Geometry-Aware Supertagging with Heterogeneous Dynamic Convolutions

Anonymous ACL submission

Abstract

The syntactic categories of categorial grammar formalisms are structured units made of smaller, indivisible primitives, bound together by the underlying grammar's category formation rules. In the trending approach of constructive supertagging, neural models are increasingly made aware of the internal category structure, which in turn enables them to more reliably predict rare and out-of-vocabulary categories, with significant implications for grammars previously deemed too complex to find practical use. In this work, we revisit constructive supertagging from a graph-theoretic perspective, and propose a framework based on heterogeneous dynamic graph convolutions, 016 aimed at exploiting the distinctive structure 017 of a supertagger's output space. We test our approach on a number of categorial grammar datasets spanning different languages and grammar formalisms, achieving substantial improvements over previous state of the art scores.

1 Introduction

Their close affinity to logic and lambda calculus has made categorial grammars a standard tool of the trade for the formally-inclined NLP practitioner. Modern flavors of categorial grammar, despite their (sometimes striking) divergences, share a common architecture. At its core, a categorial grammar is a formal system consisting of two parts. First, there is a *lexicon*, a mapping that assigns to each word a set of *categories*. Categories are quasi-logical formulas recursively built out of atomic categories by means of category forming operations. The in-034 ventory of category forming operations at the minimum has the ability to express linguistic functionargument structure. If so desired, the inventory can be extended with extra operations, e.g. to handle syntactic phenomena beyond simple concatenation, or to express additional layers of grammatical infor-040 mation. The second component of the grammar is

a small set of *inference rules*, formulated in terms of the category forming operations. The inference rules dictate how categories interact and, through this interaction, how words combine to form larger phrases. Parsing thus becomes a process of deduction comparable (or equatable, depending on the grammar's formal rigor) to program synthesis, providing a clean and elegant syntax-semantics interface. 042

043

044

045

046

047

051

052

056

057

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

077

078

079

In the post-neural era, these two components allow differentiable implementations. The fixed lexicon is replaced by *supertagging*, a process that contextually decides on the most appropriate supertags (i.e. categories), whereas the choice of which rules of inference to apply is usually deferred to a parser further down the processing pipeline. The highly lexicalized nature of categorial grammars thus shifts the bulk of the weight of a parse to the supertagging component, as its assignments and their internal make-up inform and guide the parser's decisions.

In this work, we revisit supertagging from a geometric angle. We first note that the supertagger's output space consists of a sequence of trees, which has as of yet found no explicit representational treatment. Capitalizing on this insight, we employ a framework based on heterogeneous dynamic graph convolutions, and show that such an approach can yield substantial improvements in predictive accuracy across categories both frequently and rarely encountered during a supertagger's training phase.

2 Background

The supertagging problem revolves around the design and training of a function tasked with mapping a sequence of words $\{w_1, \ldots, w_n\}$ to a sequence of categories $\{c_1, \ldots, c_n\}$. Existing supertagging architectures differ in how they implement this mapping, with each implementation choice boiling down to (i) which of the temporal and structural dependencies within and between the input and out-

152

153

154

155

157

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

133

134

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

103

087 088

090

084

put are taken into consideration, and (ii) how these dependencies are materialized.

Earlier work would utilize solely occurrence counts from a training corpus to independently map word n-grams to their most likely categories, and then attempt to filter out implausible sequences via rule-constrained probabilistic models (Bangalore and Joshi, 1999). The shift from sparse feature vectors to distributed word representations facilitated integration with neural networks and improved generalization on the mapping domain, extending it to rare and previously unseen words (Lewis and Steedman, 2014). Later, the advent of recurrent neural networks offered a natural means of incorporating temporal structure, widening the input receptive field through contextualized word representations on the one hand (Xu et al., 2015), but also permitting an auto-regressive formulation of the output generation, whereby the effect of a category assignment could percolate through the remainder of the output sequence (Vaswani et al., 2016).

More recently, the so-called constructive paradigm seeks to explore the previously ignored structure "internal" to categories. By inspecting their formation rules, Kogkalidis et al. (2019) equates categories to CFG derivations, and views a category sequence as the concatenation of their flattened depth-first projections. The goal sequence is now incrementally generated on a symbol-bysymbol basis using a transformer-based seq2seq model; a twist which provides the decoder with the means to construct novel categories on demand, bolstering co-domain generalization. The decoder's global receptive field, however, comes at the heavy price of quadratic memory complexity, which also bodes poorly with the elongated output sequences, leading to a slowed down inference speed. Expanding on the idea, Prange et al. (2021) explicates the categories' tree structure, embedding symbols based on their tree positions and propagating contextualized representations through tree edges, using either residual dense connections or a tree-structured GRU. This adaptation completely eliminates the burden of learning how trees are constructed, instead allowing the model to focus on what trees to construct, leading to drastically improved performance. Simultaneously, since the decoder is now token-separable, it permits construction of categories for the entire sentence in parallel, speeding up inference and reducing the network's memory footprint. In the process, however, it loses the ability to model interactions between auto-regressed nodes belonging to different trees, morally reducing the task once more to sequence classification (albeit now with a dynamic classifier).

Despite their common goal of accounting for syntactic categories in the zipfian tail, there are tension points between the above two approaches. In providing a global history context, the first breaks the input-to-output alignment and hides the categorial tree structure. In opting for a tree-wise bottomup decoding, the second forgets about meaningful inter-tree output-to-output dependencies. In this paper, we seek to resolve these tension points with a novel, unified and grammar-agnostic supertagging framework based on heterogeneous dynamic graph convolutions. Our architecture combines the merits of explicit tree structures, strong autoregressive properties, near-constant decoding time, and a memory complexity that scales with the input, boasting high performance across the full span of the frequency spectrum and surpassing previously established benchmarks on all datasets considered.

