

# Toward Explainable Uncertainty

Alan Fern and Tom Dietterich

School of Electrical Engineering and Computer Science  
Oregon State University

# Goal 1: Uncertainty Aware ML Systems

- Design machine learning systems that “know what they know”  
[Li, Littman, Walsh ICML’08]
  - Provide guarantees on predictions
  - Allow systems to abstain and/or produce ambiguous predictions
- Achieve this in:
  - Closed Worlds = Known Unknowns
  - Open Worlds = Unknown Unknowns
- Why?
  - Safe and Trustworthy AI
  - End User Acceptability
  - Computational Efficiency – use more complex model if simpler model is uncertain

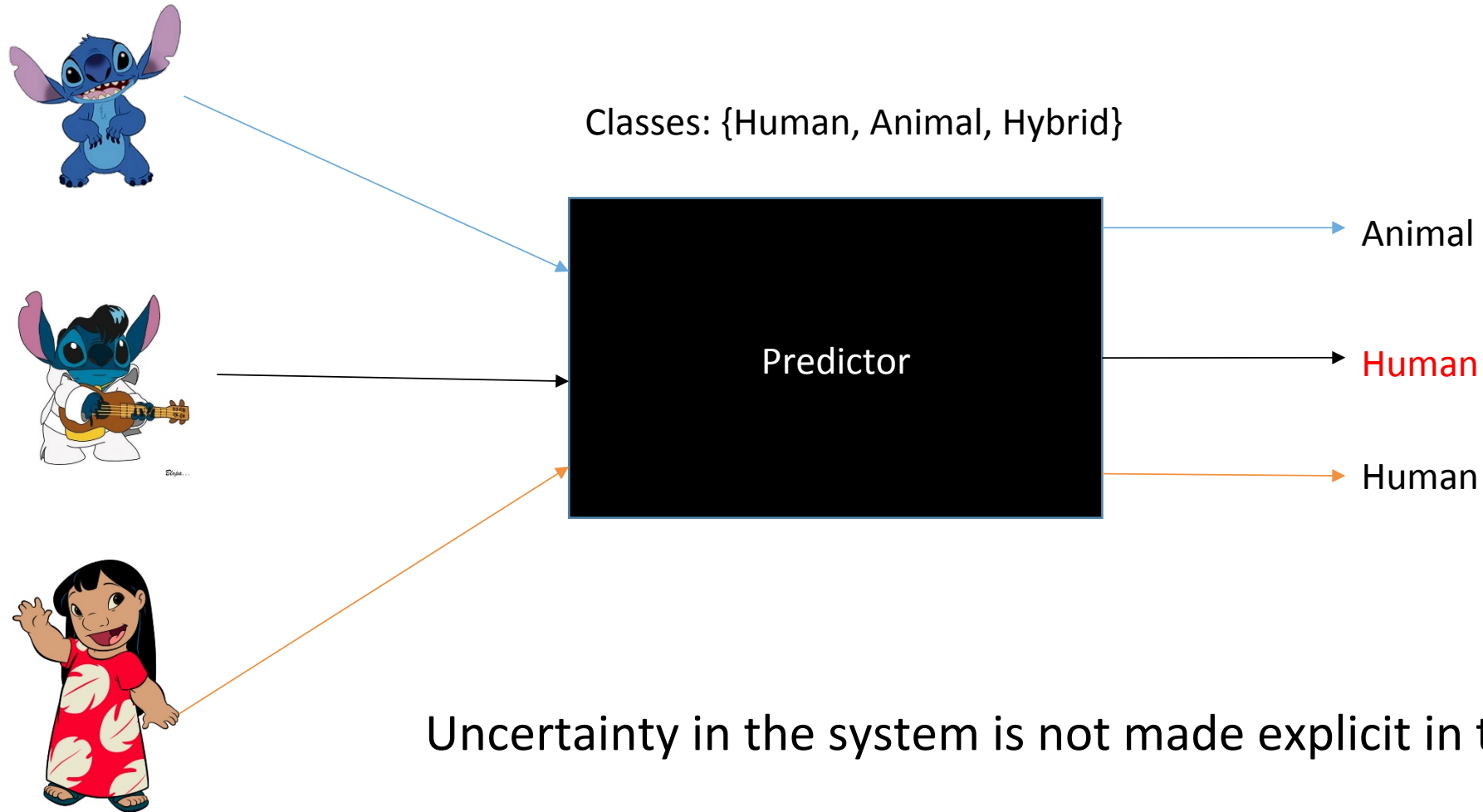
# Goal 2: Transparent Uncertainty in ML Systems

- Design machine learning systems that can “explain their uncertainty”
  - Give insight into why they abstained or produced ambiguous answer
- Achieve these goals in:
  - Closed Worlds = Known Unknowns
  - Open Worlds = Unknown Unknowns
- Why?
  - Basis for feedback to learning systems
  - Basis for investigating anomalies
  - Mechanism for building trust

# Outline

- Conformal Prediction for Uncertainty Aware Classification
  - Empirical performance in closed worlds
  - Empirical performance in open worlds
  - **Not effective in open worlds** → **Suggests integrating with anomaly detection**
  
- Explanations for Anomaly Detection
  - What is an anomaly explanation?
  - How to compute explanations?
  - How to evaluate explanations?

# Standard Classification

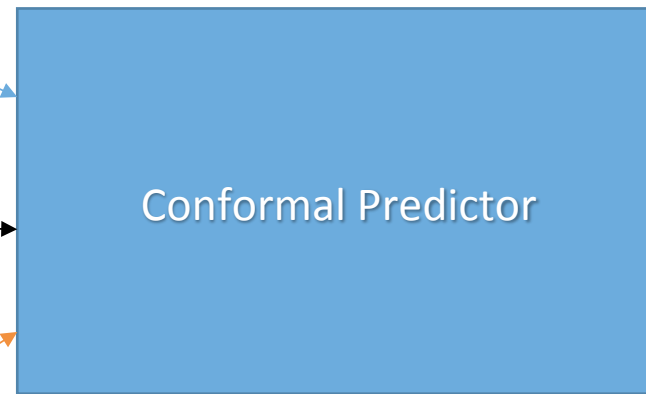
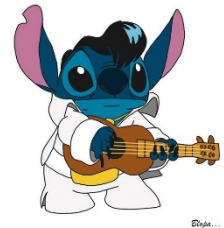


# Conformal Prediction [Vovk et al., 2005]

Most basic type of explanation of uncertainty



Classes: {Human, Animal, Hybrid}



{**Animal**, Hybrid}

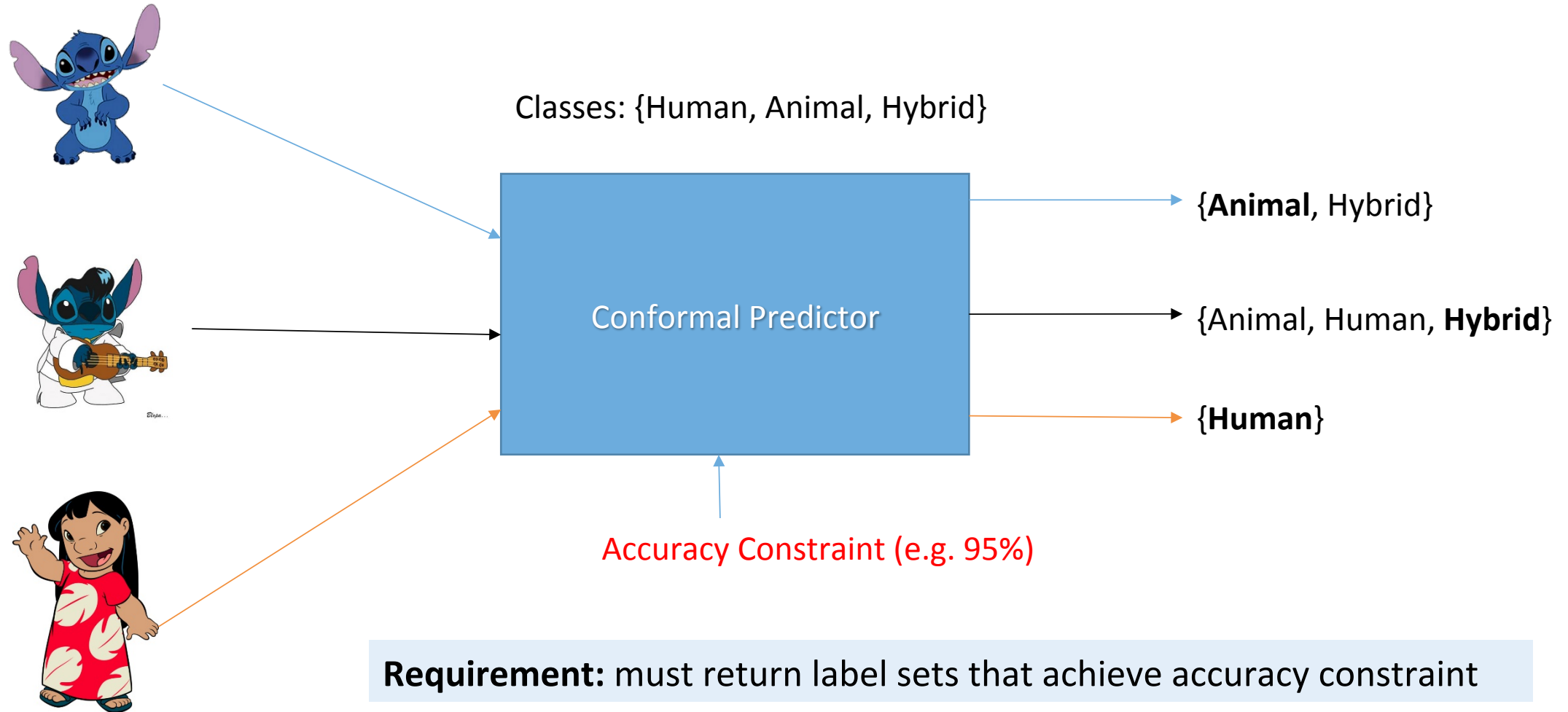
{Animal, Human, **Hybrid**}

{**Human**}



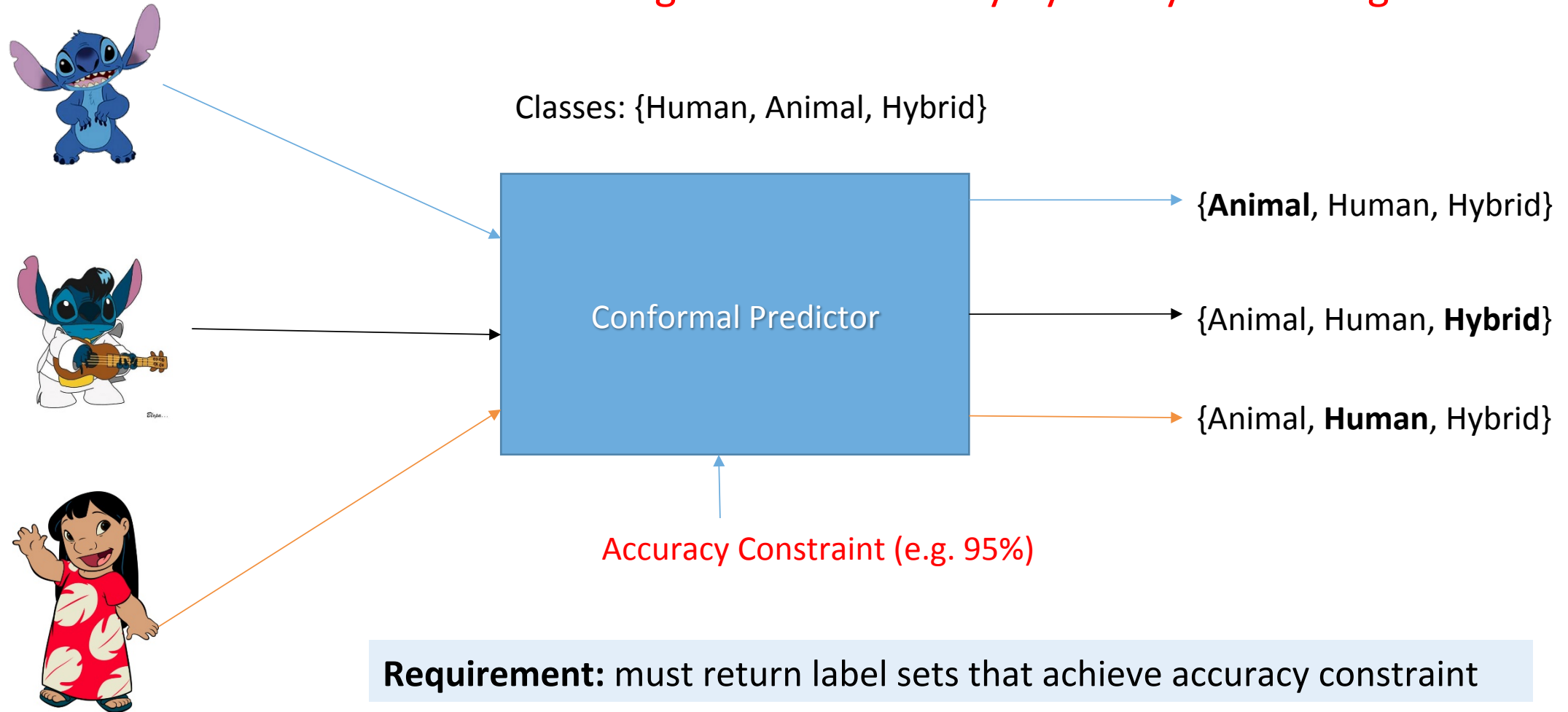
- Conformal predictors output **sets of labels**.
- Label set is correct if it contains true label.

