

# Rashomon Capacity: Measuring Predictive Multiplicity in Classification

Hsiang Hsu and Flavio P. Calmon

School of Engineering and Applied Science, Harvard University

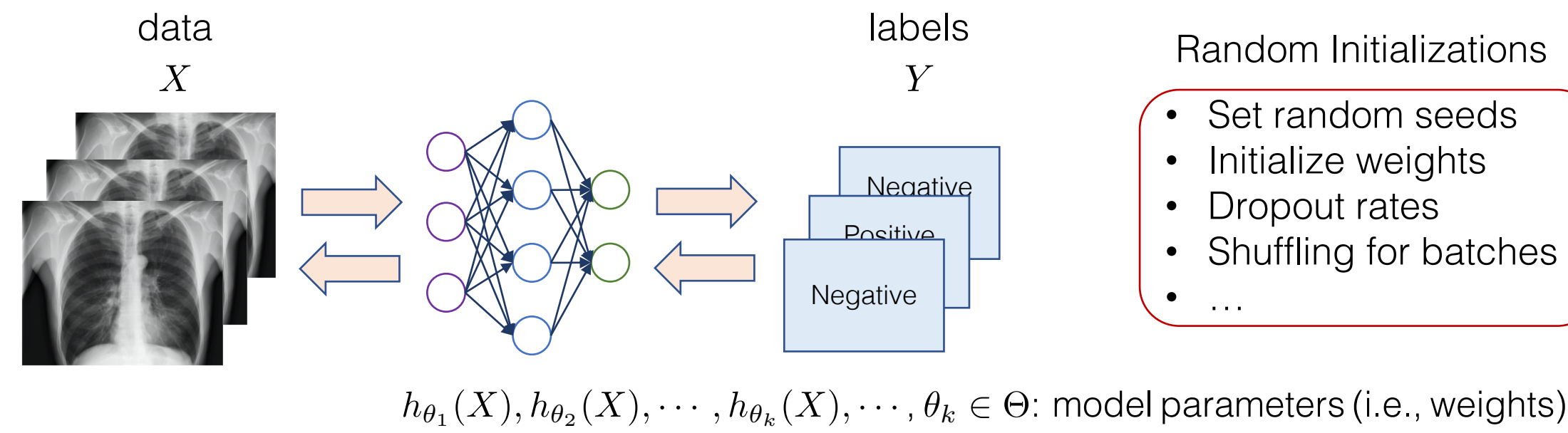


## Take-Away

1. Rashomon Capacity is a computationally tractable metric for predictive multiplicity
2. Disclosing predictive multiplicity is critical for individual-level decision-making, and current practice does not report it!

