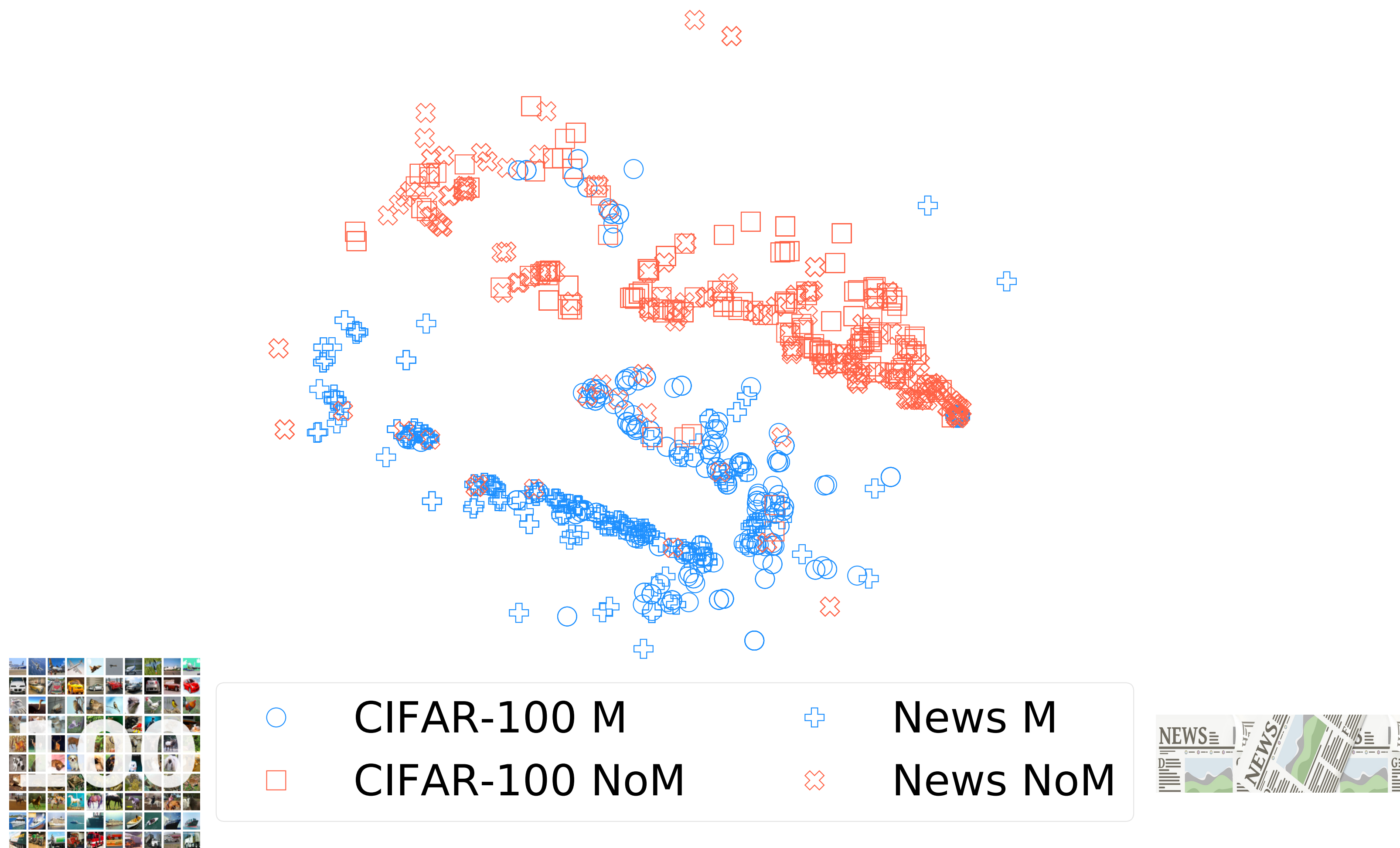


# Performance

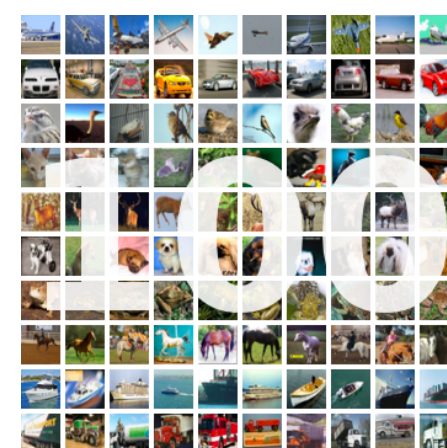
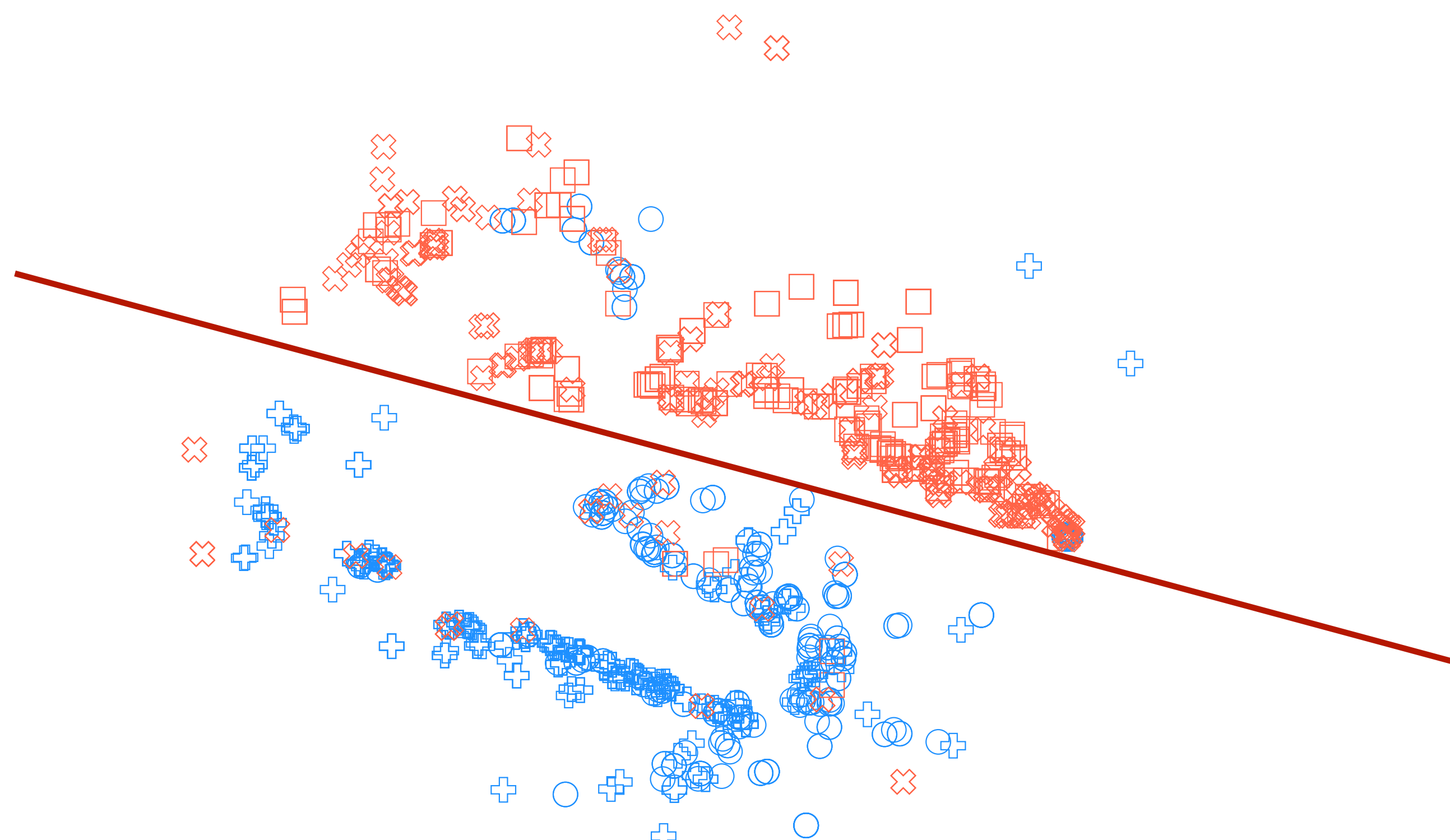




