

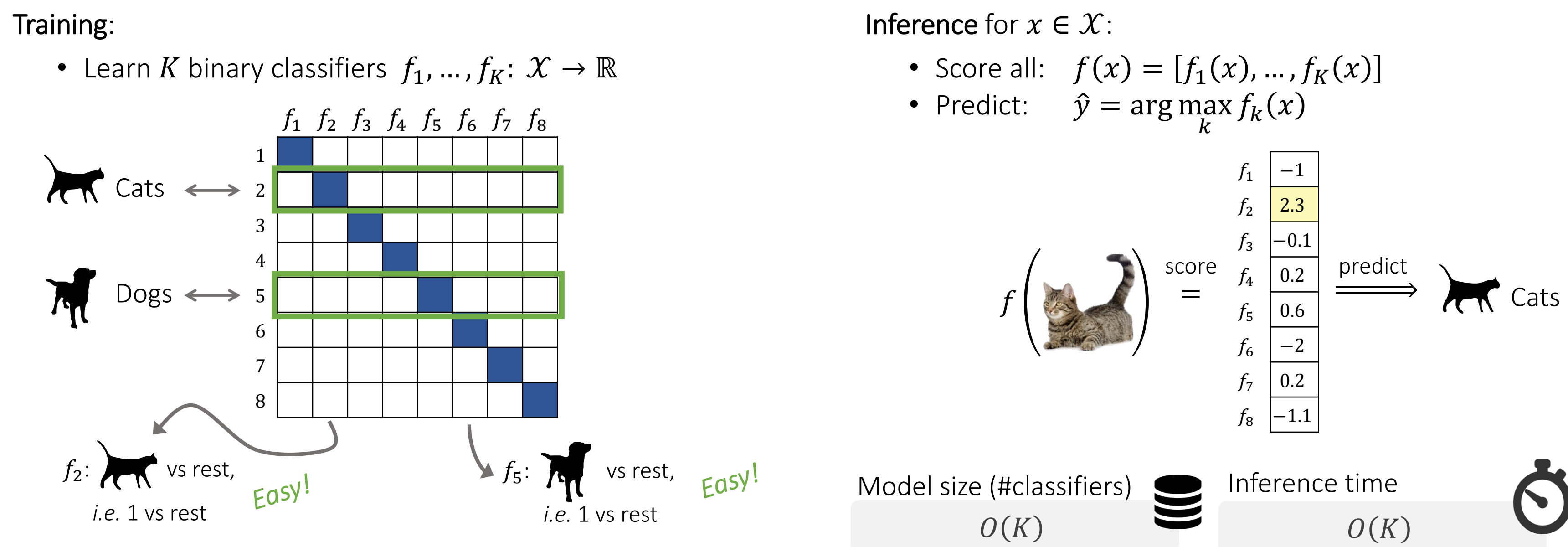
Efficient Loss-Based Decoding on Graphs for Extreme Classification