3 Methodology

3.1 Preliminaries

We will formulate supertagging as an iterative graph completion process. Our representation format is that of a heterogeneous graph, consisting of two types of nodes and edges. Lexical nodes (words) are vectors of dimensionality d_w ; the sentential structure can be encoded as a set of labeled edges in the cartesian product of words, $\mathcal{E}^w \in \mathbb{Z}^{s imes s}$, such that $\mathcal{E}^w_{i,j} = j-i$ is the distance between words j and i in the input sequence. Each lexical item *i* is associated to a binary branching tree T_i , representing the category assigned to the word in question. We will denote with $T_{i,k}$ the k-th node of the i-th tree, where enumeration starts from 1 for the root, and is inductively defined by following along the tree's breadth-first traversal. Tree nodes are represented as vectors of dimensionality d_n , the totality of n nodes in the graph then being a matrix $N \in \mathbb{R}^{n \times d_n}$; tree edges are converted into a sparse connectivity matrix $\mathcal{E}^n \in \mathbb{N}^{s \times n}$, where $\mathcal{E}_{i,\nu}^n$ is k if node ν occurs at position k within tree *i*, and zero otherwise. We will denote with $\mathcal{N}_{i,t}$ the depth-t neighborhood of a lexical item i as the set of nodes $[T_{i,2^t} \dots T_{i,2^{t+1}}]$, i.e. the set of nodes that are occupying depth t of tree i. To tell the two kinds of nodes apart, and for reasons that will

185

191

211

223

224

227

229

become clear in what is to come, we will refer to lexical nodes as states, and tree nodes as just nodes.

3.2 Breadth-First Parallel Decoding

Categories being trees, they are amenable to 186 breadth-first traversal. A sequence of categories can then be viewed as multiple sequences of primitives, sequence t corresponding to the concatenation of primitive symbols occurring along all trees 190 at depth t. Combined with ancestry information and a series of initial seeds (i.e. word vectors), it 192 is straightforward to use each sequence as a statetracking vector from which to predict the next se-194 quence of elements, essentially implementing a tree 195 unfolding function à la treeRNN. As mentioned 196 earlier, this fails to account for "horizontal" interactions between nodes occupying the same depth, and 198 also, as a result, "diagonal" interactions between 199 nodes living in different trees, practically reducing the task to a separable classification applied independently across words. A seemingly innocu-202 ous solution would be to employ a second form of recurrence, as in Alvarez-Melis and Jaakkola (2017); this would however once more break the input-to-output alignment and necessitate explicit handling of unoccupied tree positions, while inserting a multiplicative time complexity factor scaling 208 constantly with the input.¹ Our approach seeks to mitigate this by foregoing horizontal interactions, 210 but reinstating intra-tree interactions across depths. We do so by repurposing the initial sequence of lexical vectors, from input seeds to recurrent state-213 214 tracking vectors that arbitrate the decoding process across both sequence length and tree depth. In the 215 absence of a localized tree unfolding function, and 216 aiming to properly capture the "regularly irregular" 217 structure of the output space, we turn our attention towards structure-aware dynamic graph convolu-219 tions.

In high level terms, the process can be summarized as an iteration of three alternating stages of message passing rounds. Initially, state vectors are supplied from an arbitrary encoder network, and a fringe consisting of s unlabeled nodes is instantiated in alignment with the input sequence. From then on, and until a fix-point is reached:

1. Each state vector receives feedback in a manyto-one fashion from the last decoded nodes lying directly above it (initially none), yielding tree-contextual states.

2. The updated state vectors exchange messages with one another in a many-to-many fashion, yielding tree-and-sequence-contextual states.

232

233

234

235

236

237

238

239

240

241

242

243

244

245

247

248

249

250

251

252

253

254

256

257

258

259

260

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

3. The final states project class weights to their respective fringe nodes in a one-to-many fashion; depending on the arity of the decoded symbols, a next masked fringe is constructed with appropriate node positions and state-tonode edge indices; the process terminates when the next fringe is empty.

For a visual example, please refer to Appendix A.

3.3 Architecture

We now move on to detail the individual blocks that together make up the network's pipeline.

3.3.1 Node Embeddings

State vectors are temporally dynamic; they are initially supplied by an external encoder, and are then updated through a repeated sequence of three message passing rounds, described in the next subsections. Tree node embeddings, on the other hand, are not subject to temporal updates, but instead become dynamically "revealed" by the decoding process. They are computed on the basis of (i) their primitive symbol and (ii) their position within a tree.

Primitive symbol embeddings are obtained from a standard embedding table $W_e : S \to \mathbb{R}^{d_n}$ that contains a distinct vector for each symbol in the set of primitives S. When it comes to embedding positions, we are presented with a number of options. It would be straightforward to fix a vocabulary of positions, and learn a distinct vector for each. Such an approach would however lack elegance, as it would impose an ad-hoc bound to the shape of trees that can be encoded (contradicting the constructive paradigm), while also failing to account for the compositional nature of trees. We thus opt for a path-based approach, inspired by and extending the idea of Shiv and Quirk (2019). We note first that paths over binary branching trees form a semigroup, i.e. they consist of two primitives (namely a left and a right path), and an associative noncommutative binary operator that binds two paths together into a single new one. The archetypical example of a semigroup is matrix multiplication; we therefore instantiate a tensor $P \in \mathbb{R}^{2 \times n_d \times n_d}$ encoding each of the two path primitives as a linear map over symbol embeddings. From the above we can derive a function p that converts positions to linear

