

# Interpretable Neural Predictions with Differentiable Binary Variables

Joost Basting   Wilker Aziz   Ivan Titov

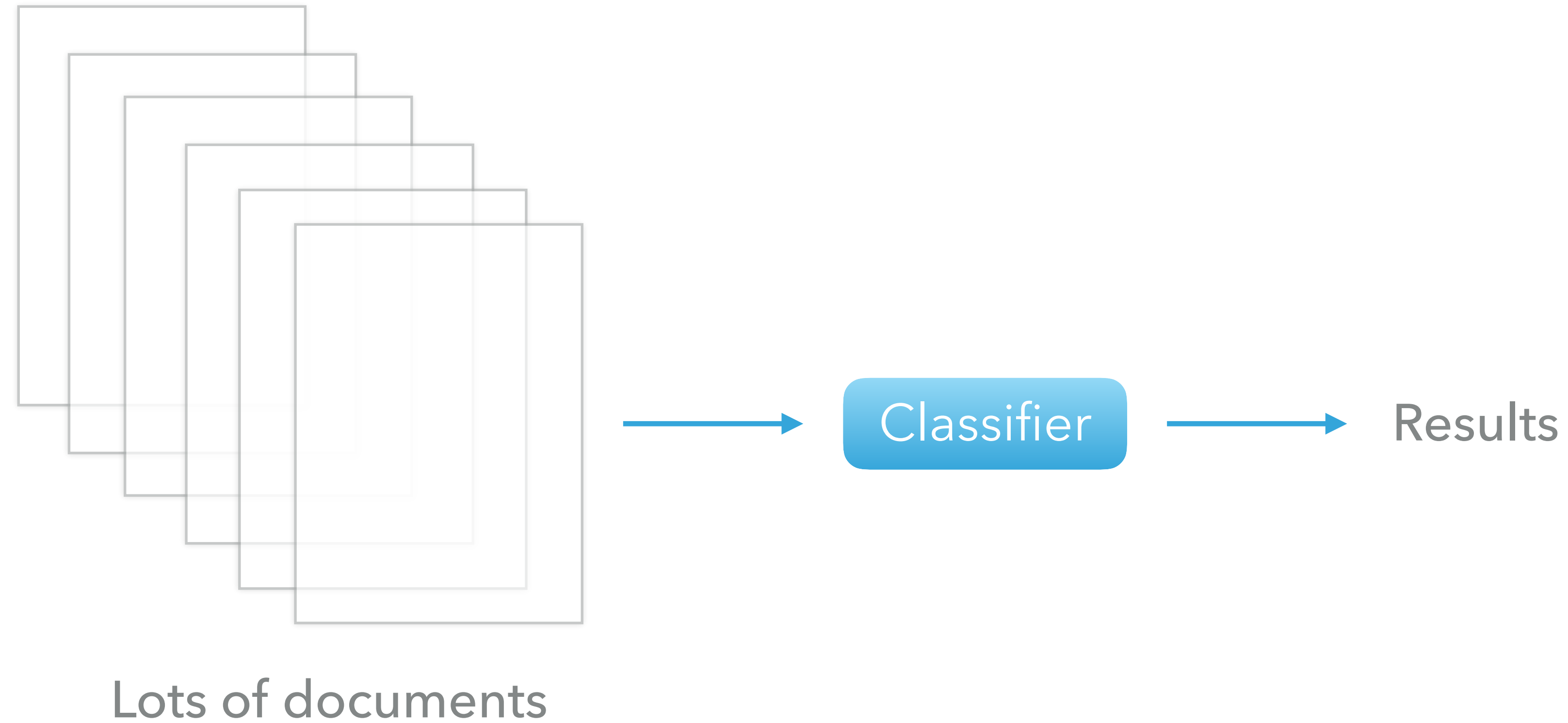
ILLC, University of Amsterdam  
ILCC, University of Edinburgh

[basting.github.io](https://basting.github.io)

ACL, 30 July 2019

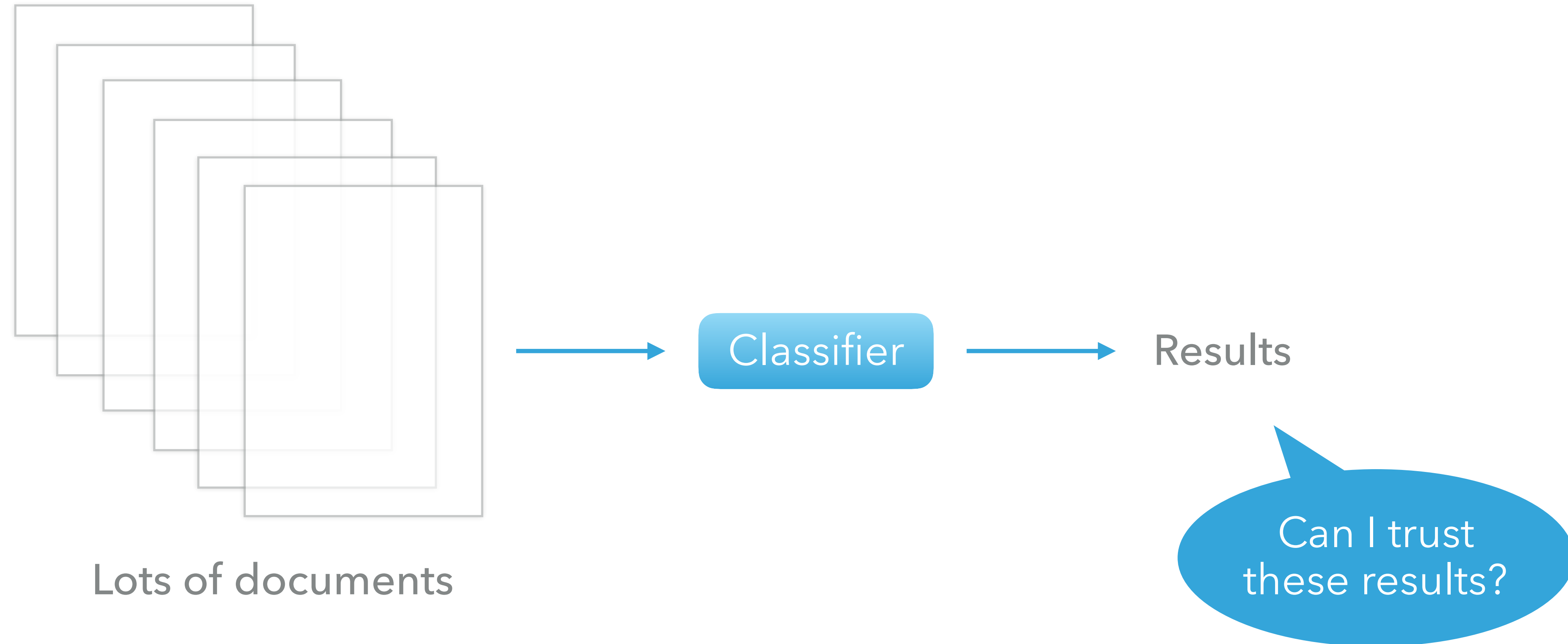
# Interpreting neural networks is difficult

2



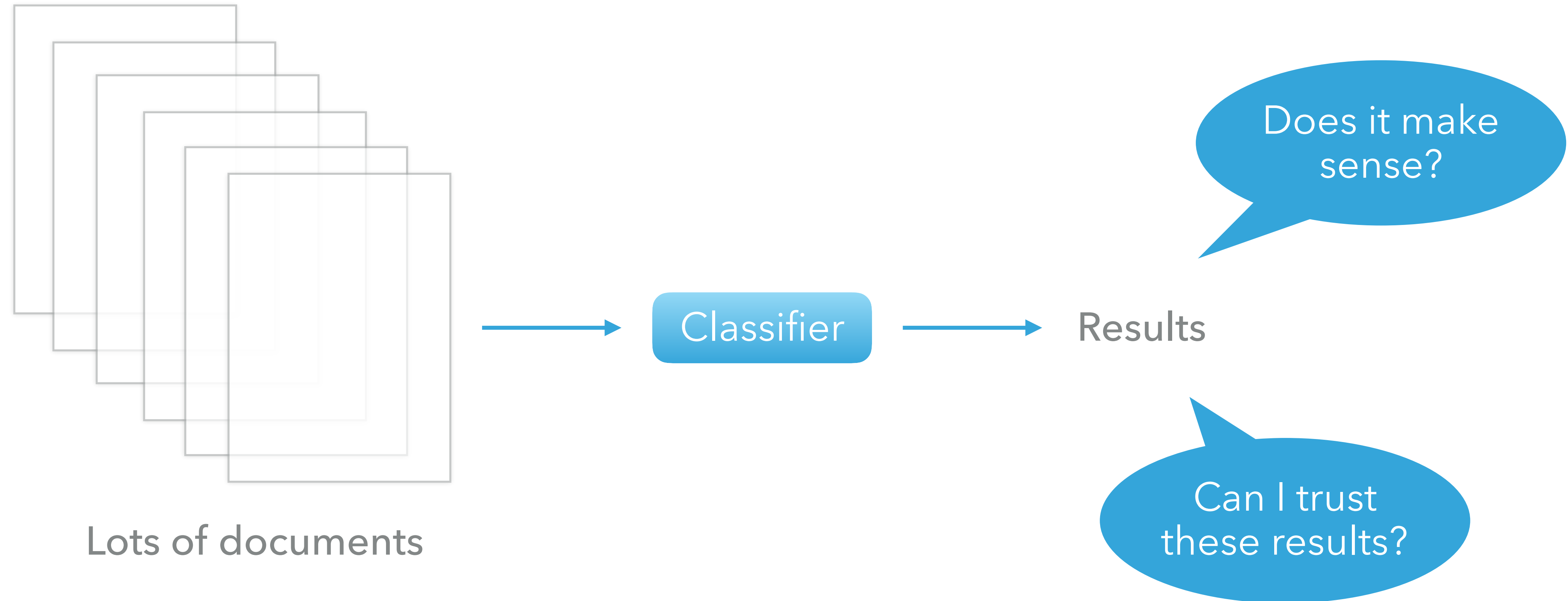
# Interpreting neural networks is difficult

2



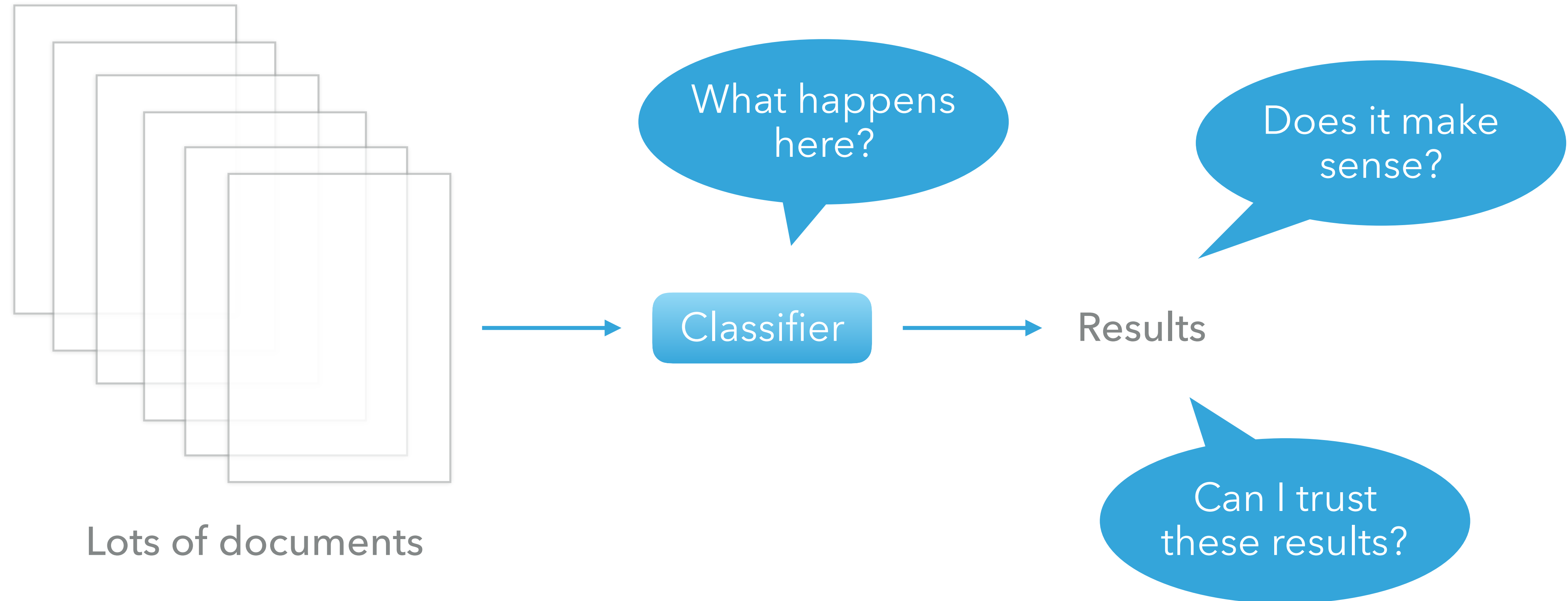
# Interpreting neural networks is difficult

2



# Interpreting neural networks is difficult

2



pours a dark amber color with decent head that does not recede much . it 's a tad too dark to see the carbonation , but fairs well . smells of roasted malts and mouthfeel is quite strong in the sense that you can get a good taste of it before you even swallow .



