# The Best Defense is a Good Offense: Adversarial Augmentation against Adversarial Attacks (A⁵)
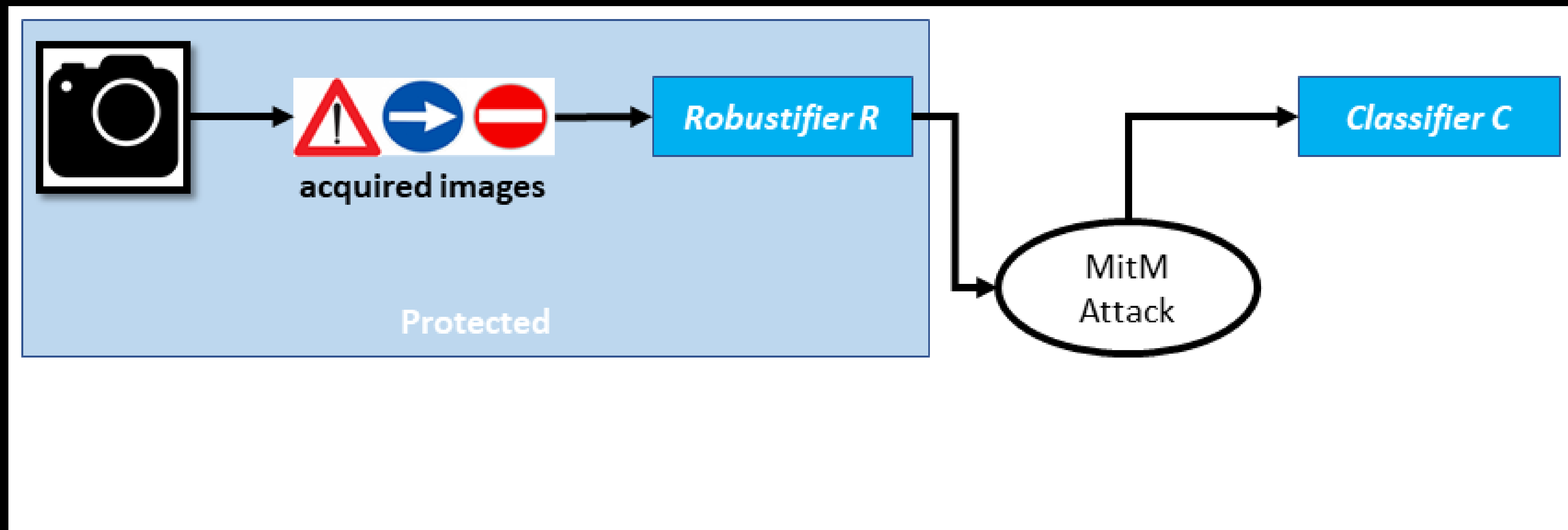
Iuri Frosio, Jan Kautz
NVIDIA