# Conformal Prediction: Accuracy



# Conformal Prediction: Accuracy

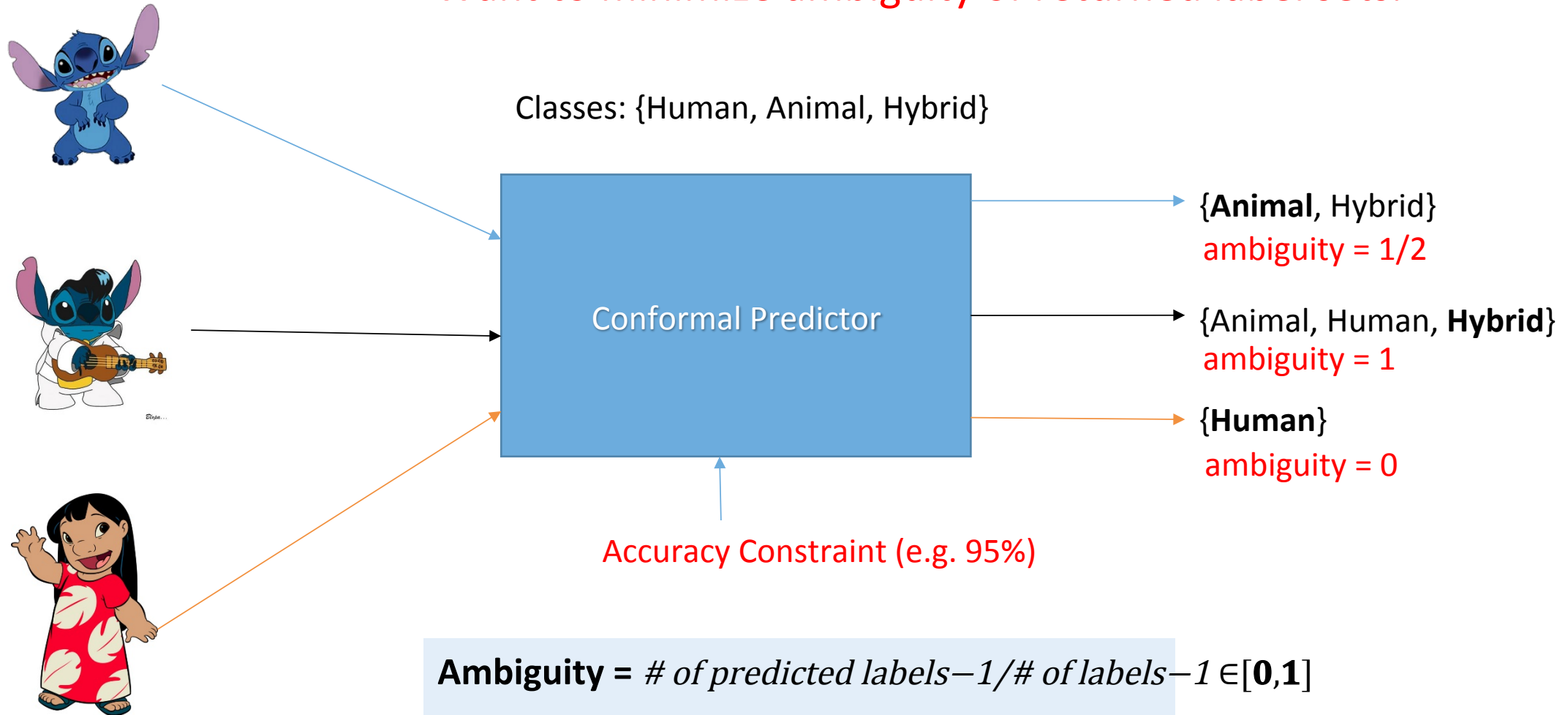
But we can get 100% accuracy by always returning all labels.





# Conformal Prediction: Ambiguity

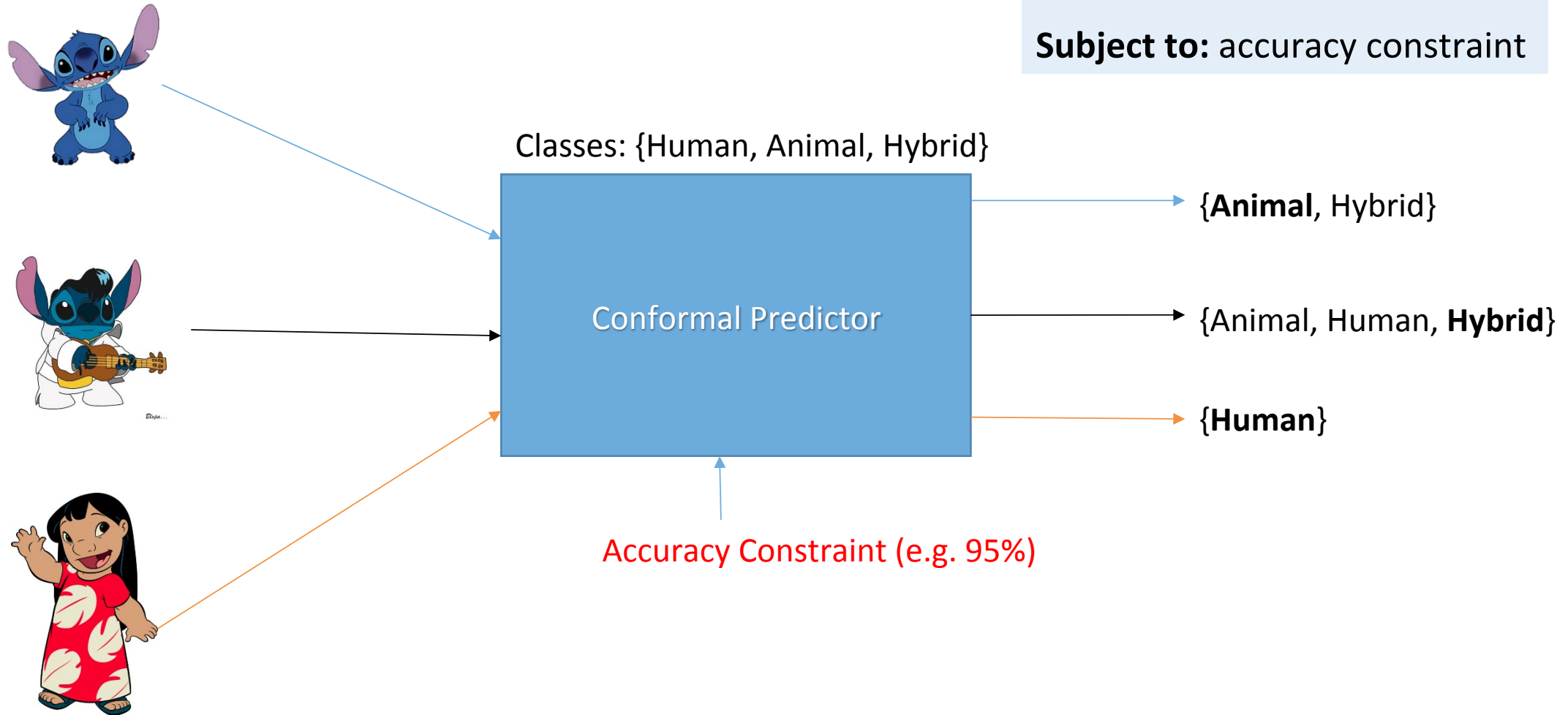
Want to minimize ambiguity of returned label sets.



