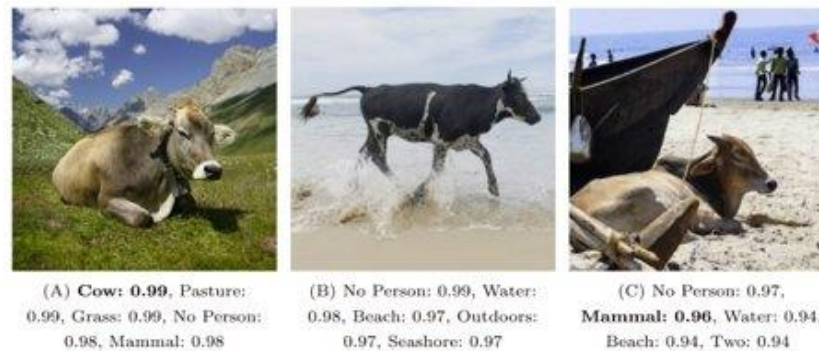


# The Role of Causality in Developing Robust AI Systems

A presentation by Dren Fazlija

# Robustness and Privacy

- Robustness: Decreasing sensitivity towards input changes



Source: [Beery2018]

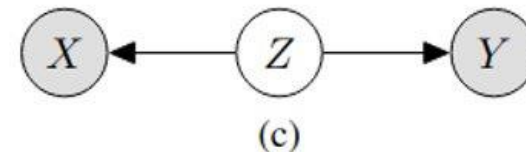
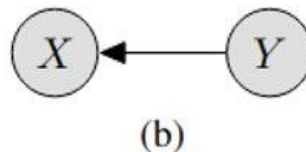
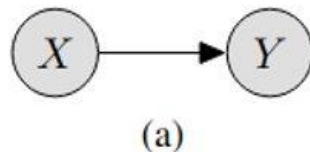
- Privacy: Defending against privacy-evasive attacks
- Causal solutions for both areas overlap significantly
  - Robustness: Methods for centralized learning setting
  - Privacy: Similar methods for decentralized/federated learning setting

# Statistical Machine Learning

- We assume that our data is independent and identically distributed (IID)
- Allows one to infer the performance of models solely through training data
  - Empirical Risk Minimization
- Very unlikely that training data covers all statistical properties of real-world inference data
- Susceptible to distributional shifts caused by unseen data

