# BANYAN: IMPROVED REPRESENTATION LEARNING WITH EXPLICIT STRUCTURE

Anonymous authors

Paper under double-blind review

# ABSTRACT

We present Banyan, a model that efficiently learns semantic representations by leveraging an inductive bias towards explicit hierarchical structure. Although typical transformer-based models excel at scale, they struggle in low-resource settings. Recent work on models exploiting explicit structure has shown promise as efficient learners in resource-constrained environments. However, these models have yet to demonstrate truly competitive performance. Banyan bridges this gap, significantly improving upon prior structured models and providing, for the first time, a viable alternative to transformer embeddings for under-represented languages. We achieve these improvements through two key innovations 1) A novel entangled tree structure that resolves multiple constituent structures into a single shared one, explicitly incorporating global context. 2) Diagonalized message passing functions that increase the influence of the inductive bias. Our final model has just 14 non-embedding parameters yet is competitive with baselines many orders of magnitude larger. Banyan outperforms its structured predecessors and competes with large unstructured models across various semantic tasks in multiple languages. Notably, it excels in low-resource settings, highlighting its potential for efficient and interpretable NLP in resource-constrained environments. These results underscore the value of appropriate inductive biases in capturing semantic relationships and open new avenues for efficient, interpretable NLP models.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

## 1 INTRODUCTION

Semantic representations of text are important for many NLP applications such as retrieval augmented generation (Lewis et al., 2020), question answering, and summarisation (Abdalla et al., 2023; Wang et al., 2022). They are also useful for clustering and organising textual data when labelled training sets are not available. At the time of writing such representations are primarily generated by large scale transformer models (Vaswani et al., 2017). These models are incredibly effective, but training them usually requires scale, both in terms of data and compute.

An alternative approach is to take inspiration from linguistics/cognitive science and explicitly incorporate structured compositions. Put simply, composition states that all you need to understand 040 the semantics of a whole are the meanings of its parts and the structure that dictates how they fit 041 together (Chomsky, 1956; Crain & Nakayama, 1987; Pallier et al., 2011; de Marneffe et al., 2006). 042 This is a very efficient principle, because novel utterances can broken down into familiar parts 043 using systematic rules, rather than having to store the meaning of each utterance individually. It 044 is thought that this principle lets humans generalise from (comparatively) little data and makes us efficient learners (Fodor & Pylyshyn, 1988; Lake et al., 2016; Ito et al., 2022; Wiedemer et al., 2023). In order to explicitly incorporate an inductive bias of this kind we need to change the modelling 046 process somewhat. Rather than keeping all the information flow internal, models must now learn 047 representations for the atomics, operate on (and/or learn) a discrete graph that dictates the mode of 048 combination, and learn functions that control information flow through such a graph. Models of this kind have demonstrated improved language modelling perplexity at cognitively plausible scales (Hu et al., 2021; 2022); better systematic generalisation (Sartran et al., 2022; Murty et al., 2023); and, 051 importantly for this paper, the ability to efficiently acquire semantics (Opper et al., 2023). 052

053 Opper et al. (2023) introduce a model called the Self-StrAE, which learns to representations which have to explicitly model compositional semantics. This demonstrated very promising performance

054 while requiring minimal resources. Both in terms of data and model size. Opening the door to 055 investigate whether more compute efficient solutions can be found for learning semantic representa-056 tions. This would be particularly useful for low resourced languages where relying on scale is not a 057 generally feasible solution. However, the Self-StrAE, while promising, still lags behind large scale 058 pre-trained transformers, even in langauges which fall outside of standard pre-training corpora. In this paper we introduce a model called Banyan, which significantly improves performance over that of (Opper et al., 2023), while simultaneously achieving even greater resource efficiency. We achieve 060 this by changing the form of the structure optimised to a graph that models global relations between 061 nodes which we call an entangled trees, as well as a message passing regime based on diagonal 062 functions which reduces parameters while producing more expressive representations. Our model, 063 Banyan, achieves competitive performance with transformer based baselines, and for the first time 064 represents a low cost yet viable alternative for producing representations for low resource languages, 065 measured using semantic textual similarity (STS) tasks. By leveraging cognitively inspired inductive 066 biases we can achieve performance comparable or better than large scale pre-trained LLMs but with 067 only 14 non-embedding parameters.

068 069

070

# 2 BACKGROUND AND RELATED WORK

Banyan is a graph neural network, specifically a recursive neural network (RvNN), that learns both structure and representations. Before detailing the model, we unpack these terms in this section.

Recursive Neural Networks: Like their recurrent cousins, recursive neural networks operate by repeatedly applying a function to update a the network state in an ordered fashion. However, rather than utilising temporal ordering (i.e. over a sequence), RvNNs operate according to some hierarchical structure, typically this given as input and most often it is a binary tree. They can be applied bottom-up (traversing from leaves to root) or top-down (from root to leaves) or both. First popularised by Socher et al. (2011; 2013), they have inspired numerous successor frameworks which differ in how the recursive function is defined. These successors include the Tree-LSTM (Tai et al., 2015), IORNN (Le & Zuidema, 2014; Ji & Eisenstein, 2015) and also Banyan.

082 Learning Structure: RvNNs typically require structure as input, sometimes such structure is 083 available or can be obtained using existing tools, but generally this is quite a limiting factor, because it limits model flexibility. A solution is to incorporate a mechanism within the model that is 084 able to induce the structure during the recursive computation. Prior approaches include the use 085 of differentiable chart parsing (Drozdov et al., 2019; 2020; Hu et al., 2021; 2022), beam search (Ray Chowdhury & Caragea, 2023), continuous relaxation (Chowdhury & Caragea, 2021; Soulos 087 et al., 2024), or reinforcement learning (Havrylov et al., 2019). While successful these solutions 880 can suffer from memory issues and hyperparameter sensitivity. In this paper we adopt the approach 089 of Opper et al. (2023) and utilise representation similarity to dictate merge order. This is both computationally inexpensive and surprisingly effective. 091

Semantic Representations of Text: Systems like Word2Vec (Mikolov et al., 2013) and GLoVe 092 (Pennington et al., 2014) model the semantics of words using the distributional hypothesis (Harris, 093 1954). This hypothesis states that the context a word is used in defines its meaning. Consequently, 094 representations are learned by setting a context window of some fixed size and then using that to 095 predict the missing word. This approach proved very effective for a long time, but words don't all 096 have one meaning - it changes in context. A natural solution is to use transformers. Initially smaller 097 (relative to today) encoder only models produced poor representations (Reimers & Gurevych, 2019). 098 However, at time of writing transformers with optional contrastive finetuning (Gao et al., 2021) have become the model of choice for producing semantic representations. 099

100 Semantic Representation Learning through Structure: Transformer embeddings are more success-101 ful than static word embeddings because they are allow for flexible contextualisation. The meaning 102 of a word can change in context and will more strongly influenced by certain neighbours rather than 103 others. A transformer can model this phenomenon by routing information between specific tokens via 104 its attention. An alternative approach is to use structure, where rather than attention dictating routing, 105 it is done by an explicit graph. Early work in this area focused on lexical semantics, using dependency parses to determine a more focused context window (Levy & Goldberg, 2014; Vashishth et al., 2019). 106 More recent work by Opper et al. (2023) use constituency parses in order to learn embeddings at both 107 the word and the sentence level. They introduce two variants of their model. StrAE: which takes

structure as input, and Self-StrAE: which learns its own structure with the representations. This later
 model, the Self-StrAE, is the starting point from which we build Banyan, and will be outlined in
 more detail in the subsequent section.