# Sounds Magic, Why?



# Sounds Magic, Why?



○	CIFAR-100 M	+	News M
□	CIFAR-100 NoM	×	News NoM



# Towards an Unsupervised Attack

---

# Towards an Unsupervised Attack

---

- Shadow model

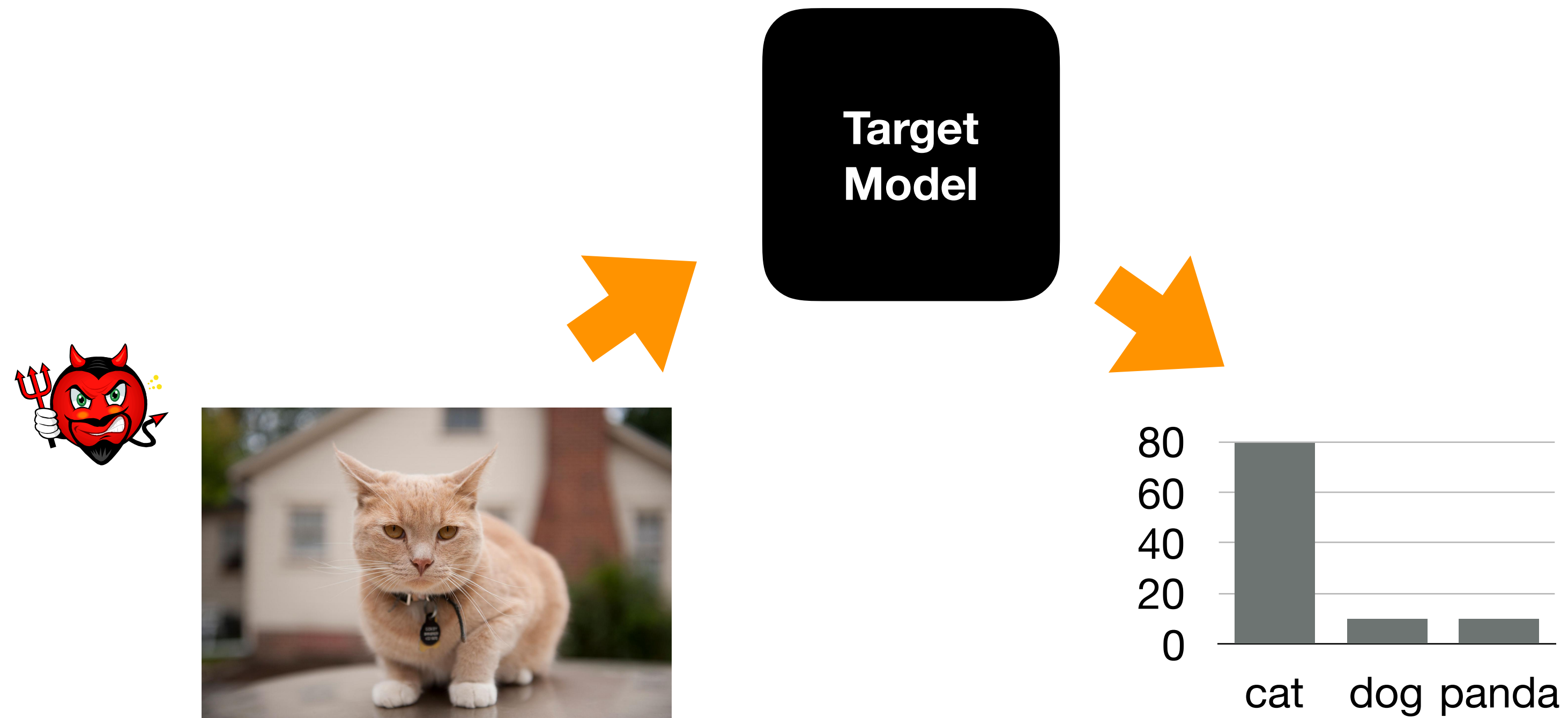
# Towards an Unsupervised Attack

---

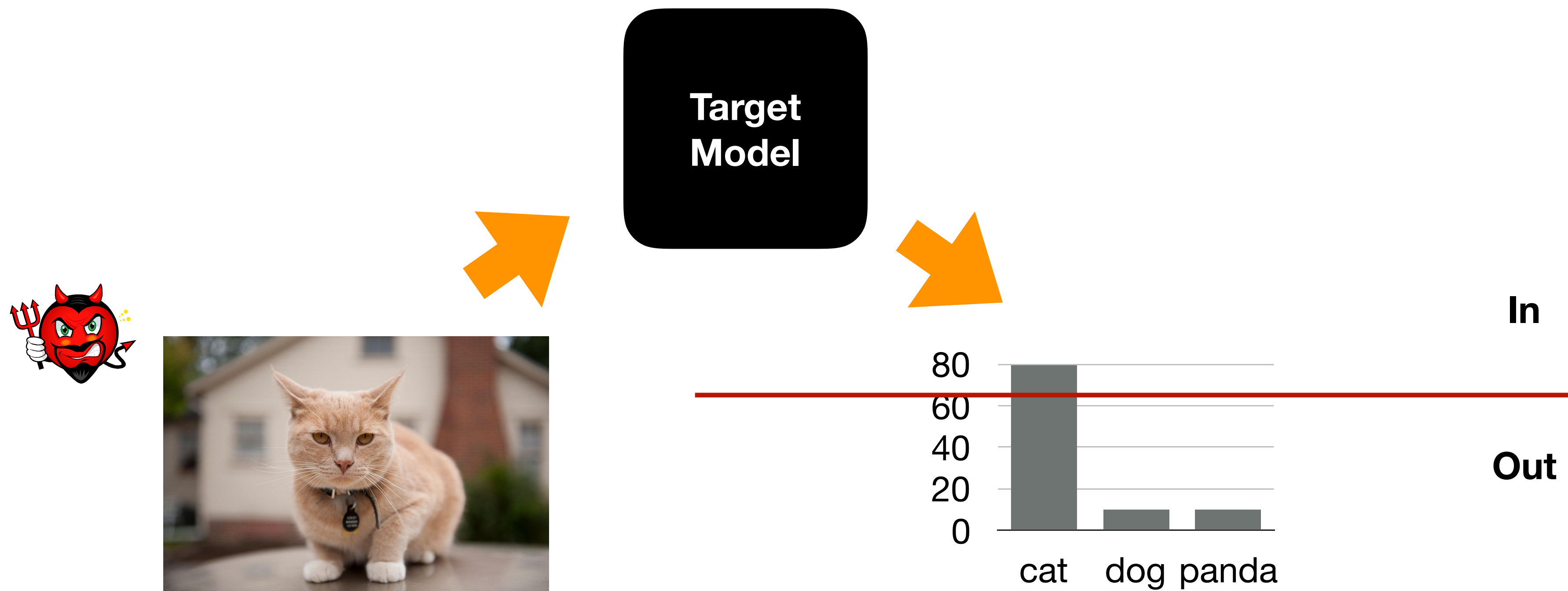
- Shadow model
- Attack model



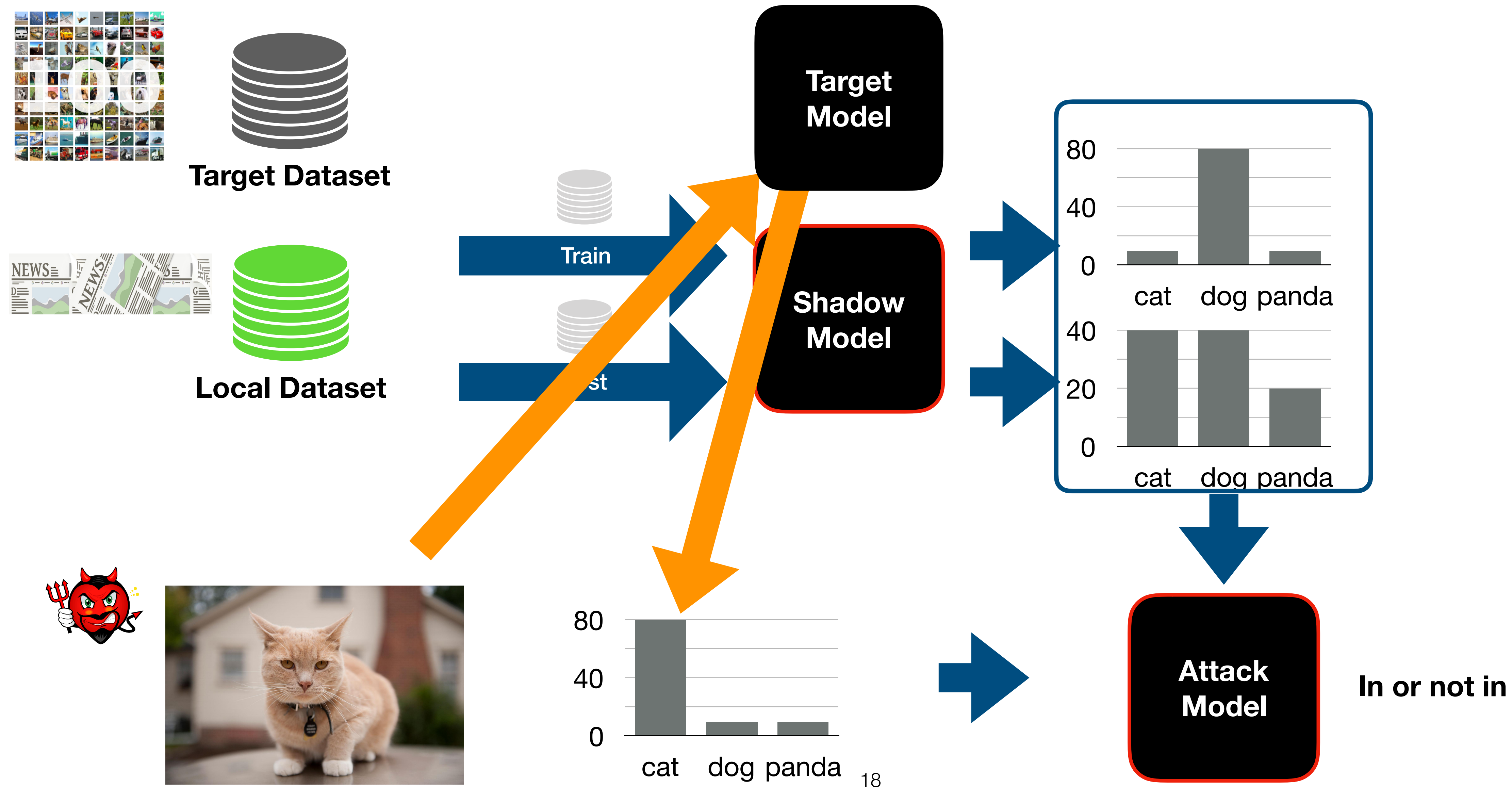
# Our Third Attack (Adversary 3)



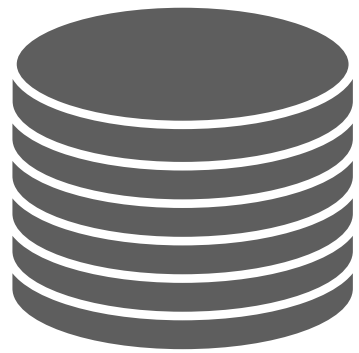
# Our Third Attack (Adversary 3)



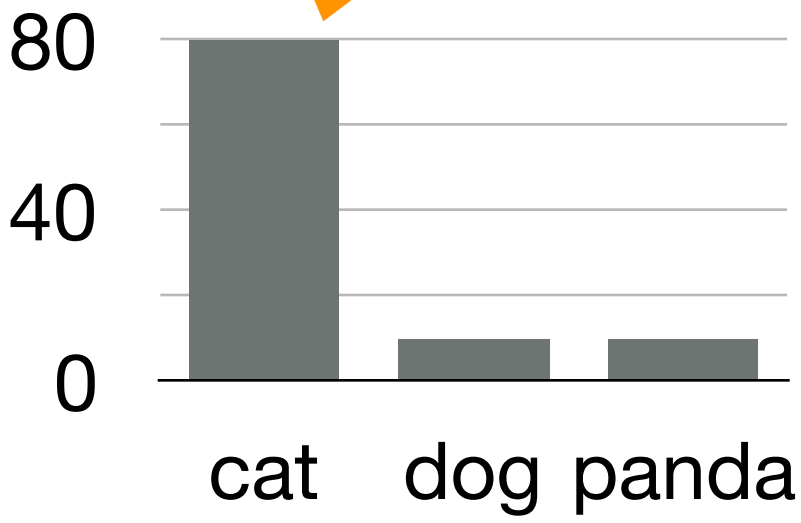
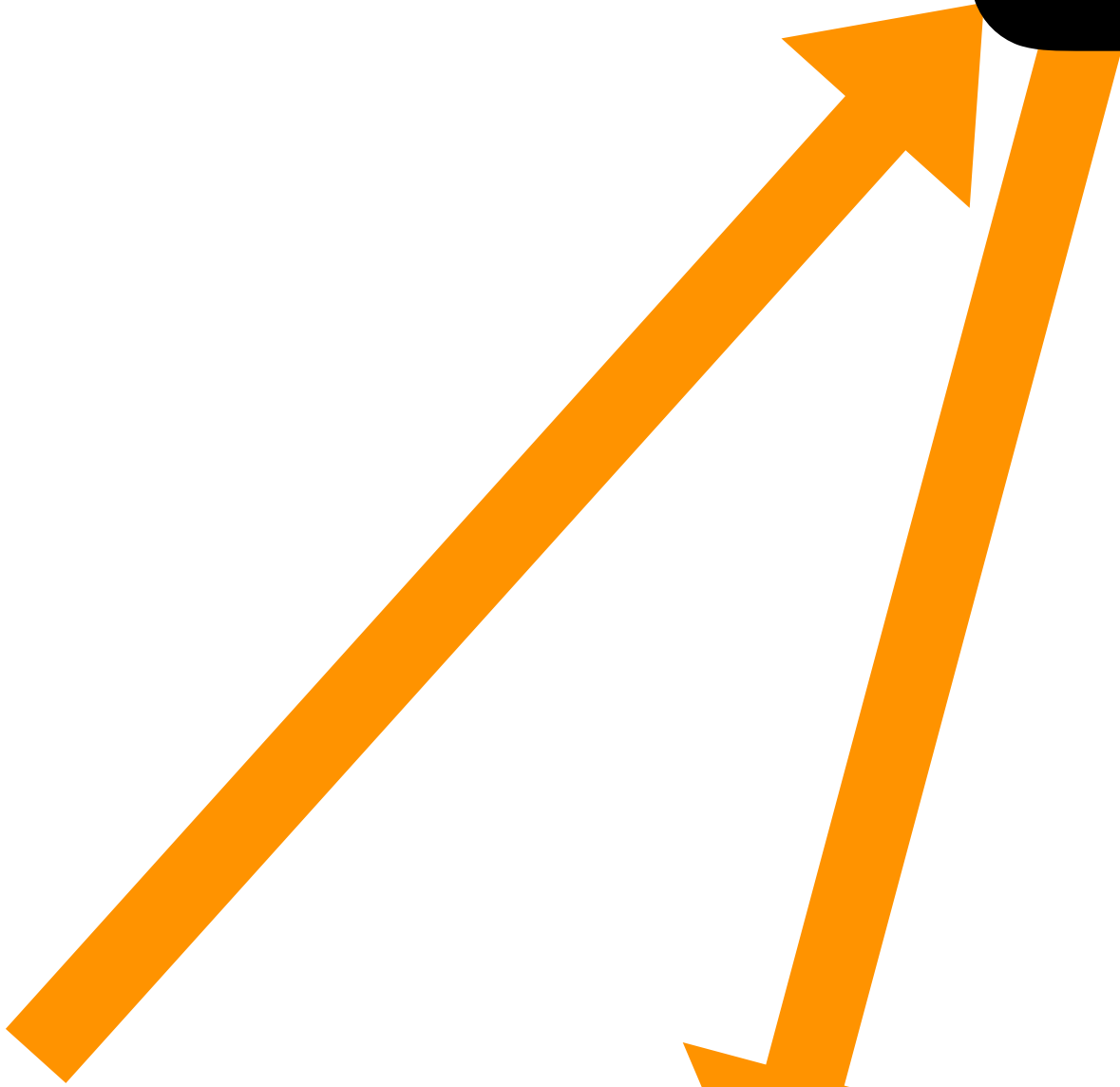
# Our Third Attack (Adversary 3)



