

Probing Syntactic Dependencies with Conditional Mutual Information and Grammatical Constraints

Anonymous ACL submission

Abstract

Unsupervised dependency parsing is a fundamental task in understanding syntactic dependency structures of natural language. Previous parameter-free methods for probing the dependency structure recover a non-trivial amount of dependencies by assuming a correlation between the syntactic dependency (a word-to-word relation) and bi-lexical dependence scores (a metric measuring one word’s influence on the other word). However, these studies assume the correlation without verifying the existence of the correlation. Furthermore, previous studies failed to utilize grammatical constraints that are beneficial to parsing performance in grammar-based unsupervised parsing methods. In this paper, we investigate the correlation between the syntactic dependency and Conditional Mutual Information (CMI) scores, a bi-lexical statistical dependence metric. We propose *delta-energy*, an unbiased estimate of the CMI, and apply it to unsupervised dependency parsing. We further assist the parsing model with three grammatical constraints. We found the delta-energy score capable of effectively separating syntactic dependencies from non-dependencies. Our unsupervised parsing model outperforms baseline parameter-free probing models in parsing performance, excelling in recovering semantically-related dependencies. The ablation study shows that the three grammatical constraints contribute to the recovery of dependencies that are semantically related and that have strong Part-Of-Speech requirements.

1 Introduction

Syntactic dependency structures are important to downstream Natural Language Processing tasks, such as Information Extraction (Tian et al., 2021; Gamallo et al., 2012), Machine Translation (Bugliarello and Okazaki, 2020; Ma et al., 2020), and Question Answering (Lyu et al., 2021). However, training a supervised dependency parser requires expensive human-annotated dependency

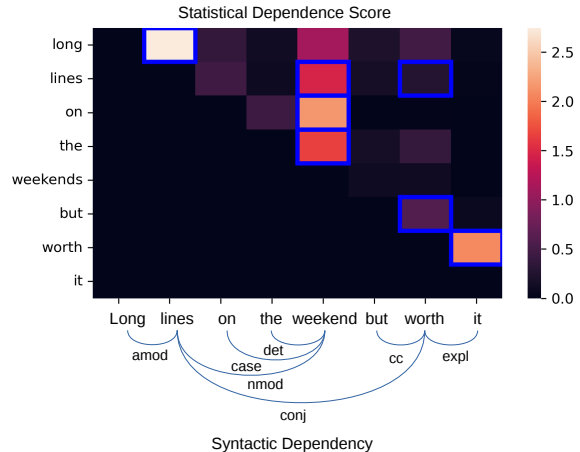


Figure 1: The correlation between syntactic dependencies and conditional mutual information scores. The lower part plots the dependency structure, and the upper part plots the conditional mutual information score. The blue box in the heatmap indicates a syntactic dependency between the corresponding words.

structures, which are only available for some languages and domains. Parameter-free probing methods (Hoover et al., 2021; Wu et al., 2020; Zhang and Hashimoto, 2021) directly extract the dependency structure from pre-trained language models without the structure annotation (a.k.a. unsupervised parsing). These methods predict the dependency structure by finding a set of dependencies that form a tree-shaped structure and that have maximum bi-lexical dependence scores.

Fig.1¹ illustrates the syntactic dependency structure and Conditional Mutual Information (CMI) scores, a bi-lexical statistical dependence metric. The syntactic dependency, represented as a word pair, encodes a word-to-word grammatical relation. For example, the dependency (“long”, “line”) with the label amod indicates that “long” serves as an

¹Conditional mutual information is a symmetric score. We study the undirected syntactic dependency to put the conditional mutual information and other dependence scores on an equal footing.

adjective modifier to “line”. On the other hand, the CMI score is a metric measuring one word’s influence on the other word. The higher the CMI score, the stronger the influence. For example, the CMI score for the word pair (“long”, “line”) is high, indicating a strong bi-lexical dependence. The above example shows a correlation between syntactic dependencies and high CMI scores. Such correlation is the cornerstone of the parameter-free probing method.

Despite the cornerstone role of the correlation, the parameter-free probing method assumes the correlation without verifying it. Whether and how much their dependence scores separate the syntactic dependency from non-dependencies (i.e., word pairs that are not connected by a syntactic dependency) remains a question. Furthermore, the probing method failed to incorporate grammatical constraints that have been shown beneficial to parsing performance by grammar-based unsupervised parsing methods (Noji et al., 2016; Naseem et al., 2010; Xu et al., 2021).

In this paper, we present a study on the correlation between the syntactic dependency and the CMI score. Our contributions are three-fold:

1. We propose *delta-energy*, an unbiased estimate of the CMI score, and derive an unsupervised parsing model from the delta-energy score. We further enhance the parsing model with three grammatical constraints: a Part-Of-Speech (POS) constraint, an adjacent-connect constraint, and a function word head constraint.
2. We verify the correlation between the syntactic dependency and the delta-energy score and show that the delta-energy score is an effective metric for separating syntactic dependencies from non-dependencies.
3. We build a state-of-the-art parameter-free unsupervised parsing model that excels at recovering semantically-related dependencies. Ablation analysis shows that the grammatical constraint has a significant contribution to the recovery of dependencies that are semantically related and that tend to have strong POS requirements.