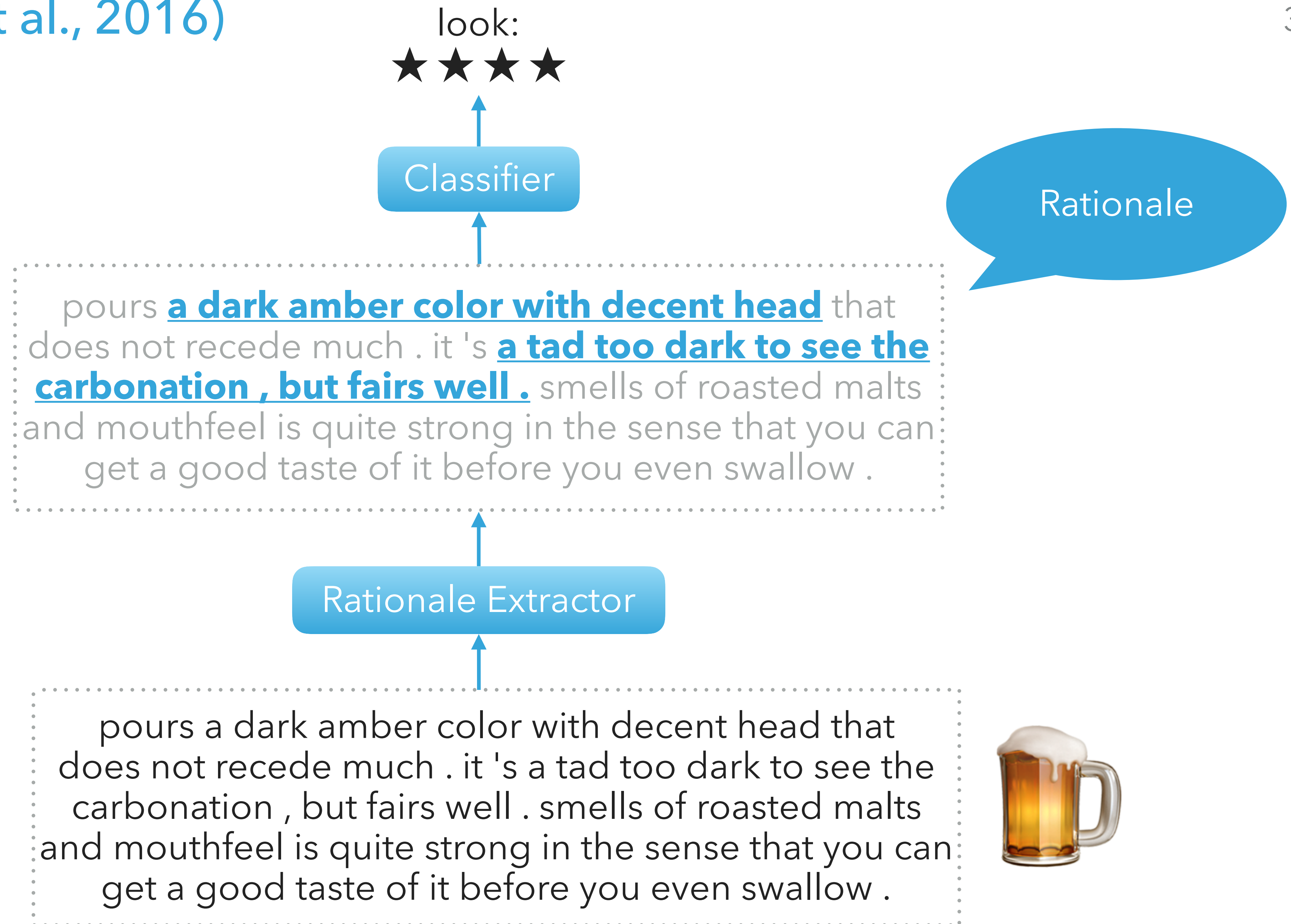
pours **a dark amber color with decent head** that does not recede much . it 's **a tad too dark to see the carbonation , but fairs well .** smells of roasted malts and mouthfeel is quite strong in the sense that you can get a good taste of it before you even swallow .

Rationale Extractor

pours a dark amber color with decent head that does not recede much . it 's a tad too dark to see the carbonation , but fairs well . smells of roasted malts and mouthfeel is quite strong in the sense that you can get a good taste of it before you even swallow .









A green speech bubble with a tail pointing towards the center of the slide.

Enough to make  
the right prediction

A rationale is a **short** and **sufficient** part of the input text.

A blue speech bubble with a tail pointing towards the center of the slide.

To be a good  
explanation

$$Z_i \mid x \sim \text{Bernoulli}(g_i(x; \phi))$$

Rationale Extractor

NN that predicts a  
sequence  $n$  of Bernoulli  
parameters

Classifier

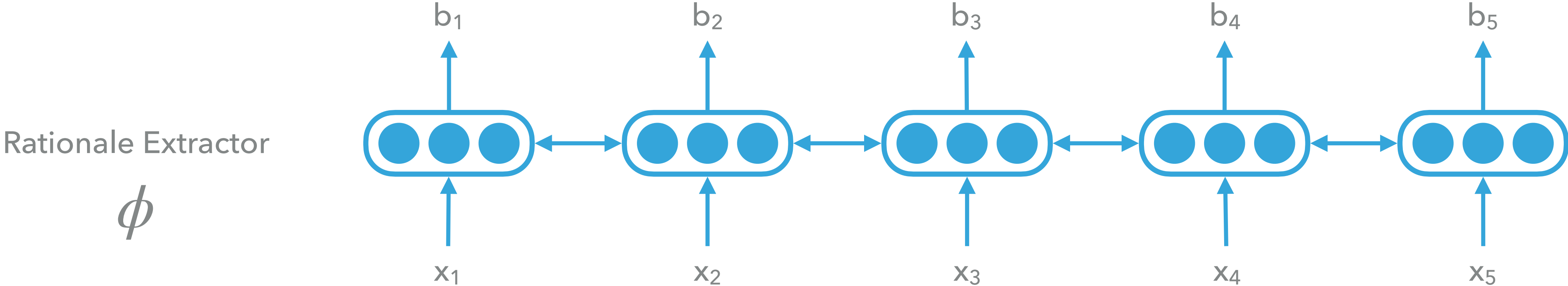
$z_i \in \{0,1\}$  can  
erase  $x_i$

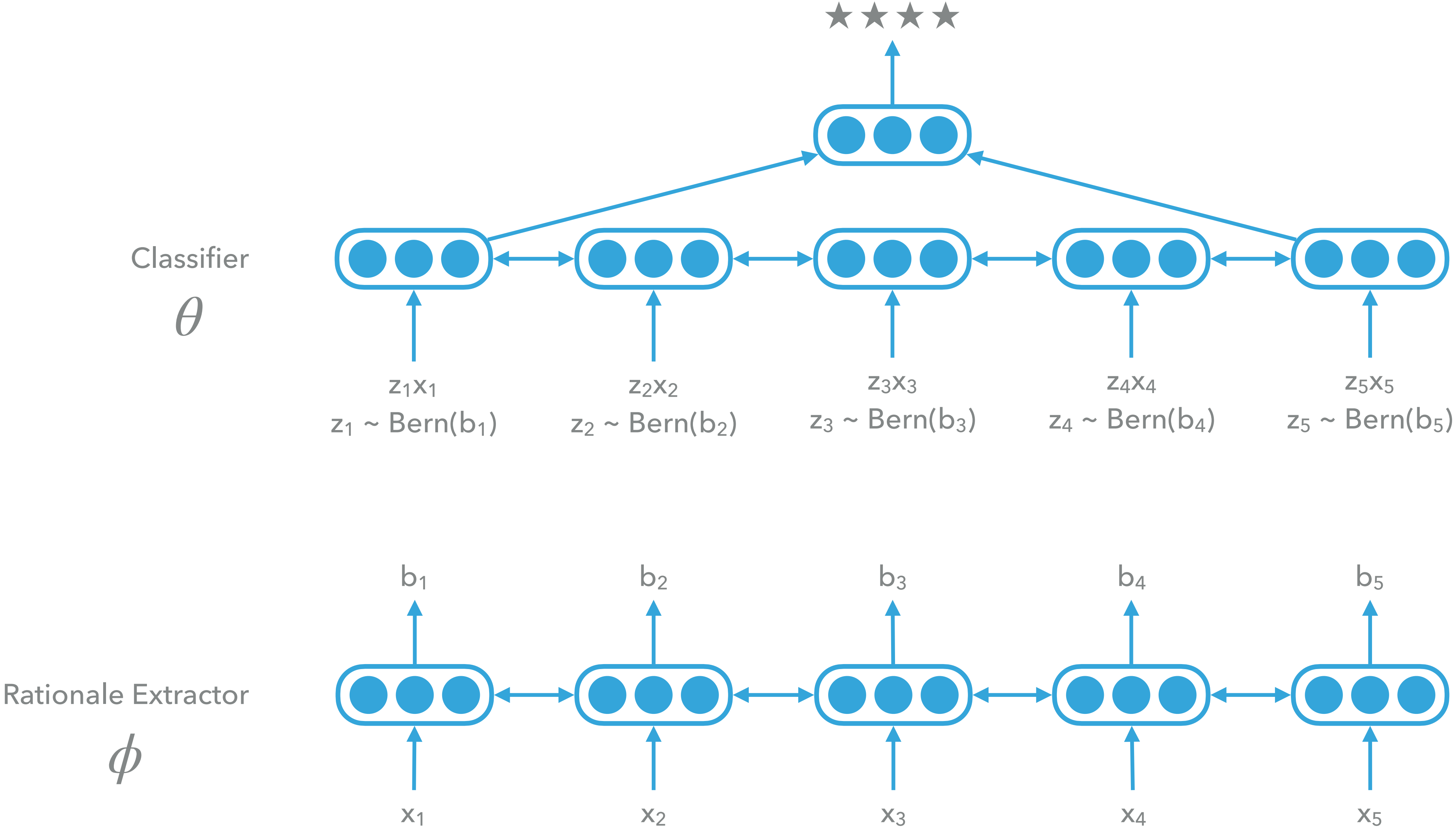
$$Y \mid x, z \sim \text{Cat}(f(x \odot z; \theta))$$

$$Z_i \mid x \sim \text{Bernoulli}(g_i(x; \phi))$$

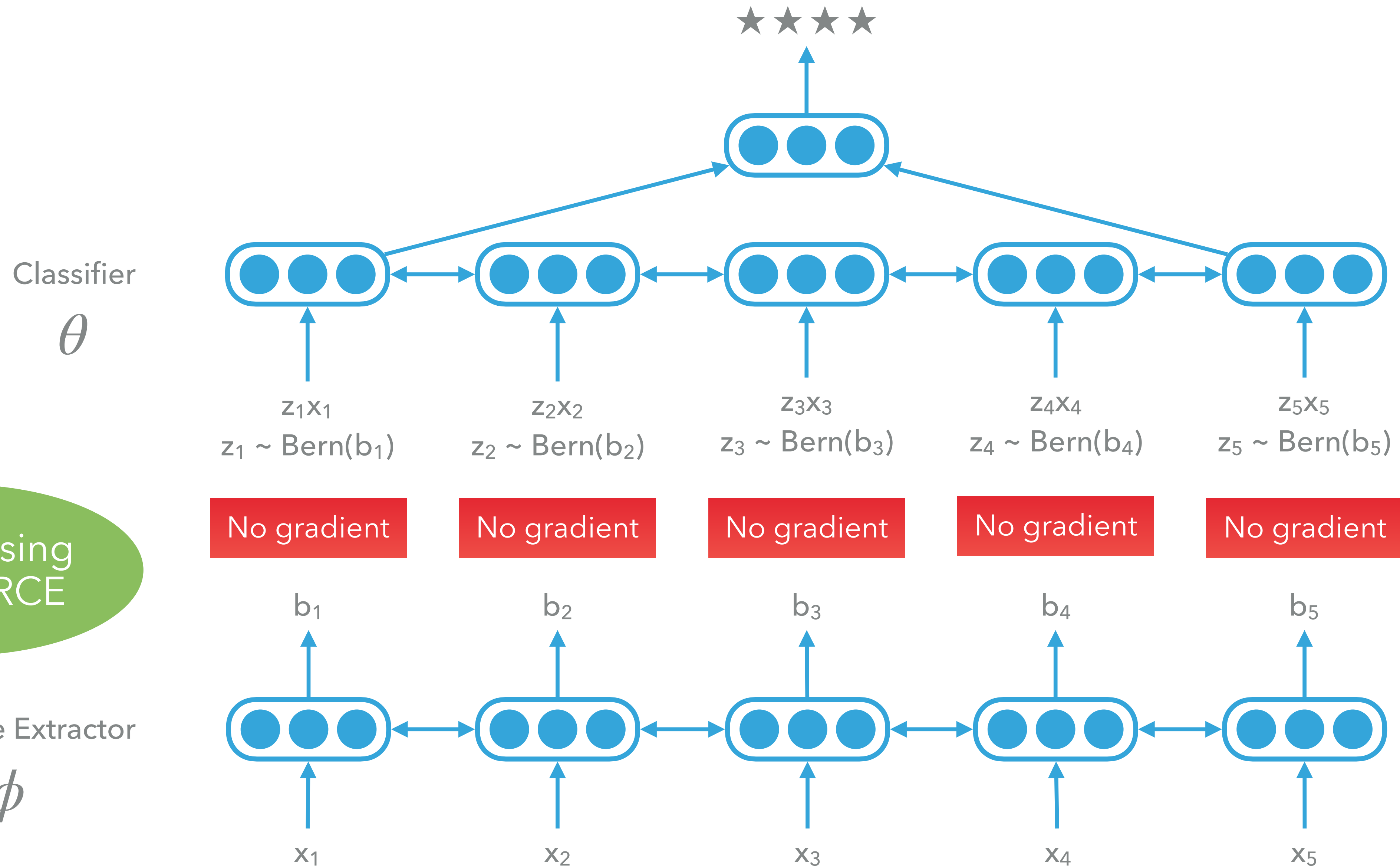
Rationale Extractor

NN that predicts a  
sequence  $n$  of Bernoulli  
parameters





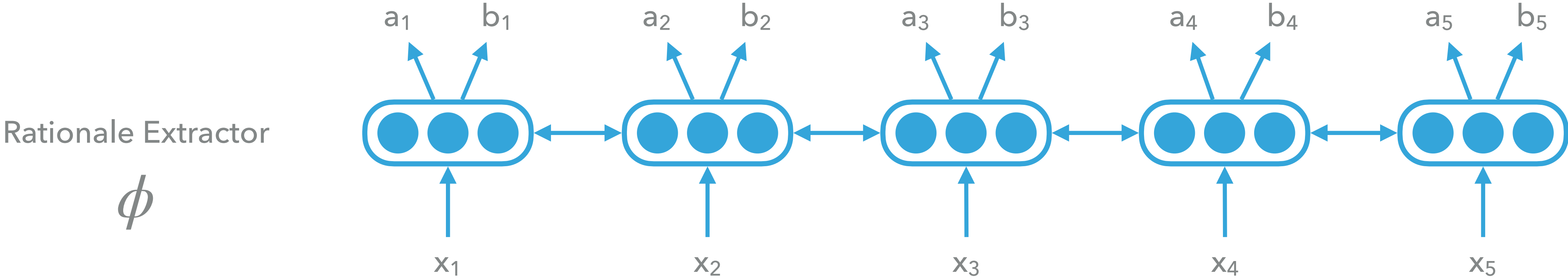
# Model of Lei et al. (2016)