# Conformal Prediction: Constrained Objective

**Minimize:** expected ambiguity

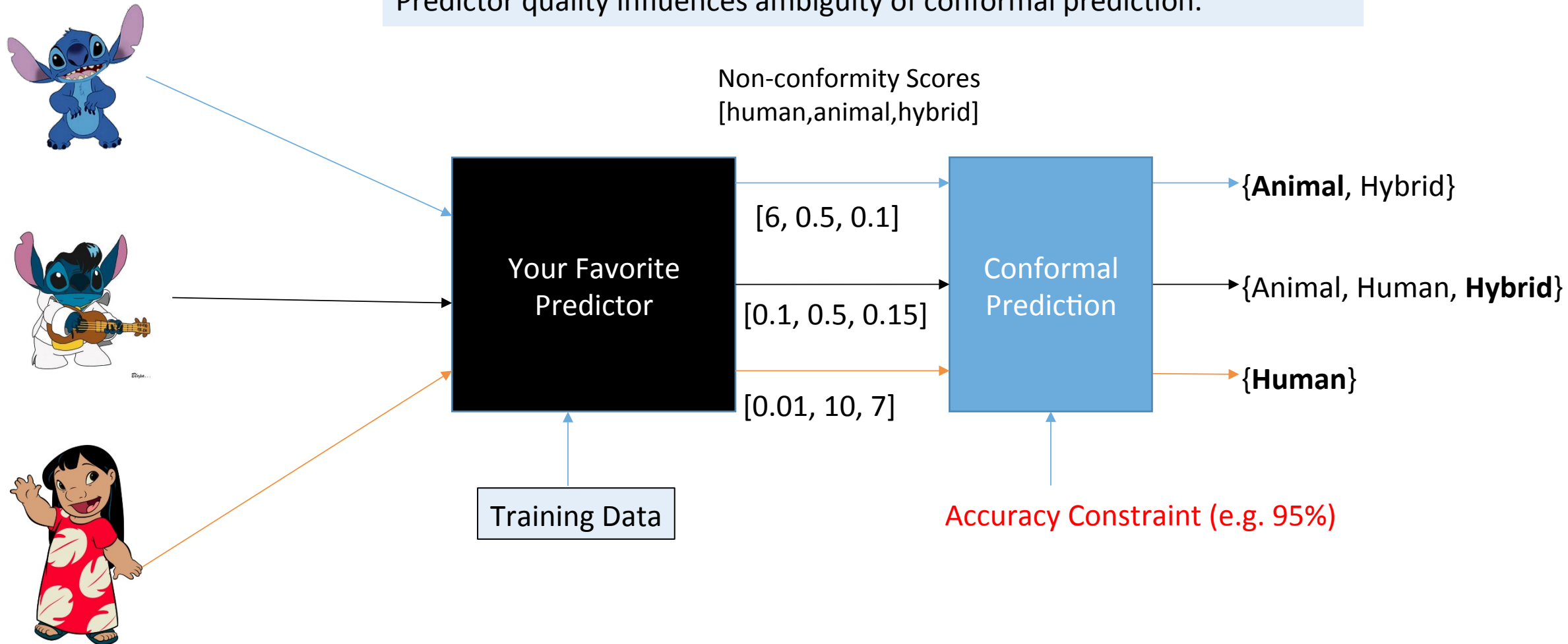
**Subject to:** accuracy constraint



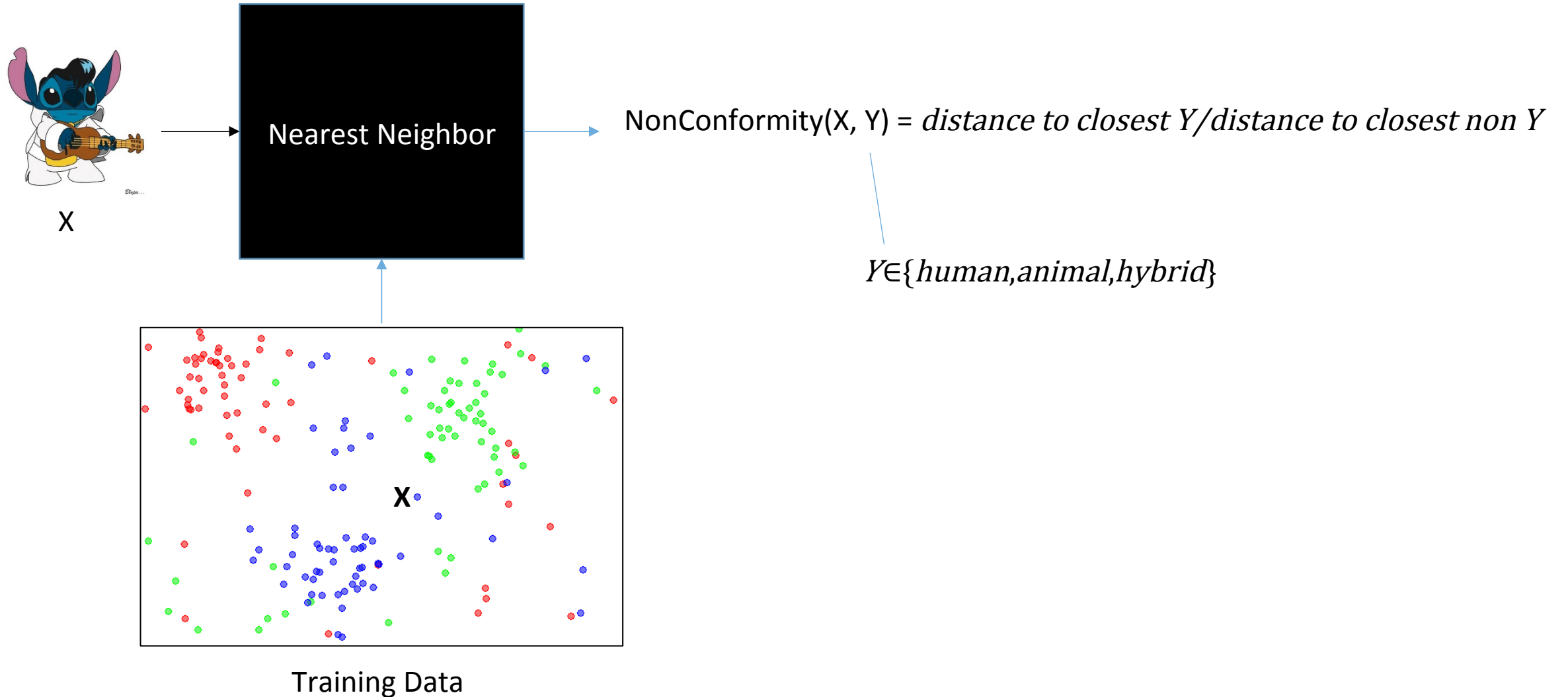
# Inside Conformal Prediction

Conformal prediction is a wrapper around any predictor that produces “non-conformity scores” over the classes relative to training data

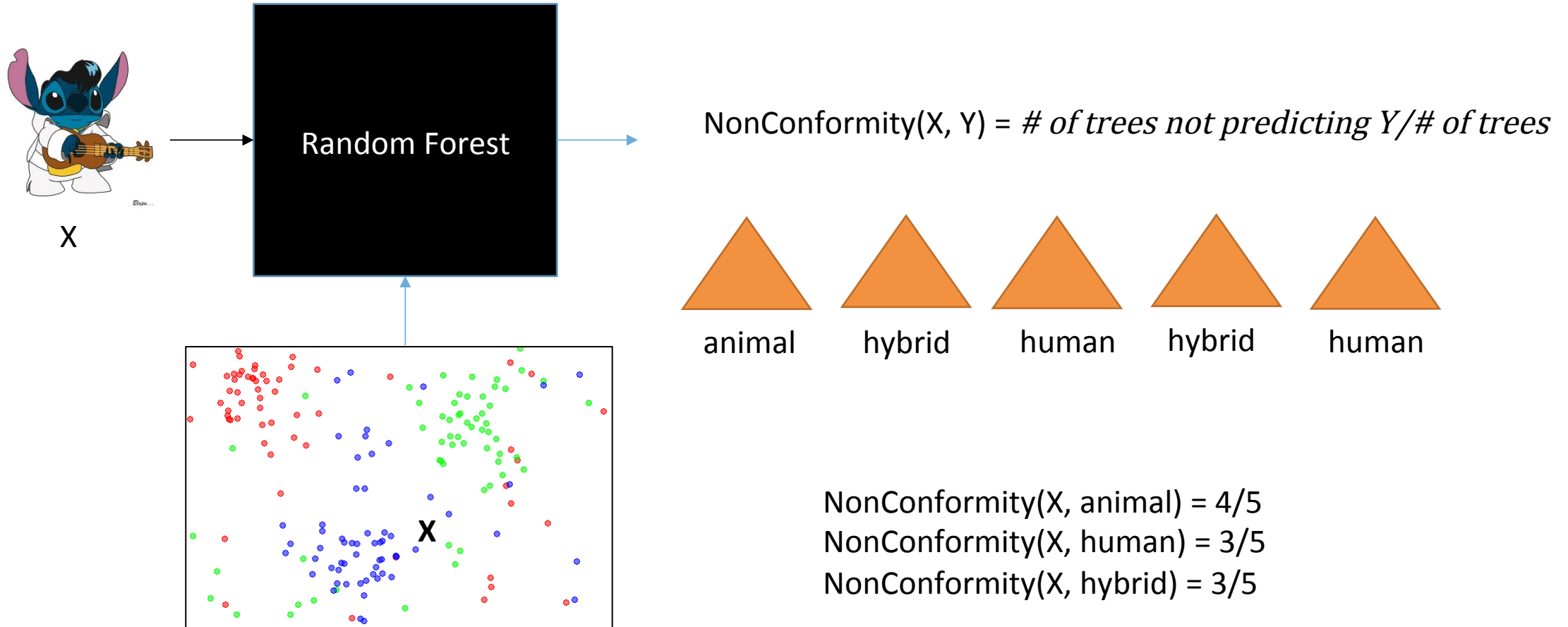
Predictor quality influences ambiguity of conformal prediction.



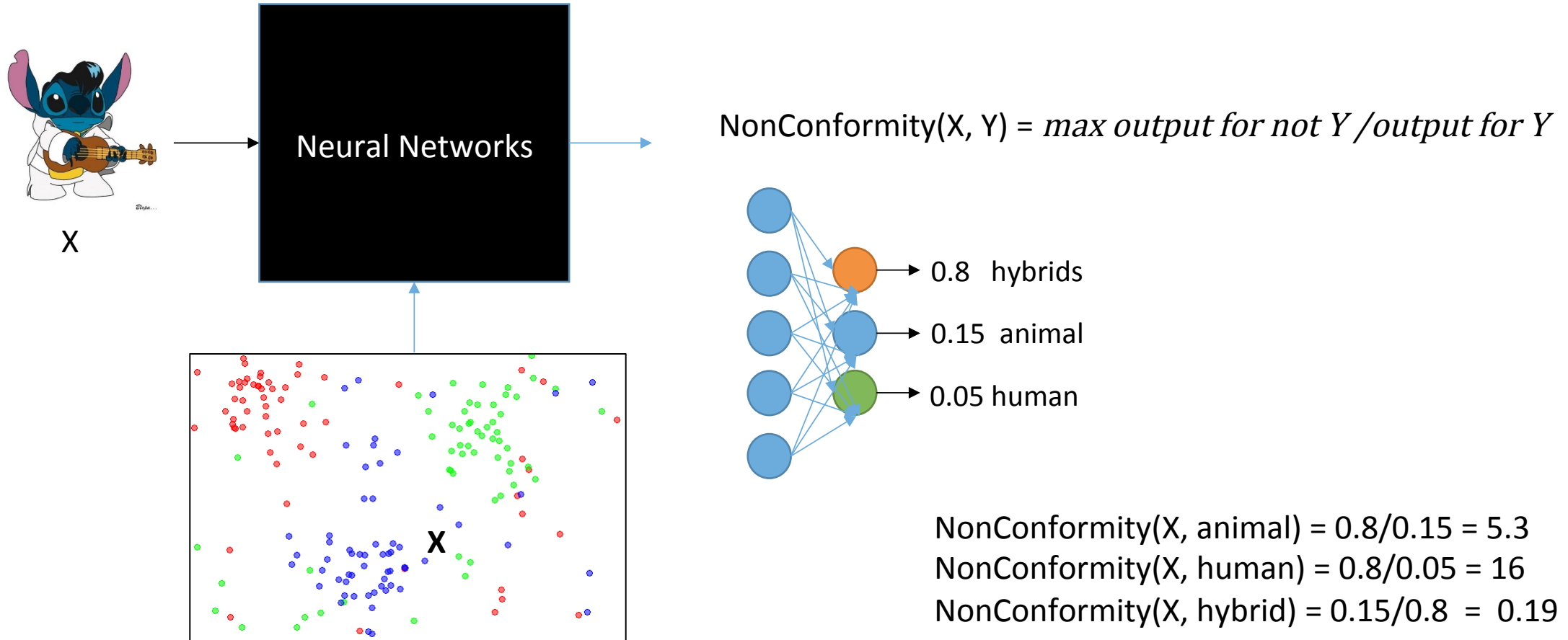
# Non-conformity Scores: Nearest Neighbor