¹Consider that at each depth t we would need to iterate over $2^t \times s$ elements.

maps, by performing consecutive matrix multipli-281 cations of the primitive weights, as indexed by the reversed binary word of a node's position; e.g. the linear map corresponding to position $12_{10} = 1100_2$ would be $p(12) = P_0 P_0 P_1 P_1 \in \mathbb{R}^{d_n \times d_n}$. We flatten the final map by evaluating it against an initial seed vector ρ_0 , corresponding to the tree root.² To 287 stabilize training and avoid vanishing or exploding weights, we model paths as unitary transformations by parameterizing the two matrices of P290 to orthogonality using the exponentiation trick on 291 skew-symmetric bases (Bader et al., 2019; Lezcano Casado, 2019). The final embedding for a symbol σ occupying position k is then given by the 294 element-wise product of its positional and content 295 embeddings $p(k)(\rho_0) \circ (W_e(\sigma)) \in \mathbb{R}^{d_n}$.

3.3.2 Node Feedback

299

302

304

312

313

314 315

317

319

321

We update states with information from the last decoded nodes using a heterogeneous messagepassing scheme based on graph attention networks (Veličković et al., 2018; Brody et al., 2021). First, we use a bottleneck layer W_b to downproject the state vector into the nodes' dimensionality. Then, given a state h_i^t and a neighborhood $\mathcal{N}_{i,t}$, we compute a self-loop score $\tilde{\alpha}_{i,t} =$ $w_a \cdot (W_b(h_i^t)||\mathbf{0})$, as well as heterogeneous scores $\tilde{\alpha}_{i,\nu,t} = w_a \cdot (h_i^t || N_{\nu})$, for each node ν in the neighborhood, where $w_a \in \mathbb{R}^{2d_n}$ a dot-product weight and $W_b(h_i^t)||N_{\nu}$ the concatenation of the downsized state vector with node ν 's position-aware embedding N_{ν} . Scores are passed through a leaky rectifier non-linearity before being normalized to attention coefficients, from which we obtain the updated states as the weighted sum of incoming messages (further processed by a shallow network W_m) and a residual connection:

$$\tilde{h}_{i}^{t} = \sum_{\mu \in \mathcal{N}_{i,t}} \alpha_{i,\nu,t} W_{m} N_{\nu} + \alpha_{i,t} h_{i}^{t}$$

States receiving no node feedback (i.e. states that have completed decoding more than one time step ago) are thus protected from updates, preserving their content. In practice, we compute attention coefficients and message vectors independently for multiple heads, but omit them from the above equations to avoid cluttering the notation.

3.3.3 State Feedback

At the end of the node feedback stage, we are left with a sequence of locally contextualized states h_i^t . Recall that, owing to our encoding of the sentential structure, states form a fully connected graph, with edges weighted by relative distances between words. We embed these distances into the encoder's vector space using an embedding table $W_r \in \mathbb{R}^{2\delta \times d_w}$, where δ the maximum allowed distance, a hyper-parameter. Edges escaping the maximum distance threshold are truncated rather than clipped, in order to preserve memory and facilitate training, leading to a natural segmentation of the sentence into (overlapping) chunks. Following standard practices, we project states into query, key and value vectors (Vaswani et al., 2017), and compute the attention scores between words i and j using relative-position weighted attention (Shaw et al., 2018):

325

326

327

329

330

331

332

333

334

335

336

337

338

341

342

343

348

350

351

352

353

355

357

358

359

360

361

362

363

365

366

367

369

$$\tilde{a}_{i,j} = d_w^{-1/2} \left(W_q \tilde{h}_i^t \circ W_r \mathcal{E}_{i,j}^w \right) \cdot W_k \tilde{h}_j^t$$

From the normalized attention scores we obtain a new set of aggregated messages:

$$m'_{i,t} = \sum_{j \in \{0..s\}} \frac{\exp(\tilde{a}_{i,j}) W_v h_j^t}{\sum_{k \in \{0..s\}} \exp(\tilde{a}_{i,k})}$$

$$347$$

Same as before, queries, keys, values, edge embeddings and attention coefficients are distributed over many heads.

Aggregated messages are passed through a swish-gated feed-forward layer (Dauphin et al., 2017; Shazeer, 2020) to yield the next sequence of state vectors:

$$h_i^{t+1} = W_3 \left(\text{swish}_1(W_1 m'_{i,t}) \circ W_2 m'_{i,t} \right)$$

where $W_{1,2}$ are linear maps from the encoder's dimensionality to an intermediate dimensionality, and vice versa for W_3 .

3.3.4 Node Prediction

Finally, from a globally-contextualized state h_i^{t+1} we need to obtain class weights for the entirety of the neighborhood $\mathcal{N}_{i,t+1}$. We start by downprojecting the state vector into the node's dimensionality using yet another shallow network W_n . The resulting feature vectors are shared across all nodes of the same tree – to discriminate between them, we gate vectors against each node's positional embedding. From the latter, we obtain class weights by matrix multiplying them against the

²In practice, paths are efficiently computed once per batch for each unique tree position during training, and stored as fixed embeddings during inference.