trained using  
REINFORCE

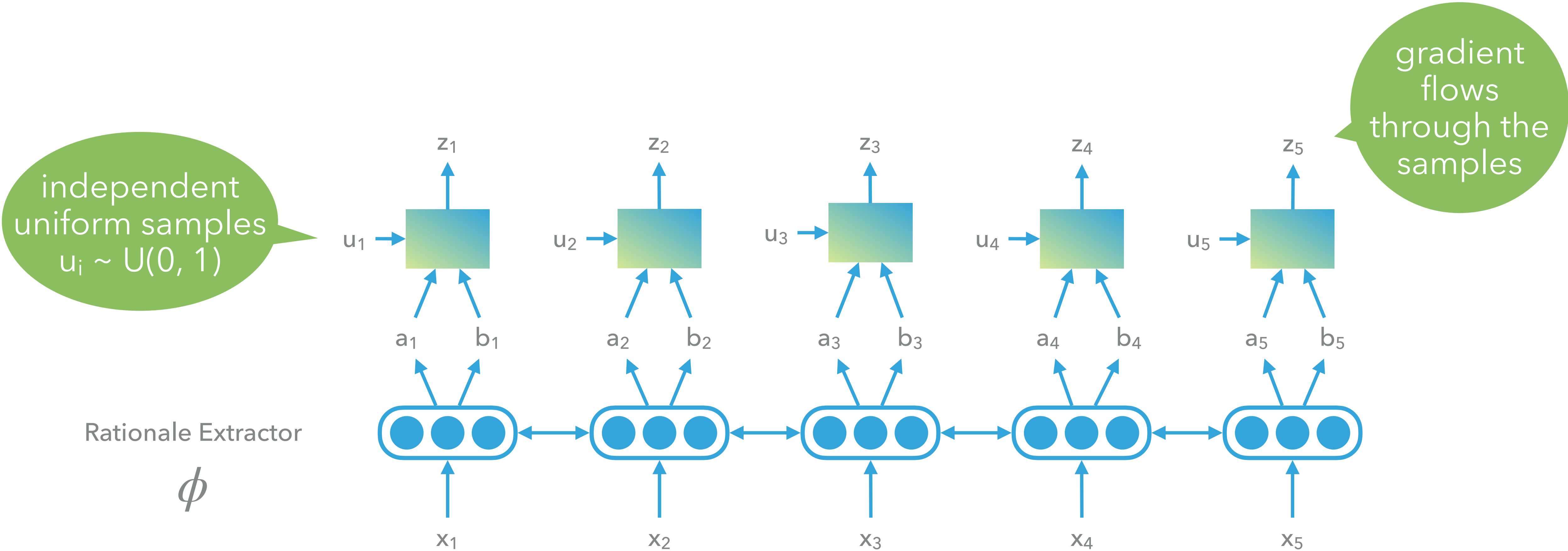
## Rationale Extractor

# Our Proposed Model

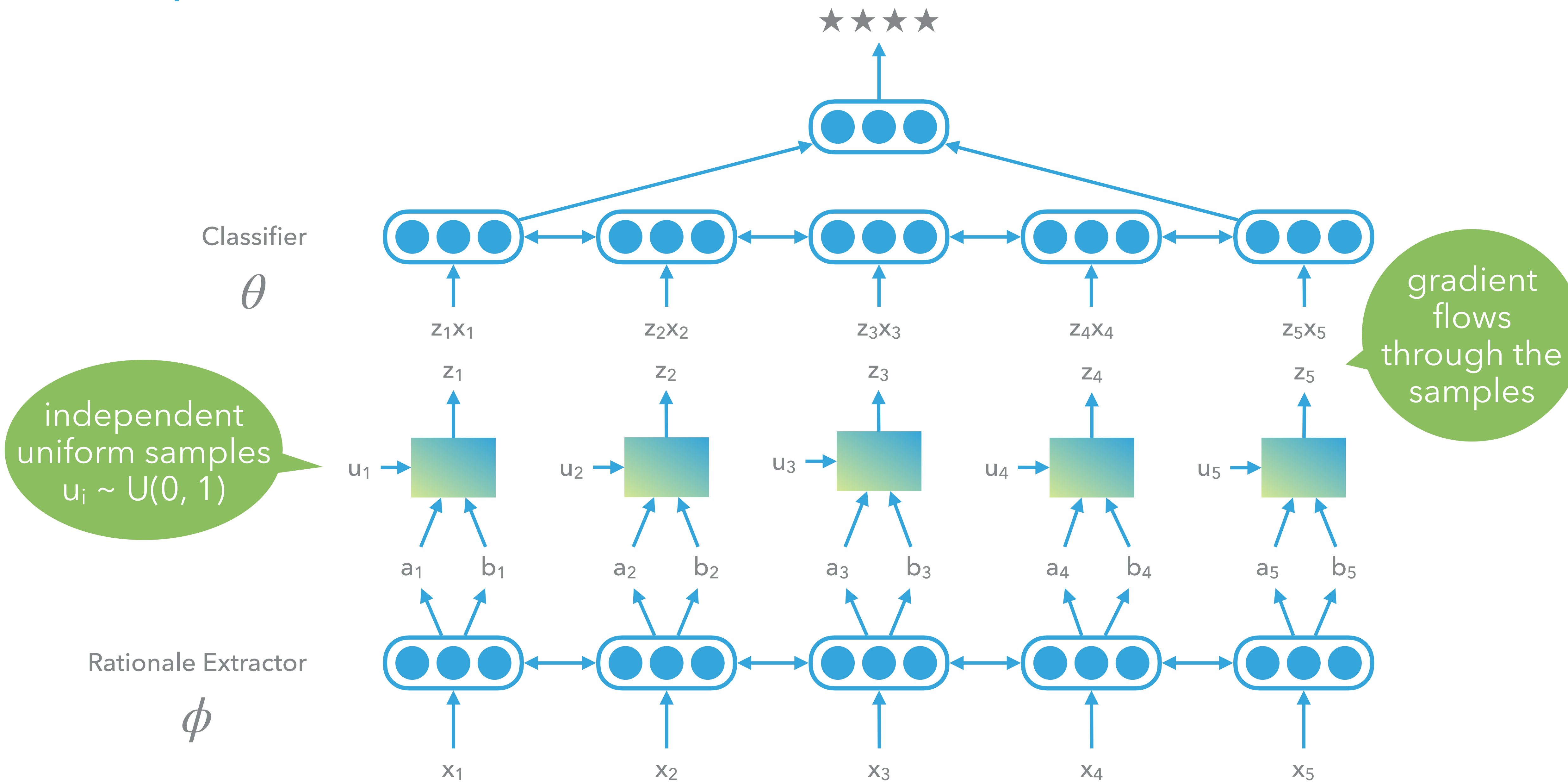




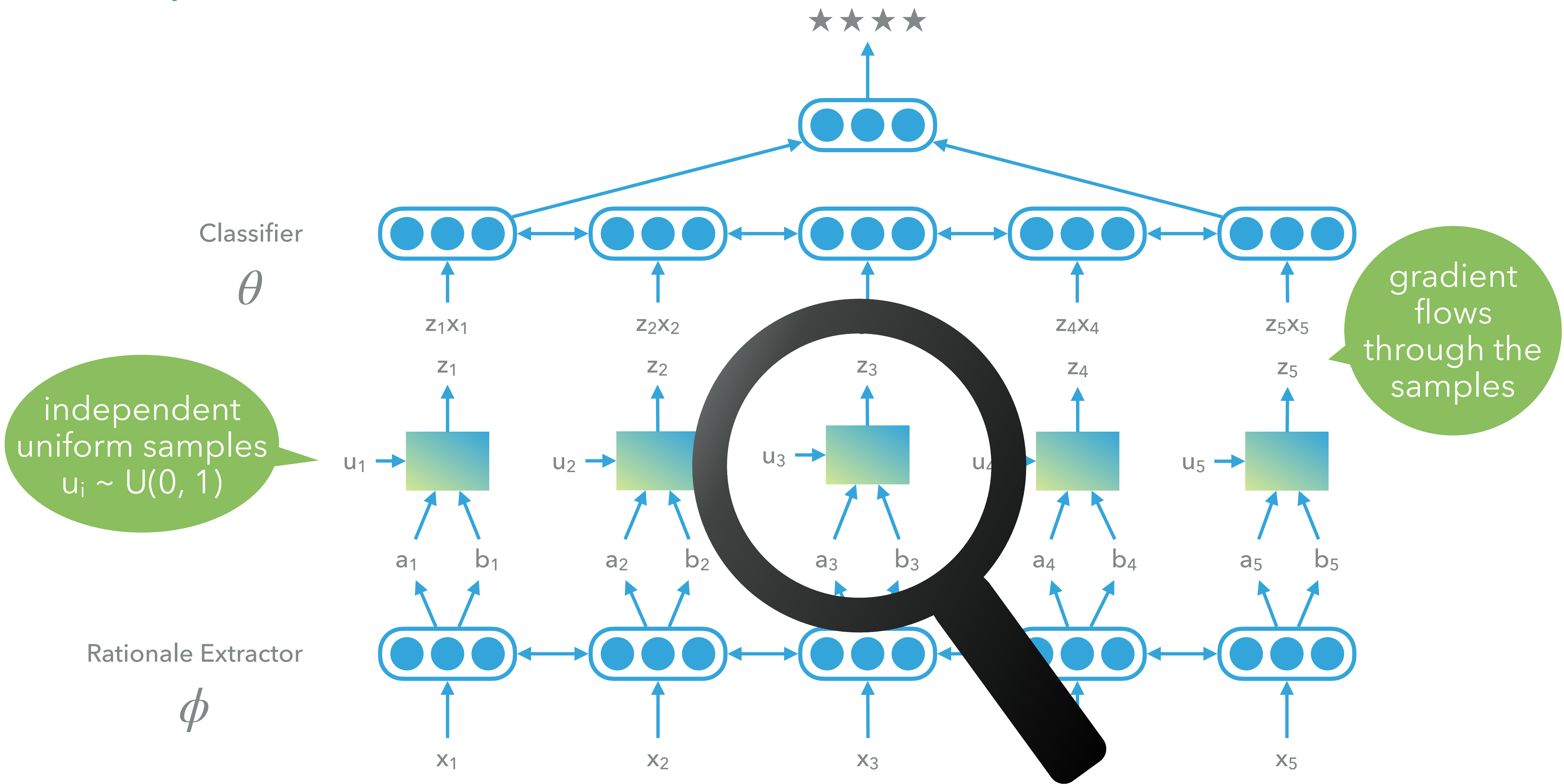
# Our Proposed Model



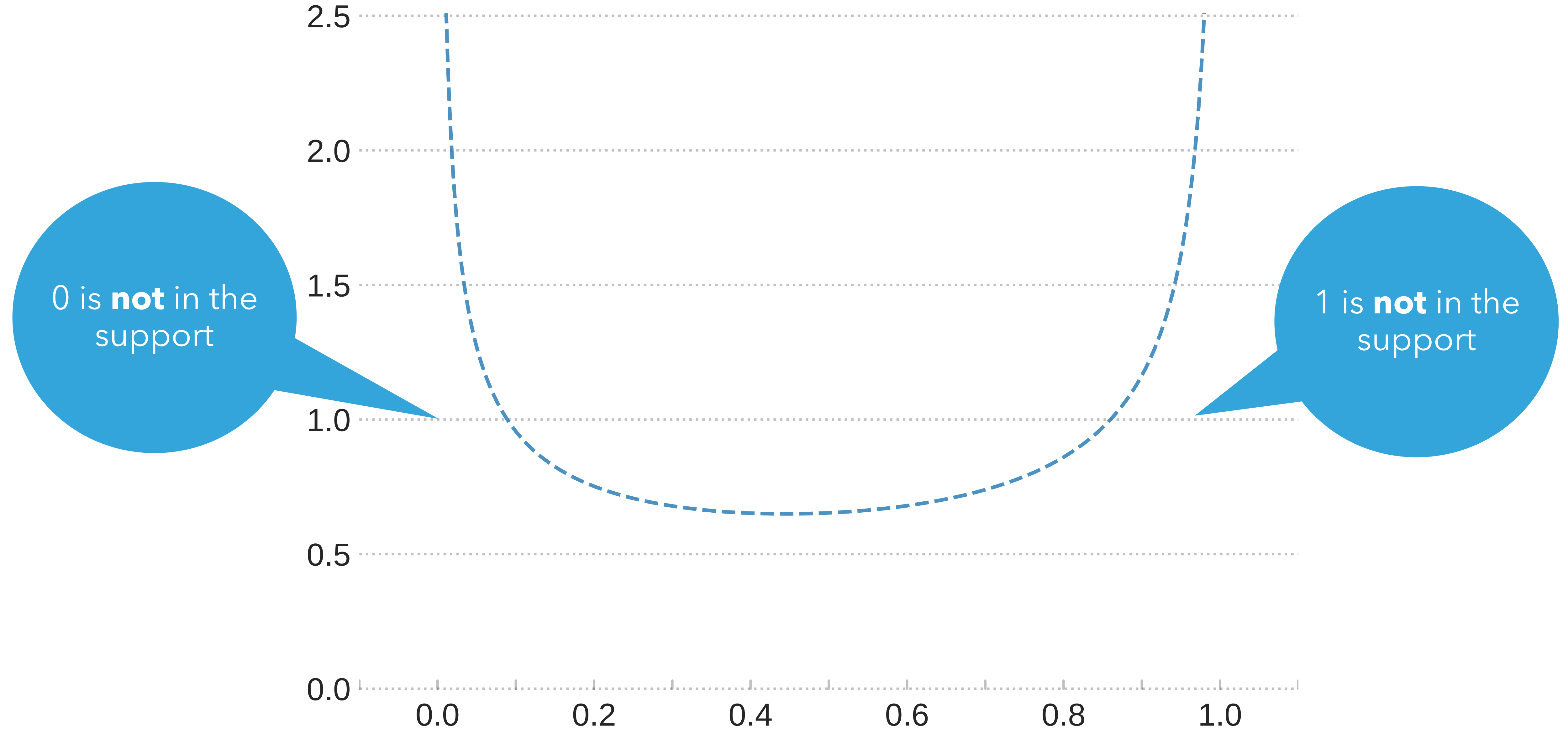
# Our Proposed Model



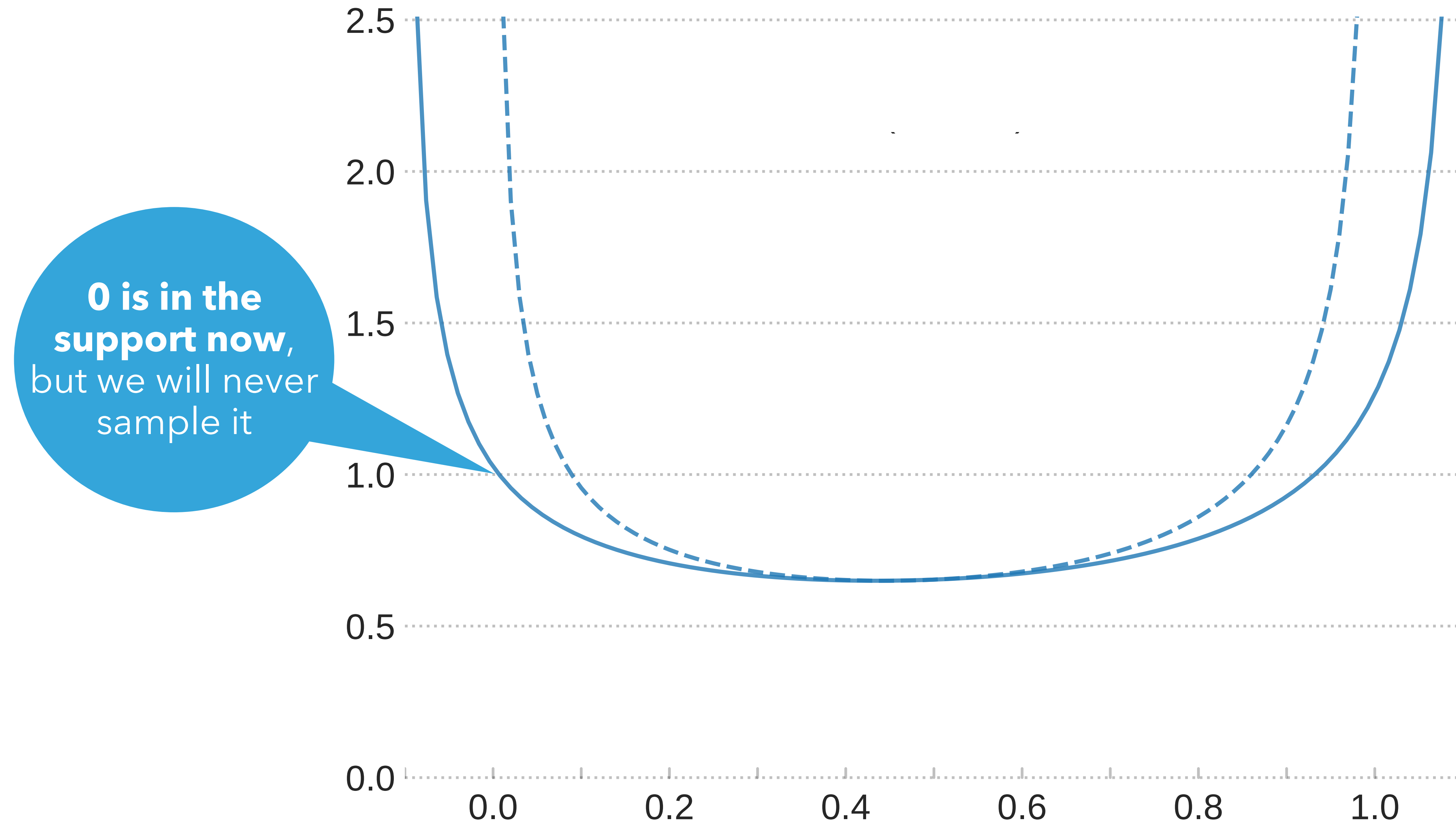