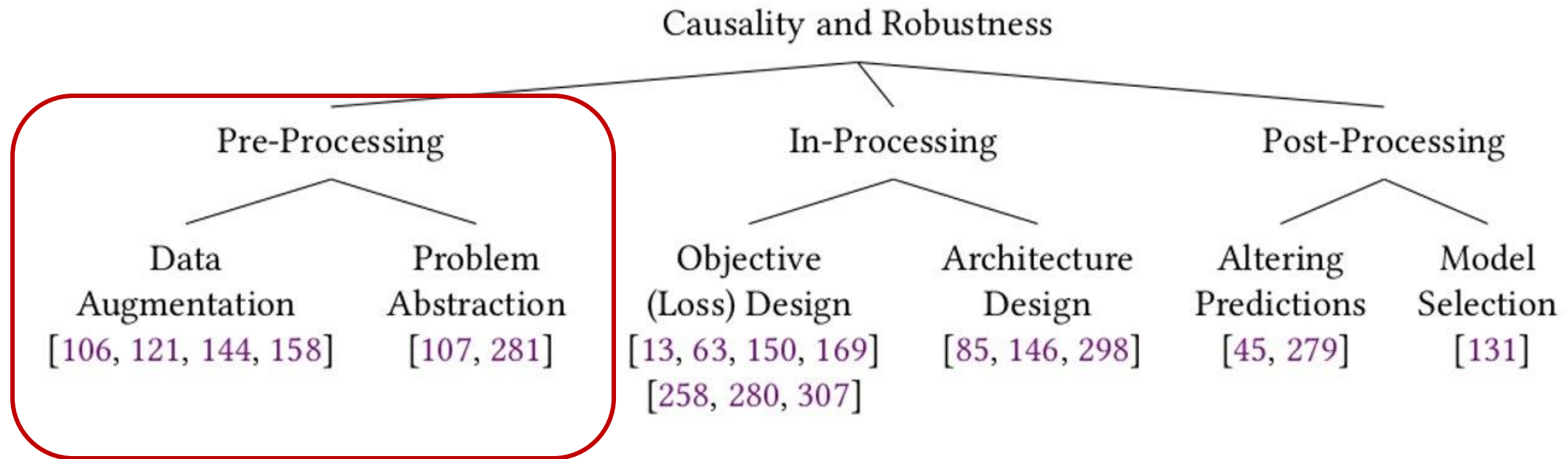
# Enhancing AI with Causality

- No definite solution for distributional shifts
- Statistical ML models are not inclined to properly understand causal relationship
  - Simply fall back on observable correlation that works best for the training data
- Causal encodings allow us to constraint this behavior
- Achievable with pre-, in- and post-processing methods



Source: [Schölkopf22]

# Overview of Causal Solutions

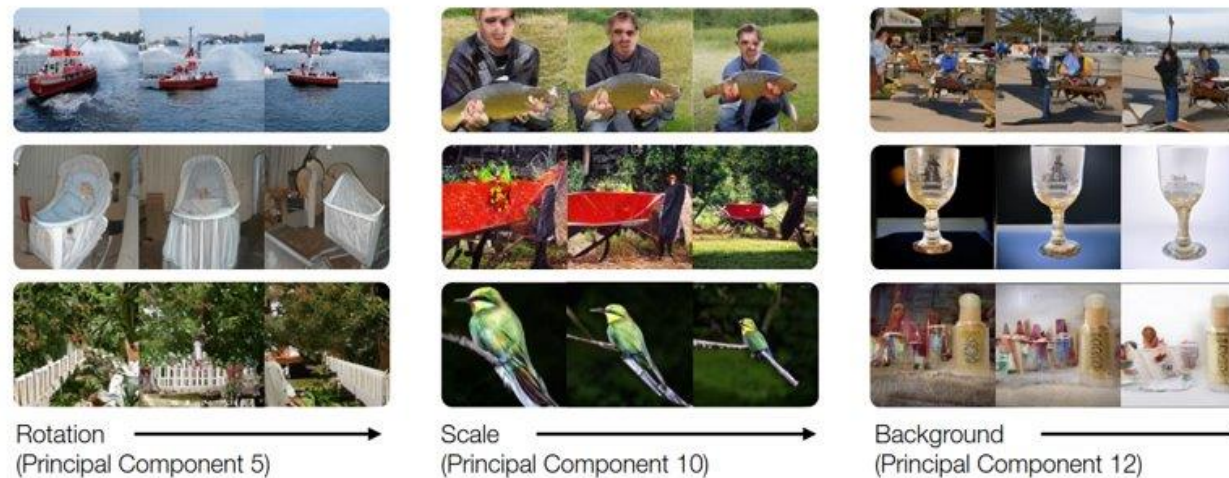


Source: [Ganguly23]

