

# Specializing Word Embeddings (for Parsing) by Information Bottleneck

Best Paper Award in EMNLP 2019  
Li and Eisner, 2019

---

**Gangwoo Kim**

Data Mining & Information Systems Lab.  
Department of Computer Science and Engineering,  
College of Informatics, Korea University

## Abstract

a **variational information bottleneck(VIB)** method to nonlinearly compress contextual embedding, keeping **only** information that helps down stream tasks

compress each embeddings to either **discrete** tag or a **continuous** vector

## Properties of Variational Information Bottleneck (VIB)

1. exploiting the *existing* information
  - <-> Fine-tuning, introducing new info.
  - less danger of overfitting / very fast
2. Stochasticity
  - <-> Dimension reduction techniques, deterministic
  - blurring unneeded capacity via randomness

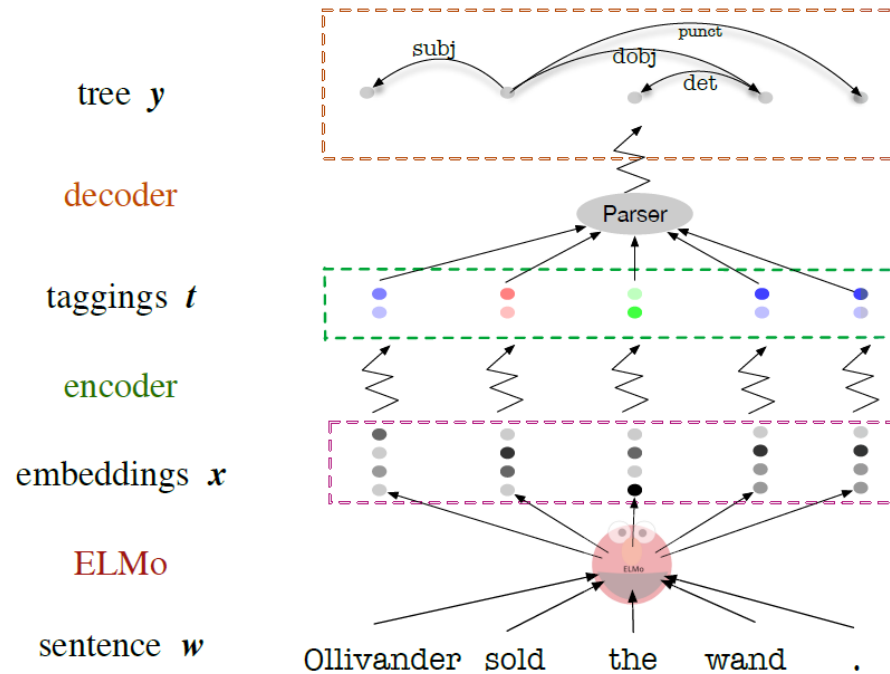


Figure 1: Our instantiation of the information bottleneck, with bottleneck variable  $T$ . A jagged arrow indicates a stochastic mapping, i.e. the jagged arrow points from the parameters of a distribution to a sample drawn from that distribution.

## Usage

The information bottleneck (IB) method originated in information theory and has been adopted in machine learning

as a training objective

N. Tishby et al, 1999, “The information bottleneck method”.  
and a framework for analyzing deep neural network

N. Tishby et al, 2015, “Deep Learning and the Information Bottleneck Principle”

**Our goal** is to learn a stochastic map from  $X$  to some compressed representation  $T$  with

- 1) leave us only the parts of  $X$  relevant to  $Y$
- 2) squeeze out any information in  $X$  irrelevant to  $Y$

both are measured by the mutual information

$$\min_{p(t|x): I(T;Y) \geq I^*} I(X;T)$$

$$\min_{p(t|x): I(T;Y) \geq I^*} I(X;T)$$

$$\mathcal{L}_{IB} = -I(Y;T) + \beta I(X;T) + \gamma \sum_{i=1}^n I(T_i; X | \hat{X}_i) \quad (2)$$

, where  $T_i$  is the tag associated with the  $i$ -th word,  $X_i$  is the ELMo embedding and  $\hat{X}_i$  is the type embedding (typical embedding)

We extend the original IB objective (1) and add terms to control the context-sensitivity of the extracted tags.