# Our Proposed Model



# Stretch-and-rectify (Louizos et al., 2018)

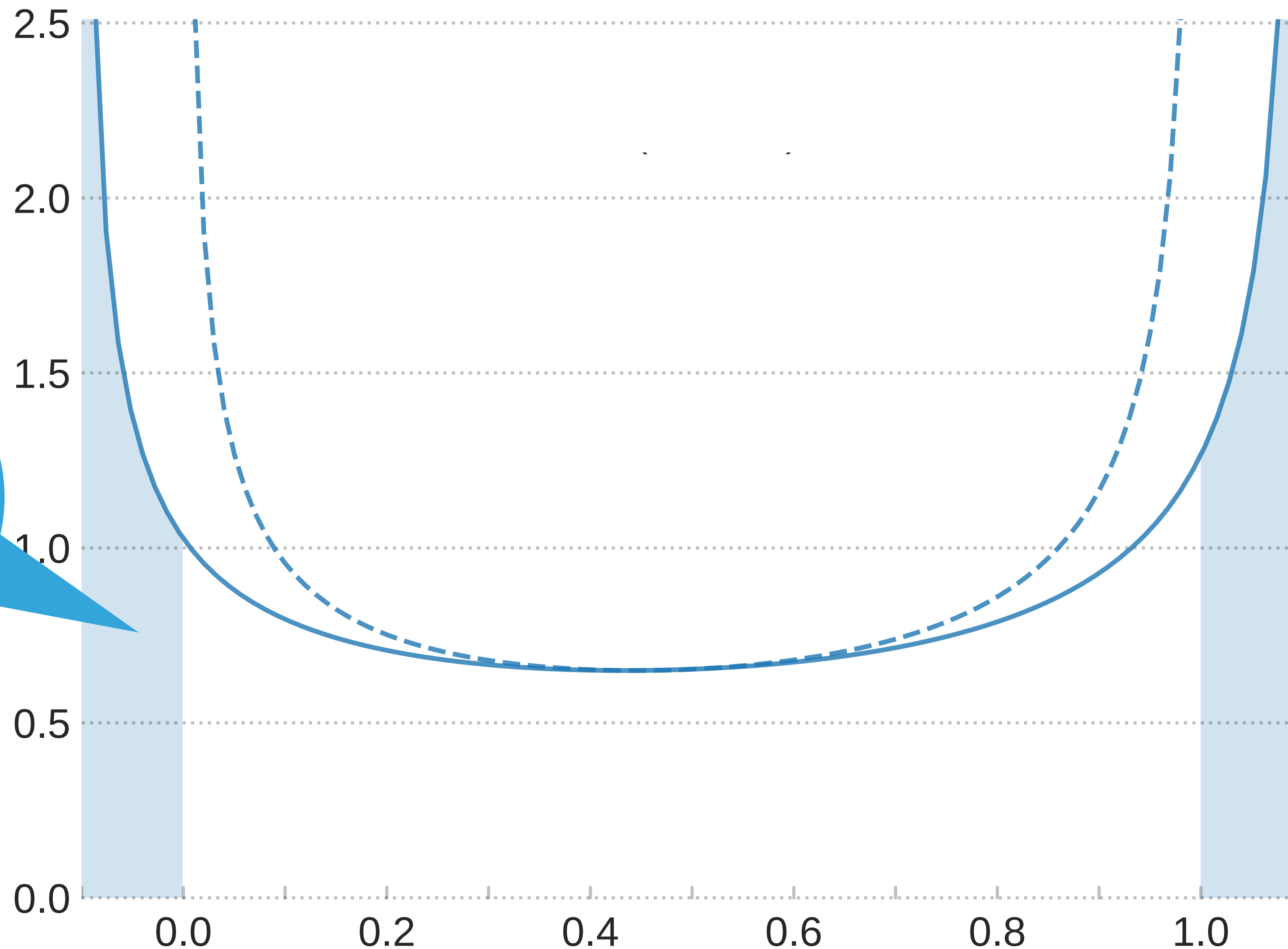


# Stretch-and-rectify (Louizos et al., 2018)



# Stretch-and-rectify (Louizos et al., 2018)

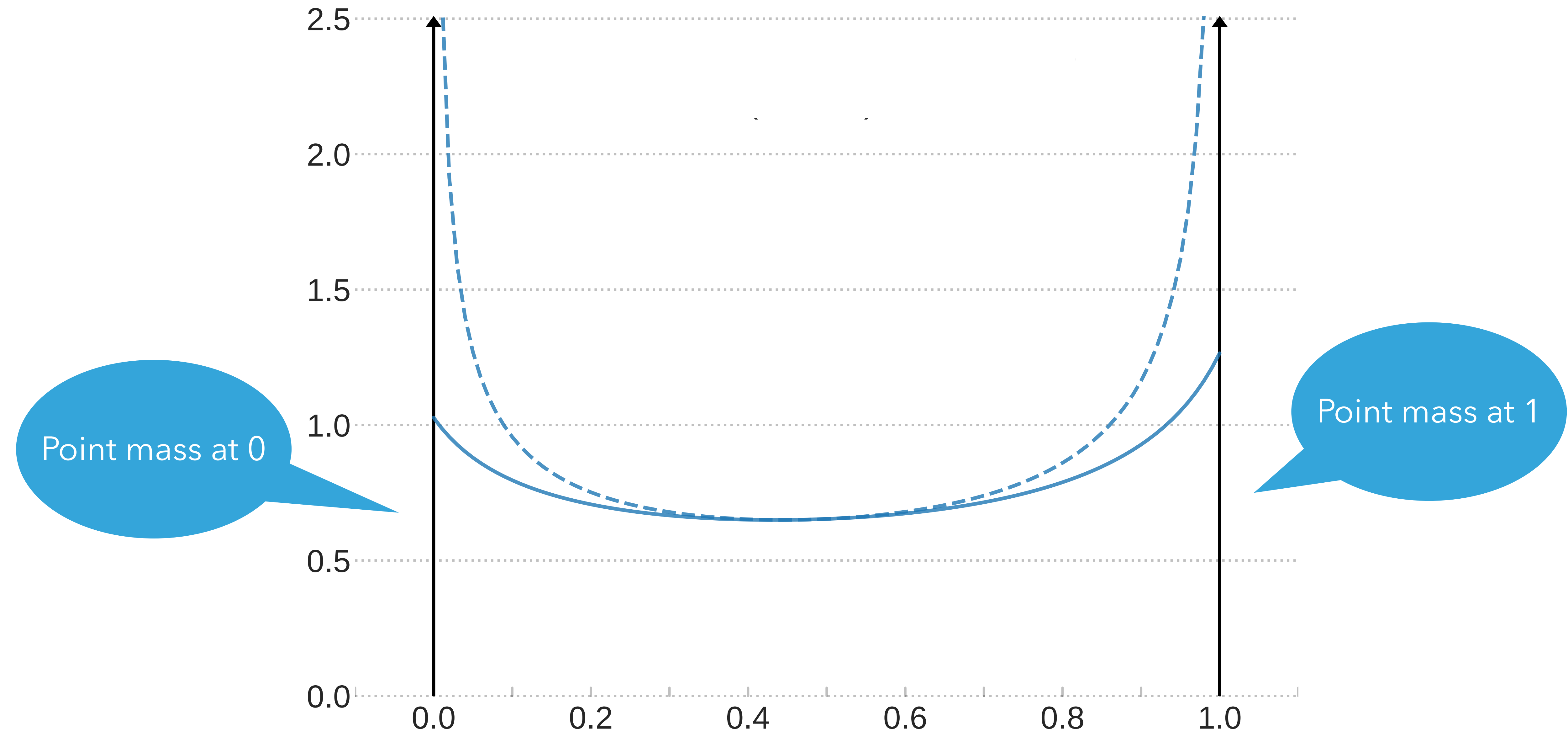
10



We **collapse**  
the shaded area  
into a **point**  
**mass**

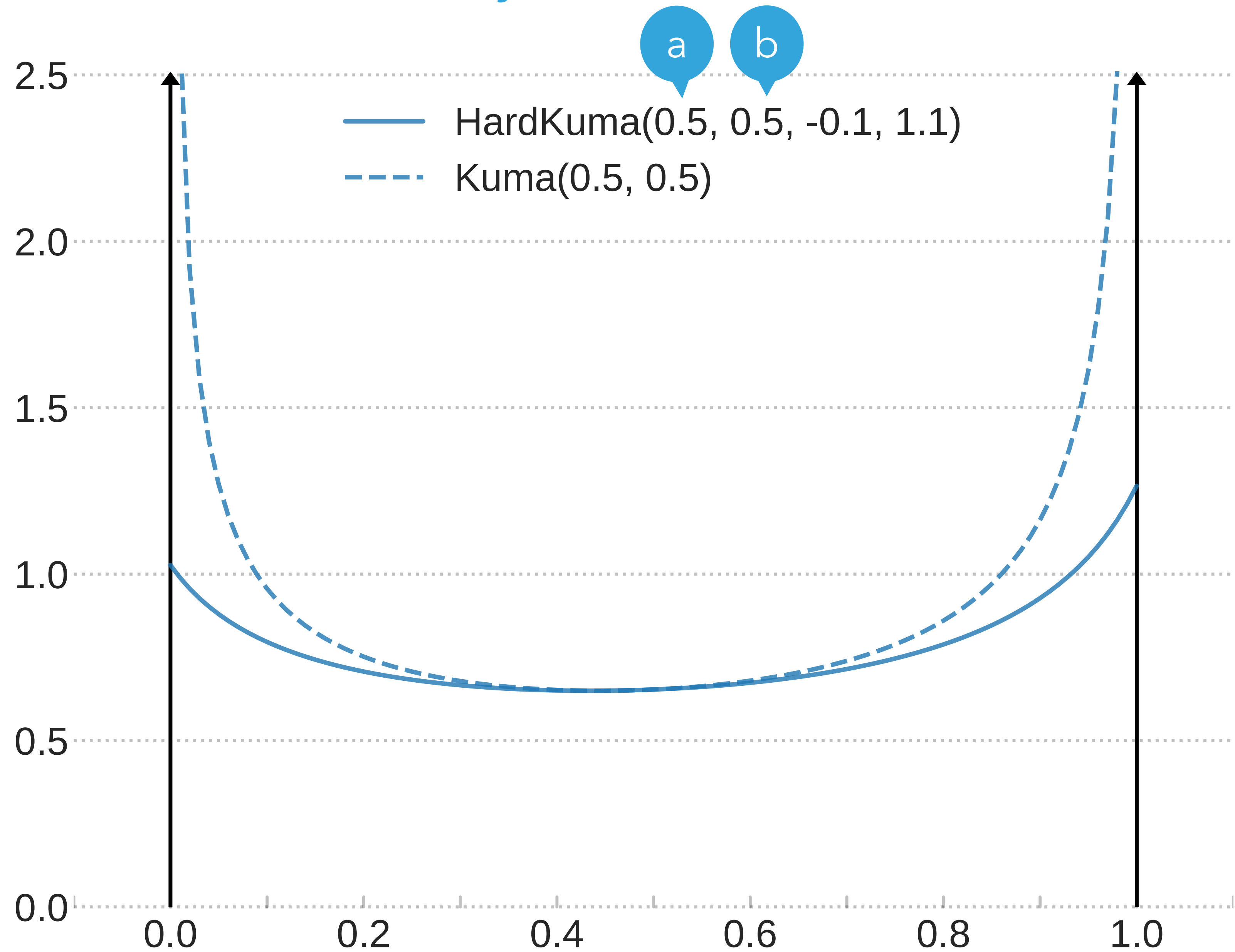
# Stretch-and-rectify (Louizos et al., 2018)