# Our Third Attack (Adversary 3)

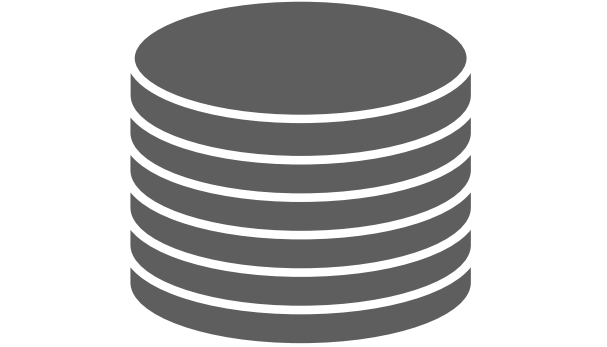


Target Dataset

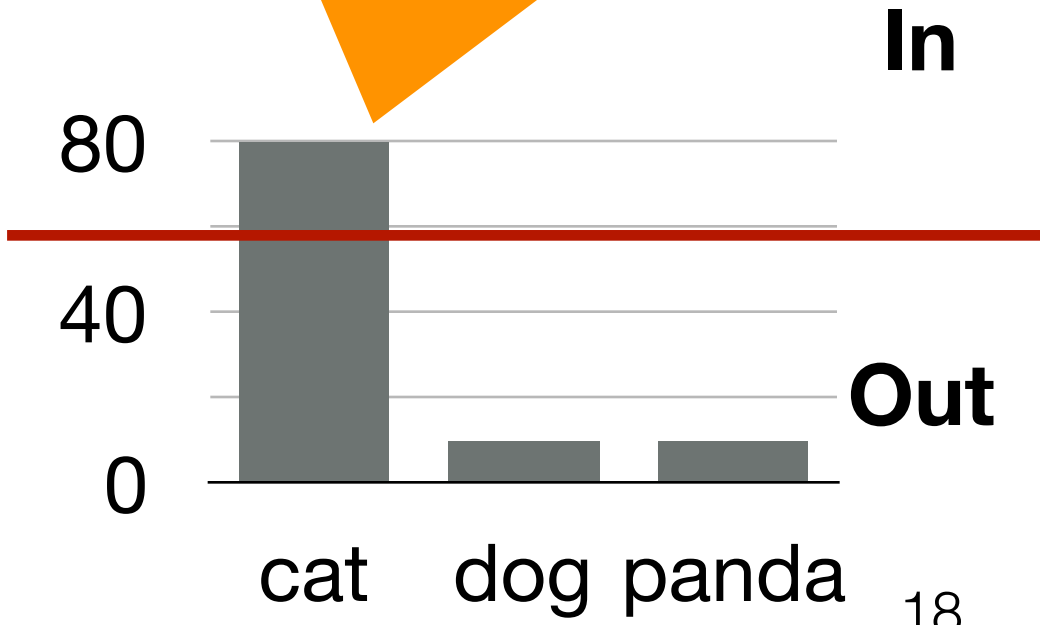
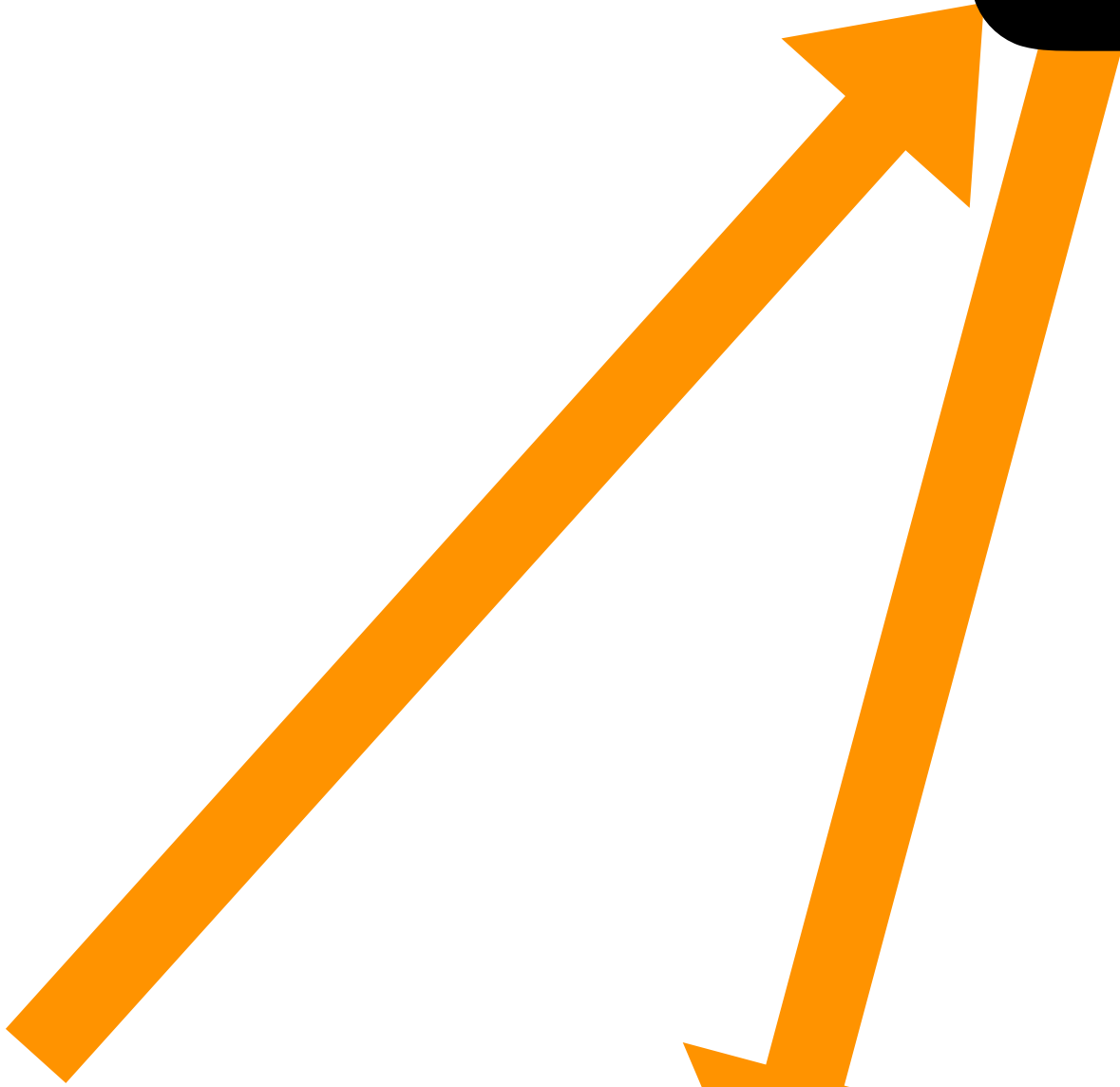