# 3 PRELIMINARY: SELF-STRAE

114 Self-StrAE involves three main components that acts 115 over a sentence  $\mathbf{w} = \langle w_n \rangle_{n=1}^N$  represented as a sequence 116 of tokens. These are: (a) a procedure to determine 117 which tokens to merge and in what order, (b) message 118 passing functions-composition and decomposition-that 119 merge and split embeddings respectively, and (c) a 120 reconstructive objective that leverages both the induced 121 structure and embeddings. While we refer the reader to 122 Opper et al. (2023) for a detailed description of these components, we briefly recap its operation to provide 123 sufficient background for the development of Banyan (§ 4). 124



Figure 1: Self-StrAE operation. Red lines indicate cosine similarity. Shared colours imply shared parameters.

(3)

126 At a high-level, Self-StrAE learns representations that both define their own structure and are in 127 turn defined by it. This is achieved by first *embedding* tokens to form an initial frontier using an 128 embedding matrix  $\Omega_{\Psi}$ . Next it then takes the adjacent tokens with the highest cosine similarity to each other (ate and doughnuts in Figure 1) and merges them into a single embedding using a 129 parametric *composition function*  $C_{\Phi}$ . This procedure is repeated until the sequence reduces to a 130 single root embedding. The resulting merge history is then treated as the induced binary tree for 131 the sentence. Self-StrAE then traverses back down the structure, recursively splitting embeddings 132 at every node using a parametric *decomposition function*  $D_{\Theta}$  to recover embeddings for the leaves. 133 Finally, the model can optionally use a *dembedding function*  $\Lambda_{\Gamma}$  to predict tokens  $\hat{w}_n$  from these leaf 134 embeddings. Figure 1 illustrates the autoencoding process. 135

Intuitively, this means that the model starts from random embeddings, and therefore an essentially
 random merge order. Throughout training, tokens which are often part of the same merges will have
 their representations drawn together, so the representation reflects what they are likely to compose
 with. The model can then leverage any regularities to better perform reconstruction. This leads the
 representations to further reflect likely compositions and consequently increases the regularity in the
 structure. Ultimately, this leads to representations which must, by virtue of the training procedure,
 reflect the compositional semantics learned by the model.

For a more formal description of the operation of the model we begin by noting that the model generates *two* sets of embeddings. One set going up from leaves to root, and another coming back down from root to leaves. We denote these  $\bar{e}$  and  $\underline{e}$  respectively We also note that an embedding is typically viewed as  $e \in \mathbb{R}^{U \times K}$  with K independent channels—of particular relevance to the composition and decomposition functions which act independently over the channels. Tokens are denoted as the vertices  $w_i \in \Delta^V$  in a V-simplex for vocabulary size V. All together, the functioning of the model is then characterised by:

112

113

125

$$\Omega_{\Psi}(w_i) = w_i \,\Psi, \quad \Psi \in \mathbb{R}^{V \times (U \star K)} \tag{1}$$

$$C_{\Phi}(\bar{e}_i, \bar{e}_{i+1}) = \operatorname{HCAT}(\bar{e}_i, \bar{e}_{i+1}) \Phi + \phi \quad \Phi \in \mathbb{R}^{2U \times U}, \phi \in \mathbb{R}^U$$
(2)

$$\Lambda_{\Gamma}(\underline{e}_i) = \underline{e}_i \, \Gamma \quad \Gamma \in \mathbb{R}^{(U \star K) \times V} \tag{4}$$

....

Given the nature of the model, a straightforward objective would be to simply reconstruct the tokens, formulated for sentence w and prediction  $\hat{\mathbf{w}}$  as  $\mathcal{L}_{CE}(\mathbf{w}, \hat{\mathbf{w}}) = -\frac{1}{N} \sum_{n=1}^{N} w_n \cdot \log \hat{w}_n$ . An alternate approach developed by Opper et al. (2023) leverages the multi-level structure of the model to define a contrastive objective over a batch of sentences  $\{\mathbf{w}_b\}_{b=1}^B$  with a total of M nodes (internal + leaves). Noting that the up and down trees share the same underlying structure (modulo reversed edges), this objective draws together corresponding up and down embeddings at a given tree position, whilst pushing away other embeddings across the batch, using the cosine similarity

 $D_{\Theta}(\underline{e}_i) = \text{HSPLIT}(\underline{e}_i \Theta + \theta) \quad \Theta \in \mathbb{R}^{U \times 2U}, \theta \in \mathbb{R}^{2U}$ 



Figure 2: Entangled trees: Example of disjoint trees being transformed into an entangled tree. We leave out denoting functions from Eqs. (1) to (4) to avoid clutter and assume them implicitly present.

metric. Denoting the pairwise similarity matrix  $A \in \mathbb{R}^{M \times M}$  between upward embeddings  $\langle \bar{e}_i \rangle_{i=1}^M$ and downward embeddings  $\langle \underline{e}_i \rangle_{i=1}^M$ , and  $A_{i\bullet}, A_{\bullet j}, A_{ij}$  the *i*<sup>th</sup> row, *j*<sup>th</sup> column, and (i, j)<sup>th</sup> entry of the matrix respectively, the objective is defined as:  $\mathcal{L}_{CO}(\bar{e}, \underline{e}) = \frac{-1}{2M} \sum_{i=1}^{M} \log (\sigma_{\tau}(A_{i\bullet}) \sigma_{\tau}(A_{\bullet i}))$  with tempered softmax  $\sigma_{\tau}(\cdot)$  (temperature  $\tau$ ) normalising over the unspecified (•) dimension.

#### 4 MODEL

186 A particularly interesting characteristic of learning with explicitly structured models such as Self-StrAE or even earlier models such as the IORNN (Le & Zuidema, 2014) is the dichotomy 187 between the upward and downward embeddings. Given their construction, the upward embeddings 188 are always *locally-contextual*: they only encapsulate the context of the span they cover. For ex-189 ample, following Fig. 1, the upward embedding  $\bar{e}$  for the span (ate doughnuts) is always the same 190 regardless of context, no matter who did the eating. In contrast, downward embeddings are always 191 globally-contextual: they must encapsulate the surrounding context by virtue of being decomposed 192 from larger spans. For our example, this implies that there are multiple downward embeddings  $e^y$  for 193 the given span, one for each  $y \in \{$ Lisa, Homer, ...  $\}$ . To learn effective embeddings then, one must 194 marginalise over these different downward embeddings to ensure that their meaning resolves over all these contexts.

196 197

166

167

169

174

175

176 177

178

179

181 182 183

185

#### 4.1 FROM TREES TO ENTANGLED TREES

199 We want to have the composition embeddings amortise 200 over all possible contexts, and simultaneously we want all 201 decompositions embeddings to resolve to the same thing. 202 The representation of an entity Lisa should encapuslate 203 everything she could possibly eat. Simulteanously if we 204 take the average of everything she could eat we should get 205 back to Lisa. Self-StrAE does not explicitly model this 206 behaviour in its structure. Decomposition embeddings of 207 the same entity only interact when we calculate the loss. On top of this, because the loss is taken over the batch, 208 they are actually treated as false negatives to each other. 209 Even though they are terms that ought not be pushed 210 away, the objective ask them to be. 211



