A Benchmark for Systematic Generalization in Grounded Language Understanding

Laura Ruis¹, Jacob Andreas², Marco Baroni³⁴ Diane Bouchacourt⁴, Brenden M. Lake⁴⁵

¹University of Amsterdam; ²Massachusetts Institute of Technology; ³ICREA; ⁴Facebook AI Research; ⁵New York University

racebook AI nesearch, new fork oniver-

Abstract

Human language users easily interpret expressions that describe unfamiliar situations composed from familiar parts ("greet the pink brontosaurus by the ferris wheel"). Modern neural networks, by contrast, struggle to interpret compositions unseen in training. In this paper, we introduce a new benchmark, gSCAN, for evaluating compositional generalization in models of situated language understanding. We take inspiration from standard models of meaning composition in formal linguistics. Going beyond an earlier related benchmark that focused on syntactic aspects of generalization, gSCAN defines a language grounded in the states of a grid world. This allows us to build novel generalization tasks that probe the acquisition of linguistically motivated rules. For example, agents must understand how adjectives such as 'small' are interpreted relative to the current world state or how adverbs such as 'cautiously' combine with new verbs. We test a strong multi-modal baseline model and a state-of-the-art compositional method finding that, in most cases, they fail dramatically when generalization requires systematic compositional rules.

1. Introduction

Human language is a fabulous tool for generalization. If you know the meaning of 'small', you can probably pick the 'small wampimuk' among larger ones, even if this is the first time you encountered wampimuks. If you know what 'walking cautiously' means, you can guess what a policeman means by 'biking cautiously' through a busy intersection (see Figure 1 for a related example). Modern deep neural networks, while achieving astounding results in many domains (LeCun et al., 2015), have not mastered comparable language-based generalization challenges, a fact conjectured to underlie their sample inefficiency and inflexibility (Lake et al., 2017; Lake & Baroni, 2018; Chevalier-Boisvert et al., 2019). Recent benchmarks have been proposed for languagebased generalization in deep networks (Johnson et al.,



Figure 1: Our gSCAN benchmark evaluates context sensitivity in situated language understanding. In these two simplified examples, the same determiner phrase "the red small circle" in the instructions will have different target referents and demand different action sequences depending on the world state and action verb. Being cautious means looking to the left and right (" $L_turn R_turn R_turn R_turn$ ") before crossing a grid line.

2017a; Hill et al., 2019), but they do not specifically test for a model's ability to perform rule-based generalization, or do so only in limited contexts. Systematic, rule-based generalization is instead at the core of the recently introduced SCAN dataset (Lake & Baroni, 2018; see also Hupkes et al., 2019 for related ideas). In a series of studies, Lake, Baroni and colleagues (Bastings et al., 2018; Loula et al., 2018; Dessì & Baroni, 2019) tested various standard deep architectures for their ability to extract general composition rules supporting zero-shot interpretation of new composite linguistic expressions (can you tell what 'dax twice' means, if you know the meaning of 'dax' and 'run twice'?). In most cases, neural networks were unable to generalize correctly. Very recent work has shown that specific architectural or training-regime adaptations allow deep networks to handle at least some of the SCAN challenges (Andreas, 2019; Lake, 2019; Nye & Tenenbaum, 2019; Russin et al., 2019; Gordon et al., 2020). However, it is unclear to what extent these proposals account for genuine compositional generalization, and to what extent they are "overfitting" to the limitations of SCAN.

SCAN simulates a navigation environment through an interpretation function that associates linguistic commands ('walk left') to sequences of primitive actions (L_TURN WALK). SCAN, however, is not grounded, in that it lacks a "world" with respect to which the agent has to interpret the commands: instead the agent must simply associate linguistic strings with fixed sequences of action symbols, essentially reducing the problem of interpretation to one of mapping syntactic strings (word sequences) to other syntactic strings (action label sequences). In real languages, by contrast, the process by which utterances are understood is both compositional and *contextual*: references to entities and descriptions of actions must be interpreted with respect to a particular state of the world. The interaction between compositional structure and contextual interpretation introduces various new types of generalization an intelligent agent might have to perform. For example, consider the meaning of size adjectives such as 'small' and 'large'. The determiner phrases 'the small bottle' and 'the large bottle' might refer to the same bottle, depending on the sizes of the bottles that surround the relevant one. We might wonder whether this and related notions of compositional generalization can be addressed using existing techniques, but SCAN's context insensitivity makes it impossible to investigate broader notions of generalization.

We introduce grounded SCAN (gSCAN), a new benchmark that, like the original SCAN, focuses on rule-based generalization, but where meaning is grounded in states of a grid world accessible to the agent. This allows modeling a much more comprehensive set of linguistic generalizations. For example, Figure 1 shows how the target referent 'the red small circle', and the action sequence required to navigate there, will change based on the state of the world. We propose a number of new linguistic generalization challenges based on gSCAN. We further test a baseline multi-modal model representative of contemporary deep neural architectures, as well as a recent method proposed to address compositional generalization in the original SCAN dataset (GECA, Andreas, 2019). Across seven different generalization splits, we show that the baseline dramatically fails on all but one split, and that GECA improves performance on an additional one only. These results demonstrate the challenges of accounting for common natural language generalization phenomena with standard neural models, and affirm gSCAN as a fruitful benchmark for developing models with more human-like compositional learning skills.

2. Related Work

Much recent work has recognized the advantages of compositional generalization for robustness and sample efficiency, and responded by building synthetic environments to evaluate aspects of this skill (Johnson et al., 2017b; Lake & Baroni, 2018; Loula et al., 2018; Bahdanau et al., 2018; Hill et al., 2019; Chevalier-Boisvert et al., 2019). Bahdanau et al. (2018) evaluate models on binary questions about object relations, generalizing to unseen object combinations after training on a small subset. Chevalier-Boisvert et al. (2019) evaluate models on tasks such as navigating to objects and moving them in a grid world, studying the influence of curriculum learning on sample efficiency. Hill et al. (2019) evaluate unseen combinations of verbs and objects in complex skill learning, demonstrating that richer grounded environments can promote generalization. Lake & Baroni and Loula et al. (2018) propose the SCAN benchmark for evaluating systematic generalization, which differs from related work by specifically testing for linguistic generalization and the acquisition of abstract compositional rules. SCAN contains a compositional set of instructions generated by a phrase-structure grammar that can unambiguously be translated into action sequences by applying an interpretation function. The data is split into training and test sets that contain systematic differences. For example, in one split models must interpret phrases that contain primitives only encountered in isolation at training time (e.g., inferring that the command 'jump twice' translates to action commands JUMP JUMP when you know that 'walk twice' translates to WALK WALK and 'jump' translates to JUMP). SCAN however lacks grounding, which substantially limits the variety of linguistic generalizations it can examine.

