

Estimating Accuracy from Unlabeled Data

A Bayesian Approach

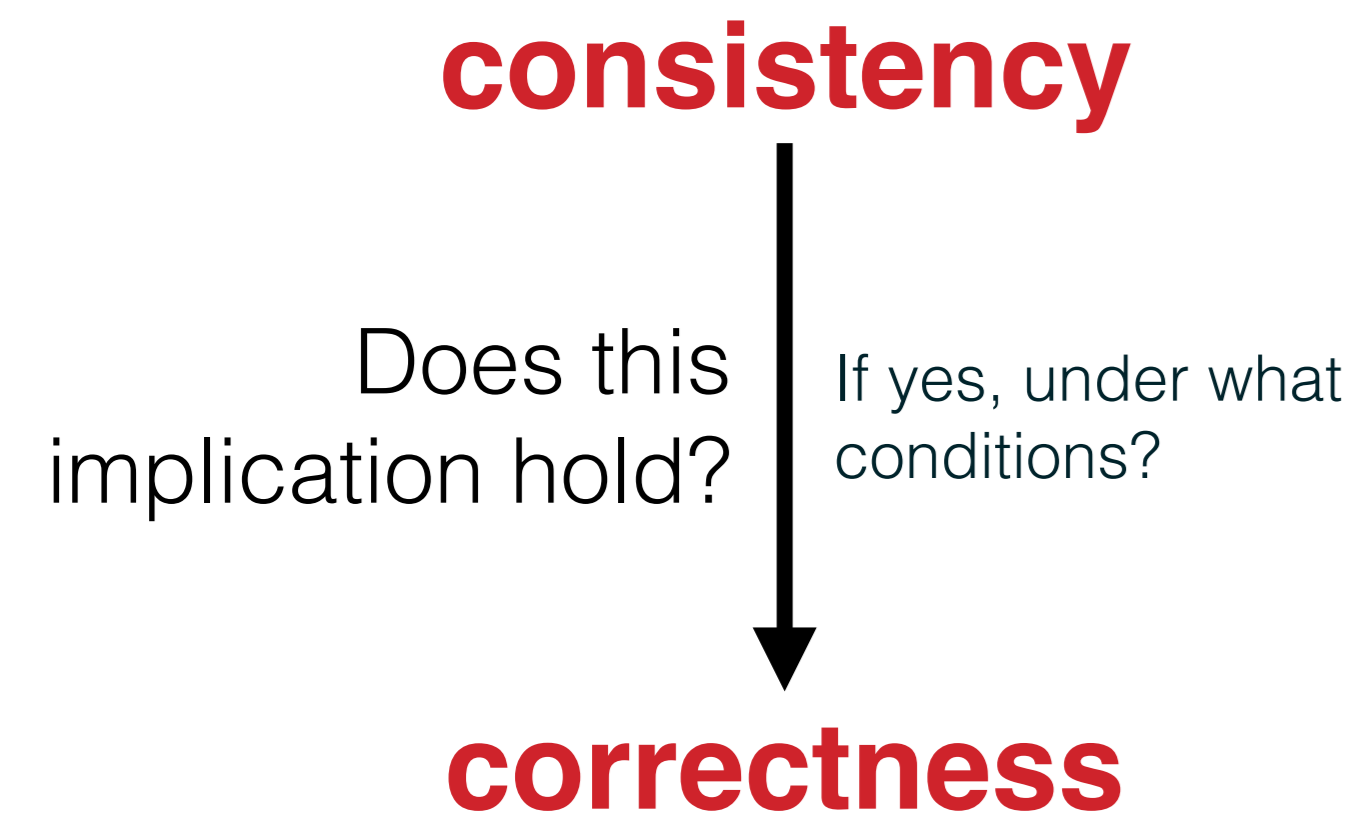
Emmanouil Antonios Platanios
e.a.platanios@cs.cmu.edu