# Our Third Attack (Adversary 3)

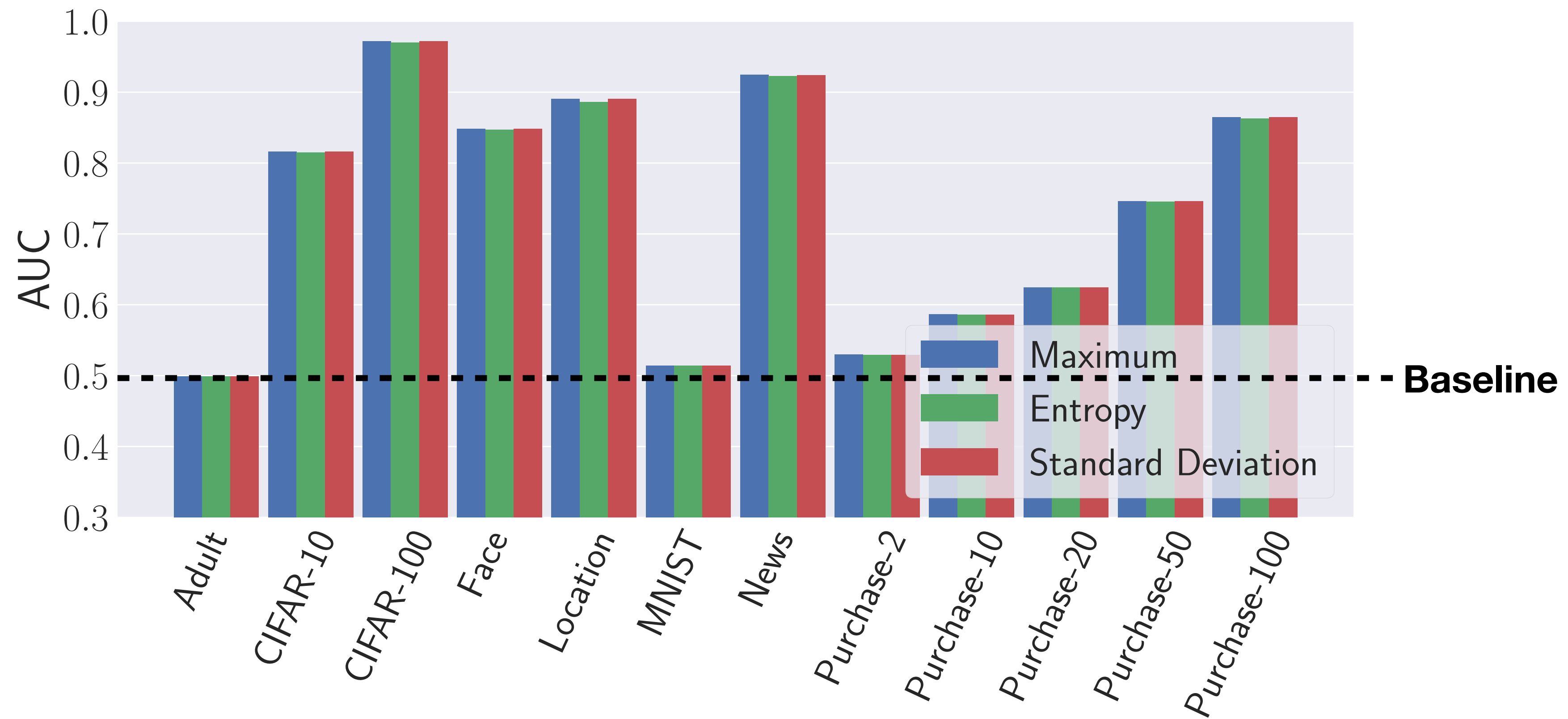


Target Dataset

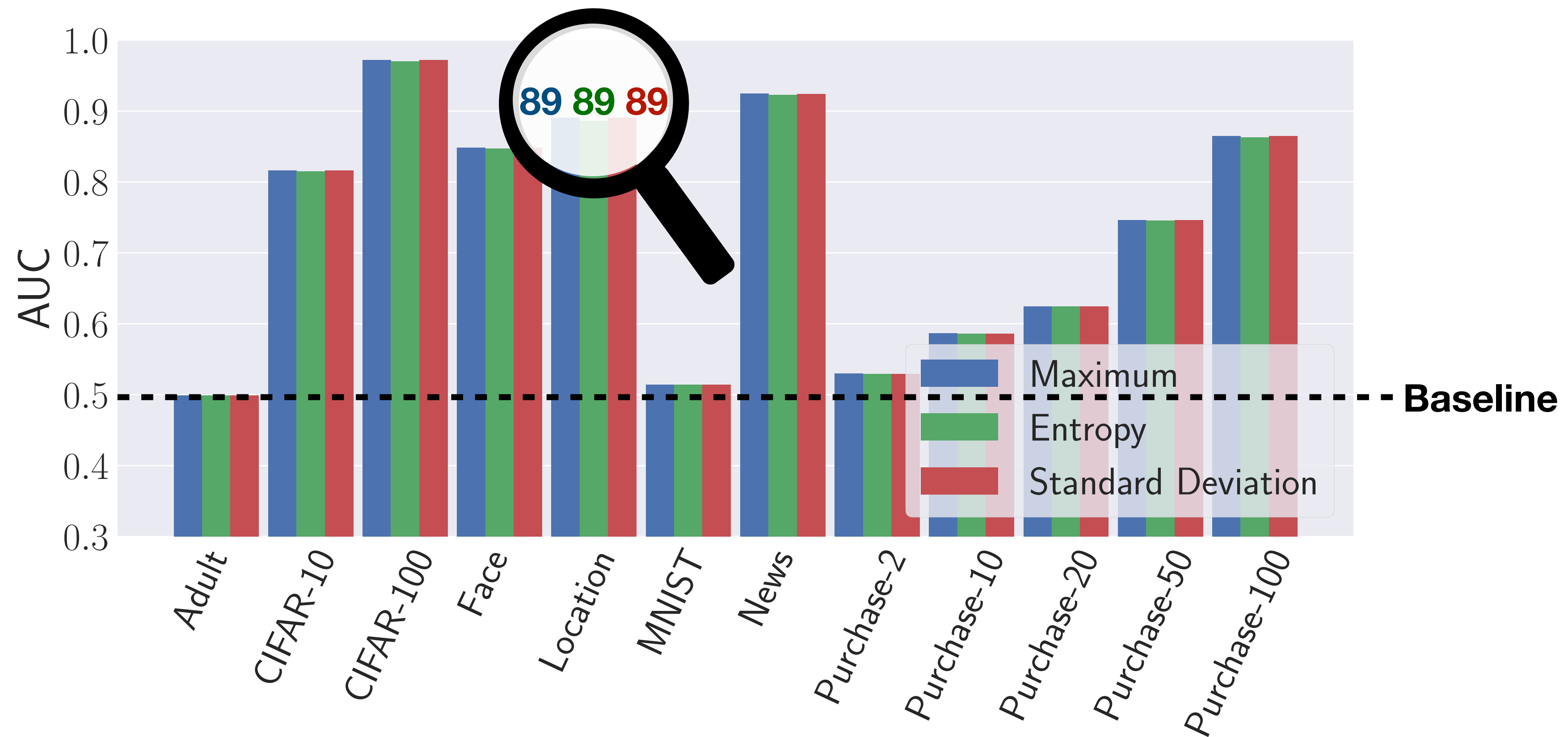




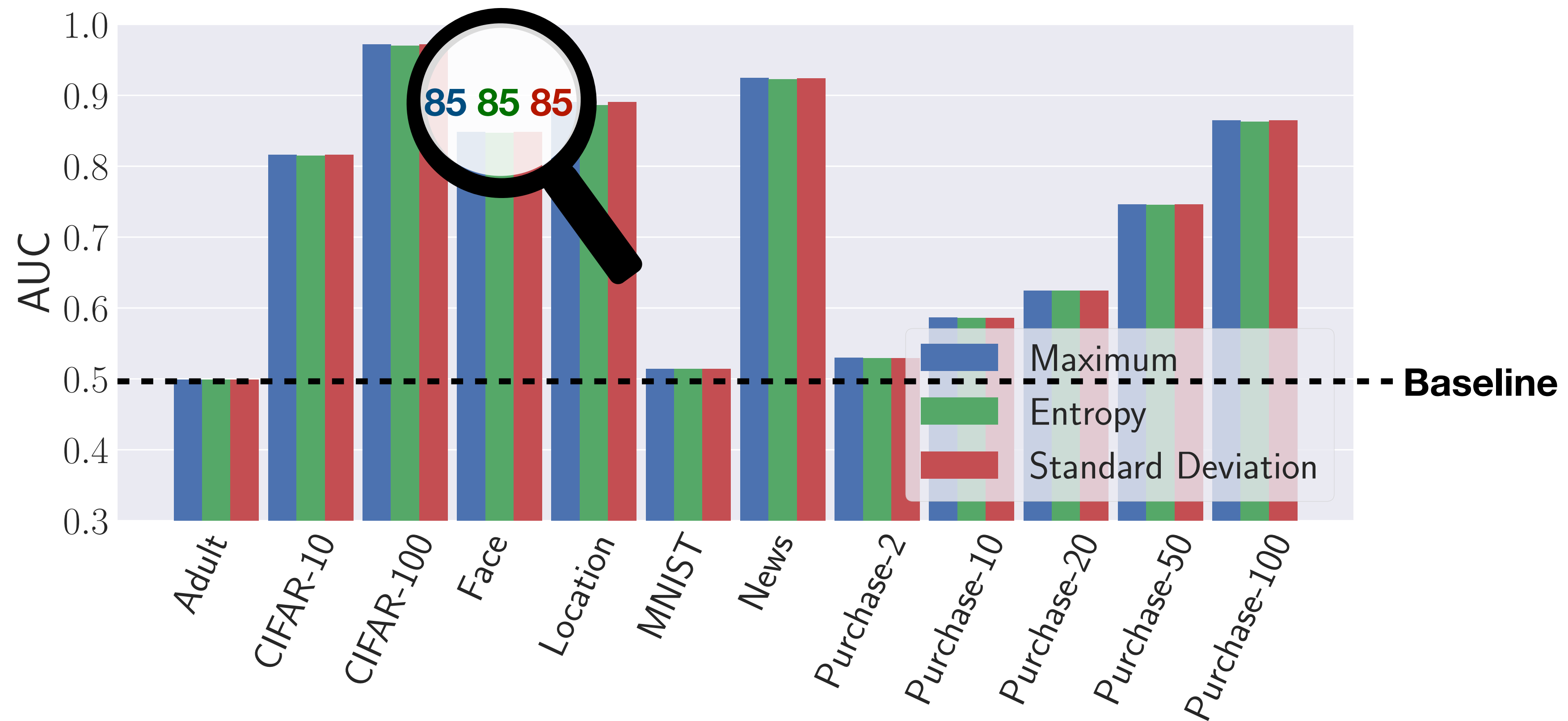
# Performance



# Performance



# Performance



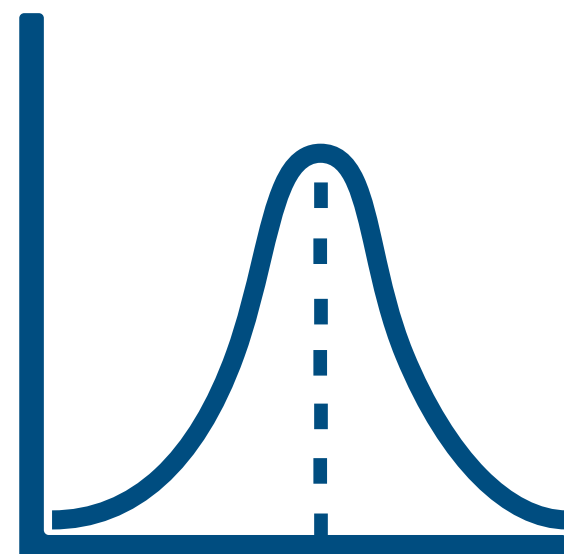
# Threshold Picking

---



**Target  
Model**

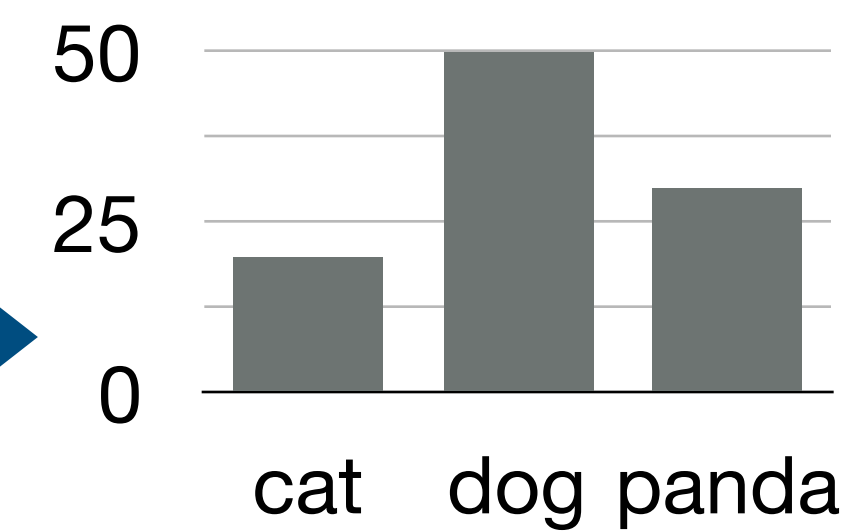
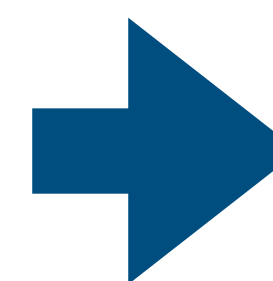
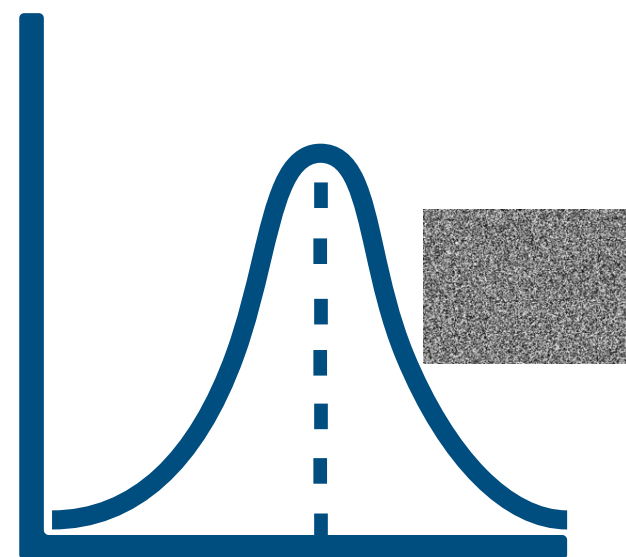
# Threshold Picking



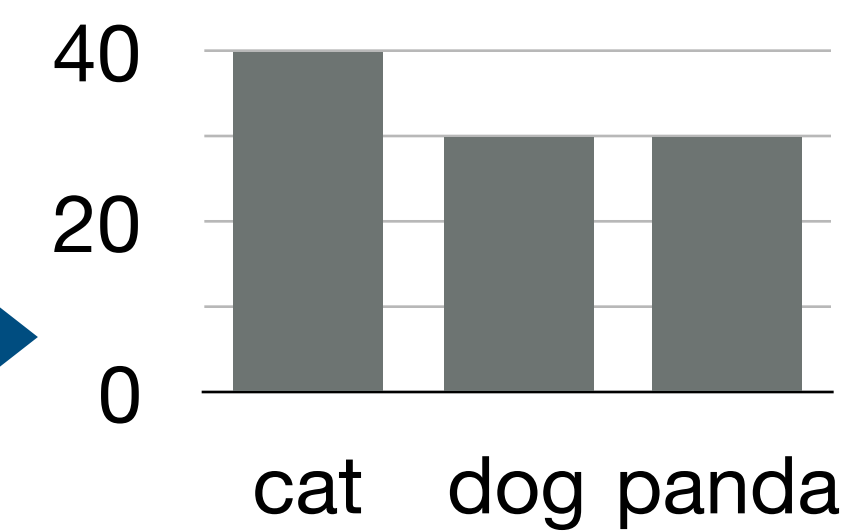
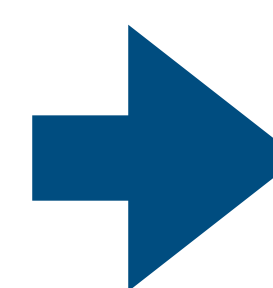
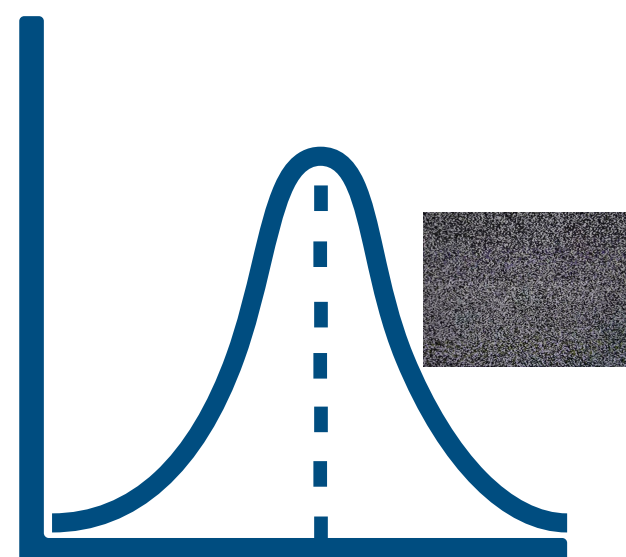
**Target  
Model**



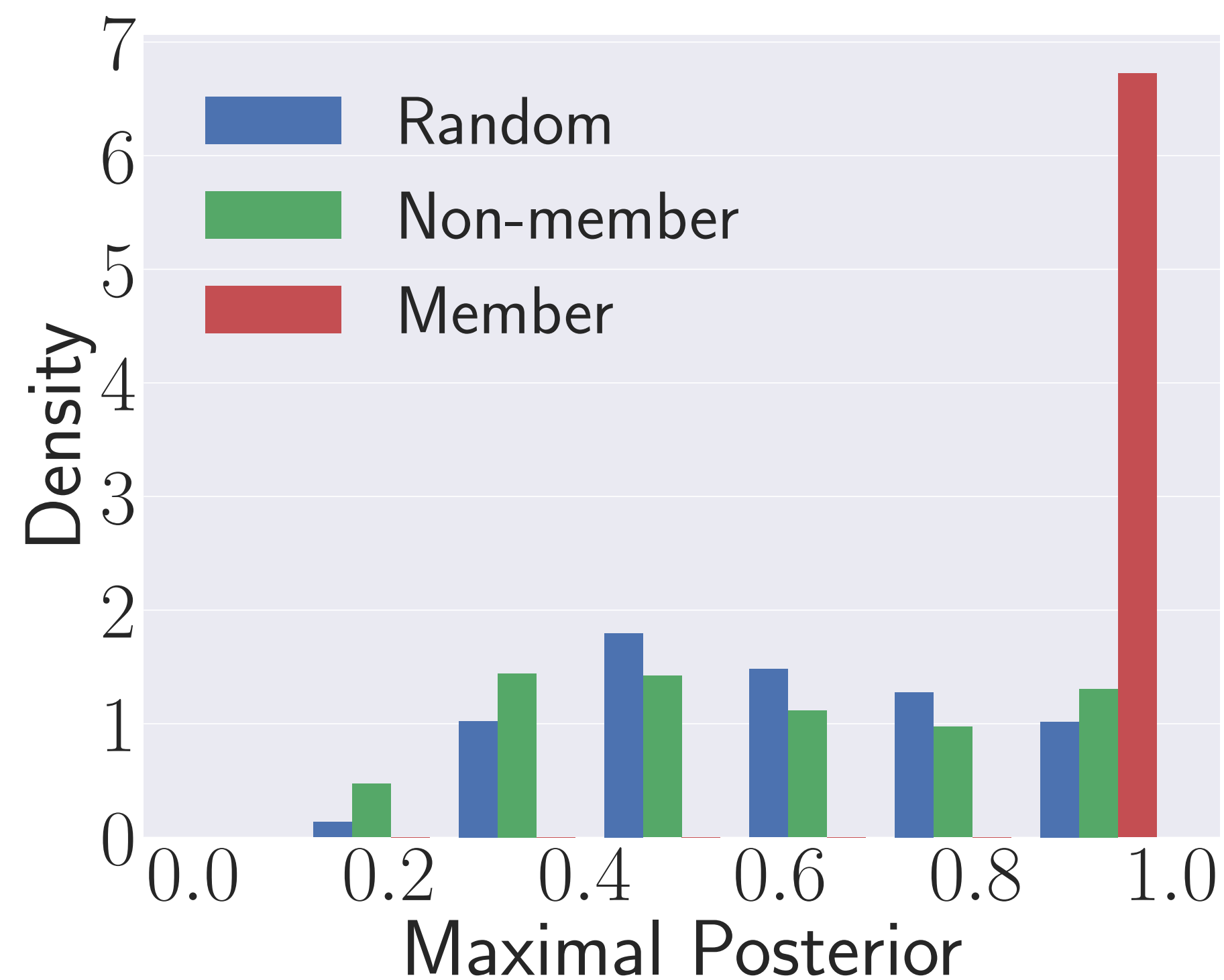
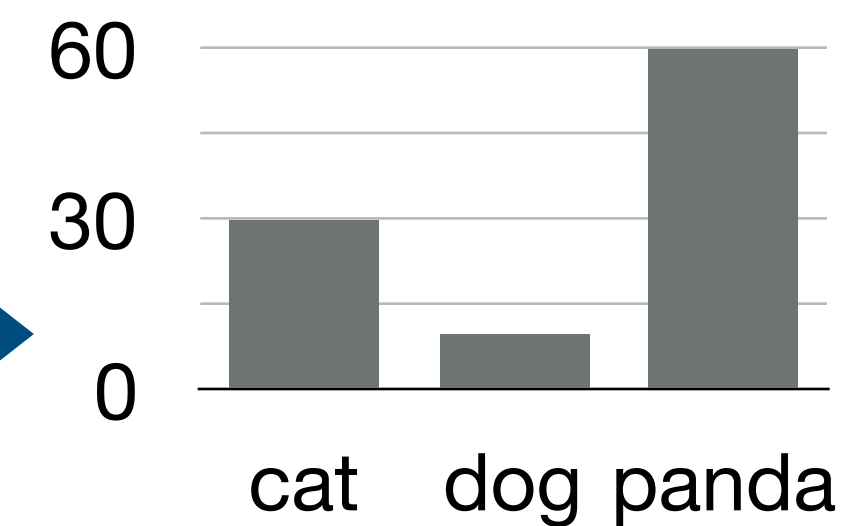
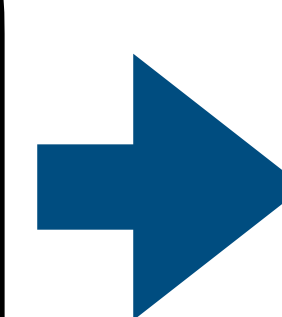
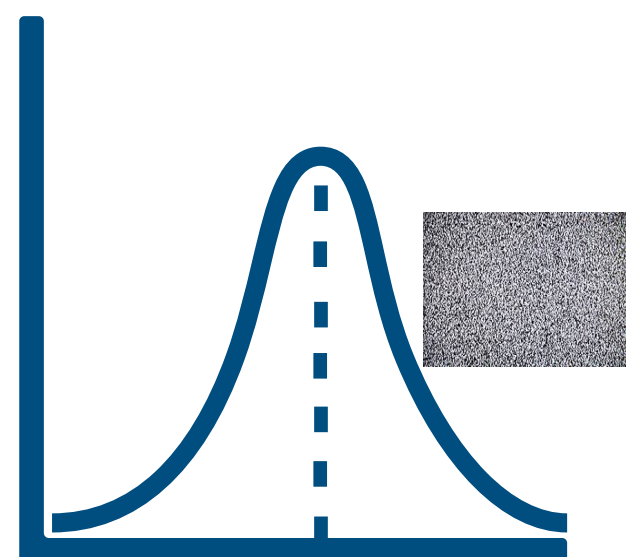
# Threshold Picking