**CVPR 2023**

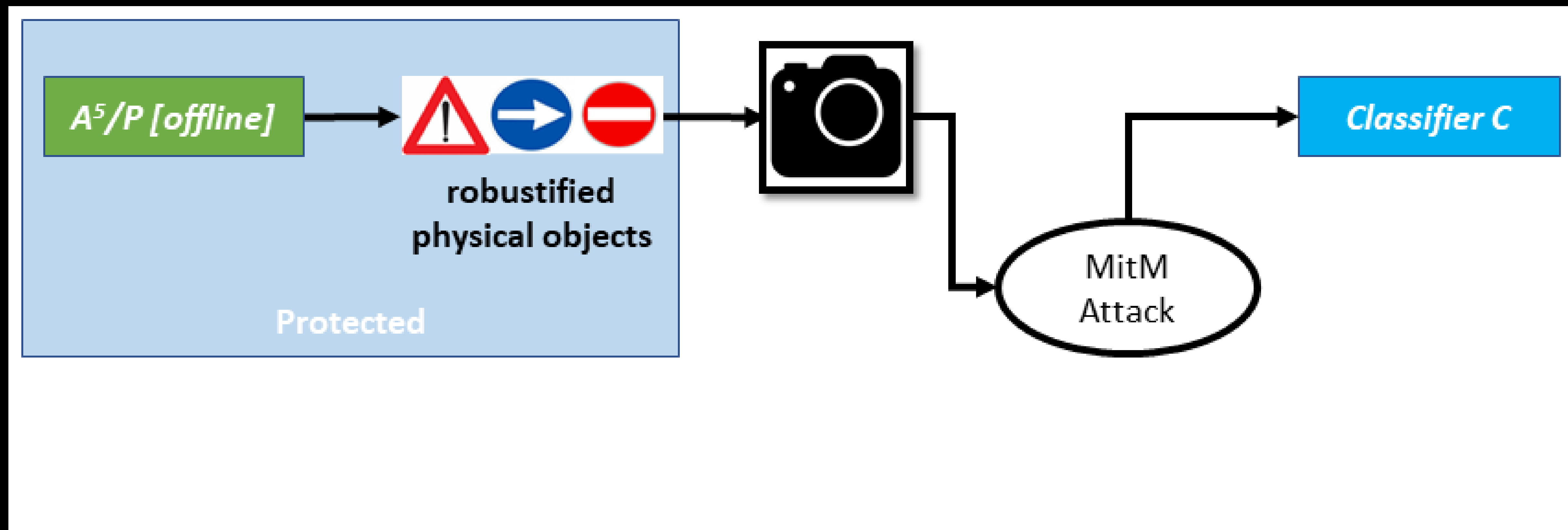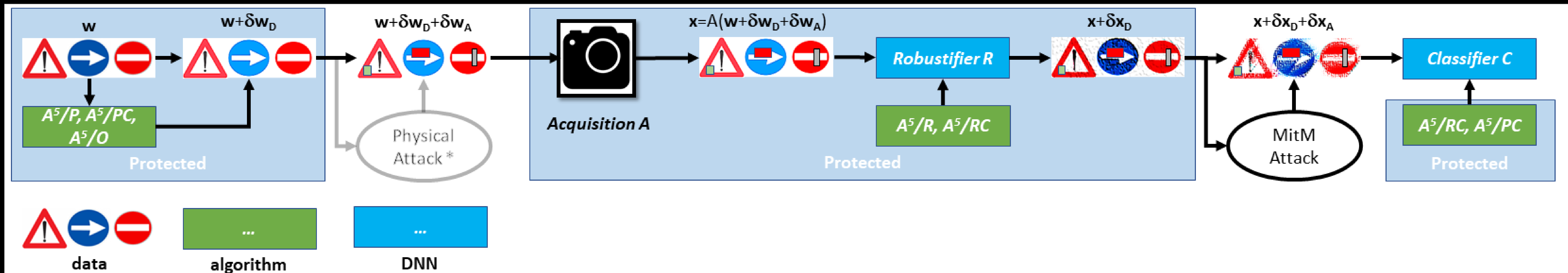**Poster ID 388**
**TUE-AM-388**