418

419

427 428 429

430 431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

transpose of the symbol embedding table (Press and Wolf, 2017):

370

371

387

397

400

401

402

403

404 405

406

407

408

409

410

411

412

413

414

415

416

417

weights_{*i*,*k*} =
$$(p(k)(\rho_0) \circ W_n h_i^{t+1}) W_e^{\top}$$

During inference, the next fringe can be easily generated with minimal structure manipulation by using the indices of binary decoded symbols to 375 extract their positions, multiply those by two (to 376 create the positions of their left children), add one (to create the positions of their right children) and finally interleave the two; in the same vein, the 379 new state indices are simply the repetition of their respective ancestor indices. We hold on to the positional embeddings of the current fringe, as they will find use in the ensuing node feedback phase 384 unaltered.

3.3.5 Putting Things Together

We compose the previously detailed components into a single layer, which acts a sequence-wide, recurrent-in-depth decoder. We insert skip connections between the input and output of the messagepassing and feed-forward layers (He et al., 2016), and subsequently normalize each using root mean square normalization (Zhang and Sennrich, 2019).

4 Experiments

We employ our supertagging architecture in a range of diverse categorial grammar datasets spanning different languages and underlying grammar formalisms. In all our experiments, we bind our model to a monolingual BERT-style language model used as an external encoder, fine-tuned during training (Devlin et al., 2018). In order to homogenize the tokenization between the one directed by each dataset and the one required by the encoder, we make use of a simple localized attention aggregation scheme. The subword tokens together comprising a single word are independently projected to scalar values through a shallow feed-forward layer. Scalar values are softmaxed within their local group to yield attention coefficients over their respective BERT vectors, which are then summed together, in a process reminiscent of a cluster-wide attentive pooling (Li et al., 2016). In cases of datalevel tokenization treating multiple words as a single unit (i.e. assigning one type to what BERT perceives as many words), we mark all words following the first with a special [MWU] token, signifying they need to be merged to the left. This effectively adds an extra output symbol to the decoder, which

is now forced to do double duty as a sequence chunker. To avoid sequence misalignments and metric shifts during evaluation, we follow the merges dictated by the ground truth labels, and consider the decoder's output as correct only if all participating predictions match, assuming no implicit chunking oracles.

4.1 Datasets

We conduct experiments on the two variants of the English CCGBank, the French TLGbank and the Dutch Æthel proofbank. A high-level overview of the datasets is presented in Table 1, and short descriptions are provided in the following paragraphs. We refer the reader to the corresponding literature for a more detailed exposition.

| | CCC original | Gbank rebank | TLGbank | Æthel |
|------------------|-----------------|------------------------|---------|-------|
| Primitives | 37 | 40 | 27 | 60 |
| Zeroary | 35 | 38 | 19 | 31 |
| Binary | 2 | 2 | 8 | 29 |
| Categories | 1323 | 1619 | 851 | 5292 |
| in train | 1286 | 1575 | 803 | 4730 |
| depth avg. | 1.94 | 1.96 | 1.99 | 1.83 |
| depth max. | 6 | 6 | 7 | 35 |
| Test Sentences | 2407 | 2407 | 1571 | 5766 |
| length avg. | 23.00 | 24.27 | 27.58 | 16.61 |
| Test Tokens | 55371 | 56395 | 44302 | 98467 |
| Frequent (100+) | 54825 | 55690 | 43289 | 95253 |
| Uncommon (10-99) | 442 | 563 | 833 | 2213 |
| Rare (1-9) | 75 | 107 | 149 | 678 |
| Unseen (OOV) | 22 | 27 | 31 | 323 |

Table 1: Bird's eye view of datasets employed and relevant statistics. Test tokens are binned according to their corresponding categories' occurrence count in the respective dataset's training set. Token counts are measured before pre-processing. Unique primitives for the type-logical datasets are counted after binarization.

CCGBank The English CCGbank (original) (Hockenmaier and Steedman, 2007) and its refined version (rebank) (Honnibal et al., 2010) are resources of Combinatory Categorial Grammar (CCG) derivations obtained from the Penn Treebank (Taylor et al., 2003). CCG (Steedman and Baldridge, 2011) builds lexical categories with the aid of two binary slash operators, capturing forward and backward function application. Some additional rules lent from combinatory logic (Curry et al., 1958) permit constrained forms of type raising and function composition, allowing categories to remain relatively short and uncomplicated while keeping parsing complexity in check. The key difference between the two versions lies in

- 448 449
- 450

their tokenization and the plurality of categories assigned, the latter containing more assignments and a more fine-grained set of syntactic primitives, which in turn make it a slightly more challenging evaluation benchmark.

French TLGbank The French type-logical tree-453 bank (Moot, 2015) is a collection of proofs ex-454 tracted from the French treebank (Abeillé et al., 455 2003). The theory underlying the resource is that 456 of Multi-Modal Typelogical Grammars (Moortgat, 457 1996); annotations are deliberately made compat-458 ible with Displacement Calculus (Morrill et al., 459 2011) and First-Order Linear Logic (Moot and Pi-460 azza, 2001) at the cost of a small increase in lexical 461 sparsity. In short, the vocabulary of operators is 462 extended with two modalities that find use in licens-463 ing or restricting the applicability of rules related 464 to non-local syntactic phenomena. To adapt their 465 representation to our framework, we cast unary 466 operators into pseudo-binaries by inserting an arti-467 ficial terminal tree in a fixed slot within them. Due 468 to the absence of predetermined train/dev/test splits, 469 we randomize them with a fixed seed at a 80/10/10 470 ratio and keep them constant between repetitions. 471