# Generative Interventions for Causal Learning

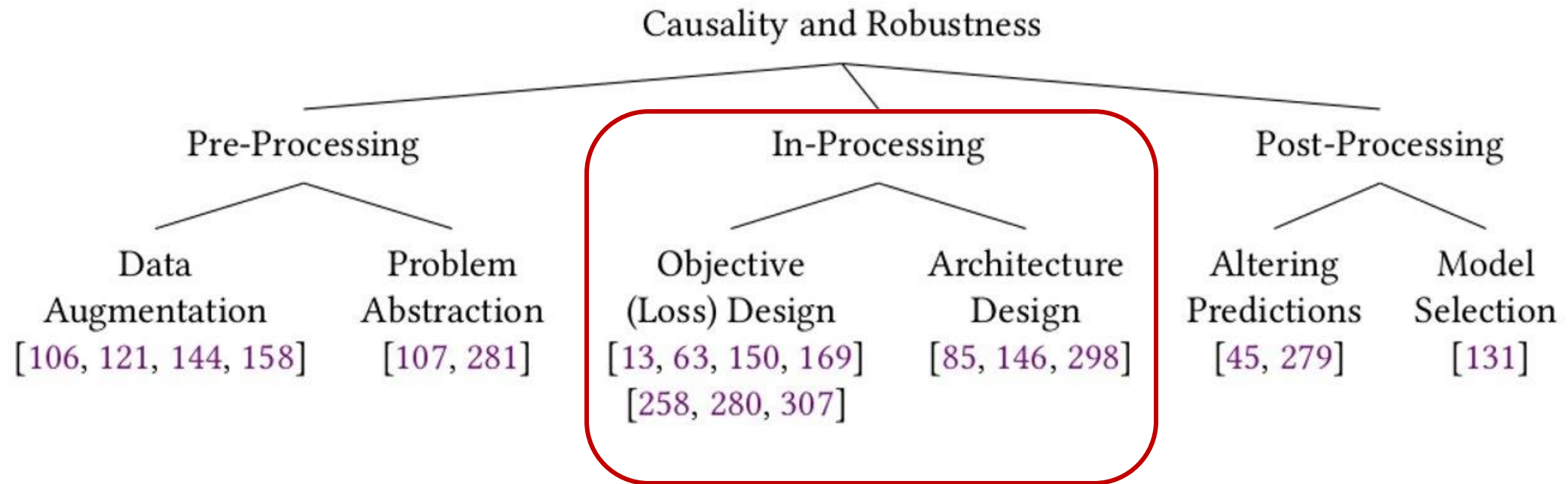
Pre-Processing Method for Robustness [Mao2021]

- Goal: Provide training data whose observable correlations better reflect causal relationships
- Simulate interventions on nuisance factors via GANs



Source: [Mao2021]

# Overview of Causal Solutions



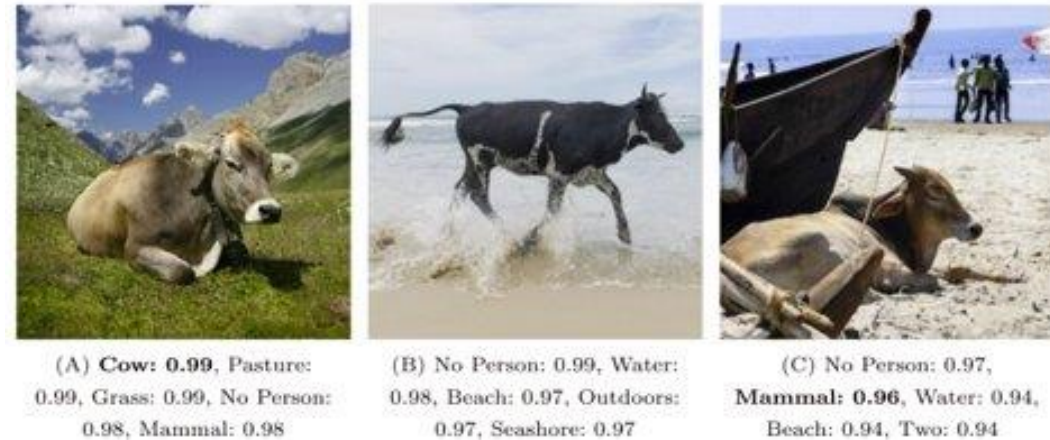
Source: [Ganguly23]



# Invariant Risk Minimization

In-Processing Method for Robustness [Arjovsky2019]