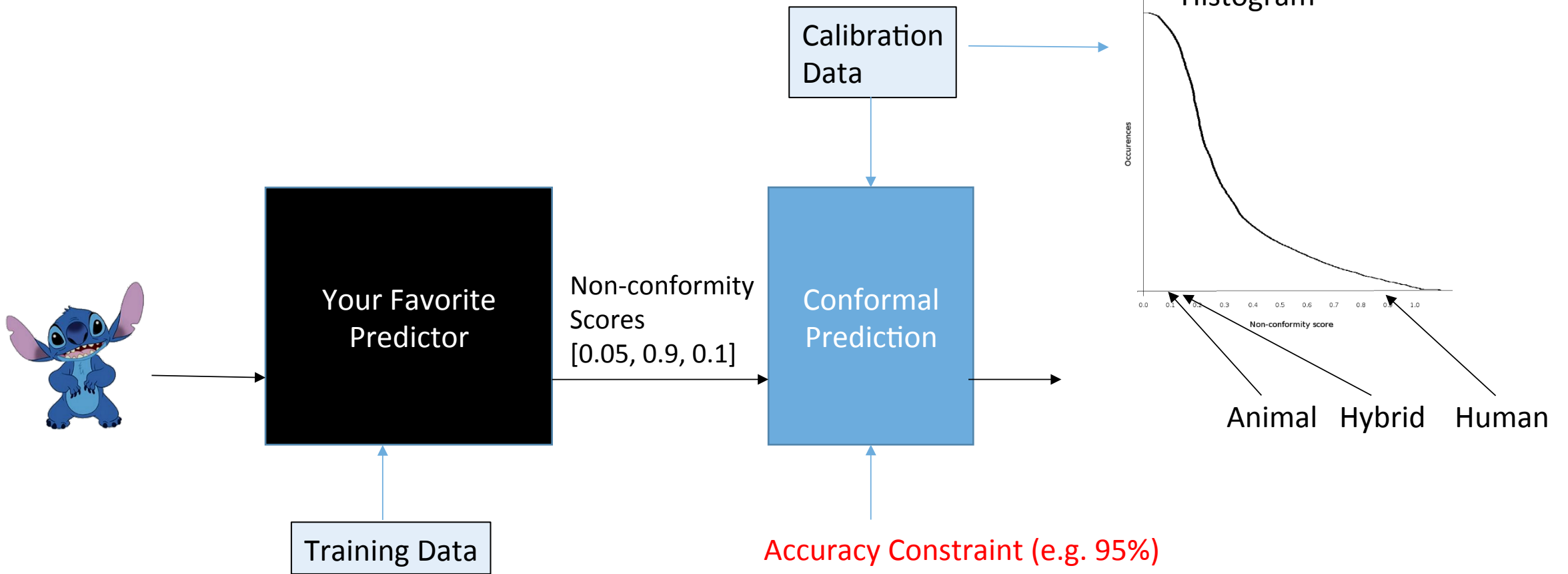
# Non-conformity Scores: Random Forest



# Conformal Prediction: Neural Networks

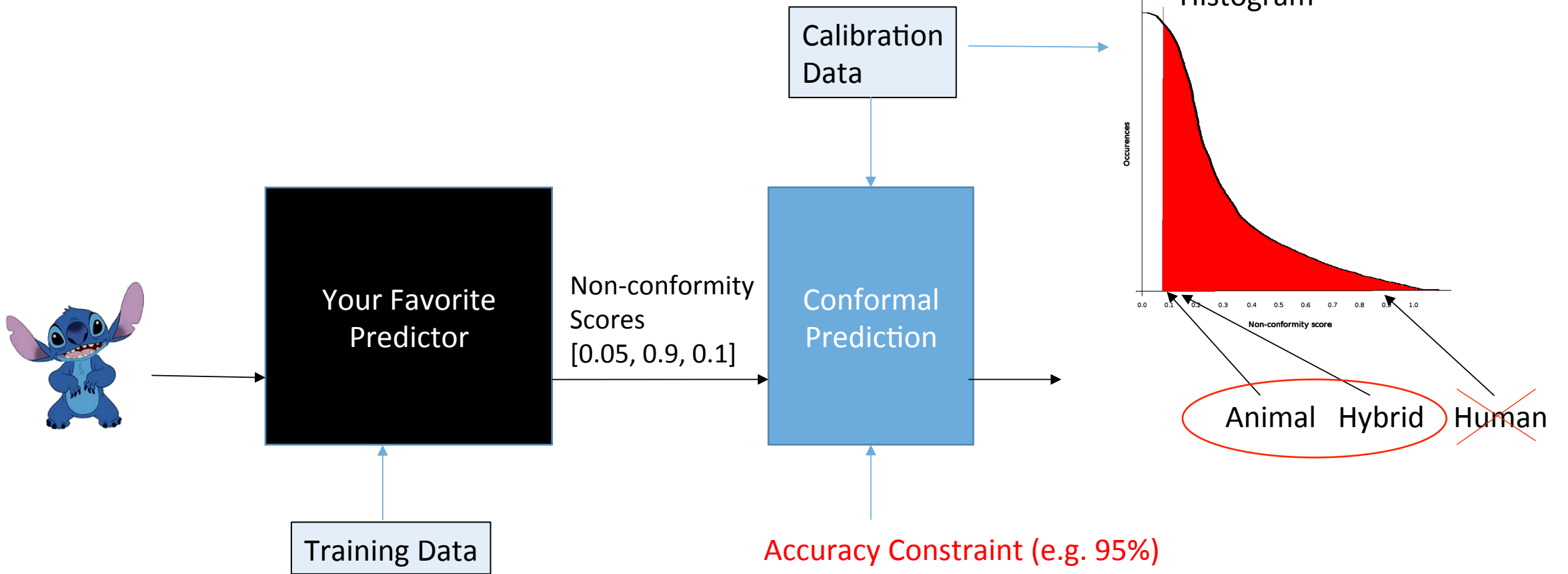


# Inside Conformal Prediction



Are at least 5% of the calibration scores weirder than this label with this example?

# Inside Conformal Prediction



Are at least 5% of the calibration scores weirder than this label with this example?



# Conformal Prediction: Empirical Evaluation

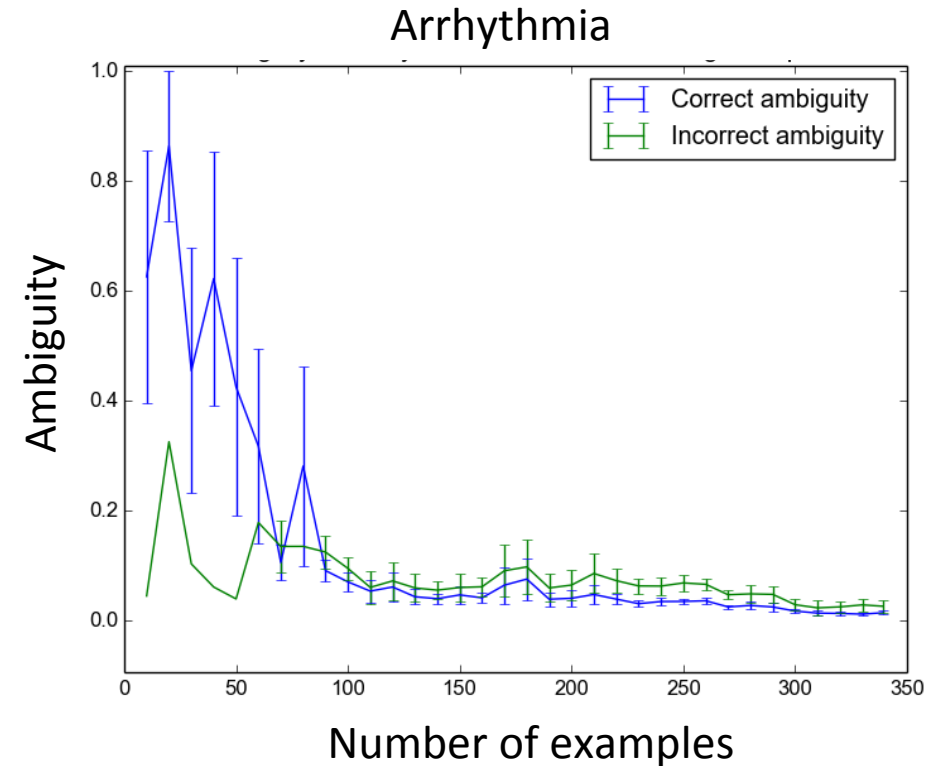
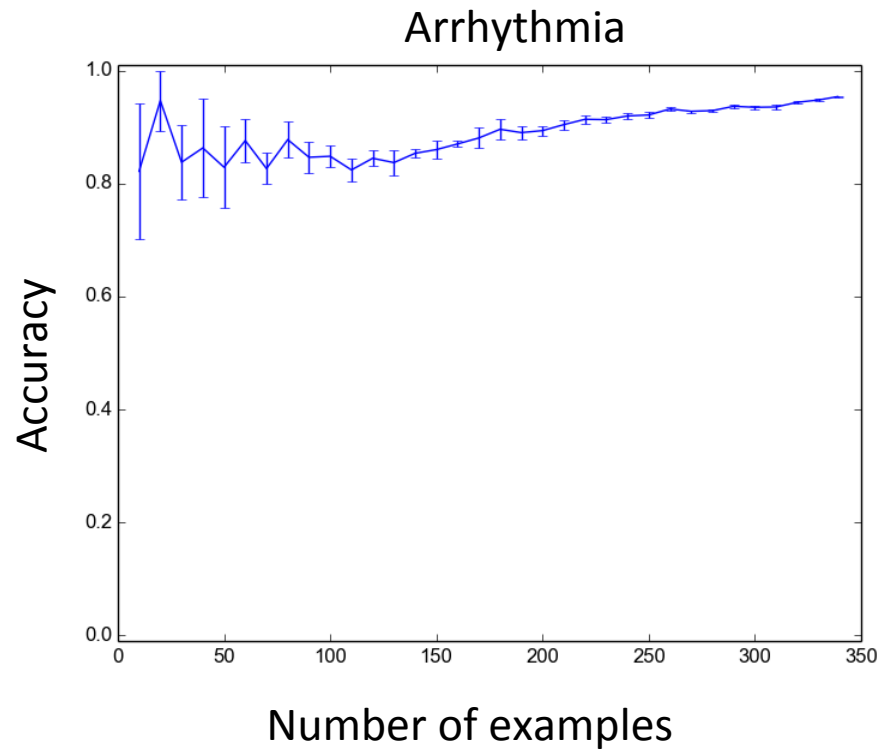
- Very few empirical evaluations of conformal prediction
  - Rarely look at ambiguity
- Most results for Nearest Neighbor -- often yields large ambiguity in our experience
- How does ambiguity vary with amount of training data in closed worlds?
- How does conformal prediction perform in open worlds?

# Closed World: Random Forest Results

- Arrhythmia: 452 data points, 13 labels, major class imbalances
- Cardiocography: 2126 instances, 10 labels, balanced classes
- Image Segmentation: 2310 instances, 8 labels, balanced classes
- Iris: 150 instances, 3 labels, balanced classes

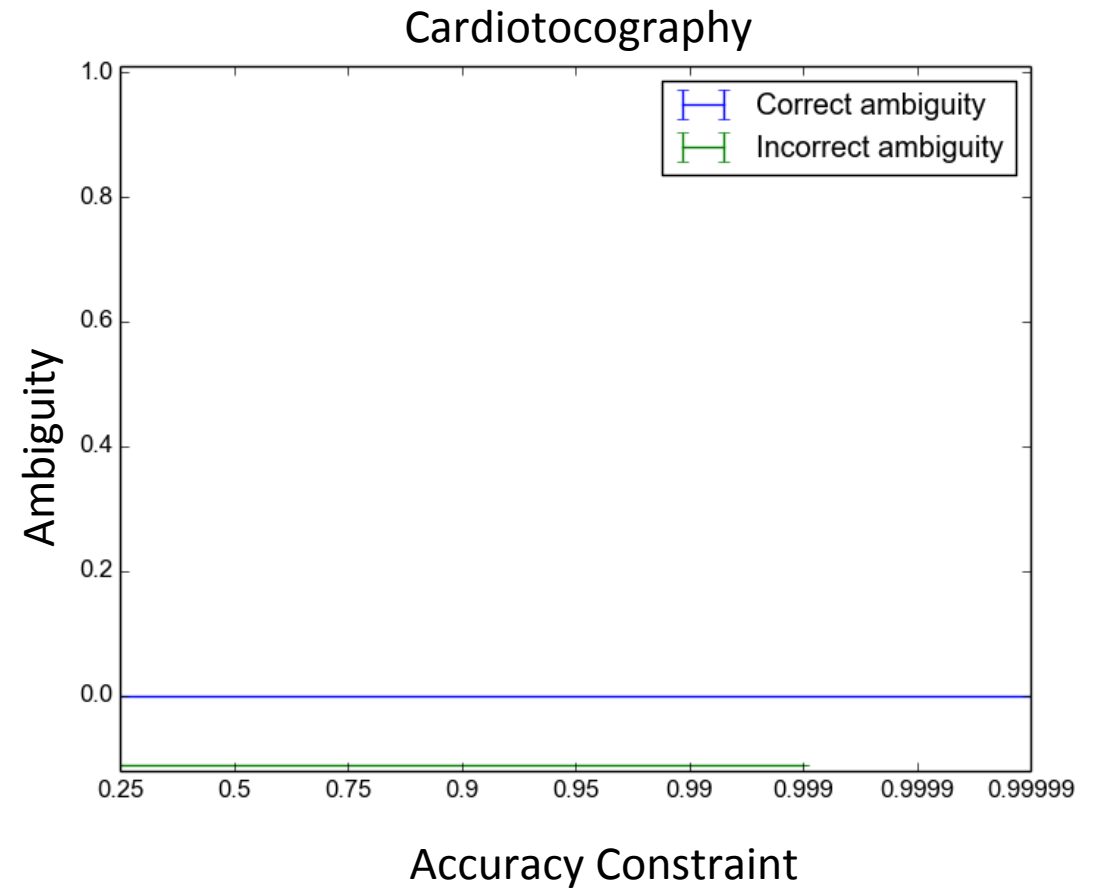
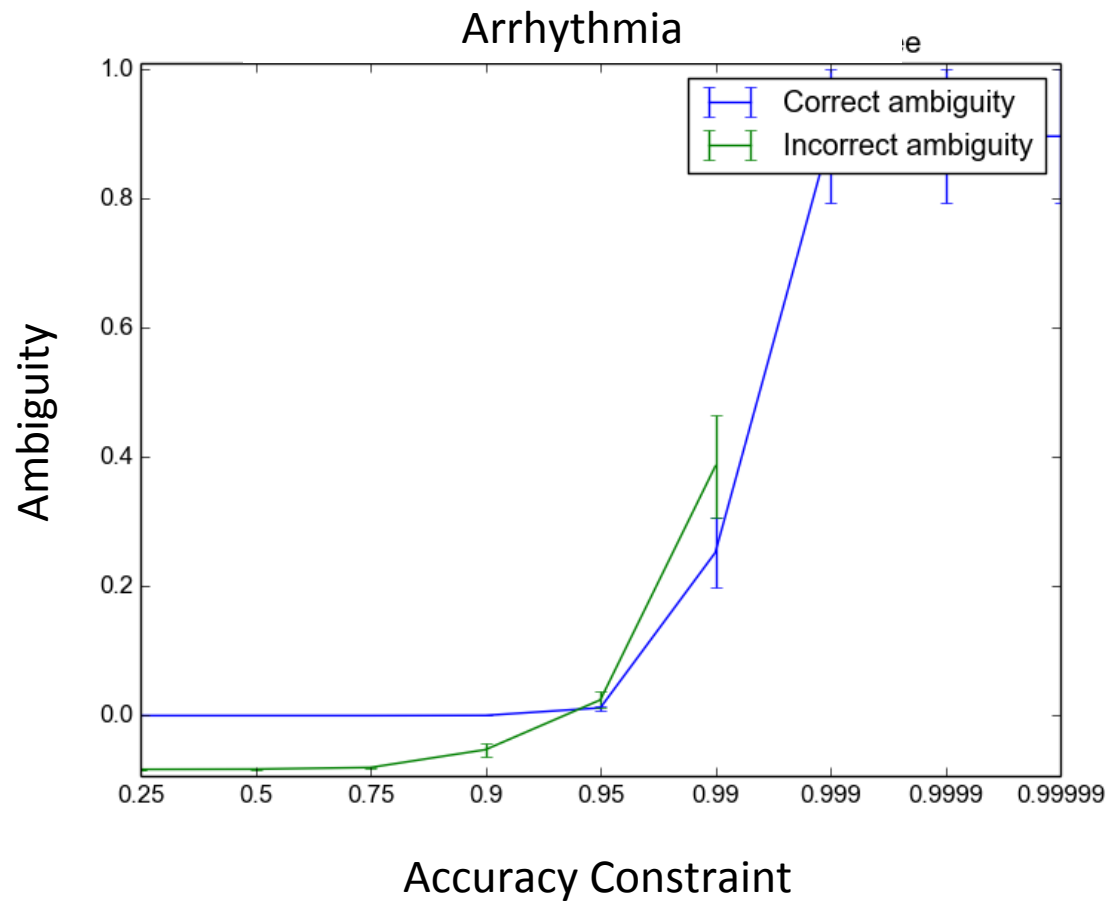
# Closed World: Random Forest