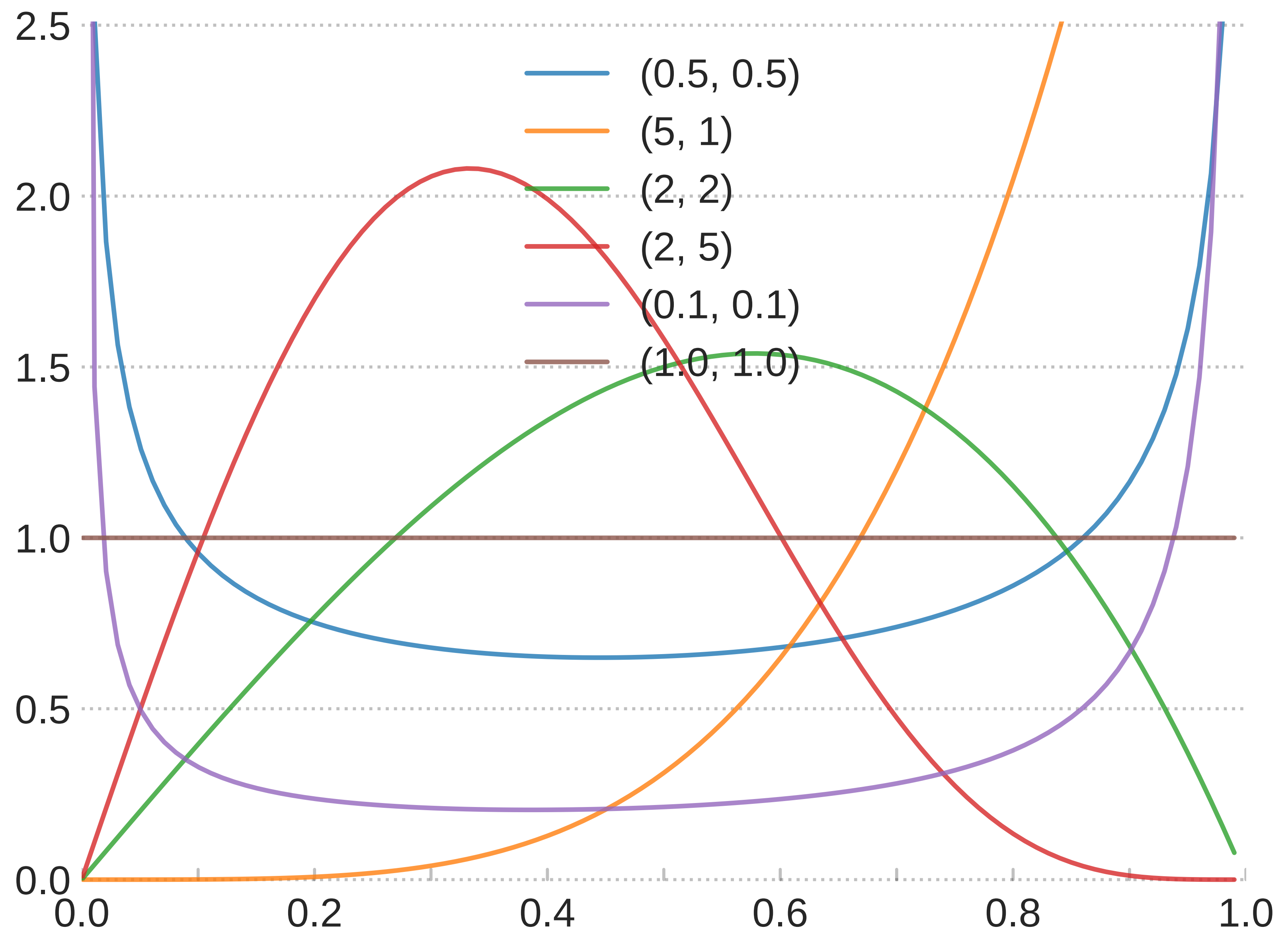
11





# In this work: Hard Kumaraswamy Distribution





- ▶ We want **short** rationales without breaking backpropagation
- ▶ Solution: relax  $L_0$  (Louizos et al., 2018)

$$L_0(z)$$

**Lei et al.**

Compute  $L_0$  for one specific assignment of  $z$

$$\mathbb{E}_{p(z|x)} [L_0(z)]$$

**Proposed model:**

$L_0$  computed for all assignments  
**on expectation**

- ▶ Baseline: **penalty for transitions** using fused lasso

$$\sum_{i=1}^{n-1} |z_i - z_{i+1}|$$

- ▶ Proposed model: compute a **relaxation** of fused lasso by computing the expected number of **zero-to-nonzero** and **nonzero-to-zero** changes:

$$\mathbb{E}_{p(z|x)} \left[ \sum_{i=1}^{n-1} \mathbb{I}[z_i = 0, z_{i+1} \neq 0] \right] + \mathbb{E}_{p(z|x)} \left[ \sum_{i=1}^{n-1} \mathbb{I}[z_i \neq 0, z_{i+1} = 0] \right]$$

# Specify target selection rate

- ▶ We want a maximum selection rate e.g. 10% of the text
- ▶ We propose a **constrained optimization** problem:

$$\min_{\phi, \theta} L(\phi, \theta) \quad \text{s.t.} \quad \mathbb{E}[L_0] < r$$

- ▶ We use Lagrangian relaxation

1. Multi-aspect sentiment analysis (BeerAdvocate, Lei et al. 2016)
  - ▶ Regression, sentiment score in  $[0,1]$
2. Stanford Sentiment (SST)
  - ▶ Classification {very negative, ..., very positive}
3. Stanford Natural Language Inference (SNLI)
  - ▶ Classification {entailment, contradiction, neutral}

Multiple aspects

Look

Smell

Taste

		Look		Smell		Taste	
		Precision	% Selected	Precision	% Selected	Precision	% Selected
Attention (Lei et al.)	Threshold	80.6	13	88.4	7	65.3	7



Multiple aspects

Look

Smell

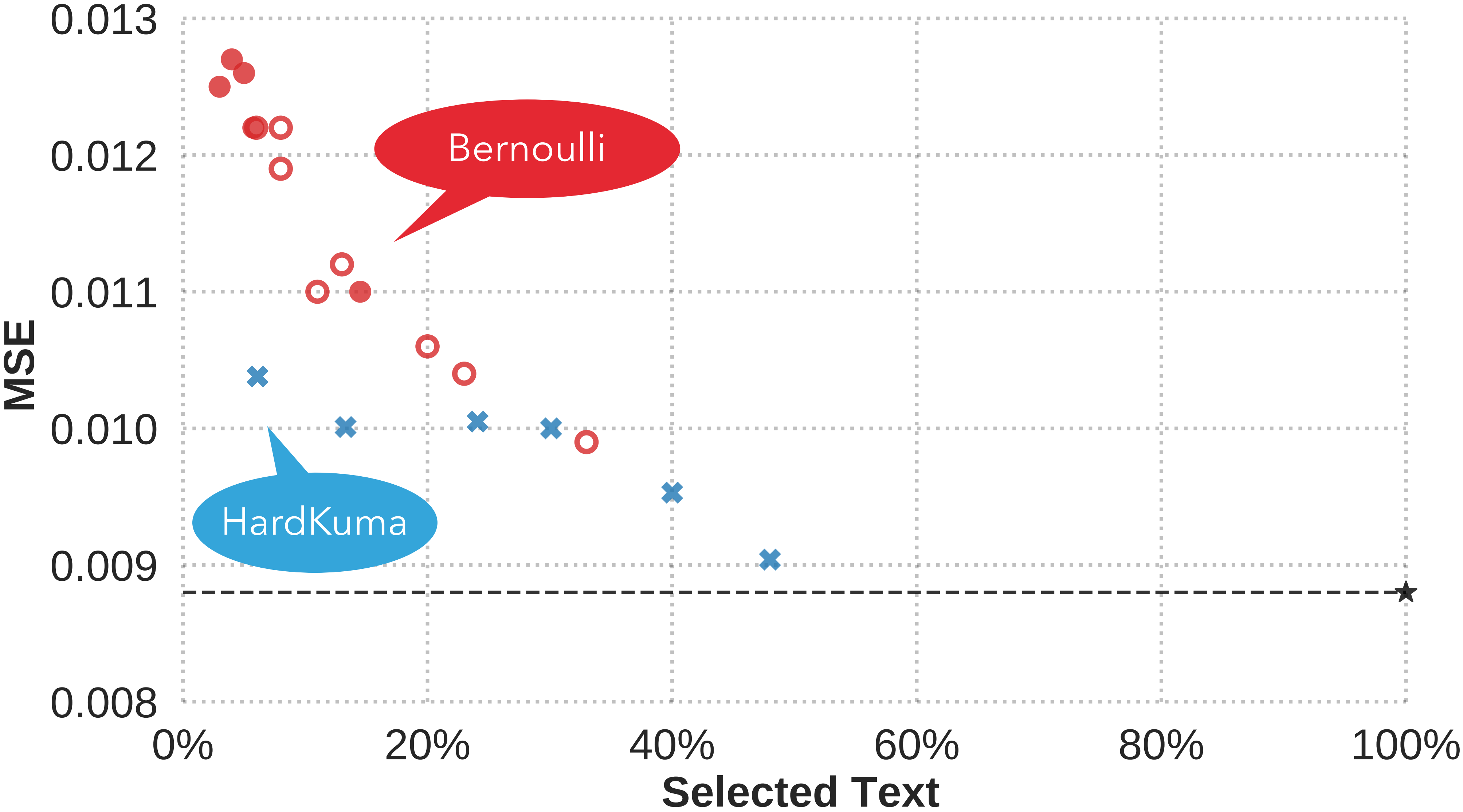
Taste

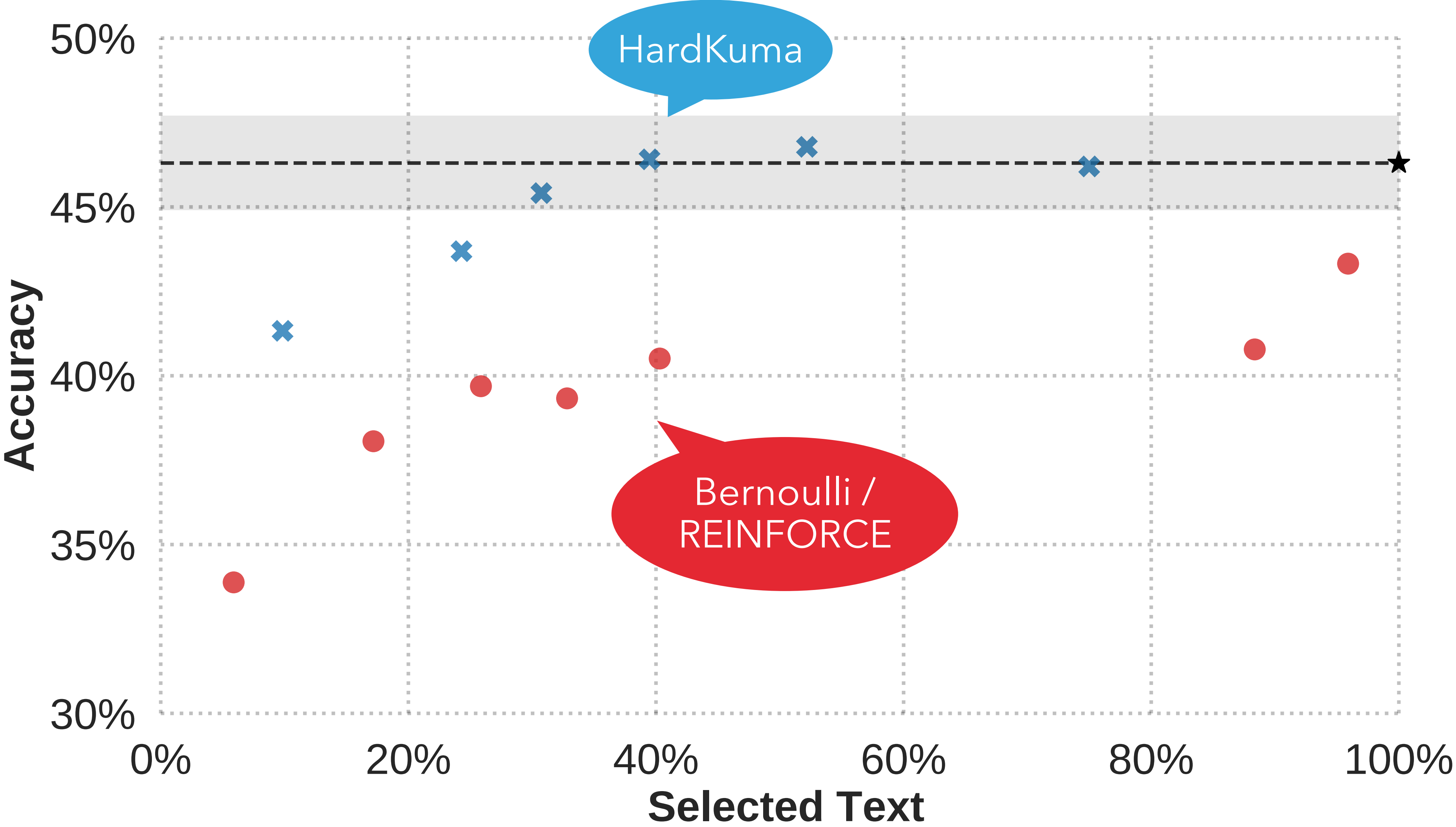
		Look		Smell		Taste	
		Precision	% Selected	Precision	% Selected	Precision	% Selected
Attention (Lei et al.)	Threshold	80.6	13	88.4	7	65.3	7
Bernoulli / REINFORCE (Lei et al.)	Tuned $\lambda$	96.3	14	95.1	7	80.2	7

Multiple aspects

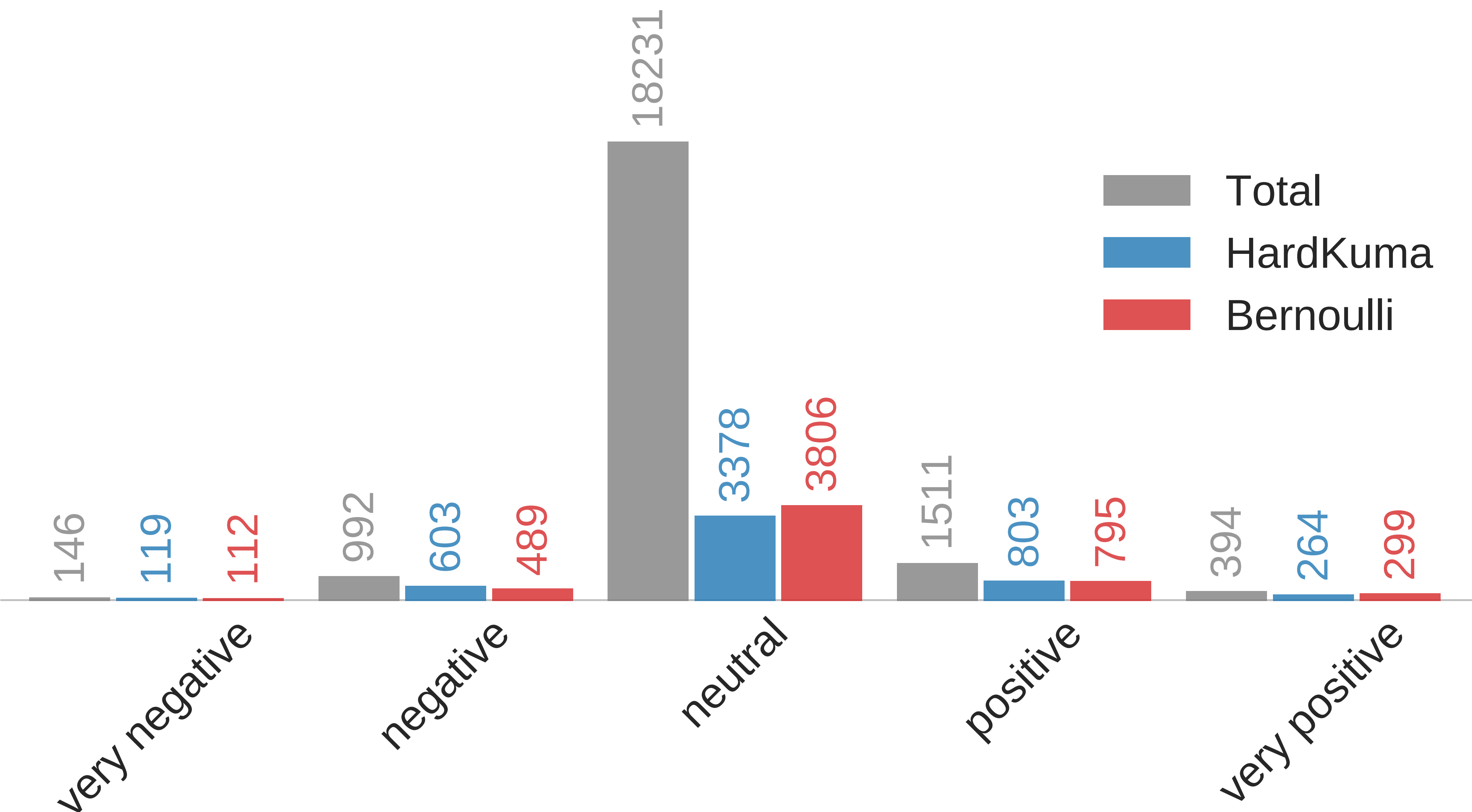
LookSmellTaste

		Precision	% Selected	Precision	% Selected	Precision	% Selected
Attention (Lei et al.)	Threshold	80.6	13	88.4	7	65.3	7
Bernoulli / REINFORCE (Lei et al.)	Tuned $\lambda$	96.3	14	95.1	7	80.2	7
HardKuma	Lagrange	98.1	13	96.8	7	89.8	7