We instantiate the variational IB estimation method for each of

the Decoder / the Token Encoder / the Type Encoder

$I(X; T)$  – the Token Encoder  $p_{\theta}(t|x)$

$$\mathbb{E}_{x,t} \left[ \log \frac{p_{\theta}(t|x)}{p_{\theta}(t)} \right] = \mathbb{E}_x \left[ \mathbb{E}_{t \sim p_{\theta}(t|x)} \left[ \log \frac{p_{\theta}(t|x)}{p_{\theta}(t)} \right] \right].$$

The outer expectation is over the true distribution of sentences  $x$   
 $\rightarrow$  an empirical estimate

To estimate the inner expectation, we could sample, drawing taggings  $t$

The troublesome term is  $\hat{p}_{\theta}(t) = \mathbb{E}_{x'} [p_{\theta}(t | x')]$ ,

$$\begin{aligned} & \overbrace{\mathbb{E}_x \left[ \mathbb{E}_{t \sim p_{\theta}(t|x)} \left[ \log \frac{p_{\theta}(t|x)}{r_{\psi}(t)} \right] \right]}^{\text{upper bound}} - \overbrace{\mathbb{E}_x \left[ \mathbb{E}_{t \sim p_{\theta}(t|x)} \left[ \log \frac{p_{\theta}(t|x)}{p_{\theta}(t)} \right] \right]}^{I(X;T)} \\ &= \mathbb{E}_x [\text{KL}(p_{\theta}(t) || r_{\psi}(t))] \geq 0 \end{aligned}$$

## **$I(T; X)$ – Two Token Encoder Architectures**

To obtain continuous tags, define  $p_\theta$  such that  $t_i$  is Gaussian-distributed

To obtain discrete tags, define  $p_\theta$  such that  $t_i$  follows a softmax distribution

$I(T; X | \hat{X})$  – the Type Encoder  $s_{\xi}(t_i | \hat{x}_i)$

We were concerned that it might not be interpretable as a tag of word  $I$  specifically  
 → no guarantee that contextual info. came from word  $i$  not its neighbors

Also we do want tag to consider context

→ use a word type embedding that does not depend on context – ELMo's level-0 embedding of word

$$\begin{aligned}
 & \overbrace{\mathbb{E}_x \left[ \mathbb{E}_{t_i \sim p_{\theta}(t_i | x)} \left[ \log \frac{p_{\theta}(t_i | x)}{s_{\xi}(t_i | \hat{x}_i)} \right] \right]}^{\text{upper bound}} - \overbrace{\mathbb{E}_x \left[ \mathbb{E}_{t_i \sim p_{\theta}(t_i | x)} \left[ \log \frac{p_{\theta}(t_i | x)}{p_{\theta}(t_i | \hat{x}_i)} \right] \right]}^{I(T_i; X | \hat{X}_i)} \\
 &= \mathbb{E}_x [\text{KL}(p_{\theta}(t_i | \hat{x}_i) || s_{\xi}(t_i | \hat{x}_i))] \geq 0
 \end{aligned}$$



$I(Y; T)$  – the Decoder  $q_\phi(y|t)$

$$\begin{aligned} & \overbrace{\mathbb{E}_{y, t \sim p_\theta} \left[ \log \frac{p_\theta(y|t)}{p(y)} \right]}^{I(Y; T)} - \overbrace{\mathbb{E}_{y, t \sim p_\theta} \left[ \log \frac{q_\phi(y|t)}{p(y)} \right]}^{\text{lower bound}} \\ &= \mathbb{E}_{t \sim p_\theta} [\text{KL}(p_\theta(y | t) || q_\phi(y | t))] \geq 0 \end{aligned}$$