Accuracy Constraint = 95%

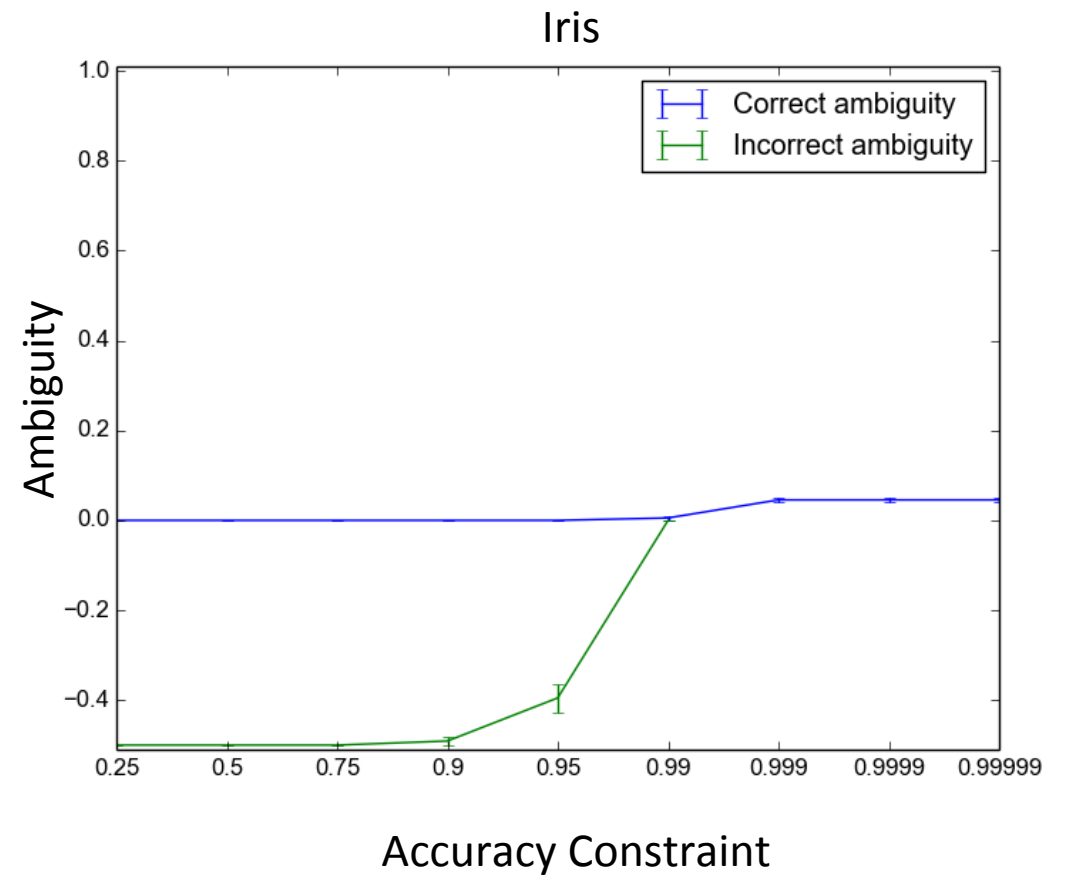
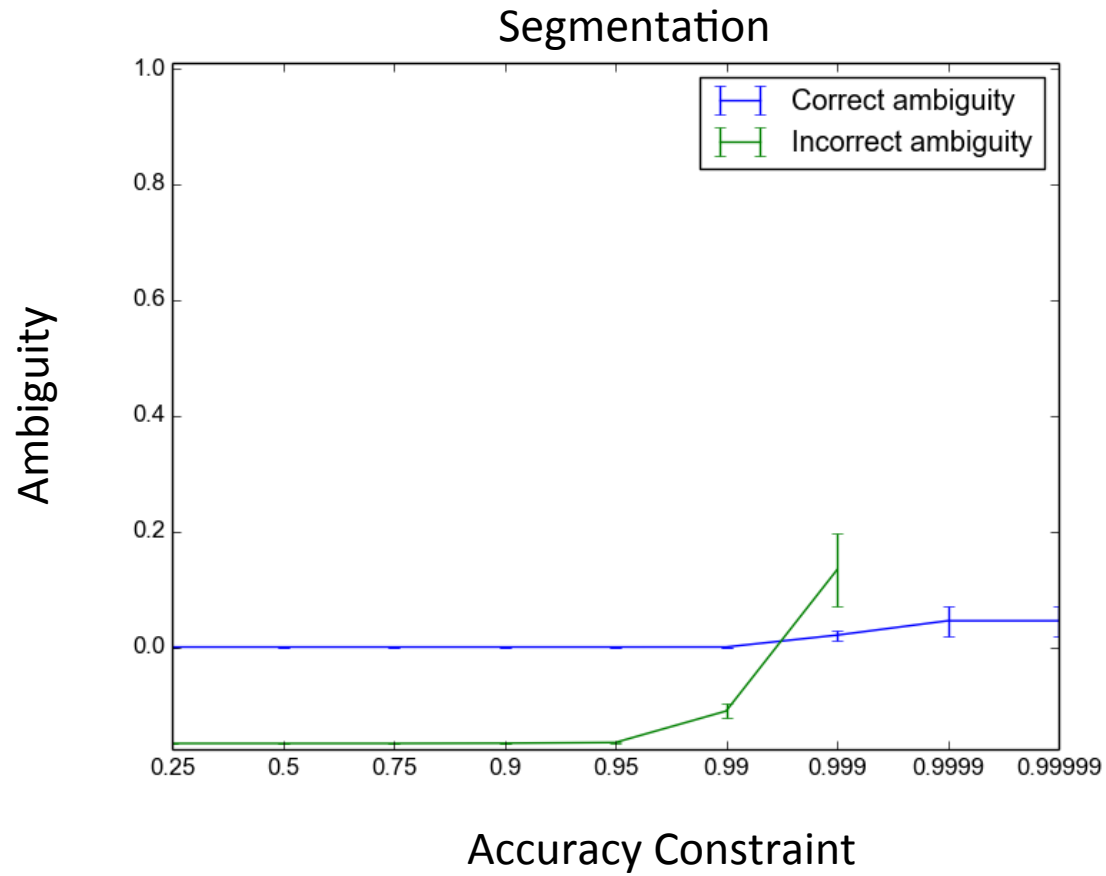


Other data sets are qualitatively similar.

# Closed World: Random Forest



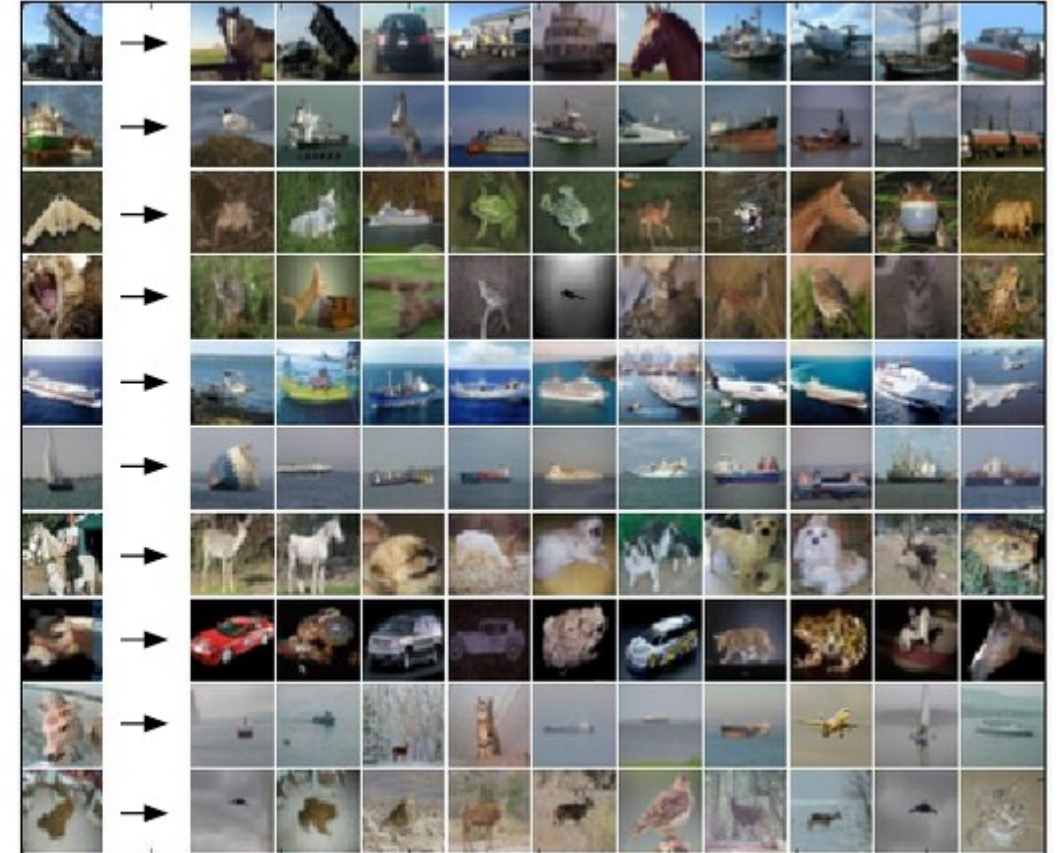
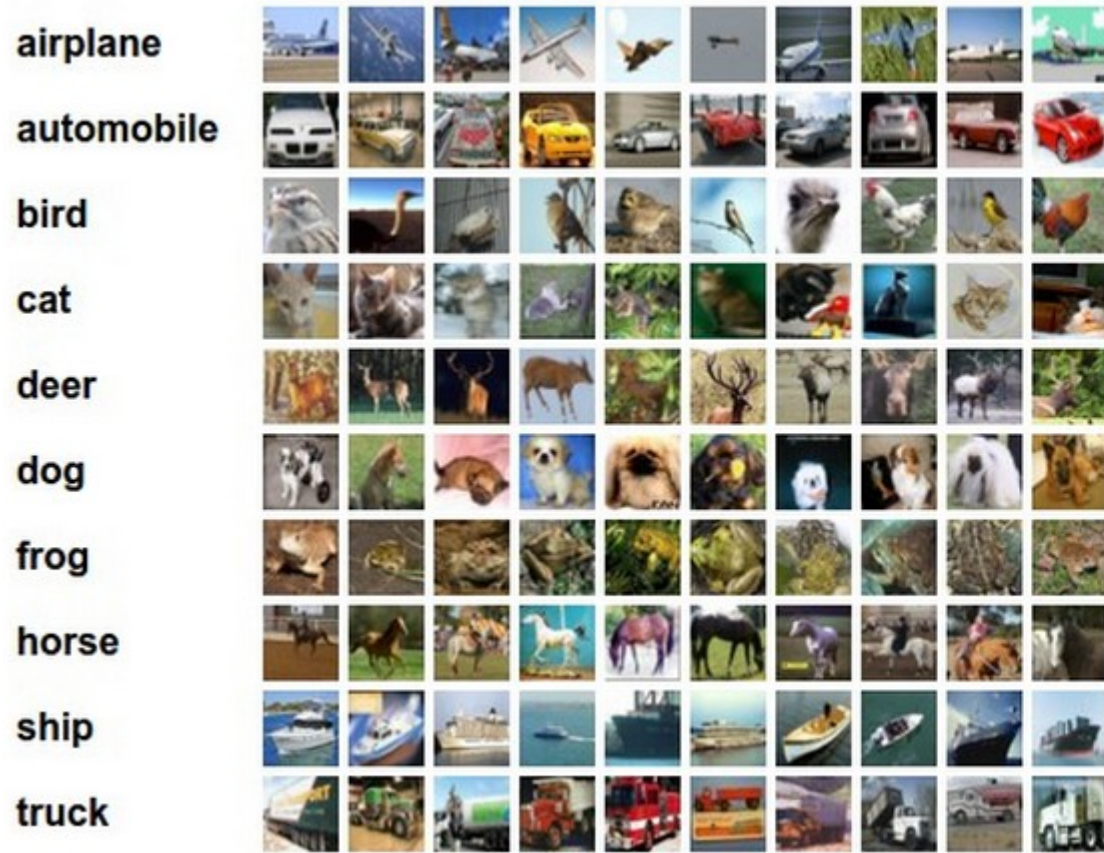
# Closed World: Random Forest