#### Algorithm 1 Banyan: Entangled Compose **Input:** Global frontier $\langle (s_n, e_n) \rangle_{n=1}^N$ , compose (o), concat ( $\diamond$ ), similarity CSIM(e, e')1: $\mathcal{A} \leftarrow \langle (s_n, e_n) \rangle_{n=1}^N$ initialise frontier 2: $(\mathcal{V}, \mathcal{E}) \leftarrow (\emptyset, \emptyset)$ ▷ initialise graph 3: while $\exists i : s_i \diamond s_{i+1} \not\in_s \mathcal{V}$ do $i^{\star} \leftarrow \arg\max_i \mathsf{CSIM}(e_i, e_{i+1})$ 4: Iocation of closest adjacent pair 5: $e_p = \circ(e_{i^{\star}}, e_{i^{\star}+1}) \quad \triangleright \text{ compute composition}$ $\begin{array}{l} \mathcal{V} \leftarrow \mathcal{V} \cup \{(s_{i^{\star}} \diamond s_{i^{\star}+1}, e_{p})\} \\ \mathcal{E} \leftarrow \mathcal{E} \cup \{p \sim i^{\star}, p \sim (i^{\star}+1)\} \end{array}$ 6: 7: 8: $\mathcal{J} \leftarrow \{j : (s_j, s_{j+1}) = (s_{i^*}, s_{i^*+1})\}$ Iocations of all occurrences of this pair $\mathcal{A} \leftarrow \mathcal{A} \setminus \{ \forall_{j \in \mathcal{J}} \ \mathcal{A}_j, \mathcal{A}_{j+1} \}$ 9: delete occurrences from those locations 10: $\mathcal{A} \leftarrow \mathcal{A} \cup_{\mathcal{J}} \{ (s_{i^*} \diamond s_{i^*+1}, e_p) \}$ > insert composition into those locations 11: **return** Graph $(\mathcal{V}, \mathcal{E})$

trees—where entangling refers to reduction of a set of disjoint tree structures into a single conjoined 214 graph structure. An example is shown in Fig. 2 with disjoint trees on the left and resulting entangled 215 tree on the right. Here, all instances of (night) and (some are born to) are captured by a single node

representing that constituent. We call our model Banyan on account of this entangling, because, like
 the tree, it can have many roots -consisting of nodes frequently reused across contexts.

Entangling: Constructing an entangled tree given a set of disjoint trees is a relatively straightforward process and is formally specified in Algorithm 1. In contrast to the agglomerative clustering employed in Self-StrAE, here we employ a global frontier spanning all leaf nodes across the given data. The key differences to the prior methods are mainly to do with constructing a graph jointly with progressing the frontier and ensuring that new nodes are never duplicated, for which we employ a node identity  $s_n$ in addition to the node embedding  $e_n$ .

Incorporating context Following the entangling of 225 trees described, the model proceeds in a similar vein 226 to Self-StrAE, by composing upwards from leaves 227 to roots (multiple roots corresponding to multiple 228 trees), and then decomposing downwards back to the 229 leaves. With entangled trees, while traversing up-230 wards each node is always composed from the same 231 two children, but on the way back down, things are 232 different as each separate context for a given node provides a different downward embedding. This is 233 shown in Fig. 3 focussing on a subgraph of the en-234



Figure 3: Upward and downward traversals for a section of the entangled tree from Fig. 2.

tangled tree from Fig. 2(right). Note that the node in question (in blue) corresponds to the span some are born to>, and has downward embeddings that incorporate context both from (endless night) and (sweet delight). This is exactly as desired, as Banyan allows explicit aggregation to derive the downward embedding that resolves over the contexts. For any upward embedding  $\bar{e}$  whose span occurs in different contexts  $y \in \mathcal{Y}$ , the corresponding downward embedding is derived by simply averaging over the different contextual down embeddings; i.e.,  $\underline{e} = 1/|\mathcal{Y}| \sum_{y} \underline{e}^{y}$ .

240 241 242

243

244

245

246

247

**Effectiveness and efficiency** Beyond the ability to explicitly incorporate context across data, entangled trees also help the contrastive objective by avoiding false negatives since they do not admit duplicate nodes by construction. Furthermore, the lack of duplicate nodes also drastically impacts the memory footprint of the model as one deals with the *set* of all nodes rather than counting each instance as its own node. These effects becomes more pronounced when entangling a larger set of instances as the likelihood of false negative and duplicates goes up together.

248 Practical estimation Given the advantages conferred by entangled trees, one would ideally want to 249 construct it over all the available data. This however is not practically feasible as the size of data 250 typically grows exponentially with time. To address this, we construct our model to estimate the 251 given objective by taking steps over *batches* of data that are of a more manageable size, noting that 252 this estimator is unbiased. To see this is the case, note that entangled trees only affects the downward 253 embeddings directly, and that batching simply means that the resolved embedding is an average over 254 samples instead of over all the data (population)—the sample mean is always an unbiased estimator 255 of the population mean.

256 257

258

## 4.2 SIMPLIFIED MESSAGE PASSING

Complementary to the development of entangled trees to incorporate context, we also explore avenues to improve the message passing with the composition (*C*) and decomposition (*D*) functions. The original formulations of these from Eqs. (2) and (3) employ concatenation and splitting along with simple single-layer linear neural networks. The authors found that these simpler formulations led to better representations than e.g., Tree-LSTM cells, because they forced the model to conform to the compression order of the structure.

But if all we need for success is to respect the compression order, then we could possibly do better
with an even simpler solution that exploits diagonalised functions (Ba et al., 2016)? These have
become a hallmark of the recent resurgence in recurrent neural networks (Peng et al., 2023; Orvieto
et al., 2023; De et al., 2024), by introducing decayed memory in the temporal dimension. Such a
parameterisation means that rather than using full matrices as our C and D functions, we instead
define them as:

271 272 273

285

$$C(\bar{e}_i, \bar{e}_{i+1}) = (\bar{e}_i \cdot \sigma(\Phi_l) + \bar{e}_{i+1} \cdot \sigma(\Phi_r)) + \phi \quad \Phi_l, \Phi_r, \phi \in \mathbb{R}^U$$
(5)

$$D(\underline{e}_i) = \left(\underline{e}_i \cdot \sigma(\Theta_l) + \theta_l, \ \underline{e}_i \cdot \sigma(\Theta_r) + \theta_r\right) \quad \Theta_l, \Theta_r, \theta_l, \theta_r \in \mathbb{R}^U$$
(6)

with sigmoid non-linearity ( $\sigma$ ) applied to parameters both for numerical stability and to make the functions enforce a decayed memory over structure depth. The repeated application of the diagonal composition function will decay the influence of nodes further down in the tree, thereby respecting the compression order of the structure. In addition, during composition parent representations can increase in magnitude as they are the sum of the two children. During decomposition child representations will, by necessity, reduce back down in magnitude towards the core. In this way the functions further mimic the information flow specified by the entangled trees.

These relatively simple changes have a pretty drastic effect, both in terms of performance (see experiments) as well as memory footprint, with parameters now reduced by a factor of U compared to the functions from Eqs. (2) and (3)

5 EXPERIMENTS:

287288 5.1 WARMUP: ENGLISH LANGUAGE EVALUATION

Goal: Having outlined Banyan, we want to test whether it can efficiently learn semantics. We start
 by evaluating on English, as there are far more test sets available than for low resource languages.

291 **Evaluation:** We want to evaluate how well Banyan is able to learn effective semantic representations. 292 Ideally we want to probe this at different levels of hierarchy, covering both the lexical and sentential 293 level. Our evaluation is unsupervised, both to directly probe the effect of pretraining with the inductive 294 bias, and because this setting has greater parity to what may be expected in a low resource domain, 295 where there are few labelled datasets. For these reasons, we turn to a series of tasks which measure 296 correlation between cosine similarity of embedding pairs for two examples and human judgements 297 of their semantic correspondence. On the word level, we use Simlex-999 (Hill et al., 2015) and 298 WordSim-S/R (Agirre et al., 2009). All tasks measure semantics, but do so on differing axes. To 299 understand this, we must first qualify the difference between semantic similarity and relatedness. Semantic similarity measures the extent to which entities act the same way. For example, 'running' 300 and 'singing' are similar as they share the role verb. Semantic relatedness measures conceptual 301 association. For example, 'singing' and 'fame' may be highly related. Simlex measures similarity at 302 the exclusion of relatedness. Wordsim S measures similarity without penalising relatedness. And 303 Wordsim R measures relatedness. On the sentence level, we use STS-12 through 16 (Agirre et al., 304 2012; 2013; 2014; 2015; 2016), the STS-B (Cer et al., 2017), SICK-R (Marelli et al., 2014) and 305 SemRel (Ousidhoum et al., 2024). Each measures slightly different aspects of sentential semantics, 306 covering similarity, relatedness, equality and entailment. A good model should do well on all of them. 307