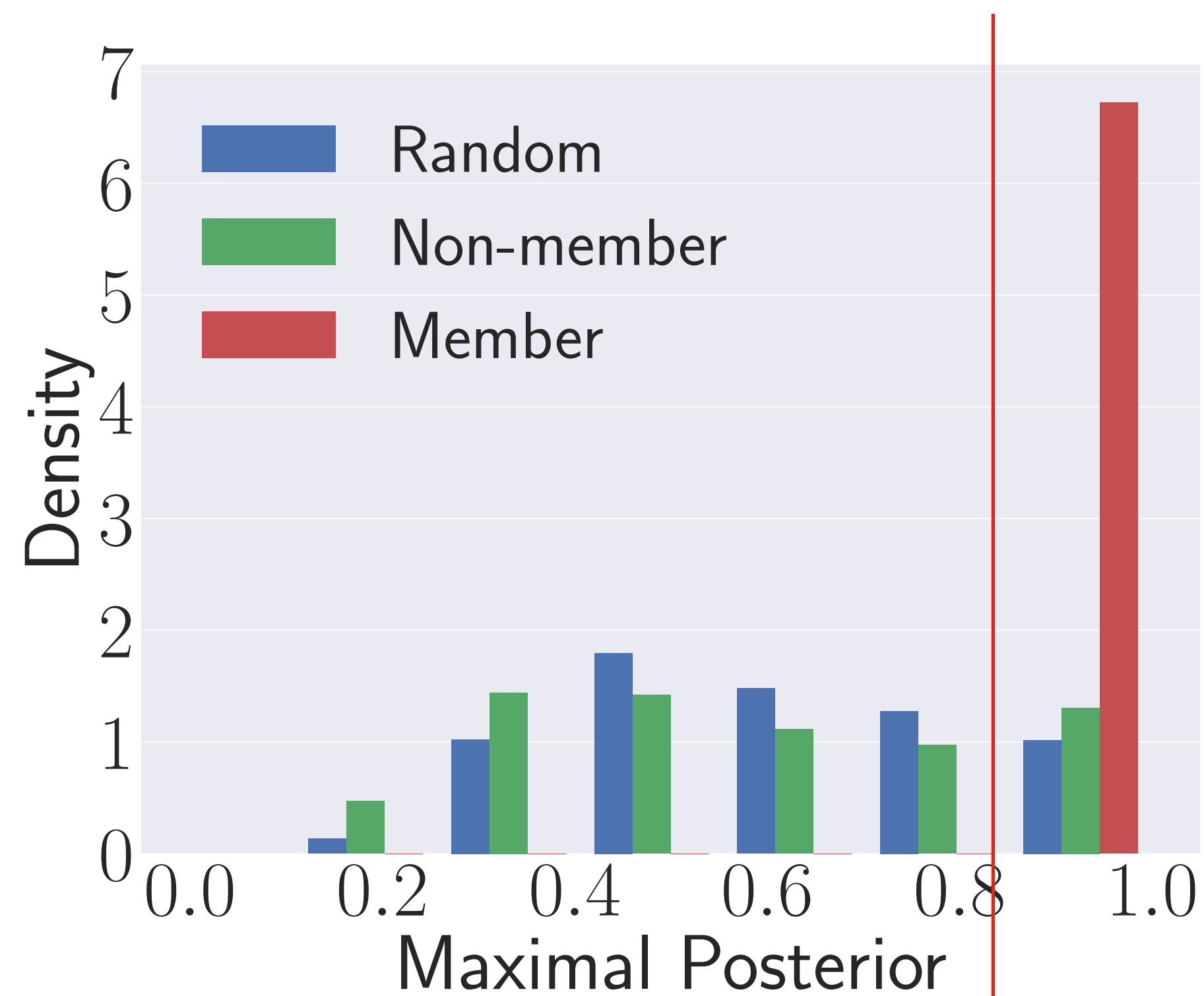
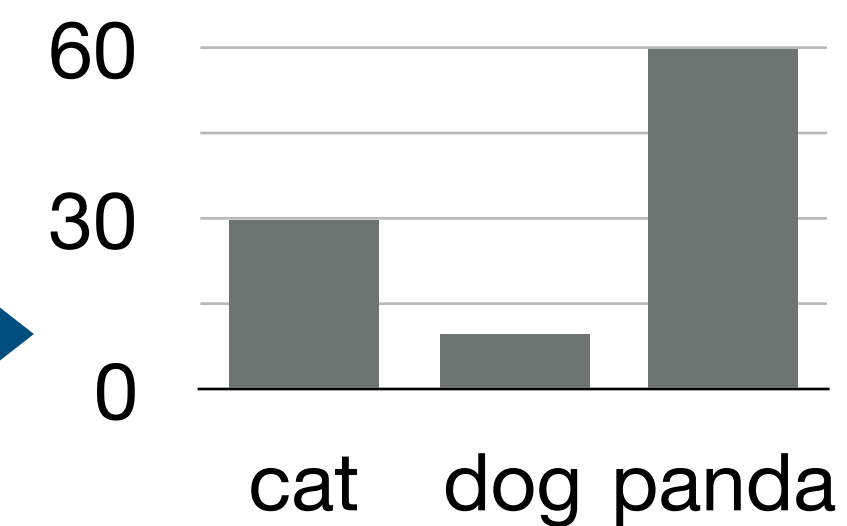
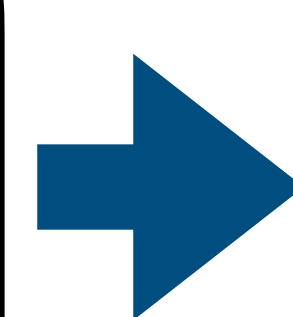
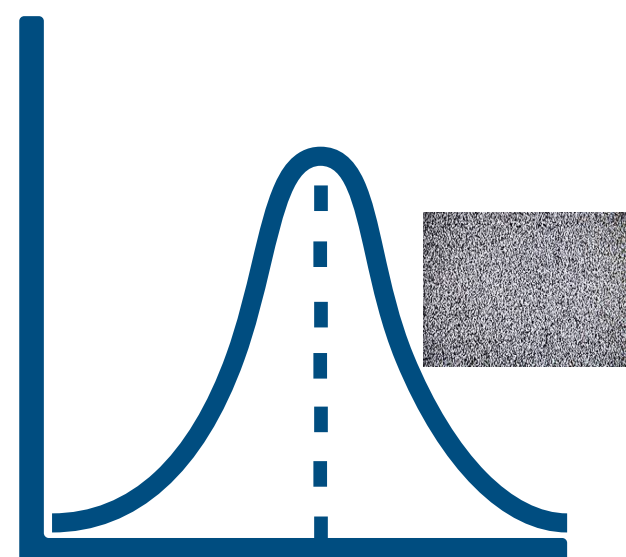
# Threshold Picking



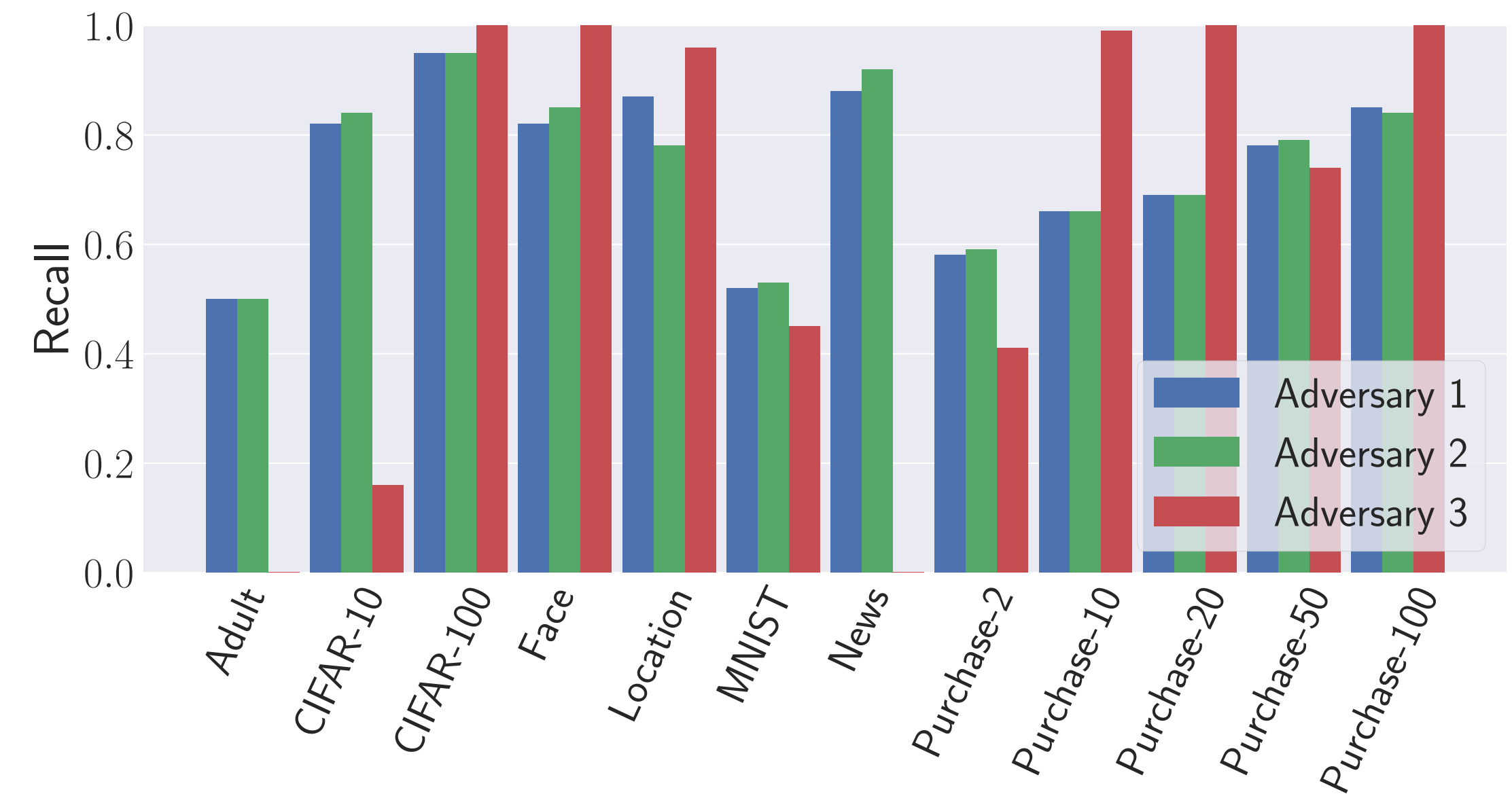
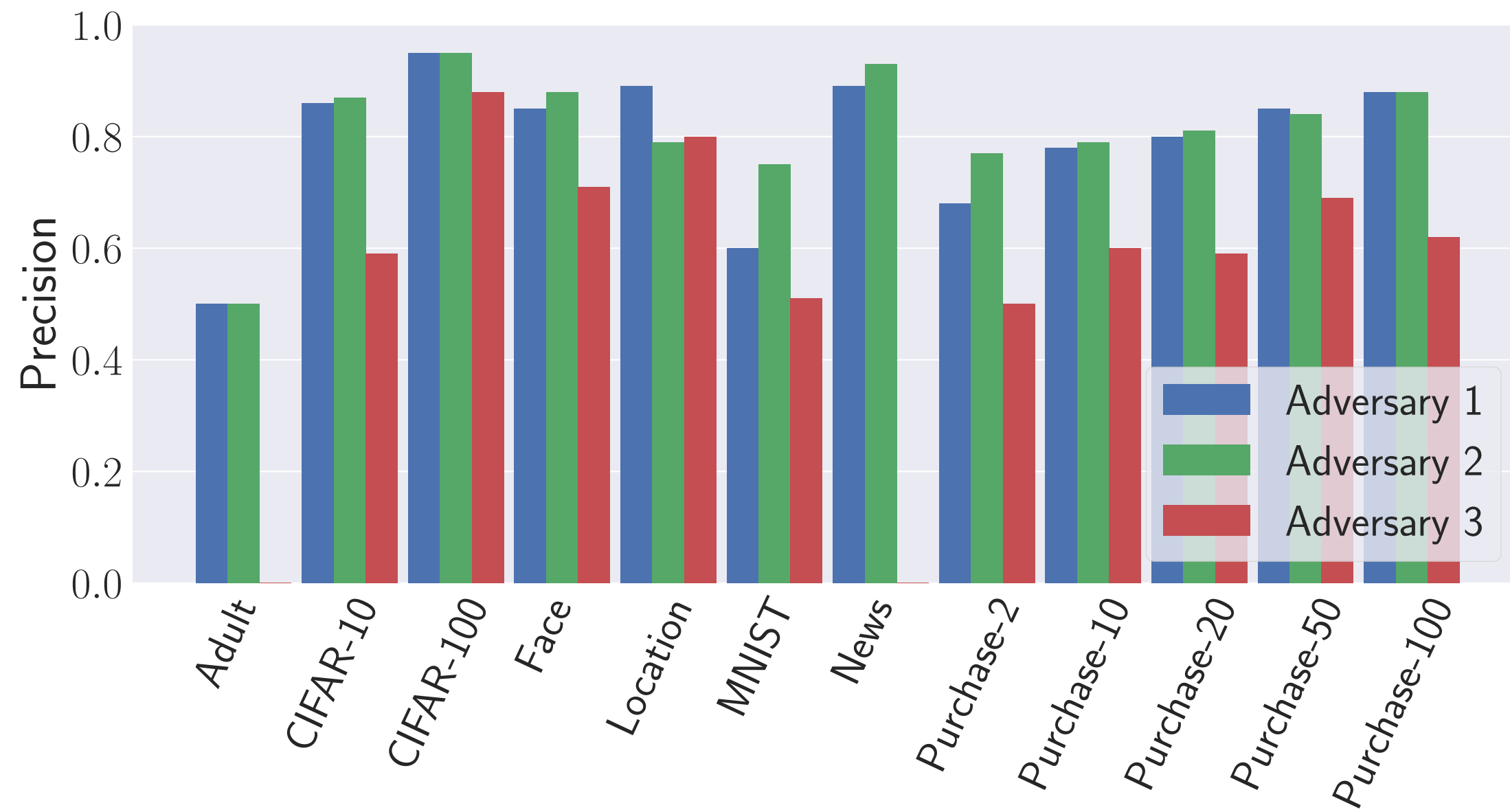
# Threshold Picking