Gordon et al. (2020) formalize the type of compositionality skills required by the SCAN dataset as equivariance to a certain group of permutations. Their model is hard-coded to be equivariant to all permutations of SCAN's verb primitives and succeeds on some of the tasks. However, the method only tackles local permutations in the input command (e.g., swapping 'walk' for 'jump') that result in local permutations in the action sequence (swapping WALK for JUMP). More realistically, in our *grounded* SCAN environment, the permutation of a word in the command could result in reference to a different target object, a different interaction with the object, or a different manner of moving through the grid. Permutations in the instruction therefore modify the correct output action sequence in a non-local manner, in their terminology. The same problem would affect the syntax-semantics separation approach proposed by Russin et al. (2019).

Another method that has successfully dealt with some of the SCAN splits is the meta-learning approach by Lake (2019). This model "learns how to learn" new primitives in meta-training episodes where words are randomly mapped to meanings. The method successfully solves some of the SCAN challenges but, in its current implementation, it is unclear how to apply it to grounded



Figure 2: Examples showing how to 'walk while spinning' and how to 'push.' On the left, the agent needs to spin around (' L_turn L_turn L_turn L_turn ') before it moves by one grid cell. On the right, it needs to push a square all the way to the wall.

semantics, where it is impossible to exploit random mapping between primitives and action symbols.

A model-agnostic method that is readily applicable to our benchmark and obtains successful results on SCAN is good-enough compositional data augmentation (GECA, Andreas, 2019). GECA identifies sentence fragments which appear in similar environments and uses those to generate novel training examples. For instance, when one sees evidence during training that 'the cat sang', 'the wug sang' and 'the cat danced' are sentences with high probability, then 'the wug danced' is also probable, but 'the sang danced' is not. The assumption here is that 'cat' and 'wug' are interchangeable, and GECA indeed helps when words are fully interchangeable, but the assumption is not always valid in more realistic, grounded language understanding.

3. The Grounded SCAN Benchmark

Our goal is to test for a broad set of phenomena in situated language understanding where humans should easily generalize, but where we expect computational models to struggle due to the systematicity of the differences between train and test examples. In this section we describe what tools we work with to achieve this goal, and in Section 5 we describe in detail how each linguistic phenomenon is tested. All code for generating the benchmark, as well as the data used in this paper is publicly available¹.

Instructions. Grounded SCAN (gSCAN) poses a simple task where an agent must execute instructions in a two-dimensional grid world with objects. We build on the formal approach of SCAN (Lake & Baroni, 2018) while evaluating a much wider range of linguistic generalizations by grounding the semantics of the input instructions. The world model allows us to examine how often an agent needs to see 'move cautiously' before applying 'cautiously' in a novel scenario, whether an agent

can identify a novel object by reasoning about its relation to other objects in the world, and whether an agent can infer how to interact with objects by identifying abstract object properties. Figure 2 shows two example commands and corresponding action sequences, which are simplified but representative of the type provided in gSCAN (actual gSCAN examples use a larger grid world with more objects). On the left, the agent must generate a sequence of target actions that leads to the circle 'while spinning.' All adverbial modifiers such as 'cautiously' or 'while spinning' require applying complex, context-sensitive transformations to the target sequence, going beyond the simple substitutions and concatenations representative of SCAN. On the right (Figure 2), the agent must push a small square, where pushing is defined to mean moving something as far as possible without hitting the wall (in this case), or another object. The agent can also be instructed to 'pull' some object, in which case it would pull the object back as far as possible. The full phrase-structure grammar for producing instructions is provided in Appendix A. The set of actions the agent can use to execute instructions is $\mathcal{A} = \{$ walk, push, pull, stay, L_turn, R_turn $\}$.

World model. Each instruction is paired with a relevant world state, presented to the agent as a tensor $\mathbf{X}_s \in \mathbb{R}^{d \times d \times c}$ for grid size d (d = 6 or 12 depending on split). The object at each grid cell is defined via one-hot encodings along three property types, namely color $\mathcal{C} = \{\text{red, green, blue, yellow}\}$, shape $\mathcal{S} = \{\text{circle, square, cylinder}\}$, and size $\mathcal{D} = \{1, 2, 3, 4\}$. Specifying the agent location and heading requires five more channels, and thus the tensorwidth is $c = 5 + |\mathcal{C}| + |\mathcal{S}| + |\mathcal{D}|$.

Each instruction also constrains the generation of the world state. For instance, each target referent from the instruction determiner phrase is ensured to be unique (only one possible target in "walk to the yellow square"). Moreover, to conform to natural language pragmatics, we ensure that if a size modifier is used in the instruction, there is always a relevant distractor object. For example, when the target referent is "the small square", we additionally place a square that is larger than the target (see Figure 2 right). For a full description of which objects we place in the world, refer to Appendix B. Further, objects of size 1 and 2 are *light*, and objects of size 3 and 4 are *heavy.* This division determines how the agent should interact with the object. If an object is light, it needs to be pushed once to move it to the next cell (executed by the action command 'push'). If an object is heavy, it needs to be pushed twice to move to the next cell ('push push').

Data splits. Equipped with this framework, we design splits with systematic differences between training and test. We distinguish two broad types of tests, compositional generalization and length generalization. Our

¹https://github.com/LauraRuis/groundedSCAN

suite of evaluations requires training just two models – one for compositional and one for length generalization – since the evaluations were designed to work with the same base training set. We also examine a "random split" with no systematic differences between training and test (Section 5A).

'Compositional generalization' evaluates combining known concepts into novel meaning (Section 5B-H). From a single training set we can evaluate how an agent handles a range of systematic generalizations, including novel object property combinations ('red square'; Section 5B,C), novel directions (a target to the south west; 5D), novel contextual references ('small yellow circle'; 5E), and novel adverbs ('pull cautiously'; 5G,H). For example, we model an analogue of the 'wampimuk' case from the introduction by holding out all examples where a circle of size 2 is referred to as 'the small circle' (Section 5E). We test whether models can successfully pick out the small circle among larger ones, even though that particular circle is referred to during training only as 'the circle' (with no other circles present) or 'the large circle' (with only smaller circles present). The shared training set across splits has more than 300k demonstrations of instructions and their action sequences, and each test instruction evaluates just one systematic difference. For more details on the number of examples in the training and test sets of the experiments, refer to Appendix C.

'Length generalization' (Section 5I) evaluates a persistent problem with sequence generation models: they fail to generalize beyond the maximum length seen during training (e.g., Graves et al., 2014). Since length generalization is entirely unsolved for SCAN, even for methods that made progress on the other splits (Lake & Baroni, 2017; Bastings et al., 2018; Russin et al., 2019; Gordon et al., 2020), we also separately generate a split to evaluate generalization to longer action sequences. We do this by using a larger grid size than in the compositional generalization splits (d = 12), and hold out all examples with a target sequence length of m > 15. During training, a model does see the full grid and all the possible instructions with enough unique world states for each instruction (see Appendix C), but they require generating less action commands than at test time. At test time the action sequences can be up to m = 47.

4. Baselines

Models are trained using supervised learning to map instructions to target sequences, given a world context. We train a multi-modal neural network baseline to generate target action sequences, conditioned on the input commands and world states (see Figure 3 for overview). The architecture is not novel and uses standard machinery (e.g., Mei et al. (2016)), but we nonetheless explain the key components below for completeness (full specification in Appendix D).

The baseline model consists of a sequence-to-sequence neural network (seq2seq) (Sutskever et al., 2014) fused with a visual encoder. The model uses a recurrent 'command encoder' to process the instructions ('Walk to the circle' in Figure 3) and a 'state encoder' to process the locations of the agent and the objects. A recurrent decoder produces the appropriate action sequence (e.g., 'walk') through joint attention over the steps of the command sequence and the spatial locations in the grid world. The input tuple $\mathbf{x} = (\mathbf{x}^c, \mathbf{X}^s)$ includes the command sequence $\mathbf{x}^c = \{x_1^c, \ldots, x_n^c\}$ and the world state $\mathbf{X}^s \in \mathbb{R}^{d \times d \times c}$, as represented over a $d \times d$ grid. The target sequence $\mathbf{y} = \{y_1, \ldots, y_m\}$ is modeled as

$$p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_{j=1}^{m} p_{\theta}(y_j \mid \mathbf{x}, y_1, \dots, y_{j-1}).$$

Command encoder. The network processes the input instruction through using a bidirectional LSTM (Hochreiter & Schmidhuber, 1997; Schuster & Paliwal, 1997) denoted as $\mathbf{h}^c = f_c(\mathbf{x}^c)$ (Figure 3). It produces a sequences of embeddings $\mathbf{h}^c = \{h_1^c, \ldots, h_n^c\}$ with a vector for each word.

State encoder. The network perceives the initial world state through a convolutional network $\mathbf{H}^s = f_s(\mathbf{X}^s)$ (Figure 3), using three kernel sizes (Wang & Lake, 2019). It produces a grid-based representation of the world state $\mathbf{H}^s \in \mathbb{R}^{d \times d \times 3c_{\text{out}}}$ with c_{out} as the number of feature maps per kernel size.

Decoder. The output decoder f_d parameterizes the distribution over action sequences based on the decoder messages, $p(\mathbf{y}|\mathbf{h}^c, \mathbf{H}^s)$. At each step, the previous output symbol y_{j-1} is embedded as $\mathbf{e}_j^d \in \mathbb{R}_e^d$, and context vectors for the command \mathbf{c}_j^c and the world state \mathbf{c}_j^s are computed from the previous decoder state \mathbf{h}_{j-1}^d using double attention (Devlin et al., 2017). The recurrent decoder operates as

$$\mathbf{h}_{j}^{d} = \mathrm{LSTM}([\mathbf{e}_{j}^{d}; \mathbf{c}_{j}^{c}; \mathbf{c}_{j}^{s}], \mathbf{h}_{j-1}^{d}),$$

which produces state \mathbf{h}_{j}^{d} based on the previous state \mathbf{h}_{j-1}^{d} and the other variables mentioned above as input. First the command context is computed as $\mathbf{c}_{j}^{c} = \text{Attention}(\mathbf{h}_{j-1}^{d}, \mathbf{h}^{c})$, attending over the input steps and producing a weighted average over \mathbf{h}^{c} . Second, while conditioning \mathbf{c}_{j}^{c} , the state context is computed as $\mathbf{c}_{j}^{s} = \text{Attention}([\mathbf{c}_{j}^{c}; \mathbf{h}_{j-1}^{d}], \mathbf{H}^{s})$, attending over grid locations and producing a weighted average over \mathbf{H}^{s} . The action emission y_{j} is then

$$p(y_j \mid \mathbf{x}, y_1, \dots, y_{j-1}) = \operatorname{softmax}(\mathbf{W}_o[\mathbf{e}_j^d; \mathbf{h}_j^d; \mathbf{c}_j^c; \mathbf{c}_j^s])$$

Training. Training optimizes cross-entropy using Adam (Kingma & Ba, 2014) with default parameters. Supervision is provided by ground-truth target sequences of which there is one for each instruction-world-state pair

Table 1: Summary of results for each split, showing exact match accuracy averaged over 3 runs and the standard deviation between the runs. Both the base-line and GECA fail on all splits except C and F.

	Exact Match (%)				
Split	Baseline	GECA			
A: Random	97.69 ± 0.22	87.6 ± 1.19			
B: Yellow squares	54.96 ± 39.39	34.92 ± 39.30			
C: Red squares	23.51 ± 21.82	78.77 ± 6.63			
D: Novel direction	0.00 ± 0.00	0.00 ± 0.00			
E: Relativity	35.02 ± 2.35	33.19 ± 3.69			
F: Class inference	92.52 ± 6.75	85.99 ± 0.85			
G: Adverb $k = 1$	0.00 ± 0.00	0.00 ± 0.00			
Adverb $k = 5$	0.47 ± 0.14	-			
Adverb $k = 10$	2.04 ± 0.95	-			
Adverb $k = 50$	4.63 ± 2.08	-			
H: Adverb to verb	22.70 ± 4.59	11.83 ± 0.31			



Figure 3: Baseline neural network. The command encoder reads the command with a biLSTM (f_c ; top left), and the state encoder reads the world state with a CNN (f_s ; top right). An LSTM decoder produces the appropriate action sequence (bottom), using joint attention over the command and the world state.

if we use the convention to always travel horizontally first and then vertically. The learning rate starts at 1e-3 and decays by 0.9 every 20,000 steps. We train for 200,000 steps with batch size 200. The best model based on the exact match on a small development set of 2,000 examples was used for the experiments. The parameters that turned out most important for performance were the kernel sizes that process the state, chosen to be 1,5 and 7 for the experiments with a grid size of 6, and 1, 5, and 13 for the grid of size 12. Both the data and the code for the model is made available², and a full specification of the parameters is given in Appendix E.

Good-enough compositional data augmentation. We use the procedure proposed by Andreas (2019) to generate novel training examples for the compositional generalization training set. GECA is run on gSCAN with essentially the same parameters as the SCAN experiments in the original paper: the context window is taken to be full sentences, with a gap size of one and a maximum of two gaps (Appendix E for details). GECA receives an input sequence consisting of the natural language command concatenated with an output sequence containing a linearized representation of the target object's feature vector. After generating augmented sequences, we re-apply these to the gSCAN dataset by modifying training examples with augmentable input sentences. Modification leaves action sequences unchanged while changing the commands and environment features. For example, an input instruction "walk to a red circle" could be modified to "walk to a red square", and the world state would be subsequently modified by replacing the target red circle with a red square.

5. Experiments

The main contribution of this work is the design of test sets that require different forms of linguistic generalization. In this section we detail both types of test sets: compositional generalization and length generalization (recall there are only two training sets, respectively). The main results averaged over three runs are summarized in Table 1, where we show the standard deviation between runs, and each split is detailed below.

A: Random split. The random split verifies that the models can learn to follow gSCAN commands when there are no systematic differences between training and test. The baseline model achieves near perfect exact match accuracy (97.69% \pm 0.22) on the 19,282 test examples, where exact match means that the entire action sequence is produced correctly. GECA performs worse (87.6% \pm 1.19) which is unsurprising since the assumption underlying the data augmentation procedure in GECA is that phrases that appear in similar environments can be permuted. This is not always correct in gSCAN (none of the verbs and adverbs can be permuted).

B, **C**: Novel composition of object properties. Here we examine whether a model can learn to recombine the properties of color and shape to recognize an unseen colored object (see Figure 4 for examples). This test is closely related to CoGenT by Johnson et al. (2017a); learning disentangled representations for object properties will be helpful here too. However, by exploiting the grounded nature of our set-up where objects are clearly distinguished from the referring expressions denoting them, we consider two separate setups, one involving *composition of references* and another involving

²https://github.com/LauraRuis/multimodal_seq2seq_gSCAN

Table 2: Exact match broken down by referred target ('Ref. Target', i.e., the target object denoted in the determiner phrase of the instruction). The \star column indicates chance performance of choosing an object uniformly at random and correctly navigating to it.

Ref. Target *		Baseline	GECA
'small red square'	8.33	13.09 ± 14.07	78.64 ± 1.10
'big red square'	8.33	11.03 ± 10.29	77.88 ± 0.95
'red square'	16.67	8.48 ± 0.90	97.95 ± 0.14
'small square'	8.33	27.13 ± 41.38	66.26 ± 13.60
'big square'	8.33	22.96 ± 32.20	68.09 ± 20.90
'square'	50	52.92 ± 36.81	95.09 ± 7.42

composition of attributes. In a first split, we hold out all data examples where a yellow square (of any size) is the target object and is referred to with the determiner phrases 'the yellow square', 'the small yellow square' or 'the big yellow square' (i.e., any phrase containing the color adjective and the shape). The training set contains examples with yellow squares as a target, but they are always referred to without a color: 'the square', 'the big square', or 'the small square', meaning the model cannot ground that target to the reference 'yellow square' (Figure 4 left). The second split never sees examples where a red square is the target in training, meaning the methods. in addition to never encountering the determiner phrase 'the red square', cannot ground a reference to this object (Figure 4 right). However, the red square does appear as a non-target background object.

The baseline shows poor performance on the 'red squares'-split that requires zero-shot generalization to the target red square $(23.51\% \pm 21.82)$; again exact match with standard deviation over runs). GECA however significantly improves performance on this split $(78.77\% \pm$ 6.63). This is precisely what GECA is designed for; permuting 'red circle' and 'yellow square' during training gives familiarity with 'red square'. Surprisingly, GECA does not improve over the baseline for the 'yellow square'split $(34.92\% \pm 39.30 \text{ for GECA and } 54.96\% \pm 39.39 \text{ for}$ the baseline). In this split, yellow squares have been seen during training as a target object, yet never referred to using their color: the referent in the instruction could have been '(small,big) square' but never '(small,big) yellow square'. We hypothesize that models overfit this pattern, and while GECA should help by generating instructions including 'yellow', their number is still too small compared to the ones of the form '(small,big) square'. In Appendix C we take a closer look at the differences in the training evidence seen by the baseline model and GECA.

When breaking down the average exact match per referred target for the 'red squares'-split, we in fact can reason about what happens (Table 2). The baseline model never sees the referred target 'red square' (except as a background object) and is unable to compositionally construct a meaningful representation of this object

(low exact match with low standard deviation). It has seen plenty evidence for other red objects, like red circles and cylinders, and other colored squares, blue or green. This result confirms results in related work that vanilla neural networks fall short in recombining familiar concepts in novel ways (Kuhnle & Copestake, 2017; Loula et al., 2018). Higher performance when only the shape is explicitly denoted (e.g. "walk to the square", when the square happens to be red) is expected because in this case the model can randomly choose one of two objects as a target (with chance as 50% correct). When only a shape is denoted, there are at most 2 objects in the world, as a result of the constraint of a unique target object combined with randomness we introduced in deciding which objects to place in the world. If an instruction only mentions 'square', there are |S| - 1 = 2possible other objects to place in the world, of which we randomly select half, resulting in 2 objects placed in the world (including the target). GECA gets a very high exact match and low standard deviation when 'the red square' or 'the square' is used, but is thrown off when a size adjective is mentioned. It seems like, again, these adjectives are too strongly grounded to the targets seen in training.

D: Novel direction. In the next experiment we examine generalizing to navigation in a novel direction. We hold out all examples from the training set where the target object is located to the south-west of the agent. The agent can train on walking in any other direction, and needs to generalize to walking to the west and then the south, or vice-versa³. Conceptually, at test time the agent needs to combine the familiar notions of walking to the south and west. Results are in the row 'novel direction' of Table 1. Both methods obtain 0 exact matches over all runs (0% correct). A closer analysis of the predictions (see Figure 8 Appendix F for 2 examples) shows that the agent usually walks all the way west (or south) and then fails to turn to the target object. The attention shows that the agent knows where to go (by attending to the correct grid cell), just not how to get there. Even though there is catastrophic failure at the task overall, the agent learned by the baseline model ends up in the correct row or column of the target $63.10\% \pm 3.90$ of the times, and the agent learned with GECA $58.85\% \pm 3.45$ of the times. This further substantiates that they often walk all the way west, but then fail to travel the distance left to the south, or vice-versa. The results on this split indicate that apparently the methods completely fail to generate target sequences that have either three occurrences of 'L_turn' (needed to walk to the west and then south for an agent that starts facing east) or two occurrences of 'R_turn' (needed to walk to the south and then west) spread over the target sequence.

 $^{^3\}mathrm{For}$ supervision, the agent walks horizontally then vertically, for evaluation both are fine.



Figure 4: Generalizing from calling an object "big square" to calling it "big yellow square" (left), and from knowing "red" and "square" to identifying a "red square" (right).

E: Novel contextual references. In natural language there are many words that can only be grounded to relative concepts. Which object one refers to when saying 'the small circle' is fully dependent on the other circles in the world state. We investigate whether a model can grasp the concept of relativity in language by considering a scenario where objects of a specific size (size 2) are never targets correctly picked by the 'small' modifier in the training phase (see example in Figure 5). At test time, the target is a circle of size 2, which is being correctly referred to as a 'small circle' (the determiner phrase may also contain a color specification). In other words, we hold out for testing all world states where the circle of size 2 is the target and the smallest circle in the world, paired with an instruction containing the word 'small'. The agent can ground the circle of size 2 to references like 'the circle', 'the green circle' or 'the big circle', but needs to generalize to that same circle being referred to as 'the small circle' at test time.

The results on this split by both methods are again substantially worse than on the random split: $35.2\% \pm$ 2.35 for the baseline and $33.19\% \pm 3.69$ for GECA. When breaking down the exact match per referred target it seems like the model is exploiting the fact that when in addition to the size modifier 'small' the color of the circle is specified, it can randomly choose between 2 circles of the specified color, as opposed to randomly choosing between any circle in the world. When generating world states for instructions containing some combination of a color, size modifier and shape in the determiner phrase (e.g. "the small red circle") during data generation we always generate 2 differently sized objects of each color-shape pair (except for the referred color-shape pair we generate 1 other). So when you recognize the color and shape in the instruction, you have a 50% chance of picking the right object. Then how to interact with it if necessary is already familiar. We observe that for data examples where the instruction specifies the color of the target in addition to the size the baseline achieves $53\%.00 \pm 1.36$ and GECA achieves $47.51\% \pm 12.59,$ suggesting the agents randomly select a circle of the specified color. Thus the obtained performance indicates a

Figure 5: Generalizing from calling an object "big" to calling it "small."



pulling heavy pushing a to square.

complete failure of genuinely understanding "small" and picking a small circle from among larger ones in arbitrary circumstances.

F: Novel composition of actions and arguments. Another phenomenon in natural language we want to model is categorization of words into classes whose entries share the same semantic properties (Pustejovsky, 1991). We study the simple case of nominal class inference, establishing two categories of nouns, that, depending on their weight (they can be light or heavy), will lead to a different interpretation of the verb taking them as patient arguments. Recall from Section 3 that pushing or pulling a heavy object over the same distance (i.e., grid cells) as a light object requires twice as many target actions of 'push' or 'pull'.

In this experiment we test a model's ability to infer the latent object class, and correctly interact with it, as illustrated in Figure 6. We hold out all examples where the verb in the instruction is 'push', and the target object is a square of size 3, meaning it is in the heavy class and needs to be pushed twice to move by one grid cell. A model should infer that this square of size 3 is 'heavy' from its extensive training experience 'pulling' this object (as it always needs to use two pull actions to move it). Note that Hill et al. (2019) similarly studies verb-noun binding.

Both methods obtain almost the exact match they get on the random split, namely $92.52\% \pm 6.75$ and $85.99\% \pm 0.85$ by the baseline and GECA respectively (compared to only examples with 'push' in the random split, for which the baseline gets an exact match of $96.64\% \pm 0.52$ and GECA of $86.72\% \pm 1.23$), and seem to be able to correctly categorize the square of size 3 in the class heavy and interact with it accordingly. This is consistent with the findings of Hill et al. (2019) with regards to generalizing familiar actions to new objects.

G, H: Novel adverbs. In the penultimate experiment we look at a model's ability to transform target sequences when an adverb modifies the verb in the input command. The adverbs all require the agent to do something (i.e., generate a particular action sequence) at some predefined interval (see Figure 1 and Figure 2 for examples). To do something *cautiously* means to look to the left and right before crossing grid lines, to do something *while spinning* requires spinning around after moving a grid cell, to do something *hesitantly* makes the agent stay put every time it moves a grid cell (with action command 'stay'), and finally to do something *while zigzagging* only applies to moving diagonally on the grid. Where normally the agent would first travel horizontally all the way and then vertically, when doing something while zigzagging the agent will alternate between moving vertically and horizontally every grid cell.

We design two experiments with adverbs. The first examines whether an agent can learn the adverb 'cautiously' from just one or a few examples and use it in different world states (few-shot learning). For instance, the agent sees a single example of an instruction with the adverb 'cautiously' during training, and needs to generalize to all other possible instructions with that adverb (while discarding examples with longer target sequences than seen during training). The second examines whether a model can generalize a familiar adverb to a familiar verb, namely 'while spinning' to 'pull'. In this experiment the agent sees ample evidence of both the tested adverb and the verb, but has never encountered them together during training. These experiments are related to the 'around right'-split introduced in SCAN by Loula et al. (2018), but in this case, the grounded meaning of each adverb in the world state changes its effect on the target sequence. Therefore we expect methods like the equivariance permutations of Gordon et al. (2020), but also GECA, to have no impact on generalization.

For the few-shot learning experiment, the models fail catastrophically when learning 'cautiously' from just one demonstration (0% correct; exact match again). We experiment with increasing the number of times 'cautiously' appears in the training set (k) for the baseline model (Table 1), but find it only marginally improves its abysmal performance. Even with as much as 50 examples of how to move cautiously, the baseline fails to generalize (4.6% correct), emphasizing the difficulty neural networks have in learning abstract concepts from limited examples. We further find that the models struggle with longer sequences, and performance quickly drops as a function of target length (Figure 9 Appendix F).

For combining spinning and pulling, both methods again fail (22.70% \pm 4.59 for the baseline and 11.83% \pm 0.31 for GECA). The methods also struggle with longer sequences and performance drops as a function of target length (Figure 10 Appendix F). Evidently, even though for the random split the methods can handle very long sequences, the networks struggle when the long sequences concern unfamiliar combinations.

I: Novel action sequence lengths. For this split we only train the baseline model, as data augmentation will not help for target sequence generation beyond the length

encountered at training. The baseline model trained on examples that require generating sequences of lengths ≤ 15 gets an exact match on a held-out test set with examples also up to length 15 of 94.98% ± 0.12 , but for the test set with target lengths of 16 it obtains an exact match of $19.32\% \pm 0.02$, for target lengths of 17 it drops to $1.71\% \pm 0.38$ and for target lengths beyond 18 the performance is below 1%. Unsurprisingly, also when this task is posed by grounded SCAN, the neural network fails to generalize to longer target sequences.

6. Conclusion

There has been important recent progress on the SCAN compositional learning benchmark through applications of meta-learning (Lake, 2019), compositional data augmentation (Andreas, 2019), permutation equivariance (Gordon et al., 2020), and syntax/semantics separation (Russin et al., 2019). Our results, on a new gSCAN (grounded SCAN) benchmark, suggest these new methods largely exploit artifacts in SCAN that are not central to the nature of compositional generalization. The gSCAN benchmark is based on the linguistic formalization of compositionality, and moves closer to formal semantics by grounding them in a world model. We trained a multi-modal baseline and a state-of-the-art method from SCAN (GECA), finding that the baseline fails on all but one split, and the state-of-the-art methods on all but two. Both methods are able to systematically generalize in the experiment where the agent needs to infer the class of an object and use that knowledge to properly interact with it (similarly to what was found by Hill et al., 2019). GECA additionally increases performance significantly on a split where the agent needs to recombine familiar concepts of color and shape into an unseen target object. However the complete failure on the subsequent splits show that fundamental advances are still needed regarding neural architectures for compositional learning.

Progress on gSCAN may come from continuing the lines of work that have made progress on SCAN. Metalearning (Lake, 2019) or equivariance permutation (Gordon et al., 2020) could support compositional generalization, if the types of generalizations examined here can be incorporated into a suitable meta-training procedure or an equivariance definition. For now, at least, future work will require highly non-trivial extensions to apply these approaches to gSCAN.

In future work, we can extend gSCAN to support reinforcement learning setup in addition to supervised learning. In addition, the gSCAN environment can be extended to supply only RGB images of the world rather than partially-symbolic state representations. These extension, however, will only make the benchmark more difficult. We hope the current challenge, as is, will stimulate many advances in compositional learning and systematic generalization.

Acknowledgments. We are grateful to Adina Williams and Ev Fedorenko for very helpful discussions, and to João Loula who did important initial work to explore compositional learning in a grid world.

References

- Andreas, J. Good-enough compositional data augmentation. CoRR, abs/1904.09545, 2019. URL http: //arxiv.org/abs/1904.09545.
- Bahdanau, D., Murty, S., Noukhovitch, M., Nguyen, T. H., de Vries, H., and Courville, A. Systematic generalization: What is required and can it be learned? *arXiv preprint*, pp. 1–16, 2018.
- Bastings, J., Baroni, M., Weston, J., Cho, K., and Kiela, D. Jump to better conclusions: SCAN both left and right. In *Proceedings of the EMNLP BlackboxNLP Workshop*, pp. 47–55, Brussels, Belgium, 2018.
- Chevalier-Boisvert, M., Bahdanau, D., Lahlou, S., Willems, L., Saharia, C., Nguyen, T. H., and Bengio, Y. BabyAI: A platform to study the sample efficiency of grounded language learning. In *Proceedings* of *ICLR*, New Orleans, LA, 2019.
- Dessì, R. and Baroni, M. CNNs found to jump around more skillfully than RNNs: Compositional generalization in seq2seq convolutional networks. In *Proceedings* of ACL, Firenze, Italy, 2019. In press.
- Devlin, J., Uesato, J., Bhupatiraju, S., Singh, R., Mohamed, A. R., and Kohli, P. RobustFill: Neural program learning under Noisy I/O. International Conference on Machine Learning (ICML), 3:1641–1658, 2017.
- Gordon, J., Lopez-Paz, D., Baroni, M., and Bouchacourt, D. Permutation equivariant models for compositional generalization in language. 2020. URL https://openreview.net/forum?id=SylVNerFvr.
- Graves, A., Wayne, G., and Danihelka, I. Neural Turing Machines. *arXiv preprint*, 2014. URL http://arxiv. org/abs/1410.5401v1.
- Hill, F., Lampinen, A., Schneider, R., Clark, S., Botvinick, M., McClelland, J. L., and Santoro, A. Emergent Systematic Generalization in a Situated Agent. arXiv preprint, 2019. URL http://arxiv. org/abs/1910.00571.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

- Hupkes, D., Dankers, V., Mul, M., and Bruni, E. The compositionality of neural networks: integrating symbolism and connectionism, 2019.
- Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, L., and Girshick, R. CLEVR: a diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of CVPR*, pp. 1988– 1997, Honolulu, HI, 2017a.
- Johnson, J., Hariharan, B., van der Maaten, L., Hoffman, J., Fei-fei, L., Zitnick, C. L., and Girshick, R. Inferring and Executing Programs for Visual Reasoning. In *International Conference on Computer Vision*, 2017b.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. CoRR, abs/1412.6980, 2014.
- Kuhnle, A. and Copestake, A. A. Shapeworld A new test methodology for multimodal language understanding. *CoRR*, abs/1704.04517, 2017. URL http://arxiv.org/abs/1704.04517.
- Lake, B. and Baroni, M. Generalization without systematicity: On the compositional skills of sequence-tosequence recurrent networks. In *Proceedings of ICML*, pp. 2879–2888, Stockholm, Sweden, 2018.
- Lake, B., Ullman, T., Tenenbaum, J., and Gershman, S. Building machines that learn and think like people. *Behavorial and Brain Sciences*, 40:1–72, 2017.
- Lake, B. M. Compositional generalization through meta sequence-to-sequence learning. In Advances in Neural Information Processing Systems, 2019.
- Lake, B. M. and Baroni, M. Still not systematic after all these years: On the compositional skills of sequence-tosequence recurrent networks. *CoRR*, abs/1711.00350, 2017. URL http://arxiv.org/abs/1711.00350.
- LeCun, Y., Bengio, Y., and Hinton, G. E. Deep learning. *Nature*, 521(7553):436-444, 2015. doi: 10. 1038/nature14539. URL https://doi.org/10.1038/ nature14539.
- Loula, J., Baroni, M., and Lake, B. Rearranging the familiar: Testing compositional generalization in recurrent networks. In *Proceedings of the EMNLP BlackboxNLP Workshop*, pp. 108–114, Brussels, Belgium, 2018.
- Mei, H., Bansal, M., and Walter, M. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *Proceedings of AAAI*, pp. 2772– 2778, Phoenix, AZ, 2016.
- Nye, M. and Tenenbaum, J. B. Learning compositional rules via neural program synthesis. 2019.

- Pustejovsky, J. The generative lexicon. Computational Linguistics, 17(4):409-441, 1991. URL https://www.aclweb.org/anthology/J91-4003.
- Russin, J., Jo, J., O'Reilly, R. C., and Bengio, Y. Compositional generalization in a deep seq2seq model by separating syntax and semantics. *CoRR*, abs/1904.09708, 2019. URL http://arxiv.org/abs/1904.09708.
- Schuster, M. and Paliwal, K. Bidirectional recurrent neural networks. *Trans. Sig. Proc.*, 45(11):2673-2681, November 1997. ISSN 1053-587X. doi: 10.1109/ 78.650093. URL http://dx.doi.org/10.1109/78. 650093.
- Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to sequence learning with neural networks. *CoRR*, abs/1409.3215, 2014. URL http://arxiv.org/abs/ 1409.3215.
- Wang, Z. and Lake, B. M. Modeling question asking using neural program generation. CoRR, abs/1907.09899, 2019. URL http://arxiv.org/abs/ 1907.09899.

	$\mathrm{VV_i} \rightarrow$	$\{walk\}$
$\mathrm{ROOT} \to \mathrm{VP}$	$\mathrm{VV}_{\mathrm{t}} \rightarrow$	$\{push, pull\}$
$\mathrm{VP} \to \mathrm{VP} \ \mathrm{RB}$	$\mathrm{RB} \rightarrow$	{while spinning,
$\mathrm{VP} \rightarrow \mathrm{VV_i}$ 'to' DP		while zigzagging,
$\mathrm{VP} \to \mathrm{VV_t}~\mathrm{DP}$		hesitantly,
$\mathrm{DP} \rightarrow `a' \mathrm{NP}$		cautiously
$\mathrm{NP} \to \mathrm{JJ} \ \mathrm{NP}$	$\rm NN \rightarrow$	$\{$ circle, square,
$\rm NP \rightarrow \rm NN$		$\operatorname{cylinder}$
	$\rm JJ \rightarrow$	$\{{\rm red}, {\rm green}, {\rm blue},$
		$big, small \}$

A. The CFG to Generate Input Commands

Where a subscript of 'i' refers to *intransitive* and of 't' to *transitive*.

B. World State Generation

We generate the world state for each instruction with two constraints: (1) there must be a unique target object for the referrent, (2) if a size modifier is present in the instruction, there must be at least one 'distractor' object present. In other words, if the target referrent is 'the small square', we additionally place a larger square than the target of any color, and if the target referrent is 'the small *yellow* square', we additionally place a larger *yellow* square. For each world generation we select at random half of the possible objects to place, where we make sure that if a size is mentioned, we always put a pair of differently sized objects for each color-shape pair. There are three different situations for which we generate a different set of objects, depending on whether a shape, color and/or size is mentioned.

1. Shape (e.g. 'the circle'):

Generate 1 randomly colored and randomly sized object of each shape that is not the target shape, randomly select half of those. See top-left of Figure 7.

- 2. Color and shape (e.g. 'the red circle'): Generate 1 randomly sized object of each color and shape pair that is not the target color and shape. See bottom-left of Figure 7.
- 3. *Size, color, and shape* (e.g. 'the small circle', 'the red small circle'):

Generate 2 randomly sized objects for each color and shape pair, making sure the size for the objects that are the same shape (and color if mentioned) as the target are smaller/larger than the target dependent on whether the size modifier is big/small. Select at random half of the pairs. See top- and bottom-right of Figure 7. We generate world states for instructions by generating all combinations of possible target objects based on the determiner phrase (e.g. 'the yellow square' gives 4 possible targets, namely a yellow square of each size), all possible relative directions between the agent and the object (e.g. the agent can be to the north-east of the target) and all possible distances between the agent and the object (e.g. if the target is to the north-east of the agent the minimum number of steps is 2 and the maximum is $2 \cdot d$, but if the target is to the north, the minimum is 1 and the maximum d). We randomly sample, for each of the combinations, possible positions of the agent and target and this gives us full dataset.



Figure 7: Four real data examples, for a grid size of 6.

C. Dataset Statistics

Table 3: Number of examples in each dataset. The first two rows for the train and test set denote data for the compositional splits, where GECA denotes the augmented set. An example is unique if it has a different input command, target commands or target object location. The third column denotes the number of unique world states the agent sees on average per input command.

		Unique Examples				
Train	Examples	Total	Per Command			
Compositional	$367,\!933$	76,033	177			
GECA	377,933	64,527	186			
Target Length	180,301	80,865	599			
Test	Examples	Unique	Per Command			
Compositional	19,282	16,381	38			
GECA	19,282	16,381	38			
Target Length	37,784	31,000	230			

If we dive a bit deeper in the augmentations to the data GECA makes, it becomes clear why performance deteriorates when this data augmentation technique is used in a dataset like grounded SCAN, and why it is not able to improve performance on the 'yellow squares'-split (Section 5B). In Table 4 we can see that although GECA correctly identifies red square target objects as missing in the training data, and augments data examples to contain target red squares, it is not able to do the same for the instruction command. Even though red squares are now seen as target objects, they are still always referred to without the color (e.g. 'the small square'). See for reference the row with blue squares as targets, which occurs in the non-augmented dataset and is therefore less affected by GECA. Also for yellow squares as target object GECA is not able to add the needed reference to the color. Additionally, when looking at the other row in Table 4, GECA also caused the references to 'red circles' to completely disappear. It seems that it is non-trivial to augment grounded SCAN with GECA to obtain improved performance on anything other but a narrow part of the test set.

Table 4: Some dataset statistics for the compositional splits. In each row the number of examples for a given target object is shown. The column 'placed' means that object was placed in the world as the target, and the column 'referred' means that it was referred to with the color (e.g. 'the red square' or 'the small yellow square')

	Non-au	igmented	Augr	nented
	Placed	Referred	Placed	Referred
Blue Squares	33,250	16,630	16,481	$5,\!601$
Red Squares	0	0	83,887	0
Yellow Squares	16,725	0	10,936	0
Red Circles	33,670	16,816	16,854	0

D. A Forward-Pass Through the Model

A full forward pass through the model can be represented by the following equations.

Encoder

Command encoder $\mathbf{h}^c = f_c(\mathbf{x}^c)$: $\forall i \in \{1, \dots, n\}$

$$\mathbf{e}_{i}^{c} = \mathbf{E}_{c}(x_{i}^{c})$$
$$\mathbf{h}_{i}^{c} = \mathrm{LSTM}_{\phi_{1}}(\mathbf{e}_{i}^{c}, \mathbf{h}_{i-1}^{c})$$

State encoder $\mathbf{H}^s = f_s(\mathbf{X}^s)$:

$$\mathbf{H}^{s} = \operatorname{ReLU}([\mathbf{K}_{1}(\mathbf{X}^{s}); \mathbf{K}_{5}(\mathbf{X}^{s}); \mathbf{K}_{7}(\mathbf{X}^{s})])$$

Where $\mathbf{E}^c \in \mathbb{R}^{|V_c| \times d_c}$ is the embedding lookup table with V_c the input command vocabulary, d_c the input embedding dimension, and \mathbf{K}_k the convolutions with kernel size k. We pad the input images such that they retain their input width and height, to enable visual attention. That

means the world state features are $\mathbf{H}^s \in \mathbb{R}^{d \times d \times 3c_{\text{out}}}$, with c_{out} the number of channels in the convolutions. We then decode with an LSTM whose initial state is the final state of the encoder LSTM. The input to the decoder LSTM is a concatenation of the previous token embedding, a context vector computed by attention over the hidden states of the encoder, and a context vector computed by conditional attention over the world state features as processed by the CNN. A forward pass through the decoder can be represented by the following equations.

Decoder $p(\mathbf{y}|\mathbf{h}^c, \mathbf{H}^s)$: $\forall j \in \{1, \dots, m\}$

$$\begin{aligned} \mathbf{h}_{0}^{d} &= \mathbf{W}_{p} \mathbf{h}_{n}^{c} + \mathbf{b}_{p} \\ \mathbf{e}_{j}^{d} &= \mathbf{E}^{d}(y_{j-1}) \\ \mathbf{h}_{j}^{d} &= \mathrm{LSTM}_{\phi_{2}}([\mathbf{e}_{j}^{d};\mathbf{c}_{j}^{c};\mathbf{c}_{j}^{s}],\mathbf{h}_{j-1}^{d}) \\ \mathbf{o}_{j} &= \mathbf{W}_{o}[\mathbf{e}_{j}^{d};\mathbf{h}_{j}^{d};\mathbf{c}_{j}^{c};\mathbf{c}_{j}^{s}] \\ \hat{\mathbf{o}}_{j} &= p_{\theta}(y_{j} \mid \mathbf{x},y_{1},\ldots,y_{j-1}) = \mathrm{softmax}(\mathbf{o}_{j}) \\ \hat{y}_{j} &= \arg\max_{V_{t}}(\hat{\mathbf{o}}_{j}) \end{aligned}$$

Textual attention $\mathbf{c}_{j}^{c} = \operatorname{Attention}(\mathbf{h}_{j-1}^{d}, \mathbf{h}^{c}):$ $\forall i \in \{1, \dots, n\}$

$$\begin{split} e_{ji}^{c} &= \mathbf{v}_{c}^{T} \tanh \mathbf{W}_{c}[\mathbf{h}_{j-1}^{d}; \mathbf{h}_{i}^{c}] \\ \alpha_{ji}^{c} &= \frac{\exp(e_{ji}^{c})}{\sum_{i=1}^{n} \exp(e_{ji}^{c})} \\ \mathbf{c}_{j}^{c} &= \sum_{i=1}^{n} \alpha_{ji}^{c} \mathbf{h}_{i}^{c} \end{split}$$

Conditional visual attention $\mathbf{c}_{j}^{s} = \operatorname{Attention}([\mathbf{c}_{j}^{c}; \mathbf{h}_{j-1}^{d}], \mathbf{H}^{s})$ $\forall k \in \{1, \dots, d^{2}\}$

$$e_{jk}^{s} = \mathbf{v}_{s}^{T} \tanh \mathbf{W}_{s}[\mathbf{h}_{j-1}^{d}; \mathbf{c}_{j}^{c}; \mathbf{h}_{k}^{s}]$$
$$\alpha_{jk}^{s} = \frac{\exp(e_{jk}^{s})}{\sum_{k=1}^{d^{2}} \exp(e_{jk}^{s})}$$
$$\mathbf{c}_{j}^{s} = \sum_{k=1}^{d^{2}} \alpha_{jk}^{s} \mathbf{H}_{k}^{s}$$

Where $\mathbf{W}_c \in \mathbb{R}^{h_d \times (h_e + h_d)}, \mathbf{v}_c \in \mathbb{R}^{h_d}, \mathbf{W}_s \in \mathbb{R}^{h_d \times (3c_{\text{out}} + h_d)}, \mathbf{v}_s \in \mathbb{R}^{h_d}, \mathbf{W}_p \in \mathbb{R}^{h_d \times h_e}, \mathbf{W}_o \in \mathbb{R}^{|V_t| \times (d_e + 3h_d)}$, with h_e the hidden size of the encoder and h_d of the decoder, c_{out} the number of channels in the encoder CNN, V_t the target vocabulary, and d_e the target token embedding dimension. All model parameters are $\theta = \{\mathbf{E}_c, \phi_1, \mathbf{K}_1, \mathbf{K}_5, \mathbf{K}_7, \mathbf{W}_c, \mathbf{v}_c, \mathbf{W}_s, \mathbf{v}_s, \mathbf{W}_p, \mathbf{E}_d, \phi_2, \mathbf{W}_o\}$.

E. (Hyper)parameters

Experiment G

Compute power used

Training a model on the data used for the experiments with a single GPU takes less than 24 hours.

Table 5:	Parameters	\mathbf{for}	$_{\rm the}$	models	not	explicitly	mentioned	ir
Section 4.								

Situation Enc.	Value	Decoder	Value
$c_{ m out}$	50	d_e	25
Dropout p on \mathbf{H}^s	0.1	LSTM Layers	1
Command Enc.	Value	h_d	100
d_c	25	Dropout p on \mathbf{E}^d	0.3
LSTM Layers	1	Training	Value
h_e	100	β_1	0.9
Dropout p on \mathbf{E}^c	0.3	β_2	0.999

Additional Parameters for the GECA model

Gap size: 1, maximum gaps: 2.



Figure 9: The exact match decreases when target length increases for the adverb split where the agent needs to generalize 'cautiously' after seeing 50 demonstrations. Note that for this experiment, the tested target lengths are *not* longer than encountered during training.

Experiment H



F. Additional Results

This section contains additional experimental results that were referred to in the main text in Section 5.

Experiment D



Figure 8: Visualized predictions from the models, where a darker grid means a higher attention weight for that cell.

Figure 10: The exact match decreases when target length increases for the adverb split where the agent needs to generalize 'while spinning' to the verb 'pull'. Note that for this experiment, the tested target lengths are *not* longer than encountered during training.