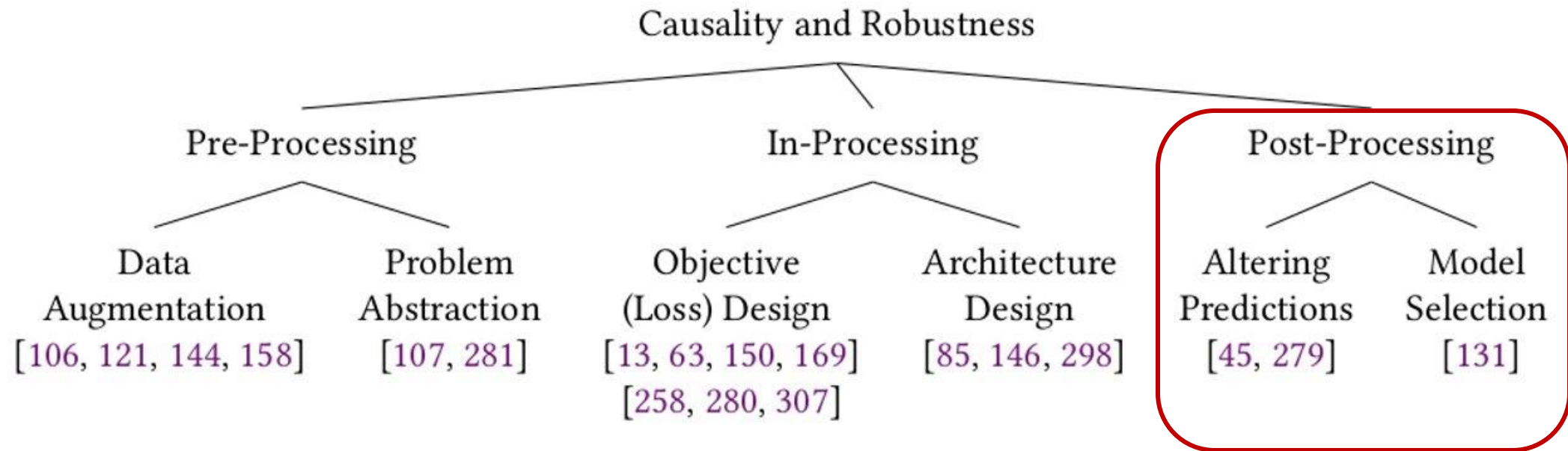
- Feature invariance relates to its causal importance
  - E.g., image background can greatly vary across data points
  - Therefore, it is not important for predicting the label
- Allows one to develop causal models without causal encodings
- Idea: Promote consistent behavior across different environments
- Successful at increasing robustness of image classifiers in the OOD setting



Source: [Beery2018]



# Overview of Causal Solutions

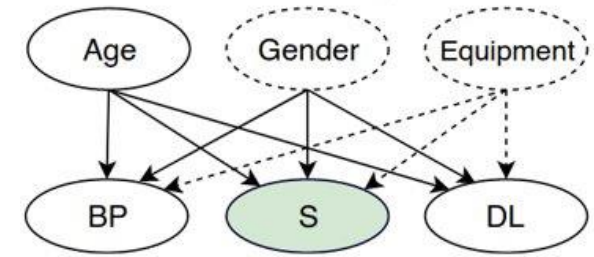


Source: [Ganguly23]