Itay Evron, Edward Moroshko, and Koby Crammer

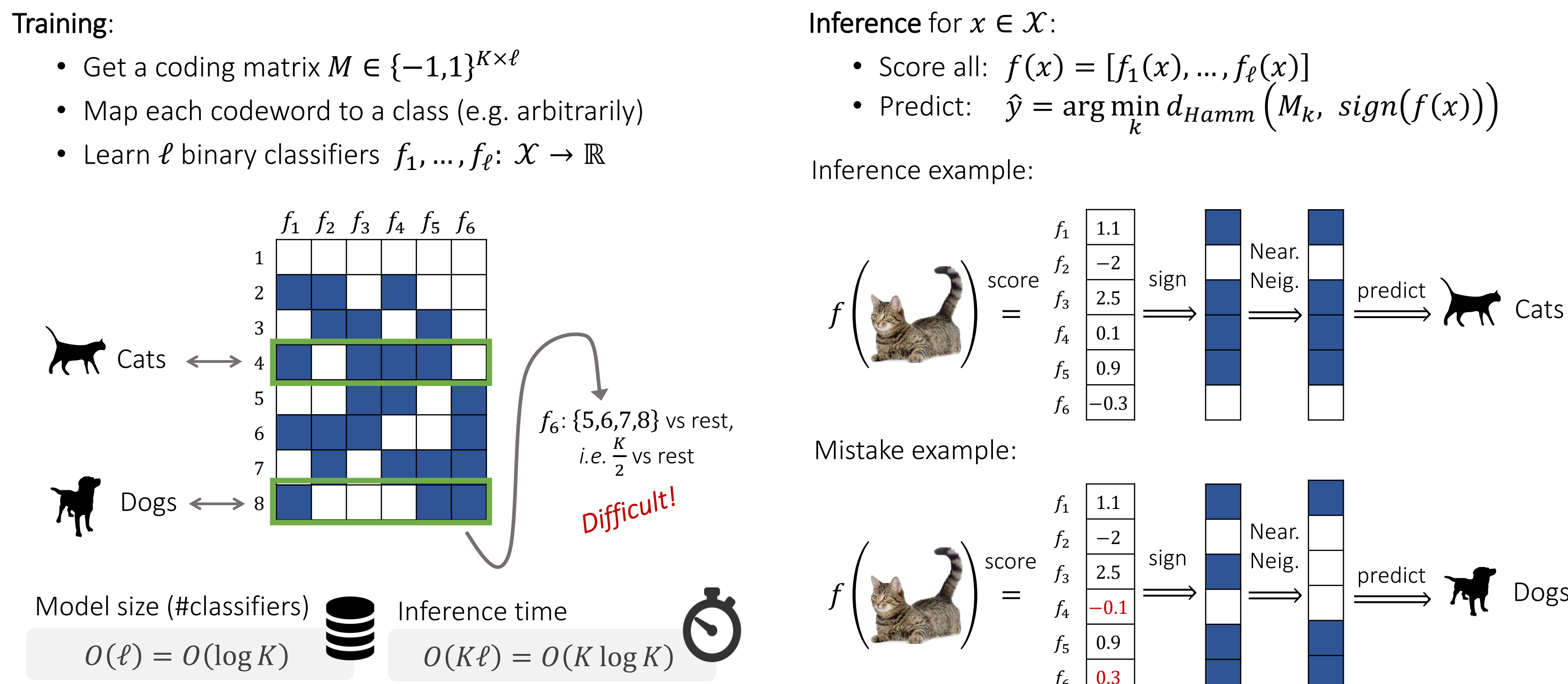
Extreme multiclass classification

- Tasks with an extremely large number of classes K .
 - Time and space complexities during training and inference become critical.
- We propose a graph-based classification scheme with time and space complexities logarithmic in K .

One vs Rest – Simple but expensive



Error Correcting Output Coding



Loss based decoding

- Instead of minimizing the Hamming distance, predict the class that minimizes some loss:

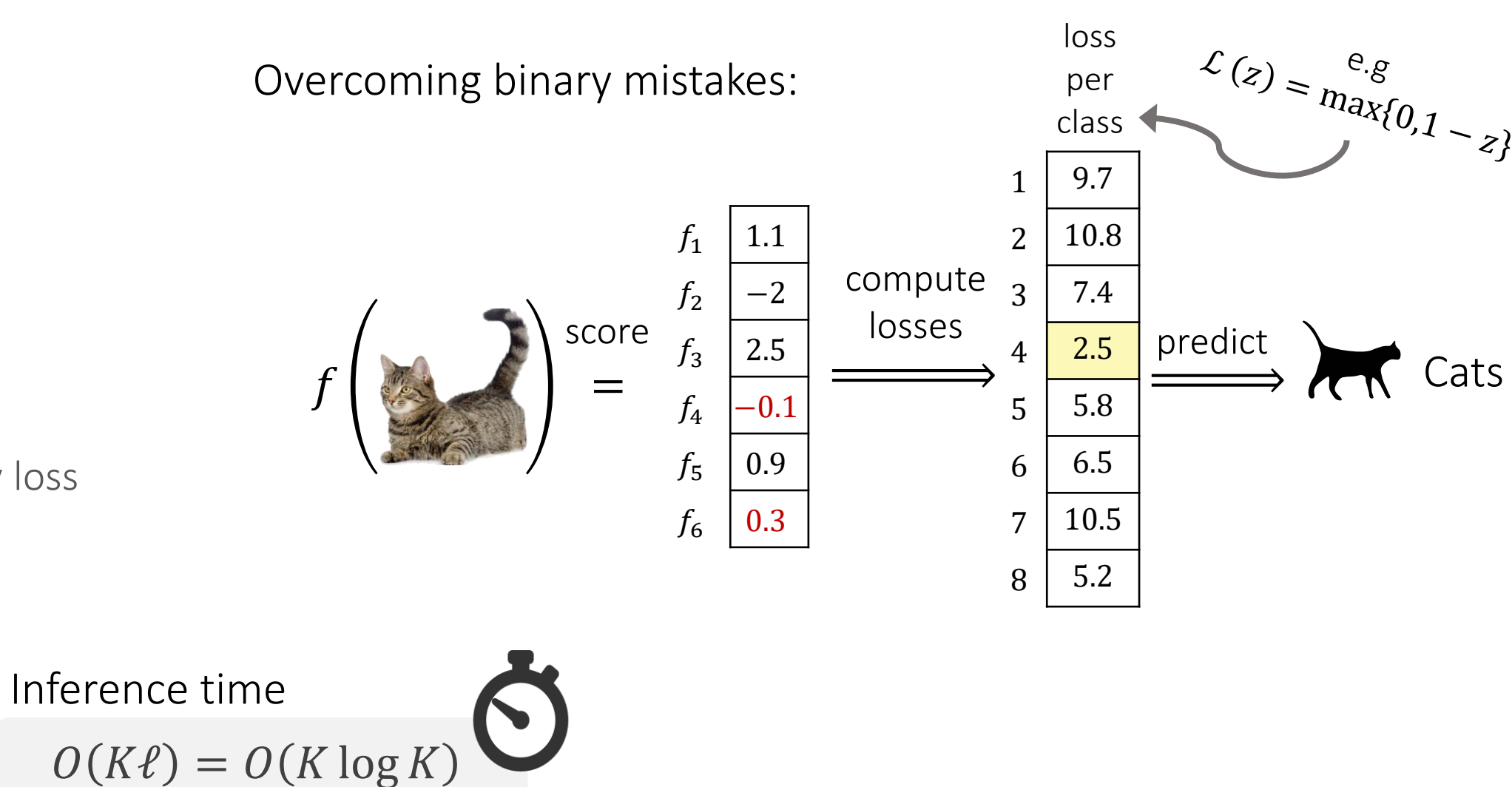
$$\hat{y} = \arg \min_k \sum_{j=1}^{\ell} \mathcal{L}(M_{k,j} \times f_j(x))$$
- An upper bound of the training multiclass error is proportional to:

$$\frac{\ell \times \varepsilon}{\rho}$$

Number of classifiers ℓ

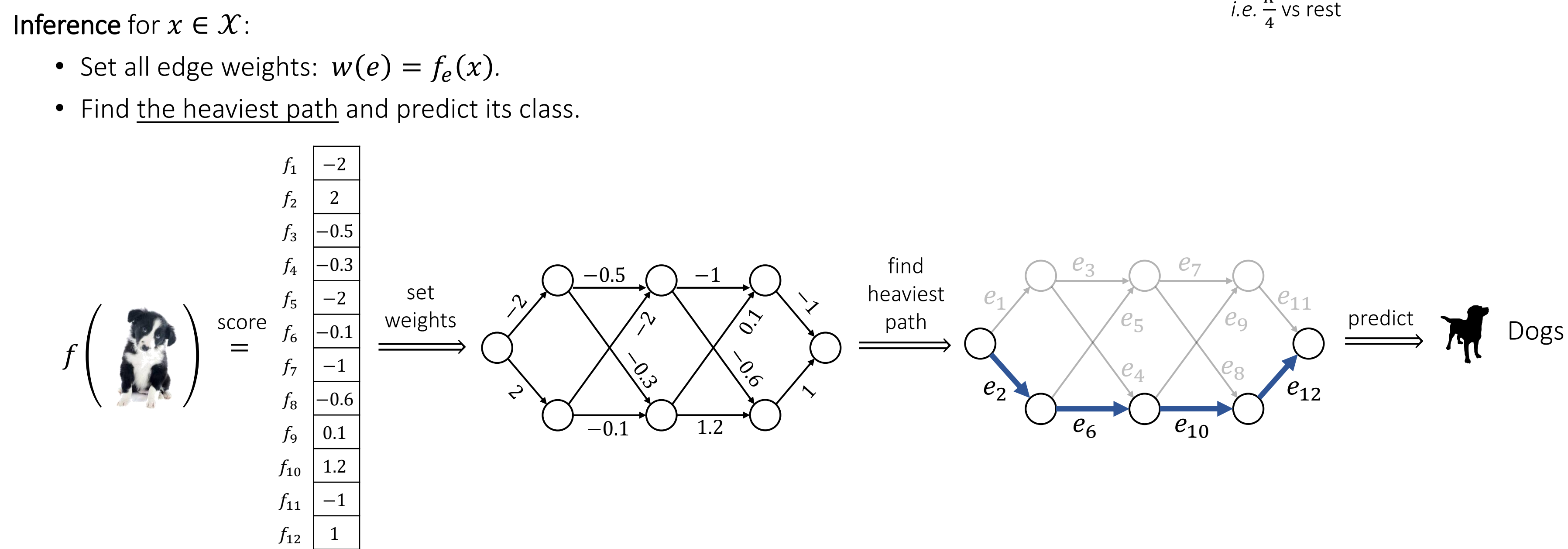
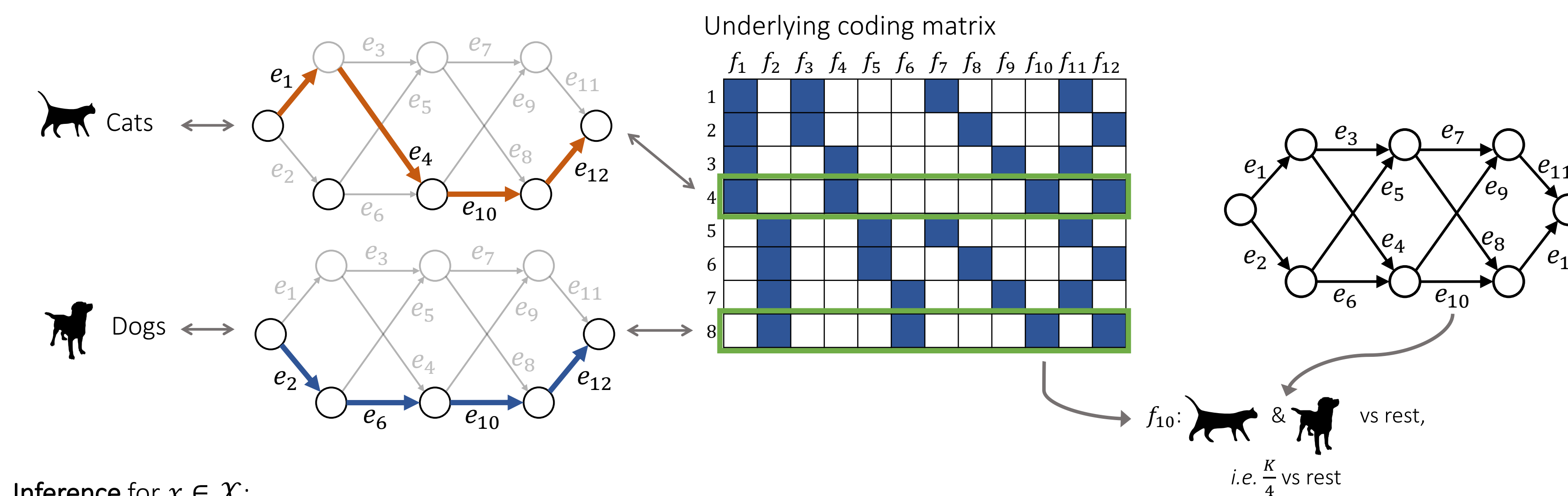
Average binary loss ε

Minimum row distance ρ
- The decoding loss function matters.