Æthel Our last experimental test bed is 472 Æthel (Kogkalidis et al., 2020a), a dataset of 473 type-logical proofs for written Dutch sentences, 474 automatically extracted from the Lassy-Small 475 476 corpus (Noord et al., 2013). Æthel is geared towards semantic parsing, which means categories 477 employ linear implication $-\infty$ as their single binary 478 operator. An additional layer of dependency infor-479 mation is realized via unary modalities, now lifted 480 481 to classes of operators distinguishing complement and adjunct roles. The grammar assigns concrete 482 instances of polymorphic coordinator types, as 483 a result containing more and sparser categories 484 (some of which distinctively tall); considering also 485 its larger vocabulary of primitives, it makes for a 486 good stress test for our approach. We experiment 487 with the latest available version of the dataset 488 (version 0.9.dev1 at the time of writing). Same 489 as before, we impose a regular tree structure, 490 this time by merging adjunct (resp. complement) 491 markers with the subsequent (resp. preceding) 492 binary operator, which makes for an unambiguous 493 and invertible representational translation. 494

4.2 Implementation

495

496

497

We implement our model using PyTorch Geometric (Fey and Lenssen, 2019), which provides a highlevel interface to efficient low-level protocols, facilitating fast and pad-free graph manipulations. We share a single hyper-parameter setup across all experiments, obtained after a minimal logarithmic search over sensible initial values. Specifically, we set the node dimensionality d_n to 128 with 4 heterogeneous attention heads and the state dimensionality d_w to 768 with 8 homogeneous attention heads. We train using AdamW (Loshchilov and Hutter, 2018) with a batch size of 16, weight decay of 10^{-2} , and a learning rate of 10^{-4} , scaled by a linear warmup and cosine decay schedule over 25 epochs. During training we provide strict teacher forcing and apply feature and edge dropout at 20% chance. Our loss signal is derived as the label-smoothed negative log-likelihood between the network's prediction and the ground truth label (Müller et al., 2019). We procure pretrained base-sized BERT variants from the transformers library (Wolf et al., 2020): RoBERTa for English (Liu et al., 2019), BERTje for Dutch (de Vries et al., 2019) and CamemBERT for French (Martin et al., 2020), which we fine-tune during training, scaling their learning rate by 10% compared to the decoder.

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

4.3 Results

We perform model selection on the basis of validation accuracy, and gather the corresponding test scores according to the frequency bins of Table 1. Table 2 presents our results compared to relevant published literature. Evidently, our model surpasses established benchmarks in terms of overall accuracy, matching or surpassing the performance of both traditional supertaggers on common categories and constructive ones on the tail end of the frequency distribution.

We observe that the relative gains appear to scale with respect to the task's complexity. In the original version of the CCGbank, our model is only slightly superior to the next best performing model (in turn only marginally superior to the token-based classification baseline), whereas in the rebank version the absolute difference is one order of magnitude wider. The effect is even further pronounced for the harder type-logical datasets, which are characterized by a longer tail, leading to performance comparable to CCGbank's for the French TLGbank (despite it being significantly smaller and sparser), and a 10% absolute performance leap for Æthel (despite its unusually tall and complex types). We attribute this to increased returns from performance

| | accuracy (%) | | | | |
|---|--------------------------------|--------------------------------|----------------------|---|----------------------------|
| model | overall | frequent | uncommon | rare | unseen |
| CCG (original) | | | | | |
| Symbol Sequential LSTM /w n-gram oracles (Liu et al., 2021) | 95.99 | 96.40 | 65.83 | 8.65 [!] | |
| Symbol Sequential LSTM (Bhargava and Penn, 2020) | 96.00 | - | - | - | \sim 5 |
| Cross-View Training (Clark et al., 2018) | 96.10 | - | - | - | n/a |
| Recursive Tree Addressing (Prange et al., 2021) | 96.09 | 96.44 | 68.10 | 37.40 | 3.03 |
| BERT Token Classification (Prange et al., 2021) | 96.22 | 96.58 | 70.29 | 23.17 | n/a |
| Auentive Convolutions (Tian et al., 2020) | 90.25 | 90.04 | /1.04 | n/a | n/a |
| Heterogeneous Dynamic Convolutions (this work) | $96.29{\scriptstyle \pm 0.04}$ | $96.61{\scriptstyle \pm 0.04}$ | 4 72.06 ±0.72 | $34.45{\scriptstyle\pm1.58}$ | $\textbf{4.55}_{\pm 2.87}$ |
| CCG (rebank) | | | | | |
| Symbol Sequential Transformer [†] (Kogkalidis et al., 2019) | 90.68 | 91.10 | 63.65 | 34.58 | 7.41 |
| Symbol Sequential LSTM [†] (Bhargaya and Penn, 2020) | 93.92 | 94.39 | 65.48 | 19.00 | 0.00 |
| TreeGRU (Prange et al., 2021) | 94.62 | 95.10 | 64.24 | 25.55 | 2.47 |
| Recursive Tree Addressing (Prange et al., 2021) | 94.70 | 95.11 | 68.86 | 36.76 | 4.94 |
| Token Classification (Prange et al., 2021) | 94.83 | 95.27 | 68.68 | 23.99 | n/a |
| Heterogeneous Dynamic Convolutions (this work) | $95.07{\scriptstyle\pm0.04}$ | 95.45±0.04 | 71.40 ±1.15 | $\textbf{37.19}{\scriptstyle \pm 1.81}$ | 3.70 ± 0.00 |
| French TLGbank | | | | | |
| ELMo & LSTM Classification (Moot, 2019) | 93.20 | 95.10 | 75.19 | 25.85 | n/a |
| BERT Token Classification [‡] | 95.93 | 96.44 | 81.39 | 47.45 | n/a |
| Heterogeneous Dynamic Convolutions (this work) | $95.92{\scriptstyle\pm0.01}$ | 96.40 ± 0.01 | 81.48±0.97 | $\textbf{55.37}{\scriptstyle \pm 1.00}$ | 7.26 ±2.67 |
| Æthel | | | | | |
| Symbol Sequential Transformer [*] (Kogkalidis et al., 2020b) | 83.67 | 84.55 | 64.70 | 50.58 | 24.55 |
| BERT Token Classification [‡] | 93.52 | 94.83 | 71.85 | 38.06 | n/a |
| Heterogeneous Dynamic Convolutions (this work) | 93.67 ±0.04 | 94.72±0.13 | 3 73.45±0.46 | 53.83 ±1.14 | 15.79±1.32 |

¹Accuracy over both bins, with a frequency-truncated training set (authors claim no difference when using the full set).