# Analysis: Word Count per Sentiment



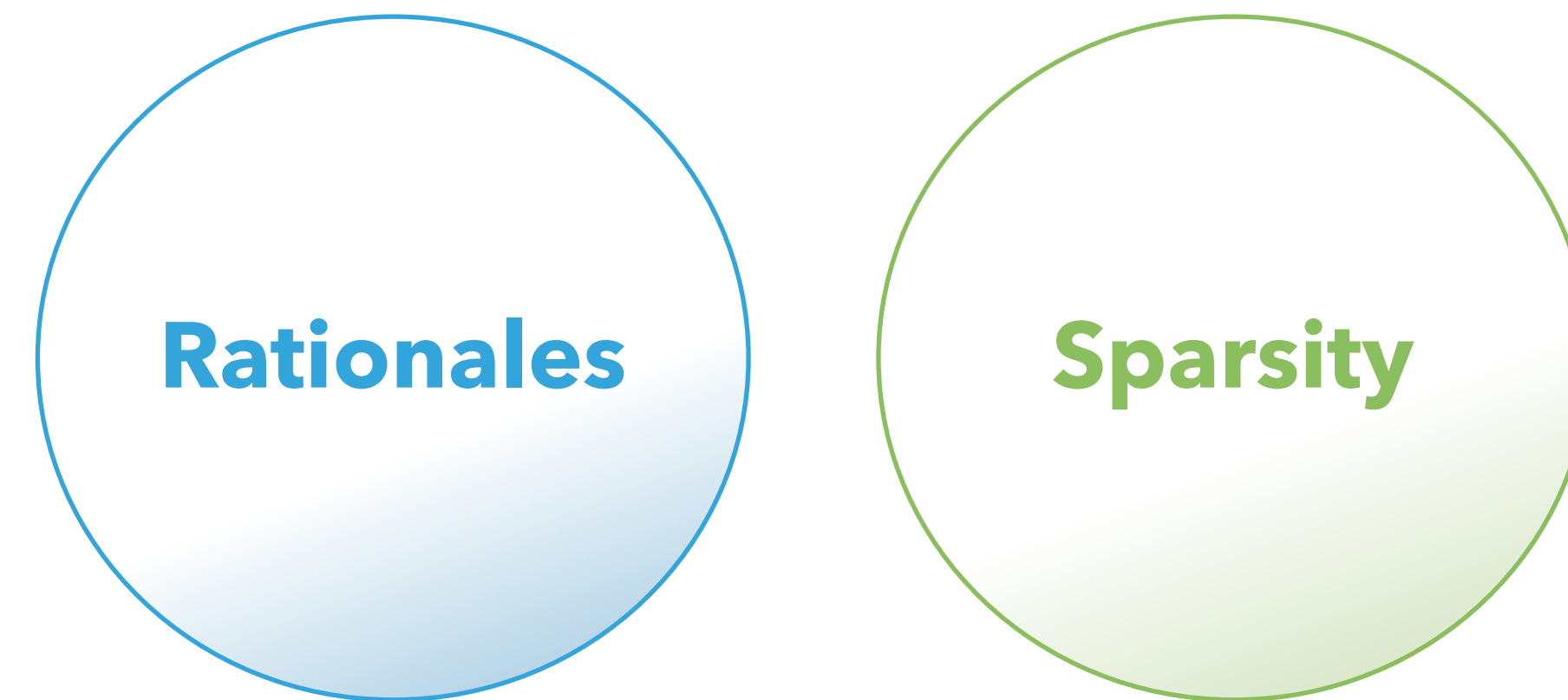
- ▶ NLI: predict {entailment, contradiction, neutral} given premise & hypothesis
- ▶ Baseline: Decomposable Attention model (Parikh et al., 2016)
- ▶ We replace Hypothesis-Premise attention with HardKuma attention

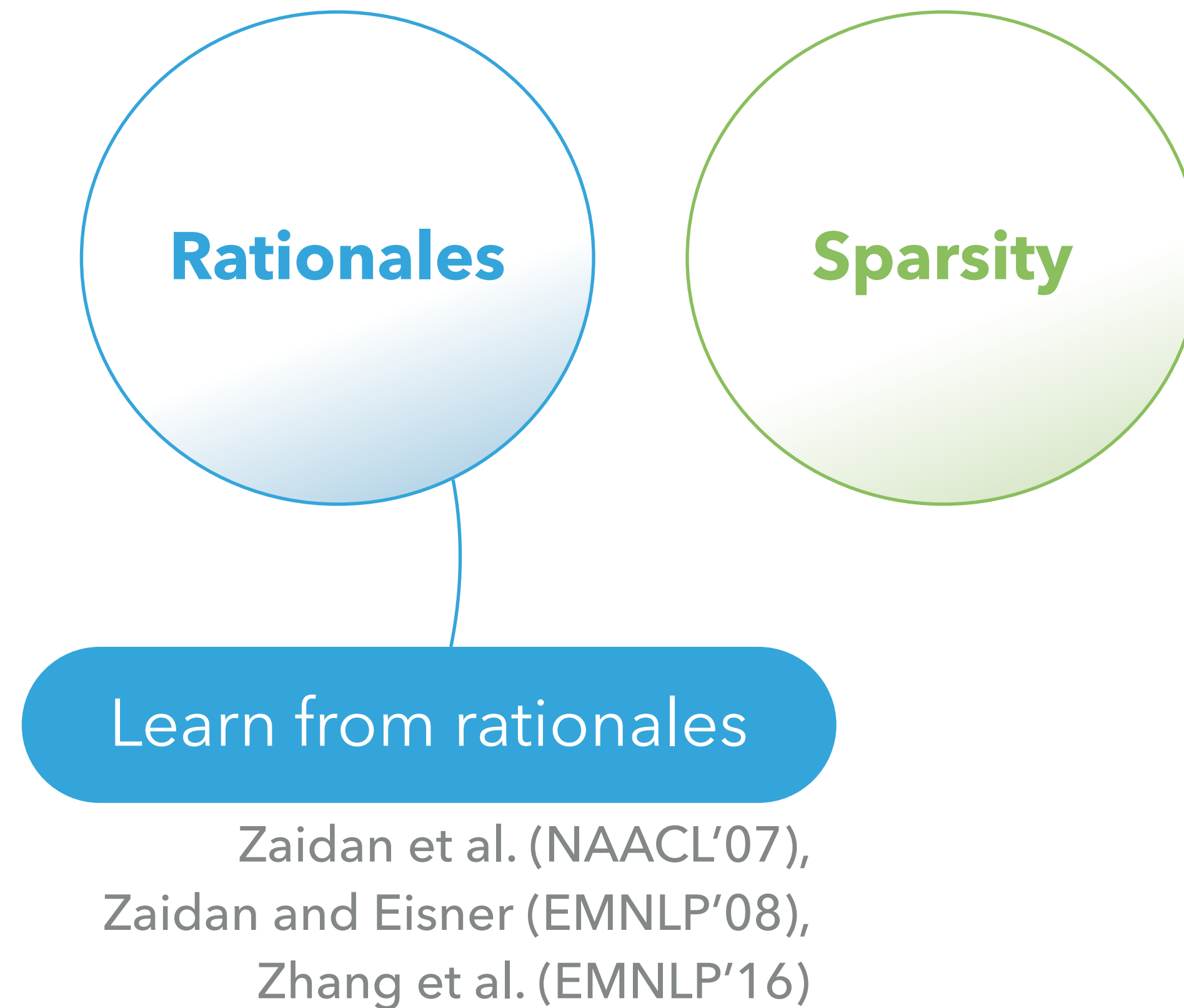
	<s>	The	man	is	walking	his	cat	.
<s>	0	0	0	0	0	0	0	0
Young	0	0	0	0	0	0	0	0
man	0	0	77	21	0	0	0	0
walking	0	0	0	0	88	0	0	0
dog	0	0	0	0	0	0	86	0

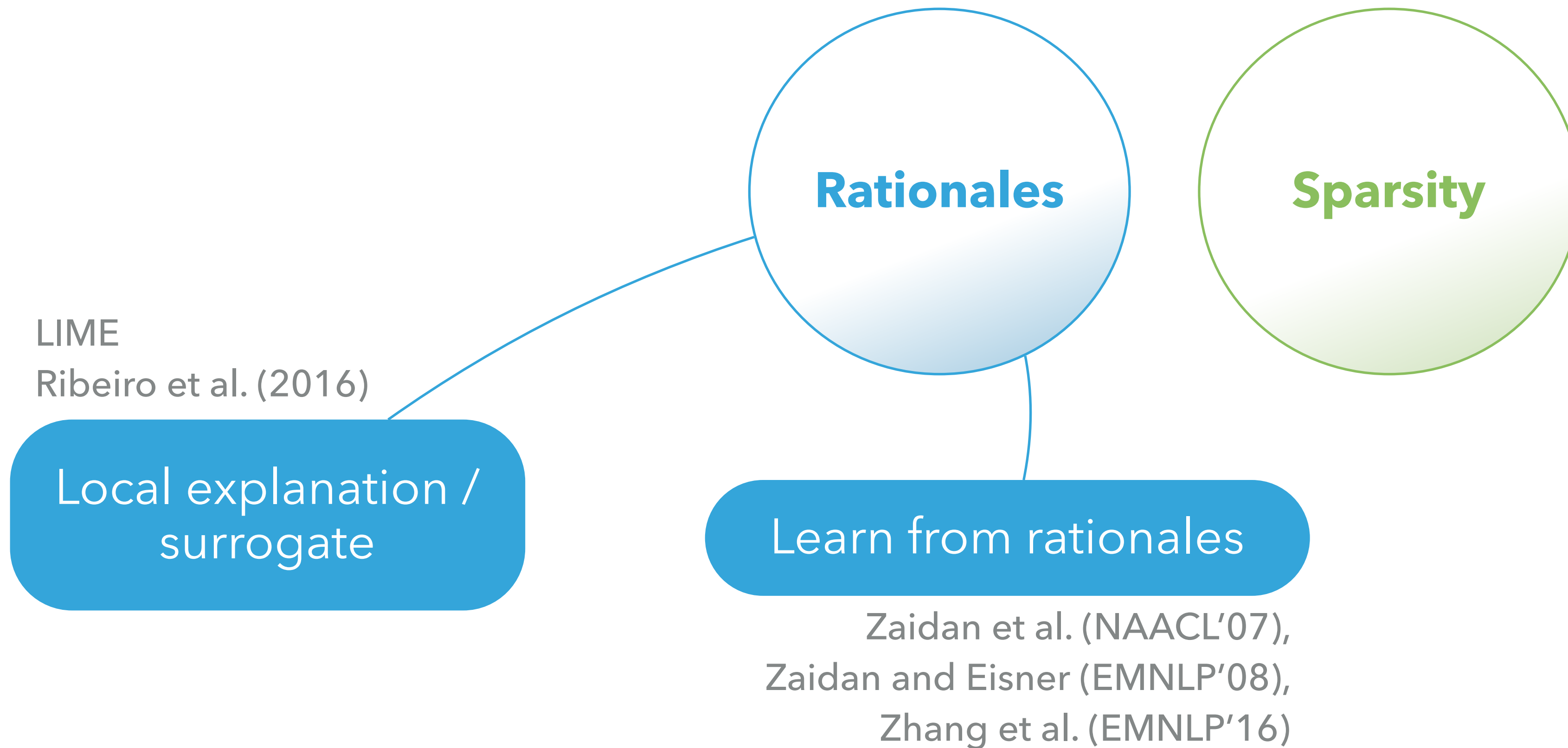
SNLI Accuracy		
	Dev	Test
LSTM (Bowman et al.)	–	80.6
DA (Parikh et al.)	–	86.3
<b>DA (reimpl.)</b>	<b>86.9</b>	<b>86.5</b>
<b>DA HardKuma</b>	<b>86.0</b>	<b>85.5</b>

Only drop **1%**  
with **8.6%** non-zero  
attention cells









this beer **pours ridiculously clear with tons of carbonation** that forms a rather impressive rocky head that settles slowly into a fairly dense layer of foam. **this is a real good lookin' beer,** unfortunately it gets worse from here ... first, **the aroma is kind of bubblegum-like and grainy,** next, the taste is sweet and grainy with an unpleasant bitterness in the finish. ... overall, the fat weasel is good for a fairly cheap buzz, but only if you like your beer grainy and bitter .

Lei et al. (EMNLP'16)  
Rationalizing Neural Predictions

Jointly train &  
learn rationales

**Rationales**

**Sparsity**

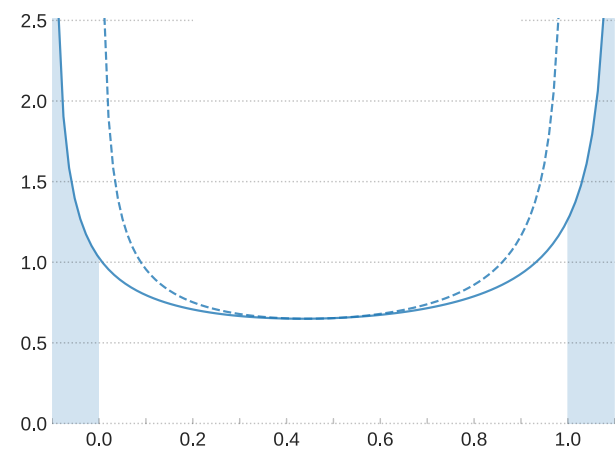
LIME  
Ribeiro et al. (2016)

Local explanation /  
surrogate

Learn from rationales

Zaidan et al. (NAACL'07),  
Zaidan and Eisner (EMNLP'08),  
Zhang et al. (EMNLP'16)

this beer **pours ridiculously clear with tons of carbonation** that forms a rather impressive rocky head that settles slowly into a fairly dense layer of foam. **this is a real good lookin' beer,** unfortunately it gets worse from here ... first, **the aroma is kind of bubblegum-like and grainy,** next, the taste is sweet and grainy with an unpleasant bitterness in the finish. ... overall, the fat weasel is good for a fairly cheap buzz, but only if you like your beer grainy and bitter .