Baselines: We compare against the Self-StrAE, GloVe embeddings (Pennington et al., 2014) and 308 a RoBERTa (Liu et al., 2019) in the medium configuration from (Turc et al., 2019). Self-StrAE 309 stands as the closest point of comparison to Banyan. Self-StrAE indicates the performance level of 310 structured representation learning lies, as well as any improvements we are able to achieve. GloVe 311 lets us compare to traditional static embeddings. This comparison probes whether our model is 312 learning anything more than just simple bag of word features. To obtain sentence embeddings, we 313 report results using both the simple average of the word embeddings and the average with filler 314 words removed following (Reimers & Gurevych, 2019). These filler words contribute little semantic 315 information and their removal has been shown to improve performance. For RoBERTa, we report results using both the standard model, and again after enhancing RoBERTa through an extra round of 316 contrastive SimCSE training (Gao et al., 2021), as a further STS baseline. In both cases, we generate 317 sentence embeddings through mean pooling. To produce static embeddings from RoBERTa to use in 318 lexical evaluation, we follow Bommasani et al. (2020) and average the contextualised representations 319 of all occurrences of the word in the training set. The RoBERTa is intended as a stronger baseline. It 320 has significantly more parameters than Banyan and is able to model meaning in context unlike GloVe. 321

**Hyperparameters and Pre-training Details:** For all models we set the embedding size to 256. For Self-StrAE we use the configuration of (Opper et al., 2023) and set embeddings as square matrices (i.e., K=16 and U=16). For Banyan we set these values to K=128 and U=2, because the more Table 1: Sentence level results for models pretrained on English. Higher is better. Results represent the average across four random initialisations. Only columns where there is no standard deviation overlap between models are bolded. Spearman's  $\rho$  is \* 100 following convention.

207										
321	Model	STS-12	STS-13	STS-14	STS-15	STS-16	STS-B	SICK-R	SemRel	Score
328	Self-StrAE	$31.98\pm0.58$	$53.88\pm0.68$	$37.73\pm0.70$	$55.23\pm0.58$	$55.55\pm0.47$	$39.53 \pm 1.61$	$51.78\pm0.29$	$50.05\pm0.92$	$46.59\pm0.43$
329	GloVe + stopword rm	$\begin{array}{c} 31.61 \pm 0.31 \\ 39.00 \pm 0.57 \end{array}$	$\begin{array}{c} 21.69 \pm 0.12 \\ 41.61 \pm 0.19 \end{array}$	$\begin{array}{c} 27.37 \pm 0.10 \\ 39.31 \pm 0.18 \end{array}$	$\begin{array}{c} 40.42 \pm 0.09 \\ 51.06 \pm 0.35 \end{array}$	$\begin{array}{c} 29.27 \pm 0.12 \\ 45.14 \pm 0.14 \end{array}$	$\begin{array}{c} 28.25 \pm 0.08 \\ 48.40 \pm 0.07 \end{array}$	$\begin{array}{c} 50.20 \pm 0.25 \\ 52.80 \pm 0.04 \end{array}$	$\begin{array}{c} 41.20 \pm 0.43 \\ 42.37 \pm 0.13 \end{array}$	$\begin{array}{c} 33.75 \pm 0.04 \\ 44.96 \pm 0.10 \end{array}$
330	RoBERTa + SimCSE	$\begin{array}{c} 42.77 \pm 1.27 \\ 50.63 \pm 1.45 \end{array}$	$\begin{array}{c} 51.70 \pm 1.30 \\ 62.23 \pm 2.51 \end{array}$	$\begin{array}{c} 45.67 \pm 1.42 \\ 54.17 \pm 2.10 \end{array}$	$\begin{array}{c} 63.97 \pm 0.81 \\ 68.77 \pm 3.00 \end{array}$	$\begin{array}{c} 59.60 \pm 0.61 \\ 66.67 \pm 1.40 \end{array}$	$\begin{array}{c} 39.97 \pm 0.95 \\ 53.53 \pm 1.18 \end{array}$	$\begin{array}{c} 52.93\pm0.23\\ \textbf{56.87}\pm\textbf{1.16}\end{array}$	$\begin{array}{c} 52.73 \pm 0.58 \\ 59.27 \pm 0.93 \end{array}$	$\begin{array}{c} 51.08 \pm 0.61 \\ 59.02 \pm 1.45 \end{array}$
551	Banyan	$51.20\pm0.007$	$\textbf{69.10} \pm \textbf{0.002}$	$\textbf{63.20} \pm \textbf{0.004}$	$\textbf{73.20} \pm \textbf{0.002}$	$66.60\pm0.002$	$61.50 \pm 0.002$	$55.50\pm0.003$	$\textbf{61.60} \pm \textbf{0.002}$	$\textbf{62.70} \pm \textbf{0.001}$
332										

333 independent channels we allowed the better the model seemed to perform. We refer the reader to the 334 Appendix for ablations. We also note that because we can perform this reduction in channel size, the 335 number of non-embedding parameters for Banyan drops to just 14, as these are directly proportional to 336 U. We trained Self-StrAE and Banyan for 15 epochs (circa 15k steps and sufficient for convergence) 337 using the Adam optimizer (Kingma & Ba, 2015), with a learning rate of 1e-3 for Banyan and 1e-4 for Self-StrAE using a batch size of 512. We applied dropout of 0.2 on the embeddings and 0.1 338 on the composition and decomposition function outputs. The temperature hyper-parameter for the 339 Self-StrAE was set to 0.2. To process the graphs we used DGL (Wang et al., 2020). The GloVe 340 baseline was trained for 15 epochs with a learning rate of 1e-3, and a window size of 10. We used the 341 official C++ implementation. RoBERTa medium was trained for 200,000 steps, (10% of which were 342 used for warmup). We used a learning rate of 5e-5, and a linear schedule. Positional embeddings 343 are relative key-query. The configuration for RoBERTa medium is 8 layers, 8 attention heads and 344 2048 dimensional feedforward layers. We used the Transformers library to implement and train the 345 model (Wolf et al., 2020). For SimCSE training, we used the default parameters and the official 346 implementation for unsupervised RoBERTa training from Gao et al. (2021). As our pre-training corpus 347 we selected the WikiText-103 benchmark dataset (Merity et al., 2016). The RoBERTa and GloVe baselines are trained on the full corpus (103 million tokens), representing the upper-middle end of 348 the level of data scale that might be available for a language. Whereas we trained Self-StrAE and 349 Banyan on a uniform subsample of 10 million tokens, representing the lower end of how many tokens 350 might be available, because these explicit structure models are supposed to be efficient learners. 351

Results: Results are shown in Tables 1 and 2.
On both the word level and sentence level
Banyan does much better than Self-StrAE. We
ablate the reasons for this in more detail later in
the manuscript. Both models suffer on SimLex
because they need to model both similarity and
relatedness as the latter dictates merge (related

**Results:** Results are shown in Tables 1 and 2. Table 2: Word level results analogous to Table 1.

			•	
Model	Simlex	Wordsim-S	Wordsim-R	Score
Self-StrAE	$13.80\pm0.41$	$54.38\pm0.78$	$52.85 \pm 1.27$	$40.34\pm0.66$
GloVe	$27.47\pm0.25$	$62.53\pm0.42$	$51.00\pm0.56$	$47.00\pm0.38$
RoBERTa	$\textbf{29.23} \pm \textbf{0.64}$	$61.97 \pm 2.38$	$46.00\pm2.13$	$45.73 \pm 1.71$
Banyan	$16.57\pm0.02$	$\textbf{63.25} \pm \textbf{0.03}$	$\textbf{69.00} \pm \textbf{0.01}$	$\textbf{49.61} \pm \textbf{0.02}$

concepts often compose together). However, the important thing to note is that the structured models 359 effectively transfer the same performance from the word level to the sentence level. They can take 360 advantage of composition, and transfer the meaning of the parts to understanding the meaning of the 361 whole. The GloVe baseline is good on the word level, but does not generalise to the sentence level as 362 well as the transformer, even when we give it stopword removal. It cannot transfer semantic knowl-363 edge seamlessly to different levels of complexity. Banyan can, and is able to achieve comparable or 364 better performance than the SimCSE RoBERTa despite being much smaller and exposed to 10x less pre-training data. This means we have a structured model that remains efficient and cheap, and also effective at representation learning. 366

367 368

369

## 5.2 MULTILINGUAL EVALUATION:

Goal: From the results in English we know that Banyan is an efficient learner: it can produce good
 representations without requiring large-scale data or compute. This implies potential use for under
 represented communities, whose languages are not well covered by current NLP approaches. Now
 we have to test that.