# Threshold Picking

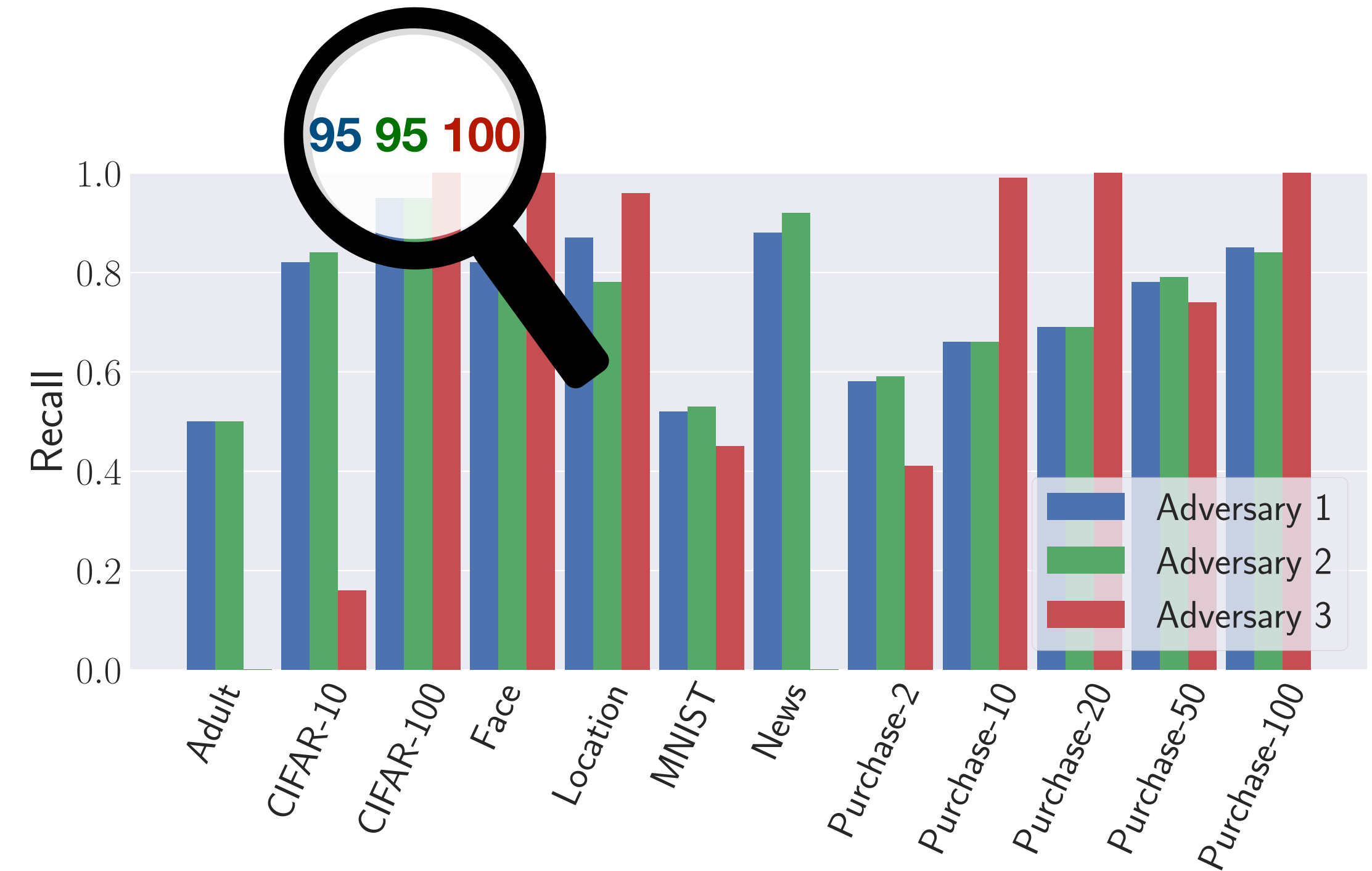
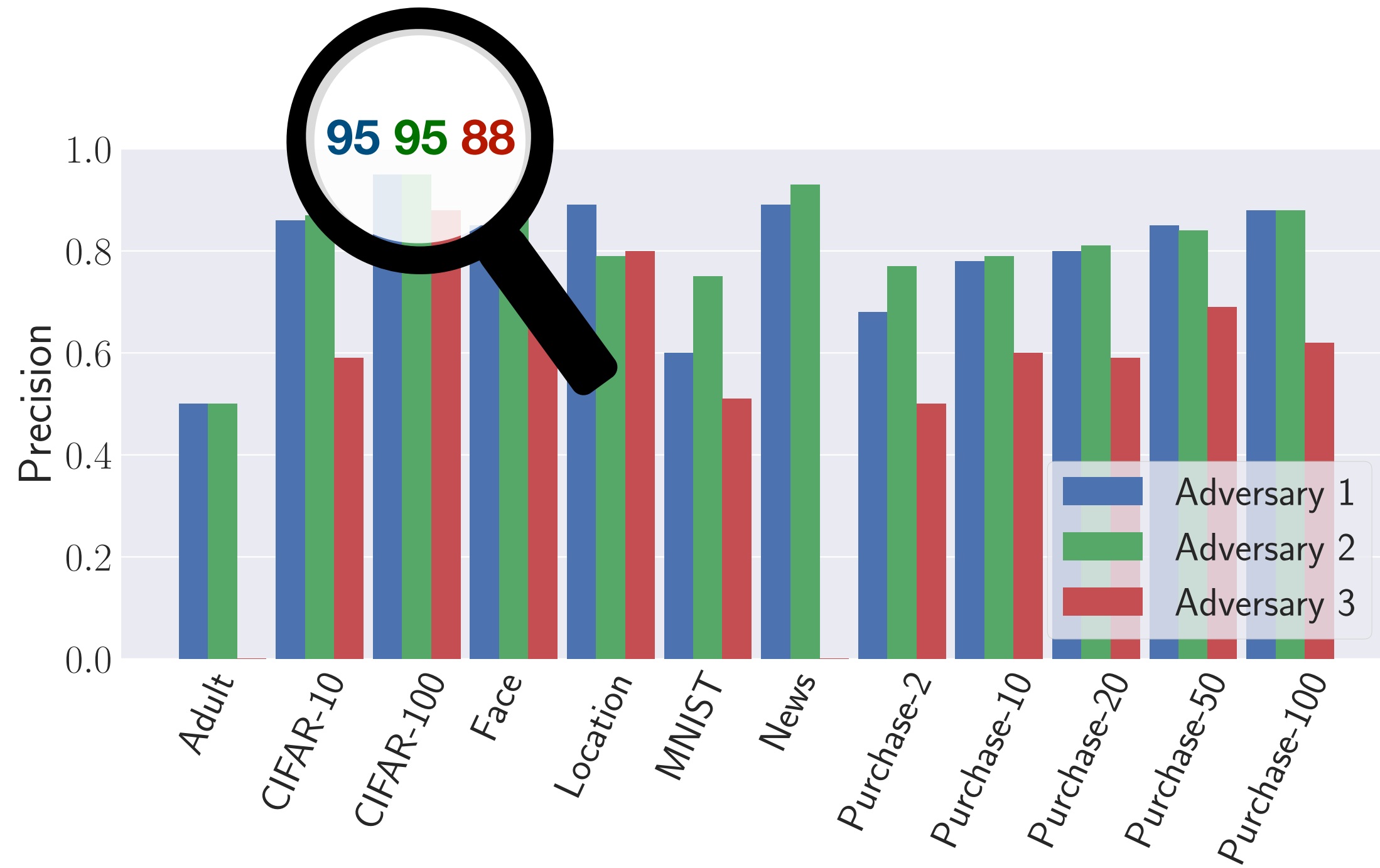