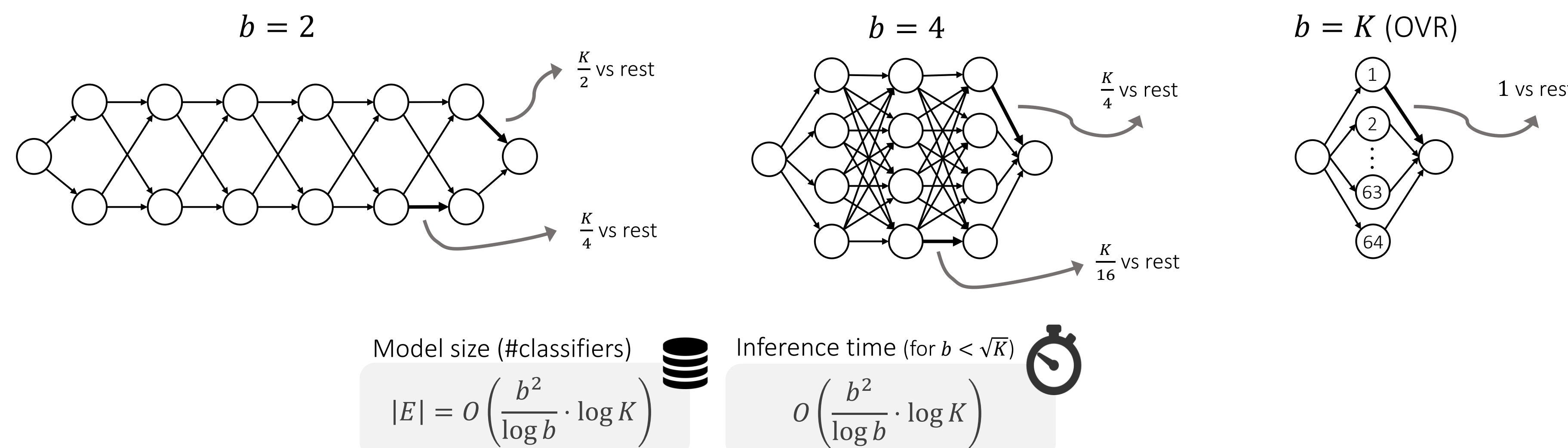
Our model – Wide-LTLS

- Based on LTLS [Jasinska and Karampatziakis 2016].
- Build a trellis graph with exactly K paths.
- Map each path to a class (e.g. arbitrarily).
- For each edge e (in parallel):
 - Train a binary classifier $f_e: \mathcal{X} \rightarrow \mathbb{R}$ on the entire training set:
 - Separate classes (=paths) that use this edge, from the classes that do not.