Evaluation: Learning semantic representations for low resource languages remains an ongoing
 challenge in NLP. A core problem is not just the lack of training data, but also the lack of evaluation
 datasets. Recent work by Ousidhoum et al. (2024) has sought to address this issue, providing
 semantic relatedness test sets for several low resource Asian and African languages. These test sets
 are evaluated the same as before, comparing the cosine similarity between model embeddings for

Table 3: Multilingual Results. Banyan performance is taken over four random seeds. Baselines
marked with † have been finetuned on supervised semantic similarity datasets. FT denotes unsupervised finetuning using masked language modelling on the same corpora as Banyan.

381	Model	Indonesian	Arabic	Telugu	Marathi	Mor. Arabic	Kinyarwanda	Hausa	Afrikaans	Spanish	Amharic	Hindi	Score
382	XLM-R Llama-3.1 (8B) Mistral Nemo	46.7 53.4 50.7	31.6 31.1 20.1	46.3 65.6 57	55.7 63.4 52.3	17.4 19.4 15.1	13.2 19.7 16.3	4.1 6.1 1.8	56.2 65.4 58.3	68.9 66.7 66.2	57.3 64.1 53.2	52.7 61.7 55.8	40.92 46.96 40.62
384	MiniLM-L12† Paraphrase XLM-R†	39 46.1	16.1 61	34.8 58.1	39.5 79.6	13.5 7.1	35 43.2	32.7 22.5	74.1 76.8	58.8 71.7	9.6 64.6	43.8 52	36.08 52.97
	XLM-R (FT)	47.9	33.6	68.8	75.1	21.6	19.4	14.6	72.6	72.8	59.6	57.6	49.41
385	Banyan	$44.17\pm1.11$	$43.20\pm1.82$	$\textbf{71.13} \pm \textbf{0.91}$	$67.67 \pm 0.64$	$\textbf{52.00} \pm \textbf{2.25}$	$\textbf{46.1} \pm \textbf{0.32}$	$\textbf{43.7} \pm \textbf{1.21}$	$\textbf{78.68} \pm \textbf{0.30}$	$60.95\pm0.76$	$\textbf{66.18} \pm \textbf{0.46}$	$61.83\pm0.6$	57.78

two sequences with human judgements of their semantic match. As before, evaluation is zero-shot
unsupervised. Allowing us to evaluate Banyan on Indonesian, Arabic, Telugu, Marathi, Moroccan
Arabic, Kinyarwanda, Hausa, Afrikaans, Spanish, Amharic and Hindi. These represent a spectrum in
terms how well resourced they are. For example, Spanish and Hindi are reasonably well represented,
while Moroccan Arabic and Kinyarwanda have extremely little training data.

**Baselines:** We select XLM-R (Conneau et al., 2019): a transformer encoder trained on 2TB of 392 multilingual data. Llama 3.1 8B (Dubey et al., 2024): a decoder only LLM trained on 15 trillion 393 tokens. Mistral Nemo 12B: a decoder only LLM designed with multi-lingual capacities in mind. In 394 addition we also compare against two specialised embedding models from the sentence transformers 395 range (Reimers & Gurevych, 2019): Mini-LM-L12-V2 and Paraphrase-XLM-R-Multilingual-V1. 396 These are pre-trained transformer encoders that have been finetuned on supervised datasets designed 397 to produce high quality semantic representations. The baselines we select here are emblematic of 398 the kind of models one might reach for in order to embed a corpus. For all models we use mean 399 pooling to produce the sentence representation following Reimers & Gurevych (2019); Li & Li 400 (2024). Finally, for parity we include an XLM-R baseline which is finetuned on the same corpora.

401 Banyan Pre-training and Hyperparameters: For Afrikaans, Spanish and Amharic we obtained 402 corpora from Leipzig Corpora Collection<sup>1</sup> (Goldhahn et al., 2012). For Amharic we utilised a 403 MiT licenced pre-training set of 1 million sequences available on the Huggingface hub at this link. 404 Kinyarwanda and Hausa data was sourced from Opus (Nygaard & Tiedemann, 2003). Each dataset 405 consists of roughly 10 million tokens. We utilise a pre-trained BPE tokenizer for each language 406 from the BPEMB Python package (Heinzerling & Strube, 2018). Though the package also provides 407 pre-trained embeddings, we solely use the tokenizer and learn embeddings from scratch. For the 408 model hyperparameters we keep all the settings from the experiments on English. For XLM-R we finetune for up to 100k steps with early stopping, using a linearly scheduled learning rate of 5e-5 with 409 10 percent of stepping serving as warmup. XLM-R runs at batch size 128 across 4xA40 45gb cards. 410

411 **Results:** See Table 3. In Spanish, a well resourced language with high coverage, the transformer 412 baselines almost all outperform Banyan. However, as languages become lower resourced the picture 413 changes, and Banyan outperforms or is comparable to the baselines. This even includes the multilin-414 gual XLM-R that has undergone supervised training to produce better representations. While finetuning XLM-R improves performance the amount of benefit it provides is not uniform and is insufficient 415 to prove viable in the very low resource cases. Banyan is able to learn competitive representations 416 consistently across languages, unsupervised and with very little data, meaning it provides a viable 417 alternative for producing embeddings cheaply and efficiently for low resource languages. 418

419

431

391

# 420 5.3 EFFICIENCY

421 Alongside its embedding matrix, Banyan has two central 422 components: the composition and decomposition func-423 tions. We diagonalise these functions so that they are both 424 easier to compute and have fewer parameters than standard 425 weight matrices,  $(2U \text{ rather than } 2U \times U)$ , achieving a fur-426 ther order of magnitude reduction in parameters compared 427 with the already minimal Self-StrAE.



428
429
429
420
420
420
420
421
421
422
423
424
425
426
427
428
429
429
429
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420
420

Figure 4: Total #nodes in entangled trees vs sentential trees as batch size grows.

<sup>&</sup>lt;sup>1</sup>For Spanish and Hindi we select the mixed corpus and uniformly subsample to reduce size to  $\approx 10M$  tokens.

Table 4: Ablations of modelling changes made for Banyan. Higher is better. Results represent the average across four random initialisations. Only columns where there is no standard deviation overlap between models are bolded. Spearman's  $\rho$  is \* 100 following convention.

