Revisiting the Compositional Generalization Abilities of Neural Sequence Models

Anonymous ACL submission

Abstract

Compositional generalization is a fundamental trait in humans, allowing us to effortlessly combine known phrases to form novel sentences. Recent works have claimed that standard seq-005 to-seq models severely lack the ability to compositionally generalize. In this paper, we focus on one-shot primitive generalization as introduced by the popular SCAN benchmark. We demonstrate that modifying the training distribution in simple and intuitive ways enables standard seq-to-seq models to achieve nearperfect generalization performance, thereby showing that their compositional generalization abilities were previously underestimated. We perform detailed empirical analysis of this phenomenon. Our results indicate that the gener-017 alization performance of models is highly sensitive to the characteristics of the training data which should be carefully considered while designing such benchmarks in future.

1 Introduction

021

037

According to the *principle of compositionality*, the meaning of a complex expression (e.g., a sentence) is determined by the meaning of its individual constituents and how they are combined. Humans can effectively recombine known parts to form new sentences that they have never encountered before. Despite the unprecedented achievements of standard seq-to-seq networks such as LSTMs and Transformers in NLP tasks, previous work has suggested that they are severely limited in their ability to generalize compositionally (Lake and Baroni, 2018; Furrer et al., 2020).

Problem Statement. Our work relates to a central challenge posed by compositional generalization datasets such as SCAN (Lake and Baroni, 2018) and Colors (Lake et al., 2019), which we refer to as *one-shot primitive generalization*: The dataset consists of *input-output sentence* pairs (e.g. 'walk twice \rightarrow WALK WALK'); input sentences



Figure 1: Overview of the SCAN generalization task (left) and our approach (right) that enables standard neural sequence models to generalize compositionally.

041

043

044

045

047

049

051

053

060

are formed from primitive words ('walk') and function words ('twice') and are generated by a CFG; output sentences are obtained by applying an interpretation function. Crucially, there is a systematic difference between the train and test splits¹: While the former has a *single* example of an *isolated primitive* (e.g., the primitive definition 'jump \rightarrow JUMP' in SCAN), the latter consists of compositional sentences with this isolated primitive (e.g. 'jump twice \rightarrow JUMP JUMP'). See Fig. 1, left.

A model with the right inductive bias should generalize on the test data after having seen compositional expressions with other primitives during training. The need for such inductive bias is justified via psychological experiments (Lake et al., 2019) indicating that humans do have the ability to generalize on such tasks. Previous works have suggested that seq-to-seq models lack the appropriate inductive bias necessary to generalize on this task since they achieve near-zero accuracy on both

¹We use the term *systematicity* in the rest of the paper to refer to this difference.

074

084

091

101

102

103

105

106

107

061

SCAN and Colors benchmarks. This has led to the development of many specialized architectures (Li et al., 2019; Gordon et al., 2020; Chen et al., 2020), learning procedures (Lake, 2019; Conklin et al., 2021) and data augmentation methods (Andreas, 2020; Guo et al., 2020) to solve the task.

Contributions. The primary claim of our paper is that, contrary to prior belief, neural sequence models such as Transformers and RNNs do have an inductive bias² to generalize compositionally which can be enabled using the right supervision. (i) We show that by making simple and intuitive changes to the training data distribution, standard seq-to-seq models can achieve high generalization performance even with a training set of size less than 20% of the original training set. In particular, if we incorporated examples with more novel primitives in the training set without necessarily increasing the size of the training set (see right part of Fig. 1), then the generalization performance of standard seq-to-seq models improves and reaches near-perfect score after a certain point. Our results also exemplify the importance of the training distribution apart from architectural changes and demonstrate that providing the right supervision can significantly improve the generalization abilities of the models. (ii) We investigate the potential cause behind the improvement in generalization performance and observe that the embedding of the isolated primitive becomes more similar to other primitives when the train set has higher number of primitives and their use cases. (iii) To understand the phenomenon better, we characterize the effect of different distributions, model capacity and transferability, and show that the parameters of the experimental setting play a crucial role while evaluating the generalization abilities of models.