2 Background

2.1 Conditional Mutual Information

In this paper, we measure bi-lexical statistical dependence scores using Conditional Mutual Information (Eq. 1). Given a sentence $x = (x_1, \dots, x_n)$, CMI measures one word’s (X_i) influence on the other word’s (X_j) distribution under side information c . The side information can include various types of information, contextual or grammatical. For example, the two words’ context $x_{-ij} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_{j-1}, x_{j+1}, \dots, x_n)$ can be the side information. CMI computes the expected log probability ratio between the joint probability $p(x_i, x_j|c)$ and the product of the marginal probabilities $p(x_i|c)p(x_j|c)$. CMI measures the distance between the joint distribution from the marginal product distribution. The higher the CMI score, the further the joint and the marginal product distribution are, the stronger the statistical dependence is between X_i and X_j .

$$I(X_i; X_j|c) := \mathbb{E}_{x_i, x_j|c} \left[\log \frac{p(x_i, x_j|c)}{p(x_i|c)p(x_j|c)} \right] \quad (1)$$

2.2 Extracting Syntactic Dependency via Measuring Bi-Lexical Dependence

Hoover et al. (2021); Wu et al. (2020) measure the bi-lexical dependence score of word pairs (X_i, X_j) under the context x_{-ij} . They define the dependence score as the difference between the informed probability $p(x_j|x_i, x_{-ij})$ and the null probability $p(x_j|X_i = [\text{MASK}], x_{-ij})$ given by a masked language model. The informed probability measures the probability of x_j given the content of x_i , while the null probability measures the probability of x_j without knowing the content of x_i . The null probability can serve as an approximation to the marginal probability $p(x_j|x_{-ij})$ (Xu et al., 2020). After obtaining the dependence score for all word pairs, the probing method selects a dependency tree that maximizes the sum of the dependence scores using Maximum Spanning Tree algorithms (Prim, 1957; Eisner, 1997; Edmonds, 1967).

Hoover et al. (2021) uses a log ratio of the informed and the null probability (pmi, Eq. 2) as the dependence score. The pmi score is a single-point estimate of the CMI score, and the single-point estimation is well-known to have a high estimation variance. Moreover, the pmi score uses the null probability to approximate the marginal probability, which adds further bias to the estimation. The

two issues of the pmi score make it an unreliable estimate of the CMI score. On the other hand, Wu et al. (2020) uses the Euclidean distance between the embedding generating the informed probability $e_{\text{informed_prob}}$ and the embedding generating the null probability $e_{\text{null_prob}}$ (Eq. 3). This method operates in the embedding space instead of the probabilistic space and is not directly comparable with our method.

$$I_{ij}^{\text{pmi}} := \log \frac{p(x_j|x_i, x_{-ij})}{p(x_j|X_i = [\text{MASK}], x_{-ij})} \quad (2)$$

$$I_{ij}^{\text{pert}} := \|e_{\text{informed_prob}} - e_{\text{null_prob}}\|^2 \quad (3)$$

2.3 Sampling from Language Models

Monte Carlo estimation is the standard approach to estimating the CMI score reliably. Goyal et al. (2022) proposes a Metropolis-Hastings (MH) algorithm to sample from language models. Starting with a sentence x with arbitrary content on the word X_i and X_j , the MH method iteratively performs the following steps over the two words:

1. samples a proposal word x'_i from a proposal distribution $q(x'_i|x)$. In this case, the proposal distribution would be a mask language model distribution with X_i set to [MASK].
2. computes the acceptance probability based on the proposal probability and the target probability $p(x')$. Here, p is the language model distribution of choice and $x' = (x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)$
3. accepts or rejects the proposal word according to the acceptance probability. If accepted, $x \leftarrow x'$. Otherwise, $x \leftarrow x$.

The (x_i, x_j) samples produced by the MH algorithm are guaranteed to converge to the target distribution $p(x_i, x_j|x_{-ij})$ as long as the target distribution is irreducible and aperiodic (Besag, 2004). Common language model distributions satisfy the irreducibility and aperiodicity conditions (Goyal et al., 2022).

Nonetheless, the convergence speed of the MH algorithm can be slow in practice. Multi-try MH algorithms (Martino, 2018) mitigate the slow convergence problem by independently proposing n samples and accepting the sample with the highest probability in the target distribution. This approach enables the multi-try MH algorithm to explore the high-probability region of the target distribution

more efficiently than the original MH algorithm, leading to faster convergence.

3 Method

Our method consists of two stages: inducing CMI scores and decoding syntactic dependency structure from the CMI score. We incorporate the POS constraint in the induction stage and incorporate the adjacent-connect and the function word head constraint in the decoding stage.

3.1 Inducing CMI Scores

We define the bi-lexical dependence score as the CMI between two words X_i and X_j under side information c . The side information includes, mandatorily, the context x_{-ij} and, optionally, the POS constraint y_i and y_j for X_i and X_j . We compute the CMI with Eq. 4 and use a causal language model (CLM) distribution for p . Our experiments show that CLMs provide higher-quality samples than masked language models. Eq. 4 is an equivalent form of Eq. 1 that is more suitable for the sampling-based estimation of the CMI score. The first term is the expected probability of samples from the joint distribution $X_i X_j|c$, whereas the second term is the expected probability of samples from the marginal product distribution $X_i|c \otimes X_j|c$.

$$I_{ij}(x) = \mathbb{E}_{(x_i, x_j) \sim X_i X_j|c} [\log p(x_i, x_j, c)] - \mathbb{E}_{(x'_i, x'_j) \sim X_i|c \otimes X_j|c} [\log p(x'_i, x'_j, c)] \quad (4)$$