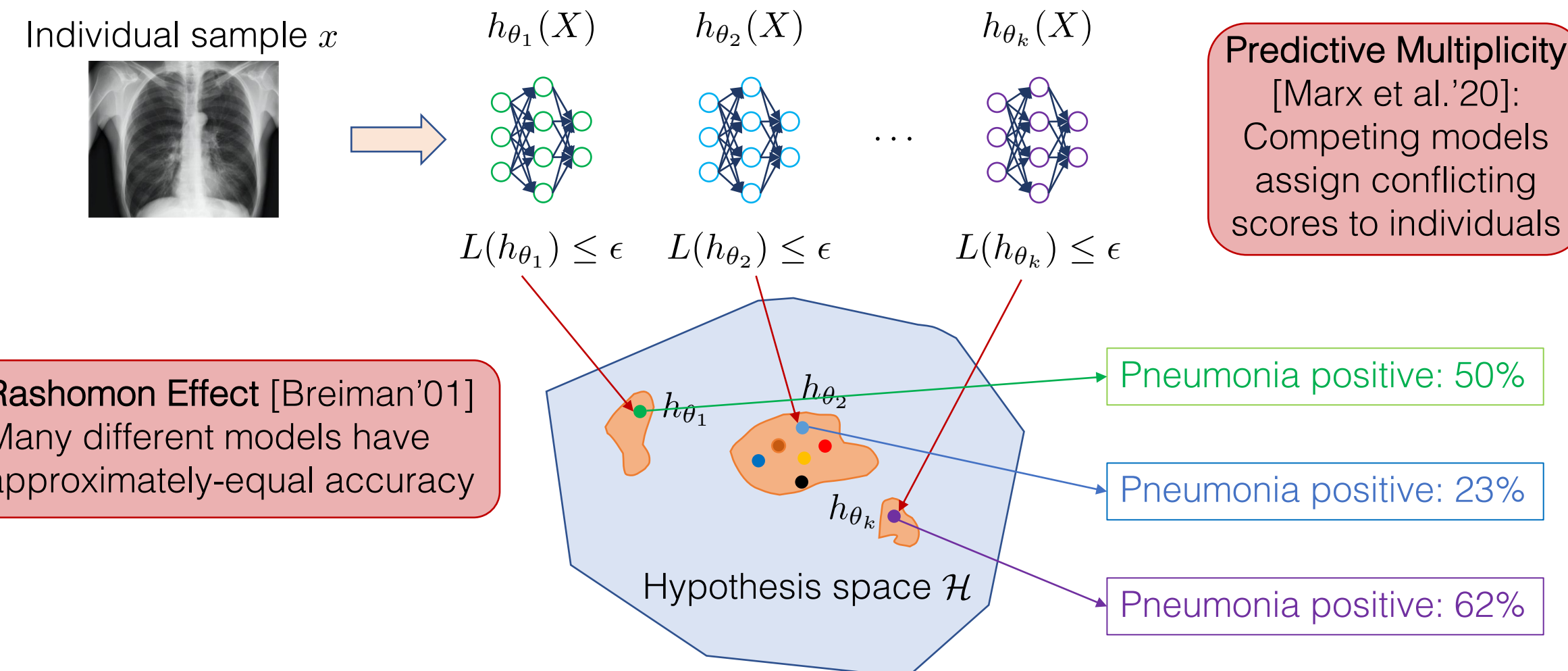
## Rashomon Effect and Predictive Multiplicity