[†]Numbers from Prange et al. (2021).

[‡]Our replication.

548

549

550

551

552

554

556

558

559

560

562

564

565

566

*Model trained and evaluated on an older dataset version and tree sequences spanning less than 140 nodes in total.

Table 2: Model performance across datasets and compared to recent studies. Numbers are taken from the papers cited unless otherwise noted. For our model, we report averages and standard deviations over 6 runs. Bold face fonts indicate (within standard deviation of) highest performance.

in the rare and uncommon bins; there is a synergistic effect between the larger population of these bins pronouncing even minor improvements, and acquisition of rarer categories apparently benefiting from the plurality of their respective bins in a self-regularizing manner.

Finally, to investigate the relative impact of each network component, we conduct an ablation study where message passing components are removed from their network in their entirety. Removing the state feedback component collapses the network into a token-wise separable recurrence, akin to a graph-featured RNN without a hidden-to-hidden affine map. Removing the node feedback component turns the network into a Universal Transformer (Dehghani et al., 2018) composed with a dynamically adaptive classification head. Removing both is equatable to a 1-to-many contextualized token classification that is structurally unfolded in depth. Our results, presented in Table 3, verify first a positive contribution from both components, indi-

| | -sf | -nf | -sf-nf |
|----------------|-------|-------|--------|
| CCG (original) | -0.05 | -0.01 | -0.08 |
| CCG (rebank) | -0.12 | -0.04 | -0.07 |
| French TLGbank | -0.13 | -0.14 | -0.23 |
| Æthel | -0.24 | -0.12 | -0.37 |

Table 3: Absolute difference in overall accuracy when removing the state and node feedback components (averages of 3 repetitions).

cating the importance of both information sharing axes. In three out of the four datasets, the relative gains of incorporating state feedback outweigh those of node feedback, and are most pronounced in the case of Æthel, likely due to its positionally agnostic types. With the exception of CCGrebank, relinquishing both kinds of feedback largely underperforms having either one, experimentally affirming their compatibility. 569

570

571

572

573

574

575

576

577

578

579

5 Related Work

Our work bears semblance and owes credit to various contemporary lines of work. From the architec-

tural angle, we perceive our work as an application-581 specific offspring of weight-tied architectures, dy-582 namic graph convolutions and structure-aware self-583 attention networks. The depth recurrence of our decoder is inspired by weight-tied architectures (Dehghani et al., 2018; Bai et al., 2019) and their graph-586 oriented variants (Li et al., 2016), which model neu-587 ral computation as the fix-point iteration of a single layer against a structured input, thus allowing for a dynamically adaptive computation "depth" - albeit 590 with a constant parameter count. Analogously to 591 structure-aware self-attention networks (Zhu et al., 592 2019; Cai and Lam, 2020) and graph attentive net-593 works (Veličković et al., 2018; Yun et al., 2019; 594 Ying et al., 2021; Brody et al., 2021), our decoder employs standard query/key and fully-connected attention mechanisms injected with structurally biased representations, either at the edge or at the node level. Finally, akin to dynamic graph ap-599 proaches (Liao et al., 2019; Pareja et al., 2020), our decoder forms a closed loop system that autoregressively generates its own input, in the process becoming exposed to subgraph structures that drastically differ between time steps.

605 From the application angle, our proposal is a refinement of and a continuation to recent advances in categorial grammar supertagging. Similar to 607 the transition from words to subword units (Sennrich et al., 2016), constructive supertaggers seek to bolster generalization by disassembling syntactic categories into smaller indivisible units, thereby 611 incorporating structure at a finer granularity scale. 612 The original approach of Kogkalidis et al. (2019), 613 later adopted by Bhargava and Penn (2020), employed seq2seq models to directly translate an in-615 put text to a flattened projection of a categorial 616 sequence, demonstrating that the correct prediction 617 of categories unseen during training is indeed feasi-618 619 ble. Prange et al. (2021) improved upon the process through the explicit accounting of the tree structure embedded within categorial types, while Liu et al. 621 (2021) explored the orthogonal approach of employing a transition-based "parser" over individual 623 categories. Outside the constructive paradigm, Tian 624 et al. (2020) employed graph convolutions over sen-625 tential edges built from static, lexicon-based preferences. Our approach is a bridge between prior 627 works; our modeling choice of structure-aware 628 graph convolutions boasts the merits of explicit 629 sentential and tree-structured edges, a structurally constrained, valid-by-construction output space, 631

favorable memory and time complexities, partial auto-regressive context flows, end-to-end differentiability with no vocabulary requirements, and minimal rule-based structure manipulation. 632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

6 Conclusion

We have proposed a novel supertagging methodology, where both the linear order of the output sequence and the tree-like structure of its elements is made explicit. To represent the different information sources and their disparate sizes and scales, we turned to heterogeneous graph attention networks. To capture the auto-regressive dependencies between different trees, we formulated the task as a dynamic graph completion process, aligning each subsequent temporal step with a higher order tree node neighborhood and predicting them in parallel across the entire sequence. We tested our methodology on four different datasets spanning three languages and as many grammar formalisms, establishing new state of the art scores in the process. Through our ablation studies, we showed the importance of incorporating both intra- and inter-tree context flows, to which we attribute our system's performance.