Avinava Dubey
akdubey@cs.cmu.edu

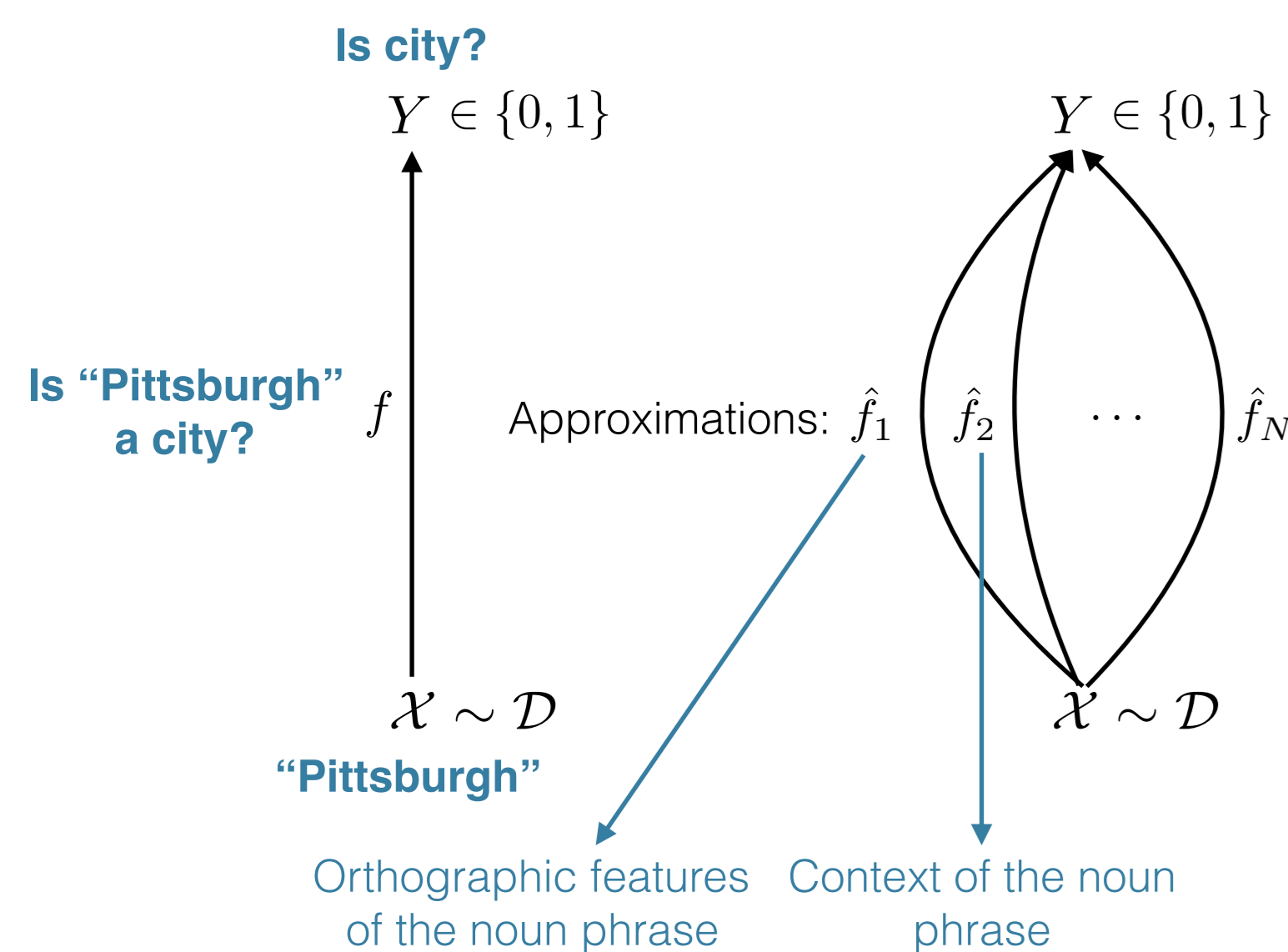
Tom Mitchell
tom.mitchell@cs.cmu.edu

1. Problem

Using **only unlabeled data** we can measure **consistency** but not **correctness**. Therefore:



There exists a **binary function** f that we do not know. Instead, we have a set of function approximations to that function and **we want to know how accurate they are**.