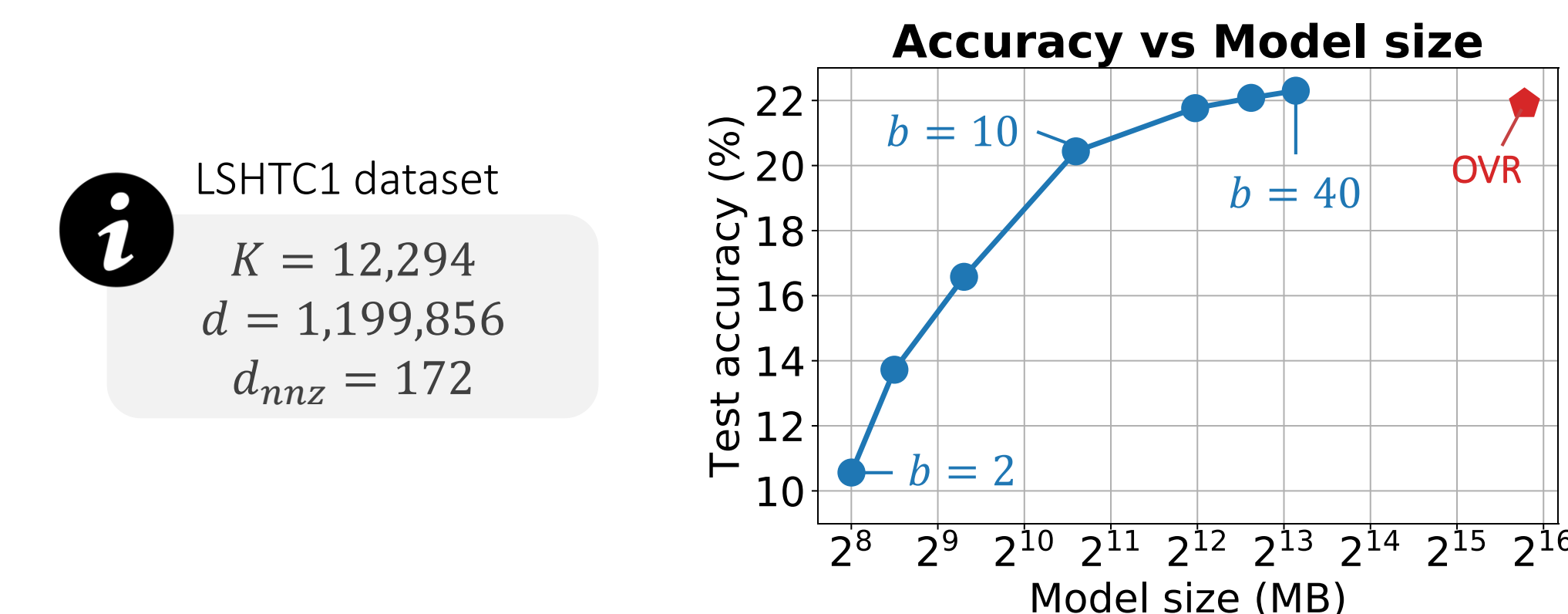
Increasing the graph width

- The same number of classes can be represented with different graph widths b .
- For instance, the following graphs all have $K = 64$ paths (=classes):



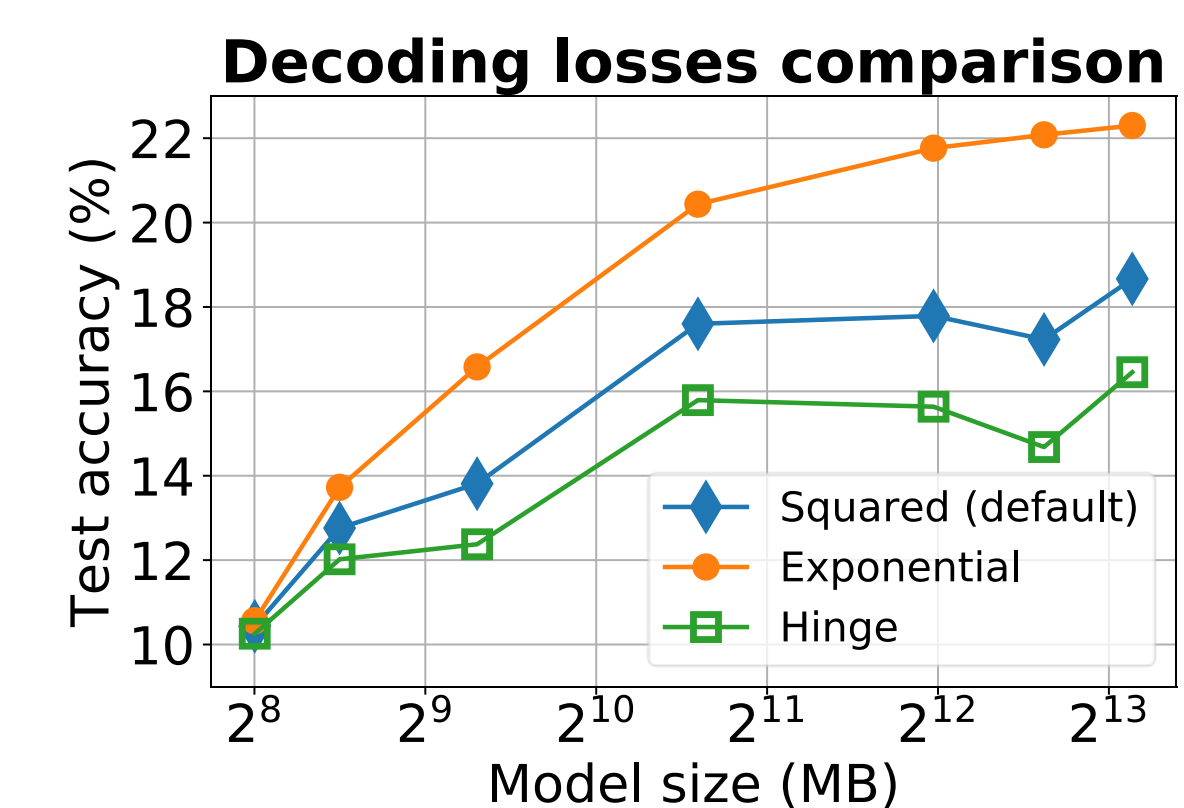
Graph width controls performance

- Our model offers a tradeoff between accuracy and model size.



W-LTLS as loss-based decoding

- We prove that W-LTLS performs loss-based decoding with the squared loss $\mathcal{L}(z) = (1-z)^2$.
- We show how to generalize W-LTLS to any loss function \mathcal{L} , and perform loss based decoding in time logarithmic in K .
- The decoding loss function matters!
 - The loss function can be chosen quickly after training.



Wider graph – Easier binary problems

- The subproblems are $\frac{K}{b^2}$ vs rest, thus get easier.

