An Information-Theoretic Parameter-Free Bayesian Framework for Probing Labeled Dependency Trees from Attention Scores

Anonymous ACL submission

Abstract

Figuring out how neural language models 2 comprehend syntax acts as a key to reveal-3 ing how they understand languages. We 4 systematically analyzed methods of ex-5 tracting syntax from models, namely prob-6 ing, and found five limitations yet widely 7 exist in previous probing practice. We pro-8 posed a method that can directly extract la-9 beled dependency trees from attention 10 scores without training any network, while 11 being able to calculate the mutual infor-12 mation (MI) in a mathematical-rigorous 13 way. Compared with previous approaches, 14 our method has a much simpler model, 15 while being able to probe more complex de-16 pendency trees, providing much more fine-17 grained information about model explana-18 tion at the same time. We demonstrated our 19 method's effectiveness by systematically 20 comparing it with a great many competitive 21 baselines, and gained informative conclu-22 sions, shedding light on our method's ex-23 planation potential. Our code is included in 24 the "software" materials of the openreview 25 system to keep anonymity, and we'll make 26 them publicly available upon publication. 27

Introduction 28

³⁰ (LLMs) have left the world with deep impressions. ³¹ This process is accompanied by confusion, since ⁶⁹ even larger to fit in the dimensionality, inevitably 32 LLMs are usually trained on simple next-token-33 prediction LM tasks, while languages have com-71 34 plex hierarchical structures, known as syntax, and 72 bitter tradeoff: previous methods are putting their 35 recent brain science has proved that humans im-³⁶ plicitly build up syntax structures while reading ⁷⁴ Since vector-based hidden states have completely 37 (Lopopolo et al., 2020; Dotan et al., 2022; Fallon et 75 different modalities compared with dependency ³⁸ al., 2024). Thus, figuring out whether language ⁷⁶ trees, a trainable mapping network is necessary.

³⁹ models have learned to implicitly comprehend syn-40 tax is the key to revealing the essentials of their in-41 telligence.

There are already methods for extracting syntax 42 43 (or other linguistics concepts) from model internal 44 states, called probing methods. A common practice 45 of probing is to train a supervised classifier net-⁴⁶ work on top of model states (Hewitt and Manning, 47 2019; Pimentel et al., 2020; Müller-Eberstein et al., 48 2022) to predict dependency syntax trees, or di-49 rectly take some model states as evidence for syn-50 tax (Htut et al., 2019).

Despite the insights they gave, it is obvious that 51 ⁵² previous probing methods are **explaining by unex-**⁵³ plainibility: Most of them are introducing external 54 trainable networks to extract syntax, ranging from simple linear mappings (Liu et al., 2019) to deep 56 MLPs (Hewitt and Liang, 2019; Voita and Titov, 57 2020; Pimentel et al., 2020a) or pseudo attention 58 heads (Pimentel et al., 2022). This is causing a ⁵⁹ trade-off: Linear mappings are simple and explain-60 able, but have limited expressivity. Deeper networks can fit any co-relationships, but a deep prob-62 ing network is unexplainable itself, so it's natural 63 to raise the doubt on whether the extracted syntax 64 structures really come from the probed LM, or just 65 the strong probes have learned to unconditionally 66 predict them. Moreover, since modern LLMs have 29 Recent advancements in Large Language Models 67 larger hidden dimensions compared with pre-68 trained models, the trainable networks have to be 70 making them more unexplainable.

> If we dive deeper, we might find clues about this 73 attention mainly on contextualized hidden states. 77 Using hidden states is also a primary cause for the 78 aforementioned concern of did the probe learn the

79 task: Contextualized embeddings embed abundant 131 trap of directly using attention scores as depend-⁸⁰ semantics, so even though the probed LM knows ¹³² ency probabilities. We systematically compared ⁸¹ nothing about syntactics, it's still possible that the ¹³³ with a series of strong baselines, even with those ⁸² deep probing model learns it. As an example, if we ¹³⁴ methods requiring far more complex probe net-⁸³ see two words, *eat* and *breakfast*, even without any ¹³⁵ works, and achieved state-of-the-art head im-⁸⁴ context, there're still good reasons for us to believe ¹³⁶ portance estimation and tree-constructing perfor-⁸⁵ that *breakfast* acts as the object of the verb *eat*. This ¹³⁷ mance. We further derived informative conclusions ⁸⁶ is exactly the case of this concern.

87 ⁸⁸ the better choice? Maybe we should *put attention* ¹⁴⁰ elegant way, while offering vast possibilities for the ⁸⁹ on attention. Attention is the only component that ¹⁴¹ upcoming conclusion-intensive research thanks to ⁹⁰ involves inter-token relationships (MLP and ¹⁴² its fine-grained MI and probabilities functions. ⁹¹ add/norm are applied token-wise), while depend-⁹² ency syntactics is exactly inter-word relationships. ¹⁴³ 2 93 Moreover, attention maps are topologically con-⁹⁴ sistent with syntax trees: Attention scores are ma-¹⁴⁴ Just after the born of deep contextualize embed-95 trices, while dependency trees can also be de- 145 dings (Peters et al., 2018) and transformer-based ⁹⁶ scribed as adjacency matrices. Even the matrix size ¹⁴⁶ pre-trained models (Devlin et al., 2019), research-97 can be the same (ignoring sub-word tokenization, which is also overridable by simple indexing).

Unfortunately, despite those nice consistencies, 99 to our best knowledge, there're only few works fo-100 cusing on probing attention (Clark et al., 2019; Htut 101 et al., 2019; Vig and Belinkov, 2019; Ravishankar et al., 2021), only being able to extract inferior or even incomplete dependency trees. There seems to be a contradiction. Why is it? Our opinion is that, due to those consistencies, researchers are over-107 trusting attention scores, which is right another ¹⁰⁸ limitation. Since attention scores are softmax-normalized, constituting a probability distribution across tokens, they tend to directly use attention scores as probabilities of a dependency relationship 111 112 between two words. However, attention scores def-113 initely do not only have this single functionality of syntax, so filtering out highly syntactical attention 115 heads together with transformations on attention scores is necessary. 116

Based on our analysis, we proposed our method 117 ¹¹⁸ of Information-theoretic Parameter-free Bayesian 119 Probing (IPBP): Instead of training supervised networks, we chose to directly estimate the multivari-120 ate probabilistic distributions between attention scores and dependency relationships. With those distributions, we're able to estimate mutual infor-124 mation (MI) in a mathematically rigorous way, obtaining a good metric for each head's individual importance for each dependency label. We further de-127 signed a novel decoding algorithm incorporating 128 the estimated MI and Bayesian posteriors, being 178 Belinkov, 2019; Ravishankar et al., 2021), yielding 129 able to efficiently reconstruct *labeled* dependency

138 on the estimated MI and distributions. In a word, If hidden states are not yet good enough, what's 139 our method is addressing the two limitations in an

Related Work

147 ers have started to investigate whether or not lin-¹⁴⁸ guistics properties are embedded in these models 149 (Conneau et al., 2018; Liu et al., 2019; Tenny et al., 150 2019; Hewitt and Manning, 2019). Then arguments 151 began in this area. The frontline of these arguments 152 is about what probing model can we use to prevent 153 it from learning the task itself. While early practices 154 and preliminary works suggested strictly-linear 155 probes (Alain and Bengio, 2017; Hewitt and Man-156 ning, 2019; Liu et al., 2019), Hewitt and Liang, ¹⁵⁷ 2019 proposed control tasks that penalizes models 158 being ability to learn the task itself, and had at-159 tempts on several Deep MLPs. Furthermore, Pi-160 mentel et al., 2020b admitted this trade-off and 161 took probing as an accuracy-complexity two-goal 162 optimizing problem, and most radically, Pimentel 163 et al. 2020a, insisted that probes should be as deep 164 and complex as possible since they used them as 165 estimations of \mathcal{V} -Information (Xu et al., 2020). 166 Apart from disputes, there are also alternative the-167 ories proposed by the researchers, like the code-de-168 scription-length theory by Voita and Titov, 2020 169 and the architectural bottleneck principle by Pi-170 mentel et al., 2022. These theories can be seen as 171 patches under the supervised probing context since 172 they're also addressing the complexity vs. accuracy 173 tradeoff.