3.1.1 Sampling with the POS Constraint

Directly applying the MH algorithm to the CLM could not incorporate the POS constraint into the CMI score because the CLM predicts the next word based on its preceding context without considering any POS constraint. We incorporate the POS constraint by applying a mask over the CLM's output distribution. Rijkhoff (2007) points out that POS is a set of words with the same grammatical properties. By the definition, we can impose a POS constraint y_i for the word X_i by masking out words that do not have the grammatical property specified by the POS. We translate the idea into Eq.5. Eq.5 defines the conditional distribution $p(x_i|x_j, x_{-ij}, y_i)$ by renormalizing the CLM probability of a word $p(x_i|x_j, x_{-ij})$ with the total probability of all words satisfying the POS constraint $\sum_{X_i} p(x_i|x_j, x_{-ij}) \mathbb{1}_{(Y(X_i)=y_i)}$

$$\begin{aligned}
& p(x_i|x_j, x_{-ij}, y_i) \\
& := \frac{p(x_i|x_j, x_{-ij})\mathbb{1}_{(Y(X_i)=y_i)}}{\sum_{X_i} p(x_i|x_j, x_{-ij})\mathbb{1}_{(Y(X_i)=y_i)}} \quad (5)
\end{aligned}$$

3.1.2 Estimating with the POS constraint

Estimating the CMI requires computing the joint probability of the sentence (x_i, x_j, x_{-ij}) and the POS constraint (y_i, y_j) . Unfortunately, one can not compute the joint probability straightforwardly as the CLM does not model the POS constraint. We propose *delta-energy* (Eq.6), an unbiased estimate of the CMI score, to overcome this issue. Compared to the CMI score, the delta-energy score eliminates the POS constraint y_i and y_j inside the expectation, enabling straightforward computation of the probability using CLM. We prove that the elimination is safe such that the delta-energy score is equivalent to the CMI score when the side information c contains only the POS constraint and the remaining context (Appendix A.1)

$$\begin{aligned}
I_{ij}^{DE}(x) = & \mathbb{E}_{(x_i, x_j) \sim X_i X_j | c} [\log p(x_i, x_j, x_{-ij})] \\
& - \mathbb{E}_{(x'_i, x'_j) \sim X_i | c \otimes X_j | c} [\log p(x'_i, x'_j, x_{-ij})] \quad (6)
\end{aligned}$$

3.2 Decoding Syntactic Dependency from Delta-Energy Scores

We apply Prim’s algorithm (Prim, 1957) to decode an undirected dependency tree from the symmetric delta-energy score. We additionally apply two grammatical constraints at this stage: the adjacent-connect constraint and the function-word head constraint.

The adjacent-connect constraint is a default strategy when a word is not statistically dependent on the rest of the sentence (i.e., $\forall j, I_{ij}(x) \approx 0$). In that case, we default the word to be connected with its right neighbor, inspired by the high parsing performance of a trivial baseline that connects adjacent words (Klein and Manning, 2004). We set a threshold τ such that a word is automatically connected to its right neighbor if the accumulative delta-energy score between the word and the rest of the sentence is below τ .

The function-word head constraint (Noji et al., 2016) prevents function words from being a syntactic head in the decoded structure. In the context of undirected dependencies, the constraint prevents function words from having more than one connection to other words. We enforce the constraint by

gradually decreasing delta-energy scores related to the function word that violates the constraint. As we will see in Section 4, the two constraints are effective in improving the parsing performance.

4 Experiment

4.1 Experiment Setup

We use three datasets for experiments: EWT-10, WSJ-10 (Klein and Manning, 2004), and PUD. Among the three, the EWT-10 and the WSJ-10 dataset contain sentences shorter than 10 words (excluding punctuations) from the English Web Treebank (Bies, Ann et al., 2012) and the Penn Treebank (Marcus, Mitchell P. et al., 1999) respectively. The main reason for using the EWT-10 and the WSJ-10 datasets is the high computational cost of the delta-energy estimation. For example, running the delta-energy estimation on the development section of EWT-10 takes 48 GPU hours on a single A100 GPU. In addition, the WSJ-10 dataset is widely used for unsupervised dependency parsing (Klein and Manning, 2004; Cohen and Smith, 2009). The PUD dataset contains the full English section of the Parallel Universal Dependency treebank (Zeman et al., 2018). The EWT-10 and the PUD dataset contain dependencies annotated in the universal dependency format (Nivre et al., 2020) while the WSJ-10 contains dependencies annotated in the Stanford dependency format (de Marneffe and Manning, 2008). We use the development section of the EWT-10 dataset (i.e., EWT-DEV-10) to analyze the correlation between the syntactic dependency and the delta-energy score and to evaluate the parsing model derived from the delta-energy score. We use the test section of the EWT-10 dataset, the WSJ-10 dataset, and the PUD dataset to evaluate the models’ parsing performance on universal dependencies, on Stanford dependencies, and on long sentences, respectively. The parsing performance is measured in Unlabelled Undirected Attachment Score (UUAS) (Nivre and Fang, 2017) due to the symmetricity of the delta-energy score. We compute the UUAS score for syntactic dependencies that connect actual words for all experiments (i.e., we exclude the root dependency from the evaluation).

We use the bert-large-cased model (Devlin et al., 2019) for the proposal distribution and the gpt2-large model (Radford et al., 2019) for the target distribution when sampling for the EWT-10 and the WSJ-10 datasets. For the PUD dataset, we

Dependence Scores	P-Value (\downarrow)	Cohen's d (\uparrow)
delta-energy	0.00E+00	1.14
pmi	2.38E-214	0.69
perturbed-masking	0.00E+00	1.11

Table 1: P-value and Cohen's d value for the three scores in separating the syntactic dependency from the non-dependency. The p-value indicates whether the score can separate the two dependencies, and the d value indicates the separation effect. A low p-value and a high d value indicate a good separation score.

alternatively use the opt-125m model (Zhang et al., 2022) for the target distribution to speed up the sampling process. We take one sample for every 12 sampling steps to avoid correlation between subsequent samples. In total, we collect 128 samples for every word pair. We limit words that can be sampled from the bert-large-cased model to the vocabulary of the bert model to reduce the computational and implementation complexity. We use the POS tag provided in the dataset for implementing the POS constraint and the function word head constraint. We run each experiment once because our method is parameter-free and also because of the high computational cost.

We use three baselines for analyses: the pmi baseline (Hoover et al., 2021), the perturbed-masking baseline² (Wu et al., 2020), and the adjacent-connect baseline (Klein and Manning, 2004).

4.2 Correlation between Syntactic Dependencies and Dependence Scores

The core question we seek to answer in this paper is: whether and to what degree does the syntactic dependency correlate with bi-lexical dependence scores? We answer this question by studying whether and how much the dependence score can separate the syntactic dependency from the non-dependency. We compare the delta-energy score with the dependence score derived from the pmi and the score derived from the perturbed-masking baseline. The experiment shows that the delta-energy score can separate and separates the syntactic and the non-dependency well.

The first column in Table 1 shows the p-value for the t-test with a null hypothesis that the syntactic dependency and the non-dependency have the same mean dependence score. All three scores

²we use the released code for the experiment but corrected an implementation bug. See Appendix.A.3

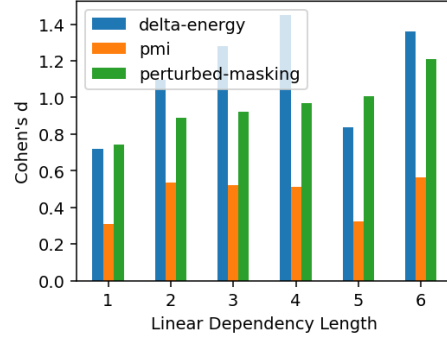


Figure 2: Cohen's d value by linear dependency lengths

have a p-value of 0 or close to 0, suggesting that all dependence scores can separate the syntactic and the non-dependency as two statistical populations.

The second column shows Cohen's d value, which measures the separation effect of the dependence score. The delta-energy score has the highest d value of 1.14³ among the three dependence scores. The perturbed-masking score has a medium d value of 1.11, and the pmi score has the lowest d value of 0.69. This result indicates that the delta-energy score is the best score for separating the dependency and the non-dependency group. The high d value suggests that the delta-energy model could perform better than the two baseline models in parsing performance, as we will see in the next section.

However, syntactic dependencies are not uniformly distributed across all dependency lengths. The syntactic dependency, on average, has shorter lengths than the non-dependency. The discrepancy creates a concern that the above analysis is not only measuring the separation of the syntactic dependency and the non-dependency but also the effect of the short and the long dependency. To counteract the concern, Fig. 2 breaks down the d value by the linear dependency length (i.e., the number of words between the word pair). The delta-energy score has the highest d value for most dependency lengths. The result reinforces the observation derived from Table 1 that the delta-energy score is the best score in separating the syntactic dependency from the non-dependency.

4.3 Parsing Performance of the Delta-Energy Model

Table 2 shows the parsing performance of the delta-energy model and the baseline parsing models on

³A d value of 0.5 indicates a medium effect, 0.8 a large effect, and 1.2 a very large effect (Sawilowsky, 2009)

Dependence Score	UUAS
delta-energy	0.631
pmi	0.559
perturbed-masking	0.586
adjacent-connect	0.497

Table 2: UUAS scores of the delta-energy and the baseline models on the EWT-DEV-10 dataset

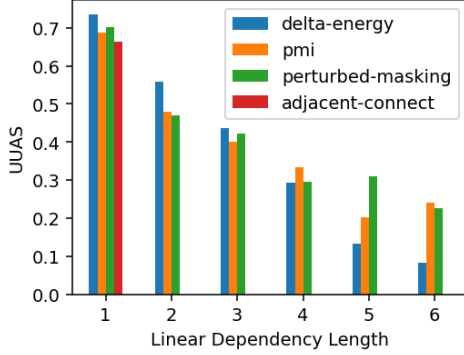


Figure 3: UUAS scores for syntactic dependencies of different lengths on the EWT-DEV-10 dataset

the EWT-DEV-10 dataset. The delta-energy model performs the best, leading the perturbed-masking model (the second-best model) by 0.046. The delta-energy, the perturbed-masking, and the pmi model outperform the adjacent-connect baseline by a large margin. The result confirms that the delta-energy score is the best score for separating the syntactic dependency from the non-dependency. The high performance of the delta-energy, the perturbed-masking, and the pmi model indicates that one can recover a non-trivial amount of syntactic dependencies by measuring the bi-lexical dependence score.

Fig. 3 plots the UUAS scores of the delta-energy model and the baseline models for recovering the syntactic dependency of different lengths. The delta-energy model performs the best for the syntactic dependency with lengths up to 3, performs similarly to the perturbed-masking model for the syntactic dependency with a length of 4, and performs the worst for the syntactic dependency with lengths 5 and 6. The result reveals the source of the delta-energy model’s improvement: the short-length dependency, which makes up the majority of the syntactic dependency. For example, the EWT-DEV-10 dataset has 484 dependencies of lengths greater or equal to 4 while containing 4036 dependencies of lengths less than 4.

Fig. 4 shows the UUAS scores for relations

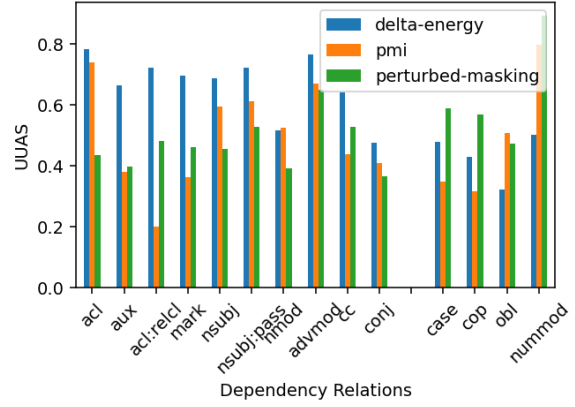


Figure 4: UUAS of relations where the performance difference between the delta-energy and the perturbed-masking model is greater than 0.1. The left part plots the relations where the delta-energy model performs better, and the right part plots the relations where the perturbed-masking model performs better.

where the performance difference between the delta-energy and the perturbed-masking model is more than 0.1. The delta-energy model outperforms the perturbed-masking model in recovering the semantically-related dependencies while underperforming in recovering functionally-related dependencies. The result indicates that the delta-energy model is more sensitive to semantically-related dependencies than functionally-related dependencies.

4.4 Ablation Study

model	UUAS (% loss)	Precision (% loss)
+UPOS, +ADJC, +FNWH	0.6313 (0.00%)	0.6332 (0.00%)
+UPOS, +ADJC, -FNWH	0.6064 (-3.94%)	0.6064 (-4.23%)
+UPOS, -ADJC, +FNWH	0.5518 (-12.59%)	0.5934 (-6.29%)
+UPOS, -ADJC, -FNWH	0.5319 (-15.75%)	0.5712 (-9.79%)
-UPOS, +ADJC, +FNWH	0.5937 (-5.96%)	0.5966 (-5.78%)
-UPOS, +ADJC, -FNWH	0.5544 (-12.18%)	0.5544 (-12.44%)
-UPOS, -ADJC, +FNWH	0.5403 (-14.41%)	0.58 (-8.40%)
-UPOS, -ADJC, -FNWH	0.4992 (-20.93%)	0.535 (-15.51%)

Table 3: Ablation analysis for the delta-energy model. +UPOS indicates the use of the POS constraint, +ADJC indicates the adjacent-connect constraint, and +FNWH indicates the function word head constraint.

Table 3 presents an ablation study for the three grammatical constraints. The UPOS, ADJC, and FNWH represent the POS, the adjacent-connect, and the function word head constraint, respectively. The +/- sign indicates whether the model uses the constraint. The Table shows that removing the POS or the FNWH constraint equally decreases the UUAS and the precision score. On the other

hand, removing the ADJC constraint decreases the UUAS score more than the precision score. This is because, in some cases, the delta-energy score measures a 0 dependence score between one word and the rest of the sentence. The 0 dependence score creates an orphan problem in that the word is statistically disconnected from the rest of the sentence, resulting in an underprediction of the syntactic dependency. The ADJC constraint mitigates the orphan problem by forcibly connecting the word to its right neighbor. The ablation study with the UUAS score indicates that all grammatical constraints benefit the parsing performance and that the adjacent-connect constraint is important in resolving the orphan problem.

Fig. 5 analyzes which relation the grammatical constraint helps the most. The figures plot the UUAS for the relation where removing the respective grammatical constraint causes more than 0.1 loss in UUAS. Fig. 5a shows that the ADJC constraint improves performance for a wide range of dependencies. Fig. 5b shows that the POS constraint improves performance for dependencies with strong POS requirements. For example, the conj relation requires two words to have the same POS tag. The parataxis relation also includes cases where the two words have the same POS tag (Nivre et al., 2020). Fig. 5c shows that the FNWH constraint improves performance for semantically-related dependencies. The nsubj, obl, and advcl dependencies connect words with their semantic arguments. The nmod and acl relation connect words with their modifiers. These dependencies contribute to the semantics of the sentence. The result suggests that grammatical constraints are important for decoding syntactic dependencies from language models.

4.5 Comparison with State-of-the-Arts in Unsupervised Parsing

Type	Models	EWT-TEST-10	WSJ-10	PUD
Parameter-free Probing Models	delta-energy	0.615	0.592	0.525
	perturbed-masking	0.591	0.584	0.507
	mlmbias	0.352	0.586	0.495
Grammar-based Models	dmv	0.611	0.597	0.484
	lcdmv	0.659	0.614	0.554

Table 4: Comparison of the delta-energy model with unsupervised parsing models. The best score for the parameter-free probing models is in bold.

Table 4 compares the delta-energy model with two parameter-free probing models (perturbed-masking and mlmbias (Zhang and Hashimoto,

2021)) and two parametric grammar-based models (dmv (Klein and Manning, 2004) and lcdmv (Noji et al., 2016))⁴. The dmv model has the same grammatical constraint as the delta-energy model, while the lcdmv model has an additional constraint that the sentence can not have a deep recursive center-embedding structure (Noji et al., 2016). Table 4 shows that the delta-energy model performs the best among the parameter-free probing models. Compared to the dmv model, the delta-energy model performs better on the PUD dataset and performs similarly on the EWT-TEST-10 and the WSJ-10 datasets. The better performance on the PUD dataset highlights the strength of the delta-energy model in comparison with the dmv model, a grammar-based model using the same grammatical constraint. Nonetheless, the delta-energy model falls behind the lcdmv model because the lcdmv model has access to additional grammatical constraints. The result again reinforces the importance of the grammatical constraint in recovering syntactic dependencies from language models.

5 Related Works

5.1 Parameter-free Probing Methods

The pmi score (Hoover et al., 2021) measures the bi-lexical dependence score using the log-ratio between the informed and the null probability given by the BERT model. Besides the estimation problem mentioned in Section 2, the pmi score failed to utilize the POS information. In comparison, our method can utilize the POS information as a constraint and improves parsing performance with the information.

The perturbed-masking score (Wu et al., 2020) measures the bi-lexical dependence score using the Euclidean distance of the embedding that generates the informed and the null probability. Despite the simple approach, the perturbed-masking score performs well in recovering the syntactic dependency. However, the perturbed-masking score operates in the embedding space, making it difficult to establish a direct connection between the syntactic dependency and the language modeling objective. In contrast, our delta-energy score operates in the probabilistic space and, consequently, can establish a more direct connection with the language modeling objective. Furthermore, the perturbed masking score cannot utilize the POS information like the

⁴We use the code released by Noji et al. (2016) for the dmv and the lcdmv model

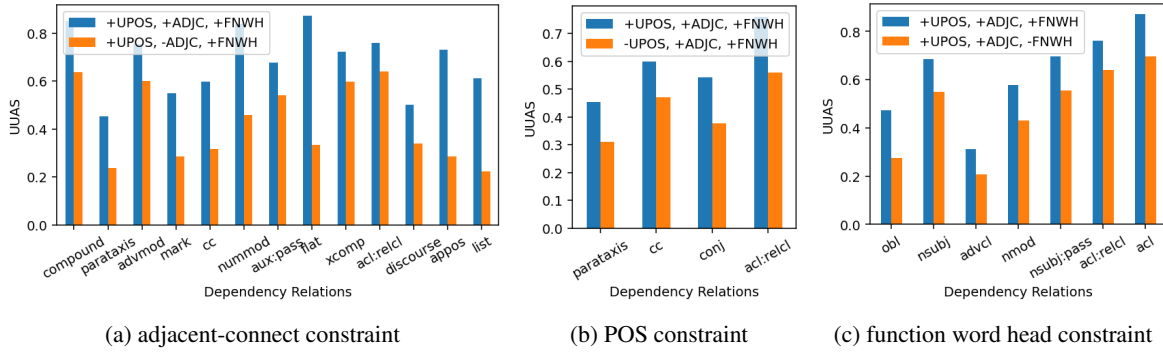


Figure 5: Ablation analysis by dependency relations where removing the respective constraint causes more than 0.1 loss in UUAS score

pmi score, whereas our delta-energy score can utilize the POS information for better performance.

Zhang and Hashimoto (2021) measures the bilexical dependence score using their formulation of the “conditional mutual information”. However, their formulation has two theoretical issues. Firstly, their formulation has an upper bound of 0, in contrast to the widely-known CMI, which is a strictly non-negative metric. Secondly, two statistically independent variables can obtain the maximum value under their formulation. The two issues disqualify their formulation as a valid dependence score. We present the proof in Appendix A.2.

5.2 Parametric Grammar-based Methods

The Grammar-based parametric method (Klein and Manning, 2004; Noji et al., 2016) induces grammar by maximizing the likelihood of observed sentences. While the method can theoretically avoid the data availability problem of lacking dependency annotations, most studies assume the POS information (Han et al., 2020). Effectively, the grammar-based method is still constrained by the availability of the POS information. On the other hand, our method can utilize the POS information as supplementary information. Moreover, the grammar-based method requires a special initialization (Klein and Manning, 2004; Yang et al., 2020) or grammatical constraints (Noji et al., 2016) to induce grammar successfully. As shown in Sec.4, our method can benefit from the constraint but does not require the constraint to extract the dependency properly.

6 Conclusions

In this paper, we studied the correlation between syntactic dependencies and CMI scores derived from causal language models and the application

of the CMI score on unsupervised parsing. We proposed delta-energy, an unbiased estimate of the CMI score that allows the incorporation of POS constraints. We verified that syntactically connected words are more statistically dependent under causal language model distributions. The delta-energy score is the best metric for separating syntactic dependencies from non-dependencies. We found that the unsupervised parsing model induced by the delta-energy score outperforms baseline models by a large margin. The delta-energy model outperforms baseline models in recovering semantically-related dependencies but underperforms in recovering functionally-related dependencies. Our ablation study shows that the POS, the adjacent-connect, and the function word head constraint benefit the parsing performance. The POS constraint contributes to the recovery of dependencies with strong POS requirements. The adjacent-connect constraint boosts the performance in recovering a wide range of dependencies. The function word head constraint significantly contributes to recovering semantically related dependencies. The result indicates the importance of grammatical constraints in extracting syntactic dependencies from language models. The delta-energy model performs strongly against state-of-the-art parameter-free probing models and matches the performance of the grammar-based parametric model using similar grammatical constraints.

7 Limitations

A concern we have is the high computational cost of the MH algorithm. At every sampling step, the MH algorithm has to evaluate the probability of the sentence $x' = (x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)$ with the proposal sample x'_i . Since we collect 128 sam-

622 ples and take one sample for every 12 sampling
 623 steps, we have to evaluate the sentence probability
 624 1536 times to estimate the CMI score for a word
 625 pair. The above is for the constant factor of the
 626 computational complexity. The total computational
 627 complexity for a sentence with n words is $O(n^5)$
 628 considering that the CLM model has a computa-
 629 tional complexity of $O(n^3)$ ($O(n^2)$ complexity for
 630 one pass through the transformer model and $O(n)$
 631 steps to obtain the probability for each word in
 632 the sentence). The high computational cost pre-
 633 vents us from conducting a large-scale multilingual
 634 experiment for languages that have dependency
 635 annotations (Nivre et al., 2020) available.

636 References

- 637 Julian Besag. 2004. [An introduction to markov chain](#)
 638 [monte carlo methods](#). In *Mathematical Foundations*
 639 *of Speech and Language Processing*, pages 247–270.
 640 Springer New York.
- 641 Bies, Ann, Mott, Justin, Warner, Colin, and Kulick, Seth.
 642 2012. [English web treebank](#).
- 643 Emanuele Bugliarello and Naoaki Okazaki. 2020. [En-](#)
 644 [hancing machine translation with dependency-aware](#)
 645 [self-attention](#). In *Proceedings of the 58th Annual*
 646 *Meeting of the Association for Computational Lin-*
 647 *guistics*, pages 1618–1627, Online. Association for
 648 Computational Linguistics.
- 649 Shay Cohen and Noah A. Smith. 2009. [Shared logis-](#)
 650 [tic normal distributions for soft parameter tying in](#)
 651 [unsupervised grammar induction](#). In *Proceedings*
 652 *of Human Language Technologies: The 2009 An-*
 653 *ual Conference of the North American Chapter of*
 654 *the Association for Computational Linguistics*, pages
 655 74–82, Boulder, Colorado. Association for Computa-
 656 tional Linguistics.
- 657 Marie-Catherine de Marneffe and Christopher D. Man-
 658 ning. 2008. [The Stanford typed dependencies repre-](#)
 659 [sentation](#). In *Coling 2008: Proceedings of the work-*
 660 *shop on Cross-Framework and Cross-Domain Parser*
 661 *Evaluation*, pages 1–8, Manchester, UK. Coling 2008
 662 Organizing Committee.
- 663 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
 664 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
 665 [deep bidirectional transformers for language under-](#)
 666 [standing](#). In *Proceedings of the 2019 Conference of*
 667 *the North American Chapter of the Association for*
 668 *Computational Linguistics: Human Language Tech-*
 669 *nologies, Volume 1 (Long and Short Papers)*, pages
 670 4171–4186, Minneapolis, Minnesota. Association for
 671 Computational Linguistics.
- 672 Jack Edmonds. 1967. [Optimum branchings](#). *Journal*
 673 *of Research of the National Bureau of Standards*
 674 *Section B Mathematics and Mathematical Physics*,
 675 71B(4):233.
- Jason Eisner. 1997. [Bilexical grammars and a cubic-](#)
 676 [time probabilistic parser](#). In *Proceedings of the Fifth*
 677 *International Workshop on Parsing Technologies*,
 678 pages 54–65, Boston/Cambridge, Massachusetts,
 679 USA. Association for Computational Linguistics. 680
- Pablo Gamallo, Marcos Garcia, and Santiago Fernández-
 681 Lanza. 2012. [Dependency-based open information](#)
 682 [extraction](#). In *Proceedings of the Joint Workshop*
 683 *on Unsupervised and Semi-Supervised Learning in*
 684 *NLP*, pages 10–18, Avignon, France. Association for
 685 Computational Linguistics. 686
- Kartik Goyal, Chris Dyer, and Taylor Berg-Kirkpatrick.
 687 2022. [Exposing the implicit energy networks behind](#)
 688 [masked language models via metropolis–hastings](#).
 689 In *The Tenth International Conference on Learning*
 690 *Representations, ICLR 2022, Virtual Event, April 25-*
 691 *29, 2022*. OpenReview.net. 692
- Wenjuan Han, Yong Jiang, Hwee Tou Ng, and Kewei Tu.
 693 2020. [A survey of unsupervised dependency pars-](#)
 694 [ing](#). In *Proceedings of the 28th International Con-*
 695 *ference on Computational Linguistics*, pages 2522–
 696 2533, Barcelona, Spain (Online). International Com-
 697 mittee on Computational Linguistics. 698
- Jacob Louis Hoover, Wenyu Du, Alessandro Sordoni,
 699 and Timothy J. O’Donnell. 2021. [Linguistic depen-](#)
 700 [dencies and statistical dependence](#). In *Proceedings*
 701 *of the 2021 Conference on Empirical Methods in Nat-*
 702 *ural Language Processing*, pages 2941–2963, Online
 703 and Punta Cana, Dominican Republic. Association
 704 for Computational Linguistics. 705
- Dan Klein and Christopher Manning. 2004. [Corpus-](#)
 706 [based induction of syntactic structure: Models of](#)
 707 [dependency and constituency](#). In *Proceedings of*
 708 *the 42nd Annual Meeting of the Association for*
 709 *Computational Linguistics (ACL-04)*, pages 478–485,
 710 Barcelona, Spain. 711
- Chenyang Lyu, Lifeng Shang, Yvette Graham, Jennifer
 712 Foster, Xin Jiang, and Qun Liu. 2021. [Improving](#)
 713 [unsupervised question answering via summarization-](#)
 714 [informed question generation](#). In *Proceedings of the*
 715 *2021 Conference on Empirical Methods in Natural*
 716 *Language Processing*, pages 4134–4148, Online and
 717 Punta Cana, Dominican Republic. Association for
 718 Computational Linguistics. 719
- Nianzu Ma, Sahisnu Mazumder, Hao Wang, and Bing
 720 Liu. 2020. [Entity-aware dependency-based deep](#)
 721 [graph attention network for comparative preference](#)
 722 [classification](#). In *Proceedings of the 58th Annual*
 723 *Meeting of the Association for Computational Lin-*
 724 *guistics*, pages 5782–5788, Online. Association for
 725 Computational Linguistics. 726
- Marcus, Mitchell P., Santorini, Beatrice, Mary Ann
 727 Marcinkiewicz, and Taylor, Ann. 1999. [Treebank-3](#). 728
- Luca Martino. 2018. [A review of multiple try MCMC](#)
 729 [algorithms for signal processing](#). *Digit. Signal Pro-*
 730 *cess.*, 75:134–152. 731

A Appendix

A.1 Equivalence of Delta-Energy and CMI

Proposition 1. I_{ij}^{DE} is equivalent with I_{ij} when the side information c contains only the POS information (y_i, y_j) and the remaining context x_{-ij} .

We first look at the definition of CMI. Let $c = (x_{-ij}, y_i, y_j)$

$$I_{ij}(x) \quad (7)$$

$$:= \mathbb{E}_{X_i X_j | c} \left[\log \frac{p(x_i, x_j | c)}{p(x_i | c)p(x_j | c)} \right] \quad (8)$$

$$= \mathbb{E}_{X_i X_j | c} \left[\log \frac{p(x_i, x_j, c)}{p(x_i | c)p(x_j, c)} \right] \quad (9)$$

$$= \mathbb{E}_{\substack{(x_i, x_j) \sim X_i X_j | c \\ (x'_i, x'_j) \sim X_i | c \otimes X_j | c}} \left[\log \frac{p(x_i, x_j, c)}{p(x'_i | c)p(x'_j, c)} \right] \quad (10)$$

$$= \mathbb{E}_{\substack{(x_i, x_j) \sim X_i X_j | c \\ (x'_i, x'_j) \sim X_i | c \otimes X_j | c}} \left[\log \frac{p(x_i, x_j, c)}{p(x'_i, x'_j, c)} \right] \quad (11)$$

$$= \mathbb{E}_{\substack{(x_i, x_j) \sim X_i X_j | c \\ (x'_i, x'_j) \sim X_i | c \otimes X_j | c}} \left[\log \frac{p(x_i, x_j, x_{-ij}) \mathbb{1}_{Y(x_i, x_j) = y_i, y_j}}{p(x'_i, x'_j, x_{-ij}) \mathbb{1}_{Y(x'_i, x'_j) = y_i, y_j}} \right] \quad (12)$$

$$= \mathbb{E}_{\substack{(x_i, x_j) \sim X_i X_j | c \\ (x'_i, x'_j) \sim X_i | c \otimes X_j | c}} \left[\log \frac{p(x_i, x_j, x_{-ij})}{p(x'_i, x'_j, x_{-ij})} \right] \quad (13)$$

$$= I_{ij}^{DE}(x) \quad (14)$$

The key to the proof lies in that the samples always satisfy the condition $Y(x_i, x_j) = y_i, y_j$ and $Y(x'_i, x'_j) = y_i, y_j$. Consequently, the indicator function will always return 1 and enables us to safely remove the POS information inside the expectation.

A.2 Theoretical Issues of Zhang and Hashimoto (2021)

They proposed a formulation of ‘‘conditional mutual information’’ (Eq.15)

$$I_{ij}^{ZH}(x) = \mathbb{E}_{X_i X_j | x_{-ij}} \left[\log p(x_i | x_j, x_{-ij}) - \log \mathbb{E}_{X_j | x_i, x_{-ij}} p(x_i | x_j, x_{-ij}) \right] \quad (15)$$

We prove the following propositions

Proposition 2. The upper bound of I_{ij}^{ZH} is 0.

Proof.

$$(15) = \mathbb{E}_{X_i | x_{-ij}} \left[\mathbb{E}_{X_j | x_i, x_{-ij}} \log p(x_i | x_j, x_{-ij}) - \log \mathbb{E}_{X_j | x_i, x_{-ij}} p(x_i | x_j, x_{-ij}) \right] \quad (16)$$

$$\leq \mathbb{E}_{X_i | x_{-ij}} \left[\mathbb{E}_{X_j | x_i, x_{-ij}} \left[\log p(x_i | x_j, x_{-ij}) - \mathbb{E}_{X_j | x_i, x_{-ij}} \log p(x_i | x_j, x_{-ij}) \right] \right] \quad (17)$$

$$= 0 \quad (18)$$

□

Proposition 3. Two statistically independent can reach the maximum value of 0 under I_{ij}^{ZH} .

Proof. Let the two random variables be defined over a two-value set $X_i, X_j = \{0, 1\}$. Each value has a probability of 0.5. Consequently, we have the joint and the marginal probability as shown in the following table.

X_i		0	1
X_j	Prob	0.5	0.5
0	0.5	0.25	0.25
1	0.5	0.25	0.25

$$I^{ZH}(X_i; X_j) = (2 * 0.5) \left[(0.5 * 2) \log 0.5 \right. \quad (19)$$

$$\left. - \log(0.5 * 2 * 0.5) \right] \quad (20)$$

$$= 0 \quad (21)$$

□

A.3 Implementation bug in Wu et al. (2020)

Datasets	Released	Corrected
EWT-DEV-10	0.626	0.586
EWT-TEST-10	0.581	0.591
WSJ-10	0.572	0.584
PUD	0.501	0.507

Table 5: UUAS of the released implementation vs. our corrected implementation on four datasets

Wu et al. (2020) applies a softmax normalizing the dependence score I_{ij} based on all dependence scores I_i . related to the word X_i . However, we found that their implementation produces a normalized score with a value far greater than 1, which

892 is impossible as the softmax function produces re-
893 sults in the range of $[0, 1]$. We fix the bug using
894 the implementation provided in the `scipy` (Virta-
895 nen et al., 2020) package. Table 5 compares the
896 performance of the released implementation and
897 our corrected implementation. We see that, except
898 for the EWT-DEV-10 dataset, the corrected imple-
899 mentation outperforms the released implementa-
900 tion. This suggests that our implementation does
901 not artificially lower the method’s performance in
902 any way.