Consistency definition:

Given **unlabeled input data**, X_1, \dots, X_S , we observe the **sample agreement rates**:

$$\hat{a}_{\mathcal{A}} = \frac{1}{S} \sum_{s=1}^S \mathbb{I} \left\{ \hat{f}_i(X_s) = \hat{f}_j(X_s), \forall i, j \in \mathcal{A} : i \neq j \right\}$$

Correctness definition:

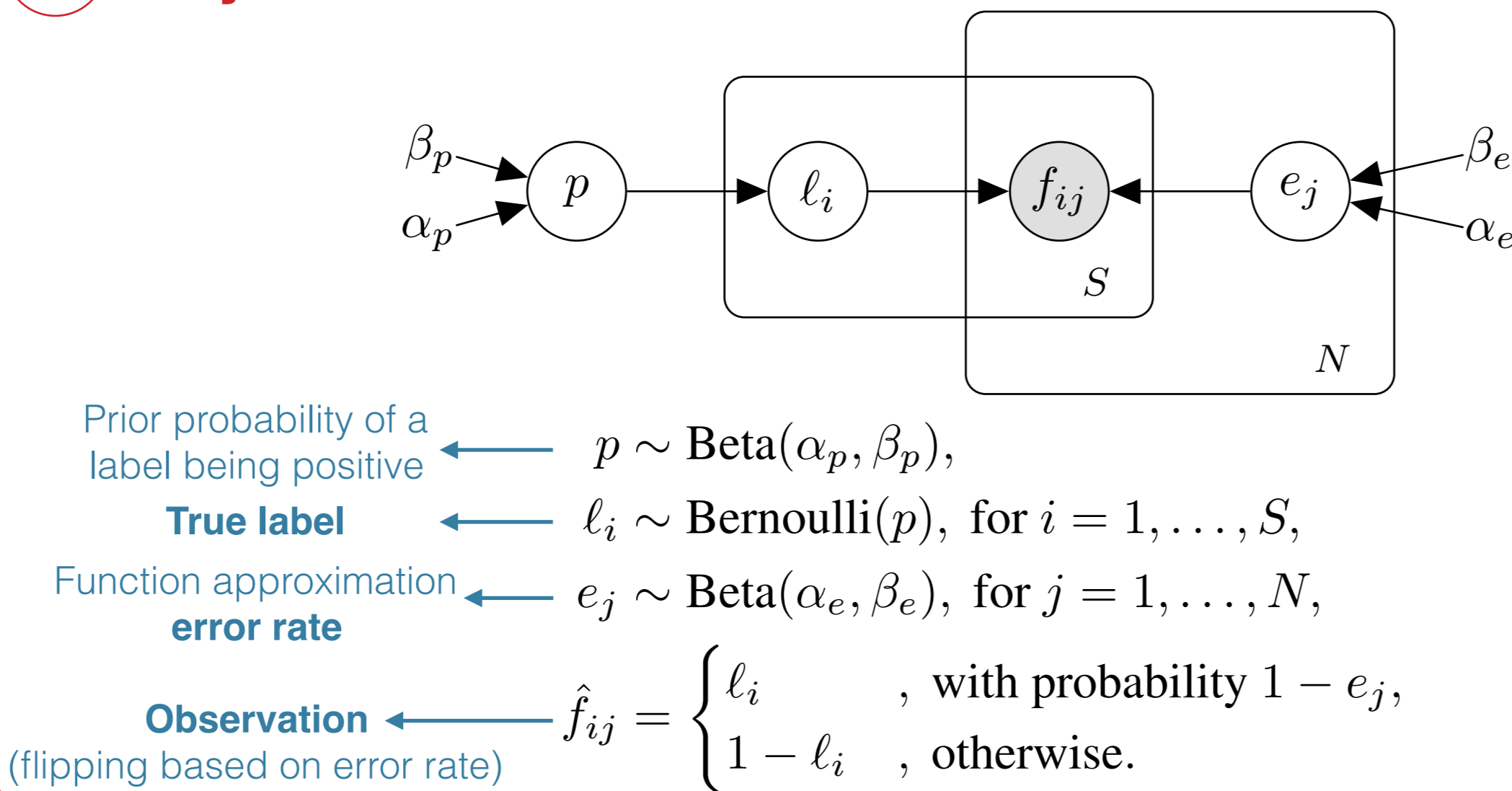
Error Rate: The probability over $\mathbb{P}(\mathcal{X}) = \mathcal{D}$ of disagreeing with the correct output label.