Apart from supervised probes, there're also na-174 175 ïve parameter-free probes, mainly based on extract-176 ing dependency trees (or partial dependency arcs) 177 from attention scores (Clark et al. 2019, Vig and 179 not yet good enough probing performances. If we 130 trees, while preventing us from dropping into the 180 take a broader view, we'll also find parameter-free 181 explaination methods for more general-purpose 182 concept in deep learning research (Mu and Andreas, 227 183 2020; Antverg and Belinkov, 2021). Together with 228 lationships {nsubj, dobj, ...} and f(l, a), P(l), 184 some supervised works (Radford et al., 2019; 229 f(a) is short for the density value of joint distribu-Lakretz et al., 2019; Dalvi et al., 2019), these meth- $_{230}$ tion $f(L, A_{b,h})$ at $L = l, A_{b,h} = a$, density of mar-186 ods, also called neuron analysis methods, were sys-187 tematically evaluated by a recent work (Fan et al., 188 2024). In section 4, we'll systematically compare 189 our methods with the principles of these strong 190 baselines.

3 **IPBP Methodology** 191

¹⁹³ method into key points in the first section, and then ²³⁹ 194 introduce the details.

195 3.1 **Key Aspects Analysis**

196 Given a sentence $X = x_1 x_2 \dots x_n$ and an arbitrary ¹⁹⁷ token pair (x_i, x_i) , there might or might not be a ¹⁹⁸ dependency arc from x_i to x_j . We define $l^{[i][j]}$ as ²⁴⁴ Now assume we have a dataset \mathcal{D} consisting of a ¹⁹⁹ the variable for which kind of dependency exists ²⁴⁵ series of *<sentence, dependency tree>* pairs. We from x_i to x_j . $l^{[i][j]}$ can be a specific dependency ²⁴⁶ also have a model with \mathcal{B} blocks and \mathcal{h} attention ²⁰¹ type like nsubj, or ϕ when there's no dependency 202 arc from x_i to x_i . If the sentence is sent into a trans-203 former-based LM, there will also be a series of attention scores from x_i to x_j , namely $a_{b,h}^{[i][j]}$, which $a_{251}^{[250]}$ tween token pairs having dependency l. 205 stands for the attention score of the h-th attention 252 ²⁰⁶ head from the *b*-th transformer block. If we take ²⁵³ sentence $X \in D$, we feed $X = x_1 \dots x_n$ into the 207 the dataset as a series of token pairs, and get the 254 model, and for any token pair $\langle x_i, x_j \rangle$ $(i, j \in$ 208 observation $l^{[i][j]}$ and $a_{b,h}^{[i][j]}$ with respect to each 255 $\{1 \dots n\}$), we'll have its dependency relationship ²⁰⁹ token pair, we'll get two co-occurring dataset-wide ²⁵⁶ $l^{[i][j]}$ and a series of attention scores $a_{b,h}^{[i][j]}$ for any variables, L and $A_{b,h}$ which stand for the depend- $_{257} b \in \{1 \dots b\}$ and $h \in \{1 \dots h\}$. We'll add each at-211 spect to any token pair. 212

Therefore, the goal of our probing can be di-213 214 vided into two:

- MI Estimation: estimating mutual infor- 262 215 216 $MI(L; A_{b,h}).$ 217
- Tree Reconstruction: A method of deriving 218 a full dependency tree based on attention 219 scores $A_{h,h}$ 220

221 ²²² $A_{b,h}$ is continuous, the joint distribution is a mix-223 ture distribution, and the formula of MI is slightly 224 different from its classical discrete or continuous 225 ones, shown as follows:

226 MI
$$(L; A_{b,h}) = \sum_{l \in \mathcal{L} \cup \{\phi\}} \int f(l, a) \log \frac{f(l, a)}{P(l)f(a)} da$$

Where \mathcal{L} stands for the set of all dependency re-²³¹ ginal distribution $f(A_{b,h})$ at $A_{b,h} = a$, and scalar 232 probability P(L = l).

233 Moreover, the second goal can be regarded as a ²³⁴ Bayesian inference process taking $A_{b,h}$ as evi- $_{235}$ dence and L as hypothesis. The posterior distribu-236 tions ($f(L = l | A_{h,h} = a)$) are required for tree re-237 construction. Therefore, the key to achieving these ¹⁹² To foster understanding, we'll first break our ²³⁸ two goals is those probabilistic distributions.

> The above formulation intuitively explains our 240 method in a nutshell, in the following sections ²⁴¹ we'll dive deep into how we can infer these distri-242 butions from the dataset.

243 3.2 **Getting the Distributions**

247 heads within each block. We first initialize a series ²⁴⁸ of attention score sets $\mathcal{A}_{b,h;l}$ where $b \in \{1 \dots \mathcal{B}\}$, 249 $h \in \{1 \dots h\}$ and $L \in \mathcal{L} \cup \{\phi\}$. $\mathcal{A}_{h,h;l}$ means all 250 possible attention scores of attention head b, h be-

Then we'll iterate over the dataset. For a specific ency type and head (b, h)'s attention score with re-₂₅₈ tention score $a_{b,h}^{[i][j]}$ to the corresponding attention ²⁵⁹ score set $\mathcal{A}_{b,h;l^{[i][j]}}$. After iteration, all attention 260 score sets will have all possible attention scores for any token pair $\langle x_i, x_i \rangle$ in any sentence $X \in \mathcal{D}$.

After gaining the attention values, we'll estimate mation (MI) between L and every $A_{b.h}$: 263 those key probabilities. The most intuitive one ²⁶⁴ might be P(L = l), since we can simply take the 265 empirical probability on the dataset as the approxi-²⁶⁶ mate value, which is $\hat{P}(L = l) = \frac{|\mathcal{A}_{b,h;l}|}{\sum_{l' \in \mathcal{L} \cup \{\phi\}} |\mathcal{A}_{b,h;l'}|}$

 $_{267}$ where b, h can be any value and this equation Specifically, since L is a discrete variable while ²⁶⁸ means the proportion of token pairs with depend- $_{269}$ ency *l* among all possible token pairs from the da-270 taset. The tricky ones are the continuous probabili-271 ties. Since we already have abundant attention 272 score samples, we'll use Kernel Density Estimation 273 (KDE) to estimate the continuous ones. Specifi-(1) 274 cally, for every possible $\mathcal{A}_{b,h;l}$, we can regard it as

275 the observation of the attention variable $A_{b,h}$ under 319 coarse-grained, and not helpful for discovering the 276 the circumstance of L = l. The samples in $\mathcal{A}_{b,h;l}$ 320 different functionality of every attention head. 277 follow the conditional density of $f(A_{b,h}|L=l)$. ³²¹ Therefore, we should tweak Equation 1 to make the 278 We use the Gaussian kernel and take a specific 322 MI formulation fit this specialist assumption. The $_{279}$ bandwidth *B* (See Appendix A). Therefore, the $_{323}$ new formulation is as follows: ²⁸⁰ kernel density $\hat{f}(A_{b,h}|L = l)$ can be estimated as: 324

$$\frac{1}{|\mathcal{A}_{b,h;l}| \cdot B} \sum_{i=1}^{|\mathcal{A}_{b,h;l}|} \frac{1}{\sqrt{2\pi \cdot \sigma_{\mathcal{A}_{b,h;l}}}} \exp\left(-\frac{x_0 - \mathcal{A}_{b,h;l}^{(i)}}{B}\right)^{32}$$

283 and $\mathcal{A}_{h,h:l}^{(i)}$ means the *i*-th value of $\mathcal{A}_{b,h:l}$.

with $A_{b,h}$ as evidence and L as hypothesis, then the $\frac{1}{332}$ the possibility of dependency relationships other 287 estimated $\hat{f}(A_{b,h}|L)$ s can be seen as the *likelihood* 333 than *l*, and can be estimated using $1 - \hat{P}(l)$. 288 densities. Applying the Bayesian theorem, we'll 289 get the following equation:

290
$$f(L|A_{b,h}) = \frac{f(A_{b,h}|L)P(L)}{f(A_{b,h})} = \frac{f(A_{b,h},L)}{f(A_{b,h})}$$
291 (3)

292 $_{293} \hat{f}(A_{b,h}|L)$ act as the catalyst of the whole process $_{340}$ pendency relationship $l \in \mathcal{L} \cup \{\phi\}$, we filter out a ²⁹⁴ of IPBP. We can multiply $\hat{f}(A_{b,h}|L)$ with the ³⁴¹ series of attention heads highly responsible for l, $\hat{P}(L)$ s, which are already gained by empirical prob- 342 constituting the head set \mathcal{H}_l . We then infer the pos-296 abilities, to get the joint densities $\hat{f}(A_{b,h}, L)$. By 343 sibilities of dependency arcs of l based on the pos-297 summing over all possible Ls, we'll get the esti- 344 teriors of heads from \mathcal{H}_l , and use MI_{binary} to bal-²⁹⁸ mated marginal density $\hat{f}(A_{b,h})$. Now that each ³⁴⁵ ance between posteriors conditioned on each head 299 term in the above equation is available, we can get 346 from \mathcal{H}_l , forming the overall possibility for a de-³⁰⁰ the estimated posterior probability $\hat{f}(L|A_{b,h})$. By ³⁴⁷ pendency arc with relation *l*. Finally, we use a de-³⁰¹ now, the probabilities required for the whole IPBP ³⁴⁸ coding algorithm to build the dependency tree ³⁰² process, as mentioned in Section 3.1, are all set.

303 3.3 **Estimating MI**

305 to our two main goals: MI estimation and Tree Re- 353 also differs, setting a fixed threshold for all possible 306 construction. However, maybe we should reex- 354 ls will favor those relationships with larger MI val-307 amine the MI formulation in Equation 1: the 355 ues. Therefore, an adaptive threshold conditioning ³⁰⁸ MI($L; A_{b,h}$) in Equation 1 measures how much ³⁵⁶ on the relation l is necessary. Remind that mutual ³⁰⁹ common information head (b, h) has about *every* ³⁵⁷ information is upper-bounded by the individual en-310 possible dependency relationship. However, there 358 tropies of each random variable, in our case, ³¹¹ might not be such an all-around attention head that ³⁵⁹ $H(\mathbf{1}_{\{l\}}(L))$ and $H(A_{b,h})$ (where $\mathbf{1}_{\{l\}}(L)$ is the in-³¹² is responsible for *every* possible dependency type, ³⁶⁰ dicator function meaning whether or not L equals ³¹³ but more probable that some heads are responsible ³⁶¹ to l). Since $A_{b,h}$ is continuous, and the entropy an-314 for certain dependency labels. This kind of special- 362 alogs of continuous variables (variational entropies) 315 ist head is also the assumption of preliminary at- 363 are known as inferior analogs, possibly non-posi-³¹⁶ tention-analyze work like (Htut et al., 2019). Even ³⁶⁴ tive, making it unable to act as an upper bound, we ³¹⁷ though there do exist such versatile heads, an MI ³⁶⁵ choose to estimate $H(\mathbf{1}_{\{l\}}(L))$ as: 318 corresponding to all dependency types is still too

$$MI_{\text{binary}}(l; A_{b,h}) = \int f(l, a) \log \frac{f(l, a)}{P(l)f(a)} da + \int f(\neg l, a) \log \frac{f(\neg l, a)}{P(\neg l)f(a)} da$$
(4)

In that equation, $f(\neg l, a)$ is short for the density value of $f(\neg l, A_{b,h})$ at $A_{b,h} = a$, where ⁽²⁾ ₃₂₈ $f(\neg l, A_{b,h})$ stands for the joint density between Where $\sigma_{\mathcal{A}_{b,h;l}}$ is the standard deviation of $\mathcal{A}_{b,h;l}$, $\frac{320}{329}$ any *L* other than *l* and $A_{b,h}$. In practice, this joint ³³⁰ density can be gained by marginalizing $\hat{f}(A_{h,h}, L)$ Again, if we take the view of Bayesian inference, ₃₃₁ over all possible Ls except for l. $P(\neg l)$ stands for

Getting Highly Syntactical Heads 334 **3.4**

335 By now, having posterior distributions $\hat{f}(L|A_{h,h})$ 336 and MI_{binary} feasible for estimating the independ-) 337 ent importance of each dependency type, the road 338 towards the goal of Tree Reconstruction is clear. This gives us inspirations: the likelihoods 339 The basic idea of our approach is: for every de-349 based on these overall possibilities.

To filter out \mathcal{H}_l , it's natural to set a threshold on 350 ³⁵¹ MI_{binary}($l; A_{b,h}$). However, since for different de-³⁰⁴ With all these distributions, we're able to proceed ³⁵² pendency relationships, the magnitudes of MI_{binary}

 $\hat{H}(\mathbf{1}_{\{l\}}(L)) = \hat{P}(L)\log\hat{P}(L) + \hat{P}(\neg L)\log\hat{P}(\neg L)$ (5) ⁴¹² This is approximately equivalent to the Loga-366 367

³⁶⁹ form [0, 1] scale. An overall threshold can be set ³⁷⁰ for this proportion, resulting in \mathcal{H}_l for every *l*.

Tree-Reconstruction Algorithm 371 3.5

³⁷² After getting \mathcal{H}_l s, another problem occurs: As 373 mentioned before, previous probing practices ³⁷⁴ mainly aim at building unlabeled trees. Even those 375 supervised dependency parsing methods (Dozat and Manning, 2017; Tian et al., 2022) are training 377 separate networks for predicting arcs, and then predicting labels for those predicted arcs. Therefore, 378 these methods are operating on a simple probability space with only probabilities on the existence of de-381 pendency arcs. What's more, in their methods, 429 ity of non-existence. The overall probability there's only one network responsible for predicting 430 $P(x_i, x_j; l)$, meaning the probability of an arc of l 383 necessary that we design a decoding algorithm that 433 lated as follows: 385 not only balances each posterior but also consti-386 tutes a uniform probability space.

We first make an assumption that the overall 435 388 389 possibility of dependency arcs is independently 436 390 conditioned on each head in \mathcal{H}_l (otherwise the problem might be too complex). An ideal resort for ³⁹² balancing each posterior is to treat the prediction of ⁴³⁸ valid probability space. By now, the two problems ³⁹³ dependency arcs as a *voting problem*: for depend-⁴³⁹ introduced by *multi-head* and *multi-label* are all ³⁹³ dependency arcs as a voting problem. For depend-³⁹⁴ ency l, each head $\langle b_i, h_i \rangle \in \mathcal{H}_l$ can be seen as a ⁴⁴⁰ solved. We're just one step towards building the ⁴⁴¹ tree, that is, the decoding algorithm utilizing the ⁴⁴¹ uce, that is, the account of the second probabilities. Specifically, following previ-396 399 thirds) voting the arc belongs to l. However, due to 446 might refer to Appendix A for implementation de-400 the non-discrete weights, the problem cannot be ef- 447 tails of our methods, like hyperparameters and our ficiently dynamically programmed, resulting in a 448 GPU-optimized KDE and integral methods. search space of $\mathcal{O}(2^{|\mathcal{H}_l|})$, which will be rather in-403 efficient during inference. Instead, we relax this 449 3.6 404 voting problem to an easy-computing while ra- 450 Since MI estimation is a small hot topic in statistics, 405 tional form: We take the geometric mean of the 451 in case of re-inventing wheels, we've done research 406 posteriors. Specifically, let $GP_{\mathcal{H}_{i}}(x_{i}, x_{j}; l)$ be the 452 on related methods. We found two methods sharing 407 geometrically-averaged probability of an arc of l 453 (minor) principles with our method: The first one 408 between tokens x_i and x_j conditioned on heads in 454 (Moon et al., 1995) is a method estimating MI be-409 \mathcal{H}_l . In logarithm space, the geometric mean is:

$$410 \log \operatorname{GP}_{\mathcal{H}_{l}}(x_{i}, x_{j}; l) = \frac{\sum_{\langle b_{k}, h_{k} \rangle} \operatorname{MI}_{\operatorname{binary}}(l; A_{b_{k}, h_{k}}) \cdot f(L=l; A_{b_{k}, h_{k}})}{\sum_{\langle b_{m}, h_{m} \rangle \in \mathcal{H}_{l}} \operatorname{MI}_{\operatorname{binary}}(l; A_{b_{m}, h_{m}})}$$

$$411 \tag{6}$$

413 rithmic Opinion Pooling technique widely adopted If we divide the MI with the entropy, the result- 414 in Bayesian inference, thus acting as a reasonable ³⁶⁸ ing proportions $\frac{MI_{binary}(l;A_{b,h})}{\hat{H}(\mathbf{1}_{\{l\}}(L))}$, $\forall l$ will be in a uni-⁴¹⁵ approximation when the number of experts (in our 416 case, heads in \mathcal{H}_{l}) is relatively large. However, the ⁴¹⁷ problem of Logarithmic Pooling is that, if we sum 418 over all probabilities of each possible dependency 419 relationship (in our MI_{binary} case, l and $\neg l$), it is 420 not guaranteed to be 1, recalling the second prob-421 lem of *uniform probability space*. To resolve this, 422 we build a larger multivariate probability space of $_{423} \{0,1\}^{|\mathcal{L}|+1}$. We take the voting process of the de-⁴²⁴ pendency between x_i and x_i as $|\mathcal{L}| + 1$ independ- $_{425}$ ent votes, the ℓ -th voting votes for the existence of 426 the ℓ -th dependency from $|\mathcal{L}|$, using the ⁴²⁷ GP_{\mathcal{H}_i} ($x_i, x_j; l$) in Equation 6 as the probability of 428 existence, and $1 - GP_{\mathcal{H}_{i}}(x_{i}, x_{i}; l)$ as the probabilprobabilities, our method, on the other hand, has a $_{431}$ between tokens x_i and x_i conditioned on all highly bunch of posterior probabilities. Therefore, it's $_{432}$ responsible heads $\mathcal{H}_1 \cup ... \cup \mathcal{H}_{|\mathcal{L}|} \cup \mathcal{H}_{\phi}$, is calcu-

 $P(x_i, x_i; l) =$

$$\operatorname{GP}_{\mathcal{H}_l}(x_i, x_j; l) \cdot \prod_{l' \in |\mathcal{L}| + \{\phi\} - \{l\}} [1 - \operatorname{GP}_{\mathcal{H}_l}(x_i, x_j; l)]$$
(7)

While the probability of not arc between x_i and $_{437} x_i$ is $1 - \sum_{l \in |\mathcal{L}|} P(x_i, x_i; l)$, thus resulting in a probability of a dependency arc of l can be seen as 443 ous supervised dependency parsing works, we're the probability of a series of heads with total 444 using the Eisner dynamic programming algorithm weights larger than a proportion (like half or two- 445 (Eisner, 1996) as the decoding algorithm. Readers

The Novelty of IPBP

455 tween two observations within a time series using 456 KDE. They're doing three individual KDEs, with 457 one multivariate one. While it's a known issue that 458 KDE quickly becomes inferior when variables be-459 come more than one, known as the *dimensionality* 460 *curse*, their method is inevitably introducing errors

461 (and also unapplicable to our attention-dependency 506 462 mixed-joint distribution setting). We, instead, dex- 507 lient score. Following the original authors, we set 463 terously circumvented the curse and made the least 508 it to the top 99.5% value among values in $\mathcal{A}_{h,h,l}$. ⁴⁶⁴ number of estimations possible (limited to 1) by ex- ₅₀₉ 465 ploiting mixed-joint distribution and Bayesian the- 510 refers to a series of methods performing correlation 466 orems.

467 tribution (Gao et al., 2017). However, they use a 513 formulation of these methods: 468 kNN-like algorithm to estimate point-wise mutual 469 information (PMI) and average it over the dataset. 470 515 Their method didn't provide any valid probability 471 472 distributions, thus offering no possibility of tree re-516 473 construction, and also providing less chance for da-474 taset-level or visualization-based explanations.

Experiments 475 4

Baselines 476 4.1

⁴⁷⁸ pare our method with a series of probing as well as ⁵²² the network. When $\lambda_1 = 1$, $\lambda_2 = 0$, this equation ⁴⁷⁹ neuron analysis baselines. Corresponding to the ⁵²³ becomes Lasso (Radford et al., 2019), when $\lambda_1 =$ 480 two sub-tasks introduced in Section 3.1, we first in- $_{524}$ 0, $\lambda_2 = 1$, it becomes Ridge (Lakretz et al., 2019), 481 troduce a series of head-selection baselines, where $_{525}$ and λ_1 , $\lambda_2 = 1$ corresponds to ElasticNet (Dalvi et 482 we replace the estimated MI with other criteria, and 526 al., 2019). We use ElasticNet as a representative. 483 keep the tree-construction algorithm unchanged. 527 After gaining the trained W_{θ} , we use the weight en-484 485 methods. This is better for illustrating the individ- $_{530}$ LFF(*l*; *b*, *h*). 486 ual contributions of each corresponding submodule. 531 487 488 several strong neuron analysis methods evaluated 533 tional probabilities, and use the mean logarithm 489 by a recent paper (Fan et al., 2024): 490

491 a parameter-free method, which gets the correla- 536 the state-of-the-art entropy estimation algorithm, 492 493 tion scores by calculating mean values with respect 537 used by previous methods also taking information-494 to different concepts alongside the dataset. In our 538 theoretic perspectives (Pimentel et al., 2020a, Pi-

⁴⁹⁶ PL(*l*; *A*_{*b*,*h*}) =
$$\sum_{l' \in \mathcal{L} + \{\phi\} - \{l\}} \left| \bar{\mathcal{A}}_{b,h;l} - \bar{\mathcal{A}}_{b,h;l'} \right|$$
 (8)

Where \mathcal{A}_{\dots} denotes the mean value of a spe-497 ⁴⁹⁸ cific attention score set. Note that despite its sim-⁵⁴ plicity, this method is evaluated as the method that 54 499 is most consistent with others by Fan et al., 2024, 500 thus most robust. 501

IoU (Mu and Andreas, 2020): This method uses 502 Jaccard Similarity as a correlation criterion. In our 503 ⁵⁰⁴ implementation, we use the following form:

505 IoU
$$(l; A_{b,h}) = \frac{|\mathcal{A}_{b,h,l} \cap [\tau, +\infty)|}{|\mathcal{A}_{b,h,l}| + \sum_{l' \in \mathcal{L} + \{\phi\} - \{l\}} |\mathcal{A}_{b,h,l'} \cap [\tau, +\infty)|}$$
 (9)

Where τ is a threshold serving as selecting a sa-

The Linear Feedforward Family: This method ⁵¹¹ ranking by training a supervised linear network W_{θ} . Another one is also focusing on mixed joint dis- 512 Specifically, the equation below gives a uniform

$$W_{\theta} =$$

514

$$\underset{W_{\theta} \in \Theta}{\operatorname{argmin}} \left[\sum_{X \in \mathcal{D}} \sum_{x_{i}, x_{j} \in X \times X} \log P_{\theta} \left(l = l^{[i][j]} \middle| a_{1, 1 \dots \hat{\mathcal{H}}, \hat{h}}^{[i][j]} \right) + \lambda_{1} \|\theta\|_{1} + \lambda_{2} \|\theta\|_{2} \right]$$
(10)

Where W_{θ} is a matrix of shape $\mathcal{B}\mathcal{h} \times (|\mathcal{L}| + 1)$, 517 518 and $a_{1,1\dots,\ell,\hbar}^{[i][j]}$ denotes the concatenation of atten-519 tion scores between x_i and x_j for all attention 520 heads, and $P_{\theta}\left(l^{[i][j]} | a_{1,1\dots \mathcal{C}, h}^{[i][j]}\right)$ stands for the 477 In this section, we're going to systematically com- 521 probability of the ground-truth label estimated by We'll also compare the tree-construction algorithm $_{528}$ try mapping attention score of head (b, h), to the with common practices of previous attention-based $_{529}$ probability of relation l as the correlation value,

V-Information: Xu et al., 2020 proposed to use For head-selection baselines, we'll start from 532 a trainable network as an approximation of condi-534 probabilities as approximations of conditional en-Probeless (Antverg and Belinkov, 2021): This is 535 tropies based on the law of large number. This is 495 situation, we use the following instead of MI_{binary}: 539 mentel et al., 2022). Specifically, in our case, we $H_{\mathcal{V}}$ 540 use max $H_{\mathcal{V}}(l|A_{,,\cdot}) - H_{\mathcal{V}}(l|A_{b,h})$ as replacement ⁵⁴¹ of $MI_{\text{binary}}(l; A_{b,h})$, and as the equation shows:

Where $MLP_{b,h;l}(\cdot)$ are deep MLPs individually trained using head $\langle b, h \rangle$ to predict label *l*.