Other than architectural adjustment and optimizations, several interesting ideas present themselves as promising research avenues. First, it is worthwhile to consider adaptations of our framework to either allow an efficient integration of more "exotic" context pathways, e.g. sibling node interactions, or alter the graph's decoding order altogether. On a related note, for formalisms faithful to the linear logic roots of categorial grammars, it seems reasonable to anticipate that the goal graph can be compactified by collapsing primitive nodes of opposite polarity according to their interactions, unifying the tasks of supertagging and parsing with a single end-to-end framework. Finally, and despite its success, our methodology is not without limitations. Crucially (and like all decoders that perform multiple assignments concurrently) our model trades inference speed for an incompatibility with local greedy algorithms like beam search - finding ways to reconcile the two is a pressing matter.

Practice aside, our results pose further evidence that lexical sparsity, historically deemed the categorial grammar's curse, might well just require a change of perspective to tame and deploy as the answer to the very problem it poses.

References

- 690 694 701
- 710 711 716 717 718 719 720 721

- 724
- 725

729

- 733

- Anne Abeillé, Lionel Clément, and François Toussenel. 2003. Building a treebank for French. In Treebanks, pages 165–187. Springer.
- David Alvarez-Melis and Tommi S. Jaakkola. 2017. Tree-structured decoding with doubly-recurrent neu-In 5th International Conference ral networks. on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Philipp Bader, Sergio Blanes, and Fernando Casas. 2019. Computing the matrix exponential with an optimized Taylor polynomial approximation. Mathematics, 7(12):1174.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. 2019. Deep equilibrium models. Advances in Neural Information Processing Systems, 32.
- Srinivas Bangalore and Aravind K. Joshi. 1999. Supertagging: An approach to almost parsing. Computational Linguistics, 25(2):237-265.
- Aditya Bhargava and Gerald Penn. 2020. Supertagging with CCG primitives. In Proceedings of the 5th Workshop on Representation Learning for NLP, pages 194-204, Online. Association for Computational Linguistics.
- Shaked Brody, Uri Alon, and Eran Yahav. 2021. How attentive are graph attention networks? arXiv preprint arXiv:2105.14491.
- Deng Cai and Wai Lam. 2020. Graph transformer for graph-to-sequence learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 7464–7471.
- Kevin Clark, Minh-Thang Luong, Christopher D. Manning, and Quoc Le. 2018. Semi-supervised sequence modeling with cross-view training. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 1914-1925, Brussels, Belgium. Association for Computational Linguistics.
- Haskell Brooks Curry, Robert Feys, William Craig, J Roger Hindley, and Jonathan P Seldin. 1958. Combinatory Logic, volume 1. North-Holland Amsterdam.
- Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. 2017. Language modeling with gated convolutional networks. In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 933-941.
- Wietse de Vries, Andreas van Cranenburgh, Arianna Bisazza, Tommaso Caselli, Gertjan van Noord, and Malvina Nissim. 2019. BERTje: A Dutch BERT model. arXiv preprint arXiv:1912.09582.

Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. 2018. Universal transformers. In International Conference on Learning Representations.

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

762

763

764

765

766

767

768

769

770

771

773

775

776

777

778

779

780

781

782

783

784

785

786

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Matthias Fey and Jan E. Lenssen. 2019. Fast graph representation learning with PyTorch Geometric. In ICLR Workshop on Representation Learning on Graphs and Manifolds.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770-778.
- Julia Hockenmaier and Mark Steedman. 2007. CCGbank: a corpus of CCG derivations and dependency structures extracted from the Penn Treebank. Computational Linguistics, 33(3):355-396.
- Matthew Honnibal, James R Curran, and Johan Bos. 2010. Rebanking CCGbank for improved np interpretation. In Proceedings of the 48th annual meeting of the association for computational linguistics, pages 207–215.
- Konstantinos Kogkalidis, Michael Moortgat, and Tejaswini Deoskar. 2019. Constructive type-logical supertagging with self-attention networks. In Proceedings of the 4th Workshop on Representation Learning for NLP (RepL4NLP-2019), pages 113-123, Florence, Italy. Association for Computational Linguistics.
- Konstantinos Kogkalidis, Michael Moortgat, and Richard Moot. 2020a. ÆTHEL: Automatically extracted typelogical derivations for Dutch. In Proceedings of the 12th Language Resources and Evaluation Conference, pages 5257-5266, Marseille, France. European Language Resources Association.
- Konstantinos Kogkalidis, Michael Moortgat, and Richard Moot. 2020b. Neural proof nets. In Proceedings of the 24th Conference on Computational Natural Language Learning, pages 26-40, Online. Association for Computational Linguistics.
- Mike Lewis and Mark Steedman. 2014. Improved CCG parsing with semi-supervised supertagging. Transactions of the Association for Computational Linguistics, 2:327-338.
- Mario Lezcano Casado. 2019. Trivializations for gradient-based optimization on manifolds. Advances in Neural Information Processing Systems, 32.
- Yujia Li, Richard Zemel, Marc Brockschmidt, and Daniel Tarlow. 2016. Gated graph sequence neural networks. In Proceedings of ICLR'16.