2 Enabling Generalization by Providing the Right Supervision

Setup. We will focus on the SCAN and Colors datasets.³ Both of these datasets have a single *iso-lated primitive*. We refer to all other primitives as *example primitives*. The original training set of SCAN has 13.2k examples while the test set has 7.7k examples. Colors has just 14 training examples and 8 test examples. More details on these datasets can be found in Appendix B.



Figure 2: Generalization performance (\uparrow) on SCAN and Colors improves with higher number of example primitives in the train set.

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

Adding More Primitives. Both the SCAN and Colors training set have three example primitives. We modify the training set such that the set of example primitives present in the dataset is higher. To do so, we add new primitives to the language which are simply random words (e.g., 'swim', 'clap', etc.) that have the same semantics and follow the same grammar rules as other existing primitives (see Fig. 1 for illustration). The new primitives act as example primitives in our training set. For SCAN, we control the size of the training set such that it is at most the size of the original dataset.⁴ To generate the training set, we randomly sample the examples from the new grammar and add one separate example with the isolated primitive. The test set is untouched and remains the same.

Main Observation. Fig. 2 shows the generalization performance of Transformers and LSTM based seq-to-seq models. We observe that there is a clear trend of improvement in compositional generalization as we increase the number of example primitives and their use cases. It is also surprising to see that on SCAN, Transformers perform at par with some recently proposed specialized architectures (Li et al., 2019; Gordon et al., 2020) and even better than certain architectures (Russin et al., 2019).

Implication. Since the training set contains only one example with the isolated primitive and the test set is untouched, one-shot primitive generalization is preserved. Hence our results clearly show that standard neural sequence models have 'some' inductive bias required to generalize on such outof-distribution tasks even if it is not as strong as that of specialized architectures designed primarily

²However, note that this inductive bias is not as strong as that of specialized architectures designed for these tasks.

³Results on COGS (Kim and Linzen, 2020) can be found in Appendix D.

⁴The training set size |T| is kept fixed by discarding original examples and adding (|T|/#primitives) examples per primitive. Because of extremely small data size, we cannot do this for Colors while also trying to illustrate our idea.



Figure 3: Measuring the distance of embedding of *iso-lated primitive* with embeddings of example primitives for learned Transformer and LSTM models as we increase the number of example primitives in SCAN.

to solve these tasks. Our results are in contradiction to previously suggested limitations of standard seqto-seq models in terms of primitive generalization (Lake and Baroni, 2018; Furrer et al., 2020; Baroni, 2020). While it is important to develop architectures with better compositional generalization abilities, we wish to highlight that synthetic benchmarks such as SCAN require a model with very strong inductive bias and tend to underestimate the generalization abilities of baseline models.

143

144

145

146

147

148

149

150

152

153

154

155

156

157

159

161

162

165

166

167

168

170

171

173

174

175

176

177

Other findings. We find that model capacity critically influences the generalization performance of models. We explore this in more detail in Appendix C. While we have shown that these models can generalize from one-shot exposure to primitive definitions, our results also hold for the more general case where the one-shot exposure of the primitive is in a sentence⁵ (e.g. 'jump twice \rightarrow JUMP JUMP').

Prior Work. Note that our work is unrelated to previous works that propose data augmentation approaches for compositional generalization tasks (Andreas, 2020; Guo et al., 2020; Akyürek et al., 2021). (1) The datasets created by some of these augmentation methods do not preserve the systematic differences between train and test sets, while our datasets do. 6 (2) The objective of these works was to devise a method to improve the performance on compositional generalization whereas the focus of our work is not to develop a general method; rather we want show that baseline seq-to-seq models are capable of generalizing compositionally even without breaking systematicity. (3) These methods add additional data resulting in datasets of larger sizes whereas we control for data size.



⁶We discuss this in more detail in the Appendix G.



Figure 4: Visualizing the *t*-SNE reduced embeddings of isolated primitive (\blacktriangle), example primitives (\blacksquare) and non-primitives (\bullet) from a learned Transformer model as we increase number of example primitives in SCAN.

2.1 Embedding of Isolated Primitive

Our results raise the question: Why do Transformers and LSTMs generalize better when the training data has more example primitives? Compositional generalization in our setting requires a model to learn to apply the same rules to the isolated primitive as it does to the other example primitives. Thus, we analyze the change in the learned embedding of the isolated primitive (such as 'jump') with respect to other primitives in different settings.