Under each head-selection setting, for fair com-547 548 parison, we set a limit of the total number of syntactical heads $\sum_{l \in \mathcal{L} \cup \{\phi\}} |\mathcal{H}_l|$ of 2000.

For the tree-construction alternative, we use 551 Raw attention score: Under this setting, we're still ⁵⁵² using the estimated MI as head importance criteria,

⁵⁵³ while for a specific head $\langle b, h \rangle \in \mathcal{H}_l$, we use the ⁶⁰¹ 4.3 IPBP Structural Alternatives statention score $a_{b,h}^{[\cdot][\cdot]}$ instead of the posterior 602 Apart from comparing with previous methods, $f(A_{h,h},L)$ in the reconstruction algorithm. This 603 we're also curious about our model's designs. ⁵⁵⁶ simple intuitive is the common underlying princi-⁶⁰⁴ Therefore, we propose two alternative structures: ⁵⁵⁷ ple of previous works focusing on attention (Clark ⁶⁰⁵ 558 et al., 2019; Vig and Belinkov, 2019; Ravishankar 606 samples exhibit a long-tail characteristic: most 559 et al., 2021). We found that due to the absence of 607 samples come from A_{ϕ} , since most pairs of words ⁵⁶⁰ our estimated posteriors, if $\sum_{l \in \mathcal{L} \cup \{\phi\}} |\mathcal{H}_l|$ reaches ⁶⁰⁸ don't have dependency arc in between. \mathcal{A}_{ϕ} might ⁵⁶¹ 2000, the scores of all heads will be rather noisy. ⁶⁰⁹ be noisy, consisting of various non-syntactic inter-⁵⁶² Therefore, we choose to select top-k heads based ⁶¹⁰ token relationships, and MI estimations based on 563 on MI for each label. We did a grid search and 611 samples in \mathcal{A}_{ϕ} might be affected by this long tail ⁵⁶⁴ found the top-8 settings have ideal performance.

Model. Dataset and Metrics 565 4.2

567 open llama 7b is a decoder-based LLM consisting 616 also calculated a more syntactical MI, namely $_{568}$ of $3\overline{2}$ layers and 32 attention heads within each $_{617}$ MI_{pos}, with the following formulation: ⁵⁶⁹ layer. Compared with pre-trained language models 570 like BERT (Devlin et al., 2019), open llama 7b 618 571 might consist of attention heads with rather varied 572 functionalities, offering more insights under the 620 contemporary LLM research context. 573

574 ⁵⁷⁴ Specifically, open_nama_10 is a decoder based ⁵⁷⁵ model having triangular-masked attention scores. ⁶²² by $\hat{P}_{pos}(L = l) = \frac{\mathcal{A}_{b,h;l}}{\sum_{l' \in \mathcal{L}} \mathcal{A}_{b,h;l'}}$ and $\hat{f}_{pos}(L, A_{b,h}) =$ 576 In implementation, we cache the Key Values of each attention head and use them to re-calculate the ${}^{623} \hat{f}(A_{b,h}|L)\hat{P}_{pos}(L)$. During implementation, we'll 577 unmasked attention scores. While our reconstruc- 624 use a balance factor α and calculate the mixed MI tion is inevitably introducing "useless" attention ⁶²⁵ $MI_{mix}(\cdot;\cdot) = \alpha MI_{binary}(\cdot;\cdot) + (1 - \alpha) MI_{pos}(\cdot;\cdot)$ 579 ⁵⁸⁰ scores, we think that it is still necessary for two rea- ⁶²⁶ Arc First: Unlike previous methods, we're disons: 1. Making compromises to the decoder struc- 627 rectly obtaining labeled dependency trees, bypass-581 ⁵⁸² ture will hinder our method from applying to non-⁶²⁸ ing the process of dependency arc predicting. 583 decoder models (Chung et al., 2024; Zeng et al., 629 We're curious about whether it's a good choice. 2024), thus less universal. 2. As sentences become 630 Under this setting, instead of estimating 584 ⁵⁸⁵ longer, the softmax-normalized scores will be di- ⁶³¹ $\hat{f}(A_{b,h}|L=l)$, we'll directly estimate the unla-⁵⁸⁶ luted. This is more serious for triangular attention ₆₃₂ beled likelihoods $\hat{f}(A_{b,h}|L \in \mathcal{L})$ and $\hat{f}(A_{b,h}|L =$ since it has rows of varying lengths. While softmax $_{633} \phi$), and calculate the corresponding multivariate is not bijective, using cached QK to reconstruct the 634 probabilities together with corresponding MI val-588 unnormalized scores is inevitable. 589

Following previous supervised dependency 590 parsing works (Tian et al., 2022), we use Universal 591 ⁵⁹² Dependencies (UD) 2.9 (Zeman et al., 2021), as da-⁵⁹³ taset, with 39832 sentences in the training set and ⁵⁹⁴ 1700 sentences in the validation set. UD 2.9 is an 595 English treebank covering texts from multiple 596 sources like literature, news articles, spoken lan-⁵⁹⁷ guages, etc., with diverse morphological and gram-598 matical features. We also use labeled attachment 599 scores (LAS), and unlabeled attachment scores 600 (UAS) as metrics.

Positive MI: We noticed that the attention score ⁶¹² noisy distribution. Other score sets $\mathcal{A}_1, \dots \mathcal{A}_{|\mathcal{L}|}$ are 613 having approximately the same magnitudes and 614 their corresponding token pairs are guaranteed to ⁵⁶⁶ We're using open llama 7b¹ as our probed model. ⁶¹⁵ have any dependency relationship. Therefore, we

$$\mathsf{MI}_{\mathsf{pos}}(L; A_{b,h}) = \sum_{l \in \mathcal{L}} \int f_{\mathsf{pos}}(l, a) \log \frac{f_{\mathsf{pos}}(l, a)}{P_{\mathsf{pos}}(l) f_{\mathsf{pos}}(a)} \mathrm{d}a$$
(12)

In that equation, $P_{pos}(\cdot)$, $f_{pos}(\cdot, \cdot)$ actually stand Specifically, open_llama_7b is a decoder-based ⁶²¹ for conditional possibilities when $l \neq \phi$, estimated

635 ues. We'll compare UAS to check the quality of re-636 constructed unlabeled trees.

637 Transposed: Sometimes, we're unsure whether 638 the attended token acts as a dependency head, or a 639 dependant. So we'll let $l^{[i][j]}$ correspond to $a^{[j][i]}_{b,h}$ 640 and repeat the whole IPBP process in this setting.

641 **4.4 Result and Analysis**

642 Results are shown in Table 1. We can see that our 643 method is overperforming all competitive baselines, 644 including the state-of-the-art conditional entropy

¹https://github.com/openlm-research/open llama

Method	UAS	LAS
Probeless	34.8	20.9
IoU	38.3	26.6
ElasticNet	41.9	31.3
V-Information	41.3	20.9
Raw Score	32.3	16.6
IPBP	<u>49.1</u>	<u>30.6</u>
IPBP (transposed)	42.6	28.0
$IPBP + MI_{pos}$	<u>49.9</u>	<u>34.8</u>
IPBP (arc only)	36.5	N/A

Table 1: Results of our IPBP and different baselines.