- 790
- 799

- 804
- 807
- 810
- 813 814

816 817 818

- 819
- 821 822 823

825 826

827

831

832

833 834

836 838

841

- Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Will Hamilton, David K Duvenaud, Raquel Urtasun, and Richard Zemel. 2019. Efficient graph generation with graph recurrent attention networks. Advances in Neural Information Processing Systems, 32.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.
- Yufang Liu, Tao Ji, Yuanbin Wu, and Man Lan. 2021. Generating CCG categories. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 13443–13451.
- Ilya Loshchilov and Frank Hutter. 2018. Fixing weight decay regularization in adam.
- Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. 2020. CamemBERT: a tasty French language model. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7203-7219, Online. Association for Computational Linguistics.
- Michael Moortgat. 1996. Multimodal linguistic inference. JoLLI, 5(3/4):349-385.
- Richard Moot. 2015. A type-logical treebank for French. Journal of Language Modelling Vol, 3(1):229-264.
- Richard Moot. 2019. Reconciling vectors with proofs for natural language processing. Compositionality in formal and distributional models of natural language semantics, 26th Workshop on Logic, Language, Information and Computation (WoLLIC 2019). Retrieved from https://richardmoot. github.io/Slides/WoLLIC2019.pdf.
- Richard Moot and Mario Piazza. 2001. Linguistic applications of first order intuitionistic linear logic. Journal of Logic, Language and Information, 10(2):211-232.
- Glyn Morrill, Oriol Valentín, and Mario Fadda. 2011. The displacement calculus. Journal of Logic, Language and Information, 20(1):1–48.
- Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. 2019. When does label smoothing help? Advances in neural information processing systems, 32.
- Gertjan van Noord, Gosse Bouma, Frank Van Eynde, Daniël de Kok, Jelmer van der Linde, Ineke Schuurman, Erik Tjong Kim Sang, and Vincent Vandeghinste. 2013. Large scale syntactic annotation of written Dutch: Lassy. In Essential Speech and Language Technology for Dutch, pages 147–164. Springer.

Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao Schardl, and Charles Leiserson. 2020. EvolveGCN: Evolving graph convolutional networks for dynamic graphs. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 5363–5370.

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

894

895

896

- Jakob Prange, Nathan Schneider, and Vivek Srikumar. 2021. Supertagging the long tail with tree-structured decoding of complex categories. Transactions of the Association for Computational Linguistics, 9:243-260.
- Ofir Press and Lior Wolf. 2017. Using the output embedding to improve language models. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, pages 157-163, Valencia, Spain. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715-1725, Berlin, Germany. Association for Computational Linguistics.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-attention with relative position representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 464-468, New Orleans, Louisiana. Association for Computational Linguistics.
- Noam Shazeer. 2020. GLU variants improve transformer. arXiv preprint arXiv:2002.05202.
- Vighnesh Shiv and Chris Quirk. 2019. Novel positional encodings to enable tree-based transformers. Advances in Neural Information Processing Systems, 32.
- Mark Steedman and Jason Baldridge. 2011. Combinatory categorial grammar. In Robert Borsley and Kersti Börjars, editors, Non-Transformational Syntax: Formal and Explicit Models of Grammar, pages 181-224. Wiley-Blackwell.
- Ann Taylor, Mitchell Marcus, and Beatrice Santorini. 2003. The Penn treebank: an overview. Treebanks, pages 5-22.
- Yuanhe Tian, Yan Song, and Fei Xia. 2020. Supertagging Combinatory Categorial Grammar with attentive graph convolutional networks. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6037–6044, Online. Association for Computational Linguistics.

- 904 905 906 907 908 909 910
- 911 912 913 914 915 916 917
- 918 919
- 923 925
- 926 927 928
- 929
- 930 931
- 935

936

938

939

- 942
- 943

947

950

951

953

Jie Zhu, Junhui Li, Muhua Zhu, Longhua Qian, Min Zhang, and Guodong Zhou. 2019. Modeling graph structure in transformer for better AMR-to-text generation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 5459-5468, Hong Kong, China. Association for Computational Linguistics.

- Ashish Vaswani, Yonatan Bisk, Kenji Sagae, and Ryan Musa. 2016. Supertagging with LSTMs. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 232-237, San Diego, California. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems, 30.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In International Conference on Learning Representations.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38-45, Online. Association for Computational Linguistics.
- Wenduan Xu, Michael Auli, and Stephen Clark. 2015. CCG supertagging with a recurrent neural network. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 250-255, Beijing, China. Association for Computational Linguistics.
- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do transformers really perform badly for graph representation? Advances in Neural Information Processing Systems, 34.
- Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. 2019. Graph transformer networks. Advances in neural information processing systems, 32.
- Biao Zhang and Rico Sennrich. 2019. Root mean
- square layer normalization. Advances in Neural Information Processing Systems, 32.

954

A Visualization of the decoding process



(a) State vectors independently receive auto-regressive feedback from their last decoded respective fringe in a many-to-one fashion.



(b) The tree-contextual states exchange messages with one another in a many-to-many fashion.



(c) The final states project class weights to their respective fringe nodes in a one-to-many fashion; depending on the arity of the decoded symbols, a next masked fringe is constructed.

Figure 1: Visualization of one step of the decoding process for an (abstract) example sequence, focusing on the central tree T_i and starting from the partially decoded output at step 1. Node content is intentionally left unspecified so as not to add grammar-specific overhead, but tree structure is assumed fixed and given by binary nodes $T_{i,1} T_{i,3}$ $T_{i,7}$ and $T_{i+1,1}$ (rest zeroary). The computations prescribed by each subfigure take place in parallel across all nodes, trees & sentences in the batch.

955

. . .

. . .