Lei et al. (EMNLP'16)  
Rationalizing Neural Predictions

Jointly train & learn rationales

Louizos et al. (ICLR'18)  
Sparsifying parameters

Stretch & rectify



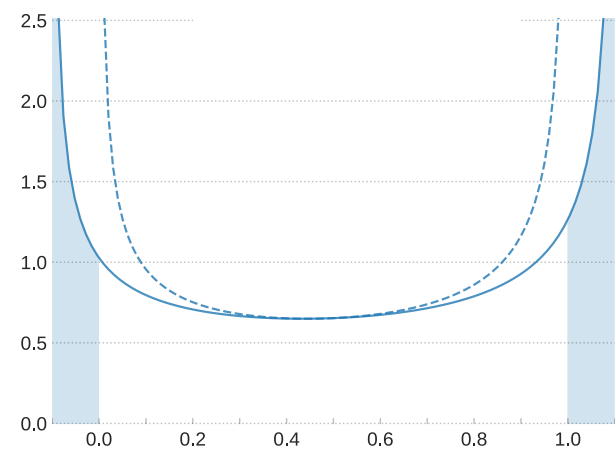
LIME  
Ribeiro et al. (2016)

Local explanation / surrogate

Learn from rationales

Zaidan et al. (NAACL'07),  
Zaidan and Eisner (EMNLP'08),  
Zhang et al. (EMNLP'16)

this beer **pours ridiculously clear with tons of carbonation** that forms a rather impressive rocky head that settles slowly into a fairly dense layer of foam. **this is a real good lookin' beer,** unfortunately it gets worse from here ... first, **the aroma is kind of bubblegum-like and grainy,** next, the taste is sweet and grainy with an unpleasant bitterness in the finish. ... overall, the fat weasel is good for a fairly cheap buzz, but only if you like your beer grainy and bitter .



Lei et al. (EMNLP'16)  
Rationalizing Neural Predictions

Jointly train & learn rationales

**This work**

Louizos et al. (ICLR'18)  
Sparsifying parameters

Stretch & rectify

**Rationales**

**Sparsity**

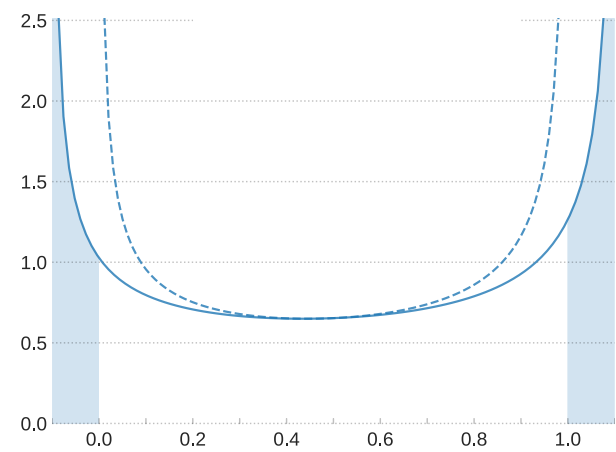
LIME  
Ribeiro et al. (2016)

Local explanation / surrogate

Learn from rationales

Zaidan et al. (NAACL'07),  
Zaidan and Eisner (EMNLP'08),  
Zhang et al. (EMNLP'16)

this beer **pours ridiculously clear with tons of carbonation** that forms a rather impressive rocky head that settles slowly into a fairly dense layer of foam. **this is a real good lookin' beer,** unfortunately it gets worse from here ... first, **the aroma is kind of bubblegum-like and grainy,** next, the taste is sweet and grainy with an unpleasant bitterness in the finish. ... overall, the fat weasel is good for a fairly cheap buzz, but only if you like your beer grainy and bitter .



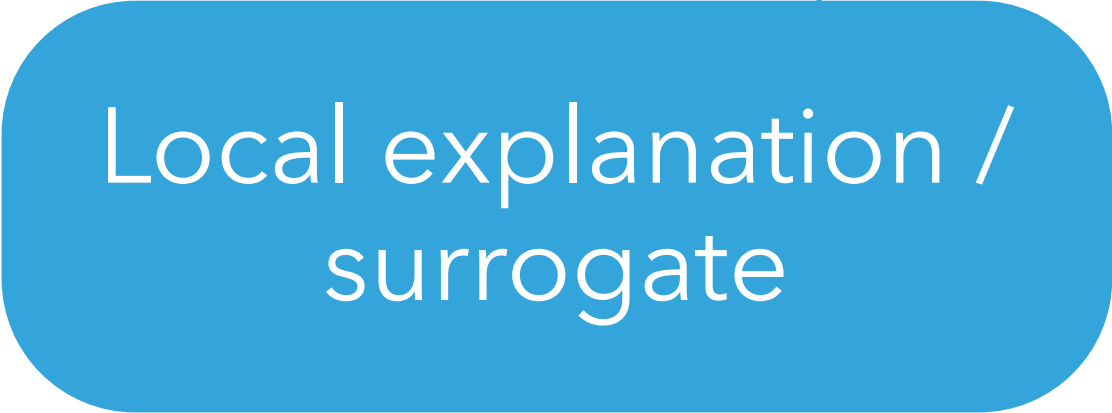
Lei et al. (EMNLP'16)  
Rationalizing Neural Predictions



Louizos et al. (ICLR'18)  
Sparsifying parameters



LIME  
Ribeiro et al. (2016)



Voita et al. (ACL'19)



Zaidan et al. (NAACL'07),  
Zaidan and Eisner (EMNLP'08),  
Zhang et al. (EMNLP'16)



- ▶ Differentiable approach to extractive rationales
  - ▶ Stretch and rectify using HardKuma
  - ▶ Support for binary outcomes
- ▶ Objective to specify the percentage of selected text
- ▶ Future work: interpretable QA / fact checking
- ▶ Code online: [github.com/bastings](https://github.com/bastings)
  - ▶ DIY: add a HardKuma layer to your classifier!



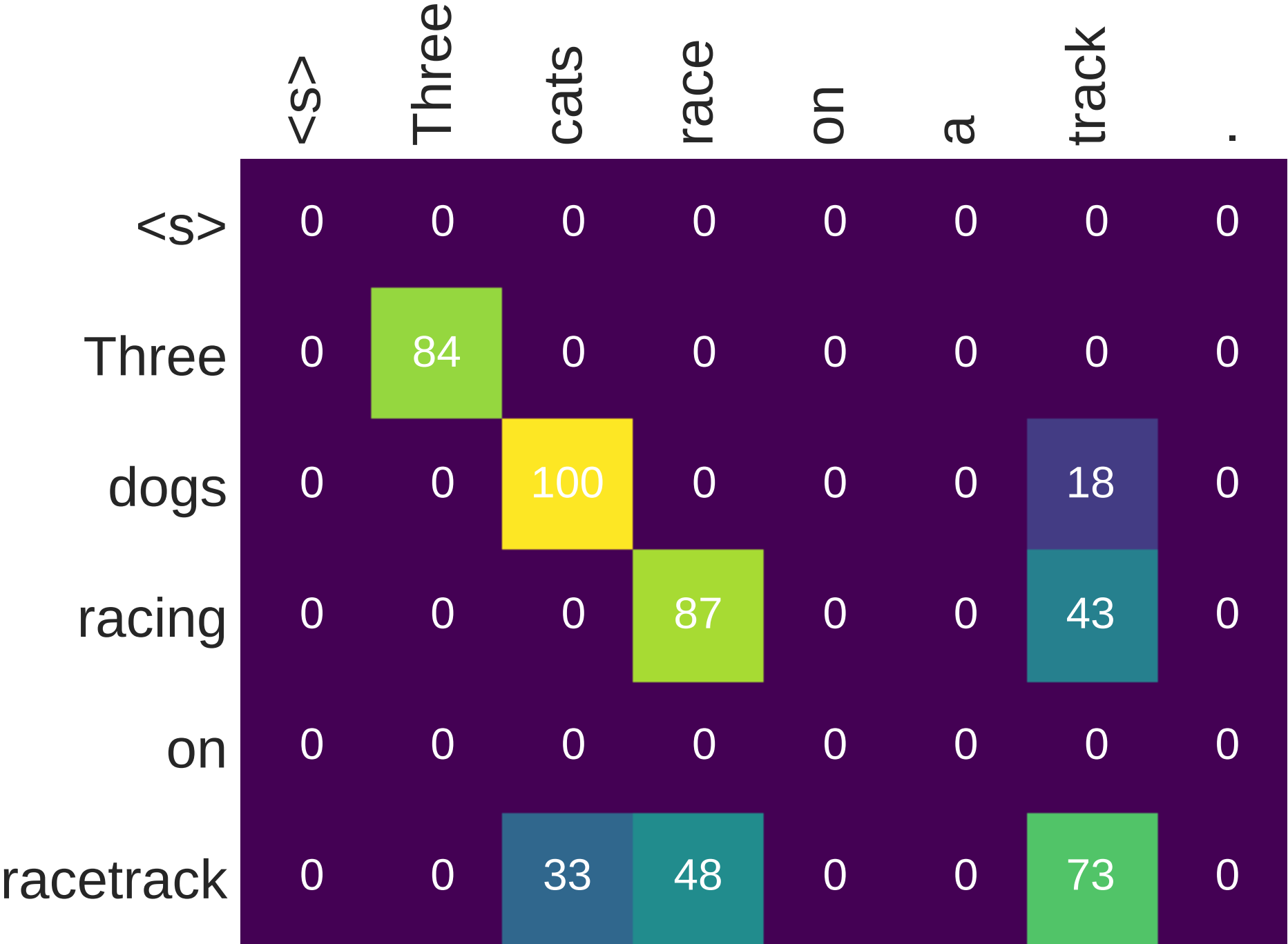
Thank you!

Code online @ [github.com/bastings](https://github.com/bastings)

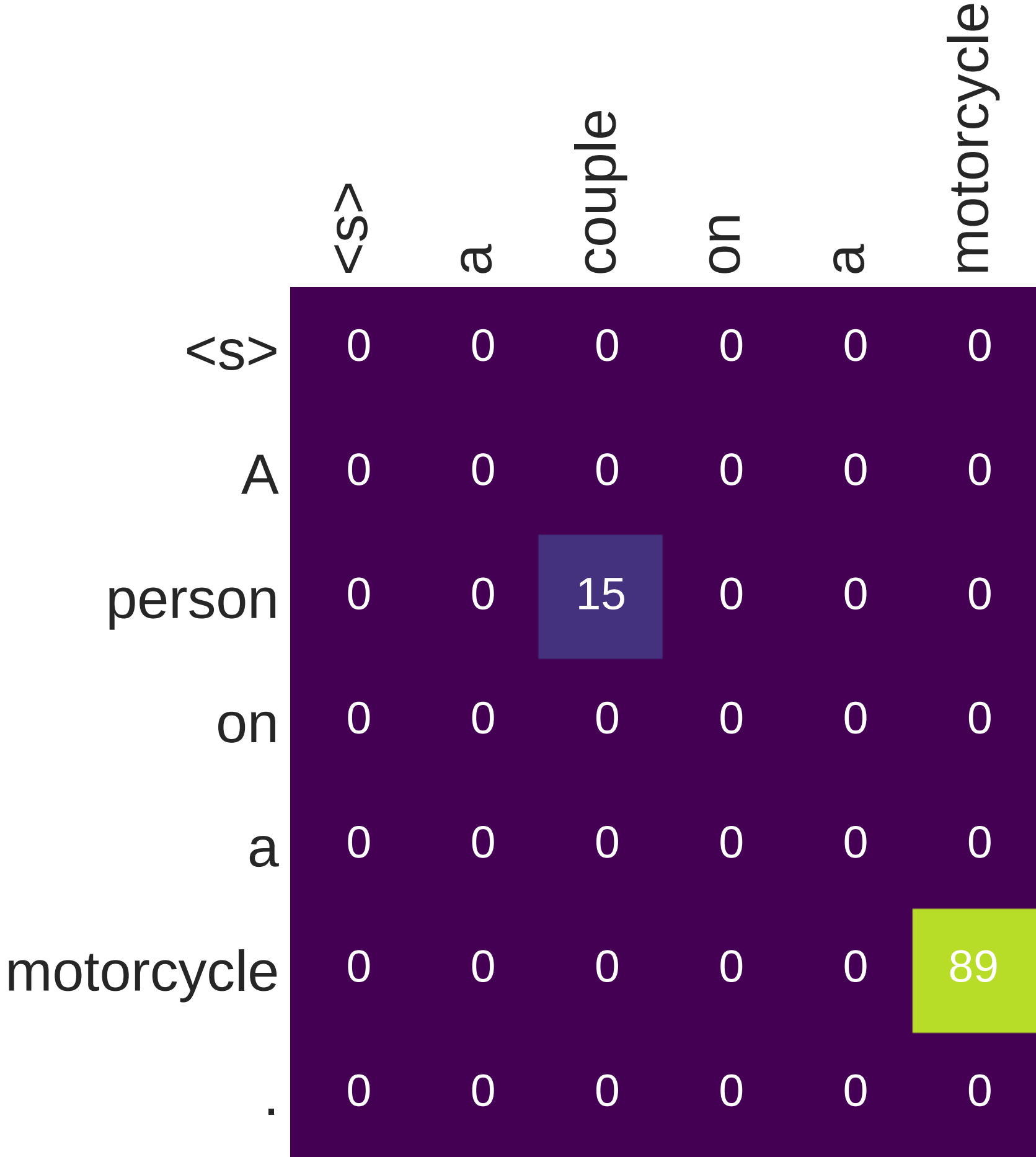


Check out our new NMT toolkit for novices Joey NMT at [github.com/joeynmt](https://github.com/joeynmt)