# Adversarial Augmentation against Adversarial Attacks (A$^5$)

A full framework for PREEMPTIVE, CERTIFIED protection

# Adversarial Augmentation against Adversarial Attacks (A$^5$)

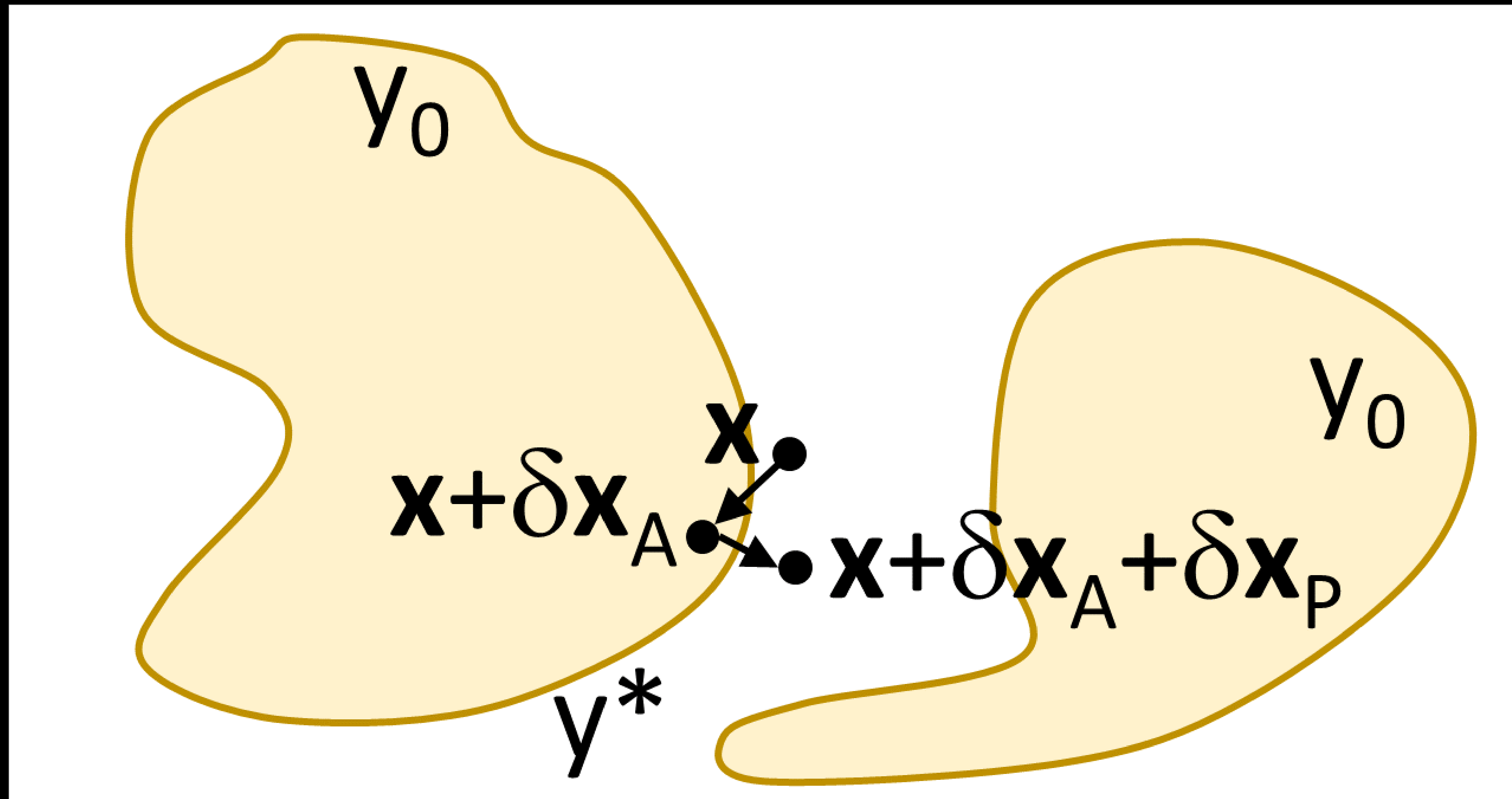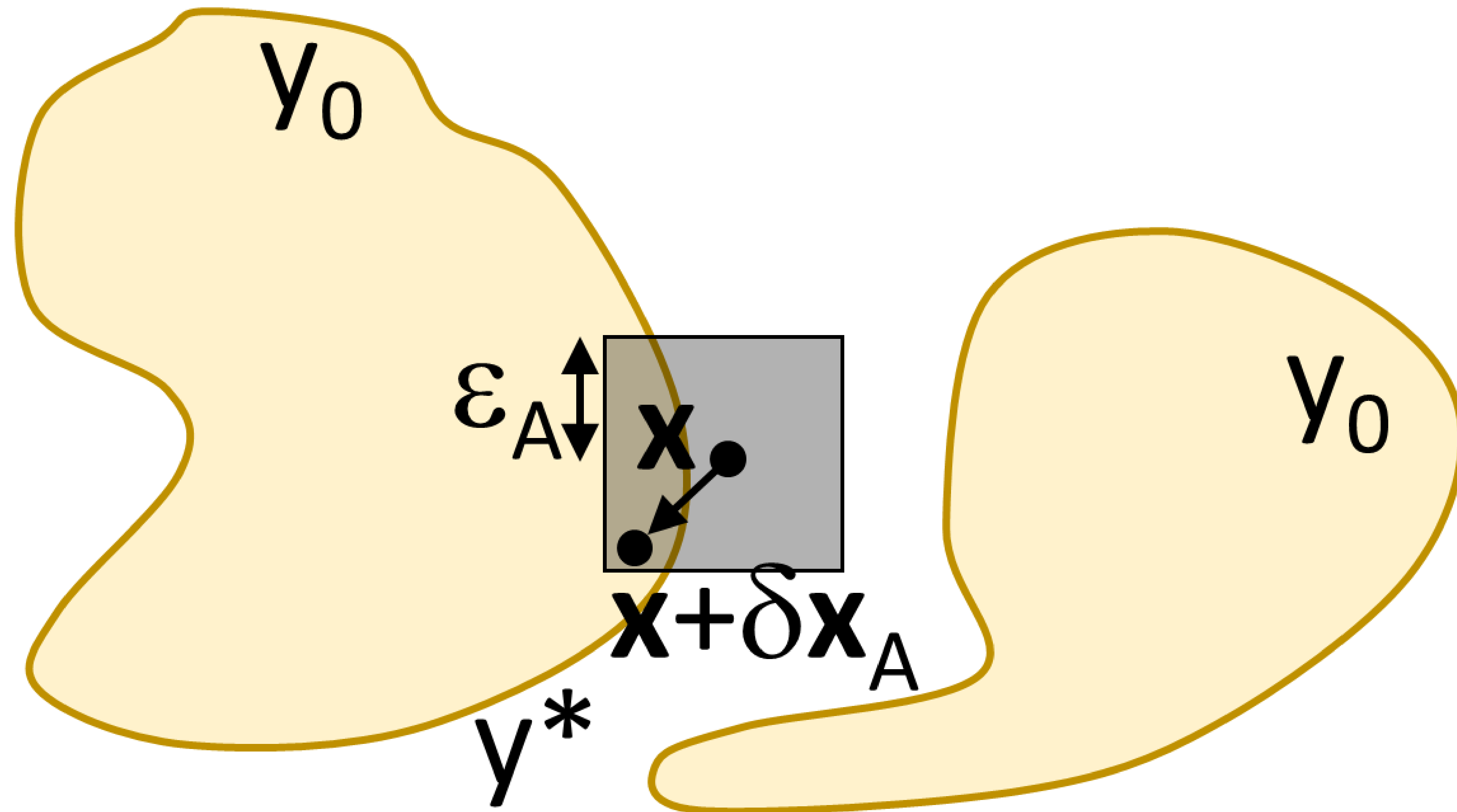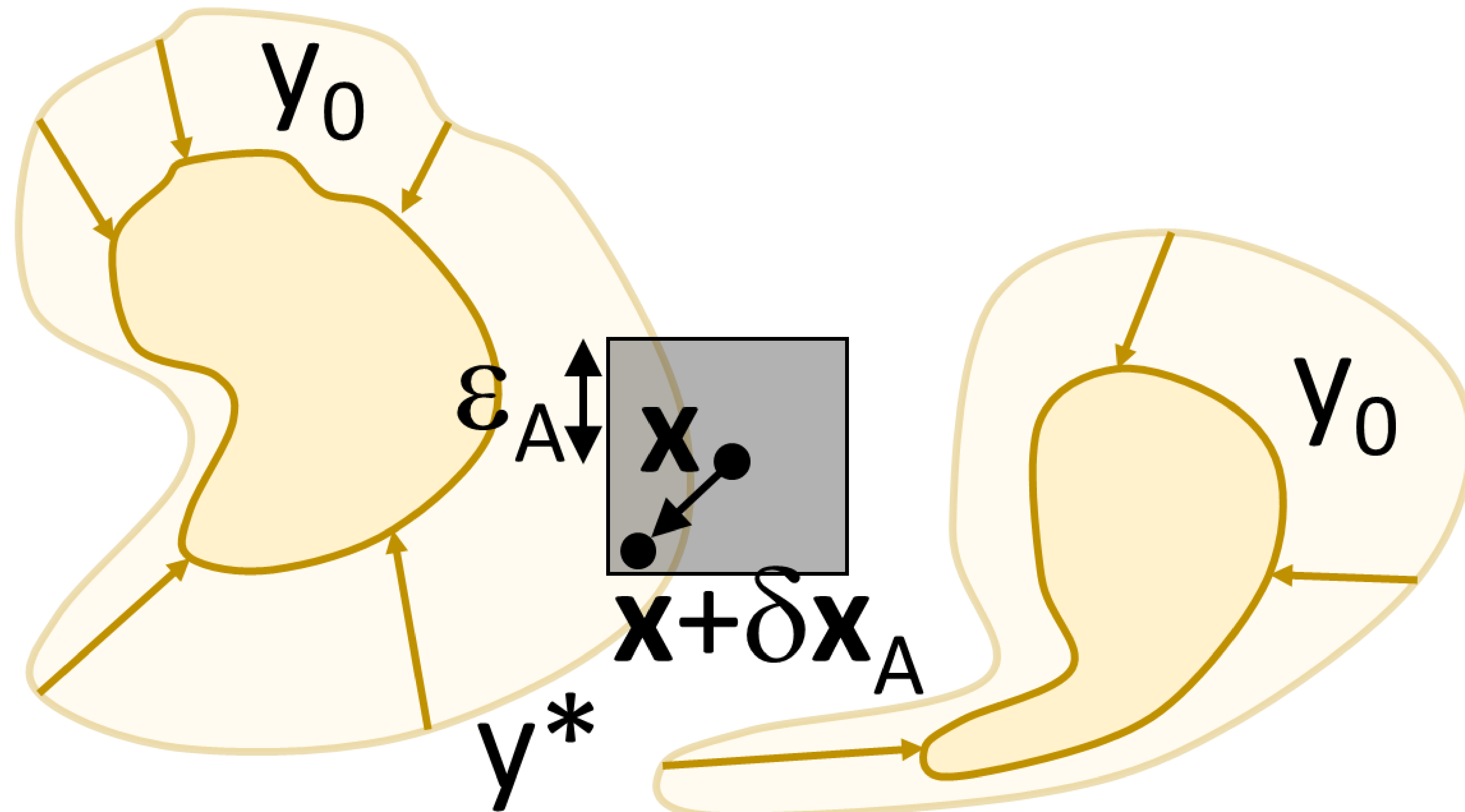A full framework for PREEMPTIVE, CERTIFIED protection