$$e_{\mathcal{A}} = \mathbb{P}_{\mathcal{D}} \left(\bigcap_{i \in \mathcal{A}} [\hat{f}_i(X) \neq Y] \right)$$

Error Event ← $E_{\mathcal{A}}$

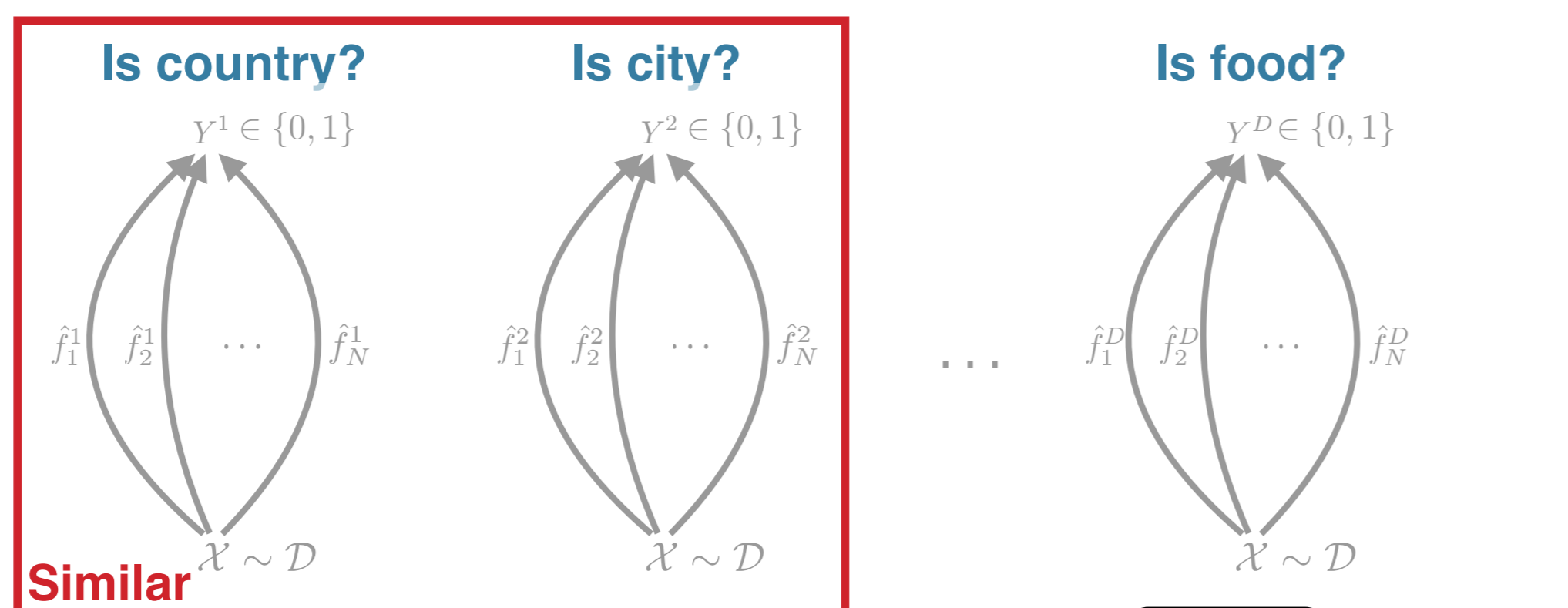