# Example: Contradiction



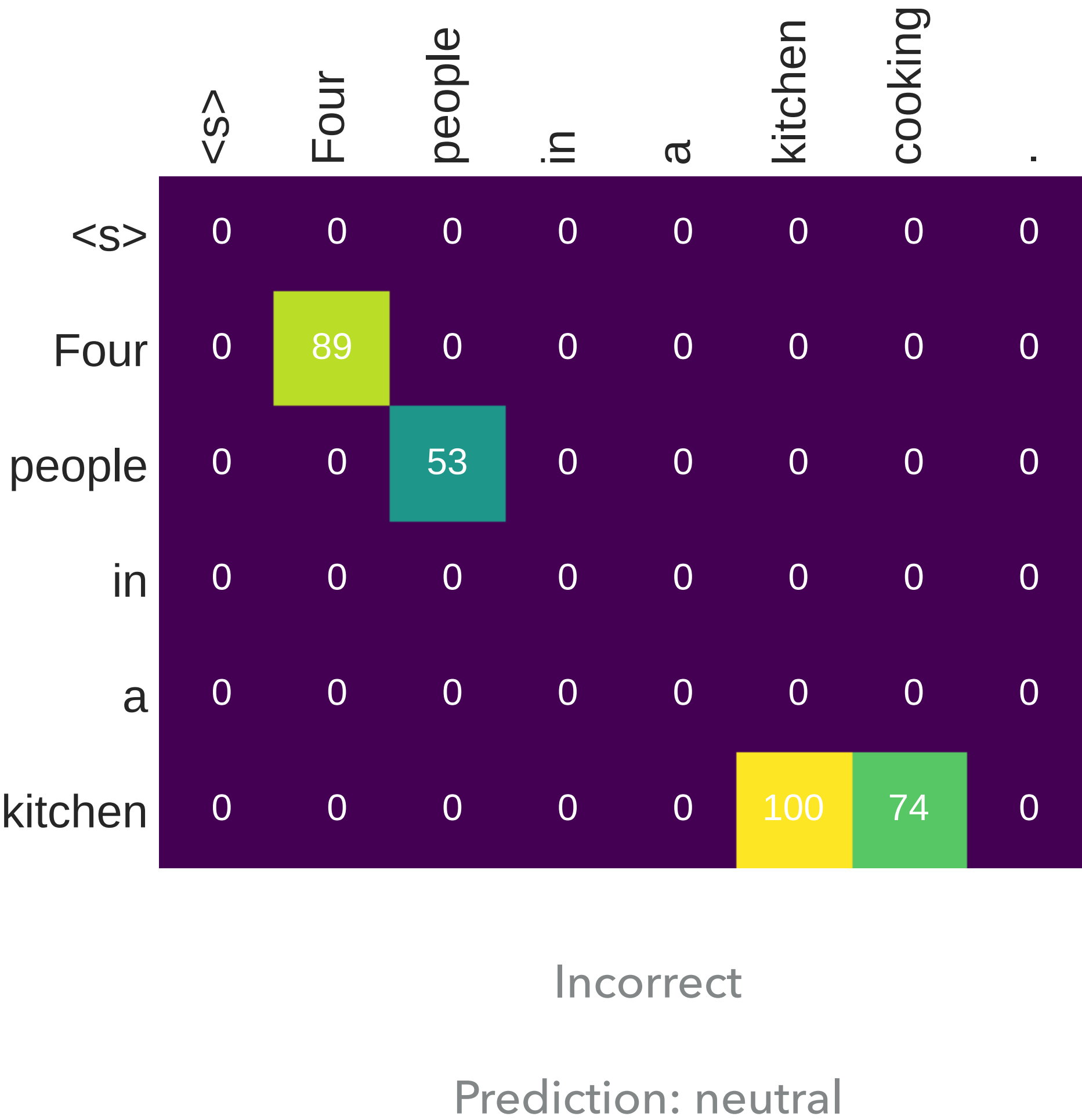
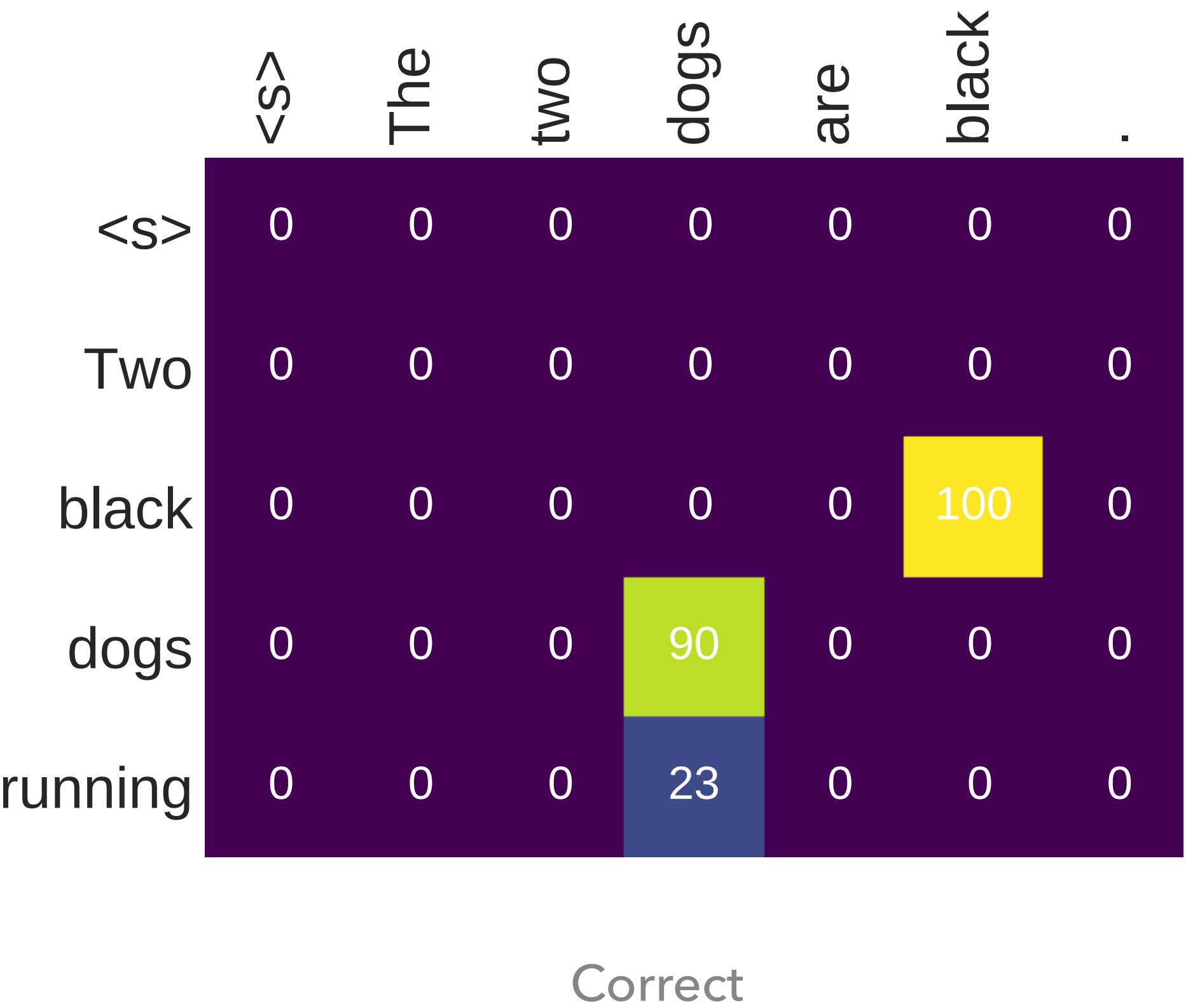
Correct



Incorrect

Prediction: entailment

# Example: Entailment



# Example: Neutral

	<S>	They	are	in	the	desert	.
<S>	0	0	0	0	0	0	0
People	0	0	0	0	0	0	0
walking	0	0	0	0	0	0	0
through	0	0	0	0	0	0	0
dirt	0	0	0	0	0	81	0
.	0	0	0	0	0	0	0

Correct

	<S>	A	dog	found	a	bone
<S>	0	0	0	0	0	0
A	0	0	0	0	0	0
dog	0	0	89	13	0	12
gnawing	0	0	0	0	0	47
on	0	0	0	0	0	0
a	0	0	0	0	0	0
bone	0	0	12	14	0	76
.	0	0	0	0	0	0

Incorrect

Prediction: entailment

$$\begin{aligned}\log P(y \mid x) &= \log \mathbb{E}_{P(z|x,\phi)} [P(y \mid x, z, \theta)] \\ &\stackrel{\text{JL}}{\geq} \mathbb{E}_{P(z|x,\phi)} [\log P(y \mid x, z, \theta)] \\ &= \mathcal{E}(\phi, \theta)\end{aligned}$$



We maximise this lower bound on the log-likelihood

The diagram illustrates the components of a baseline regularizer. The equation is presented as a sum of three terms, each highlighted in a colored box. Callouts in speech bubbles explain each term: the first term is the loss for sufficient rationales, the second is the L0 norm for short rationales, and the third is the fused lasso for coherent rationales.

$$\min_{\phi, \theta} L(\phi, \theta) + \lambda_0 \sum_{i=1}^n z_i + \lambda_1 \sum_{i=1}^{n-1} |z_i - z_{i+1}|$$

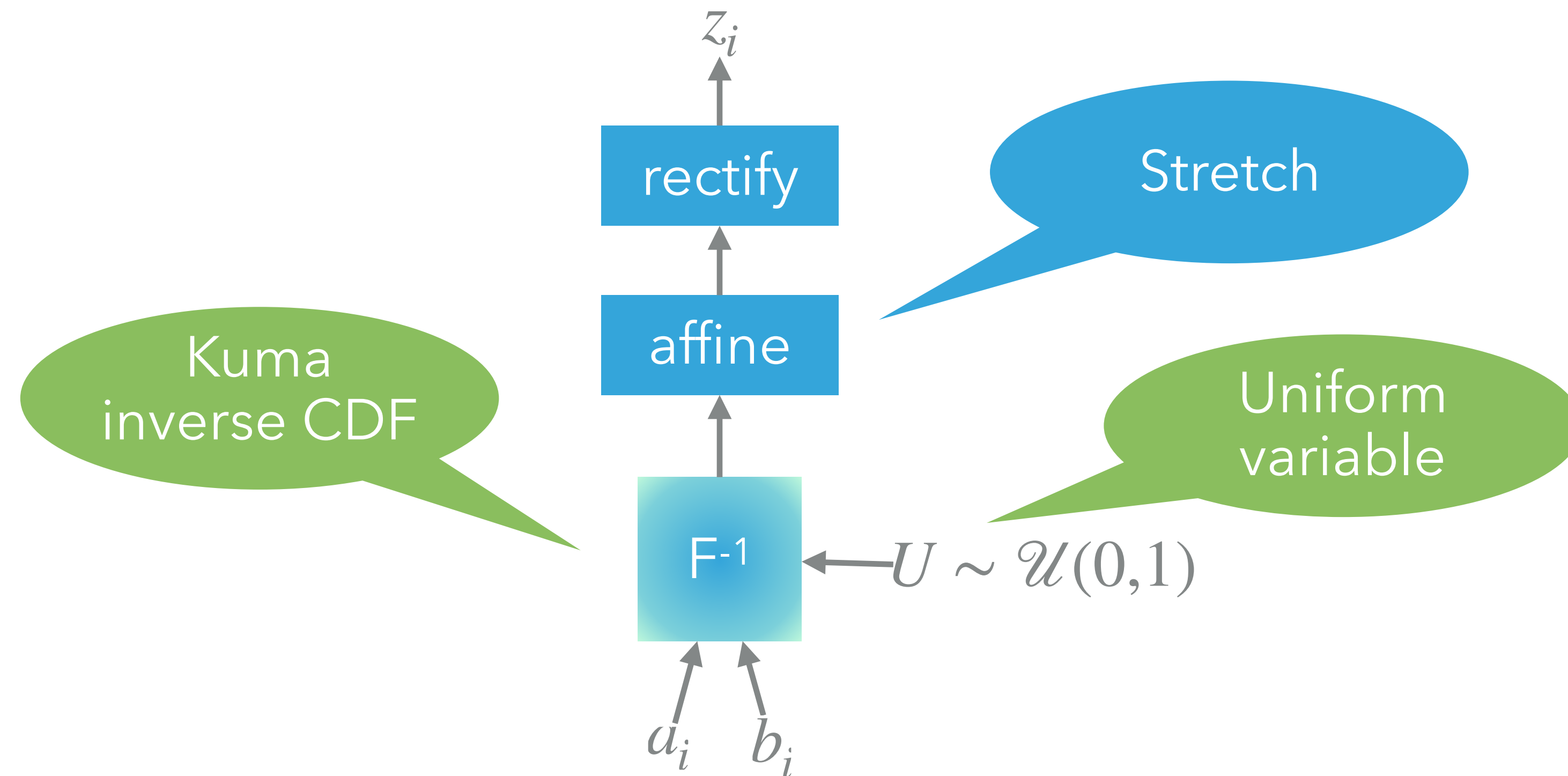
Loss for **sufficient** rationales

$L_0$  for **short** rationales

Fused lasso for **coherent** rationales

# How to get a HardKuma sample from a uniform variable

32



Kuma  
inverse CDF

$$F_K^{-1}(u; a, b) = \left(1 - (1 - u)^{\frac{1}{b}}\right)^{\frac{1}{a}} \quad u \in [0, 1]$$

Sample  
using uniform  
random source  $U$

$$F_Z^{-1}(U; a, b) \sim \text{Kuma}(a, b) \quad U \sim \mathcal{U}(0, 1)$$



- ▶ We **stretch** the support of the Kuma to  $(l, r)$ :

$$F_T(t; a, b, l, r) = F_K\left(\frac{(t - l)}{(r - l)}; a, b\right)$$

- ▶ And define a **rectified** random variable:

$$H \sim \text{HardKuma}(a, b, l, r)$$

- ▶ by passing a Kuma sample  $t$  through a hard sigmoid:

$$T \sim \text{Kuma}(a, b, l, r) \quad h = \min(1, \max(0, t))$$

Support  
in **closed** interval  
[0, 1]

- ▶ Sampling  $h=0$  means sampling any  $t \in (l,0]$

- ▶ with mass under Kuma:

$$\mathbb{P}(H = 0) = F_K \left( \frac{-l}{r-l}; a, b \right)$$

- ▶ Sampling  $h=1$  means sampling any  $t \in [1,r)$

- ▶ with mass under Kuma:

$$\mathbb{P}(H = 1) = 1 - F_K \left( \frac{1-l}{r-l}; a, b \right)$$

$$\begin{aligned}\mathbb{E}_{p(z|x)} [L_0(z)] &\stackrel{\text{ind}}{=} \sum_{i=1}^n \mathbb{E}_{p(z_i|x)} [\mathbb{I}[z_i \neq 0]] \\ &= \sum_{i=1}^n 1 - \mathbb{P}(Z_i = 0) ,\end{aligned}$$

- ▶ We can also compute a relaxation of fused lasso by computing the expected number of **zero-to-nonzero** and **nonzero-to-zero** changes:

$$\mathbb{E}_{p(z|x)} \left[ \sum_{i=1}^{n-1} \mathbb{I}[z_i = 0, z_{i+1} \neq 0] \right] + \mathbb{E}_{p(z|x)} \left[ \sum_{i=1}^{n-1} \mathbb{I}[z_i \neq 0, z_{i+1} = 0] \right]$$
$$\stackrel{\text{ind}}{=} \sum_{i=1}^{n-1} \mathbb{P}(Z_i = 0)(1 - \mathbb{P}(Z_{i+1} = 0)) + (1 - \mathbb{P}(Z_i = 0))\mathbb{P}(Z_{i+1} = 0)$$

$$U \sim \mathcal{U}(0,1)$$

$$F_X^{-1}(u) \sim X$$

# Why are gradients possible?

39

- ▶ We consider the case where we need derivatives of a function  $L(u)$  of the underlying uniform variable  $u$ , as when we compute reparameterized gradients in variational inference. By chain rule:

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial h} \times \frac{\partial h}{\partial t} \times \frac{\partial t}{\partial k} \times \frac{\partial k}{\partial u}$$

Depends on differentiable observation model

Derivative for hardsigmoid  
0 for  $t < 0$ ,  
1 for  $0 < t < 1$ ,  
0 for  $t > 1$ ,  
undef for  $t = \{0, 1\}$

r-l

Depends on Kuma inverse CDF, no challenge