195										
433	Model	STS-12	STS-13	STS-14	STS-15	STS-16	STS-B	SICK-R	SemRel	Score
436	Standard Trees	$31.98 \pm 0.58$	$53.88 \pm 0.68$	$37.73\pm0.70$	$55.23\pm0.58$	$55.55\pm0.47$	$39.53 \pm 1.61$	$51.78\pm0.29$	$50.05\pm0.92$	$46.59\pm0.43$
	+ diag functions	$35.13 \pm 0.33$	$56.05 \pm 0.24$	$40.58 \pm 0.05$	$58.83 \pm 0.10$	$56.78 \pm 0.21$	$44.10 \pm 0.14$	$53.35 \pm 0.17$	$52.65 \pm 0.17$	$49.68 \pm 0.06$
437	++ CE loss	$47.10 \pm 1.04$	$61.85\pm1.44$	$58.60 \pm 1.34$	$70.45\pm0.57$	$62.45\pm0.70$	$59.50\pm0.53$	$\textbf{59.00} \pm \textbf{0.26}$	$60.33\pm0.26$	$59.91 \pm 0.54$
438	Entangled Trees	$38.98 \pm 0.39$	$61.75\pm0.14$	$43.65\pm0.46$	$58.21 \pm 0.41$	$55.29 \pm 0.23$	$46.15\pm0.71$	$53.93 \pm 0.16$	$52.53\pm0.09$	$51.31\pm0.13$
400	+ diag functions	$44.15 \pm 0.002$	$62.80 \pm 0.002$	$48.30 \pm 0.001$	$64.60 \pm 0.002$	$60.30 \pm 0.001$	$49.80 \pm 0.002$	$55.14 \pm 0.001$	$57.70 \pm 0.001$	$55.23 \pm 0.001$
439	++ CE loss	$\textbf{51.20} \pm \textbf{0.007}$	$\textbf{69.10} \pm \textbf{0.002}$	$\textbf{63.20} \pm \textbf{0.004}$	$\textbf{73.20} \pm \textbf{0.002}$	$\textbf{66.60} \pm \textbf{0.002}$	$\textbf{61.50} \pm \textbf{0.002}$	$55.50\pm0.003$	$\textbf{61.60} \pm \textbf{0.002}$	$\textbf{62.70} \pm \textbf{0.001}$

440 441

Fig. 4). This is because the number of reused constituent nodes also grows as batch size increases, 442 and entangled trees capture the set of all constituents, which consequently does not grow as drastically. 443 In practical terms, because entangled trees requires fewer nodes, and each node requires two distinct 444 embeddings ( $\bar{e}$  and  $\underline{e}$ ) to be held for it, reducing the number of nodes required leads to radical 445 improvements in memory efficiency. Put together, these changes mean that we can train Banyan very 446 quickly as we can use large batches and its small number of parameters ensure quick convergence. 447 On a single Nvidia A40 GPU with a batch size of 1024, Banyan pretrains from scratch in under 50 448 minutes, meaning that the total cost of pretraining a Banyan model sits at around 30 cents<sup>2</sup>. Free-tier 449 Google Colab users can achieve similar results in about two hours with a smaller batch size. Inference can also be performed on CPU on typical laptops, because the model is so small. Combined with 450 its data efficiency, we believe this provides a promising alternative for low resource languages and 451 communities. 452

453

455

481

454 5.4 ABLATIONS

We have shown that Banyan is more effective than its Self-StrAE predecessor, but what is the impact
of the different modelling changes we made? To test this we ablate our results from our first set of
experiments on English (see Table 4).

459 The simplest positive impacts to see are from the introduction of the diagonalised composition and 460 decomposition functions. These are sigmoided scalar values with which we multiply embeddings. Therefore they act similarly to the fast weights of Ba et al. (2016), decaying in the influence of 461 embeddings further down in the tree on the root representation. This means that the embeddings 462 produced by the model are restricted to conform to the compression order dictated by the structure, 463 and we know from Opper et al. (2023), that the more we can enforce this constraint the better our 464 representations will end up. Secondly, such simple message passing functions bias the representa-465 tion space towards informative separability. There has to be some signal from which to perform 466 reconstruction, and all the pressure is now on the representations. 467

Switching the objective to cross entropy over the vocabulary rather than the contrastive objective used by Opper et al. (2023) also yields significant benefits. This is likely because the contrastive loss is supposed to be beneficial because it enforces a pressure for representations to be uniformly distributed in space (Wang & Isola, 2020). However, our other modelling changes already push towards this quality. While it is a shame because the contrastive loss is conceptually elegant. It is known to have problems and eventually lead to shortcut solutions (Robinson et al., 2021). Therefore having a more robust objective grounded in data, like the cross entropy over the vocabulary, is actually quite nice.

Finally, changing to entangled trees is also beneficial. The effect is more pronounced before switching to the cross entropy objective, as it removes the issue of false negatives as discussed in Section 4. However, it also is beneficial beyond this. Entangling explicitly allows the model to take advantage of shared constituency structure between complex sequences, because it combines the information from all incoming parent messages. The fact that performance improves using the cross entropy objective shows that explicitly allowing the model to take advantage of such systematicity is useful.

483	Model	Banyan	Self-StrAE	RoBERTa(M)	All-MiniLM-L12-V2	XLM-R	Llama 3.1	Mistral Nemo
484	Params	14	1072	$\approx 10 \mathrm{M}$	≈21M	$\approx 85 M$	$\approx 8B$	$\approx 12B$
485								

<sup>2</sup>Current cloud computing costs sourced from: https://www.runpod.io/pricing

# <sup>486</sup> 6 CONCLUSION, LIMITATIONS AND FUTURE WORK

We introduce Banyan, a Self-Structuring AutoEncoder. Banyan's focus on global, entangled structure and simplified message passing exploits the benefits of structured compositions inherent in language data. It is more effective and efficient than prior work from which we draw three central conclusions.

Firstly, explicitly modelling structured compositions is an effective inductive bias. Table 5 shows
the parameters for the structured models versus the baselines. The structured models are far smaller,
with tens or thousands of parameters instead of millions or billions. And nonetheless, Banyan is still
competitive across several metrics, indicating we have found an efficient learning procedure.

Secondly, we have not yet fully exploited the potential of the inductive bias. Banyan still relies on
 greedy agglomerative clustering to induce structure. This is effective, but sub-optimal. Future work
 could make the structure induction procedure parametric and learnable. The type of structure models
 are exposed to impacts the quality of learnt semantic representations (Opper et al., 2023). So if *how* we induce structure improves, the model should learn significantly better representations.

Thirdly, while this paper focuses on recursive neural networks for the purposes of efficiency and low resource applicability, the method could in principle be applied to representations from pre-trained transformer models. Firstly, the transformers attention essentially defines a soft (fully connected) graph between tokens, which could serve as a more flexible basis for constructing Banyan's discrete structures. Moreover, the entangled tree structure essentially serves as a map of the conceptual associations learned by a model, and could provide an interesting probe into the representation space of pre-trained LLMs.

Finally, good and cheap embedding models are useful for many applications. For example, the digital humanities need to organise corpora of ancient languages, making it easier for researchers to access texts they need. But these corpora are small, and these languages are unlikely to be present in pretraining corpora of larger models. Banyan provides an efficient solution for producing representations for both these use cases and low resource languages and under represented communities more generally. To conclude, Banyan addresses the problem of efficient learning in low-resource settings.