178

179

181

182

183

184

185

187

188

189

190

191

193

194

195

196

197

198

199

200

201

202

203

204

205

In particular, we compare the average distance with other primitives before and after adding certain number of primitives to training data (this is the same setting as mentioned earlier in this section). We find that as we increase the number of example primitives in the training set, the embedding of the isolated primitive gets closer to the example primitives (Fig. 3) in terms of Euclidean, Manhattan and Cosine Distances. If the embedding of the isolated primitive is closer to other primitives, then the model is more likely to operate over it in a similar fashion and apply the same rules as it does over the other primitives.

This phenomenon is also illustrated in *t*-SNE plots (Fig. 4) of the learned embeddings where the embedding of the isolated primitive seems closer to example primitives when there are more example primitives in the dataset. Hence, a possible reason behind improved generalization performance could be the difference in the learned embeddings.⁷

⁷More fundamental reasons for difference in learned embeddings, such as learning dynamics, are beyond our scope.



Figure 5: Measuring the generalization performance of Transformer on different types of train set distributions of the SCAN dataset.

3 Exploring the Impact of Training Distributions

208

209

210

211

212

213

214

215

216

217

219

221

222

233

239

240

242

243

244

In this section, we analyze the influence of different training distributions on the generalization performance of the model. In the previous experiments, the data generating distribution was uniform over all possible samples. Here, we alter the distribution by varying the number of examples for each primitive. We experiment with linearly, quadratically and exponentially increasing probability distribution functions. For instance, in the quadratically increasing case, a training set with 10 primitives will have one primitive with 1 example, the next one with 4, another one with 9 examples and so on. The general idea is that all the example primitives do not have equal representation in the training data. Upon training the models on different distributions, we observed that the models generalize well even with fewer number of example primitives when their distribution is linearly or quadratically increasing (Fig. 5a). On the other hand models struggle to generalize when the distribution is very skewed (exponential). In that case, most primitives $(\sim 60\%)$ appear in only one sentence in the training data and some primitives ($\sim 10\%$) appear in over hundred sentences each. The failure to generalize on such data implies that extra primitives must be added as part of multiple sentences; just adding the definition or a single example for each primitive does not help the model to leverage it.

We then try to characterize the relationship between the number of example primitives and the amount of data required for the model to generalize well on the test data, when the example primitives are uniformly distributed. We create different training sets by varying the total number of example primitives #primitives; for each example primitive, we draw #examples number of samples uniformly from the CFG. Fig. 5b shows the generalization performance of Transformers for each of these training sets. The size of each training set is the product of the row and column values (#primitives \times #examples). As expected, the upper-right triangle has higher scores indicating that the sample requirement decreases as we add more primitives to the dataset. Surprisingly, the top-left cell indicates that Transformers can achieve high performance even with 2k training examples which is less than **20%** of the original SCAN training set. 247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

283

285

286

287

288

289

290

292

293

294

Transfer. We wish to check whether the inductive bias that is enabled when a model is trained on more number of example primitives can be transferred to a scenario where the number of example primitives is limited. We create a *pretraining* set with 50 example primitives uniformly distributed, each of them having 200 examples. The *finetuning* set is the original SCAN training set and the test set is the original SCAN test set. The model is first trained from scratch on the pretraining set and then finetuned on the finetuning set. We find that if we allow all the parameters of a Transformer Model to be updated during the finetuning phase on the original SCAN training set, then the model generalizes very poorly. On the other hand, when we freeze the weights of the encoder and decoder after the pretraining phase, and only allow the embedding and output layers to be updated then the model generalizes near-perfectly on the test set. Our hypothesis is that in the latter setting, the task becomes simpler for the model since it only has to align the embeddings of the primitives in the finetuning phase with the embeddings of the primitives seen during the pretraining phase. This experiment also indicates that the previously learned rules during pretraining can help a model to compositionally generalize on novel primitives.

4 Conclusion

While it is essential to make progress in building architectures with better compositional generalization abilities, we showed that the generalization performance of standard seq-to-seq models (often used as baselines) is underestimated. A broader implication of our experiments is that although systematicity must be preserved when designing such benchmarks, it is imperative to carefully explore different parameters associated with the experimental setup to draw robust conclusions about a model's generalization abilities.