2. Approach

1 Bayesian Error Estimation

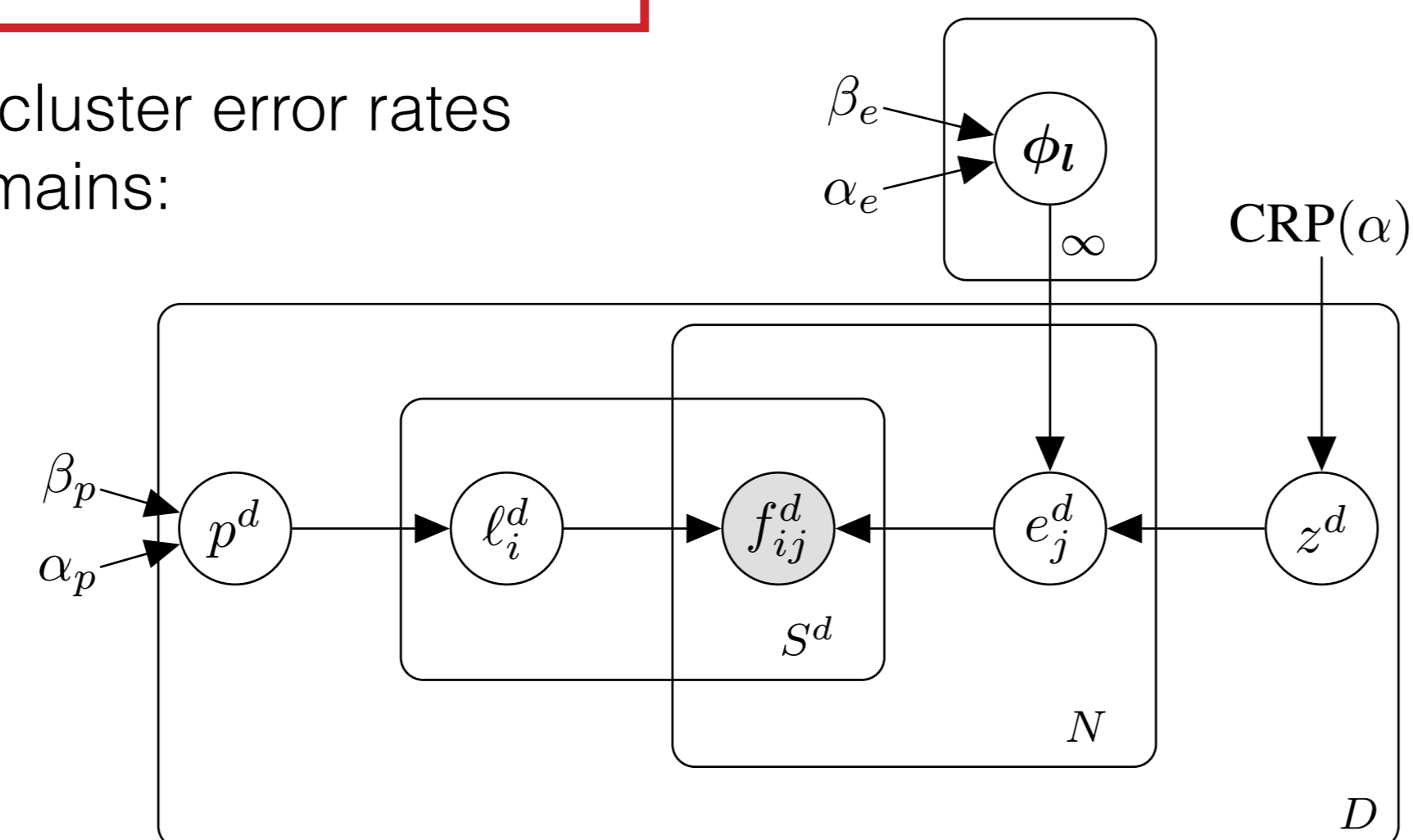


2 Coupled Bayesian Error Estimation

What about multiple domains?