## Architecture

This parser uses a Bi-LSTM to extract features from compressed tags or vectors and assign scores to each tree edge

During training, the decoder compute only an approximation to  $q_\phi(y|t)$  for the golden tree  $y$

# 3 Training and Inference



## Objective Function

$$\mathcal{L}_{IB} = -I(Y; T) + \beta I(X; T) + \gamma \sum_{i=1}^n I(T_i; X | \hat{X}_i) \quad (2)$$

$$\mathbb{E}_{x,y} \left[ \mathbb{E}_{t \sim p_{\theta}(t|x)} [-\log q_{\phi}(y|t)] + \beta \text{KL}(p_{\theta}(t|x) || r_{\psi}(t)) + \gamma \sum_{i=1}^n \text{KL}(p_{\theta}(t_i | x) || s_{\xi}(t_i | \hat{x}_i)) \right] \quad (3)$$

## Details

We must apply the re-parametrization trick to backpropagate  
We use the Gumbel-softmax variant for discrete  $t$

## Objective Function

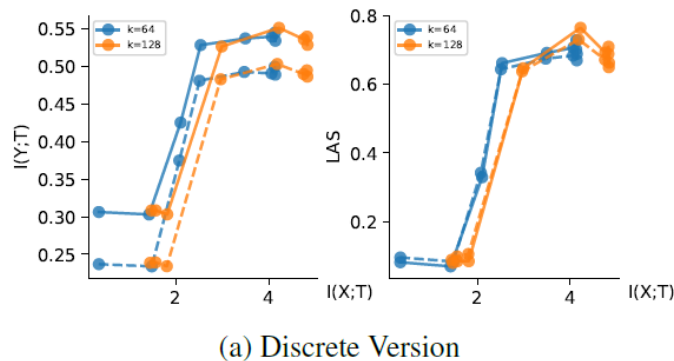
We will examine our compressed tags on a subset of **Universal Dependencies(UD)**, a collection of dependency treebanks across 76 languages using the same POS tags and dependency labels

We alternate between improving the model on even epochs and the variational distributions on odd epochs

1. We show the relationship between  $I(Y;T)$  and  $I(X;T)$  on English
2. Across 9 languages, we study how our automatic tags correlates with gold POS tags
3. We also show how decreasing  $\beta$  gradually refines the automatic discrete tag set

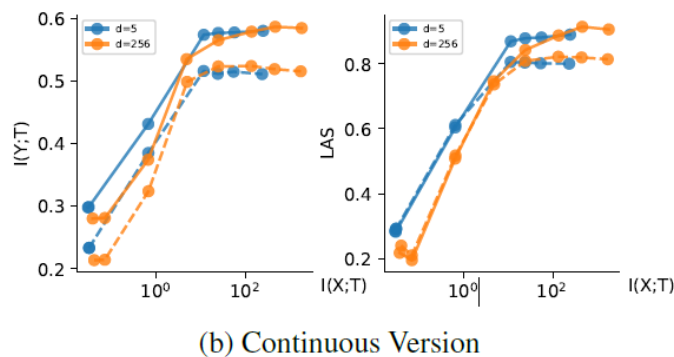
1. We show the relationship between  $I(Y;T)$  and  $I(X;T)$  on English
2. Across 9 languages, we study how our automatic tags correlates with gold POS tags
3. We also show how decreasing  $\beta$  gradually refines the automatic discrete tag set

# 1. We show the relationship between $I(Y;T)$ and $I(X;T)$ on English



Diminishing returns

after some point the additional info. does not contribute much to predicting  $Y$



The simpler, the better

at each level of  $I(X, T)$ , low-dimensional one perform on par with high-dimensional one