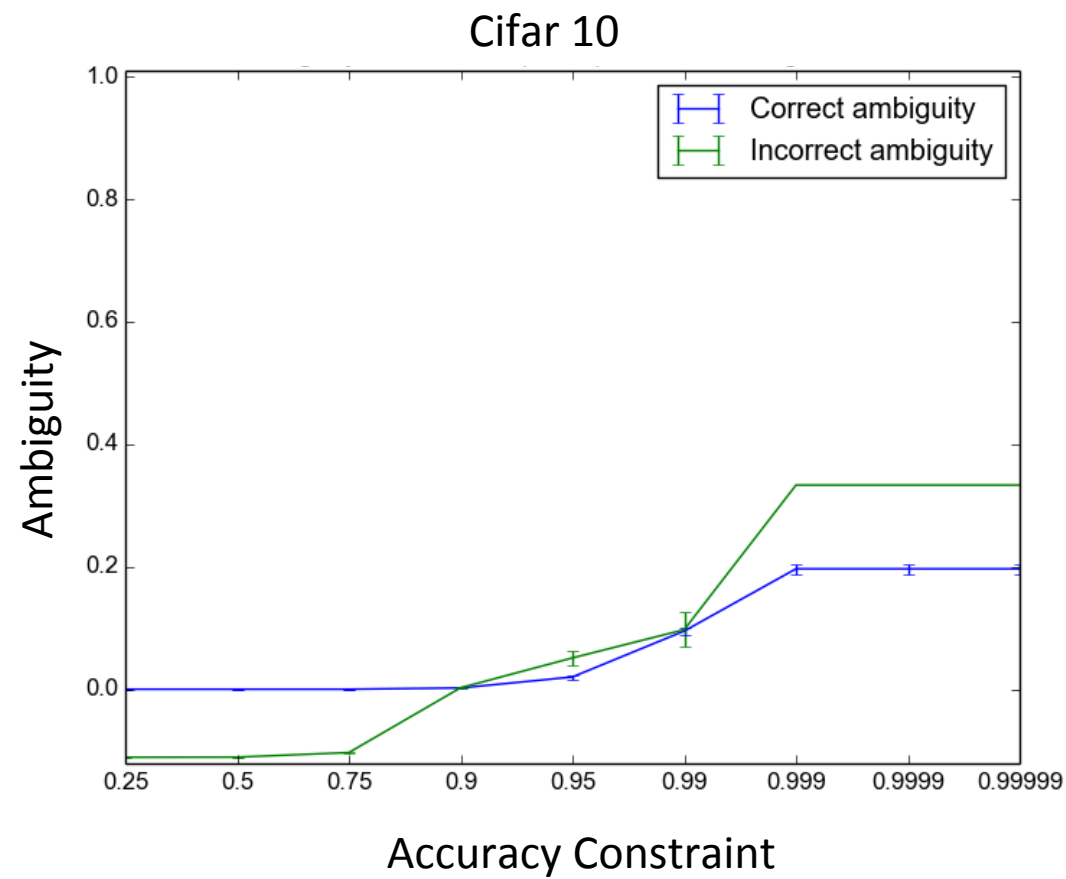
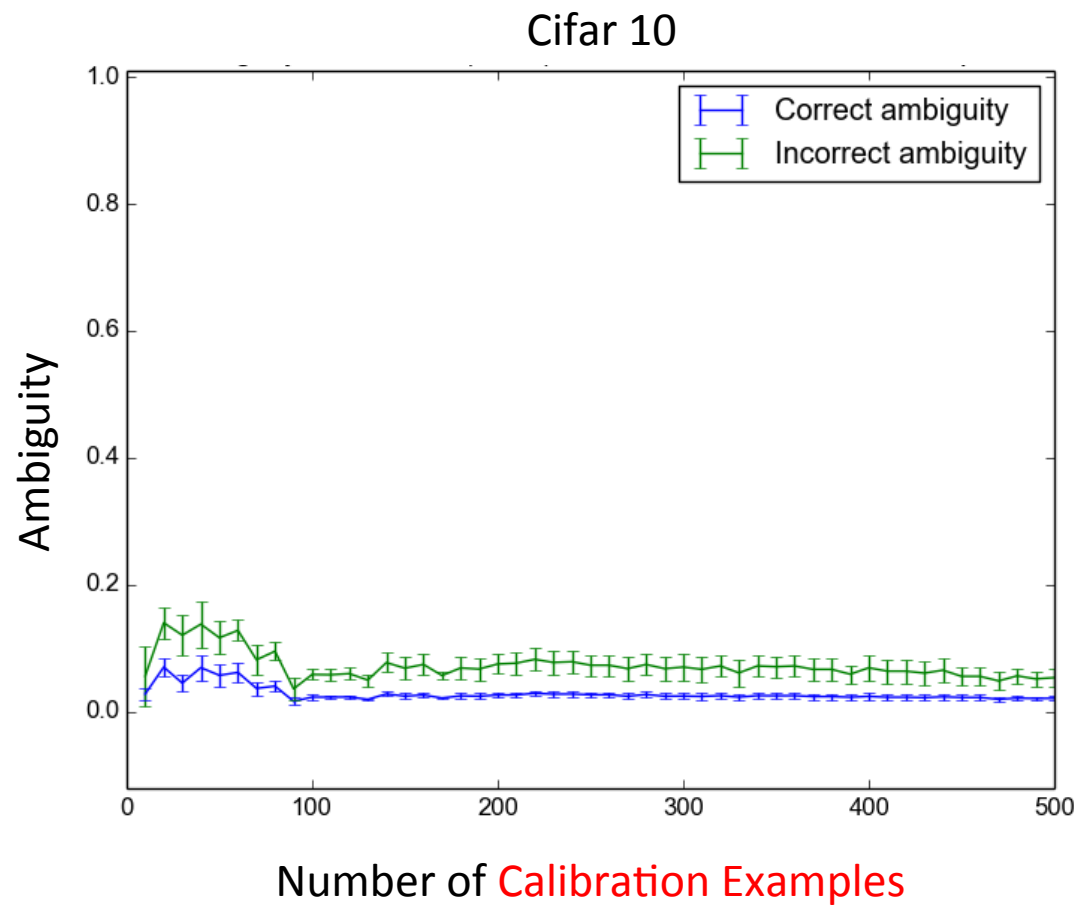
# Closed World: Random Forest

- Overall, we see “ideal behavior” on these data sets.
- Close to 0 ambiguity with small amount of data.

# Closed World: Convolutional Network: Cifar 10



# Closed World: Deep Net

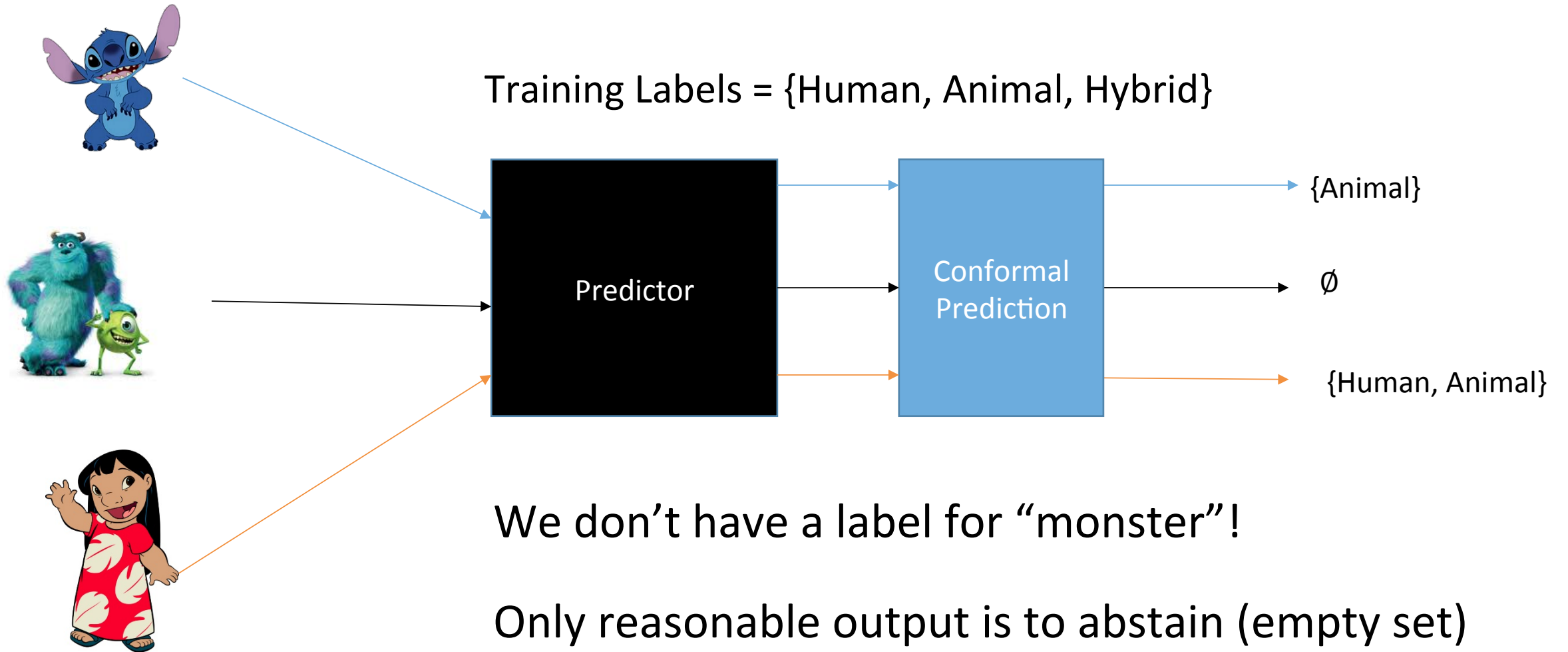




# Closed World Observations

- Overall, we see “ideal behavior” on these data sets for both random forest and convolutional network.
- Close to 0 ambiguity with small amount of data.
- Neural network very rarely abstains (negative ambiguity) compared to the random forest