# Comparing All Attacks





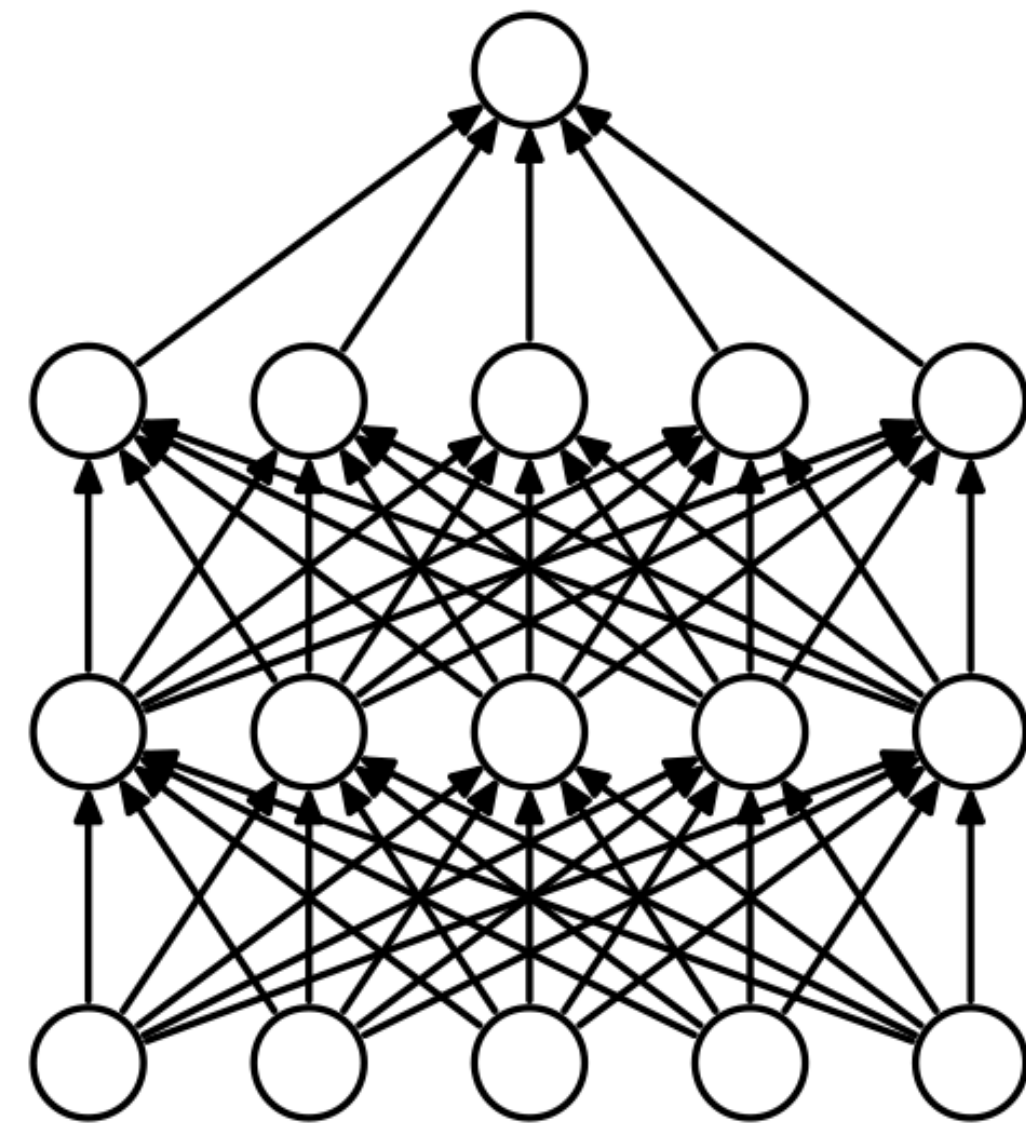
# Comparing All Attacks



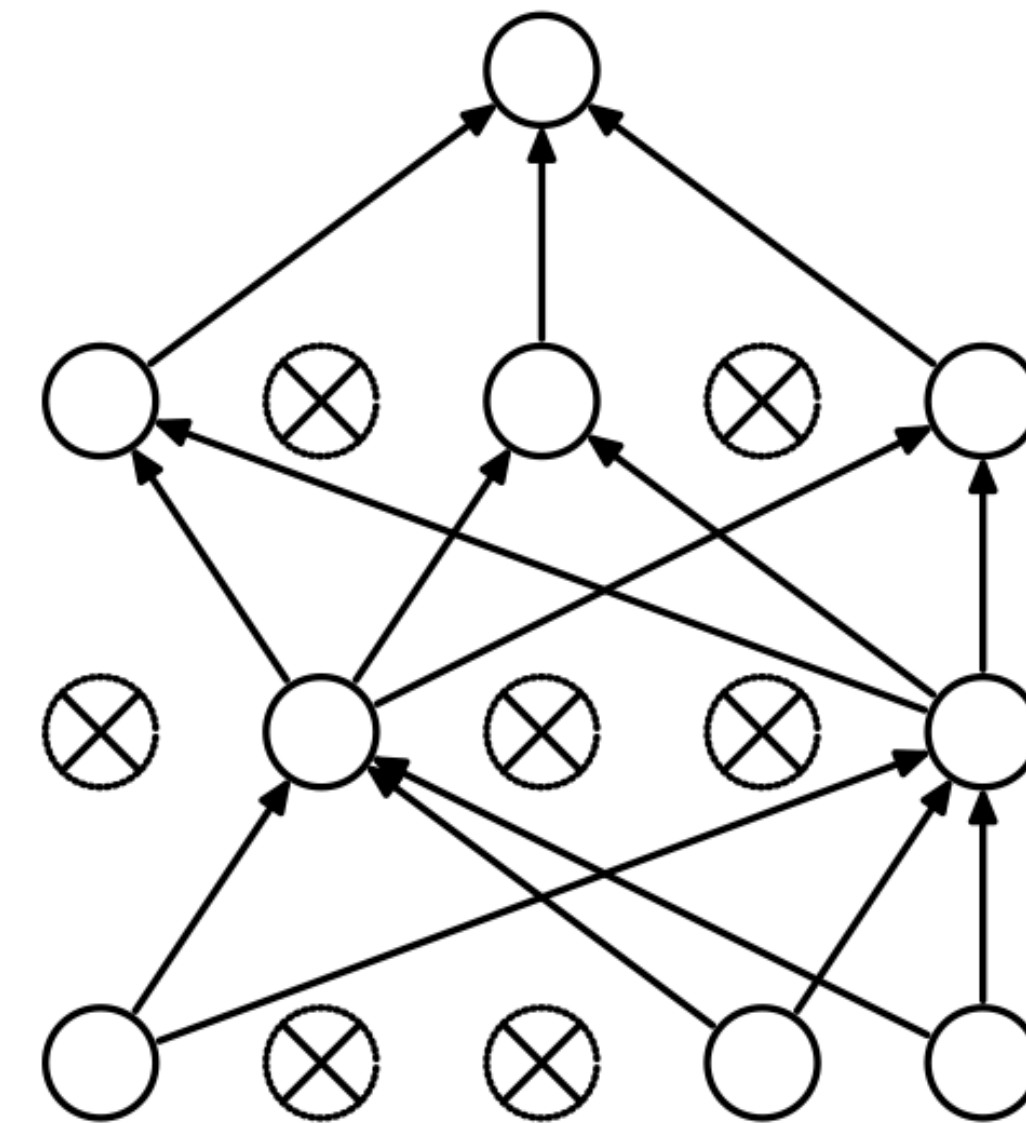
# Defense

---

- Dropout



(a) Standard Neural Net

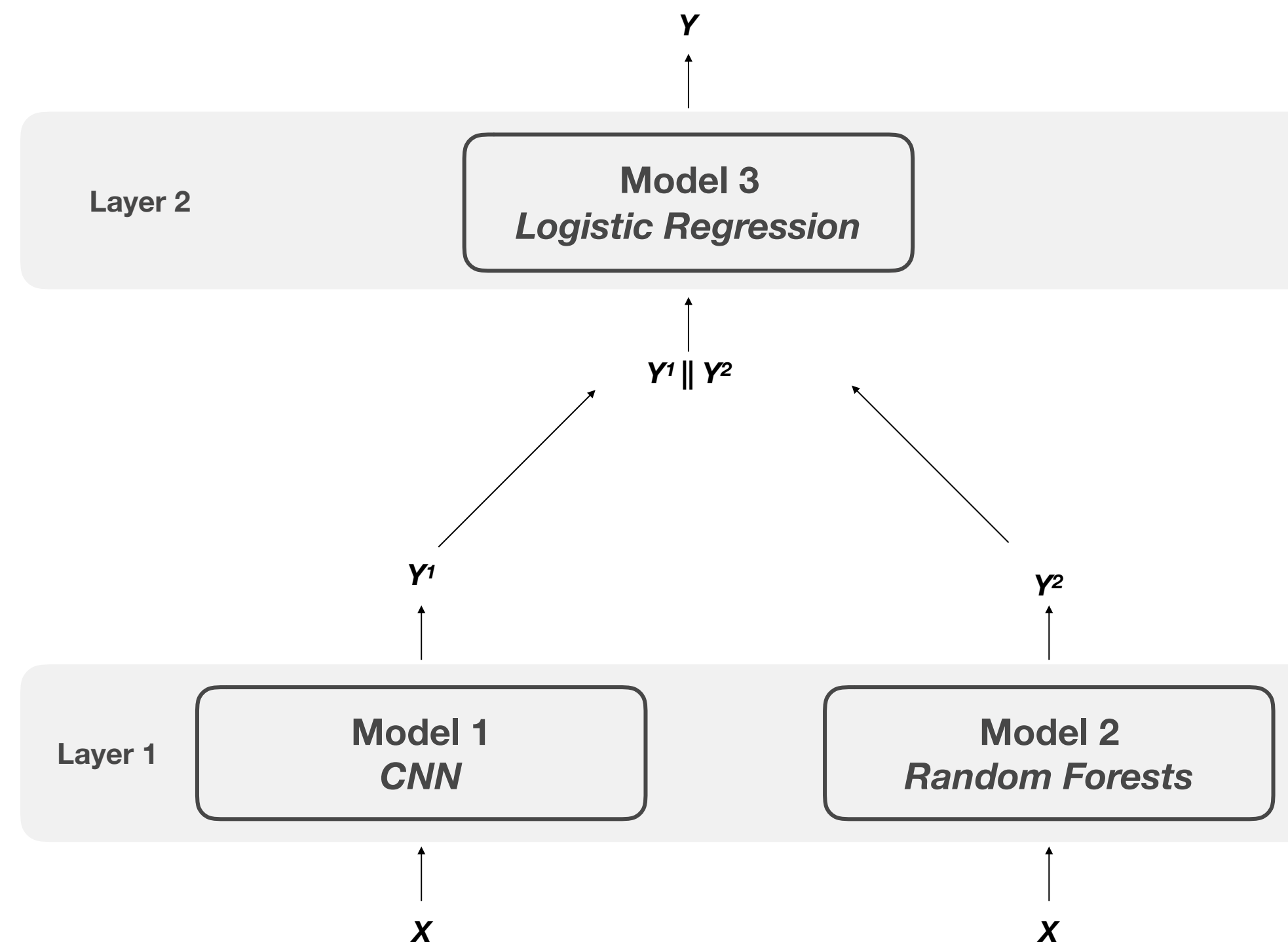


(b) After applying dropout.

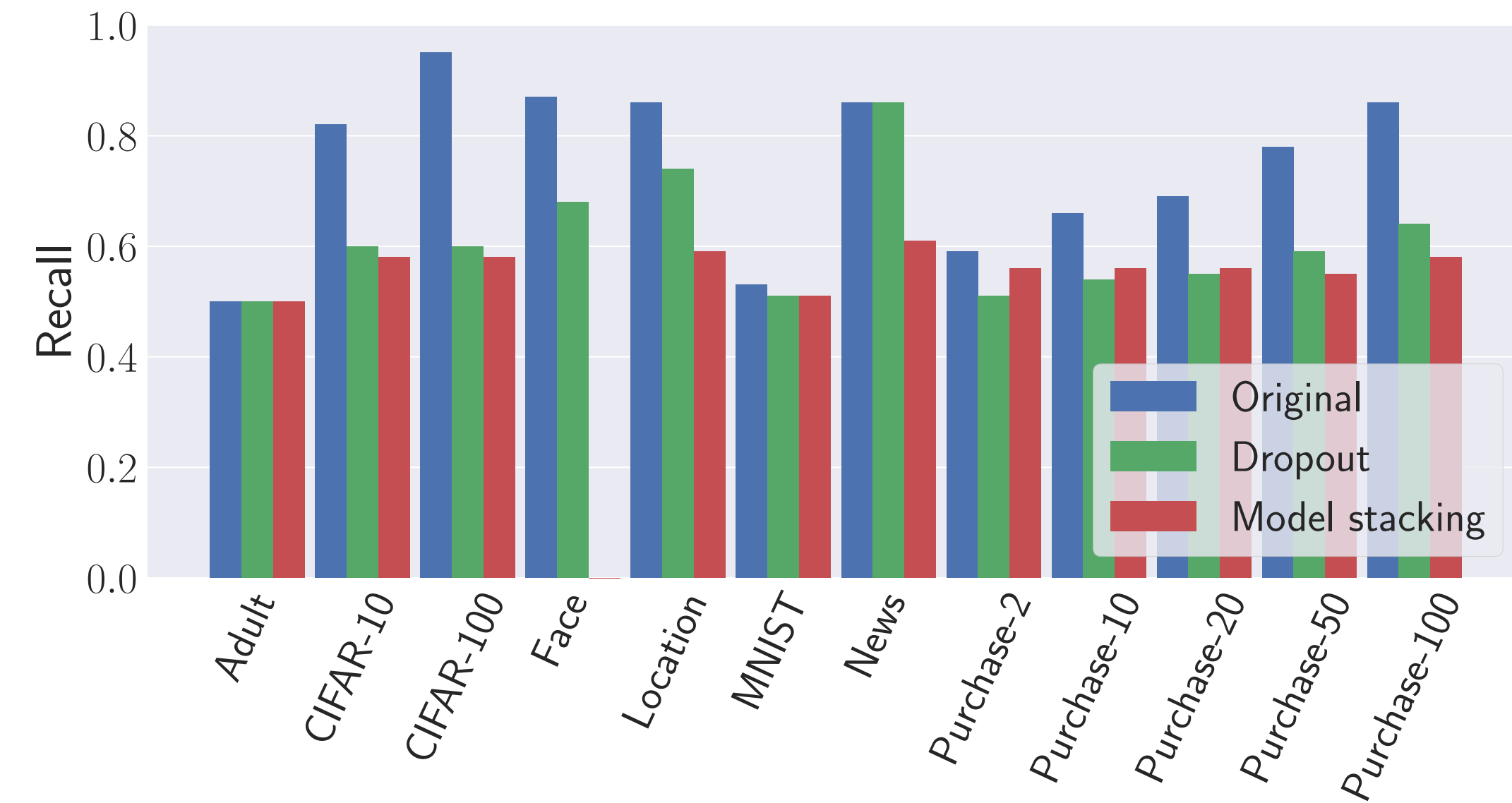
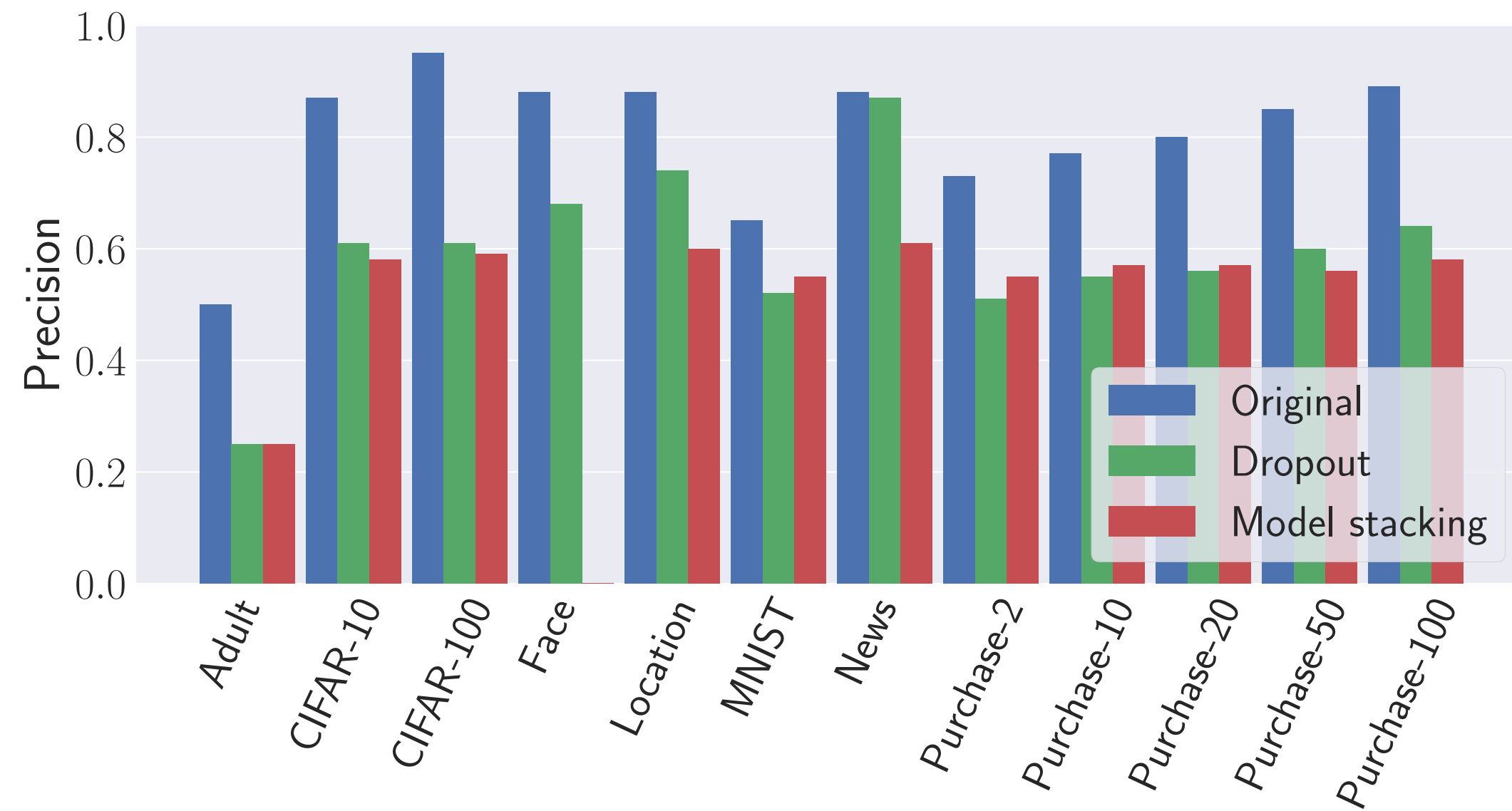
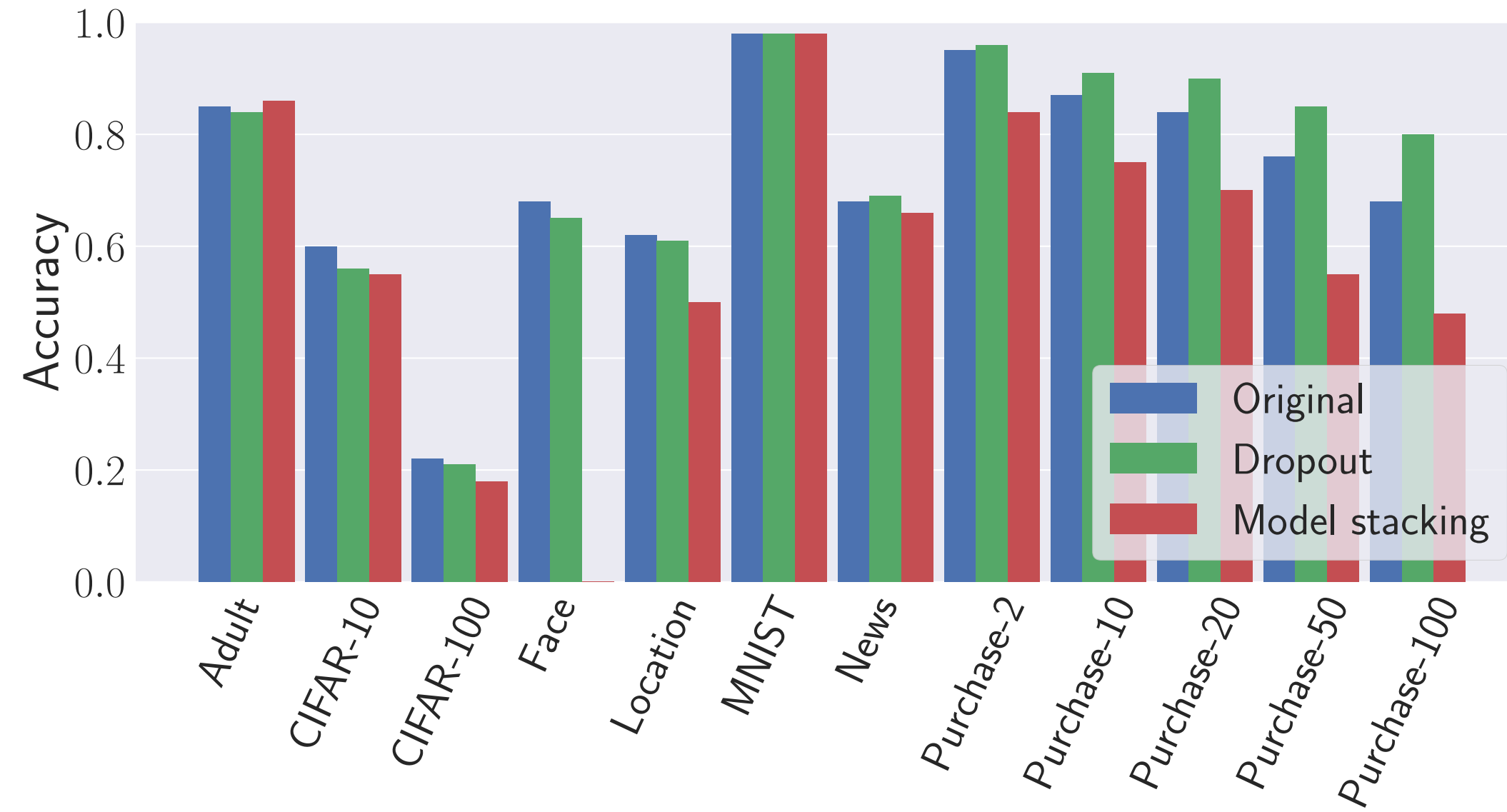
Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting", JMLR 2014