## 2. Across 9 languages, we study how our automatic tags correlates with gold POS tags

### Continuous Version

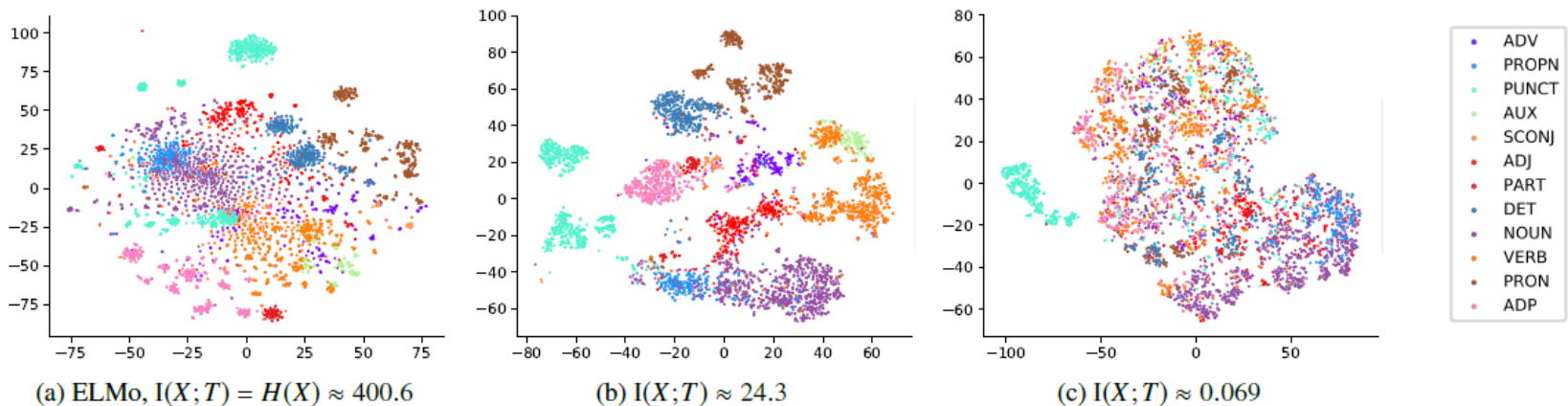


Figure 3: t-SNE visualization of VIB model ( $d = 256$ ) on the projected space of the continuous tags. Each marker in the figure represents a word token, colored by its gold POS tag. This series of figures (from left to right) shows a progression from no compression to moderate compression and to too-much compression.

- Note that the gold POS tags were not used in training
- An interesting observation is that NOUN and PROPN
- heavier compression squeezes out information

## 2. Across 9 languages, we study how our automatic tags correlates with gold POS tags

### Continuous Version

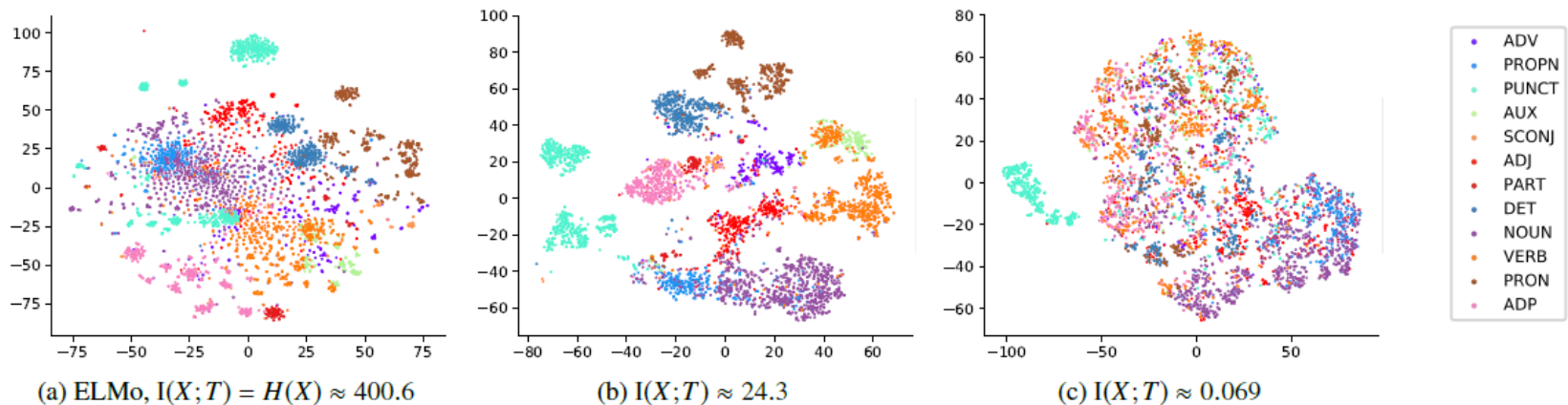


Figure 3: t-SNE visualization of VIB model ( $d = 256$ ) on the projected space of the continuous tags. Each marker in the figure represents a word token, colored by its gold POS tag. This series of figures (from left to right) shows a progression from no compression to moderate compression and to too-much compression.

- Note that the gold POS tags were not used in training
- An interesting observation is that NOUN and PROPN
- heavier compression squeezes out information



### 3. We also show how decreasing $\beta$ gradually refines the automatic discrete tag set Discrete Version

Deterministic Annealing, a method that gradually decrease  $\beta$  during training

→ results hierarchical structure refelects properties of English syntax

i.e. POS → (anaphors / Possesive pronouns)

# 6 Engineering Evaluation



Learning how to compress ELMo's tags for a given task is a fast alternative to fine-tuning all the ELMo parameters.

→ With on a single GPU, it is able to train on 10K sents. in 100 sec per epoch

Training our compression method does improve our generalization performance

Models	Arabic	Hindi	English	French	Spanish	Portuguese	Russian	Chinese	Italian
Iden	0.751	<b>0.870</b>	0.824	0.784	0.808	0.813	0.783	0.709	<b>0.863</b>
PCA	0.743	<b>0.866</b>	0.823	0.749	0.802	0.808	0.777	0.697	0.857
MLP	0.759	<b>0.871</b>	0.839	0.816	<b>0.835</b>	0.821	0.800	0.734	<b>0.867</b>
VIBc	<b>0.779</b>	<b>0.866</b>	<b>0.851</b>	<b>0.828</b>	<b>0.837</b>	<b>0.836</b>	<b>0.814</b>	<b>0.754</b>	<b>0.867</b>
POS	0.652	0.713	0.712	0.718	<b>0.739</b>	<b>0.743</b>	<b>0.662</b>	0.510	0.779
VIBd	<b>0.672</b>	<b>0.736</b>	<b>0.742</b>	<b>0.723</b>	<b>0.725</b>	0.710	<b>0.651</b>	<b>0.591</b>	<b>0.781</b>

Table 2: Parsing accuracy of 9 languages (LAS). Black rows use continuous tags; gray rows use discrete tags (which does worse). In each column, the best score for each color is boldfaced, along with all results of that color that are not significantly worse (paired permutation test,  $p < 0.05$ ). These results use only ELMo layer 1; results from all layers are shown in Table 3 in the appendix, for both LAS and UAS metrics.

task labels

a protected attribute

classifier

attacker

a representation

document

Goal : We want decision  $y = f(x)$  to be  
oblivious to  $z$

## Attacker

After the classifier  $c(h(x))$  is fully trained,  
we use the encoder  $h \rightarrow z$

## Definition

A protected attribute has *leaked*  
if we can train a classifier  $h \rightarrow z$

A protected attribute has *guarded*  
if we cannot train it

## Corpus - Twitter messages

### 1. DIAL

- Task : binary emoji-based **sentiment** and binary **tweet-mention** prediction  
Sentiment : Positive v Negative  
Mention : conversational v non-conversational
- Protected : race of authors  
Race : Light(Standard American English) v Dark (Else)

### 2. PAN 16

- Task : to classify a given **dependency relation** between two tokens
- Protected : Age and gender  
Age : (18-34) v (35+)  
gender : male v female

We collected 160K for training and 10K for development from each.  
Each split is balanced with respect to both the main and the protected labels.