# Causal Model Selection

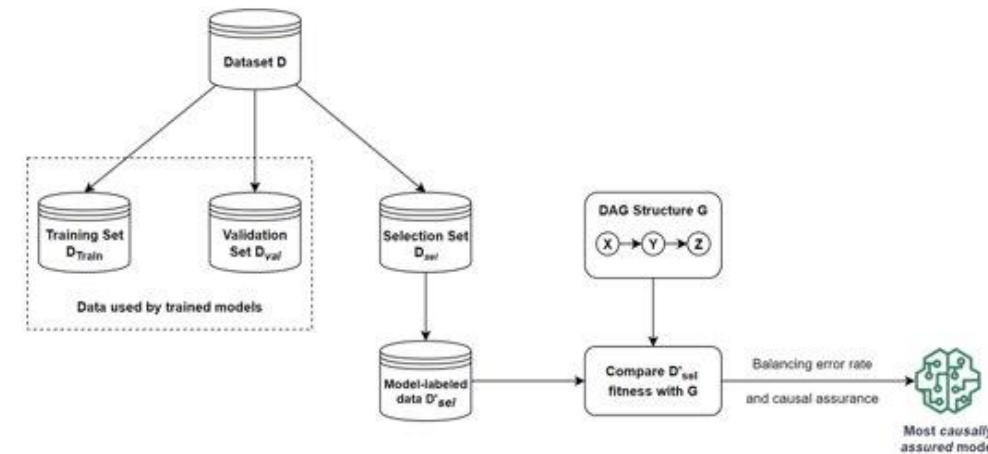
Post-Processing Method for Robustness [Kyono2019]

- Builds on causal invariance assumption
- Goal: pick trained model that best reflects the causality intrinsic to the domain
- Let each model overwrite the labels of the data points with their own prediction
- Resulting dataset should then reflect the very same relationships
- Successfully picks robust models for real-world tabular datasets



(d) Powerlifting dataset

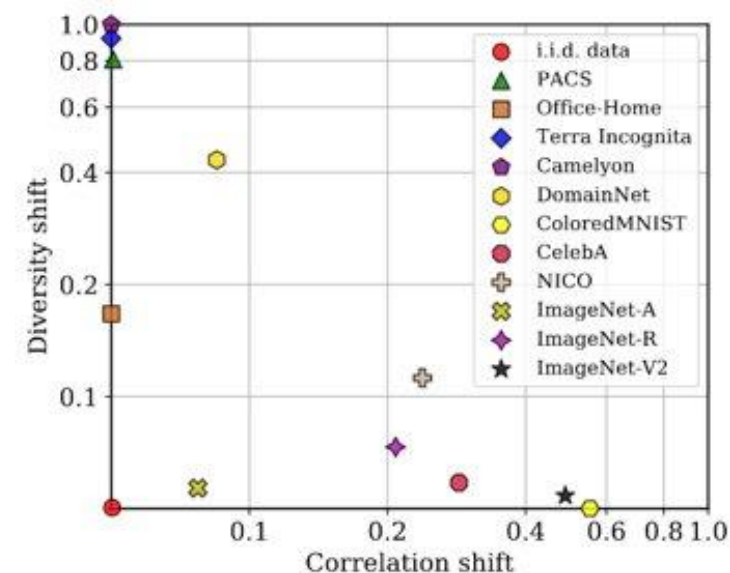
Source: [Kyono2019]



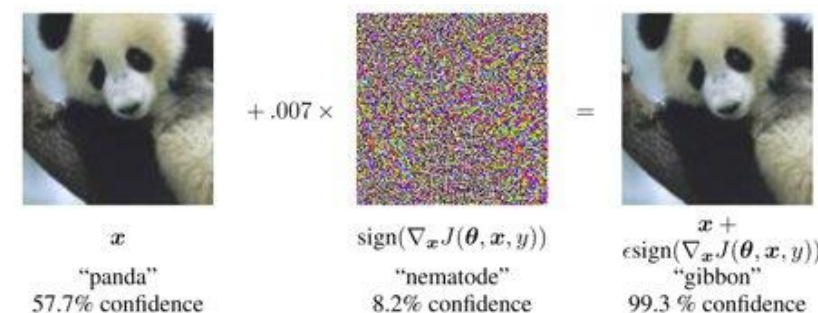
Inspired by [Kyono2019]

# Future Work

- Categorize OOD learning abilities of causal solutions
- Explore related fields of Causal ML
  - E.g., Neurosymbolic AI or Object-Centric Learning
- Further explore Adversarial Machine Learning
  - Most solutions are designed for natural OOD data
  - Interesting subarea: Certified robustness



Source: [Ye2022]



Source: [Goodfellow2014]