646 principles with our method while requiring a much 697 The conclusion is consistent with our intuitions and 647 more computational budget. In fact, our imple- 698 worth hypothesizing, but it was never justified by ⁶⁴⁸ mented \mathcal{V} -Information MLP is optimized using ⁶⁹⁹ previous works, which either lack good MI-like cri-649 several tricks (see Appendix B), while during its 700 teria or focus on unlabeled trees. Thanks to the 652 drawn that supervised methods may still fall behind 705 these top-10 labels, we calculated the Pearson cor-656 657 ready selecting attention heads based on our esti- 710 lation. 659 660 mated MI, it still has a great performance gap with our posterior-based method, further justifying the 711 5 661 necessity of our posterior-based algorithm.

For structure alternatives, we notice that incor-663 porating MI_{pos} will give performance benefits, 664 665 shedding light on potential improvements to our 666 methods. The transposed setting will still capture a 667 relatively smaller portion of dependencies. Last but not least, by comparing with our arc-based baseline, *high-quality* trees, and also *transparent for expla-*669 we'll find we're actually at a triangular balance, we 670 probed for more accurate, also labeled trees, while 671 choosing a more straightforward method, with no 672 need for individual arc probing.

673 4.5 **Further Analysis**

675 grained analysis of our reconstructed trees and estimated MI values. Instead of listing up MI and doing trivial analyses, we decide to provide two intriguing and informative conclusions, giving inspira-678 tions to upcoming works. 679

The first conclusion is that decoder models 681 adaptively capture look-back/ahead dependen-682 cies: Since the masked decoder attention can only

look back, there are good reasons that dependencies 683 that also look back (pointing to front words) can be well captured. What makes it more intriguing is 685 that dependencies looking ahead might also be captured in a look-back manner. We draw this conclusion by comparing the top-10 most well-reconstructed labels between original IPBP and the transposed alternative. We find that there're more lookahead dependencies (5 of 10) under the transposed setting compared with the original setting (3 of 10). The second conclusion is that model layers corre-693 ⁶⁹⁴ spond to tree layers to some extent: lower layers 695 are for local/phrasal dependencies, while higher 645 estimation method, V-Information, which shares 696 layers are for global/sentence-wide dependencies. training, we still find that the trained MLPs are rel- 701 fine-grained MI, we can calculate the MI-weighted atively good at detecting arcs while having poor 702 layer indices for each label, where smaller performances on labeling. This aligns with its low 703 weighted indices indicate dependencies having LAS in tree-construction results. An insight can be 704 more lower-layer heads responsible for it. Among statistical ones, especially when the data is long- 706 relation coefficient ρ between the weighted layer tailed or low-dimensional. Moreover, even though 707 index and average depth (maximum distance to leaf the head-selection settings for the raw-score 708 nodes) of each dependency label, getting a result of method is specifically tuned, and the method is al- 709 0.69 with p=0.03 for a null hypothesis of no corre-

Conclusion

712 We proposed a method that can estimate MI and 713 reconstruct labeled dependency trees without intro-714 ducing any trainable networks. Indeed, our method 715 is achieving an "impossible triangle": it has simpler 716 architectures requiring negligible computation 717 budgets, while producing more complicated and 719 nation, meaning that researchers can get fine-720 grained head-level MI estimation, and a bunch of 721 intuitive probability functions, without worrying 722 about did my network furtively learnt the task? 723 Through comparing with a series of competitive 724 baselines, we ensured its effectiveness, and then 674 Like previous probing works, we'll also do fine- 725 made two informative conclusions based on our es-726 timated MI and reconstructed trees. The number of 727 conclusions is limited due to content limit, and 728 since our method is providing an analytical back-729 bone, we strongly appeal to future research for 730 fine-grained analysis on those estimated MI values 731 and distributions.

732 Limitations

733 Despite its efficiency, our method still has several 734 shortages: The most important one is that, to pre-785 735 vent the problem from being too complicated and 786 736 bounded by the curse of dimensionality, our 787 737 method does not consider the multivariate case, 788 738 taking an assumption that all attention heads are in- 789 790 739 dependent. Moreover, as mentioned in Section 4.2, 740 the introduction of "useless" attention scores is also 791 792 741 noticeable, meaning that the density estimations ⁷⁴² might contain noises. Lastly, our method is only ap-⁷⁹³ F 743 plicable to discrete-continuous mixtures, where all 794 744 probed concepts are discrete labels, but not appli-745 cable to multivariate continuous joint distributions. 797

746 Ethical Considerations

747 Since our method is an explanation method, read- 800 r48 ers might exploit our method to perform syntactical 749 attacks, like getting poorly captured dependency 750 labels and designing specific prompts to confuse 751 models. For models that are put into use in produc-752 tion environments, this might cause unexpectable 808 753 effects. 807

754

755 References

Guillaume Alain and Yoshua Bengio. 2017. Under- 811 756

standing intermediate layers using linear classifier 812 757

- probes. In The 5th International Conference on 758
- Learning Representations. 759

760 Omer Antverg and Yonatan Belinkov. 2021. On the Pit- 815

falls of Analyzing Individual Neurons in Language 816 761

Models. In International Conference on Learning 762

- Representations. 763
- 764 Hyung Won Chung, Le Hou, Shayne Longpre, Barret 819 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi 820 765 Wang, Mostafa Dehghani, Siddhartha Brahma, Al-766 bert Webson, Shixiang Shane Gu, Zhuyun Dai, Mi-767 822 rac Suzgun, Xinyun Chen, Aakanksha Chowdhery, 768 Alex Castro-Ros, Marie Pellat, Kevin Robinson, 769 Dasha Valter, Sharan Narang, Gaurav Mishra, Ad-770 ams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, 825 Yimin Fan, Fahim Dalvi, Nadir Durrani, and Hassan 771 Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja- 826 772 cob Devlin, Adam Roberts, Denny Zhou, Quoc V. 827 773 Le and Jason Wei. 2024. Scaling instruction-fine- 828 774 tuned language models. Journal of Machine Learn-775 ing Research, 25(70), 1-53. 776 830 777 Kevin Clark, Urvashi Khandelwal, Omer Levy, and 831 Christopher D. Manning. 2019. What Does BERT 832 778
- Look at? An Analysis of BERT's Attention. In Pro-779
- ceedings of the 2019 ACL Workshop BlackboxNLP: 780
- 834 Analyzing and Interpreting Neural Networks for 781 835

NLP, pages 276-286, Florence, Italy. Association for Computational Linguistics.

782

783

- 784 Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single \$&!#* vector: Probing sentence embeddings for linguistic properties. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics
- ahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. 2019. What is one grain of sand in the desert? analyzing individual neurons in deep nlp models. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 33, No. 01, pp. 6309-6317). 798
- 799 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171-4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 808 Dror Dotan and Nadin Brutmann. 2022. Syntactic chunking reveals a core syntactic representation of 809 multi-digit numbers, which is generative and auto-810 matic. Cognitive research: principles and implications, 7(1), 58.
- 813 Timothy Dozat and Christopher D. Manning. Deep Biaffine Attention for Neural Dependency Parsing. 2017. In International Conference on Learning Representations. 2017.
- 817 Jason M. Eisner. 1996. Three New Probabilistic Models for Dependency Parsing: An Exploration. In COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics.
- 821 Jacqueline Fallon and Liina Pylkkänen. 2024. Language at a glance: How our brains grasp linguistic structure from parallel visual input. Science Advances, 10(43), eadr9951.
 - Sajjad. 2024. Evaluating neuron interpretation methods of nlp models. Advances in Neural Information Processing Systems, 36.
- 829 Weihao Gao, Sreeram Kannan, Sewoong Oh, and Pramod Viswanath. 2017. Estimating mutual information for discrete-continuous mixtures. Advances in neural information processing systems, 30.
- 833 John Hewitt and Percy Liang. 2019. Designing and Interpreting Probes with Control Tasks. In Proceedings of the 2019 Conference on Empirical Methods