# 540 REFERENCES

- Mohamed Abdalla, Krishnapriya Vishnubhotla, and Saif Mohammad. What makes sentences semantically related? a textual relatedness dataset and empirical study. In Andreas Vlachos and Isabelle Augenstein (eds.), *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 782–796, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.eacl-main.55. URL https://aclanthology.org/2023.eacl-main.55.
- Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Paşca, and Aitor Soroa. A
   study on similarity and relatedness using distributional and WordNet-based approaches. In *North American Chapter of the Association for Computational Linguistics (NAACL)*, pp. 19–27, 2009.
- Eneko Agirre, Mona Diab, Daniel Cer, and Aitor Gonzalez-Agirre. Semeval-2012 task 6: A pilot on semantic textual similarity. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics - Volume 1: Proceedings of the Main Conference and the Shared Task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*, SemEval '12, pp. 385–393, USA, 2012. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. \*SEM 2013 shared task: Semantic textual similarity. In Mona Diab, Tim Baldwin, and Marco Baroni (eds.), Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity, pp. 32–43, Atlanta, Georgia, USA, June 2013. Association for Computational Linguistics. URL https://aclanthology.org/ S13-1004.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei
  Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. SemEval-2014 task 10: Multilingual
  semantic textual similarity. In Preslav Nakov and Torsten Zesch (eds.), *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pp. 81–91, Dublin, Ireland,
  August 2014. Association for Computational Linguistics. doi: 10.3115/v1/S14-2010. URL
  https://aclanthology.org/S14-2010.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei
  Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz Uria, and
  Janyce Wiebe. SemEval-2015 task 2: Semantic textual similarity, English, Spanish and pilot on
  interpretability. In Preslav Nakov, Torsten Zesch, Daniel Cer, and David Jurgens (eds.), *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pp. 252–263, Denver,
  Colorado, June 2015. Association for Computational Linguistics. doi: 10.18653/v1/S15-2045.
  URL https://aclanthology.org/S15-2045.
- Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Rada Mihalcea, German Rigau, and Janyce Wiebe. SemEval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pp. 497–511, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/S16-1081. URL https://aclanthology.org/S16-1081.
- Jimmy Ba, Geoffrey Hinton, Volodymyr Mnih, Joel Z. Leibo, and Catalin Ionescu. Using fast weights to attend to the recent past, 2016. URL https://arxiv.org/abs/1610.06258.
- Rishi Bommasani, Kelly Davis, and Claire Cardie. Interpreting Pretrained Contextualized Representations via Reductions to Static Embeddings. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4758–4781, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.431. URL https://aclanthology.org/2020.acl-main.431.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. SemEval-2017 task
   1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Workshop on Semantic Evaluation (SemEval)*, pp. 1–14, 2017.
  - N. Chomsky. Three models for the description of language. *IRE Transactions on Information Theory*, 2(3):113–124, 1956. doi: 10.1109/TIT.1956.1056813.

- Jishnu Ray Chowdhury and Cornelia Caragea. Modeling hierarchical structures with continuous recursive neural networks. *CoRR*, abs/2106.06038, 2021. URL https://arxiv.org/abs/2106. 06038.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek,
   Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un supervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116, 2019. URL
   http://arxiv.org/abs/1911.02116.
- Stephen Crain and Mineharu Nakayama. Structure dependence in grammar formation. *Language*, 63 (3):522–543, 1987.
- Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Desjardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and Caglar Gulcehre. Griffin: Mixing gated linear recurrences with local attention for efficient language models, 2024.
- Marie-Catherine de Marneffe, Bill MacCartney, and Christopher D. Manning. Generating typed dependency parses from phrase structure parses. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*, Genoa, Italy, May 2006. European Language
  Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2006/
  pdf/440\_pdf.pdf.
- Andrew Drozdov, Patrick Verga, Mohit Yadav, Mohit Iyyer, and Andrew McCallum. Unsupervised
   latent tree induction with deep inside-outside recursive autoencoders. *CoRR*, abs/1904.02142,
   URL http://arxiv.org/abs/1904.02142.
- Andrew Drozdov, Subendhu Rongali, Yi-Pei Chen, Tim O'Gorman, Mohit Iyyer, and Andrew Mc-Callum. Unsupervised parsing with S-DIORA: Single tree encoding for deep inside-outside recursive autoencoders. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4832–4845, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.392.
   //aclanthology.org/2020.emnlp-main.392.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 625 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, 626 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston 627 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, 628 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris 629 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton 630 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David 631 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 632 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip 633 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, 634 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, 635 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, 636 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu 637 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph 638 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, 639 Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz 640 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence 641 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas 642 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, 643 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, 644 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan 645 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, 646 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, 647 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit

Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, 649 Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia 650 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, 651 Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, 652 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, 653 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 654 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, 655 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, 656 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen 657 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe 658 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya 659 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex 660 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei 661 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew 662 Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, 664 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt 665 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao 666 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon 667 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide 668 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, 669 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily 670 Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix 671 Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank 672 Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, 673 Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid 674 Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-675 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste 676 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, 677 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, 678 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik 679 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly 680 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, 681 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, 682 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria 683 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, 684 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle 685 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, 687 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia 688 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro 689 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, 690 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, 691 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan 692 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara 693 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh 694 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 696 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, 699 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, 700 Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang,

Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang,
Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait,
Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The Ilama 3 herd
of models, 2024. URL https://arxiv.org/abs/2407.21783.

- Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical analysis.
   *Cognition*, 28(1-2):3–71, 1988.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6894–6910, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL https: //aclanthology.org/2021.emnlp-main.552.
- Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. Building large monolingual dictionaries at the Leipzig corpora collection: From 100 to 200 languages. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pp. 759–765, Istanbul, Turkey, May 2012. European Language Resources Association (ELRA). URL http: //www.lrec-conf.org/proceedings/lrec2012/pdf/327\_Paper.pdf.
- 723 Zellig S. Harris. Distributional structure. *Word*, 10:146–162, 1954.

715

724

753

- Serhii Havrylov, Germán Kruszewski, and Armand Joulin. Cooperative learning of disjoint syntax and semantics. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *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)*, pp. 1118–1128, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1115.
  URL https://aclanthology.org/N19-1115.
- Benjamin Heinzerling and Michael Strube. BPEmb: Tokenization-free Pre-trained Subword Embeddings in 275 Languages. In Nicoletta Calzolari (Conference chair), Khalid Choukri, Christopher
  Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga (eds.), *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*(*LREC 2018*), Miyazaki, Japan, May 7-12, 2018 2018. European Language Resources Association
  (ELRA). ISBN 979-10-95546-00-9.
- Felix Hill, Roi Reichart, and Anna Korhonen. SimLex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695, 2015.

Xiang Hu, Haitao Mi, Zujie Wen, Yafang Wang, Yi Su, Jing Zheng, and Gerard de Melo. R2D2:
Recursive transformer based on differentiable tree for interpretable hierarchical language modeling.
In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 4897–4908, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.379. URL https://aclanthology.org/2021.acl-long.379.

- Xiang Hu, Haitao Mi, Liang Li, and Gerard de Melo. Fast-R2D2: A pretrained recursive neural network based on pruned CKY for grammar induction and text representation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2809–2821, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main. 181. URL https://aclanthology.org/2022.emnlp-main.181.
- Takuya Ito, Tim Klinger, Douglas H. Schultz, John D. Murray, Michael W. Cole, and Mattia Rigotti.
   Compositional generalization through abstract representations in human and artificial neural networks, 2022. URL https://arxiv.org/abs/2209.07431.