# Adversarial Augmentation against Adversarial Attacks ($A^5$)
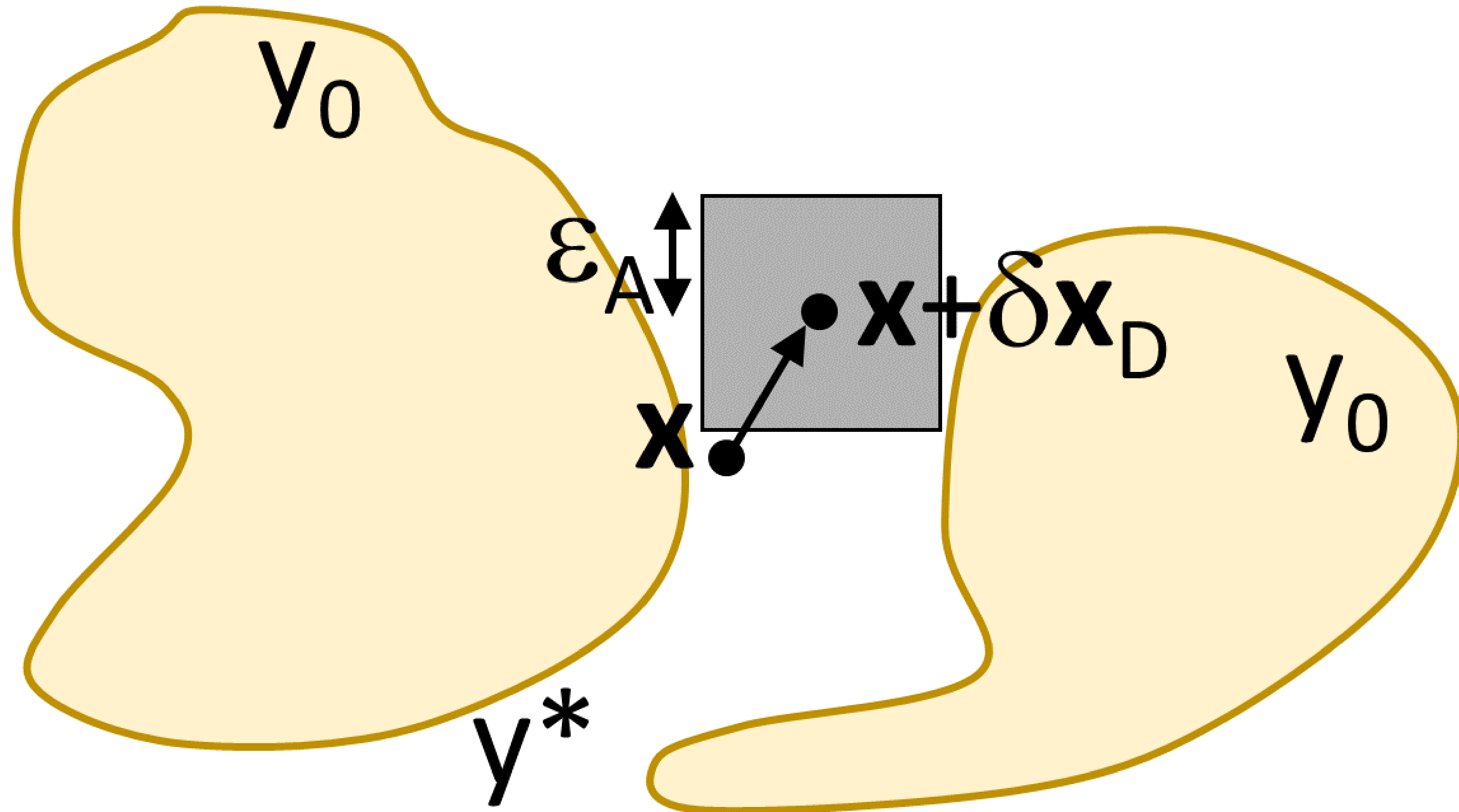
## A full framework for PREEMPTIVE, CERTIFIED protection