# Open World: Conformal Prediction

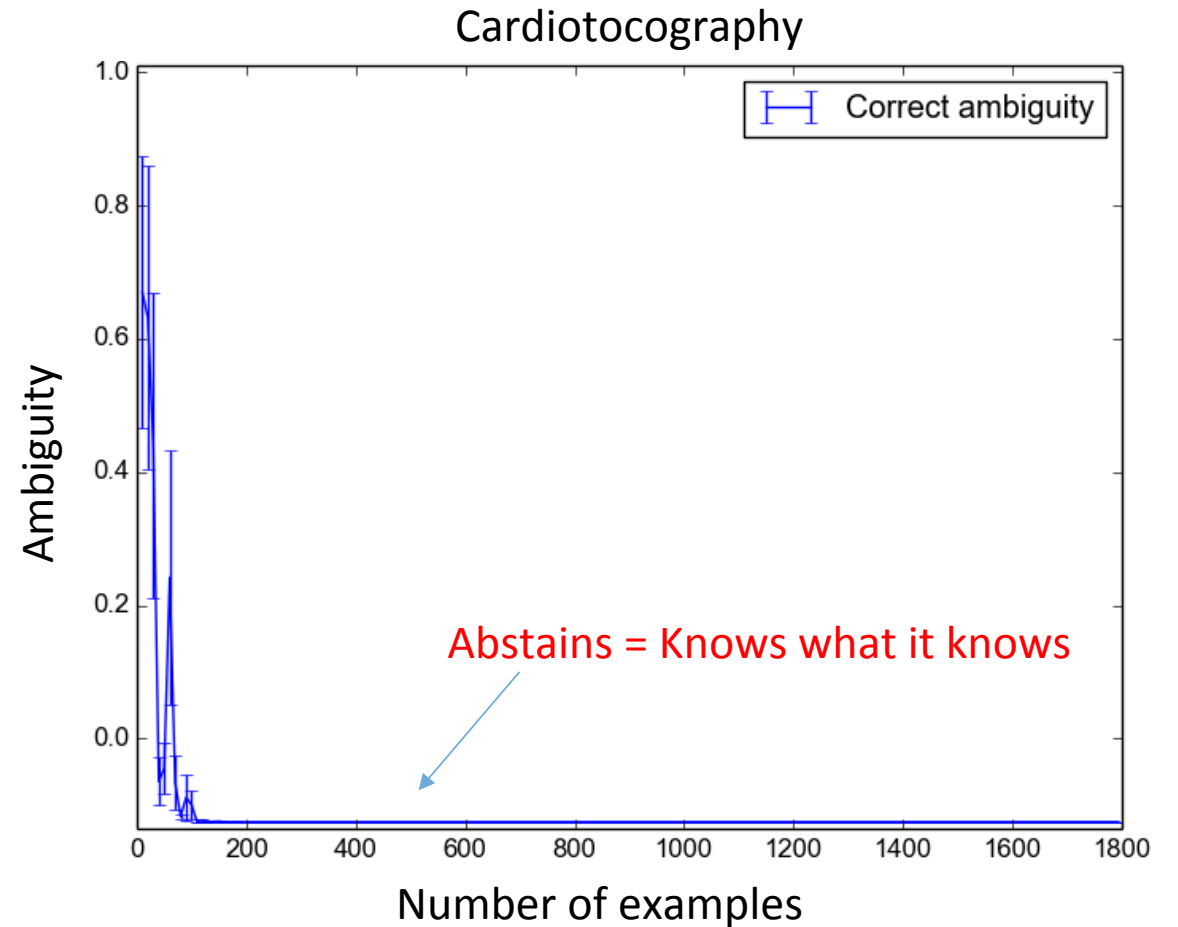
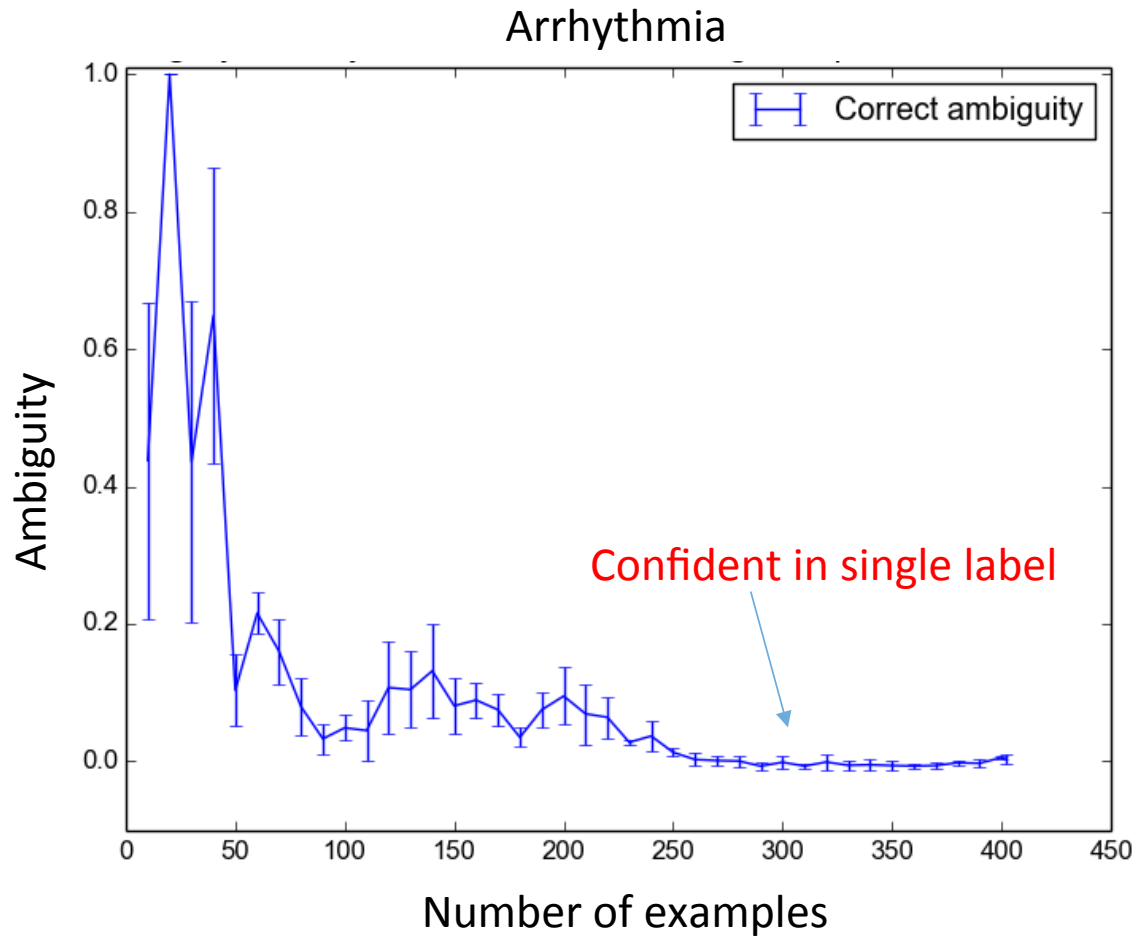


# Open World Experiments

- Feed novel classes to conformal predictor
- **Random Forest** : withheld a label from each training set
- **Convolutional Network** : feed it images that have nothing to do with  
Cifar 10

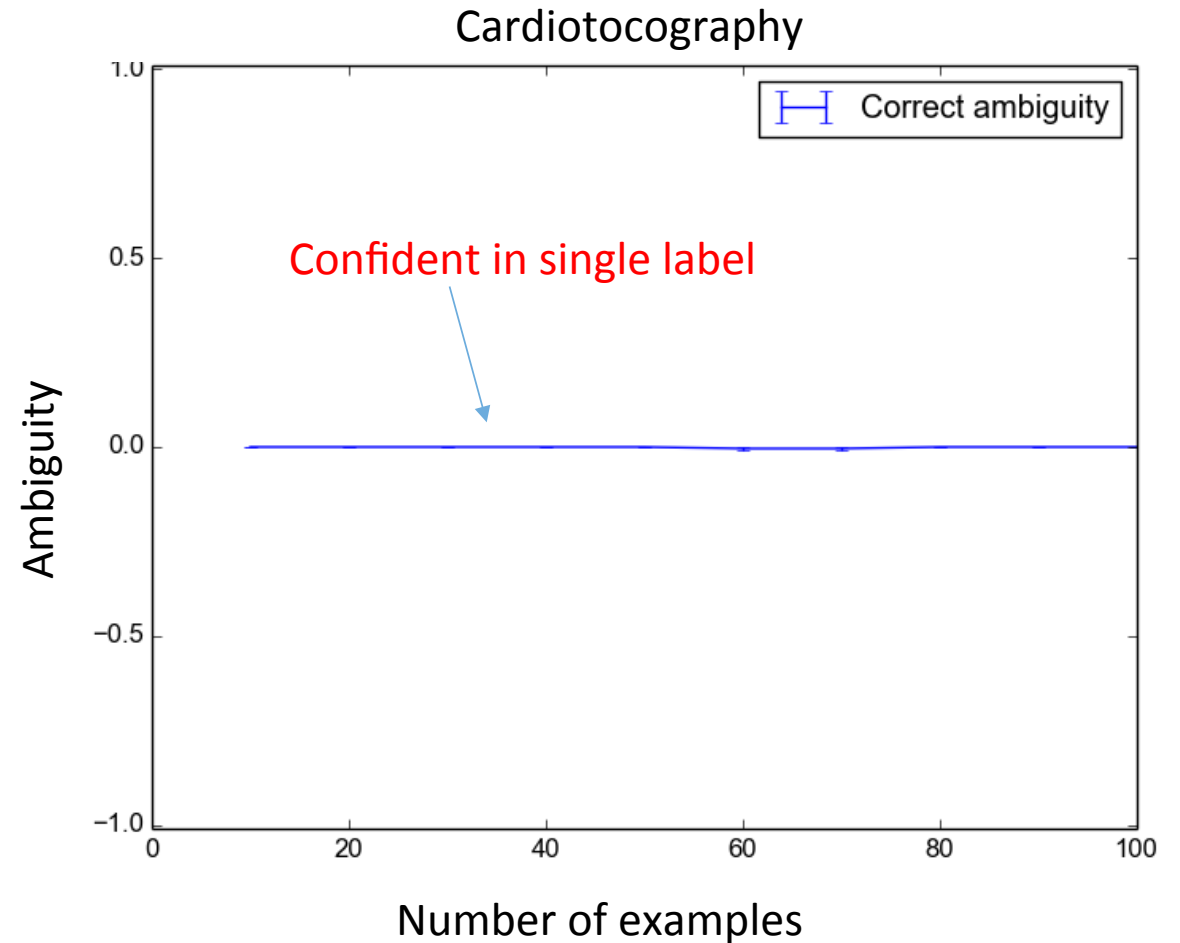
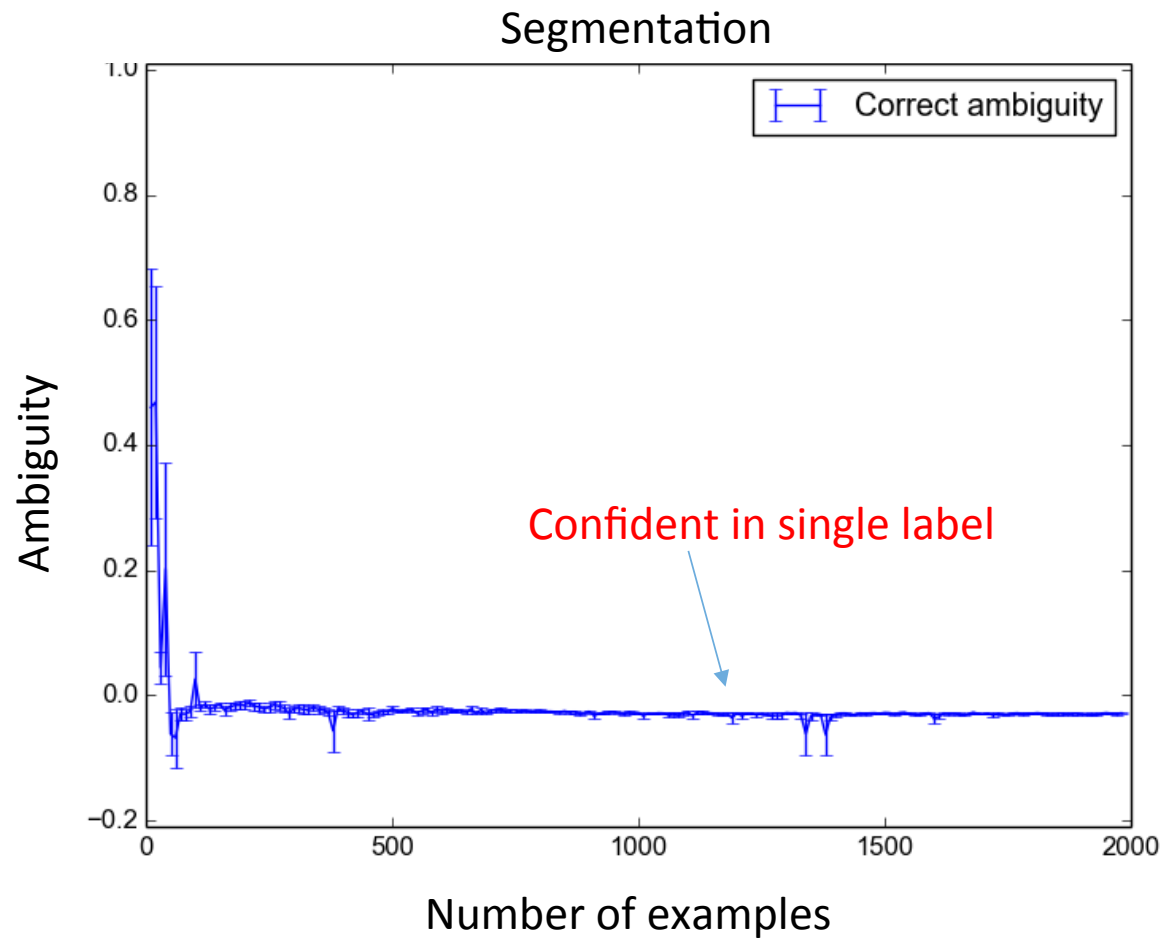
# Open World: Random Forest