# Defense

- Dropout
- Model Stacking



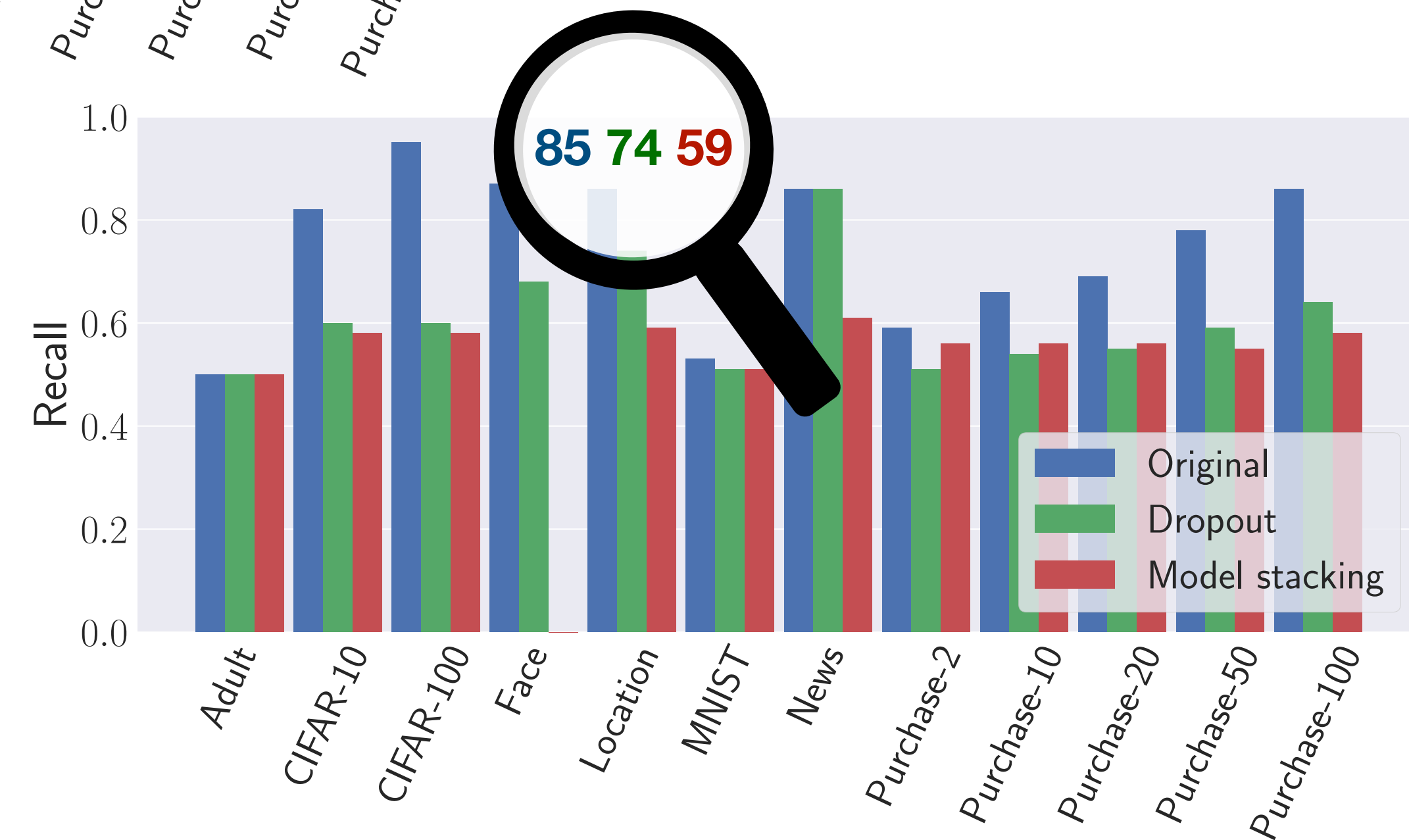
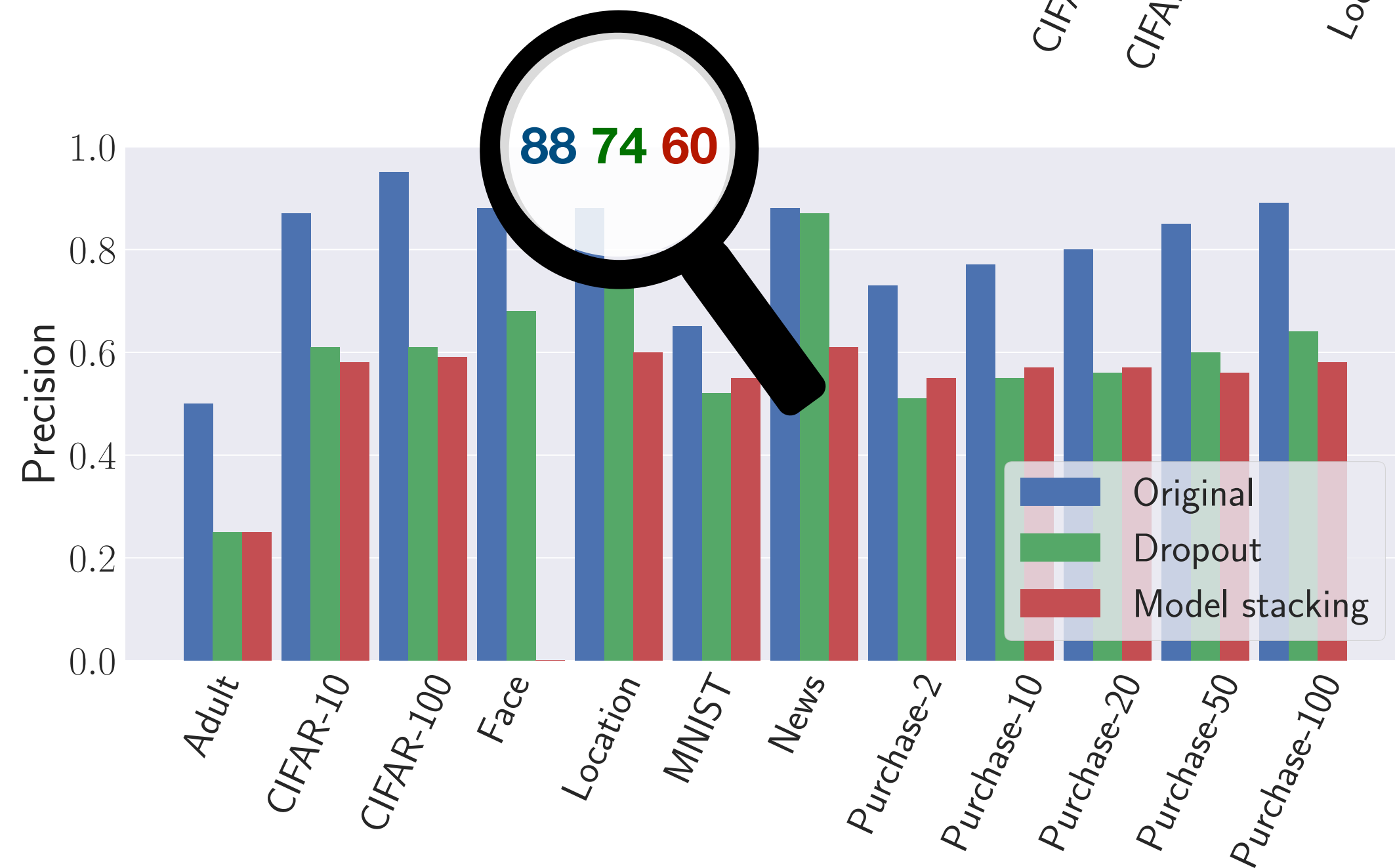
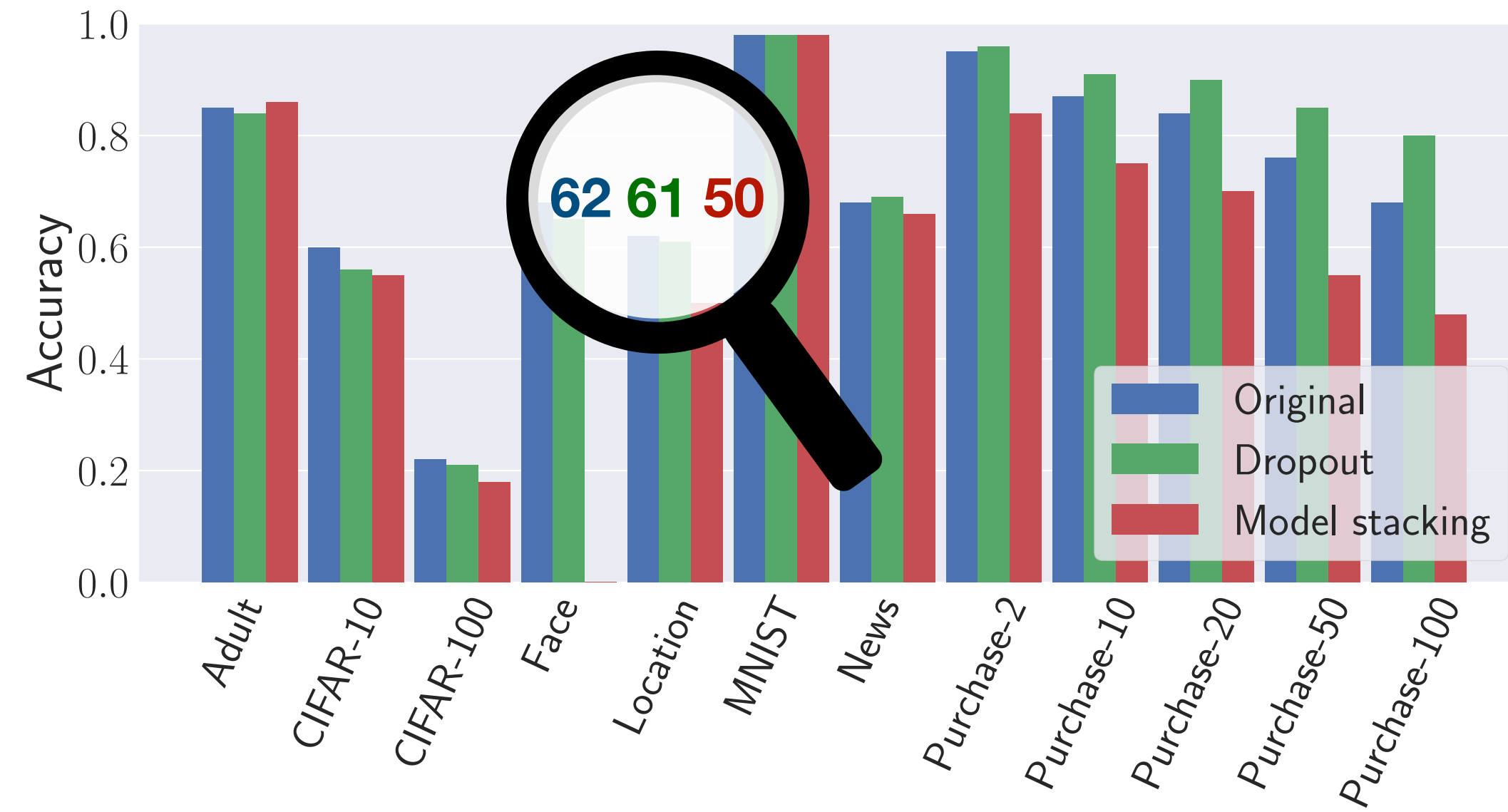
# Performance





# Performance

	Accuracy	Precision	Recall
<b>Dropout</b>	<b>2 %</b>	<b>16 %</b>	<b>13 %</b>
<b>Model stacking</b>	<b>19 %</b>	<b>32 %</b>	<b>31 %</b>



# Conclusion

---

# Conclusion

---

- Membership inference attack simpler

# Conclusion

---

- Membership inference attack simpler
- Overfitting is a common enemy



# Conclusion

---

- Membership inference attack simpler
- Overfitting is a common enemy
- Defenses against membership inference

# Conclusion

---

- Membership inference attack simpler
- Overfitting is a common enemy
- Defenses against membership inference

**Thank you for your attention!  
Questions?**

**ahmed.salem@cispa.saarland**  
**<https://ahmedsalem2.github.io/>**  
**@AhmedGaSalem**