References

297

298

300

301

302

305

307

310

311

313

314

315

316

317

319

320

321

322

323

330

331

332

338

339

341

345

347

351

- Ekin Akyürek, Afra Feyza Akyürek, and Jacob Andreas. 2021. Learning to recombine and resample data for compositional generalization. In *International Conference on Learning Representations*.
 - Jacob Andreas. 2020. Good-enough compositional data augmentation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7556–7566, Online. Association for Computational Linguistics.
 - Anonymous. 2022. From SCAN to real data: Systematic generalization via meaningful learning. In Submitted to The Tenth International Conference on Learning Representations. Under review.
 - Marco Baroni. 2020. Linguistic generalization and compositionality in modern artificial neural networks. *Philosophical Transactions of the Royal Society B*, 375(1791):20190307.
 - Paul Bloom. 2000. *How Children Learn the Meanings* of Words. MIT Press.
 - Xinyun Chen, Chen Liang, Adams Wei Yu, Dawn Song, and Denny Zhou. 2020. Compositional generalization via neural-symbolic stack machines. In *Advances in Neural Information Processing Systems*, volume 33, pages 1690–1701. Curran Associates, Inc.
 - Henry Conklin, Bailin Wang, Kenny Smith, and Ivan Titov. 2021. Meta-learning to compositionally generalize. 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), pages 3322–3335, Online. Association for Computational Linguistics.
 - Róbert Csordás, Kazuki Irie, and Juergen Schmidhuber.
 2021. The devil is in the detail: Simple tricks improve systematic generalization of transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 619–634, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Daniel Furrer, Marc van Zee, Nathan Scales, and Nathanael Schärli. 2020. Compositional generalization in semantic parsing: Pre-training vs. specialized architectures.
- Jonathan Gordon, David Lopez-Paz, Marco Baroni, and Diane Bouchacourt. 2020. Permutation equivariant models for compositional generalization in language. In *International Conference on Learning Representations*.
- Demi Guo, Yoon Kim, and Alexander Rush. 2020. Sequence-level mixed sample data augmentation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5547–5552, Online. Association for Computational Linguistics.

Najoung Kim and Tal Linzen. 2020. COGS: A compositional generalization challenge based on semantic interpretation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9087–9105, Online. Association for Computational Linguistics. 353

354

356

357

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

389

390

391

392

393

396

397

398

399

400

401

402

403

- Brenden Lake and Marco Baroni. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2873–2882. PMLR.
- Brenden M Lake. 2019. Compositional generalization through meta sequence-to-sequence learning. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc.
- Brenden M. Lake, Tal Linzen, and Marco Baroni. 2019. Human few-shot learning of compositional instructions.
- Yuanpeng Li, Liang Zhao, Jianyu Wang, and Joel Hestness. 2019. Compositional generalization for primitive substitutions. 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 4293–4302, Hong Kong, China. Association for Computational Linguistics.
- Qian Liu, Shengnan An, Jian-Guang Lou, Bei Chen, Zeqi Lin, Yan Gao, Bin Zhou, Nanning Zheng, and Dongmei Zhang. 2020. Compositional generalization by learning analytical expressions. In *Advances in Neural Information Processing Systems*, volume 33, pages 11416–11427. Curran Associates, Inc.
- Santiago Ontañón, Joshua Ainslie, Vaclav Cvicek, and Zachary Fisher. 2021. Making transformers solve compositional tasks.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc.
- Jake Russin, Jason Jo, Randall C. O'Reilly, and Yoshua Bengio. 2019. Compositional generalization in a deep seq2seq model by separating syntax and semantics.

A Implementation Details

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

We use 8 NVIDIA Tesla P100 GPUs each with 16
GB memory to run our experiments. All models are implemented in PyTorch (Paszke et al., 2019).
We do not use any pretrained models and all embeddings are learnt from scratch. Parameters are updated using Adam Optimization. All results are an average of 5 different runs with random seeds. The dataset-specific hyperparameters used for each model are shown in Table 1.

B Primitive Generalization Datasets

In this paper, we show results on three datasets that evaluate primitive generalization.

SCAN (Lake and Baroni, 2018) is a supervised sequence-to-sequence semantic parsing task wherein the natural language input command has to be transformed to the corresponding set of actions. The complete dataset consists of all the commands (a total of 20,910) generated by a phrasestructure grammar and the corresponding sequence of actions, produced according to a semantic interpretation function. The benchmark consists of 4 splits: random, add jump, turn left and length. We work on the 'add jump' split which was designed to test primitive generalization. In this split, the test set (size: 7706) is made up of all the compositional sentences with the primitive 'jump' (which we refer to as the *isolated primitive*). The train set (size: $13,204^8$) has just one example of the isolated primitive (i.e. the primitive definition 'jump \rightarrow JUMP') and other examples demonstrating the definitions and compositions of the three other primitives (which we refer to as the *example primitives*). Table 2 illustrates the task.

Colors (Lake et al., 2019) is a sequence-tosequence task that was designed to measure human inductive biases. Apart from the challenge of primitive generalization, this dataset poses an additional challenge of low-resource learning for neural sequence models. The train set has just 14 examples that are either primitive definitions of the four primitives or examples with compositions of the three example primitives and three operations (concatenation, repetition and wrapping). The test set has 8 examples⁹ with compositions of the isolated



Figure 6: Measuring the generalization performance of a Transformer of varying capacity across increasing number of primitives in the SCAN train set.



Figure 7: Measuring the generalization performance of an LSTM of varying capacity across increasing number of primitives in the Colors train set.

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

primitive ('zup'). Table 2 illustrates the task.

COGS (Kim and Linzen, 2020) is a semantic parsing task of mapping english natural language sentences to their corresponding logical forms. Apart from primitive generalization, COGS also evaluates other types of systematic generalization such generalizing to higher depths or generalizing to novel syntactic structures. The size of the train set is 24,155 and that of the test set is 21,000.

C How model capacity affects compositional generalization

We analyze the relationship between the model capacity and the number of example primitives in the training set. We vary the number of primitives as per the description in Section 2. We evaluate the generalization performance of the models while gradually increasing the number of parameters by increasing the size of its embeddings and intermediate representations. For each experiment, we exhaustively finetune the rest of the hyperparameters (eg. dropout, learning rate, batch size ...etc) to select the best model. We observe a general trend (refer to Fig. 6 and Fig. 7) where models start to overfit and have poor generalization performance

⁸The dataset released by (Lake and Baroni, 2018) is of size 14,670 which has many repetitions of the 'jump \rightarrow JUMP' primitive definition. In this work, we remove all these repetitions since they do not significantly help in generalization.

⁹The original dataset has two additional examples which

evaluate length generalization. Since we focus only on primitive generalization, we do not evaluate on these.

	SCAN		COLORS		COGS	
Hyperparameters	Transformer	LSTM	Transformer	LSTM	Transformer	LSTM
Embedding Size	[64, 128 , 256]	[64, 128 , 256]	[16, 32 , 64]	[16, 32 , 64]	[384 , 512]	[64, 128 , 256]
Hidden/FFN Size	[256 , 512]	[64, 128]	[16, 32 , 64]	[16, 32, 64]	[512 , 1024]	[128, 256 , 512]
Heads	[2 , 4]	N/A	[4 , 8]	N/A	[2 , 4]	N/A
Number of Layers	[2 , 3]	[1, 2]	[2 , 3]	[1, 2]	[2 , 3]	[1, 2]
Learning Rate	[3e-4, 5e-4 , 8e-4]	[5e-3, 8e-3 , 1e-2]	[8e-4 , 1e-3]	[5e-3, 8e-3 , 1e-2]	[3e -4, 5e -4, 8e-4]	[5e-3, 8e-3 , 1e-2]
Batch Size	[128 , 256]	[128, 256]	[1 , 2]	[1 , 2]	[128 , 256]	[128 , 256]
Dropout	[0.1 , 0.2]	[0.1 , 0.2]	[0.1 , 0.2]	[0.1 , 0.2]	[0.1 , 0.2]	[0.1 , 0.2]
Epochs	150	150	150	150	150	150
Avg Time/Epoch	30	40	2	3	60	80

Table 1: Different hyperparameters and the values considered for each of them in the models. The best hyperparameters for each model for all the datasets (with maximum number of primitives of all the settings studied in this paper) are highlighted in bold. Average Time/Epoch is measured in seconds.

TRAIN:	
jump	JUMP
run after run left	LTURN RUN RUN
run	RUN
look left twice and look opposite right	LTURN LOOK LTURN LOOK RTURN RTURN LOOK
Test:	
jump twice after look	LOOK JUMP JUMP
turn left after jump twice	JUMP JUMP LTURN
jump right twice after jump left twice	LTURN JUMP LTURN JUMP RTURN JUMP RTURN JUMP

Table 2:	An illustration	of the	primitive	generalization	task in SCAN.
----------	-----------------	--------	-----------	----------------	---------------

TRAIN:	
dax	RED
lug	BLUE
wif	GREEN
zup	YELLOW
lug fep	BLUE BLUE BLUE
lug blicket wif	BLUE GREEN BLUE
dax kiki lug	BLUE RED
TEST:	
zup fep	YELLOW YELLOW YELLOW
zup blicket lug	YELLOW BLUE YELLOW
zup kiki dax	RED YELLOW

Table 3: An illustration of the primitive generalizationtask in Colors.

as we increase the model size. Note that all these models are able to achieve near-perfect accuracies on the SCAN random split. This shows that carefully controlling the model size is important for achieving compositionally generalization since all these model achieve near-perfect accuracies on random splits. On such small datasets, larger models might simply memorize the input-output mappings in the train set¹⁰. We also find that as we increase the number of example primitives, the models are

473

474

475

476

477

478

479

480

481

482



Figure 8: Decrease in generalization performance on our COGS primitive generalization test set with a decrease in the percentage of example primitives and their use cases present in the train set.

less susceptible to overfitting and achieve relatively better generalization performance.

483

484

485

486

487

488

489

490

491

492

493

494

D Removing Primitives hurts Generalization on COGS

Unlike SCAN and Colors, both of which have a single isolated primitive and only three example primitives, COGS has 3 isolated primitives - a verb, a common noun and a proper noun which are supported by 80 verbs, 40 common nouns and 20 proper nouns as example primitives. We hypothesize that this high number of example primitives might be one of the reasons behind the high perfor-

¹⁰Such memorization has been cited as a potential reason why models fail at compositional generalization (Conklin et al., 2021).

COMPLEXI	TY SENTENCE
1	jump twice
2	jump thrice and look
3	run twice after jump opposite left
4	jump around left and walk opposite left twice

Table 4: Sentences of varying complexities featuring the isolated primitive 'jump'.

mance of Transformers on COGS (Csordás et al., 2021; Ontañón et al., 2021), as far as primitive generalization is concerned.

495

496

497

498

499

502

503

507

508

510

511

512

513

514

515

516

517

518

519

520

521

523

524

525

526

530

To validate our hypothesis, we systematically reduce the number of example primitives in COGS and evaluate the model. The test set of COGS focusing on primitive generalization consists of 5000 examples. If we directly start removing the primitives from the train set, we risk having outof-vocabulary tokens in the test set. Hence we select a portion of the test set of size 1218 which exludes 129 example primitives. We will hold this test set fixed and vary the percentage of the 129 example primitives to be inserted in the train set. For each example primitive, samples are drawn uniformly from the original COGS train set. Note that even though the number of example primitives and their use cases will vary in the train set, we control the total train set size to be always 2500 for fair evaluation.

The results of our experiment can be seen in Fig. 8. We see a clear trend of decrease in generalization performance as we decrease the number example primitives and their use cases. This is in tandem with the results shown in Section 2 and further validates the idea that providing more example primitives and their use cases helps neural sequence models generalize on the primitive generalization task. Our results help explain that the gap in performance of neural sequence models on primitive generalization tasks in COGS and primitive generalization tasks in SCAN or Colors is at least partially caused by the difference in the number of example primitives and their use cases in these datasets.

E Implicit Word Learning

531 Drawing analogy from human vocabulary acquisi-532 tion (Bloom, 2000), our primitive generalization 533 setting corresponds to the case when a child is 534 explicitly explained the meaning of a word. But 535 children can learn word meaning from implicit us-536 age. In our setting this would translate to using



Figure 9: Measuring the similarity of the embedding of *isolated primitive* with the embeddings of example primitives for learned Transformer and LSTM models as we increase the number of example primitives in the Colors train set.

a primitive in a more complex construction, say 'jump twice \rightarrow JUMP JUMP' instead of the original 'jump \rightarrow JUMP'. It would be interesting to evaluate how well seq-to-seq models learn the meanings of words from a single sentence and whether they learn to use that word compositionally with other words. 537

538

539

540

541

542

543

544

545

546

547

548

549

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

We consider the 'add jump' split in SCAN. Instead of providing the 'jump \rightarrow JUMP' primitive definition in the train set, we instead provide one compositional sentence featuring 'jump'. We vary the complexity of this sentence as shown in Table 4. Similar to the case of providing only the primitive definition, we observe that models are unable to generalize and achieve near-zero accuracies.

We now wish to see whether the presence of more number of primitives and their sentences in the train set helps a model generalize in this scenario (like it did for primitive definitions as shown in Section 2). We consider the setup of having 100 primitives and their sentences in the train set (Section 2) apart from the one sentence with the word 'jump'. We find that models are able to achieve near-perfect generalization accuracies.

This shows that our idea holds more generally: Adding more primitives and their sentences helps a model effectively learn the meaning of a new primitive, whether specified explicitly via a primitive definition or implicitly in a sentence.

F Details of Experimental Setups and Other Results

F.1 Embedding of Isolated Primitive

We scale the embedding vectors to unit L2-norm for calculating the euclidean distance and unit L1norm for calculating the manhattan distance. For



Figure 10: Visualizing the *t*-SNE reduced embeddings of isolated primitive (\blacktriangle), example primitives (\blacksquare) and non-primitives (\bullet) from a learned LSTM model as we increase the number of example primitives in the Colors train set.

9



Figure 11: Measuring the generalization performance of LSTM on different types of train set distributions of the SCAN dataset.

Colors dataset as well, we compare the average distance with other primitives before and after adding primitives to the training data. We again find that as we increase the number of example primitives in the training set, the embedding of the isolated primitive ('zup') gets closer to the example primitives (refer to Fig. 9) in terms of Euclidean, Manhattan and Cosine Distances.

We additionally show the t-SNE plots of the learned embeddings for the LSTM model on the Colors dataset (Fig. 10).

F.2 Distributions

In Section 3, we showed results of the Transformer model on various train set distributions of the SCAN dataset. We also experimented with the LSTM model, the results of which can be found in Fig. 11. We see the same trend as we saw for Transformers.

G A Note on Other Data Augmentation Methods

Applying data augmentation methods such as GECA (Andreas, 2020) on SCAN will lead to addition of augmented training examples containing combinations of the isolated primitive 'jump'.

Concurrent to this work, Anonymous (2022) proposed a data augmentation method based on the theory of meaningful learning. Similar to our work, they also augment the train set by adding more primitives (e.g. 'jump_0', 'jump_1', ..., 'jump_n'). However, compared to our work, their setup is completely different: The new primitives that they add to the train set are all still mapped to the output token of an example primitive 'jump' (i.e. 'JUMP'). Their train set has examples showing compositions of 'jump' while their test set evaluates for novel compositions of the newly added primitives. We argue that their setup cannot be considered oneshot primitive generalization since now the model can see the output token 'JUMP' in composition with other words. We claim that this familiarity with the output token enables a model to generalize well on the test data even if the newly added primitives are only presented one-shot in the train set. Indeed, Lake and Baroni (2018) also suggested that the reason why models are able to do well on the 'turn left' split of SCAN is because the train set consists of many examples that have the output token 'LTURN' used compositionally.

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

To validate our claim, we propose a simple experiment. In the original SCAN 'add jump' split, we map 'jump \rightarrow WALK' instead of 'jump \rightarrow JUMP' for all examples (primitive definitions as well as compositional sentences) in both the train and test sets. In this setup, even though the input word 'jump' is seen only once at train time, it's mapping 'WALK' is used compositionally in many examples. On evaluating a Transformer model on this split, we found that it achieves near-perfect accuracies. This shows that providing compositional examples with the output token of the isolated primitive not only breaks systematicity, but is the main reason behind the high performance of models in that setting.

573

574

591

593

594

595