Ambiguity for just novel classes

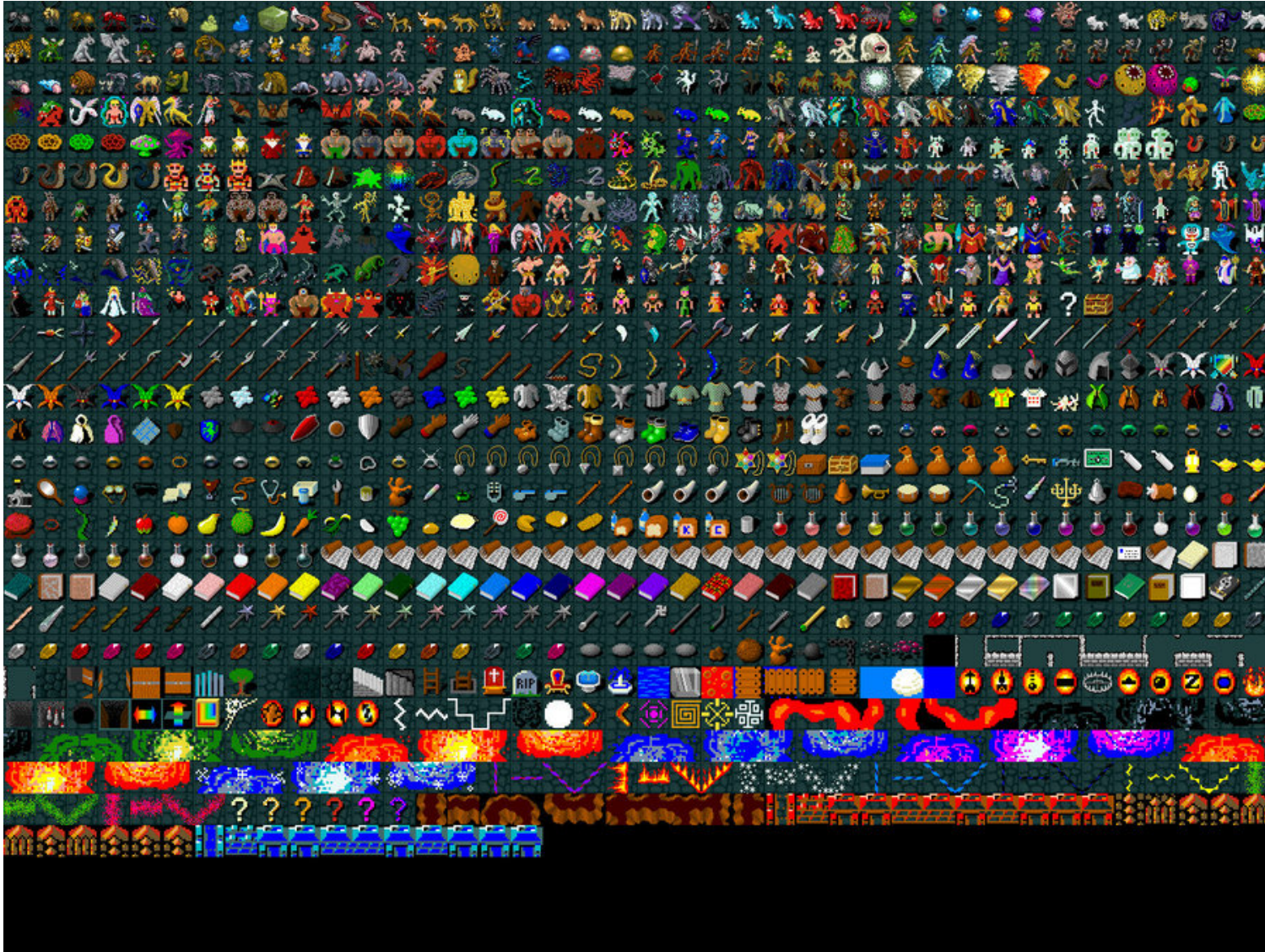


# Open World: Random Forest

Ambiguity for just novel classes



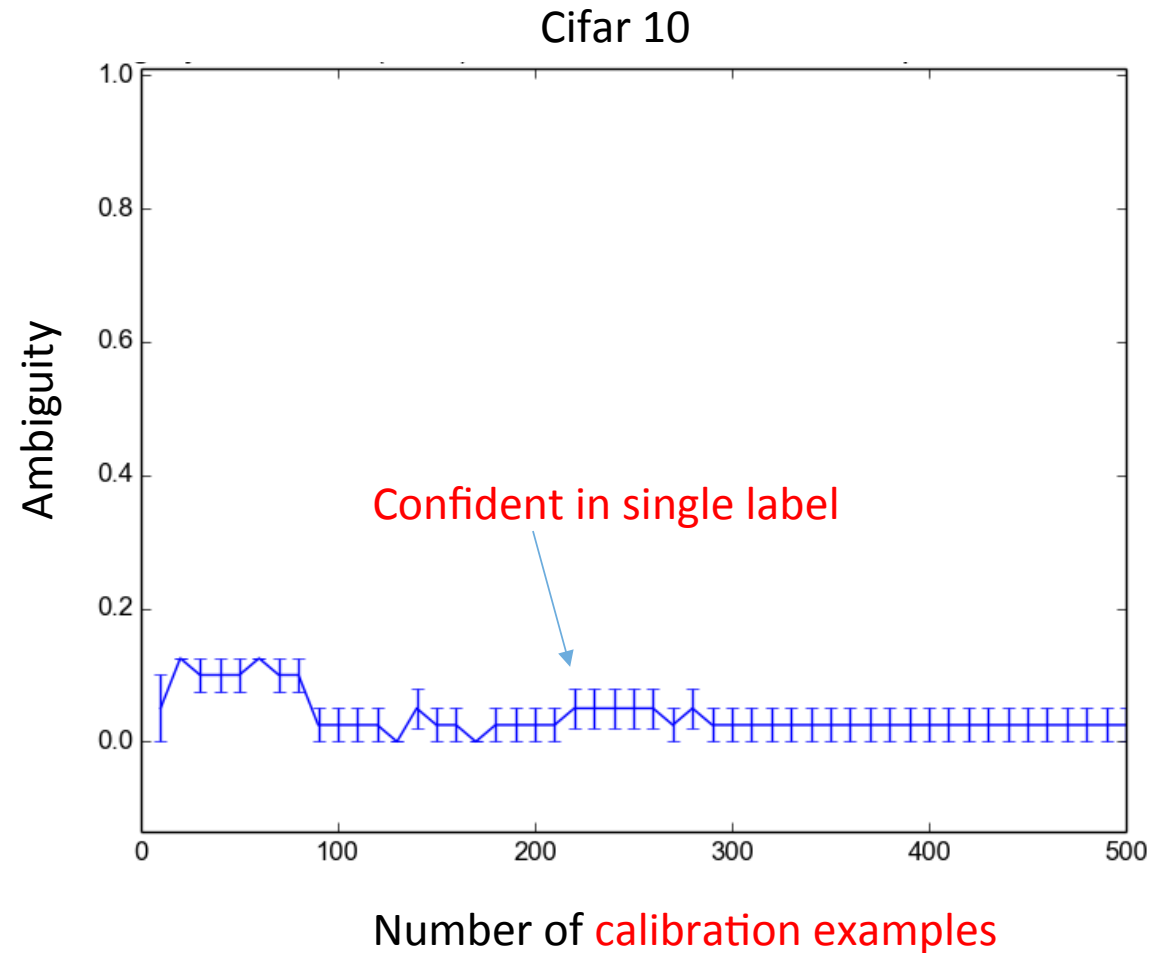
# Open World: Convolutional Network



Nethack sprite sheet  
images

# Open World: Convolutional Network

Ambiguity for just novel Nethack images

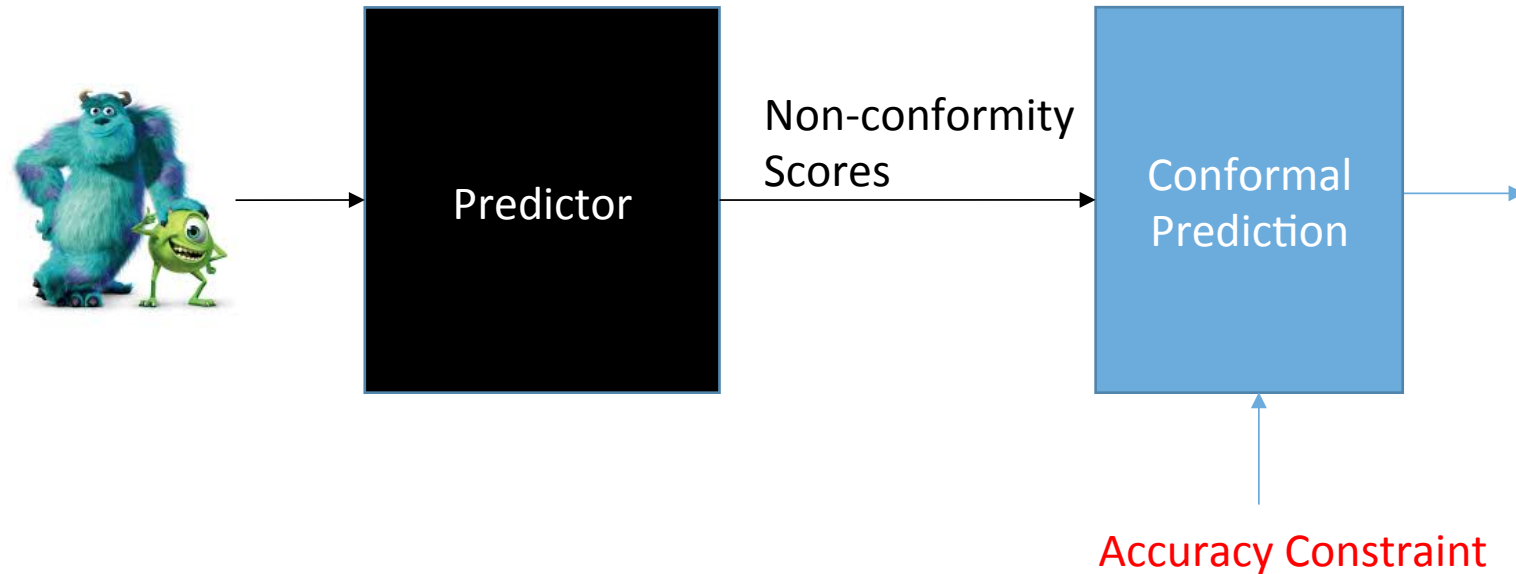


# Open World Observations

- In all but one case there was practically no abstention
- The theory of conformal prediction does not address the issue of open worlds
- Appears that standard conformal prediction on its own is not sufficient for open worlds



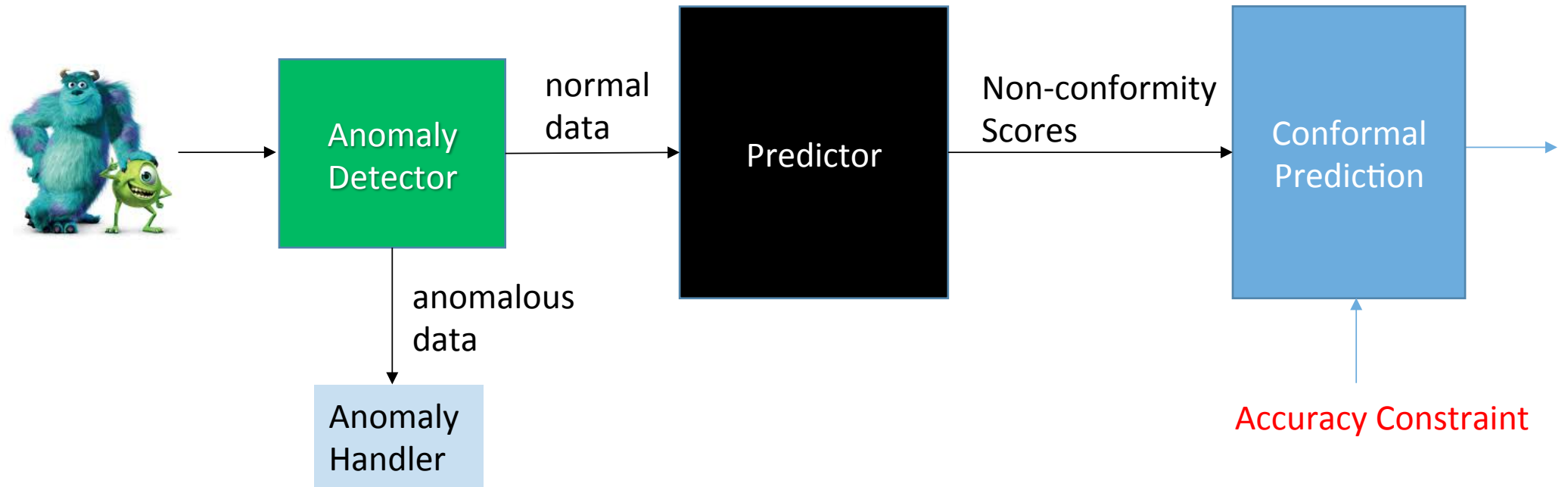
# Next Steps for Open Worlds



New algorithms for training predictors.

**Goal:** yield reliable abstention for novel classes.

# Next Step for Open Worlds



How to select anomaly threshold?

Can we provide any guarantees in open worlds?