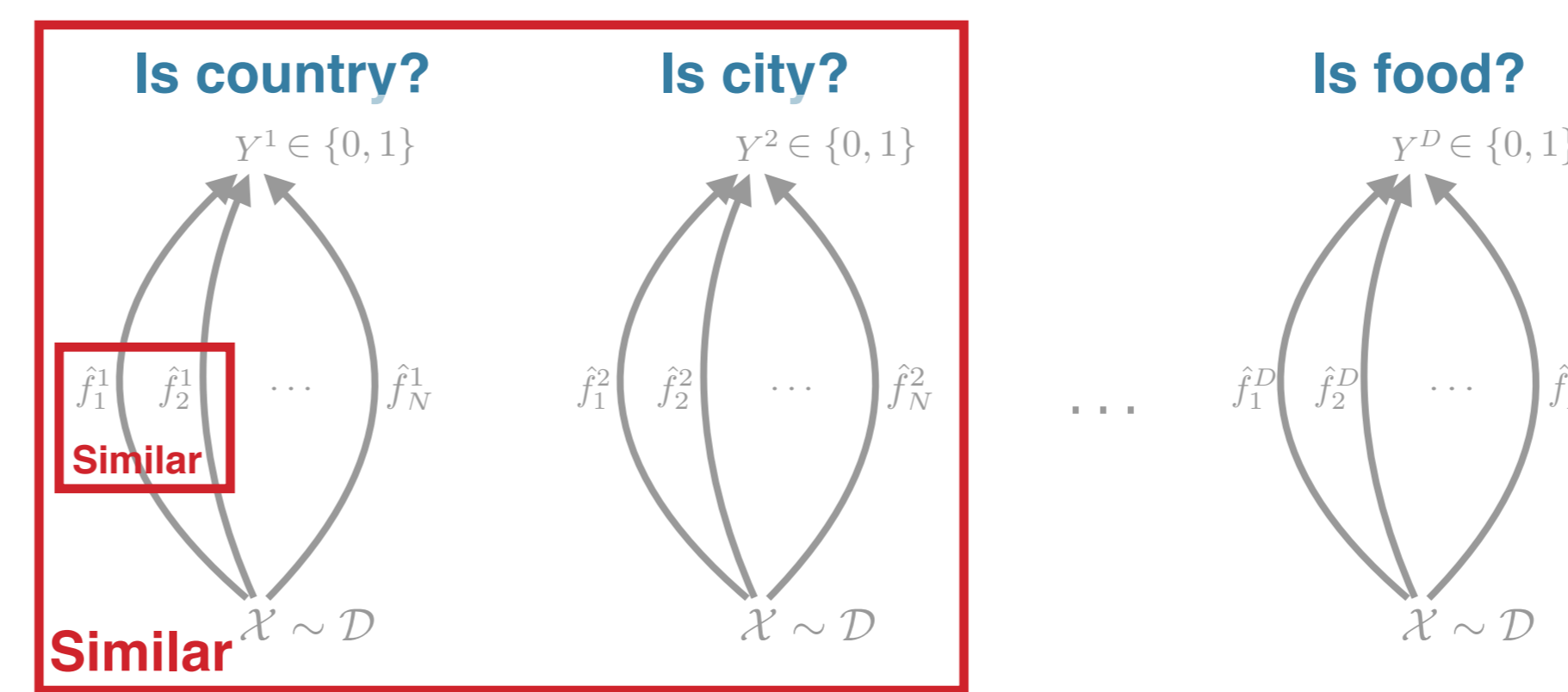


We can cluster error rates over domains:

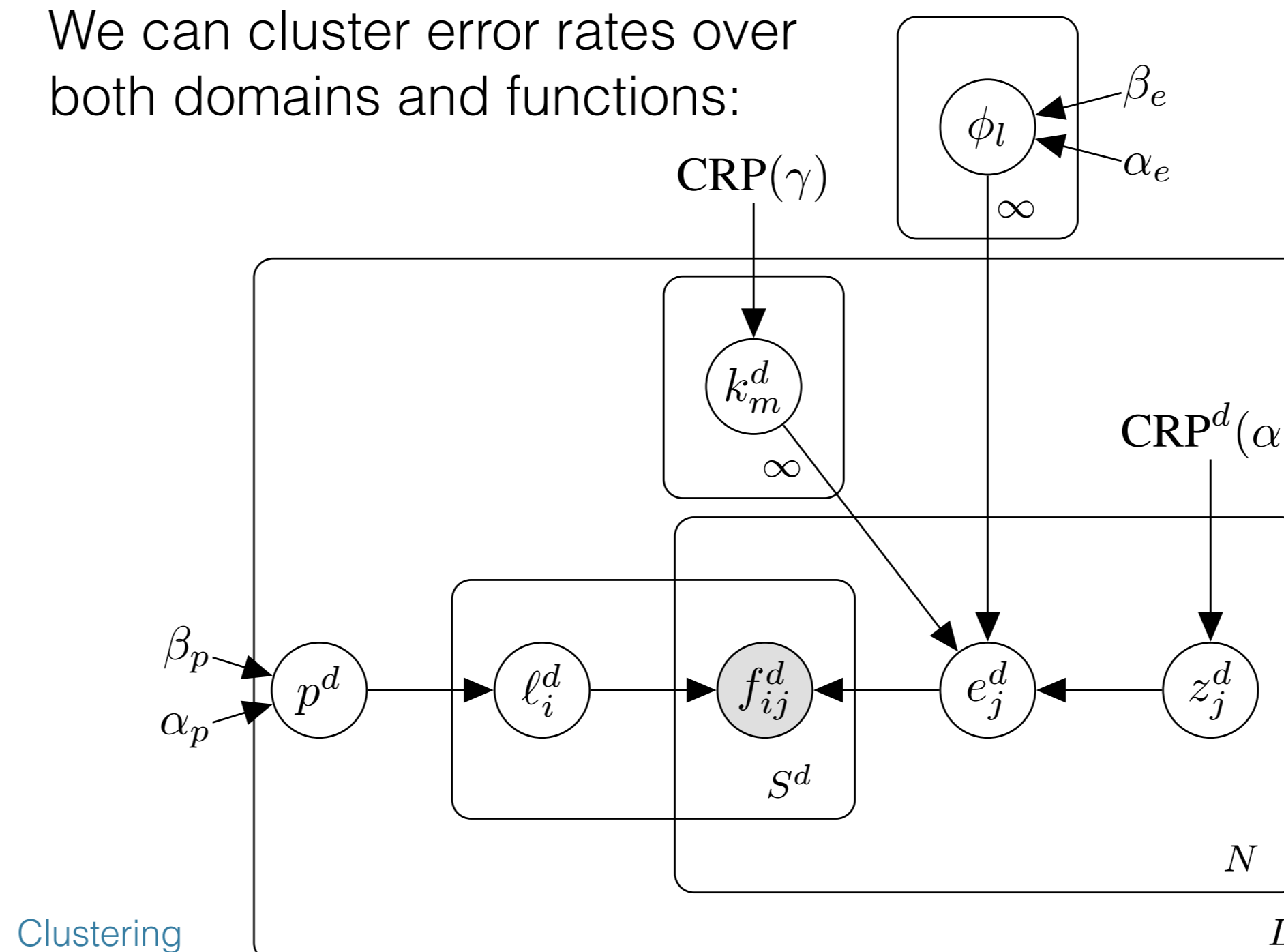


$p^d \sim \text{Beta}(\alpha_p, \beta_p)$, for $d = 1, \dots, D$,
 $l_i^d \sim \text{Bernoulli}(p^d)$, for $i = 1, \dots, S^d$, and $d = 1, \dots, D$,
 $z^d \sim \text{CRP}(\alpha)$, for $d = 1, \dots, D$,
 $e_j^d = [\phi_{z^d}]_j$, for $j = 1, \dots, N$, and $d = 1, \dots, D$,
 $\hat{f}_{ij}^d = \begin{cases} l_i^d & \text{with probability } 1 - e_j^d, \\ 1 - l_i^d & \text{otherwise.} \end{cases}$

3 Hierarchical Coupled Bayesian Error Estimation



We can cluster error rates over both domains and functions:



$p^d \sim \text{Beta}(\alpha_p, \beta_p)$, for $d = 1, \dots, D$,
 $l_i^d \sim \text{Bernoulli}(p^d)$, for $i = 1, \dots, S^d$, and $d = 1, \dots, D$,
 $\phi_l \sim \text{Beta}(\alpha_e, \beta_e)$, for $l = 1, \dots, \infty$,
 $k_m^d \sim \text{CRP}(\gamma)$, for $d = 1, \dots, D$, and $m = 1, \dots, \infty$,
 $z_j^d \sim \text{CRP}^d(\alpha)$, for $d = 1, \dots, D$, and $j = 1, \dots, N$,
 $e_j^d = \phi_{k_{z_j^d}^d}$, for $j = 1, \dots, N$, and $d = 1, \dots, D$,
 $\hat{f}_{ij}^d = \begin{cases} l_i^d & \text{with probability } 1 - e_j^d, \\ 1 - l_i^d & \text{otherwise.} \end{cases}$

Our methods implicitly use agreement rates in order to estimate function error rates. We are using the agreement between the function outputs and the true underlying labels in order to infer both the error rates of our functions and those labels, jointly.