# Bound computation

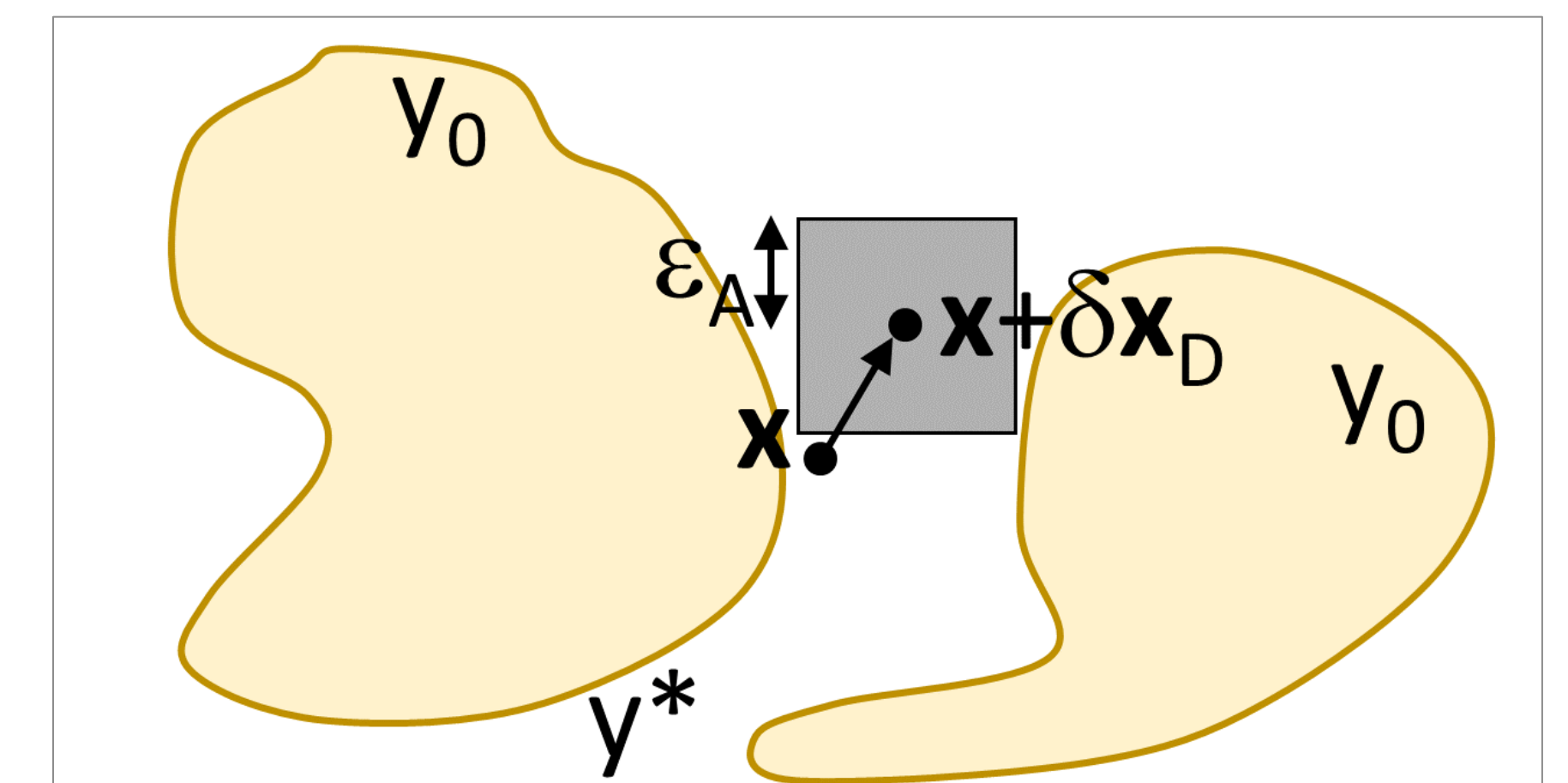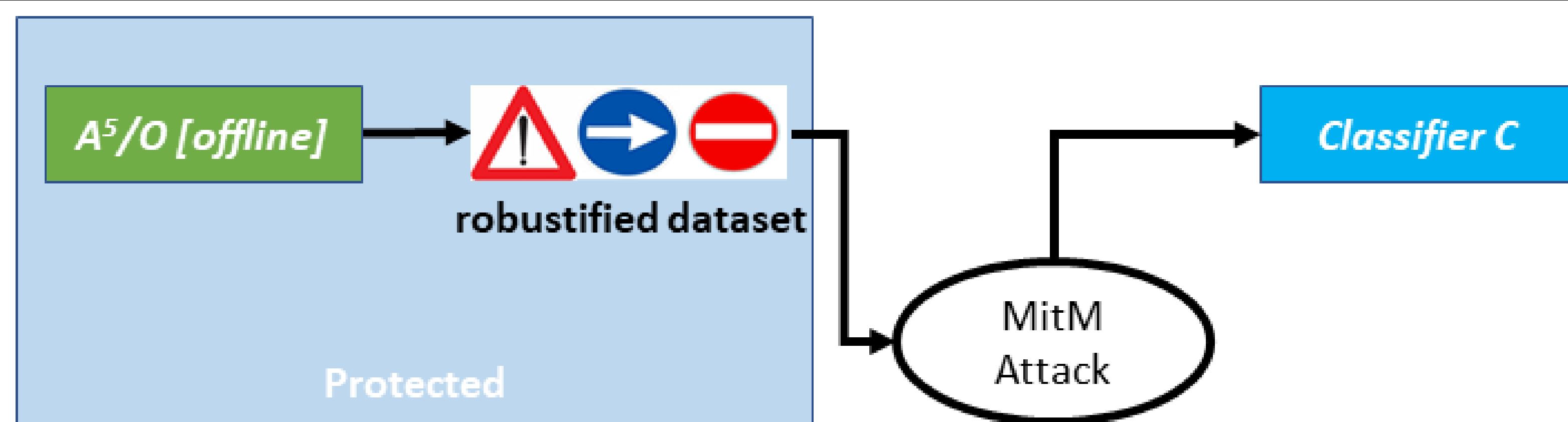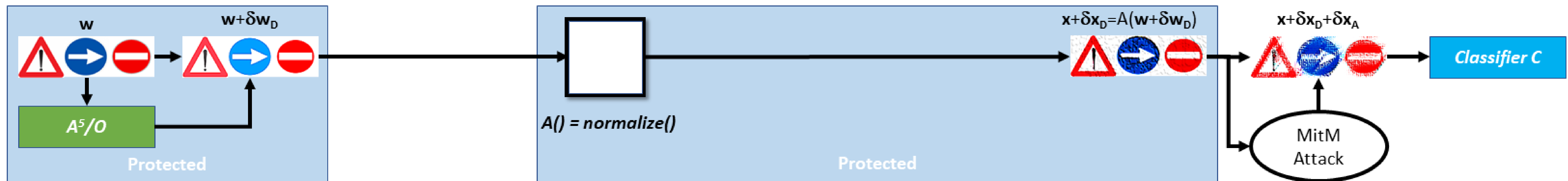