### Training Time



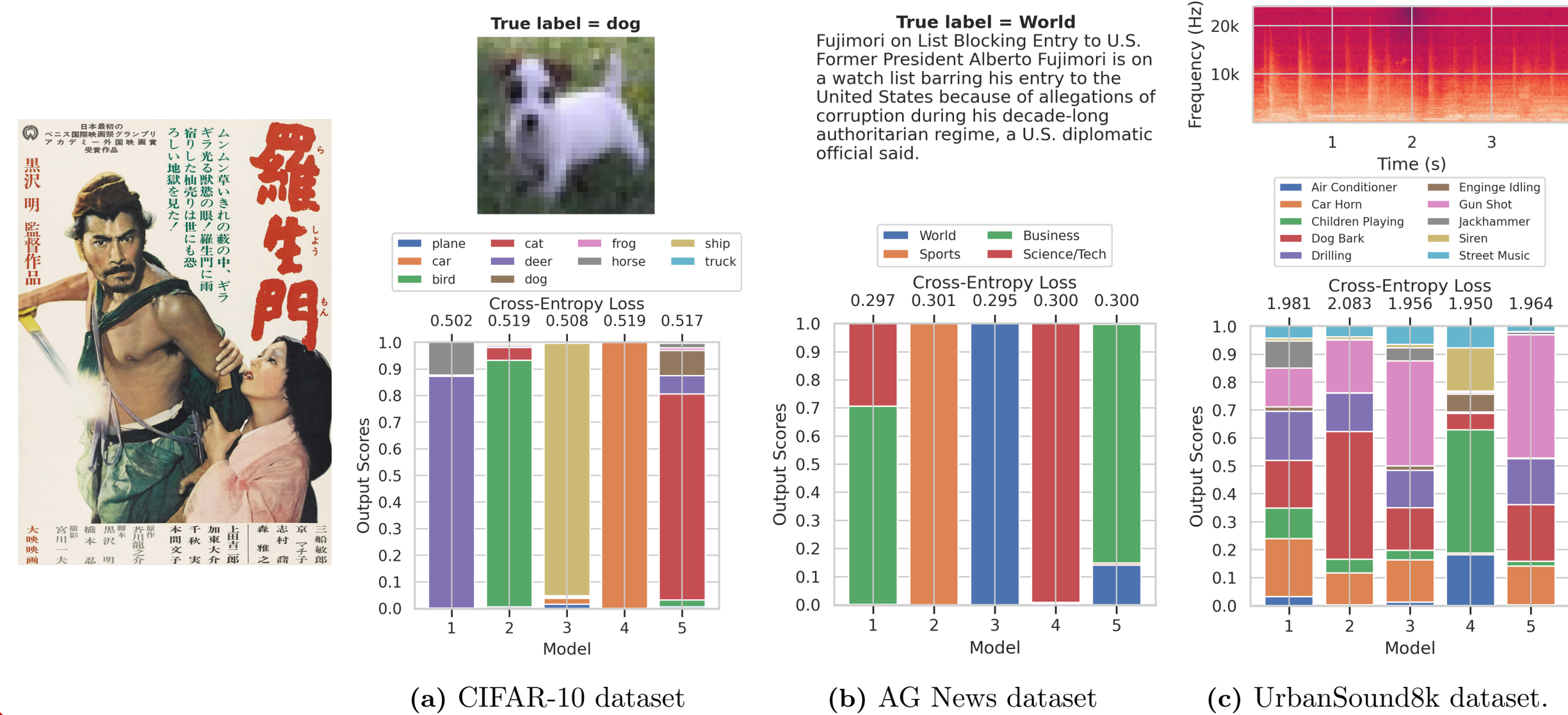
### Test Time

test loss:  $L(h_\theta) = \mathbb{E}[\ell(h_\theta(X), Y)]$



**Rashomon Effect [Breiman'01]**  
Many different models have approximately-equal accuracy

Rashomon set:  $\mathcal{R}(\mathcal{H}, \epsilon) \triangleq \{h_\theta \in \mathcal{H}; L(h_\theta) \leq \epsilon\}$

