# Generating Demonstrations for In-Context Compositional Generalization in Grounded Language Learning

Anonymous ACL submission

## Abstract

In-Context-learning and few-shot prompting are viable methods compositional output generation. However, these methods can be very sensitive to the choice of support examples used. Retrieving good sup-005 ports from the training data for a given 007 test query is already a difficult problem, but in some cases solving this may not even be enough. We consider the setting of 009 grounded language learning problems where 011 finding relevant supports in the same or similar states as the query may be diffi-013 cult. We design an agent which instead generates possible supports inputs and targets current state of the world, then uses them in-context-learning to solve the test query. We show substantially improved per-017 formance on a previously unsolved compositional generalization test without a loss of performance in other areas. The approach is general and can even scale to instructions expressed in natural language.

## 1 Introduction

026

040

041

042

043

It is thought that a compositional understanding of language and the world (so-called *compositional generalization*). around is something that enables efficient learning in both humans (Chomsky, 1957; Tenenbaum, 2018) and machines (Sodhani et al., 2021; Jang et al., 2021). However, a long line of work and many different datasets show that Deep Learning approaches do not always achieve such compositional generalization. Some solutions to address this deficiency include modular architectures, data augmentation, and sparsity. A recent line of work concerns in-context learning (ICL). Instead of providing a query and asking for the target directly, a few examples of query-target pairs (supports) are also provided. Recent work indicates that supports covering the elements of the query can help enable compositional generalization even if neither shows the desired behaviour exactly (Gupta et al., 2023). A follow up question is how to find examples for each query. Most prior work suggests retrieval from the training data (Pasupat et al., 2021).

However, in the Grounded Language Learning case, retrieval approaches might not be sufficient to make compositional generalization by ICL work well. The expected outputs are conditional not only on the query, but also on the state of the world. Therefore, searching for nearby examples in the input space is problematic. Using the query alone means that it is unlikely that *state-relevant* examples will be retrieved. There might not be query-covering examples in the same state from the training data. Using *similar* states is also challenging because small changes in the state can result in large changes to the target sequence. Searching for nearby examples in the *output space* (Zemlyanskiy et al., 2022) is more promising, but it also relies on being able to find state-relevant covering outputs. It is difficult to make a retrieval-based strategy that works well in all cases.

045

046

047

050

051

052

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

087

Instead of retrieval, we suggest that generation of the supports based on the state might be a better alternative. We contribute the following:

- We confirm that in-context learning is a useful method for unlocking output compositional generalization in the grounded language learning context.
- We show that support selection for in-context learning is a crucial piece of the puzzle and that retrieval from the training set is not enough to get the best performance due to the challenge of the query state being potentially unobserved in the retrieval examples.
- We propose a new method, **DemoGen**, to generate the necessary supports which show different instructions and targets of which the query instruction requires a composition of. Our experiments on gSCAN show that using in-context learning with these supports unlocks superior compositional generalization performance.
- We extend the gSCAN dataset to naturallanguage like instructions to show further that that **DemoGen** method can scale well to natural-language like instructions as well.



Figure 1: Generating demonstrations on gSCAN with **DemoGen** for an ICL Transformer (Figure 2). The **Instruction Generator** takes as input the current state and  $I_q$  and produces similar instructions  $I_1, ... I_n$  likely to occur in the same state, sorted by likelihood (parens). A **Bootstrap Transformer** trained on the training data generates the corresponding actions  $A_1...A_n$  in that state. Some instructions are more helpful than others. Instructions in green,  $I_{1,3,6,8,13,16}$  show both the correct object in  $I_q$  and also either one of the verb or adverb. Instructions in yellow,  $I_5$  show the correct object, but an irrelevant verb and adverb combination. Instructions in red,  $I_{18,21,22}$  show a different object to the target one. Actions in grey  $A_{13,16,18,21,22}$  show an incorrect target sequence. As long as the instructions and actions in green are included in the support set, a sufficiently powerful model can use them and ignore the other supports.

#### 2 Background

093

100

101

102

103

104

105

106

109

110

111

112

113

114

115

116

117

118

119

## 2.1 Compositional Generalization and Grounded Language Learning

Compositional Generalization refers to the ability of a system to learn the rule for how solutions to sub-problems may be combined in some way, then apply the rule to unseen combinations of known sub-problem solutions. It can be seen in both the inductive and productive sense. In the inductive sense, the system must produce some known symbol in response to a unseen combination of known query inputs. In the *productive* sense, the system must produce some unseen combination of known output symbols. The capability of Deep Learning to perform compositional generalization has been studied extensively. Early experiments showed the challenge of doing so on both RNNs (Lake and Baroni, 2018) and Transformers (Hupkes et al., 2020) and many datasets have been created to demonstrate the problem, both with synthetic and "realistic" natural language data (Bastings et al., 2018; Kim and Linzen, 2020; Keysers et al., 2020; Li et al., 2021; Yin et al., 2021; Finegan-Dollak et al., 2018). As more datasets become available, so do approaches to handle the compositional generalization problem. Most approaches generally fall into some combination of data augmentation (Andreas, 2020; Li and McClelland, 2022; Chen et al., 2022b; Qiu et al., 2022; Akyürek et al., 2021), neural module networks (Andreas et al., 2016b; Buch et al., 2021; D'Amario et al., 2021; Ruis and Lake, 2022) and meta-learning (Lake, 2019; Conklin et al., 2021).

The field of Grounded Language Learning is natural fit to study the problems of both inductive and productive compositional generalization. We can test inductive generalization by placing the agent in a state with a novel combination of input sym-

bols. Productive generalization can be tested by giving instructions that require generating some novel combination of outputs conditioned on the state. While the former problem is extensively explored by related work, the latter has received less attention and therefore the focus of this work. 125

126

127

128

129

130

131

133

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

158

159

160

#### 2.2 In-context Learning

Meta-learning and ICL are promising approaches for compositional generalization in sequence generation tasks. In this paradigm, a few *support inputs* and corresponding *support outputs* for a given *query* sequence are provided and the task is to predict the correct *target* sequence (Lake, 2019; Conklin et al., 2021). This has been popularized by the notion of ICL in large language models, where a few examples of the input-output pairs as well as a query are given as part of a *prompt*, then the target is predicted autoregressively (Brown et al., 2020; Min et al., 2022a), which has been shown to enable compositional generalization in sequence generation (Chen et al., 2022a; Logeswaran et al., 2020).

#### 2.3 Support Selection for ICL

ICL methods are sensitive to the choice of support sets used. Mitchell et al. (2021) found that selecting supports that were not relevant to the task at hand degraded performance when using sequence based meta-learning with SCAN. As we also show in our experiments, ICL approachs with a poorly chosen procedure for selecting supports may be worse on all tasks compared to when no ICL is used at all.

Different approaches have been proposed for finding good examples. Many methods try to pick good examples from the training data, for example by using a similarity index (Pasupat et al., 2021), or with a metric that takes into account diversity and local structure coverage (Levy et al., 2022; Gupta



Figure 2: The model architecture for sequence-to-sequence ICL. Each support state  $S_1, ..., S_n$ , support instruction  $I_1, ..., I_n$  and corresponding support targets  $A_1, ..., A_n$ , as well as the query state  $S_q$  and query instruction  $I_q$  are used as inputs to a Transformer Encoder (along with positional encoding). Right-shifted query targets  $a_q^1, ..., a_q^n$  are used as inputs to a Transformer Decoder. Both the support targets and query targets use the same random permutation on every training step.

et al., 2023; Ye et al., 2023). Such retrieval is potentially problematic, because getting relevant output supports requires that the retrieved inputs are evaluated in the same or a very similar state, which can increase the complexity of the retrieval problem.

There are also generative approaches to create the support examples, for example subproblem decomposition (Yang et al., 2022), chain-of-thought (Wei et al., 2022), least-to-most-prompting (Zhou et al., 2022) asking for diverse examples (Yu et al., 2023). These approaches can get very impressive results on ungrounded compositional generalization benchmarks, but they have their own requirements including reliance on information in large language models or special helper prompts about the input structure. A hybrid of generation and retrieval is GandR Zemlyanskiy et al. (2022) which first guesses the output using a helper model and retrieves examples based on output similarity. Our work extends on the generated-example paradigm with the idea of *generating* support instructions for a query state, then *solving* those support instructions using a "bootstrap" model. We explain in Section 3.2 why this is important in the grounded language learning setting.

# 3 Method

161

162

163

164

165

166

168

171

172

173

174

175

176

178

179

180

181

184

185

186

187

188

189

191

192

193

194

196

197

198

In this section, we describe a method we call **De-moGen**. The method is designed to work with datasets where there is both an instruction and a state in the input.

## 3.1 In-Context Learning

ICL can be realized with a large-context encoderdecoder Transformer (see Figure 2). For an initial state  $S_q$  and instruction  $I_q$ , the model is trained to generate a sequence of targets  $A = a_1^Q, ..., a_m^Q$ using a set of support inputs  $I_1, ..., I_n$  and the corresponding support outputs  $A_1, ..., A_n$ .

The entire set of support states  $S_1, ..., S_n$ , support instructions  $I_1, ..., I_n$  and corresponding support targets  $A_1, ..., A_n$ , along with the query state  $S_q$  and query instruction  $I_q$  are passed as one big sequence to the Transformer Encoder, using sinecosine positional encoding in (Vaswani et al., 2017). Right-shifted query targets are used as inputs to the Transformer Decoder with causal masking. 202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

We do not use a pre-trained model and train only on each problem's own training set to eliminate the risk of having pre-trained on the test set. To ensure that we still get in-context learning from the ICL Transformer, we use the technique of permuting the symbol-index mapping of the support and query targets on every training step, similar to Chan et al. (2022).

#### 3.2 Support Set Generation

Choosing the support inputs  $I_1, ..., I_n$  and outputs  $A_1, ..., A_n$  for the ICL model is not a trivial problem. **DemoGen** generates the support sets using two models trained on the training data - an **Instruction Generator** and **Bootstrap Transformer**.

Instruction Generator Support inputs are generated by a BART-like masked language model (Lewis et al., 2020). The model is trained to reconstruct a partially masked sentence. It is trained on a balanced dataset of all the instructions in the training data to ensure that inputs occurring less often have a reasonable chance of being sampled. To generate support inputs, some percentage of the tokens (including padding tokens) in the query  $I_q$ (in this work, 20%) are randomly masked and then the entire input is reconstructed by autoregressive decoding. This process is repeated  $k \ge n$  times, to form  $I_1, \ldots, I_k$ . We deduplicate the samples and remove  $I_q$  from  $I_1, ..., I_k$ . We also *filter* the supports by the use of a *scoring model*. The scoring model estimates probability that a generated support is in-distribution, conditioned on any relevant context. The score is the length-normalized log-likelihood of generated support inputs. We assume that conditionally in-distribution supports are more likely to be solveable by the **Bootstrap Transformer** below. We take the top n by score to get  $I_1, ..., I_n$ .

	SCAN	$\operatorname{COGS}$	ReaSCAN	RTFM	DescWorld	MetaWorld	gSCAN	gSCAN-NL
State-cond.			$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Inductive		$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Productive	$\checkmark$	$\checkmark$					$\checkmark$	$\checkmark$
Natural Lang.		$\checkmark$		$\checkmark$	$\checkmark$			$\checkmark$
Fewshot						$\checkmark$	$\checkmark$	$\checkmark$

Table 1: Comparison of benchmark datasets. Only gSCAN and gSCAN-NL (see Section 4.1) tests productive compositional generalization in a state-conditioned setting and gSCAN-NL extends the input instructions to something more like natural language.

Bootstrap Transformer Support outputs  $A_1, ..., A_n$ are generated from the  $(S_1, I_1), \dots, (S_n, I_n)$ pairs by an Autoregressive Transformer model trained on the same training data. Examples of the generated support inputs and outputs are shown in Figure 1.

243

244

245

246

247

248

251

257

261

264

267

271

272

273

274

275

276

277

279

281

284

287

Generating both the support inputs and outputs 249 has a few interesting advantages. Compared to retrieving on inputs, we can generate examples which we know will be relevant to the current state and also generate examples which might not be found in the training data for a given query. Compared to retrieving based on the predicted output, we 256 can generate a greater diversity of supports which would be valid in the state, as opposed to fetching the same output over and over again in many dif-258 ferent states. The only assumption we make is that the model used to generate the support targets is capable of inductive compositional generalization, but not necessarily productive compositional gener-262 alization. In practice, this is already true with the 263 Transformer architecture (Qiu et al., 2021; Sikarwar 265 et al., 2022). One challenge with generating the supports is that our support generator might come up with support inputs that are either not relevant or not solvable in the current state. We show in 268 the experiments that the presence of irrelevant sup-269 ports is not a problem as long as the other useful supports are also present.

#### 4 Experiments

#### 4.1Dataset

We examine which dataset would be appropriate to evaluate DemoGen on. Since we know that incontext learning helps specifically when it comes to productive compositional generalization, we want the dataset to test this case. We also limit out dataset search to the state-conditioned setting, where it makes sense to generate demonstrations conditioned on the state. To really test our method, we also want a dataset using instructions in the form of natural language as well. The result of the dataset consideration is found in Table 1. We considered well-known compositional generalization and grounded language learning datasets. SCAN (Lake and Baroni, 2018), COGS (Kim and Linzen,

2020), CogNiTIon (Li et al., 2021), CFQ (Keysers et al., 2020) and SMCalFlow-CS (Yin et al., 2021) test productive generalization, but are not stateconditioned. ReaSCAN (Wu et al., 2021), gSCAN-RS (Qiu et al., 2021), RTFM (Zhong et al., 2020), BabyAI (Chevalier-Boisvert et al., 2019), ALFRED (Shridhar et al., 2020) and DescribeWorld (Weir et al., 2023) are state-conditioned but mainly test inductive generalization. MetaWorld (Yu et al., 2019) tests productive generalization, but in the few-shot learning context where examples are already given. gSCAN (Ruis et al., 2020) is the only dataset which tests *productive* generalization in a state-conditioned setting, however it uses very simplistic instructions. Based on this criteria, we choose to evaluate on gSCAN, but also extend it by rewriting the instructions using an LLM to resemble natural language.

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

gSCAN is a Minigrid-based environment with a single training data set and 8 out-of-distribution test splits covering various compositional generalization scenarios. An agent receives an instruction with a target object, a verb to apply to that object and an adverb which affects both navigation and the verb. About 360,000 demonstrations of navigating to various objects and performing some task on them with various adverbs are provided as a training set. A success happens when the agent performs the expected sequence of actions exactly. The input and action vocabularies are small and the instructions constructed using a simple grammar. Typically the instructions follow the form "[verb] a [size] [color] [object] [adverb]", where [size], [color] and [adverb] are sometimes omitted. The in-distribution split is 100% solvable by deep learning.

More challenging are the eight out-of-distribution test splits. Splits B, C, E, F require inductive generalization, for example identifying a "red square" as a goal in split C and a size-3 object being "small" in relation to other objects in split E. Extensions to gSCAN such as ReaSCAN (Wu et al., 2021) and Relational Splits (gSCAN-RS) (Qiu et al., 2021) test further such scenarios. Splits D, G and H require productive generalization at testing-time. Split D requires navigating to an object that is south-west of the agent, which in practice requires the production of LTURN(2) ...  $LTURN^1$ . Split H requires composing a the verb "pull" with the adverb "while spinning", which requires the production of novel fragments LTURN(4) PULL. Split G is a few-shot learning split for a new behaviour "cautiously".

	Parses	Words	<b>Zipf</b> $\alpha$	RMSE
gSCAN	18	18	1.99	0.11
NL-gSCAN	1550	859	1.29	0.01

Table 2: Linguistic properties of the baseline (gSCAN) and paraphrased datasets (NL-gSCAN)

Natural Language Instructions We also ex-339 tend the gSCAN dataset such that the instructions 340 are less formulaic and more like natural language. 341 342 By prompting the openai-gpt3.5 model with 25 different examples of paraphrases for an instruc-343 tion, we can generate paraphrases of all the other instructions in the dataset. To validate that the paraphrased dataset looks more like natural language, we estimate the  $\alpha$  parameter (closer to 1.0) 348 meaning more natural) for a Zipf distribution using maximum likelihood estimation using the method in (Clauset et al., 2009) and also calculate the number of unique parses with spaCy. The paraphrased data has an  $\alpha$  of 1.29 vs 1.99 along with a better fit, and there is a greater diversity of both words (18 vs 859) and syntax parses (18 vs 1550). The target object description was retained in approximately 99% of cases. Examples of paraphrases and further analysis given in Appendix L. Paraphrased 357 instructions are also shown in Figure 1.

Related Work on gSCAN Various approaches to gSCAN including graph networks (Gao et al., 2020), linguistic-assisted attention (Kuo et al., 2021), symbolic reasoning (Nye et al., 2021), auxiliary tasks (Jiang and Bansal, 2021; Hein and Diepold, 2022), modular networks (Heinze-Deml and Bouchacourt, 2020; Ruis and Lake, 2022), logic programming (Yang et al., 2023) and data augmentation (Setzler et al., 2022; Ruis and Lake, 2022) have been proposed. These approaches tend to make some trade-off between performance and generalizability. Transformers have been shown to work well on on the *inductive* category of splits (Qiu et al., 2021) as well as on ReaSCAN and gSCAN-RS (Sikarwar et al., 2022), but there is no general approach which works well on the *productive* category. In this work, we aim to show that an ICL approach along with a support generation strategy that does not assume too much about the problem is a feasible general approach at least for problems like the one in Split H.

#### 4.2 What makes for good supports?

We first explore what sort of supports work well for gSCAN. These methods are based on pre-existing knowledge of the dataset. When we perform experiments with the ICL Transformer, we use the architecture described in Section 3.1 trained to 300,000 steps with batch size 128, hidden size of 512, 8 attention heads, 12 layers and 16 supports per query example. Training was run for 300,000 iterations over 10 seeds. We perform evaluation every 5000 steps on a random subsample of the validation data and the best by split-A (in-distribution) performance are reported. Detailed information on hyperparmeters is given in Appendix C

381

383

384

385

386

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

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

Heuristic Select the best instructions and outputs for a given state which show; 1) going to the same object, 2) showing the target verb in combination with other adverbs, 3) showing the target adverb in combination with other verbs. Note that the generated supports might contain test-set input-output pairs, meaning that we assume extra knowledge not available to the learning agent. The heuristic can be seen as an upper bound on we could expect from an optimal demonstration generator.

**Random Instructions (RD)** Support instructions are selected randomly, without the use of the heuristic described above. Instructions can be about any object in the same state, not just the target one.

**Other States (OS)** We generate the same instructions as in the Heuristic approach but demonstrations are in states different to the query state. Such states are extracted from the training data. The sampled states are also included in the supports and used by the ICL Transformer. If the training data does not contain a state with the same instruction as the one generated by the expert, that instruction is not included in the support set.

Table 4 shows the coverage of the supports over the query according to some hand-written metrics. Heuristic gets full coverage in every category. If we demonstrate random instructions in the same state (RD), only show demonstrations describing the same object 16% of the time (1). If we pick known good instructions for the query demonstrated in different states (OS) then we often describe the correct object, but the outputs look very different to the query, because the starts (2) or finishes (3)in a different position and the agent-target distance is often different (4). This is also reflected in the ICL Transformer performance in Table 3, where demonstrations of relevant instructions in different states show a very wide performance gap and demonstrations in the same state with randomly chosen instructions perform better, but still overall worse than the Heuristic. This supports the idea

338

359

362

363

365

369

370

374

375

378

379

<sup>&</sup>lt;sup>1</sup>In this work, where an action or subsequence is repeated n times, we use the notation (ACT1 ACT2)(n)

	No	ICL	Algori	thmi	с	Retr	ieval	C	eneration	
	$\mathbf{TF}$	$\mathbf{FT}$	Heuristic	$\mathbf{R}\mathbf{D}$	$\mathbf{OS}$	CovR	$\mathbf{GandR}$	DemoGen	DG-NP	DG-NF
Α	1.0	1.0	1.0	0.77	0.99	$0.99\pm.01$	$0.99\pm.01$	$1.0 \pm .01$	$0.94\pm.06$	$0.96\pm.02$
В	1.0	1.0	1.0	0.62	0.0	$0.98\pm.01$	$0.88\pm.05$	$1.0 \pm .01$	$0.92\pm.05$	$0.92\pm.02$
С	0.96	1.0	1.0	0.66	0.2	$0.83\pm.30$	$0.92\pm.03$	$0.98\pm.02$	$0.72\pm.27$	$0.85\pm.03$
D	0.01	0.16	0.50	0.0	0.0	$0.0 \pm .00$	$0.0\pm.00$	$0.03\pm.02$	$0.0\pm.00$	$0.01\pm.01$
Ε	1.0	1.0	1.0	0.59	0.0	$0.99\pm.01$	$0.99\pm.01$	$0.99\pm.01$	$0.92\pm.09$	$0.86\pm.03$
$\mathbf{F}$	1.0	1.0	1.0	0.75	0.99	$0.99\pm.01$	$0.99\pm.01$	$0.99\pm.01$	$0.92\pm.08$	$0.95\pm.01$
G	0.0	0.0	0.0	0.0	0.0	$0.0 \pm .00$	$0.0\pm.00$	$0.0 \pm .00$	$0.0\pm.00$	$0.0\pm.00$
Η	0.22	0.22	0.86	0.15	0.0	$0.56 \pm .10$	$0.17\pm.01$	$0.8\pm.05$	$0.18\pm.02$	$0.62 \pm .2$

Table 3: Success rates on gSCAN for different splits (A–H). Numbers are  $\pm$  standard deviation over 10 seeds, measured after 300,000 steps. Variances are shown only for retrieval and generation experiments and are negligible on other experiments. Heuristics, Retrieval and Generation all use ICL Transformer as the architecture, with supports generated by each method. TF is a Transformer baseline and FT is the same Transformer fine-tuned on generated demonstrations from DemoGen. Best non-oracle results bolded.

	RD	OS	CR	GR	DG
(1) Desc. Obj.	0.16	1.00	0.33	0.68	0.33
(2) Agent Pos.	1.00	0.03	1.00	0.08	1.00
(3) Tgt. Pos.	0.16	0.03	0.39	0.08	0.44
(4) Same Diff.	0.16	0.02	0.39	0.09	0.44
(5) Tgt. Obj.	0.16	0.19	0.27	0.14	0.44
(6) Verb & $(1)$	0.16	0.43	0.88	0.15	1.00
(7) Advb & $(1)$	0.16	0.33	0.78	0.51	0.88
(8) $(6)$ & $(7)$	0.16	0.19	0.70	0.00	0.88
(9) (4) & (8)	0.16	0.00	0.62	0.00	0.88

Table 4: Fraction of supports matching criteria from on each generation method on Split H. Omitted is **Heuristic**, which is 1.0 in every category. (6)-(8) are calculated based on whether any support in a query's supports match that criteria. Other splits are shown in Appendix F

that our support selection procedure needs to find demonstrations that both cover what is requested in the query and also do so in the same state as the query.

#### 4.3 Retrieval vs Generation

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

In the real world, we don't have access to a heuristic function to generate good supports. Instead we have to come up with them using the data we are already given. We can either try to retrieve good supports from the dataset or try to generate them somehow. We compare the following state-of-theart retrieval methods tested on other productive compositional generalization problems and compare them to DemoGen. Further details of implementations are given in Appendix E and D

451 Coverage Retrieval (CR, CovR) Supports are
452 retrieved using a similarity index on states and in453 structions, then chosen based on query coverage
454 similar to Gupta et al. (2023). Instructions are
455 encoded with TF-IDF and states are flattened, one-

hot encoded, then projected along their 320 principal components. The influence of the state and instructions on encoding similarity is balanced by multiplying instruction vectors by the ratio of the state vector norm to the instruction vector norm, contatenating and renormalizing. For each query input and state, we find the 128 nearest neighbours, then rank them descending by their one and twogram coverage. Examples from the retrievals are chosen greedily such that all the one-grams and two-grams in the query are covered maximally. 456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

**GandR (GR)** Supports are retrieved using the Generate-and-Retrieve strategy (Zemlyanskiy et al., 2022). In that strategy, a "helper" model trained on the training data makes an initial guess for the outputs of a given query in a state, even if that query is out of distribution. Then examples for later in-context learning are fetched by similarity of their output sequence to the guessed output sequence. In our implementation, similar to CovR, both query instructions and *outputs* are TF-IDF encoded and retrieved from a similarity index. 128 examples are chosen and similar again to CovR, we greedily pick examples from the 128 nearest output neighbours covering the query input to avoid picking the same (non-covering) instruction many times.

**DemoGen (DG)** Our generation strategy as described in Section 3.2. 2048 instructions are sampled from the language model, deduplicated, and ranked to get the top 16 instructions and corresponding support targets for the query state. A Transformer with the same architecture as given in Section 3.1 is used as the Bootstrap model.

## 4.4 Retrieval Methods vs Generation

The main challenge for retrieval methods is that the supports inputs and outputs for some test splits may not exist in the training data. In gSCAN, we also found that most states don't have close

-	$\mathbf{TF}$	CovR	$\operatorname{Gand} \mathbf{R}$	DemoG
Α	$1.0 \pm .00$	$0.98 \pm .03$	$0.94\pm.01$	$0.99\pm.00$
В	$0.99\pm.00$	$0.93 \pm .08$	$0.92\pm.06$	$0.96\pm.00$
$\mathbf{C}$	$0.99\pm.03$	$0.68 \pm .37$	$0.9\pm.04$	$0.97\pm.00$
D	$0.08\pm.16$	$0.0 \pm .00$	$0.0\pm.00$	$0.01\pm.01$
Е	$0.98\pm.03$	$0.95\pm.08$	$0.89\pm.01$	$0.98\pm.00$
$\mathbf{F}$	$1.0 \pm .00$	$0.88 \pm .11$	$0.94\pm.02$	$0.98\pm.00$
$\mathbf{G}$	$0.0 \pm .00$	$0.0 \pm .00$	$0.0\pm.00$	$0.0\pm.00$
Η	$0.19\pm.03$	$0.44 \pm .23$	$0.17 \pm .00$	$0.59\pm.06$

Table 5: Success rates for a non-ICL Transformer (TF) retrieval baselines and DemoGen on NL-gSCAN. Best results bolded.

near neighbours. An average example's nearest neighbour had a hamming similarity of  $0.74\pm0.107$ (i.e., 10 of 36 cells would be different in the nearest neighbour). Detailed similarity analysis is given further in Appendix A.1. This is also reflected in the properties of what gets retrieved. In Table 4, the distance between the agent and the target object is often different in the query versus the supports (4) and there are fewer demonstrations showing the same exact same target object (5). They also do not always have both (8) the correct verb (6) and adverb (7) in the retrieved supports. On GandR the adverb can significantly change the outputs, so supports with the same verb (but without the adverb) are not selected. For both methods there are even fewer cases where there is least one demonstration of both the correct verb and of the adverb on on the same path (9).

494

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

512 Deficiencies in query coverage aside, these baselines are still stronger on Split H than many previously 513 published results. CovR retrieves examples that 514 515 are very close to the query state like and gets a 516 success rate of 56% on gSCAN Split H and 44% on NL-gSCAN Split H with high variance, However on 517 Split C, CovR loses performance compared to the 518 baseline and has high variance between seeds on both datasets. The other inductive generalization splits on NL-gSCAN also have small but not negligi-521 ble loss compared to a non-ICL Transformer when 522 using CovR to retrieve the supports. GandR gets 17% on Split H, but retains good performance on 524 525 the other splits. However we lose about 10 points on splits B and C compared to the Transformer 526 baseline on both datasets and also about 5 points 527 on Split F of NL-gSCAN.

529 Generating the Supports How does generating 530 the supports with DemoGen compare? In Table 4 531 we see that the generated instructions cover the dif-532 ferent aspects of the instruction and they are made 533 in the same state. This means that the agent start-534 ing position is preserved (2), the path between the 535 starting the target position (between supports and 536 target) is better preserved (4) and, crucially, both the correct verb (6) and adverb (7) are present in the demonstration in combination with the correct object. Demonstrating the right things also has an impact on performance. **DemoGen**, gets 80% on productive generalization Split H for gSCAN and even 59% for the more challenging NL-gSCAN. Performance also remains good on the the inductive generalization splits for both datasets. We provide a summary and detailed comparison to prior work on gSCAN in Appendix B. Aside from (Hein and Diepold, 2022), a specialized architecture with some additional supervision, ours is the best result on Split H.

On Splits D and G, performance on retrieval methods and DemoGen is still not good. The reason is they require generation of a pattern that won't be seen in the outputs in any permutation of the labels. In the case of Split D, it requires LTURN(2) WALK(n) LTURN(1) WALK(n). Only 6% of the data matches this pattern in any index-label permutation. In the case of split G, (LTURN RTURN(3) LTURN WALK)(n) is required. Only 0.0001% of training data matches that up to a permutation. In contrast, Split H requires (LTURN(4) PULL(n)), and there are many examples from the "push a [size] [color] [object]" set of instructions matching that up to a permutation.

Comparing retrieval and generation, we see that retrieval is a good start for finding good supports, but in the state-conditioned setting they still don't live up to what appears to be possible if we selected known good supports using the Heuristic method in Section 4.2. On the other hand, support generation can get very close to the heuristic method and without any prior knowledge of the problem. Generating the supports for the state also doesn't cause a loss of performance on the inductive generalization splits.

4.5 Ablations and Further Questions

	Valid	Correct	C & V	$\mathbf{C} \mid \mathbf{V}$
Α	0.79	0.70	0.70	0.88
В	0.73	0.64	0.64	0.88
$\mathbf{C}$	0.61	0.50	0.50	0.83
D	0.65	0.24	0.24	0.36
Ε	0.78	0.66	0.66	0.84
$\mathbf{F}$	0.73	0.63	0.63	0.87
G	0.79	0.72	0.72	0.91
Η	0.79	0.56	0.56	0.71

Table 6: DemoGen supports: Fraction of valid instructions, correct targets, correct and valid (C & V) and correct given valid (C  $\mid$  V) on synthetic data by split, according to an oracle function.

**Support Quality** Ideally, support should comprise *valid* support inputs (eg, tasks that are actually solveable in a state) and they should be *correct*  537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

564

565

566

567

568

569

570

571

572

573

574

578 enough to facilitate ICL. We investigated this on supports generated by our method and reported the 579 results in Table 6. On average, about 77% of generated support inputs are valid. A support output is 582 *correct* if it matches what an oracle generator would have generated for the corresponding instruction 583 584 and state. 50% of the support pairs were both correct and valid. The number is clearly lower on splits where a Transformer is not able to solve the task well. For example on Split H, there may be "pull an [object] while spinning" in the support inputs, where [object] is not the target object.

590

592

593

594

598

601

610

611

612

613

614

615

618

619

621

622

623

625

Permutations Our ICL Transformer uses a different symbol-index mapping on each training step. On gSCAN, the sequence "WALK(5) RTURN WALK(5)" would be translated into "RTURN(5) LTURN RTURN(5)" under the permutation WALK  $\rightarrow$  RTURN,  $RTURN \rightarrow LTURN$ . One concern is the possibility that a query target with the same symbols for pull ... while spinning is generated after permutation during training, however the probability of this happening is very low. We measured that for a single pass through the training data, approximately 3% of the generated support instructions matched pull ... while spinning, 0.3% of the permuted query outputs matched PULL actions followed by four LTURN instructions, and their intersection was 0.001% of all sampled supports.

Architectural Ablations We also compare the effect of various ablations on gSCAN success rate in Table 3. Fine-Tuning (FT) on the supports generated by DemoGen improves performance marginally on Split D, but not Split H, which shows the importance of using in-context learning for productive generalization. Removing the permuter block (DG-NP) reduces performance to a similar level of not using ICL at all. Remvoing Filtering (DG-NF) reduces average Split C and split H performance drops by about 13 and 20 points respectively with higher variance. We also tried other variants of the Transformer architecture, including RoFormer (Su et al., 2021), Universal Transformer (Dehghani et al., 2019) and Perceiver (Jaegle et al., 2022), which all had similar results compared to a regular Transformer.

Criteria	Success Rate
Remove Same Object Remove Same Adverb Remove Same Verb	$\begin{array}{c} 0.67 \pm 0.17 \\ 0.3 \pm 0.16 \\ 0.21 \pm 0.04 \end{array}$

Table 7: DemoGen Split H success rate when 16 supports are chosen, excluding specified supports.

Ablations on Supports We also examine how important it is to have the right demonstrations at inference time. Figure 3 shows the effect of scaling



Figure 3: Performance of a ICL Transformer trained using Demogen and different numbers of demonstrations at evaluation time on each split. Performance is averaged over 10 different initializations.

626

627

628

629

630

631

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

the number of inference-time demonstrations on performance of DemoGen. With 4 demonstrations and less, exact match performance suffers quite a lot, and the best performance is found with around 12 demonstrations. Additionally, we examine how DemoGen performs on Split H when demonstrations matching certain criteria are removed from the support set. Removing those matching the same object makes a 13 point impact on success rate. Bigger changes come from removing those matching the same adverb (50 points) or verb (59 points). Learning with permutations alone is not enough its also important that the supports cover the types of output behaviour that are found in the target.

# 5 Conclusion

In-Context Learning can help improve performance on challenging compositional generalization problems, but the choice of support examples is crucial to its performance. In the grounded-language learning case, we showed that retrieval may not be enough to get good supports. We demonstrate that generated supports better cover what is required for productive generalization with our support analysis and ablation studies.

We proposed **DemoGen**, a method for sampling support inputs from an autoregressive language model conditioned on the query state, then and solving them using a bootstrap model. When **DemoGen** used with in-context learning, our method outperforms both the best non-retrieval architectures with non-specific architectures and also other strong retrieval based baselines on the challenging Split H of gSCAN, while retaining good performance on other splits. Our method is general and also works well even if the instructions resemble natural language.

## 6 Limitations

662

688

690

704

705 706

710

714

715

716

In this section, we discuss the limitations of our work.

First, on the dataset and evaluation. gSCAN is a
synthetic and with quite simple instructions. We
wanted to evaluate on instructions that were challenging like natural language, but we did not have
the resources to crowdsource annotations for every
data point in gSCAN. Therefore, we relied on commercial large language models to generate similar
instructions instead. These instructions aren't a
substitute for exactly human-generated language,
but they are a good approximation.

675 In this work we decided to dive deep into evaluation 676 on gSCAN instead of evaluating on a broader set of datasets. The main reason for this is first that 677 we are not aware of any other datasets which test 679 the output-sequence level compositional behaviour generalization demanded by for example gSCAN Split H. The second reason is that gSCAN is a diagnostic dataset with output sequence rules which are not noisy and easy to understand for humans. This means that we can more precisely measure the properties of the generated supports and their 686 effectiveness with respect to performance on the problem. 687

> Another limitation of this work is that supports need to be generated at test time for the test set. In this work, we pre-generated the supports for the test set, though a real-time application of this work on unseen examples would need to run the generation process, which could make inference time much longer. There are also other methods to improve the performance of the support input and support output procedure, for example quantization (Dettmers et al., 2022), KV-caching, early stopping, etc.

## 7 Ethics

We used commercial large language models to generate paraphrases of the inputs to test the scalability of our method to natural language data in Section 4.1. These commercial large language models come with their own range of documented ethical issues, such as the capability to amplify harmful biases and misinformation, labour exploitation in training, energy consumption and permission to use web-scale training data. There is also an economic ethical aspect, where the use of the large language model displaces humans who may have been willing to perform the labelling. For our usecase, it was by many orders of magnitude cheaper to use the large language model than crowd-sourced labelling at a fair wage. On the other hand, we believe that there are better uses of human time than paraphrasing hundreds of thousands of examples of simple navigation problems for the purpose of producing a single research paper.

717

718

719

720

721

723

724

725

726

727

728

729

730

731

732

733

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

Our work covers the foundational issue of compositional generalization in grounded language learning, so we don't expect direct applications of it to have the potential to cause social harm. However, the work should be adapted with care. In particular, it is important that the model generating the supports for ICL is actually generating supports which are useful for generating the downstream problem. Generating outputs to a problem with generated wrong input-output pairs is likely to result in even more wrong outputs. Our work shouldn't be deployed in safety critical situations, but instead should be seen as a step towards achieving better data-driven compositional generalization.

# 8 Code and Resources

Our project code can be found at https: //acl-2024-demogen-submission.s3.eu-north-1. amazonaws.com/demogen\_code\_submission.zip.

The paraphrased gSCAN dataset referred to in Section 4.1 can be found at https: //acl-2024-demogen-submission.s3.eu-north-1. amazonaws.com/dataset-paraphrased.txt.

# 9 Computational Resource Usage and Reproducibility Requirements

Experiments were run on our internal GPU cluster. Running a ICL experiment to 300,000 iterations takes about 3 days on a MI250x GPU. For 6 different experiment runs with 10 seeds each, the total compute time is about 330 GPU-days, though the experiments can be run in parallel. The number of GPU-days we used to produce this work was much higher, because of tweaks to the experimental conditions, debugging, restarting failed jobs, etc.

# 10 Bibliographical References

## References

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, ACL 2020, Online, July 5-10, 2020.

Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016a. Learning to compose neural networks for question answering. In NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016.

884

885

Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016b. Neural module networks. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

770

773

774

775

779

781

786

789

790

791

793

795

804

805

807

810

811

812

814

Jasmijn Bastings, Marco Baroni, Jason Weston, Kyunghyun Cho, and Douwe Kiela. 2018. Jump to better conclusions: SCAN both left and right. In Proceedings of the Workshop: Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2018, Brussels, Belgium, November 1, 2018.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are fewshot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.

Shyamal Buch, Li Fei-Fei, and Noah D. Goodman. 2021. Neural event semantics for grounded language understanding. Transactions of the Association for 798 Computational Linguistics, 9.

Stephanie Chan, Adam Santoro, Andrew K. 799 Lampinen, Jane Wang, Aaditya Singh, Pierre H. Richemond, James L. McClelland, and Felix Hill. 2022. Data distributional properties drive emergent in-context learning in transformers. In NeurIPS.

Devendra Singh Chaplot, Kanthashree Mysore Sathyendra, Rama Kumar Pasumarthi, Dheeraj Rajagopal, and Ruslan Salakhutdinov. 2018. Gatedattention architectures for task-oriented language grounding. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018.

Wei-Lin Chen, Cheng-Kuang Wu, and Hsin-816 Hsi Chen. 2023. Self-icl: Zero-shot in-817 context learning with self-generated demonstrations. arXiv:2305.15035.

Yanda Chen, Ruigi Zhong, Sheng Zha, George 819 Karypis, and He He. 2022a. Meta-learning via 820 language model in-context tuning. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022.

825 Zining Chen, Weiqiu Wang, Zhicheng Zhao, Aidong Men, and Hong Chen. 2022b. Bag of tricks for out-827 of-distribution generalization. arXiv:2208.10722.

Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. 2019. BabyAI: A platform to study the sample efficiency of grounded language learning. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.

Noam Chomsky. 1957. Syntactic Structures.

Aaron Clauset, Cosma Rohilla Shalizi, and Mark E. J. Newman. 2009. Power-law distributions in empirical data. SIAM Rev., 51(4).

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, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021.

Vanessa D'Amario, Tomotake Sasaki, and Xavier Boix. 2021. How modular should neural module networks be for systematic generalization? In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual.

Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. 2019. Universal transformers. In International Conference on Learning Representations.

Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P. Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing discrete text prompts with reinforcement learning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, Abu Dhabi, United Arab Emirates, volume abs/2205.12548.

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Llm.int8(): 8-bit matrix multiplication for transformers at scale. arXiv:2208.07339.

Andrew Drozdov, Nathanael Schärli, Ekin Akyürek, Nathan Scales, Xinying Song, Xinyun Chen, Olivier Bousquet, and Denny Zhou. 2022. Compositional semantic parsing with large language models. arXiv:2209.15003.

Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui Zhang, and Dragomir R. Radev. 2018. Improving text-to-sql evaluation methodology. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers.

Tong Gao, Qi Huang, and Raymond J. Mooney. 2020. Systematic generalization on gSCAN with language conditioned embedding. In Proceedings of

944

945

946

the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, AACL/IJCNLP 2020, Suzhou, China, December 4-7, 2020.

886

895

898

900

901

902

903

904

905

906

907

908

924

926

Divyansh Garg, Skanda Vaidyanath, Kuno Kim, Jiaming Song, and Stefano Ermon. 2022. LISA: Learning interpretable skill abstractions from language. In Advances in Neural Information Processing Systems.

Prasoon Goyal, Raymond J. Mooney, and Scott Niekum. 2021. Zero-shot task adaptation using natural language. arXiv:2106.02972.

Shivanshu Gupta, Matt Gardner, and Sameer Singh. 2023. Coverage-based example selection for incontext learning. arXiv:2305.14907.

Alice Hein and Klaus Diepold. 2022. A minimal model for compositional generalization on gscan.
In Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2022, Abu Dhabi, United Arab Emirates (Hybrid), December 8, 2022.

909 Christina Heinze-Deml and Diane Bouchacourt.
910 2020. Think before you act: A simple baseline
911 for compositional generalization. arXiv:2009.13962,
912 2009.13962.

Felix Hill, Andrew K. Lampinen, Rosalia Schneider, Stephen Clark, Matthew Botvinick, James L.
McClelland, and Adam Santoro. 2020. Environmental drivers of systematicity and generalization in a situated agent. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.

920 Matthew Honnibal and Ines Montani. 2017. spaCy
921 2: Natural language understanding with Bloom
922 embeddings, convolutional neural networks and in923 cremental parsing. To appear.

Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. 2020. Compositionality decomposed: How do neural networks generalise? J. Artif. Intell. Res., 67.

928Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste929Alayrac, Carl Doersch, Catalin Ionescu, David930Ding, Skanda Koppula, Daniel Zoran, Andrew931Brock, Evan Shelhamer, Olivier J Henaff, Matthew932Botvinick, Andrew Zisserman, Oriol Vinyals, and933Joao Carreira. 2022. Perceiver IO: A general archi-934tecture for structured inputs & outputs. In Inter-935national Conference on Learning Representations.

936 Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kap937 pler, Frederik Ebert, Corey Lynch, Sergey Levine,
938 and Chelsea Finn. 2021. BC-Z: zero-shot task gen939 eralization with robotic imitation learning. In 5th
940 Annual Conference on Robot Learning, 8-11 Novem941 ber 2021, London, UK.

942 Yichen Jiang and Mohit Bansal. 2021. Inducing943 transformer's compositional generalization ability

via auxiliary sequence prediction tasks. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021.

Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. 2020. Measuring compositional generalization: A comprehensive method on realistic data. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.

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 2020, Online, November 16-20, 2020.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In Proceedings of the Thirty-Sixth Conference on Neural Information Processing Systems, volume 35.

Yen-Ling Kuo, Boris Katz, and Andrei Barbu. 2021. Compositional networks enable systematic generalization for grounded language understanding. In Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021.

Brenden M. Lake. 2019. Compositional generalization through meta sequence-to-sequence learning. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada.

Brenden M. 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, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, volume 80.

Brenden M. Lake, Tal Linzen, and Marco Baroni. 2019. Human few-shot learning of compositional instructions. In Proceedings of the 41th Annual Meeting of the Cognitive Science Society, CogSci 2019: Creativity + Cognition + Computation, Montreal, Canada, July 24-27, 2019.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021.

Itay Levy, Ben Bogin, and Jonathan Berant. 2022. Diverse demonstrations improve in-context compo-

sitional generalization. In 61st Annual Meeting of 1004 the Association of Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Mar-1005 1006 jan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. 1008 BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual 1010 Meeting of the Association for Computational Lin-1011 1012 quistics, ACL 2020, Online, July 5-10, 2020.

Yafu Li, Yongjing Yin, Yulong Chen, and Yue Zhang. 2021. On compositional generalization of 1014 1015 neural machine translation. In Proceedings of the 1016 59th Annual Meeting of the Association for Com-1017 putational Linguistics and the 11th International 1018 Joint Conference on Natural Language Processing, 1019 ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021. 1020

1021

1022

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1050

1055

Yuxuan Li and James McClelland. 2022. Systematic generalization and emergent structures in transformers trained on structured tasks. In *NeurIPS* '22 Workshop on All Things Attention: Bridging 1024 Different Perspectives on Attention.

Lajanugen Logeswaran, Yao Fu, Moontae Lee, and 1026 1027 Honglak Lee. 2022. Few-shot subgoal planning with language models. In Proceedings of the 2022 1028 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, 1031 United States, July 10-15, 2022. 1032

Lajanugen Logeswaran, Ann Lee, Myle Ott, Honglak Lee, Marc'Aurelio Ranzato, and Arthur Szlam. 2020. Few-shot sequence learning with transformers. In Workshop on Meta-Learning, NeurIPS 2020, virtual, volume abs/2012.09543.

Sewon Min, Mike Lewis, Luke Zettlemover, and Hannaneh Hajishirzi. 2022a. Metaicl: Learning to learn in context. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022.

So Yeon Min, Devendra Singh Chaplot, Pradeep Ku-1045 1046 mar Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov. 2022b. FILM: Following instruc-1047 tions in language with modular methods. In Inter-1048 1049 national Conference on Learning Representations.

Eric Mitchell, Chelsea Finn, and Christopher D. Manning. 2021. Challenges of acquiring compositional inductive biases via meta-learning. In AAAI Workshop on Meta-Learning and MetaDL Challenge, MetaDL@AAAI 2021, virtual, February 9, 2021, volume 140.

Maxwell Nye, Michael Henry Tessler, Joshua B. 1056 Tenenbaum, and Brenden M. Lake. 2021. Improv-1057 ing coherence and consistency in neural sequence 1058 models with dual-system, neuro-symbolic reason-1059 1060 ing. In Advances in Neural Information Processing 1061 Systems.

Panupong Pasupat, Yuan Zhang, and Kelvin Guu. 2021. Controllable semantic parsing via retrieval augmentation. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021.

1063

1065

1066

1067

1068

1069

1070

1071

1073

1074

1075

1076

1078

1079

1080

1081

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1098

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

Linlu Qiu, Hexiang Hu, Bowen Zhang, Peter Shaw, and Fei Sha. 2021. Systematic generalization on gSCAN: What is nearly solved and what is next? In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021.

Linlu Qiu, Peter Shaw, Panupong Pasupat, Pawel Krzysztof Nowak, Tal Linzen, Fei Sha, and Kristina Toutanova. 2022. Improving compositional generalization with latent structure and data augmentation. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022.

Laura Ruis, Jacob Andreas, Marco Baroni, Diane Bouchacourt, and Brenden M. Lake. 2020. A benchmark for systematic generalization in grounded language understanding. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.

Laura Ruis and Brenden M. Lake. 2022. Improving systematic generalization through modularity and augmentation. In Proceedings of the 44th Annual Conference of the Cognitive Science Society.

Matthew Setzler, Scott Howland, and Lauren A. Phillips. 2022. Recursive decoding: A situated cognition approach to compositional generation in grounded language understanding. Arxiv:2201.11766.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. ALFRED: A benchmark for interpreting grounded instructions for everyday tasks. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pages 10737–10746. Computer Vision Foundation IEEE.

Ankur Sikarwar, Arkil Patel, and Navin Goyal. 1116 2022. When can transformers ground and com-1117 pose: Insights from compositional generalization 1118 benchmarks. In Proceedings of the 2022 Confer-1119 ence on Empirical Methods in Natural Language
Processing, Abu Dhabi, United Arab Emirates.

1122Shagun Sodhani, Amy Zhang, and Joelle Pineau.11232021. Multi-task reinforcement learning with1124context-based representations. In Proceedings of the112538th International Conference on Machine Learn-1126ing, ICML 2021, 18-24 July 2021, Virtual Event.

1127Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen,1128and Yunfeng Liu. 2021. Roformer: Enhanced1129transformer with rotary position embedding.1130arXiv:2104.09864.

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

Josh Tenenbaum. 2018. Building machines that learn and think like people. In Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS 2018, Stockholm, Sweden, July 10-15, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA.

1143Jason Wei, Xuezhi Wang, Dale Schuurmans,1144Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,1145Quoc V Le, and Denny Zhou. 2022. Chain of1146thought prompting elicits reasoning in large lan-1147guage models. In Advances in Neural Information1148Processing Systems.

Nathaniel Weir, Xingdi Yuan, Marc-Alexandre 1149 Côté, Matthew Hausknecht, Romain Laroche, Ida 1150 Momennejad, Harm Van Seijen, and Benjamin 1151 Van Durme. 2023. One-shot learning from a demon-1152 stration with hierarchical latent language. In Pro-1153 ceedings of the 2023 International Conference on 1154 Autonomous Agents and Multiagent Systems, AA-1155 MAS '23, page 2388–2390, Richland, SC. Inter-1156 national Foundation for Autonomous Agents and 1157 Multiagent Systems. 1158

1159Zhengxuan Wu, Elisa Kreiss, Desmond Ong, and1160Christopher Potts. 2021. ReaSCAN: Compositional1161reasoning in language grounding. In Thirty-fifth1162Conference on Neural Information Processing Sys-1163tems Datasets and Benchmarks Track (Round 1).

1164Jingfeng Yang, Haoming Jiang, Qingyu Yin, Dan-<br/>qing Zhang, Bing Yin, and Diyi Yang. 2022. SE-<br/>QZERO: Few-shot compositional semantic parsing<br/>with sequential prompts and zero-shot models. In<br/>Findings of the Association for Computational Lin-<br/>guistics: NAACL 2022.

1170 Zhang Yang, Adam Ishay, and Joohynung Lee. 2023.
1171 Coupling large language models with logic programming for robust and general reasoning from text.
1173 In Findings of the 61th Annual Meeting of the Association for Computational Linguistics.

1175Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu,1176and Lingpeng Kong. 2023. Compositional exemplars for in-context learning. In International Con-1177plars for in-context learning. In International Con-

ference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research, pages 39818–39833. PMLR.

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1223

Pengcheng Yin, Hao Fang, Graham Neubig, Adam Pauls, Emmanouil Antonios Platanios, Yu Su, Sam Thomson, and Jacob Andreas. 2021. Compositional generalization for neural semantic parsing via spanlevel supervised attention. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021.

Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. 2019. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In 3rd Annual Conference on Robot Learning, CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings, volume 100 of Proceedings of Machine Learning Research, pages 1094–1100. PMLR.

Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. Generate rather than retrieve: Large language models are strong context generators. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* 

Yury Zemlyanskiy, Michiel de Jong, Joshua Ainslie, Panupong Pasupat, Peter Shaw, Linlu Qiu, Sumit Sanghai, and Fei Sha. 2022. Generate-and-retrieve: use your predictions to improve retrieval for semantic parsing.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv:2210.03493*.

Victor Zhong, Tim Rocktäschel, and Edward Grefenstette. 2020. Rtfm: Generalising to new environment dynamics via reading. In *International Conference on Learning Representations*.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed H. Chi. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv:2205.10625*.

# A Details of the gSCAN Dataset

1224

Statistics on the gSCAN dataset are reproduced in Table 8 for the reader's convenience.

	Num. Examples	Length $\pm$ std.
Train	367933	$14.35 \pm 10.07$
А	19282	$13.35\pm8.87$
В	18718	$13.95\pm9.72$
$\mathbf{C}$	37436	$14.07\pm9.78$
D	88642	$17.96 \pm 10.78$
Ε	16808	$13.31\pm9.48$
$\mathbf{F}$	11460	$16.50 \pm 12.40$
G	112880	$33.46 \pm 16.90$
Н	38582	$43.07\pm19.67$

Table 8: Statistics on the gSCAN dataset and test splits

1226

1227

1228

1229

1230

1231

1232

1233

## A.1 Nearest Neighbour Similarity Distribution

	1	2	4	8	16	32	64	128	256	512	1024	2048	4096	8192
$\operatorname{split}$														
a	0.743	0.740	0.737	0.734	0.731	0.725	0.719	0.712	0.702	0.692	0.684	0.681	0.670	0.663
b	0.743	0.741	0.738	0.735	0.731	0.726	0.719	0.710	0.700	0.692	0.686	0.681	0.667	0.661
с	0.759	0.756	0.754	0.750	0.747	0.742	0.738	0.731	0.722	0.710	0.702	0.699	0.690	0.682
d	0.753	0.750	0.748	0.745	0.741	0.737	0.732	0.725	0.714	0.707	0.702	0.699	0.689	0.679
$\operatorname{dev}$	0.767	0.764	0.761	0.757	0.754	0.747	0.742	0.736	0.727	0.715	0.705	0.702	0.693	0.686
е	0.668	0.668	0.668	0.666	0.663	0.659	0.655	0.649	0.638	0.626	0.617	0.614	0.606	0.599
f	0.748	0.746	0.744	0.740	0.737	0.731	0.726	0.721	0.712	0.698	0.689	0.687	0.678	0.670
g	0.767	0.764	0.761	0.757	0.754	0.747	0.741	0.735	0.726	0.714	0.705	0.702	0.693	0.685
h	0.767	0.763	0.760	0.757	0.753	0.747	0.741	0.735	0.726	0.713	0.704	0.702	0.692	0.685
$\operatorname{train}$	1.000	0.767	0.762	0.758	0.754	0.747	0.741	0.735	0.726	0.714	0.705	0.702	0.693	0.685

Table 9: Average state nearest neighbour similarity (between the shown split and the training split) for each split. X-axis is log-scale. The plots show the average hamming similarity between points in a split and their Nth nearest neighbour in the training split.

We visualize the average nth training-data nearest neighbour similarity distribution for each dataset split in Figure 4. We created the figure by taking 1000 random examples from each splits, then finding their 8192 nearest neighbours using a inner-product index over normalized one-hot encoded state representations.

In most cases, even the closest nearest neighbour state has quite many differences and these differences grow as we pick nearest neighbours further away from a training data point. This means that it is hard to find an example in the training set containing different instructions in the exact same environment layout.

# **B** Additional Comparisons

	seq2seq	GECA	FiLM	RelNet	LCGN	ViLBERT
	(Ruis et al., 2020)	(Ruis et al., 2020)	(Qiu et al., 2021)	(Qiu et al., 2021)	(Gao et al., 2020)	(Qiu et al., 2021)
Α	$97.15 \pm .46$	$87.6 \pm 1.19$	$98.83 \pm .32$	$97.38 \pm .33$	$98.6 \pm .9$	$99.95 \pm .02$
В	$30.05 \pm 26.76$	$34.92 \pm 39.30$	$94.04 \pm 7.41$	$49.44 \pm 8.19$	$99.08 \pm .69$	$99.90 \pm .06$
$\mathbf{C}$	$29.79 \pm 17.70$	$78.77 \pm 6.63$	$60.12 \pm 8.81$	$19.92 \pm 9.84$	$80.31 \pm 24.51$	$99.25 \pm .91$
D	$0.00 \pm .00$	$0.00 \pm .00$	$0.00 \pm .00$	$0.00 \pm .00$	$0.16 \pm .12$	$0.00 \pm .00$
Е	$37.25 \pm 2.85$	$33.19 \pm 3.69$	$31.64 \pm 1.04$	$42.17 \pm 6.22$	$87.32 \pm 27.38$	$99.02 \pm 1.16$
F	$94.16 \pm 1.25$	$85.99 \pm .85$	$86.45 \pm 6.67$	$96.59 \pm .94$	$99.33 \pm .46$	$99.98 \pm .01$
Η	$19.04 \pm 4.08$	$11.83 \pm .31$	$11.71 \pm 2.34$	$18.26 \pm 1.24$	$33.6 \pm 20.81$	$22.16 \pm .01$
	GroCoT	Planning	RD Random/RL	Modular	CMA-ES	Role-Guided
	(Sikarwar et al., 2022)	2020	(Setzler et al., 2022)	(Ruis and Lake, 2022)	(Hein and Diepold, 2022)	(Kuo et al., 2021)
А	99.9	$94.19 \pm .71$	$98.39 \pm .17$	$96.34 \pm .28$	$99.7 \pm .1$	$96.73 \pm .58$
В	99.9	$87.31 \pm 4.38$	$62.19 \pm 24.08$	$59.66 \pm 23.76$	$73.5 \pm 25.4$	$94.91 \pm 1.30$
$\mathbf{C}$	99.9	$81.07 \pm 10.12$	$56.52 \pm 29.70$	$32.09 \pm 9.79$	$99.4 \pm .4$	$67.72 \pm 10.83$
D	0.0		$43.60 \pm 6.05$	$0.00 \pm .00$	$2.2 \pm 1.5$	$11.52 \pm 8.18$
Е	99.8	$52.8 \pm 9.96$	$53.89 \pm 5.39$	$49.34 \pm 11.60$	$97.4 \pm 2.0$	$76.83 \pm 2.32$
F	99.9		$95.74 \pm .75$	$94.16 \pm 1.25$	$99.1 \pm .6$	$98.67 \pm .05$
Η	22.9		$21.95 \pm .03$	$76.84 \pm 26.94$	$98.4 \pm 1.1$	$20.98 \pm 1.98$

Table 10: Additional related work comparisons. Splits G and I are not included.



Figure 4: Average state nearest neighbour similarity (between the shown split and the training split) for each split. X-axis is log-scale. The plots show the average hamming similarity between points in a split and their Nth nearest neighbour in the training split.

In this section of the appendix, we describe in more detail other related work on gSCAN and provide the results reported by those works in Table 10 for easier comparison with our experimental results.

**Modular** A recent work by Ruis and Lake (2022). It uses a specialized decomposition into Perception, Interaction, Navigation and Transformation Modules, each of which are trained independently with their own training outputs, then connected together at test time. The modular decomposition gives a prior on how the problem should be solved (for example by decomposition into egocentric and allocentric plans). The work also describes how data augmentation can be used to improve the model, but we show the results coming from use of the modular architecture alone. This approach can get good performance on Splits G and H. Performance on other splits is either slightly improved or comparable to the baseline in Ruis et al. (2020), which is likely due to the use of a similar underlying architecture of RNNs and CNNs as feature encoders.

**Role-Guided** (Kuo et al., 2021) This approach uses linguistic priors to decompose the parsing problem and specify how sub-parsers are connected. It can achieve some level of performance on Split D and comparable performance on Split H to the Transformer.

ViLBERTis an adaptation of the ViLBERT model for gSCAN by Qiu et al. (2021) and extended on1248by Sikarwar et al. (2022). The state is first one-hot encoded, a few 2D convolution layers are applied to it.1249The state is then flattened and the channel values for each pixel are treated as vectors for each location in1250the state. Afterwards, there are several layers of cross-attention between the instruction tokens and the1251state tokens. The cross-attented representations are concatenated together and used as input to a causal1252Transformer decoder to decode the outputs.1253

1234 1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

GECA Also known as "Good Enough Compositional Augmentation" (Andreas (2020)), applied to 1254 gSCAN by Ruis et al. (2020). GECA is an augmentation method which recognizes template fragments in 1255 text, then realizes those templates with other possible substitutions. Following the example in that work, 1256 if a dataset contains "she picks the wug up in Fresno" and "she puts the wug <u>down</u> in Tempe", then the 1257 augmentation method generates samples of puts down substituted into sentences containing picks up. For 1258 example the sentence "Pat picks cats up" can be augmented to "Pat puts cats down". GECA relies on 1259 being able to identify templates containing discontiguous fragments which contain at least two tokens. 1260 In the case of SCAN, GECA might identify the fragment "jump ... JUMP ... JUMP ... JUMP" from the 1261 concatenated instruction-action pair "jump thrice JUMP JUMP" and substitute it into "walk around 1262 right thrice WALK RTURN WALK RTURN WALK RTURN" such that it is augmented into "jump around right thrice 1263 JUMP RTURN JUMP RTURN JUMP RTURN". As noted by Andreas (2020), the time and space complexity of GECA can be quite large and scales with the number of recognized templates and fragments. The results reported by Ruis et al. (2020) when using GECA in Table 10 are possibly out of date, since they were 1267 generated using an RNN architecture as opposed to a Transformer, where better performance on Splits B, C, E and F has been observed. Also, GECA was only applied to the instructions and state and not 1269 to the target commands. Possibly the reason for this is that the computational and memory complexity 1270 of GECA makes it difficult to apply the joint space of the state, instruction and target commands as in gSCAN. 1271

**CMA-ES** uses sparse hard attention with CMA-ES as the optimization algorithm as opposed to a gradient-based optimizer. The model architecture consists only of a multi-layer perceptron, predicting the next token with attention over the generated output sequence. The method requires some supervision on what the goal object is, in contrast with other approaches. Its strengths are that convergence can happen very quickly and optimization can be run on lighter hardware. The method also gets very good performance on Split H, however, as of the time of writing, the authors have not yet published their code and did not provide any analysis in their paper as to why the measured Split H performance was so good, so it is not possible to make a detailed comparison with our work.

	ViLBERT (Qiu et al., 2021)	Modular (Ruis and Lake, 2022)	Role-guided (Kuo et al., 2021)	Transformer (ours) Ours	ICL Transformer Ours
Learning Rate	0.0015	0.001	0.001	0.0001	0.0001
Embedding Dim	128	128	128	512	512
Dropout	0.1	-	-	0.1	0.1
Batch Size	128	200	200	128	128
Steps	114.96K	73K	150K	300K	300K
#params	3M			88.3M	88.3M

Table 11: Hyperparameters used in our experiments and the related work

## C Experimental Details

1272

1273

1274

1275

1276 1277

1280

1282

1283

1285

1286

1287

1288

We ran experiments to determine the performance of our approach. The Transformer blocks use an embedding size  $(d_{\text{model}})$  of 512 units and fully-connected layer size  $(d_{\text{FF}})$  of 2048 units is used. We use 12 layers for each of the encoder and decoder of the encoder-decoder transformer. The learning rate is  $10^{-5}$ , we have an effective batch size of 128, and training iteration count of 300,000. During training, dropout is not used and weight decay is set to  $10^{-3}$  with the AdamW optimizer. Beta values are left at their defaults,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Learning rate warmup is used up to step 30,000 to a peak learning rate of  $10^{-5}$ , then decayed on a log-linear schedule from steps 30,000 to 300,000 to  $10^{-6}$ . Gradient norms are clipped at 0.2 to improve training stability. We use 16-bit precision during training and make use of gradient accumulation in order to simulate large batch sizes where memory is limited.

# D Implementation of GandR for gSCAN

We make small adaptations to GandR (Zemlyanskiy et al., 2022) to adapt it to the grounded setting. The baseline transformer model makes an initial prediction for the query input, then the query input and prediction are vector-encoded and used to find other similar query-output pairs using the index, which become the support inputs and outputs used for ICL. Compared to the original, we keep the  $\alpha$  trade-off between input and target components fixed as opposed to varying it. We also don't include the state in the vector though the identity of the target object and also its distance to the agent will likely be similar as we select on the basis of input and output similarity. There is also nothing to ensure that a diversity of different instructions is sampled - only the near neighbours are sampled, even if they all correspond to a single instruction.

#### Implementation of CovR for gSCAN $\mathbf{E}$

We implement the main idea behind Set-BSR (Gupta et al., 2023) for the grounded setting. States are vector-encoded and projected using PCA into 320 dimensions. Instructions are TF-IDF encoded into vectors. Both are concatenated with each other to make a vector representation of an example. The instruction component of the vector is weighted with  $\alpha = 0.125$ . The training-set vectors are placed into an inner-product index. For performance reasons, we use a Voronoi index with 512 cells and 10 cell probes per search. For each vector in a split, we search the index for the 128 nearest neighbours, sort the neighbours in descending order according to the number of matching two-grams, one-grams and the cosine similarity to the query state. Then we pick the top k = 16 examples as the support set.

#### $\mathbf{F}$ Properties of Generated Demonstrations, other splits

Properties of Generated Demonstrations for the other splits are shown in tables below.

		Split A				
	DemoG	GandR	CovR	Expert	OS	RD
(1) Desc. Obj.	0.32	0.83	0.15	1.00	1.00	0.07
(2) Agent Pos.	1.00	0.07	1.00	1.00	0.03	1.00
(3) Tgt. Pos.	0.37	0.08	0.27	1.00	0.03	0.07
(4) Same Diff.	0.37	0.31	0.27	1.00	0.02	0.07
(5) Tgt. Obj.	0.37	0.26	0.22	1.00	0.25	0.07
(6) Verb & (5)	1.00	0.93	0.91	1.00	0.50	0.07
(7) Advb & (5)	0.75	0.93	0.77	1.00	0.38	0.07
(8) $(6)$ & $(7)$	0.75	0.93	0.73	1.00	0.23	0.07
(9) (4) & (8)	0.75	0.57	0.65	1.00	0.00	0.07
(0) (1) (0)		<u>a 11 a</u>	0.000			
		Split C	~ ~		~ ~	
	DemoG	GandR	CovR	Expert	OS	RD
<ol><li>Desc. Obj.</li></ol>	0.16	0.47	0.15	1.00	1.00	0.15
(2) Agent Pos.	1.00	0.12	1.00	1.00	0.03	1.00
(3) Tgt. Pos.	0.19	0.13	0.18	1.00	0.03	0.15
(4) Same Diff.	0.19	0.44	0.18	1.00	0.02	0.15
(5) Tgt. Obj.	0.19	0.00	0.00	1.00	0.00	0.15
(6) Verb & (5)	0.79	0.00	0.00	1.00	0.00	0.15
(7) Advb & (5)	0.41	0.00	0.00	1.00	0.00	0.15
(8) (6) & (7)	0.40	0.00	0.00	1.00	0.00	0.15
(9) (4) & (8)	0.40	0.00	0.00	1.00	0.00	0.15
		Split E				
	DomoC	CandD	Com	Ermont	08	DD
	DemoG	Gallun	Covit	Expert	05	пD
<ol><li>Desc. Obj.</li></ol>	0.22	0.89	0.07	1.00	0.00	0.00
(2) Agent Pos.	1.00	0.11	1.00	1.00	0.00	1.00
(3) Tgt. Pos.	0.27	0.12	0.22	1.00	0.00	0.00
(4) Same Diff.	0.27	0.35	0.22	1.00	0.00	0.00
(5) Tgt. Obj.	0.27	0.03	0.14	1.00	0.00	0.00
(6) Verb & (5)	0.96	0.20	0.81	1.00	0.00	0.00
(7) Advb & (5)	0.50	0.20	0.63	1.00	0.00	0.00
(8) (6) & (7)	0.50	0.20	0.60	1.00	0.00	0.00
(9) (4) & (8)	0.50	0.14	0.50	1.00	0.00	0.00
		Split G				
	DemoG	GandR	CovR	Expert	OS	RD
(1) Desc. Obi.	0.39	0.91	0.31	1.00	1.00	0.20
(2) Agent Pos	1.00	0.14	1.00	1.00	0.03	1.00
(3) Tet Pos	0.50	0.14	0.37	1.00	0.03	0.20
(4) Same Diff	0.50	0.10	0.37	1.00	0.03	0.20
(5) Tat Obj	0.50	0.30	0.94	1.00	0.02	0.20
(o) igi. Obj.			11 / 44			11.40
(6) Worb $f_r(5)$	1.00	0.22	0.021	1.00	0.51	0.20
(6) Verb & (5)	1.00	0.91	0.93	1.00	0.51	0.20
(6) Verb & (5) (7) Advb & (5) (8) (6) (- (7))	1.00	0.22 0.91 0.01	0.93	1.00	0.51	0.20
6) Verb & (5) 7) Advb & (5) 8) (6) & (7) 9) (4) % (9)	1.00 0.00 0.00	0.22 0.91 0.01 0.01	0.93 0.00 0.00	1.00 1.00 1.00	0.51 0.00 0.00	0.20 0.20 0.20

<ol><li>Desc. Obj.</li></ol>	0.26	0.00	0.00	1.00	0.00	0.00
(2) Agent Pos.	1.00	0.13	1.00	1.00	0.00	1.00
(3) Tgt. Pos.	0.32	0.15	0.29	1.00	0.00	0.00
(4) Same Diff.	0.32	0.44	0.29	1.00	0.00	0.00
(5) Tgt. Obj.	0.32	0.03	0.18	1.00	0.00	0.00
(6) Verb & (5)	1.00	0.30	0.85	1.00	0.00	0.00
(7) Advb & (5)	0.66	0.30	0.71	1.00	0.00	0.00
(8) (6) & (7)	0.66	0.30	0.69	1.00	0.00	0.00
(9) (4) & (8)	0.66	0.24	0.63	1.00	0.00	0.00
		Split D				
	DemoG	GandR	$\operatorname{CovR}$	Expert	OS	RD
(1) Desc. Obj.	0.19	0.83	0.18	1.00	1.00	0.16
(2) Agent Pos.	1.00	0.03	1.00	1.00	0.02	1.00
(3) Tgt. Pos.	0.33	0.03	0.00	1.00	0.02	0.16
(4) Same Diff.	0.33	0.00	0.00	1.00	0.00	0.16
(5) Tgt. Obj.	0.33	0.20	0.05	1.00	0.10	0.16
(6) Verb & (5)	0.99	0.89	0.42	1.00	0.25	0.16
(7) Advb & (5)	0.89	0.88	0.25	1.00	0.17	0.16
(8) (6) & (7)	0.89	0.88	0.20	1.00	0.06	0.16
(9) $(4)$ & $(8)$	0.89	0.00	0.00	1.00	0.00	0.16
		Split F				
	DemoG	GandR	$\operatorname{CovR}$	Expert	OS	RD
<ol><li>Desc. Obj.</li></ol>	0.26	0.81	0.23	1.00	1.00	0.15
<ol><li>Agent Pos.</li></ol>	1.00	0.12	1.00	1.00	0.03	1.00
(3) Tgt. Pos.	0.33	0.15	0.26	1.00	0.03	0.15
(4) Same Diff.	0.33	0.37	0.26	1.00	0.02	0.15
(5) Tgt. Obj.	0.33	0.00	0.10	1.00	0.07	0.15
(6) Verb & (5)	0.96	0.00	0.00	1.00	0.00	0.15
(7) Advb & (5)	0.60	0.00	0.62	1.00	0.29	0.15
(8) (6) & (7)	0.58	0.00	0.00	1.00	0.00	0.15
(9) (4) & (8)	0.58	0.00	0.00	1.00	0.00	0.15
		Split H				
	DemoG	GandR	CovR	Expert	OS	RD
(1) Desc. Obj.	0.33	0.68	0.33	1.00	1.00	0.16
(2) Agent Pos.	1.00	0.08	1.00	1.00	0.03	1.00
(3) Tgt. Pos.	0.44	0.08	0.39	1.00	0.03	0.16
(4) Same Diff.	0.44	0.09	0.39	1.00	0.02	0.16
(5) Tgt. Obj.	0.44	0.14	0.27	1.00	0.19	0.16
(6) Verb & (5)	1.00	0.15	0.88	1.00	0.43	0.16
(7) Advb & (5)	0.88	0.51	0.78	1.00	0.33	0.16
(8) (6) & (7)	0.88	0.00	0.70	1.00	0.19	0.16
(9) $(4)$ & $(8)$	0.88	0.00	0.62	1.00	0.00	0.16

 $\operatorname{CovR}$ 

Expert

OS

RD

#### G **Heuristic Function**

The **Heuristic** function generates relevant instructions by the use of a templating mechanism, which replaces verbs and adverbs in the sentence with other verbs and adverbs, such that the whole combination is still in distribution, but not the same as the query instruction. The rules of the system are:

- Replace "pull" with "push" and "walk to"
- Replace "walk to" with "push" and "pull" (but not if "while spinning" is the adverb)

1311

1312 1313

1315

1316

1309 1310

1300

1301 1302

1303

1304

1306

1307

1308

Word	Symbol	Action	Symbol
ʻa'	0	PULL	0
'big'	1	PUSH	1
'blue'	2	STAY	2
'cautiously'	3	LTURN	3
'circle'	4	RTURN	4
'cylinder'	5	WALK	5
'green'	6		
'hesitantly'	7		
'pull'	8		
ʻpush	9		
'red'	10		
'small'	11		
'square'	12		
'to'	13		
'walk'	14		
'while spinning'	15		
'while zigzagging'	16		

Table 12: Default mapping of words and actions to symbols

- Replace "push" with "walk to" and "pull" (but not if "while spinning" is the adverb)
  - Replace "while zigzagging" with "hesitantly", nothing and "while spinning" (but not if "push" is the verb)
  - Replace "hesitantly" with "while zigzagging", nothing and "while spinning" (but not if "push" is the verb)
    - Replace "while spinning" with "hesitantly", "while zigzagging" and nothing

Examples of what the oracle function generates for a given query instruction and environment can be found in Figure 10. Actions are generated by using the same procedure provided in Ruis et al. (2020). The instruction generated by the oracle is given to the demonstration generation procedure and a demonstration is generated by that. A demonstration can also be generated by providing the oracle-generated instruction and current state representation as input to a Transformer model trained on the provided training set.

# H Permuter Blocks

1318 1319

1320

1322

1323

1324

1325

1326 1327

1328

1329

1330

1332

The permuter block shuffles the indices mapping words to symbols in the dictionary given in Table 12. Table 13 gives an example of how the permuted sequences might look to the encoders. Essentially the individual symbols no longer hold any special meaning without reference to the demonstrations, only conditional autoregressive probabilities up to a permutation hold meaning.

# 1333 I Natural-ish Language gSCAN Dataset

The dataset is generated by extracting all of the input sentences from gSCAN and its derivatives, then using the commercial gpt3.5-turbo model from OpenAI<sup>2</sup> to generate additional paraphrases of the input sentence. The paraphrases are generated by creating four dataset specific prompts, each with an 10 examples of how one instruction in the dataset may be paraphrased, then requesting 25 additional paraphrases for a different instruction in the same dataset to be completed by the language model. The prompts are given in Appendix J. The prompts modes are described as follows:

- **Simple** Paraphrases of "Push a red square"
- **Adverb** Paraphrases of "Push a red square cautiously"
- **Relational** Paraphrases of "Push a red circle that is south east of a blue circle"

1343 **ReaSCAN** Paraphrases of "Pull the yellow square that is inside of a big red box and in the same row 1344 as a small red circle and in the same column as a small cylinder while spinning"

 $<sup>^{2}</sup>$ As of 5 May 2023

Original actions	Permutation	Encoded actions	Permuted encoding
WALK(5) RTURN WALK(5)	$PULL(0) \rightarrow 0, PUSH(1) \rightarrow 5, STAY(2) \rightarrow$	$5(5) \ 4 \ 5(5)$	$4(5) \ 3 \ 4(5)$
	2, LTURN(3) $\rightarrow$ 1, RTURN(4) $\rightarrow$ 3,		
	$WALK(5) \rightarrow 4,$		
RTURN WALK(3)	$PULL(0) \rightarrow 0, PUSH(1) \rightarrow 2, STAY(2) \rightarrow$	4 5(3)	4 1(3)
	3, LTURN(3) $\rightarrow$ 5, RTURN(4) $\rightarrow$ 4,		
	$WALK(5) \rightarrow 1,$		
LTURN(4) WALK LTURN(4)	$PULL(0) \rightarrow 4$ , $PUSH(1) \rightarrow 5$ , $STAY(2) \rightarrow$	3(4) 5 3(4) 5 3(5) 5	2(4) 1 2(4) 1 2(5) 1
WALK LTURN(5) WALK	$0, \text{ LTURN}(3) \rightarrow 2, \text{ RTURN}(4) \rightarrow 3,$	3(4) 5 3(4) 5 3(4) 5	2(4) 1 2(4) 1 2(4) 1
LTURN(4) WALK LTURN(4)	$WALK(5) \rightarrow 1,$	3(4) 5	2(4) 1
WALK LTURN(4) WALK			
LTURN(4) WALK			
LTURN WALK STAY WALK	$PULL(0) \rightarrow 3, PUSH(1) \rightarrow 0, STAY(2) \rightarrow$	$3\ 5\ 2\ 5\ 2\ 5\ 2\ 5\ 2$	$5\ 1\ 2\ 1\ 2\ 1\ 2\ 1\ 2$
STAY WALK STAY WALK	2, LTURN(3) $\rightarrow$ 5, RTURN(4) $\rightarrow$ 4,		
STAY	$WALK(5) \rightarrow 1,$		
LTURN WALK STAY WALK	$PULL(0) \rightarrow 0, PUSH(1) \rightarrow 3, STAY(2) \rightarrow$	$3\ 5\ 2\ 5\ 2$	51414
STAY	4, LTURN(3) $\rightarrow$ 5, RTURN(4) $\rightarrow$ 2,		
	$WALK(5) \rightarrow 1,$		
LTURN(4) WALK LTURN(4)	$PULL(0) \rightarrow 0, PUSH(1) \rightarrow 4, STAY(2) \rightarrow$	3(4) 5 3(4) 5 3(4) 5	$1(4) \ 2 \ 1(4) \ 2 \ 1(4) \ 2$
WALK LTURN(4) WALK	5, LTURN(3) $\rightarrow$ 1, RTURN(4) $\rightarrow$ 3,	$3(4) \ 4 \ 5 \ 3(4) \ 5 \ 3(4)$	$1(4) \ 3 \ 2 \ 1(4) \ 2 \ 1(4)$
LTURN(4) RTURN WALK	$WALK(5) \rightarrow 2,$	$5\ 3(4)\ 5\ 3(4)\ 5$	2 1(4) 2 1(4) 2
LTURN(4) WALK LTURN(4)			
WALK LTURN(4) WALK			
LTURN(4) WALK			
LTURN WALK(2) PUSH	$PULL(0) \rightarrow 1, PUSH(1) \rightarrow 0, STAY(2) \rightarrow$	$3\ 5(2)\ 1$	3 2(2) 0
	5, LTURN(3) $\rightarrow$ 3, RTURN(4) $\rightarrow$ 4,		
	$WALK(5) \rightarrow 2,$		

Table 13: Actions and possible mapping permutations generated by the permuter block.

	gSCAN	$\mathbf{RS}$	ReaSCAN
Uniq. Instrs.	430	31799	4381
Uniq. Tmpls.	-	21	658
Gen. Instrs.	12778	731377	99698
Gen. Tmpls.	-	483	14683
Prompt	Simple	Relational	ReaSCAN

Table 14: Generation properties and configuration for each of the datasets

The 10 paraphrase examples were written by ourselves - the idea is that they show how adverbs and 1345 actions can be replaced by synonyms, and also show examples of the same instruction in a different 1346 sentence ordering. For example, "push a red square" can be paraphrased as "shove the red square" or 1347 "Walk to a red square and push it". The paraphrases can also include additional verbs adverbs which are distractors, for example "grasp a red square and move it along". 1349

We generate paraphrases of instructions in gSCAN, gSCAN-RS and ReaSCAN. The default generation 1350 mode creates paraphrases for each unique instruction individually. However for gSCAN-RS and ReaSCAN, 1351 the number of unique instructions is very large, which would mean that generation would come at both a 1352 high time and monetary cost. The reason for this is the combinatorial explosion of factors; in gSCAN-RS 1353 the instructions are given as target objects with positions relative to other objects. To address this 1354 problem, we also have a "template" generation mode, which replaces the object descriptions (size, color, 1355 shape) with a template placeholder, generates paraphrases for the templates, then generates "realised" 1356 paraphrases from those templates for each of the corresponding object descriptions. This reduces the 1357 number of requests to the model from hundreds of thousands to thousands. 1358

J Prompts used for GPT3.5	1359
J.1 gSCAN Simple Prompt	1360
Here are 10 similar statements to "push a red square"	1361
1. Push the red square	1362
2. Move a red square	1363
3. Shove the red square	1364

1365	4. Go to the red square and shove it
1366	5. Go to the red square and push it
1367	6. Walk to the red square and push it
1368	7. Find a red square and push it
1369	8. Locate a red square and push it
1370	9. Get to the red square and move it along
1371	10. Walk up to the red square and then really push it
1372	Can you generate 25 similar statements for " $\{\{QUERY\}\}$ " in English?
1373	K Examples of Generated Paraphrases
1374	Examples: 367933 Unique Instructions: 430 Templates: 828 Sample Responses:
1375	<b>Original</b> walk to a yellow small cylinder hesitantly
1376	1. Walk hesitantly to a small yellow cylinder
1377	2. Approach a small yellow cylinder with hesitation
1378	3. Take cautious steps towards a small yellow cylinder
1379	4. Walk slowly towards a small yellow cylinder
1380	5. Move towards a small yellow cylinder with caution
1381	6. Walk towards a small yellow cylinder carefully
1382	7. Take hesitant steps towards a small yellow cylinder
1383	8. Approach a small yellow cylinder slowly and hesitantly
1384	9. Walk towards a small yellow cylinder with apprehension
1385	10. Move towards a small yellow cylinder hesitantly
1386	11. Walk to a small yellow cylinder with reluctance
1387	12. Take tentative steps towards a small yellow cylinder
1388	13. Walk towards a small yellow cylinder with uncertainty
1389	14. Approach a small yellow cylinder with trepidation
1390	15. Walk towards a small yellow cylinder with reservation
1391	16. Take hesitant strides towards a small yellow cylinder
1392	17. Walk to a small yellow cylinder with caution
1393	18. Move towards a small yellow cylinder with unease
1394	19. Walk towards a small yellow cylinder with doubt
1395	20. Approach a small yellow cylinder with timidity
1396	21. Walk towards a small yellow cylinder with hesitance
1397	22. Take slow steps towards a small yellow cylinder
1398	23. Walk towards a small yellow cylinder with wariness
1399	24. Move towards a small yellow cylinder with hesitancy
1400	25. Walk towards a small yellow cylinder with reluctance and caution.

Ori	ginal pull a circle hesitantly	1401
1.	Hesitantly pull a circle	1402
2.	Pull a circle with hesitation	1403
3.	Gently pull a circle	1404
4.	Pull a circle slowly and cautiously	1405
5.	Approach the circle and pull it hesitantly	1406
6.	Pull the circle with care	1407
7.	Pull the circle slowly and carefully	1408
8.	Take your time and pull the circle hesitantly	1409
9.	Pull the circle with a bit of hesitation	1410
10.	Approach the circle and pull it gently	1411
11.	Pull the circle with caution	1412
12.	Pull the circle carefully and hesitantly	1413
13.	Slowly pull the circle with hesitation	1414
14.	Pull the circle with a sense of caution	1415
15.	Pull the circle with a bit of reluctance	1416
16.	Pull the circle slowly and with care	1417
17.	Hesitate before pulling the circle	1418
18.	Pull the circle with a gentle touch	1419
19.	Pull the circle with a bit of apprehension	1420
20.	Pull the circle with a sense of uncertainty	1421
21.	Pull the circle with a bit of nervousness	1422
22.	Pull the circle with a bit of trepidation	1423
23.	Pull the circle with a sense of hesitation	1424
24.	Pull the circle with a bit of doubt	1425
25.	Pull the circle with a bit of reservation	1426
Ori	ginal pull a circle while zigzagging	1427
1.	Zigzag while pulling a circle	1428
2.	Pull a circle in a zigzag pattern	1429
3.	Carefully pull a circle while zigzagging	1430
4.	Zigzag and pull a circle simultaneously	1431
5.	Pull a circle while moving in a zigzag motion	1432
6.	With caution, pull a circle while zigzagging	1433
7.	Zigzag your way to the circle and pull it	1434
8.	Pull a circle while making zigzag movements	1435
9.	Zigzag and pull the circle with care	1436
10.	Pull a circle while navigating in a zigzag direction	1437
11.	Move in a zigzag pattern while pulling a circle	1438
12.	Pull a circle while making a zigzag path	1439

14. Pull a circle while making zigzag turns 1441 15. Carefully zigzag and pull the circle 16. Zigzag and carefully pull the circle 1443 1444 17. Pull a circle while making sharp zigzag movements 18. Zigzag and pull the circle with caution 1445 1446 19. Pull a circle while making quick zigzag motions 20. Zigzag and pull the circle slowly 1447 21. Pull a circle while zigzagging in a controlled manner 1448 22. Zigzag and pull the circle with precision 1449 23. Pull a circle while making small zigzag movements 1450 24. Zigzag and pull the circle with care and attention 1451

13. Zigzag towards the circle and pull it

1440

1453

1455

1456

1457

1458

1459

1460

1452 25. Pull a circle while zigzagging smoothly.

#### Properties of Natural-ish Language gSCAN Dataset $\mathbf{L}$

parses	words	zipf a	rmse
18	18	1.99	0.11
1550	859	1.29	0.01
234	20	1.90	0.10
9785	126	1.40	0.03
1400	35	1.26	0.04
42759	631	1.22	0.01
	parses 18 1550 234 9785 1400 42759	parses words 188 18 1550 859 234 20 9785 126 1400 35 42759 631	parses         words         zipf a           18         18         1.99           1550         859         1.29           234         20         1.90           9785         126         1.40           1400         35         1.22           42759         631         1.22

NL-) dataset



Figure 6: Word frequency distribution of NL-gSCAN and

Size Color Object 100% gSCAN 99.98%98.63% SR 100% 100% 100% ReaSCAN 100%99.99% 99.93%

Figure 7: Percentage of examgSCAN, each compared to the of each dataset and its corre-best fitting Zipf distribution ples in each training set whether sponding paraphrased (denoted probability density function, the object mentioned in the syngSCAN words are in orange and thetic dataset was also found in NL-gSCAN words are in blue exactly the same way the corre-(comprising of the larger vocab- sponding paraphrased example.

#### 1454 L.1**Linguistic Properties**

In this section we examine the linguistic properties of the dataset. The main research question is whether the instructions as paraphrased by GPT3.5 look more like natural language. Clearly, the paraphrased data has greater vocabulary complexity. But merely substituting words for synonyms would not make synthetic data appear any more natural, nor does it pose any real challenges to a learning algorithm that would need to act on the instructions. We examine two other indicia, unique parses and fit to a Zipf distribution of word frequency.

ulary).

**Parses** We compute the number of unique parses among all the instructions in each training set. A 1461 *parse* is an assignment of word-role labels, indicating the linguistic role of the token in the instruction. 1462 For example, a token may be an adjective, an adverb or some sort of connector. The parses are computed 1463 over every instruction in the training data using the spaCy package. As shown in Table 5, the number of 1464 unique parses in the paraphrased datasets are an order of magnitude larger than the number of unique 1465 parses in the synthetic datasets. This reflects the diversity of instruction structures that exist in the 1466 1467 paraphrased datasets.

**Zipfian Distribution Fit** Natural language is hypothesized to fit a Zipfian power-law distribution, 1468 where the probability of drawing a word from a corpus is inversely proportional to its frequency  $p(w) \propto \frac{1}{f_w^a}$ , 1469 where a is a parameter of the distribution which varies for different corpli. We estimate a using maximum 1470 likelihood estimation using the method in (Clauset et al., 2009) and compute the root-mean-squared error 1471 (RMSE) between the estimated probability of a word according to the estimated Zipf distribution and the 1472 empirical probability that word measured by counting word frequencies. A corpus that resembles natural 1473 language more closely will have a low RMSE to its corresponding Zipf distribution. We find that the 1474 paraphrased datasets better fit their Zipf distribution. We also visualize in both Figure 6 the ordered 1475 frequency distribution of the paraphrased gSCAN dataset and its corresponding Zip probability density 1476 function. 1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1489

1490

1491

1492

1493

1494

1495

1496

1498

1499

1501

1502

1503

1504

1505

1506

1507

1516

#### L.2**Compositional Properties**

We also examine whether the datasets maintained their compositional properties. Recall that the datasets are stratified into different splits to test different compositional generalization cases. We want to test whether these cases still hold. Clearly, in the output space, the compositional stratification still holds because we do not change the output actions. In the input space, we can only measure whether the same object is mentioned in each synthetic instruction and its corresponding paraphrased instruction, because the verbs and adverbs may be changed to a synonym or a sequence of words having a similar meaning.

As shown in Table 7, the retainment of target objects is very high, never going under 98%. We can be confident that the correct target object is mentioned in the same way in the paraphrased examples.

#### $\mathbf{M}$ Evaluation of baselines on Natural-ish gSCAN, gSCAN-SR and ReaSCAN

We evaluate current published state-of-the-art models with openly available code on the new datasets using our own re-implementation. We calculate the exact-match performance using seeds 0-9 using the same hyperparameters for each model, the details of which are specified in Appendix B. The models are briefly described below:

**ViLBERT** with Cross-Attention The ViLBERT model proposed in (Qiu et al., 2021), with only crossattention between visual and text input streams, then decoding the target action sequence autoregressively. As in (Sikarwar et al., 2022), the multi-level CNN on the grid world is replaced by adding learnable position encodings.

VilBERT with GRoCoT Self-Attention The same ViLBERT model but with the tweaks proposed in (Sikarwar et al., 2022), namely self-attention layers before cross-attention layers.

**Encoder-Decoder Transformer** A standard encoder-decoder Transformer, where the transformer input sequence is the position-encoded and embedded visual stream concatenated with the instruction, and the target output sequence are the actions, decoded autoregressively.

#### Results M.1

#### Ν Image-Based gSCAN

We also created an Image-Based gSCAN where the state inputs are images and encoded with a vision transformer with patch size 12. The results are reported in Table 9. We observed a similar boost on Split H for the NL + Img dataset as well. However, we note that the model for NL + Img appeared to be underfitting, so it is possible that with a larger model that the results could have been even better.

#### 0 Examples of generated demonstrations

We provide one-example-per-method of each support generation method on Split H in Figure 10. Examples 1508 in green are valid in the environment, relevant to the target object and correctly executed. Examples in 1509 yellow are considered "not relevant" since they concern an object with different properties than the one 1510 mentioned in the query. Examples in red are not correctly executed. Examples in grey are not valid in 1511 the environment. Note that for retrieval-based methods like GandR and Retrieval, the instruction is 1512 being solved in a different state to the query one, which is the reason why the action trajectories are both 1513 valid and correct, but look very different from each other. Up to 9 of the 16 possible supports are shown. 1514

Notice that **GandR** does not demonstrate the desired adverb "while spinning" (WALK(4)), because it is 1515 only finding near neighbours of "pull", which happen only with WALK and PUSH.

	Transformer	ViLBERT	ViLBERT(PP)			
gSCAN						
А	$1.0 \pm .00$	$1.0 \pm .00$	$1.0 \pm .00$			
В	$0.86\pm.28$	$0.94\pm.11$	$0.93 \pm .09$			
$\mathbf{C}$	$0.89\pm.16$	$0.89\pm.13$	$0.82 \pm .26$			
D	$0.01\pm.02$	$0.0\pm.01$	$0.0 \pm .00$			
Ε	$0.99\pm.02$	$0.93\pm.12$	$0.71 \pm .24$			
$\mathbf{F}$	$1.0\pm.00$	$1.0$ $\pm$ .00	$1.0 \pm .00$			
G	$0.0\pm.00$	$0.0\pm.00$	$0.0 \pm .00$			
Η	$0.19\pm.06$	$0.23\pm.01$	$0.17 \pm .06$			
	Ę	gSCAN-SR				
Ι	$1.0 \pm .00$	$1.0 \pm .00$	$1.0 \pm .00$			
II	$0.95\pm.04$	$0.93\pm.04$	$0.96\pm.02$			
III	$0.99\pm.01$	$0.96\pm.03$	$1.0 \pm .00$			
IV	$1.0\pm.00$	$1.0\pm.00$	$1.0 \pm .00$			
V	$0.46\pm.26$	$0.72 \pm .1$	$0.9 \pm .04$			
VI	$0.17\pm.18$	$0.61\pm.23$	$0.89\pm.06$			
		ReaSCAN				
IID	$0.99\pm.00$	$0.98\pm.02$	$0.97\pm.01$			
A1	$0.94\pm.02$	$0.95\pm.04$	$0.95 \pm .01$			
A2	$0.61\pm.05$	$0.52\pm.13$	$0.46\pm.07$			
B1	$0.75\pm.02$	$0.79\pm.05$	$0.75 \pm .03$			
B2	$0.54\pm.02$	$0.6\pm.09$	$0.53 \pm .05$			
C1	$0.37\pm.02$	$0.32\pm.02$	$0.64 \pm .03$			
C2	$0.27 \pm .05$	$0.22 \pm .05$	$0.22 \pm .03$			

Figure 8: The evaluation results for gSCAN, gSCAN-SR and ReaSCAN at 300,000 iterations, where performance for splits B-H is measured at the point where the model performed best on split A during training. ViLBERT is the model in (Qiu et al., 2021) and Tformer is an Encoder-Decoder Transformer. Tformer(PP) the same Transformer architecture evaluated on the paraphrased dataset.

	Transf	former	DemoGen			
	$\mathbf{NL}$	+ Img	NL	+ Img		
А	$1.0 \pm .00$	$1.0 \pm .00$	$0.99\pm.00$	$0.84\pm.01$		
В	$0.99\pm.00$	$0.93\pm.08$	$0.96\pm.00$	$0.53\pm.01$		
С	$0.99\pm.03$	$0.89\pm.16$	$0.97\pm.00$	$0.54\pm.01$		
D	$0.08 \pm .16$	$0.0\pm.00$	$0.01\pm.01$	$0.11\pm.02$		
Е	$0.98\pm.03$	$0.83\pm.22$	$0.98\pm.00$	$0.67\pm.00$		
F	$1.0 \pm .00$	$1.0\pm.00$	$0.98\pm.00$	$0.88\pm.01$		
G	$0.0 \pm .00$	$0.0\pm.00$	$0.0\pm.00$	$0.0\pm.00$		
Η	$0.19\pm.03$	$0.06\pm.05$	$0.59\pm.06$	$0.48\pm.02$		

Figure 9: Evaluation on natural language and image data. NL refers to natural language instructions, NL + Img refers to natural language instructions and patch-encoded images



## $I_1 =$ "pull a red small circle hesitantly"

- $I_2$  = "push a red big circle while spinning"  $I_3$  = "walk to a small circle hesitantly"
- $I_4 =$  "pull a circle hesitantly"
- $I_5 =$  "walk to a red circle hesitantly"
- $I_6 =$  "push a red big circle hesitantly"
- $I_7 =$  "pull a circle hesitantly"
- $I_8 =$  "pull a red small cylinder hesitantly"

Instruction Generator

- $I_9 =$  "walk to a small circle while spinning"

# (a) Support set generated by ${\bf Coverage}~{\bf Retrieval}$



 $I^Q$  = "pull a yellow cylinder while spinning"



 $I_q =$  "pull a green small circle while spinning"

Query



		$A_1 =$ "LTURN(2) WALK PULL"
le zigzagging"		$A_4 = \text{``LTURN(2)}$ WALK RTURN WALK LTURN WALK PULL(2)
		$A_{14} = $ "LTURN(2) WALK PULL
	÷	$A_{15} = $ "LTURN(2) WALK PULL
	orme	$A_{16} = $ "LTURN(2) WALK PULL
	ransf	

#### (b) Support set generated by **GandR**

Transformer

 $I_1$  = "walk to a green small circle while spinning"  $I_2$  = "push a green small circle while spinning  $I_3$  = "pull a green small circle while zigzagging  $I_4$  = "pull a green small circle hesitantly  $I_5$  = "pull a green small circle

# $$\begin{split} A_1 = ``LTURN(6) (WALK LTURN(4))(5) RTURN (WALK LTURN(4))(3) WALK" \\ A_2 = ``LTURN(6) (WALK LTURN(4))(5) RTURN (WALK LTURN(4))(3) PUSH LTURN(4) PUSH \\ A_3 = ``LTURN(2) (WALK RTURN WALK LTURN)(4) WALK PULL(2) \\ A_4 = ``LTURN(2) (WALK STAY)(5) RTURN (WALK STAY)(4) \end{split}$$

 $A_5 =$  "LTURN(2) WALK(5) RTURN WALK(4)

## (c) Support set generated by **Heuristic**

- $I_1$  = "walk to a blue small circle while spinning"  $I_2$  = "push a blue small circle while spinning
- $I_3 =$  "pull a blue small circle while zigzagging
- $I_4 = ``pull a blue small circle hesitantly }$
- $I_5=\ensuremath{``}\ensuremath{\mathsf{pull}}$  a blue small circle

Instruction Generator

- - $A_1 =$  "LTURN(4) (WALK LTURN(4))(4) RTURN (WALK LTURN(4))(3) WALK"
  - $A_2 =$  "LTURN(6) (WALK LTURN(4))(4) RTURN (WALK LTURN(4))(3) PUSH LTURN(4) PUSH
  - $A_3 =$  "LTURN WALK PULL(2)
  - $A_4 =$  "LTURN(2) (WALK STAY)(2) RTURN (WALK STAY)(4) (PULL STAY)(5)
  - $A_5 =$  "LTURN(2) WALK(4) RTURN WALK(4) PULL(10)



= "pull a blue small circle while spinnir

1	(d)	Support	set	generated	bv	Other	States
	u,	Dupport	500	generated	D y	Ounci	Duduco

Transformer

Instruction Generator	$I_1=\ensuremath{``}\ensuremath{push}$ a big blue square while zigzagging"	$A_1 = \text{``LTURN(2)}$ WALK RTURN WALK LTURN WALK RTURN WALK(2) PUSH(2)''
	$I_2 =$ "push a big blue square while spinning	$A_2 = ``LTURN(6) (WALK LTURN(4))(2) RTURN (WALK LTURN(4))(3) PUSH LTURN(4) PUSH$
	$I_3 =$ "push a small yellow circle	$A_3 =$ "LTURN(2) WALK RTURN WALK(4)
	$I_4 =$ "push a big blue cylinder"	$A_4 = $ "WALK(2) LTURN WALK(2) PUSH(2)
	$I_5=$ "walk to a small green cylinder while zigzaggir ${\begin{tabular}{c} {\bf e} \\ {$	$A_5 = $ "LTURN(2) WALK RTURN WALK(2)
	$I_6$ = "pull a big blue circle while spinning"	$A_6 =$ "LTURN(4) RTURN WALK (LTURN(4) PULL)(6) PULL
	$I_7$ = "push a big blue cylinder while spinning $\Box$	$A_7 = ``(LTURN(4) WALK)(2) LTURN(5) (WALK LTURN(4))(2) PUSH LTURN(4) PUSH$
	$I_8$ = "pull a big blue cylinder	$A_8 =$ "WALK(2) LTURN WALK(2) PULL
	$I_9 =$ "push a small yellow circle while zigzagging	$A_9 = $ "LTURN(2) WALK RTURN WALK(4)



Figure 10: Demonstrations generated on Split H for different kinds of demonstration strategies.