3. Experiments

We report the **error mean squared error (MSE_{error})** between:

- True error rates (estimated from labeled data)
- Error rates estimates from unlabeled data

the **target label mean absolute deviation (MAD_{label})** — equivalent to **accuracy**.

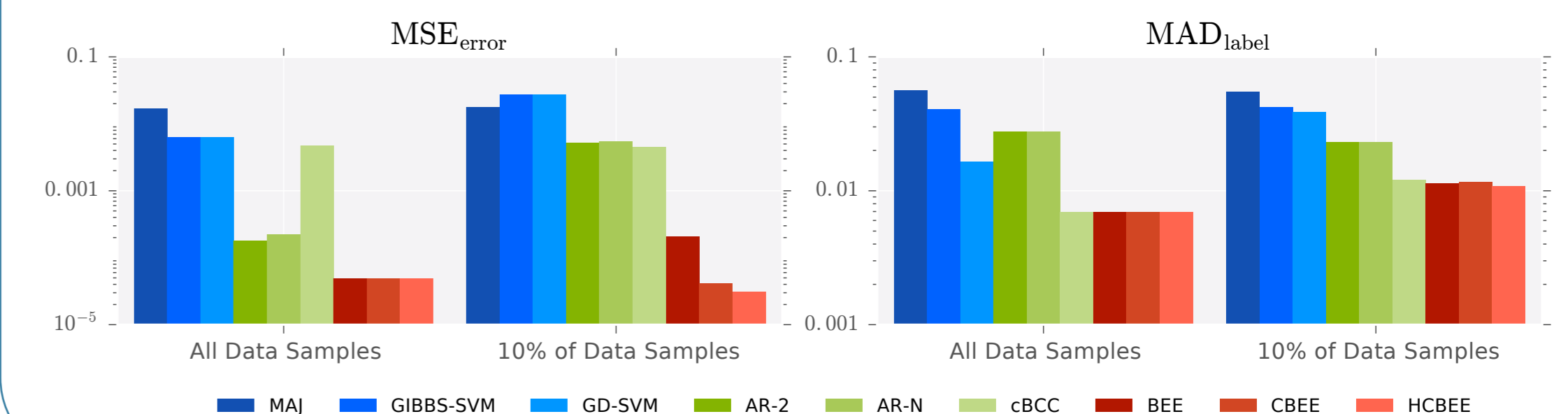
1 NELL Data Set

Task: Predict whether a noun phrase (NP) belongs to a category (e.g. "city")

4 logistic regression classifiers using different features:

- ADJ**: Adjectives that occur with the NP
- CMC**: Orthographic features of the NP
- CPL**: Phrases that occur with the NP
- VERB**: Verbs that appear with the NP

Category	# Examples
animal	20,733
beverage	18,932
bird	19,263
bodypart	21,840
city	21,778
disease	21,827
drug	20,452
fish	19,162
food	19,566
fruit	18,911
muscle	21,606
person	21,700
protein	21,811
river	21,723
vegetable	18,826



2 Brain Data Set

11 classifiers using different text representation:

- Number of letters in each word:
- Parts of speech:
- ...

Task: Find which of two 40 second long story passages corresponds to an unlabeled 40 second time series of fMRI neural activity

Passage #1	Passage #2
... They were hoping for a reason to fight Malfoy Harry had heard Fred and George Weasley complain ...
... 4 4 6 3 1 6 2 5 6 5 3 5 4 3 6 7 8 ...
... PRP VBD VBG IN DT NNP VBD VBN NNP NN TO VB NNP ...
... CC NNP NNP VB

Which passage corresponds to this fMRI recording?

1,000 labeled samples for 11 brain regions