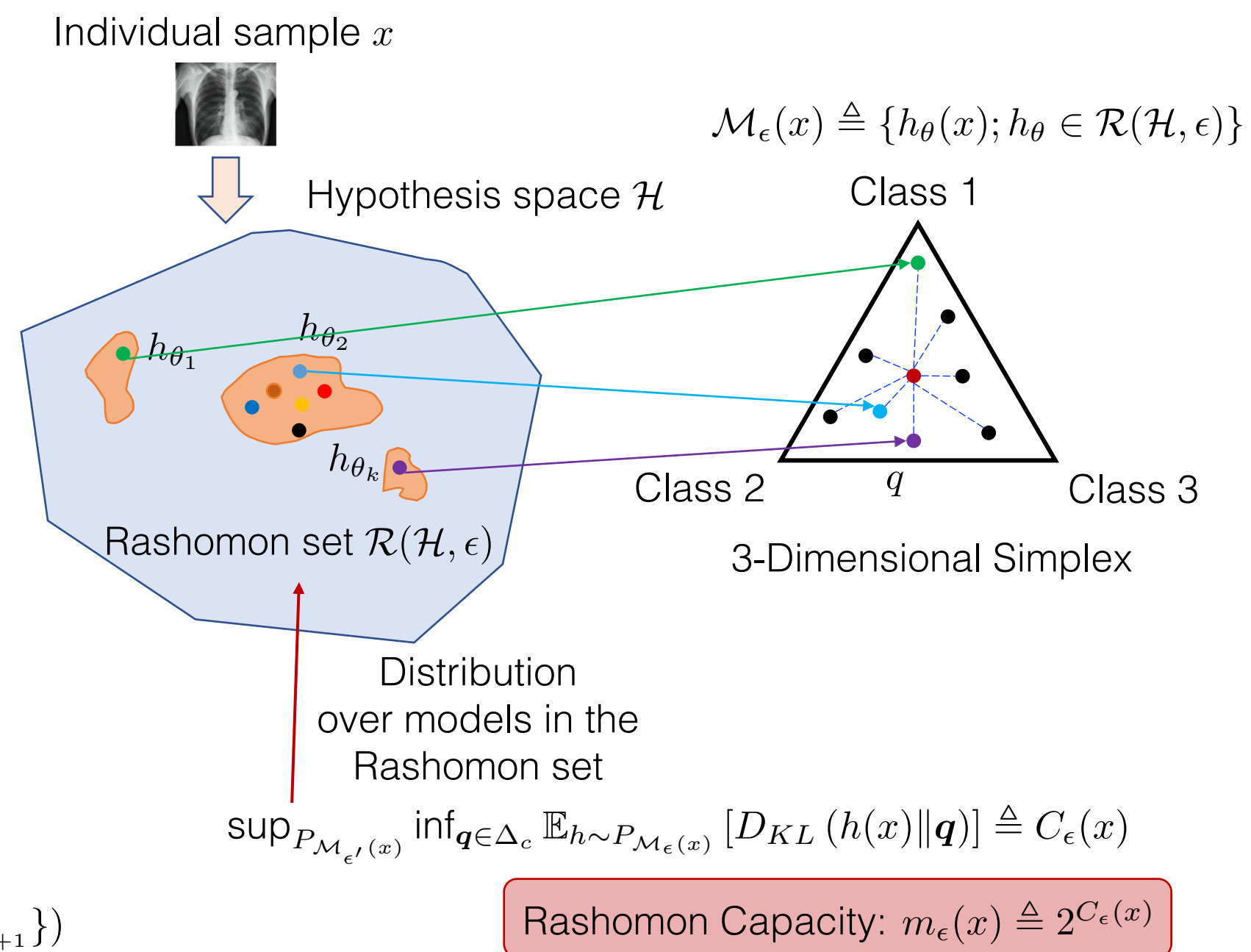


## Rashomon Capacity

1.  $h_{\theta_1}(x) = [0.05, 0.90, 0.05] \in \Delta_c$   
 $h_{\theta_2}(x) = [0.62, 0.34, 0.04] \in \Delta_c$   
 $h_{\theta_k}(x) = [0.25, 0.30, 0.45] \in \Delta_c$
2.  $h_{\theta_1}(x) = [0.05, 0.90, 0.05] \in \Delta_c$   
 $h_{\theta_2}(x) = [0.05, 0.90, 0.05] \in \Delta_c$   
 $h_{\theta_k}(x) = [0.05, 0.90, 0.05] \in \Delta_c$
3.  $h_{\theta_1}(x) = [1.00, 0.00, 0.00] \in \Delta_c$   
 $h_{\theta_2}(x) = [0.00, 1.00, 0.00] \in \Delta_c$   
 $h_{\theta_k}(x) = [0.00, 0.00, 1.00] \in \Delta_c$
4.  $h_{\theta_1}(x) = [0.05, 0.90, 0.05] \in \Delta_c$   
 $h_{theta_2}(x) = [0.62, 0.34, 0.04] \in \Delta_c$   
 $h_{\theta_k}(x) = [0.25, 0.30, 0.45] \in \Delta_c$   
 $h_{\theta_{k+1}}(x) = [0.83, 0.08, 0.09] \in \Delta_c$

### Desirable Properties

1.  $1 \leq m(x) \leq c$
2.  $m(x) = 1 \Rightarrow$  predictions from all models match
3.  $m(x) = c \Rightarrow$  there are models in the Rashomon Set; that assign each of the  $c$  classes
4. Monotonic in  $|\mathcal{R}(\mathcal{H}, \epsilon)|$



## Computational Challenges

1. How to choose  $\epsilon$ ?
  2. Approximating the Rashomon set without exhaustively searching?
- Using a reference model and approximating the true Rashomon set by a Rashomon subset

$$\tilde{\mathcal{R}}(\mathcal{H}, \epsilon') \triangleq \{h_{\theta_i} \in \mathcal{H}; L(h_{\theta_i}) \leq \hat{L}(h_{\theta^*}) + \epsilon'\}_{i=1}^K \subseteq \mathcal{R}(\mathcal{H}, \hat{L}(h_{\theta^*}) + \epsilon')$$