756 757 758	Yangfeng Ji and Jacob Eisenstein. One vector is not enough: Entity-augmented distributed semantics for discourse relations. <i>Transactions of the Association for Computational Linguistics</i> , 3:329–344, 2015. doi: 10.1162/tacl_a_00142. URL https://aclanthology.org/Q15-1024.
759 760 761	Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International Conference on Learning Representations (ICLR), San Diega, CA, USA, 2015.
762 763 764	Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. Building machines that learn and think like people. <i>CoRR</i> , abs/1604.00289, 2016. URL http://arxiv.org/abs/1604.00289.
765 766 767 768	Phong Le and Willem Zuidema. Inside-outside semantics: A framework for neural models of semantic composition. In <i>NIPS 2014 Workshop on Deep Learning and Representation Learning</i> , 2014.
769 770	Omer Levy and Yoav Goldberg. Dependency-based word embeddings. In Association for Computa- tional Linguistics (ACL)(Volume 2: Short Papers), pp. 302–308, 2014.
771 772 773 774	Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. <i>CoRR</i> , abs/2005.11401, 2020. URL https://arxiv.org/abs/2005.11401.
775 776 777	Xianming Li and Jing Li. Angle-optimized text embeddings, 2024. URL https://arxiv.org/abs/ 2309.12871.
778 779 780	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. <i>CoRR</i> , abs/1907.11692, 2019.
781 782 783 784 785	Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. A SICK cure for the evaluation of compositional distributional semantic models. In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pp. 216–223, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2014/pdf/363_Paper.pdf.
786 787 788	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. <i>CoRR</i> , abs/1609.07843, 2016. URL http://arxiv.org/abs/1609.07843.
789 790 791	Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word repre- sentations in vector space. In <i>International Conference on Learning Representations (ICLR)</i> , 2013.
792 793 794 795	Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and Christopher D. Manning. Pushdown layers: Encoding recursive structure in transformer language models, 2023. URL https://arxiv.org/ abs/2310.19089.
796 797	Lars Nygaard and Jörg Tiedemann. Opus—an open source parallel corpus. In Proceedings of the 13th Nordic Conference on Computational Linguistics, 2003.
798 799 800 801 802	Mattia Opper and Siddharth Narayanaswamy. Self-strae at semeval-2024 task 1: Making self- structuring autoencoders learn more with less. In <i>Proceedings of the 18th International Workshop</i> <i>on Semantic Evaluation (SemEval-2024)</i> , pp. 108–115. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.semeval-1.18. URL http://dx.doi.org/10.18653/v1/2024. semeval-1.18.
803 804 805 806 807 808	Mattia Opper, Victor Prokhorov, and Siddharth N. Strae: Autoencoding for pre-trained embeddings using explicit structure. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural</i> <i>Language Processing</i> , pp. 7544–7560. Association for Computational Linguistics, 2023. doi: 10. 18653/v1/2023.emnlp-main.469. URL http://dx.doi.org/10.18653/v1/2023.emnlp-main. 469.
809	Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences, 2023.

841

842

847

848

810	Nedima Ousidhoum, Shamsuddeen Hassan Muhammad, Mohamed Abdalla, Idris Abdulmumin,
811	Ibrahim Said Ahmad, Sanchit Ahuja, Alham Fikri Aji, Vladimir Araujo, Meriem Beloucif, Chris-
812	tine De Kock, Oumaima Hourrane, Manish Shrivastava, Thamar Solorio, Nirmal Surange, Kr-
813	ishnapriya Vishnubhotla, Seid Muhie Yimam, and Saif M. Mohammad. SemEval-2024 task 1:
814	Semantic textual relatedness for african and asian languages. In Proceedings of the 18th In-
815	ternational Workshop on Semantic Evaluation (SemEval-2024). Association for Computational
816	Linguistics, 2024.

- 817 Christophe Pallier, Anne-Dominique Devauchelle, and Stanislas Dehaene. Cortical representation of 818 the constituent structure of sentences. Proceedings of the National Academy of Sciences, 108(6): 819 2522–2527, 2011. doi: 10.1073/pnas.1018711108. 820
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin 821 Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemysław 822 Kazienko, Jan Kocon, Jiaming Kong, Bartlomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, 823 Ferdinand Mom, Atsushi Saito, Xiangru Tang, Bolun Wang, Johan S. Wind, Stansilaw Wozniak, 824 Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. Rwkv: 825 Reinventing rnns for the transformer era, 2023. 826
- 827 Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word 828 representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543, 829 2014.
- 830 Jishnu Ray Chowdhury and Cornelia Caragea. Beam tree recursive cells. In Andreas Krause, 831 Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett 832 (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of 833 Proceedings of Machine Learning Research, pp. 28768–28791. PMLR, 23–29 Jul 2023. URL 834 https://proceedings.mlr.press/v202/ray-chowdhury23a.html. 835
- Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-836 networks. In Empirical Methods in Natural Language Processing and International Joint Confer-837 ence on Natural Language Processing (EMNLP-IJCNLP), pp. 3982–3992, 2019. 838
- Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, and Suvrit Sra. Can 840 contrastive learning avoid shortcut solutions? CoRR, abs/2106.11230, 2021. URL https: //arxiv.org/abs/2106.11230.
- Laurent Sartran, Samuel Barrett, Adhiguna Kuncoro, Miloš Stanojević, Phil Blunsom, and Chris 843 Dyer. Transformer grammars: Augmenting transformer language models with syntactic inductive 844 biases at scale. Transactions of the Association for Computational Linguistics, 10:1423–1439, 845 2022. doi: 10.1162/tacl a 00526. URL https://aclanthology.org/2022.tacl-1.81. 846
  - Richard Socher, Jeffrey Pennington, Eric H. Huang, Andrew Y. Ng, and Christopher D. Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In Empirical Methods in Natural Language Processing (EMNLP), pp. 151–161, 2011.
- 850 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and 851 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. 852 In Empirical Methods in Natural Language Processing (EMNLP), pp. 1631–1642, 2013. 853
- 854 Paul Soulos, Henry Conklin, Mattia Opper, Paul Smolensky, Jianfeng Gao, and Roland Fernandez. 855 Compositional generalization across distributional shifts with sparse tree operations. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024. URL https: 856 //openreview.net/forum?id=f0Qunr2E0T. 857
- 858 Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from 859 tree-structured long short-term memory networks. In Association for Computational Linguistics 860 (ACL), pp. 1556–1566, 2015. 861
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: 862 The impact of student initialization on knowledge distillation. CoRR, abs/1908.08962, 2019. URL 863 http://arxiv.org/abs/1908.08962.

- 864 Shikhar Vashishth, Manik Bhandari, Prateek Yadav, Piyush Rai, Chiranjib Bhattacharyya, and Partha 865 Talukdar. Incorporating syntactic and semantic information in word embeddings using graph 866 convolutional networks. In Association for Computational Linguistics (ACL), pp. 3308–3318, 867 2019.
- 868 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information 870 Processing Systems (NeurIPS), pp. 5998–6008, 2017. 871
- 872 Bin Wang, C.-C. Jay Kuo, and Haizhou Li. Just rank: Rethinking evaluation with word and sentence 873 similarities. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long 874 Papers), pp. 6060-6077, Dublin, Ireland, May 2022. Association for Computational Linguistics. 875 doi: 10.18653/v1/2022.acl-long.419. URL https://aclanthology.org/2022.acl-long.419. 876
- 877 Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma, 878 Lingfan Yu, Yu Gai, Tianjun Xiao, Tong He, George Karypis, Jinyang Li, and Zheng Zhang. Deep 879 graph library: A graph-centric, highly-performant package for graph neural networks, 2020.
- Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through 881 alignment and uniformity on the hypersphere. CoRR, abs/2005.10242, 2020. URL https: 882 //arxiv.org/abs/2005.10242. 883
- Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Compositional 885 generalization from first principles, 2023. URL https://arxiv.org/abs/2307.05596.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, 887 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, 889 Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural 890 language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural 891 Language Processing: System Demonstrations, pp. 38-45, Online, October 2020. Association for 892 Computational Linguistics. URL https://www.aclweb.org/anthology/2020.emnlp-demos. 893 6.
- 894 895

900

901

880

#### APPENDIX Α

898 A.1 THE k AND u BALANCE

899 The change to diagonal composition functions allows us to reduce the number of total parameters while maintaining performance. This is because the number of parameters is directly proportional to channel size *u*. We show ablations for this finding in Table 6. Our findings are similar to those of 902 Opper & Narayanaswamy (2024) the smaller the channel size the better the model performs, although 903 in our case we keep things stable between seeds whereas for them when they simplified they faced 904 issues with extreme instability during training. This is thanks to the new message passing functions. 905

906 Table 6: Performance Depending on k and u values using new functions. Scores are the average of four random seeds. 907

k	и	Lex Score	STS Score
4	64	$42.9\pm0.01$	$43.5\pm0.04$
8	32	$43.2\pm0.02$	$48.6\pm0.01$
16	16	$47.02\pm0.03$	$62.2\pm0.01$
32	8	$49.2\pm0.01$	$62.9\pm0.01$
64	4	$48.7\pm0.01$	$62.9\pm0.01$
128	2	$49.61\pm0.02$	$62.7\pm0.001$
256	1	$48.7\pm0.01$	$62.9\pm0.001$