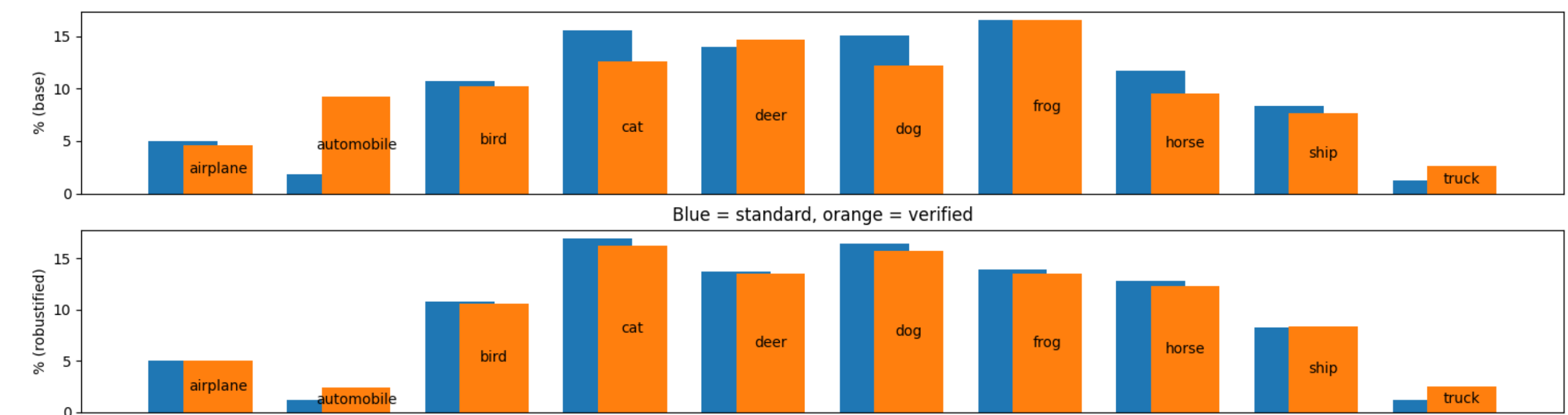# Adversarial and certified training
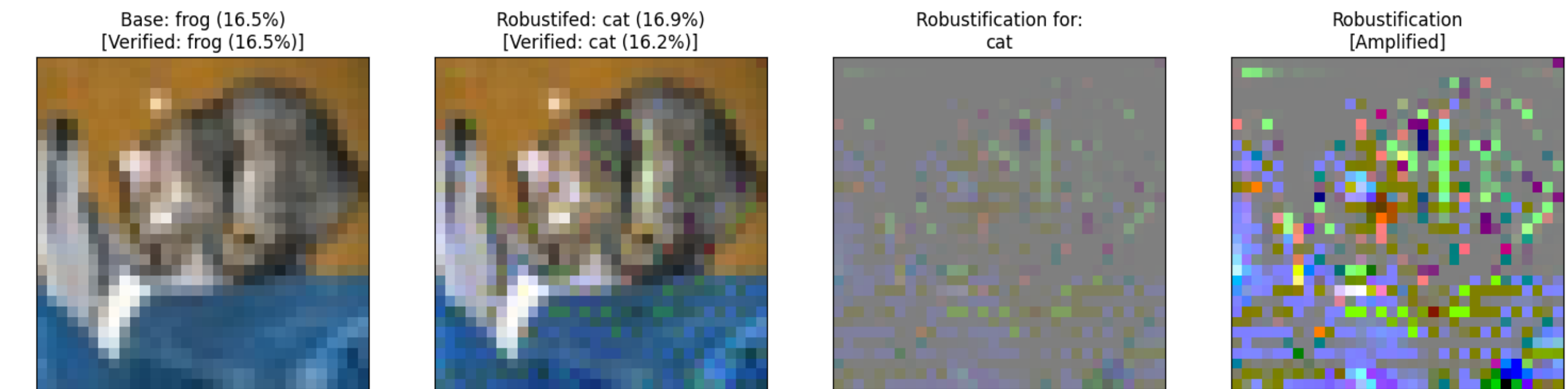
# Configurations: Offline (A⁵/O)