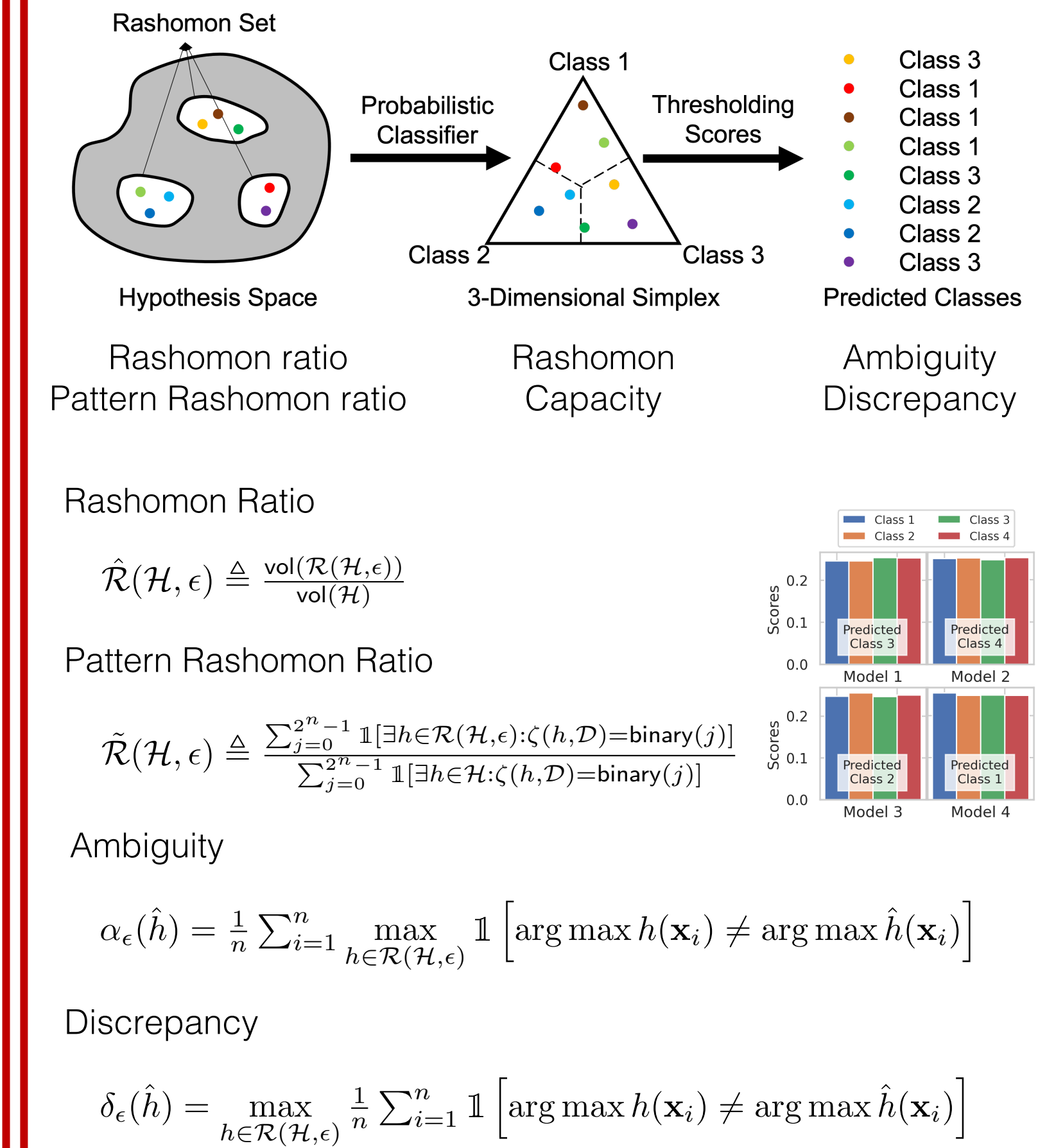
For each sample  $x$ , here are at most  $c$  models in a Rashomon subset  $\tilde{\mathcal{R}}(\mathcal{H}, \epsilon)$  whose output scores yield the same Rashomon Capacity for  $x$  as the entire Rashomon set.

### Adversarial Weight Perturbation (AWP)

$$p_k = h_{\hat{\theta}}(\mathbf{x}_i), \text{ where } \hat{\theta} = \arg \max_{\theta \in \Theta, h_\theta \in \mathcal{R}(\mathcal{H}, \epsilon)} [h_\theta(\mathbf{x}_i)]_k, \text{ for all classes } k \in [c]$$

Rashomon Capacity can be computed by the Blahut–Arimoto (BA) algorithm

## Other Metrics



## Empirical Results