814

in Natural Language Processing and the 9th Inter- 891 836

national Joint Conference on Natural Language 892 837

Processing (EMNLP-IJCNLP), pages 2733–2743, 893 838

Hong Kong, China. Association for Computational 894 839 Linguistics. 840

841 John Hewitt and Christopher D. Manning. 2019. A 896 Structural Probe for Finding Syntax in Word Repre- 897 842 sentations. In Proceedings of the 2019 Conference 898 843 of the North American Chapter of the Association 844 for Computational Linguistics: Human Language 845 900 Technologies, Volume 1 (Long and Short Papers), 846 901

pages 4129-4138, Minneapolis, Minnesota. Associ-847

ation for Computational Linguistics. 848

849 Phu Mon Htut, Jason Phang, Shikha Bordia, and Sam- 904

uel R. Bowman. 2019. Do attention heads in BERT 905 850 syntactic dependencies?. arXiv preprint 906 track 851

arXiv:1911.12246. 852 907

- 853 Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 908 Ti 2019. What Does BERT Learn about the Structure 909 854
- of Language?. In Proceedings of the 57th Annual 910 855
- Meeting of the Association for Computational Lin- 911 856

guistics, pages 3651-3657, Florence, Italy. Associa- 912 857

tion for Computational Linguistics. 858

914 Yair Lakretz, German Kruszewski, Theo Desbordes, 859

- Dieuwke Hupkes, Stanislas Dehaene, and Marco 915 Tiago Pimentel, Naomi Saphra, Adina Williams, and 860
- Baroni. 2019. The emergence of number and syntax 916 861
- units in LSTM language models. In Proceedings of 917 862
- the 2019 Conference of the North American Chapter 918 863
- of the Association for Computational Linguistics: 919 864
- Human Language Technologies, Volume 1 (Long 920 865
- and Short Papers), pages 11-20, Minneapolis, Min-866 921
- nesota. Association for Computational Linguistics. 867

868 Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Mat- 923 thew E. Peters, and Noah A. Smith. 2019. Linguistic 924 869 Knowledge and Transferability of Contextual Rep- 925 870 resentations. In Proceedings of the 2019 Conference 926 871 of the North American Chapter of the Association 927 872 for Computational Linguistics: Human Language 873 Technologies, Volume 1 (Long and Short Papers), 874 929 pages 1073-1094, Minneapolis, Minnesota. Associ-875 930

- ation for Computational Linguistics. 876
- lessandro Lopopolo, Antal Van den Bosch, Karl-877 A Magnus Petersson, and Roel M. Willems. 2020. Dis-878
- 933 tinguishing Syntactic Operations in the Brain: De-879 934
- pendency and Phrase-Structure Parsing. Neurobiol-880
- ogy of Language, 2, 152 175. 881
- Young-Il Moon, Balaji Rajagopalan, and Upmanu Lall. 937 882
- 1995. Estimation of mutual information using ker- 938 883
- nel density estimators. Physical Review E, 52(3), 884 2318. 885
- 886 Jesse Mu and Jacob Andreas. 2020. Compositional ex- 941
- planations of neurons. Advances in Neural Infor- 942 887 mation Processing Systems, 33, 17153-17163. 888 943
- 944 889 Max Müller-Eberstein, Rob van der Goot, and Barbara 945 Plank. 2022. Probing for Labeled Dependency 890

Trees. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7711-7726, Dublin, Ireland. Association for Computational Linguistics.

- 895 Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. Computational Linguistics, 19(2):313–330.
- 899 Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 2227-2237, New Orleans, Louisiana. Association for Computational Linguistics.
 - ago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. 2020a. Information-Theoretic Probing for Linguistic Structure. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4609-4622, Online. Association for Computational Linguistics.
 - Ryan Cotterell. 2020b. Pareto Probing: Trading Off Accuracy for Complexity. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 3138–3153, Online. Association for Computational Linguistics.
 - Tiago Pimentel, Josef Valvoda, Niklas Stoehr, and Ryan Cotterell. 2022. Attentional Probe: Estimating a Module's Functional Potential. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11459-11472, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 928 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. "Language are unsupervised multitask learnmodels ers." OpenAI blog 1, no. 8 (2019): 9.
- init Ravishankar, Artur Kulmizev, Mostafa Abdou, 932 Anders Søgaard, and Joakim Nivre. 2021. Attention Can Reflect Syntactic Structure (If You Let It). In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 3031-3045, Online. Association for Computational Linguistics.
- 939 Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, 940 Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das and Ellie Pavlick. 2019. What do you learn from context? Probing for sentence structure in contextualized word representations. In International Conference on Learning Representations.

935

936

902

903

913

946 Yuanhe Tian, Yan Song, and Fei Xia. 2022. Enhancing

947 Structure-aware Encoder with Extremely Limited

Data for Graph-based Dependency Parsing. In *Pro-*

949 ceedings of the 29th International Conference on

950 Computational Linguistics, pages 5438–5449,

951 Gyeongju, Republic of Korea. International Com-

⁹⁵² mittee on Computational Linguistics.

953 Jesse Vig and Yonatan Belinkov. 2019. Analyzing the

954 Structure of Attention in a Transformer Language

Model. In Proceedings of the 2019 ACL Workshop

956 BlackboxNLP: Analyzing and Interpreting Neural

⁹⁵⁷ Networks for NLP, pages 63–76, Florence, Italy. As-

⁹⁵⁸ sociation for Computational Linguistics.

959 Elena Voita and Ivan Titov. 2020. Information-Theo-

retic Probing with Minimum Description Length.

In Proceedings of the 2020 Conference on Empiri-

962 cal Methods in Natural Language Processing

(EMNLP), pages 183–196, Online. Association for
 Computational Linguistics.

965 Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stew-

art, and Stefano Ermon. 2020. A Theory of Usable

967 Information under Computational Constraints.

⁹⁶⁸ In International Conference on Learning Represen-

969 tations.

Daniel Zeman, Nivre Joakim, Abrams Mitchell, Acker mann Elia, and others. 2021. Universal Dependen-

cies 2.9. LINDAT/CLARIAH-CZ digital library at

the Institute of Formal and Applied Linguistics

974 (ÚFAL), Faculty of Mathematics and Physics,

975 Charles University.

Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, 976 Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, 977 Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, 978 Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie 979 Tang, Jing Zhang, Jingyu Sun, Juanzi Li, Lei Zhao, 980 Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie 981 Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi 982 Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, 983 Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, 984 Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, 985 Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai 986 Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, 987 Yuevan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, 988 Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu 989 Hou, Zihan Wang. 2024. Chatglm: A family of large 990 language models from glm-130b to glm-4 all 991

⁹⁹² tools. *arXiv preprint arXiv:2406.12793*.

Implementation Details of IPBP 994 A

⁹⁹⁵ In this section, we'll briefly introduce the imple-⁹⁹⁶ mentation details, like the hyperparameters and key $_{1037}$ distribution estimations are only required to be ⁹⁹⁷ algorithms we use to implement IPBP.

As shown by the source code, we use the 998 ⁹⁹⁹ PyTorch framework to implement the whole IPBP 1000 process. We're not relying on off-the-shelf pack-1041 1001 ages that have KDE functionalities like SciPy and Scikit-Learn, since their KDE implementation is CPU-based and thus too inefficient under our ex-¹⁰⁴³ take the estimated posteriors as a set of n_x discrete 1002 1003 periment settings. 1004

1005 whole long tensor $a_{b,h,l} \in \mathbb{R}^{|\mathcal{A}_{b,h;l}|}$. We calculate the polating between x_i and x_{i+1} . The interpolation-tion whole long tensor $a_{b,h,l} \in \mathbb{R}^{|\mathcal{A}_{b,h;l}|}$. We calculate the polating between x_i and x_{i+1} . The interpolation-tion whole long tensor $a_{b,h,l} \in \mathbb{R}^{|\mathcal{A}_{b,h;l}|}$. We calculate ¹⁰⁰⁷ the minimum and maximum values of $\mathcal{A}_{b,h;l}$, and ¹⁰⁴⁸ tioned in Section 3.5, like head selection and scorebuild a tensor of real numbers $X = \{x_1, x_2, \dots, x_{n_x}\}_{1049}$ weighted averaging, are all parallel-optimized, re-1009 ensuring that $x_1 < \min \mathcal{A}_{b,h;l}, x_{n_x} > \max \mathcal{A}_{b,h;l}$, 1050 sulting in being able to run inference within 5 1010 and $x_1 < x_2 < \cdots < x_{n_x}$. These discrete x values 1051 minutes on all baseline settings on a 4090 GPU. ¹⁰¹¹ serve as the points to calculate densities. Next, we¹⁰⁵² 1012 calculate the mutual differences between each 1053 values like MI_{binary}, MI_{pos}, we use the trapezoid ¹⁰¹³ point in **X** and each element in $a_{b,h,l}$, by repeating ¹⁰⁵⁴ method to estimate the integral value: as mentioned $|\mathcal{A}_{b,h;l}|$ times

1016 times, also getting a matrix $\begin{bmatrix} a_{b,h,l}^T & \cdots & a_{b,h,l}^T \end{bmatrix}^T$ of 1059 them up to get the integral values. 1017 shape $n \ge 1$ and $[a_{b,h,l}^T = a_{b,h,l}^T]^T$ of 1060 During tree reconstruction we ¹⁰¹⁷ shape $n_x \times |\mathcal{A}_{b,h;l}|$. The absolute differences of ¹⁰⁶¹ the total number of heads, i.e., $\sum_{l \in \mathcal{L} \cup \{\phi\}} |\mathcal{H}_l|$ to a ¹⁰¹⁸ the two matrices $|[X^T, ..., X^T] - [a_{b,h,l}^T, ..., a_{b,h,l}^T]|$, ¹⁰⁶² fixed value (2000), and use binary search. differences, let's say¹⁰⁶³ 1019 are the mutual

¹⁰²³ bandwidth B, with all weight w_i s equal to 1. We ¹⁰⁷⁰ following its intended usage.

1024 then calculate element-wise, following the follow-1025 ing equation:

1026
$$\frac{1}{B\sqrt{2\pi\cdot\sigma_{\mathcal{A}_{b,h;l}}}}\exp\left\{-\left(\frac{D(X^{T},\boldsymbol{a}_{b,h,l})}{B}\right)^{2}\right\}$$
 (

1027 ¹⁰²⁸ in shape $n_x \times |\mathcal{A}_{b,h;l}|$, we then calculate the row-¹⁰⁷⁵ the long-tail essential discussed in Section 4.3. This wise mean of the kernel values to get the final ker-1076 will result in many infinite V-Information values, 1029 1030 nel density values in shape n_x . Since all operations 1077 since there will be many estimated probabilities 1031 of this process are element-wise matrix operations, 1078 (for label in \mathcal{L} other than ϕ) rather close to zero. 1032 this is easily parallel-optimizable by PyTorch. As a 1079 Therefore, we apply a sample balancing technique, 1033 result, the computation time for extracting attention 1080 truncating $\mathcal{A}_{...,\phi}$ to make their numbers of samples

¹⁰³⁴ score sets (\mathcal{A}_{\dots}) and performing all kernel density 1035 estimations is within 1 hour using a single RTX 1036 4090 GPU. Since attention score allocating and 1038 done once, our method is extremely time-saving 1039 compared to most of the supervised probing meth-1040 ods.

For inferring on estimated probabilities (like in-1042 ferring on posteriors $\hat{f}(l|A_{b,h})$ in Section 3.5, we 1044 points, and an attention score in range $[x_i, x_{i+1}]$ Specifically, we take samples in $\mathcal{A}_{b,h;l}$ as a¹⁰⁴⁵ will get its corresponding posterior value by inter-

For calculating integrals, specifically, the MI 1055 in the section before, the kernel densities are de-¹⁰¹⁴ X^T for n_x times, getting a matrix $[X^T, ..., X^T]$ of ¹⁰⁵⁶ scribed by n_x points, we take the n_x points as the 1015 shape $n_x \times |\mathcal{A}_{b,h;l}|$, and repeating $a_{b,h,l}$ for n_x n_x times n_x tinter n_x times n_x times n_x times n_x times n_x

During tree reconstruction, we empirically set

We're using Universal Dependencies 2.9 ¹⁰¹⁹ are the mutual differences, let's say ¹⁰⁶³ we're using oniversal Dependencies 2.7 ¹⁰²⁰ $D(X^T, a_{b,h,l}) \in \mathbb{R}^{n_X \times |\mathcal{A}_{b,h;l}|}$. Then we calculate ¹⁰⁶⁴ (Zeman et al., 2021) as our dataset. That dataset is ¹⁰²¹ the standard deviation of $\mathcal{A}_{b,h;l}$, i.e., $\sigma_{\mathcal{A}_{b,h;l}}$, and ¹⁰⁶⁵ publicly available, using CC BY-SA 4.0 license² al-¹⁰²² take a rule-of-thumb value $\left(\frac{1}{\sum_{i=1}^{|\mathcal{A}_{b,h;l}|}w_i^2}\right)^{-\frac{1}{5}}$ ¹⁰⁶⁷ Universal Dependencies (UD) dataset is designed ¹⁰⁶⁹ cal identifications for NLP researchers, so we're

1071 B Implementation Details of Baselines

13)₁₀₇₂ For the \mathcal{V} -Information MLPs, we found that train-1073 ing on all datasets will result in a network always In order to get the kernel values, which are also 1074 predicting ϕ for all possible attention scores, due to

² https://creativecommons.org/licenses/by-nc-sa/4.0/deed.en

1081 the same as the total number of samples in other 1082 score sets $\mathcal{A}_{\cdot,\cdot;1}, \mathcal{A}_{\cdot,\cdot;2}, \dots \mathcal{A}_{\cdot,\cdot;|\mathcal{L}|}$. What's more, we 1083 also did a search on several network sizes, and 1084 found that if MLP(·) is $W_2(\operatorname{act}(W_1(\cdot)))$, where 1085 W_1 in shape 1×2 and a W_2 in shape 2×4 1086 achieves better fitting. This aligns with (Pimentel 1087 et al., 2020a) to some extent. We also use PyTorch 1088 to implement the baselines. Specifically, for Elas-1089 ticNet that requires additional training, we use 1090 AdamW optimizer, 1e - 5 for both λ_1 and λ_2 , and 1091 use a constant learning rate of 1e - 3, training for 1092 12 epochs. For the \mathcal{V} -Information MLP, since we 1093 need $\mathcal{bh} \times (|\mathcal{L}| + 1)$ individual networks for pre-1094 dicting the alternatives of binary MI, we initialize 1095 $bh \times (|\mathcal{L}| + 1)$ sets of matrices, each constituting weights of specific 1096 the а network 1097 W_1, W_2, \dots, W_n layers, with W_1 having a dimension 1098 of 1 and $W_{n \text{ lavers}}$ having a dimension of $|\mathcal{L}| + 1$. 1099 During training and inferencing, we concatenate all 1100 attention scores $a_{b,h}^{[i][j]}$ for any $b \in \{1 \dots b\}$ and 1101 $h \in \{1 \dots h\}$ into a tensor of shape $\mathcal{B}h$, and use 1102 torch.bmm to map each element of that tensor to 1103 $\mathcal{Bh} \times (|\mathcal{L}| + 1)$ probabilities (standing for the 1104 probabilities of each label conditioned on each at-1105 tention head's attention score, estimated by the var-1106 iational family). Using torch.bmm will avoid 1107 training $bh \times (|\mathcal{L}| + 1)$ networks separately, 1108 which is a disaster on computation loads, and can 1109 exploit GPU's parallel processing abilities. We use 1110 leaky relu between hidden layers and use sigmoid 1111 to form the final probabilities. We use 1e - 2 as 1112 learning rate with exponential decay (0.8 at each 1113 epoch), together with an additional warmup epoch 1114 at the beginning. The hyperparameters differ for \mathcal{V} -1115 Information since otherwise the variational family 1116 network will be more poorly trained. We also 1117 trained for 12 epochs.