# Configurations: Offline (A$^5$/O)

## Results on CIFAR10, attack 8/255

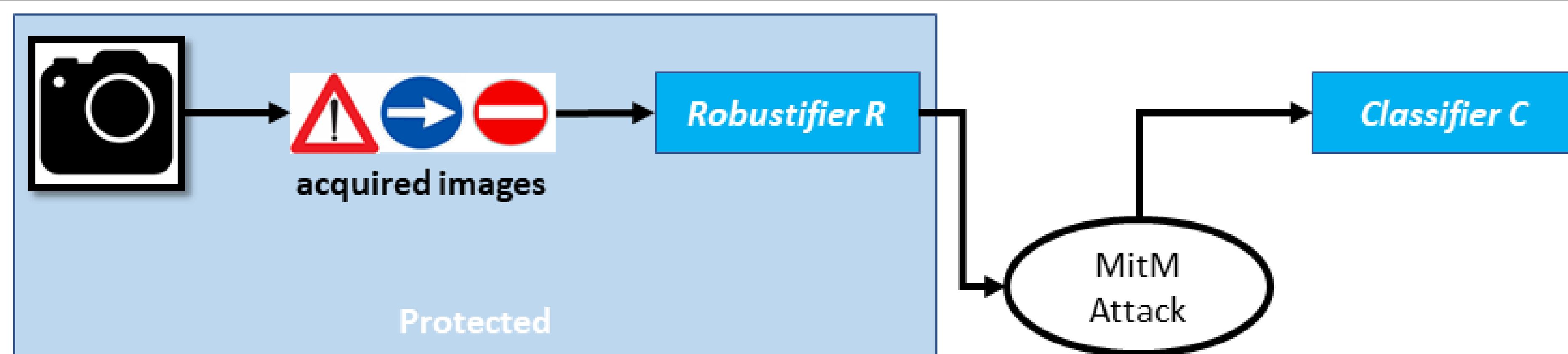| Method | Clean error | Certified error |
|--------|-------------|-----------------|
| CROWN-IBP | 54.02% | 66.94% |
| A$^5$/O | 45.68% | 49.74% |

# Configurations: Robustifier (A⁵/R)

# Configurations: Robustifier ($A^5$/R)

## Results on CIFAR10, attack 8/255

| Method | Clean error | Certified error |
|--------|-------------|-----------------|
| CROWN-IBP | 54.02% | 66.94% |
| $A^5$/O | 45.68% | 49.74% |
| $A^5$/R | 50.91% | 57.95% |

# Configurations: Robustifier and Classifier (A$^5$/RC)

## Results on CIFAR10, attack 8/255

| Method | Clean error | Certified error |
|--------|-------------|-----------------|
| CROWN-IBP | 54.02% | 66.94% |
| A$^5$/O | 45.68% | 49.74% |
| A$^5$/R | 50.91% | 57.95% |
| A$^5$/RC | 35.26% | 42.76% |

# Configurations: Physical and Classifier (A⁵/PC)

## Results on OCR

| Method | Clean error | Certified error |
|--------|-------------|-----------------|
| Vanilla | 0.89% | 100.00% |
| CROWN-IBP | 3.85% | 13.85% |
| $A^5$/P | 3.08% | 11.84% |
| $A^5$/PC | 0.73% | 4.20$ |



Augmented

Original

Defensive augmentation

After imaging with camera (simulated)

# Conclusion

- Preemptive, certifiable robustification

- Offline and on the fly for acquired images

- Offline for physical objects

- The benefit of co-training

- Technical details and more results on MNIST, FashionMNIST, TinyImage net available in our paper

- Scaling to large network architecture still problematic

- Code available for research here: