

# Deep Clustering of Text Representations for Supervision-Free Probing of Syntax

Vikram Gupta<sup>1</sup>, Haoyue Shi<sup>2</sup>, Kevin Gimpel<sup>2</sup>, Mrinmaya Sachan<sup>3</sup>

<sup>1</sup> ShareChat, India <sup>2</sup> Toyota Technological Institute at Chicago <sup>3</sup> Department of Computer Science, ETH Zurich  
vikramgupta@sharechat.co, freda@ttic.edu, kgimpel@ttic.edu, mrinmaya.sachan@inf.ethz.ch

## Abstract

We explore deep clustering of text representations for unsupervised model interpretation and induction of syntax. As these representations are high-dimensional, out-of-the-box methods like KMeans do not work well. Thus, our approach jointly transforms the representations into a lower-dimensional cluster-friendly space and clusters them. We consider two notions of syntax: *part of speech induction* (POSI) and *constituency labelling* (CoLab) in this work. Interestingly, we find that Multilingual BERT (mBERT) contains surprising amount of syntactic knowledge of English; possibly even as much as English BERT (E-BERT). Our model can be used as a supervision-free probe which is arguably a less-biased way of probing. We find that unsupervised probes show benefits from higher layers as compared to supervised probes. We further note that our unsupervised probe utilizes E-BERT and mBERT representations differently, especially for POSI. We validate the efficacy of our probe by demonstrating its capabilities as a unsupervised syntax induction technique. Our probe works well for both syntactic formalisms by simply adapting the input representations. We report competitive performance of our probe on 45-tag English POSI, state-of-the-art performance on 12-tag POSI across 10 languages, and competitive results on CoLab. We also perform zero-shot syntax induction on resource impoverished languages and report strong results.

## 1 Introduction

Contextualized text representations (Peters et al. 2018a; Devlin et al. 2019) have been used in many NLP problems such as part-of-speech (POS) tagging, syntactic parsing (Kim et al. 2020; Kitaev and Klein 2018; Zhou and Zhao 2019), coreference resolution (Lee, He, and Zettlemoyer 2018; Joshi et al. 2019) and text classification (Minnae et al. 2021; Gupta 2021) often leading to significant improvements. Recent works have shown that these representations encode linguistic information including POS (Blinkov et al. 2017), morphology (Peters et al. 2018a), and syntactic structure (Linzen, Dupoux, and Goldberg 2016; Peters et al. 2018b; Tenney, Das, and Pavlick 2019; Hewitt and Manning 2019). While there has been a lot of focus on using contextualized representations in supervised settings,

the efficacy of these representations for unsupervised learning is not well explored<sup>1</sup>. Most of the recent work in “probing” contextual representations have focused on building supervised classifiers and using accuracy to interpret these representations. This has led to a debate as it is not clear if the supervised probe is probing the model or trying to solve the task (Hewitt and Manning 2019; Pimentel et al. 2020).

Thus, we explore a new clustering-based approach to probe contextualized text representations. Our probe allows for studying text representations with relatively less task-specific transformations due to the absence of supervision. Thus, our approach is arguably a less biased way to discover linguistic structure than supervised probes (Hewitt and Manning 2019; Pimentel et al. 2020). Our work is similar in spirit to (Wu et al. 2020; Zhou and Srikumar 2021) who also investigate supervision/parameter free probing methods. We focus on two syntactic formalisms: part-of-speech induction (POSI) and constituency labelling (CoLab), and explore the efficacy of contextualized representations towards encoding syntax in an unsupervised manner. We investigate the research question: *Do contextualized representations encode information for unsupervised syntax induction? How do these perform on POSI, which has been traditionally solved using smaller context windows and morphology and span-based CoLab?*

For both formalisms, we find that naively clustering text representations does not perform well. We speculate that this is because contextualized text representations are high-dimensional and not very friendly to existing clustering approaches. Thus, we develop a deep clustering approach (Xie, Girshick, and Farhadi 2016; Ghasedi Dizaji et al. 2017; Chang et al. 2017; Yang, Parikh, and Batra 2016; Yang et al. 2017) which transforms these representations into a lower dimensional, clustering friendly latent space. This transformation is learnt jointly with clustering using a combination of reconstruction and clustering objectives. The procedure iteratively refines the transformation and the clustering using an auxiliary target distribution derived from the current soft

<sup>1</sup>Some recent work such as DIORA (Drozdo et al. 2019b,a) has explored specialized methods for unsupervised discovery and representation of constituents using ELMo (Peters et al. 2018a). (Jin et al. 2019) used ELMo with a normalizing flow model while (Cao, Kitaev, and Klein 2020) used RoBERTa (Liu et al. 2019b) for unsupervised constituency parsing.

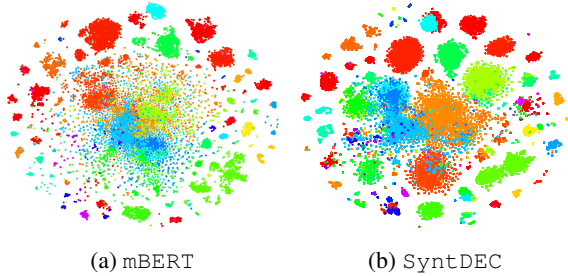


Figure 1: t-SNE visualization of mBERT embeddings (clustered using kMeans) and SyntDEC (our probe) embeddings of tokens from Penn Treebank. Colors correspond to ground truth POS tags.

clustering. As this process is repeated, it gradually improves the transformed representations as well as the clustering. We show a t-SNE visualization of mBERT embeddings and embeddings learned by our deep clustering probe (SyntDEC) in Figure 1.

We further explore architectural variations such as pre-trained subword embeddings from *fastText* (Joulin et al. 2017), a continuous bag of words (CBOW) loss (Mikolov et al. 2013), and span representations (Toshniwal et al. 2020) to incorporate task-dependent information into the latent space and observe significant improvements. It is important to note that we do not claim that clustering contextualized representations is the optimal approach for POSI as representations with short context (Lin et al. 2015), (He, Neubig, and Berg-Kirkpatrick 2018) and word-based POSI (Yatbaz, Sert, and Yuret 2012) have shown best results. Our approach explores the potential of contextualized representations for unsupervised induction of syntax and acts as an unsupervised probe for interpreting these representations. Nevertheless, we report competitive many-to-one (M1) accuracies for POSI on 45-tag Penn Treebank WSJ dataset as compared to specialized state-of-the-art approaches in the literature (He, Neubig, and Berg-Kirkpatrick 2018) and improve upon the state-of-the-art on 12 tag universal treebank dataset across multiple languages (Stratos, Collins, and Hsu 2016; Stratos 2019). We further show that our approach can be used in a zero-shot crosslingual setting where a model trained on one language can be used for evaluation in another language without using training data from the other language. We observe impressive crosslingual POSI performance, showcasing the representational power of mBERT, especially when the languages are related. Our method also achieves competitive results on CoLab, outperforming the initial DIORA approach (Drozdov et al. 2019b) and performing comparably to recent DIORA variants (Drozdov et al. 2019a) which incorporate more complex methods such as latent chart parsing and discrete representation learning. In contrast to specialized state-of-the-art methods for syntax induction, our framework is more general as it demonstrates good performance for both CoLab and POSI by simply adapting the input representations.

We further investigate the effectiveness of multilingual BERT (mBERT) (Devlin et al. 2019) for POSI across multiple languages and CoLab in English and see improve-

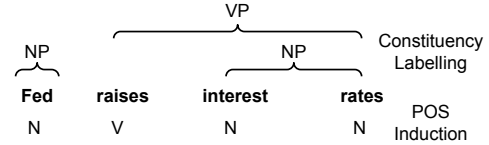


Figure 2: An illustration of POSI and CoLab formalisms.

ment in performance by using mBERT for both tasks even in English. This is in contrast with the supervised experiments where both mBERT and E-BERT perform competitively. In contrast to various supervised probes in the literature (Liu et al. 2019a; Tenney, Das, and Pavlick 2019), our unsupervised probe finds that syntactic information is captured in higher layers on average than what was previously reported (Tenney, Das, and Pavlick 2019). Upon further layer-wise analysis of the two probes, we find that while supervised probes show that all layers of E-BERT contain syntactic information fairly uniformly, middle layers lead to a better performance on the investigated syntactic tasks with our unsupervised probe.

## 2 Problem Definition

We consider two syntax induction problems in this work:

1. **Part-of-speech induction (POSI)**: determining part of speech of words in a sentence.
2. **Constituency label induction (CoLab)**: determining the constituency label for a given constituent (span of contiguous tokens).<sup>2</sup>

Figure 2 shows an illustration for the two tasks. In order to do well, both tasks require reasoning about the context. This motivates us to use contextualized representations, which have shown an ability to model such information effectively. Letting  $[m]$  denote  $\{1, 2, \dots, m\}$ , we model unsupervised syntax induction as the task of learning a mapping function  $C : X \rightarrow [m]$ . For POSI,  $X$  is the set of word tokens in the corpus and  $m$  is the number of part-of-speech tags.<sup>3</sup> For CoLab,  $X$  is the set of constituents across all sentences in the corpus and  $m$  is the number of constituent labels. For each element  $x \in X$ , let  $c(x)$  denote the context of  $x$  in the sentence containing  $x$ . The number  $m$  of true clusters is assumed to be known. For CoLab, we also assume gold constituent spans from manually annotated constituency parse trees, focusing only on determining constituent labels, following Drozdov et al. (2019a).

## 3 Proposed Method

We address unsupervised syntax induction via clustering, where  $C$  defines a clustering of  $X$  into  $m$  clusters. We define a deep embedded clustering framework and modify it to support common NLP objectives such as continuous bag of words (Mikolov et al. 2013). Our framework jointly trans-

<sup>2</sup>Note that it is not necessary for constituents to be contiguous, but we only consider contiguous constituents for simplicity.

<sup>3</sup> $X$  is distinct from the corpus vocabulary; in POSI, we tag each word token in each sentence with a POS tag.

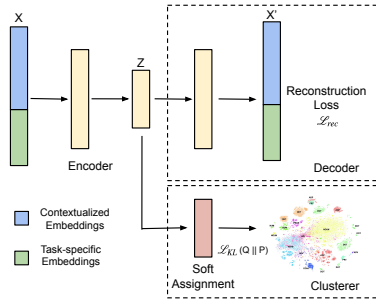


Figure 3: An illustration of our SyntDEC model.

forms the text representations into a lower-dimensions and learns the clustering parameters in an end-to-end setup.

### Deep Clustering

Unlike traditional clustering approaches that work with fixed, and often hand-designed features, deep clustering (Xie, Girshick, and Farhadi 2016; Ghasedi Dizaji et al. 2017; Jiang et al. 2016; Chang et al. 2017; Yang, Parikh, and Batra 2016; Yang et al. 2017) transforms the data  $X$  into a latent feature space  $Z$  with a mapping function  $f_\theta : X \rightarrow Z$ , where  $\theta$  are learnable parameters. The dimensionality of  $Z$  is typically much smaller than  $X$ . The datapoints are clustered by simultaneously learning a clustering  $\hat{C} : Z \rightarrow [m]$ . While  $C$  might have been hard to learn directly (due to the high dimensionality of  $X$ ), learning  $\hat{C}$  may be easier.

**Deep Embedded Clustering:** We draw on a particular deep clustering approach: Deep Embedded Clustering (DEC; Xie, Girshick, and Farhadi 2016). Our approach consists of two stages: (a) a *pretraining* stage, and (b) a *joint representation learning and clustering* stage. In the *pretraining* stage, a mapping function  $f_\theta$  is pretrained using a stacked autoencoder (SAE). The SAE learns to reconstruct  $X$  through the bottleneck  $Z$ , i.e.,  $X \xrightarrow{\text{encoder}} Z \xrightarrow{\text{decoder}} X'$ . We use mean squared error (MSE) as the reconstruction loss:

$$\mathcal{L}_{rec} = \|X - X'\|^2 = \sum_{x \in X} \|x - x'\|^2$$

The encoder parameters are used to initialize the mapping function  $f_\theta$ . In the *joint representation learning and clustering* stage, we finetune the encoder  $f_\theta$  trained in the pretraining stage to minimize a clustering loss  $\mathcal{L}_{KL}$ . The goal of this step is to learn a latent space that is amenable to clustering. We learn a set of  $m$  cluster centers  $\{\mu_i \in Z\}_{i=1}^m$  of the latent space  $Z$  and alternate between computing an *auxiliary target distribution* and minimizing the Kullback-Leibler (KL) divergence. First, a soft cluster assignment is computed for each embedded point. Then, the mapping function  $f_\theta$  is refined along with the cluster centers by learning from the assignments using an auxiliary target distribution. This process is repeated. The soft assignment is computed via the Student's  $t$ -distribution. The probability of assigning data point  $i$  to cluster  $j$  is denoted  $q_{ij}$  and defined:

$$q_{ij} = \frac{(1 + \|z_i - \mu_j\|^2/\nu)^{-\frac{\nu+1}{2}}}{\sum_{j'} (1 + \|z_i - \mu_{j'}\|^2/\nu)^{-\frac{\nu+1}{2}}}$$

where  $\nu$  is set to 1 in all experiments. Then, a cluster assignment hardening loss (Xie, Girshick, and Farhadi 2016) is used to make these soft assignment probabilities more peaked. This is done by letting cluster assignment probability distribution  $q$  approach a more peaked auxiliary (target) distribution  $p$ :

$$p_{ij} = \frac{q_{ij}^2/n_j}{\sum_{j'} q_{ij'}^2/n_{j'}} \quad n_j = \sum_i q_{ij}$$

By squaring the original distribution and then normalizing it, the auxiliary distribution  $p$  forces assignments to have more peaked probabilities. This aims to improve cluster purity, put emphasis on data points assigned with high confidence, and to prevent large clusters from distorting the latent space. The divergence between the two probability distributions is formulated as the Kullback-Leibler divergence:

$$\mathcal{L}_{KL} = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

The representation learning and clustering model is learned end-to-end.

### SyntDEC: DEC for Syntax Induction

We further modify DEC for syntax induction:

**a) CBoW autoencoders:** While DEC uses a conventional autoencoder, i.e., the input and output are the same, we modify it to support the continuous bag of words (CBoW) objective (Mikolov et al. 2013). This helps focus the low-dimensional representations to focus on context words, which are expected to be helpful for POSI. In particular, given a set of tokens  $c(x)$  that defines the context for an element  $x \in X$ , CBoW combines the distributed representations of tokens in  $c(x)$  to predict the element  $x$  in the middle. See Appendix A for an illustration.

**b) Finetuning with reconstruction loss:** We found that in the clustering stage, finetuning with respect to the KL divergence loss alone easily leads to trivial solutions where all points map to the same cluster. To address this, we add the reconstruction loss as a regularization term. This is in agreement with subsequent works in deep clustering (Yang et al. 2017). Instead of solely minimizing  $\mathcal{L}_{KL}$ , we minimize

$$\mathcal{L}_{total} = \mathcal{L}_{KL} + \lambda \mathcal{L}_{rec} \quad (1)$$

in the clustering stage, where  $\lambda$  is a hyperparameter denoting the weight of the reconstruction loss.

**c) Contextualized representations:** We represent linguistic elements  $x$  by embeddings extracted from pretrained networks like BERT (Devlin et al. 2019), SpanBERT (Joshi et al. 2020), and multilingual BERT (Devlin et al. 2019). All of these networks are multi-layer architectures. Thus, we average the embeddings across the various layers. We experimented with different layer combinations but found the average was the best solution for these tasks. We averaged the embeddings of the subword units to compute word embeddings.<sup>4</sup> For CoLab, we represent spans by concatenating the representations of the end points (Toshniwal et al. 2020).

<sup>4</sup>In our preliminary experiments, we also tried other pooling mechanisms such as min/max pooling over subwords, but average performed the best among all of them.

**d) Task-specific representations:** Previous work in unsupervised syntax induction has shown the value of task-specific features. In particular, a number of morphological features based on prefixes and suffixes and spelling cues like capitalization have been used in unsupervised POSI works (Tseng, Jurafsky, and Manning 2005; Stratos 2019; Yatbaz, Sert, and Yuret 2012). In our POSI experiments, we incorporate these morphological features by using word representations from `fastText` (Joulin et al. 2017). We concatenate `fastText` embeddings of the trailing trigram of each word with contextualized representations before passing them as input to `SyntDEC`.

## 4 Experimental Details

**Datasets:** We evaluate our approach for POSI on two datasets: 45-tag Penn Treebank Wall Street Journal (WSJ) dataset (Marcus, Santorini, and Marcinkiewicz 1993) and multilingual 12-tag datasets drawn from the universal dependencies project (Nivre et al. 2016). WSJ dataset has approximately one million words tagged with 45 part of speech tags. For multilingual experiments, we use the 12-tag universal treebank v2.0 dataset which consists of corpora from 10 languages.<sup>5</sup> The words in this dataset have been tagged with 12 universal POS tags (McDonald et al. 2013). For CoLab, we follow (Drozdov et al. 2019a) and evaluate on the WSJ test set. For POSI, as per the standard practice (Stratos 2019), we use the complete dataset (train + val + test) for training as well as evaluation. However, for CoLab, we use the train set to train our model and the test set for reporting results, following Drozdov et al. (2019a).

**Evaluation Metrics:** For POSI, we use the standard measures of many-to-one (M1; Johnson 2007) accuracy and V-Measure (Rosenberg and Hirschberg 2007). For CoLab, we use F1 score following Drozdov et al. (2019a), ignoring spans which have only a single word and spans with the “TOP” label. In addition to F1, we also report M1 accuracy for CoLab to show the clustering performance more naturally and intuitively.

**Training Details:** Similar to Xie, Girshick, and Farhadi (2016), we use greedy layerwise pretraining (Bengio et al. 2007) for initialization. New hidden layers are successively added to the autoencoder, and the layers are trained to denoise output of the previous layer. After layerwise pretraining, we train the autoencoder end-to-end and leverage the trained `SyntDEC` encoder (Section 3). K-Means is used to initialize cluster means and assignments. `SyntDEC` is trained end-to-end with the reconstruction and clustering losses. More details are in the appendix.

## 5 Part of Speech Induction (POSI)

**45-Tag Penn Treebank WSJ:** In Table 1, we evaluate the performance of contextualized representations and our probe on the 45-tag Penn Treebank WSJ dataset. KMeans clustering over the `mBERT` embeddings improves upon Brown clustering (Brown et al. 1992) (as reported by Stratos, 2019) and Hidden Markov Models (Stratos, Collins, and Hsu

Method	M1	VM
<code>SyntDEC_Morph</code>	79.5 ( $\pm 0.9$ )	73.9 ( $\pm 0.7$ )
<code>SyntDEC</code>	77.6 ( $\pm 1.5$ )	72.5 ( $\pm 0.9$ )
SAE	75.3 ( $\pm 1.4$ )	69.9 ( $\pm 0.9$ )
KMeans	72.4 ( $\pm 2.9$ )	-
Brown et al. (1992)	65.6 ( $\pm$ NA)	-
Stratos, Collins, and Hsu (2016)	67.7 ( $\pm$ NA)	-
Berg-Kirkpatrick et al. (2010)	74.9 ( $\pm 1.5$ )	-
Blunsom and Cohn (2011)	77.5 ( $\pm$ NA)	69.8
Stratos (2019)	78.1 ( $\pm 0.8$ )	-
Tran et al. (2016)	79.1 ( $\pm$ NA)	71.7 ( $\pm$ NA)
Yuret, Yatbaz, and Sert (2014)	79.5 ( $\pm 0.3$ )	69.1 ( $\pm 2.7$ )
Yatbaz, Sert, and Yuret (2012) (word-based)	80.2 ( $\pm 0.7$ )	72.1 ( $\pm 0.4$ )
He, Neubig, and Berg-Kirkpatrick (2018)	80.8 ( $\pm 1.3$ )	74.1 ( $\pm 0.7$ )

Table 1: Many-to-one (M1) accuracy and V-Measure (VM) of POSI on the 45-tag Penn Treebank WSJ dataset for 10 random runs. `mBERT` is used in all of our experiments (upper part of the table).

2016) based approach, showing that `mBERT` embeddings encode syntactic information. The stacked autoencoder, SAE (trained during pretraining stage), improves upon the result of KMeans by nearly 3 points, which demonstrates the effectiveness of transforming the `mBERT` embeddings to lower dimensionality using an autoencoder before clustering. Our method (`SyntDEC`) further enhances the result and shows that transforming the pretrained `mBERT` embeddings using clustering objective helps to extract syntactic information more effectively. When augmenting the `mBERT` embeddings with morphological features (`SyntDEC_Morph`), we improve over Stratos (2019) and (Tran et al. 2016). We also obtain similar M1 accuracy with higher VM as compared to (Yuret, Yatbaz, and Sert 2014).

**Morphology:** We also note that M1 accuracy of Tran et al. (2016) and Stratos (2019) drop by nearly 14 points in absence of morphological features, while `SyntDEC` degrades by 2 points. This trend suggests that `mBERT` representations encode the morphology to some extent.

Yatbaz, Sert, and Yuret (2012) are not directly comparable to our work as they performed word-based POSI which attaches same tag to all the instances of the word, while all the other works in Table 1 perform token-based POSI. They use task-specific hand-engineered rules like presence of hyphen, apostrophe etc. which might not translate to multiple languages and tasks. (He, Neubig, and Berg-Kirkpatrick 2018) train a POSI specialized model with Markov syntax model and short-context word embeddings and report current SOTA on POSI. In contrast to their method, `SyntDEC` is fairly task agnostic.

**12-Tag Universal Treebanks:** In Table 2, we report M1 accuracies on the 12-tag datasets averaged over 5 random runs. Across all languages, we report SOTA results and find an improvement on average over the previous best method (Stratos 2019) from 71.4% to 75.7%. We also note improvements of `SyntDEC` over SAE (70.9% to 75.7%) across languages, which reiterates the importance of finetuning representations for clustering. Our methods yield larger gains on this coarse-grained 12 tag POSI task as compared to the fine-grained 45 tag POSI task, and we hope to explore the reasons for this in future work.

**Ablation Studies:** Next, we study the impact of our choices

<sup>5</sup>We use v2.0 in order to compare to Stratos (2019).

	de	en	es	fr	id	it	ja	ko	pt-br	sv	Mean
SAE	74.8 ( $\pm 1.5$ )	70.7 ( $\pm 2.2$ )	71.1 ( $\pm 2.4$ )	66.7 ( $\pm 1.9$ )	75.4 ( $\pm 1.6$ )	66.2 ( $\pm 3.3$ )	82.1 ( $\pm 0.9$ )	65.4 ( $\pm 1.7$ )	75.1 ( $\pm 4.1$ )	61.6 ( $\pm 2.6$ )	70.9
SyntDEC	<b>81.5</b> ( $\pm 1.8$ )	<b>76.5</b> ( $\pm 1.1$ )	<b>78.9</b> ( $\pm 1.9$ )	<b>70.7</b> ( $\pm 3.9$ )	<b>76.8</b> ( $\pm 1.1$ )	<b>71.7</b> ( $\pm 3.3$ )	<b>84.7</b> ( $\pm 1.2$ )	<b>69.7</b> ( $\pm 1.5$ )	<b>77.7</b> ( $\pm 2.1$ )	<b>68.8</b> ( $\pm 3.9$ )	<b>75.7</b>
Stratos (2019)	75.4 ( $\pm 1.5$ )	73.1 ( $\pm 1.7$ )	73.1 ( $\pm 1.0$ )	70.4 ( $\pm 2.9$ )	73.6 ( $\pm 1.5$ )	67.4 ( $\pm 3.3$ )	77.9 ( $\pm 0.4$ )	65.6 ( $\pm 1.2$ )	70.7 ( $\pm 2.3$ )	67.1 ( $\pm 1.5$ )	71.4
Stratos, Collins, and Hsu (2016)	63.4	71.4	74.3	71.9	67.3	60.2	69.4	61.8	65.8	61.0	66.7
Berg-Kirkpatrick et al. (2010)	67.5 ( $\pm 1.8$ )	62.4 ( $\pm 3.5$ )	67.1 ( $\pm 3.1$ )	62.1 ( $\pm 4.5$ )	61.3 ( $\pm 3.9$ )	52.9 ( $\pm 2.9$ )	78.2 ( $\pm 2.9$ )	60.5 ( $\pm 3.6$ )	63.2 ( $\pm 2.2$ )	56.7 ( $\pm 2.5$ )	63.2
Brown et al. (1992)	60.0	62.9	67.4	66.4	59.3	66.1	60.3	47.5	67.4	61.9	61.9

Table 2: M1 accuracy and standard deviations on the 12-tag universal treebank dataset averaged over 5 random runs. mBERT is used for all of our experiments (upper part of the table). The number of epochs are proportional to the number of samples and the M1 accuracy corresponding to the last epoch is reported.

	Method	M1
E-BERT	KMeans	69.1 ( $\pm 0.9$ )
	SAE	71.6 ( $\pm 2.3$ )
	CBoW	73.8 ( $\pm 0.7$ )
	SyntDEC (SAE)	72.7 ( $\pm 1.2$ )
	SyntDEC (CBoW)	<b>74.4</b> ( $\pm 0.6$ )
mBERT	KMeans	72.4 ( $\pm 2.9$ )
	SAE	75.3 ( $\pm 1.4$ )
	CBoW	75.1 ( $\pm 0.3$ )
	SyntDEC (SAE)	<b>77.8</b> ( $\pm 1.4$ )
	SyntDEC (CBoW)	75.9 ( $\pm 0.3$ )

Table 3: Comparison of E-BERT and mBERT on the 45-tag POSI task. We report *oracle* results in this table.

on the 45-tag WSJ dataset. Table 3 demonstrates that multilingual BERT (mBERT) is better than English BERT (E-BERT) across settings. For both mBERT and E-BERT, compressing the representations with SAE and finetuning using SyntDEC performs better than KMeans. Also, focusing the representations on the local context (CBoW) improves performance with E-BERT, though not with mBERT. In the appendix, we show the impact of using different types of fastText character embeddings and note the best results when we use embeddings of the last trigram of each word.

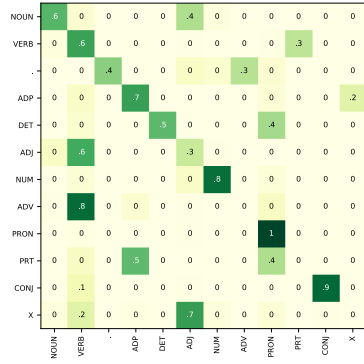
**Error Analysis:** We compared SyntDEC and KMeans (when both use mBERT) and found that SyntDEC does better on noun phrases and nominal tags. It helps alleviate confusion among fine-grained noun tags (e.g., NN vs. NNS), while also showing better handling of numerals (CD) and personal pronouns (PRP). However, SyntDEC still shows considerable confusion among fine-grained verb categories. For 12-tag experiments, we similarly found that SyntDEC outperforms KMeans for the majority of the tags, especially nouns and verbs, resulting in a gain of more than 20% in 1-to-1 accuracy. We further compare t-SNE visualizations of SyntDEC and mBERT embeddings and observe that SyntDEC embeddings show relatively compact clusters. Detailed results and visualizations (Figure 4) are shown in the appendix.

## 6 SyntDEC as an Unsupervised Probe

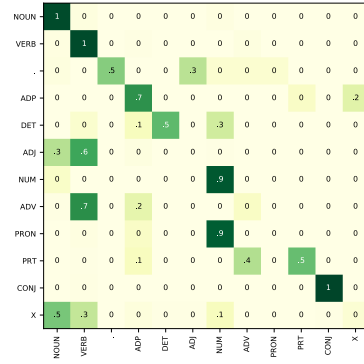
Next, we leverage SyntDEC as an unsupervised probe to analyse where syntactic information is captured in the pre-trained representations. Existing approaches to probing usually rely on supervised training of probes. However, as argued recently by (Zhou and Srikumar 2021), this can be unreliable. Our supervision-free probe arguably gets rid of any bias in interpretations due to the involvement of training data in probing. We compare our unsupervised probe to a reimplementation of the supervised shallow MLP based probe in Tenney, Das, and Pavlick (2019). Similar to their paper, we report `Expected Layer` under supervised and unsupervised settings for the two tasks in Figure 5. `Expected Layer` represents the average layer number in terms of incremental performance gains:  $E_{\Delta}[l] = \frac{\sum_{i=1}^L l * \Delta^{(l)}}{\sum_{i=1}^L \Delta^{(l)}}$ , where  $\Delta^{(l)}$  is the change in the performance metric when adding layer  $l$  to the previous layers. Layers are incrementally added from lower to higher layers. We use F1 and M1 score as the performance metric for supervised and unsupervised experiments respectively. We observe that:

1. `Expected Layer` as per the unsupervised probe (blue) is higher than the supervised probe (green) for both tasks and models showing that unsupervised syntax induction benefits more from higher layers.
2. There are larger differences between E-BERT and mBERT `Expected Layer` under unsupervised settings suggesting that our unsupervised probe utilizes mBERT and E-BERT layers differently than the supervised one. In supervised settings, both models show similar expected layers of 1.7 vs 1.8 for POSI and 3.1 vs 3.0 for CoLab. However, under unsupervised settings, the expected layer shows larger differences: 3.6 vs 5.2 for POSI and 4.3 vs 4.9 for CoLab.

In Figure 6, we further probe the performance of each layer individually by computing the F1 score for the supervised probe and the M1 score for the unsupervised probe. We observe noticeable improvement at Layer 1 for supervised POSI and Layer 1/4/6 for CoLab which also correlates with their respective `Expected Layer` values. For unsupervised settings, the improvements are more evenly shared across initial layers. Although F1 and M1 are not directly comparable, supervised performance is competitive even at



(a) mBERT: One-to-One accuracy: 54.4%



(b) SyntDEC: One-to-One accuracy: 65.9%

Figure 4: Comparison of confusion matrices of mBERT and SyntDEC for 12-tag experiments on English. One-to-one mapping is used to assign labels to clusters.

		Nearby						Distant			Mean
		en	de	sv	es	fr	pt	it	ko	id	ja
distance to en	0	0.36	0.4	0.46	0.46	0.48	0.50	0.69	0.71	0.71	-
Monolingual	76.5 ( $\pm 1.1$ )	81.5 ( $\pm 1.8$ )	68.8 ( $\pm 3.9$ )	78.9 ( $\pm 1.9$ )	70.7 ( $\pm 3.9$ )	77.7 ( $\pm 2.1$ )	71.7 ( $\pm 3.3$ )	69.7 ( $\pm 1.5$ )	76.8 ( $\pm 1.1$ )	84.7 ( $\pm 1.2$ )	75.7 -
Crosslingual	76.5 ( $\pm 1.1$ )	71.9 ( $\pm 1.5$ )	66.7 ( $\pm 1.9$ )	75.7 ( $\pm 1.4$ )	73.5 ( $\pm 1.1$ )	77.6 ( $\pm 1.1$ )	73.5 ( $\pm 1.2$ )	67.5 ( $\pm 0.9$ )	75.4 ( $\pm 1.7$ )	80.3 ( $\pm 1.3$ )	73.9 -

Table 4: POSI M1 for SyntDEC with mBERT on 12-tag universal treebank in monolingual and crosslingual settings. **Monolingual**: clusters are learned and evaluated on the same language. **Crosslingual**: clusters are learned on English and evaluated on all languages.

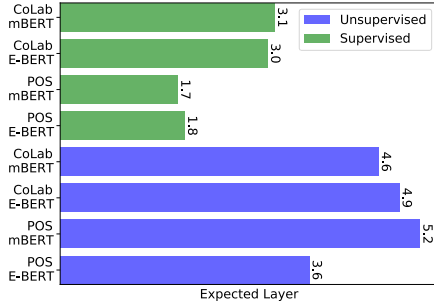
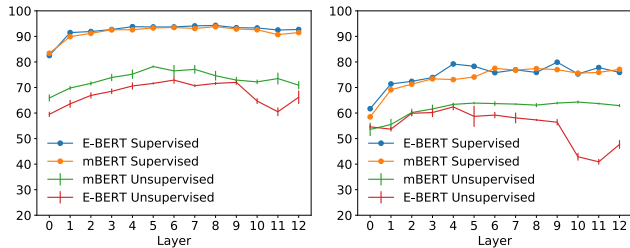


Figure 5: Expected Layer of POSI and CoLab under unsupervised SyntDEC (blue) and supervised settings (green) with E-BERT and mBERT representations.



(a) POSI

(b) CoLab

Figure 6: Comparison of M1/F1 measure for POSI and CoLab under unsupervised (SyntDEC) and supervised settings with mBERT and E-BERT representations.

higher layers while unsupervised performance drops. We present detailed results in the appendix.

## 7 Crosslingual POSI

Pires, Schlinger, and Garrette (2019); Wu and Dredze (2019); Snyder, Naseem, and Barzilay (2009) show that mBERT is effective at zero-shot crosslingual transfer. Inspired by this, we evaluate the crosslingual performance on 12-tag universal treebank (Table 4). In this task, we train the model on one language and evaluate it on another language without using training data from other language. The first row shows M1 accuracies when training and evaluating SyntDEC on the same language (monolingual). The second row shows M1 accuracies of the English-trained SyntDEC on other languages (crosslingual). In general, we find that clusters learned on a high-resource languages like English can be used for other languages. Similar to He et al. (2019), we use the distances of the languages with English to group languages as *nearby* or *distant*. The distance is calculated by accounting for syntactic, genetic, and geographic distances according to the URIEL linguistic database (Littell et al. 2017). Our results highlight the effectiveness of mBERT in crosslingual POSI. Even for Asian languages (ko, id, and ja), which have a higher distance from English, the performance is comparable across settings. For nearby languages, crosslingual SyntDEC performs well and even outperforms the monolingual setting for some languages.



Method	$F1_{\mu}$	$F1_{max}$	M1	VM
DIORA	62.5 ( $\pm 0.5$ )	63.4	-	-
DIORA <sub>CB</sub> (*)	64.5 ( $\pm 0.6$ )	65.5	-	-
DIORA <sub>CB</sub> <sup>*</sup> (*)	66.4 ( $\pm 0.7$ )	67.8	-	-
DIORA Baselines	E-BERT (**)	41.8	-	-
	ELMo (**)	58.5	-	-
	ELMo <sub>CI</sub> (**)	53.4	-	-
SyntDEC	E-BERT	60.8 ( $\pm 0.7$ )	62.7	75.4 ( $\pm 1.1$ )
	SpanBERT	61.3 ( $\pm 0.8$ )	63.3	75.9 ( $\pm 1.0$ )
	mBERT	64.0 ( $\pm 0.4$ )	64.6	79.6 ( $\pm 0.6$ )

Table 5: CoLab results on the WSJ test set using the gold parses over five random runs. Our models were trained for 15 epochs and results from the final epoch for each run are recorded. DIORA results are reported from Drozdov et al. (2019a). DIORA<sub>CB</sub> and DIORA<sub>CB</sub><sup>\*</sup> are fairly specialized models involving codebook learning (\*). We also report E-BERT and ELMo baselines from Drozdov et al. (2019a) (\*\*). We significantly outperform these previously reported E-BERT/ELMo baselines. Our results are not directly comparable to DIORA as it uses the WSJ dev set for tuning and early stopping whereas we do not.

## 8 Constituency Labelling (CoLab)

In Table 5, we report the F1 and M1 score of constituency labelling (CoLab) over the WSJ test set. We represent constituents by concatenating embeddings of the first and last words in the span (where word embeddings are computed by averaging corresponding subword embeddings). We observe improvement over DIORA (Drozdov et al. 2019b), a recent unsupervised constituency parsing model, and achieve competitive results to recent variants that improve DIORA with discrete representation learning (Drozdov et al. 2019a). Our model and the DIORA variants use gold constituents for these experiments. We compute F1 metrics for comparing with previous work but also report M1 accuracies. As with POSI, our results suggest that mBERT outperforms both SpanBERT and E-BERT for the CoLab task as well. We also note that SpanBERT performs better than E-BERT, presumably because SpanBERT seeks to learn span representations explicitly. In the Appendix (Table 7), we explore other ways of representing constituents and note that mean/max pooling followed by clustering does not perform well. Compressing and finetuning the mean-pooled representation using SyntDEC (SyntDEC\_Mean) is also suboptimal. We hypothesize that mean/max pooling results in a loss of information about word order in the constituent whereas the concatenation of first and last words retains this information. Even a stacked autoencoder (SAE) over the concatenation of first and last token achieves competitive results, but finetuning with SyntDEC improves the  $F1_{\mu}$  by nearly 4.5%. This demonstrates that for CoLab also, the transformation to lower dimensions and finetuning to clustering friendly spaces is important for achieving competitive performance.

## 9 Related Work

**Deep Clustering:** Unlike previous work where feature extraction and clustering were applied sequentially, deep clustering aims to jointly optimize for both by combining a clustering loss with the feature extraction. Various deep clustering methods have been proposed which primarily dif-

fer in their clustering approach: Yang et al. (2017) use KMeans, Xie, Girshick, and Farhadi (2016) use cluster assignment hardening, Ghasedi Dizaji et al. (2017) add a balanced assignments loss on top of cluster assignment hardening, Huang et al. (2014) introduce a locality-preserving loss and a group sparsity loss on the clustering, Yang, Parikh, and Batra (2016) use agglomerative clustering, and Ji et al. (2017) use subspace clustering. All of these approaches can be used to cluster contextualized representations, and future work may improve upon our results by exploring these approaches. The interplay between deep clustering for syntax and contextualized representations, has not previously been studied. In this paper, we fill this gap.

**Unsupervised Syntax Induction:** There has been a lot of work on unsupervised induction of syntax, namely, unsupervised constituency parsing (Klein and Manning 2002; Seginer 2007; Kim, Dyer, and Rush 2019) and dependency parsing (Klein and Manning 2004; Smith and Eisner 2006; Gillenwater et al. 2010; Spitzkovsky, Alshawhi, and Jurafsky 2013; Jiang, Han, and Tu 2016). While most prior work focuses on inducing *unlabeled* syntactic structures, we focus on inducing constituent labels while assuming the gold syntactic structure is available. This goal has also been pursued in prior work (Drozdov et al. 2019a; Jin and Schuler 2020). Compared to them, we present simpler models to induce syntactic labels directly from pretrained models via dimensionality reduction and clustering. Similar to us, (Li and Eisner 2019) also note gains for supervised NLP tasks upon reducing the representation dimension.

**Probing Pretrained Representations:** Recent analysis work (Liu et al. 2019a; Tenney et al. 2019; Jawahar, Sagot, and Seddah 2019, *inter alia*) has shown that pretrained language models encode syntactic information efficiently. Most of them train a supervised model using pretrained representations and labeled examples, and show that pretrained language models effectively encode part-of-speech and constituency information. In contrast to these works, we propose an unsupervised approach to probing which does not rely on any training data. (Zhou and Srikumar 2021) also pursue the same goals by studying the geometry of these representations.

## 10 Conclusion

In this work, we explored the problem of clustering text representations for model interpretation and induction of syntax. We observed that off-the-shelf methods like KMeans are sub-optimal as these representations are high dimensional and, thus, not directly suitable for clustering. Thus, we proposed a deep clustering approach which jointly transforms these representations into a lower-dimensional cluster friendly space and clusters them. Upon integration of a small number of task-specific features and use of multilingual representations, we find that our approach achieves competitive performance for unsupervised POSI and CoLab comparable to more complex methods in the literature. Finally, we also show that we can use the technique as a supervision-free approach to probe syntax in these representations and contrast our unsupervised probe with supervised ones.

## References

- Belinkov, Y.; Durrani, N.; Dalvi, F.; Sajjad, H.; and Glass, J. 2017. What do neural machine translation models learn about morphology? In *Proc. of ACL*.
- Bengio, Y.; Lamblin, P.; Popovici, D.; and Larochelle, H. 2007. Greedy layer-wise training of deep networks. In *Proc. of NeurIPS*.
- Berg-Kirkpatrick, T.; Bouchard-Côté, A.; DeNero, J.; and Klein, D. 2010. Painless unsupervised learning with features. In *Proc. of NAACL-HLT*.
- Blunsom, P.; and Cohn, T. 2011. A hierarchical Pitman-Yor process HMM for unsupervised part of speech induction. In *Proc. of 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*.
- Brown, P. F.; Della Pietra, V. J.; Desouza, P. V.; Lai, J. C.; and Mercer, R. L. 1992. Class-based n-gram models of natural language. *Computational linguistics*, 18(4): 467–480.
- Cao, S.; Kitaev, N.; and Klein, D. 2020. Unsupervised Parsing via Constituency Tests. *arXiv preprint arXiv:2010.03146*.
- Chang, J.; Wang, L.; Meng, G.; Xiang, S.; and Pan, C. 2017. Deep adaptive image clustering. In *Proc. of ICCV*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proc. of NAACL-HLT*.
- Droz dov, A.; Verga, P.; Chen, Y.-P.; Iyyer, M.; and McCallum, A. 2019a. Unsupervised Labeled Parsing with Deep Inside-Outside Recursive Autoencoders. In *Proc. of EMNLP-IJCNLP*.
- Droz dov, A.; Verga, P.; Yadav, M.; Iyyer, M.; and McCallum, A. 2019b. Unsupervised Latent Tree Induction with Deep Inside-Outside Recursive Auto-Encoders. In *Proc. of NAACL-HLT*.
- Ghasedi Dizaji, K.; Herandi, A.; Deng, C.; Cai, W.; and Huang, H. 2017. Deep clustering via joint convolutional autoencoder embedding and relative entropy minimization. In *Proc. of ICCV*.
- Gillenwater, J.; Ganchev, K.; Graça, J.; Pereira, F.; and Taskar, B. 2010. Sparsity in Dependency Grammar Induction. In *Proc. of ACL*.
- Gupta, V. 2021. Multilingual and Multilabel Emotion Recognition using Virtual Adversarial Training. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*. Association for Computational Linguistics.
- He, J.; Neubig, G.; and Berg-Kirkpatrick, T. 2018. Unsupervised learning of syntactic structure with invertible neural projections. In *Proc. of EMNLP*.
- He, J.; Zhang, Z.; Berg-Kirkpatrick, T.; and Neubig, G. 2019. Cross-lingual syntactic transfer through unsupervised adaptation of invertible projections. In *Proc. of ACL*.
- Hewitt, J.; and Manning, C. D. 2019. A structural probe for finding syntax in word representations. In *Proc. of NAACL-HLT*.
- Huang, P.; Huang, Y.; Wang, W.; and Wang, L. 2014. Deep embedding network for clustering. In *Proc. of International Conference on Pattern Recognition*.
- Jawahar, G.; Sagot, B.; and Seddah, D. 2019. What Does BERT Learn about the Structure of Language? In *Proc. of ACL*.
- Ji, P.; Zhang, T.; Li, H.; Salzmänn, M.; and Reid, I. 2017. Deep subspace clustering networks. In *Proc. of NeurIPS*.
- Jiang, Y.; Han, W.; and Tu, K. 2016. Unsupervised Neural Dependency Parsing. In *Proc. of EMNLP*.
- Jiang, Z.; Zheng, Y.; Tan, H.; Tang, B.; and Zhou, H. 2016. Variational deep embedding: An unsupervised and generative approach to clustering. In *Proc. of IJCAI*.
- Jin, L.; Doshi-Velez, F.; Miller, T.; Schwartz, L.; and Schuler, W. 2019. Unsupervised learning of PCFGs with normalizing flow. In *Proc. of ACL*.
- Jin, L.; and Schuler, W. 2020. The Importance of Category Labels in Grammar Induction with Child-directed Utterances. In *Proc. of International Conference on Parsing Technologies*.
- Johnson, M. 2007. Why doesn't EM find good HMM POS-taggers? In *Proc. of EMNLP-CoNLL*.
- Joshi, M.; Chen, D.; Liu, Y.; Weld, D. S.; Zettlemoyer, L.; and Levy, O. 2020. SpanBERT: Improving Pre-training by Representing and Predicting Spans. *TACL*, 8: 64–77.
- Joshi, M.; Levy, O.; Zettlemoyer, L.; and Weld, D. 2019. BERT for Coreference Resolution: Baselines and Analysis. In *Proc. of EMNLP-IJCNLP*.
- Joulin, A.; Grave, E.; Bojanowski, P.; and Mikolov, T. 2017. Bag of Tricks for Efficient Text Classification. In *Proc. of EACL*.
- Kim, T.; Choi, J.; Edmiston, D.; and Lee, S.-g. 2020. Are pre-trained language models aware of phrases? simple but strong baselines for grammar induction. *arXiv preprint arXiv:2002.00737*.
- Kim, Y.; Dyer, C.; and Rush, A. M. 2019. Compound Probabilistic Context-Free Grammars for Grammar Induction. In *Proc. of ACL*.
- Kitaev, N.; and Klein, D. 2018. Constituency Parsing with a Self-Attentive Encoder. In *Proc. of ACL*.
- Klein, D.; and Manning, C. D. 2002. A Generative Constituent-Context Model for Improved Grammar Induction. In *Proc. of ACL*.
- Klein, D.; and Manning, C. D. 2004. Corpus-Based Induction of Syntactic Structure: Models of Dependency and Constituency. In *Proc. of ACL*.
- Lee, K.; He, L.; and Zettlemoyer, L. 2018. Higher-Order Coreference Resolution with Coarse-to-Fine Inference. In *Proc. of NAACL-HLT*.
- Li, X. L.; and Eisner, J. 2019. Specializing word embeddings (for parsing) by information bottleneck. *arXiv preprint arXiv:1910.00163*.
- Lin, C.-C.; Ammar, W.; Dyer, C.; and Levin, L. 2015. Unsupervised POS Induction with Word Embeddings. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.



- Linzen, T.; Dupoux, E.; and Goldberg, Y. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *TACL*, 4: 521–535.
- Littell, P.; Mortensen, D. R.; Lin, K.; Kairis, K.; Turner, C.; and Levin, L. 2017. Uriel and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors. In *Proc. of EACL*.
- Liu, N. F.; Gardner, M.; Belinkov, Y.; Peters, M. E.; and Smith, N. A. 2019a. Linguistic Knowledge and Transferability of Contextual Representations. In *Proc. of NAACL-HLT*.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019b. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Marcus, M. P.; Santorini, B.; and Marcinkiewicz, M. A. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2): 313–330.
- McDonald, R.; Nivre, J.; Quirmbach-Brundage, Y.; Goldberg, Y.; Das, D.; Ganchev, K.; Hall, K.; Petrov, S.; Zhang, H.; Täckström, O.; Bedini, C.; Bertomeu Castelló, N.; and Lee, J. 2013. Universal Dependency Annotation for Multilingual Parsing. In *Proc. of ACL*.
- Mikolov, T.; Chen, K.; Corrado, G.; and Dean, J. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Minaee, S.; Kalchbrenner, N.; Cambria, E.; Nikzad, N.; Chenaghlu, M.; and Gao, J. 2021. Deep Learning-based Text Classification: A Comprehensive Review. *ACM Computing Surveys (CSUR)*.
- Nivre, J.; de Marneffe, M.-C.; Ginter, F.; Goldberg, Y.; Hajič, J.; Manning, C. D.; McDonald, R.; Petrov, S.; Pyysalo, S.; Silveira, N.; Tsarfaty, R.; and Zeman, D. 2016. Universal Dependencies v1: A Multilingual Treebank Collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*. European Language Resources Association (ELRA).
- Peters, M.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018a. Deep Contextualized Word Representations. In *Proc. of NAACL-HLT*.
- Peters, M.; Neumann, M.; Zettlemoyer, L.; and Yih, W.-t. 2018b. Dissecting Contextual Word Embeddings: Architecture and Representation. In *Proceedings of 2018 Conference on Empirical Methods in Natural Language Processing*.
- Pimentel, T.; Valvoda, J.; Hall Maudslay, R.; Zmigrod, R.; Williams, A.; and Cotterell, R. 2020. Information-Theoretic Probing for Linguistic Structure. In *Proc. of ACL*.
- Pires, T.; Schlinger, E.; and Garrette, D. 2019. How Multilingual is Multilingual BERT? In *Proc. of ACL*.
- Rosenberg, A.; and Hirschberg, J. 2007. V-Measure: A Conditional Entropy-Based External Cluster Evaluation Measure. In *Proc. of EMNLP-CoNLL*.
- Seginer, Y. 2007. Fast Unsupervised Incremental Parsing. In *Proc. of ACL*.
- Smith, N. A.; and Eisner, J. 2006. Annealing Structural Bias in Multilingual Weighted Grammar Induction. In *Proc. of COLING-ACL*.
- Snyder, B.; Naseem, T.; and Barzilay, R. 2009. Unsupervised multilingual grammar induction. Association for Computational Linguistics.
- Spitkovsky, V. I.; Alshawi, H.; and Jurafsky, D. 2013. Breaking Out of Local Optima with Count Transforms and Model Recombination: A Study in Grammar Induction. In *Proc. of EMNLP*.
- Stratos, K. 2019. Mutual Information Maximization for Simple and Accurate Part-Of-Speech Induction. In *Proc. of NAACL-HLT*.
- Stratos, K.; Collins, M.; and Hsu, D. 2016. Unsupervised Part-Of-Speech Tagging with Anchor Hidden Markov Models. *TACL*, 4: 245–257.
- Tenney, I.; Das, D.; and Pavlick, E. 2019. BERT Rediscovered the Classical NLP Pipeline. In *Proc. of ACL*.
- Tenney, I.; Xia, P.; Chen, B.; Wang, A.; Poliak, A.; McCoy, R. T.; Kim, N.; Van Durme, B.; Bowman, S. R.; Das, D.; et al. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. In *Proc. of ICLR*.
- Toshniwal, S.; Shi, H.; Shi, B.; Gao, L.; Livescu, K.; and Gimpel, K. 2020. A Cross-Task Analysis of Text Span Representations. In *Proc. of RepL4NLP*.
- Tran, K. M.; Bisk, Y.; Vaswani, A.; Marcu, D.; and Knight, K. 2016. Unsupervised Neural Hidden Markov Models. In *Proc. of the Workshop on Structured Prediction for NLP*.
- Tseng, H.; Jurafsky, D.; and Manning, C. 2005. Morphological features help POS tagging of unknown words across language varieties. In *Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing*.
- Wu, S.; and Dredze, M. 2019. Beto, Bentz, Becas: The Surprising Cross-Lingual Effectiveness of BERT. In *Proc. of EMNLP-IJCNLP*.
- Wu, Z.; Chen, Y.; Kao, B.; and Liu, Q. 2020. Perturbed masking: Parameter-free probing for analyzing and interpreting BERT. *arXiv preprint arXiv:2004.14786*.
- Xie, J.; Girshick, R.; and Farhadi, A. 2016. Unsupervised deep embedding for clustering analysis. In *Proc. of ICML*.
- Yang, B.; Fu, X.; Sidiropoulos, N. D.; and Hong, M. 2017. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *Proc. of ICML*.
- Yang, J.; Parikh, D.; and Batra, D. 2016. Joint unsupervised learning of deep representations and image clusters. In *Proc. of CVPR*.
- Yatbaz, M. A.; Sert, E.; and Yuret, D. 2012. Learning syntactic categories using paradigmatic representations of word context. In *Proc. of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*.
- Yuret, D.; Yatbaz, M. A.; and Sert, E. 2014. Unsupervised instance-based part of speech induction using probable substitutes. In *Proc. of COLING 2014*.
- Zhou, J.; and Zhao, H. 2019. Head-Driven Phrase Structure Grammar Parsing on Penn Treebank. In *Proc. of ACL*.
- Zhou, Y.; and Srikumar, V. 2021. DirectProbe: Studying Representations without Classifiers. In *Proc. of NAACL-HLT*.