# Learning to Ground Language to Temporal Logical Form

### Anonymous EMNLP-IJCNLP submission

### Abstract

Often times, natural language commands exhibit sequential constraints, e.g. "go through the kitchen and then into the living room." Conventional methods in classical Reinforcement Learning (RL) assume Markovian reward functions that cannot handle such non-Markovian constraints. We therefore propose to ground natural language commands to Linear Temporal Logic (LTL), but propose to do this when direct supervision is not available - utterances are labelled with trajectories but not the logical form itself. We use formal methods of LTL progression and model checking, to reward the learned latent logical forms that execute correct trajectories. These logical forms can then serve as input to planners that correctly solve the navigation task. We evaluate our framework on both goal state and path accuracy, as well as the ability to handle temporal language.

### 1 Introduction

In navigation and instruction-following tasks, an agent executes a series of actions to perform a task, in response to some task specification that defines the goal that the agent must reach. This is challenging for two reasons — first, the broad spectrum of possible language used must be grounded to a symbolic, logical representation, that can then be given to a model to find a policy that solves the task. Second,



Figure 1: Example path and instruction grounded to temporal logical form for one of the environments from (MacMahon et al., 2006). Letters indicate object positions, while colours indicate floor patterns and wall paintings, that all form propositions for LTL representations. Agents are required to correctly interpret the meaning of the instruction to navigate to the intended goal location.

the complex (e.g., temporal, sequential, conditional) constraints are non-trivial to express and achieve by conventional methods in classical Reinforcement Learning (RL), that assume a Markovian reward function. For example, commands that specify temporal ordering of events such as "only enter the living room after you have been to the kitchen" could potentially map to unbounded action sequences, unless we have adequate representation of state properties in accordance with temporal patterns over time.

We propose to ground natural language instructions to LTL task specifications that can

084 handle temporal order and furthermore allow 085 the logical formula to be structured and de-086 composed into paths in an environment. Moreover, we attempt to do this by learning from 087 demonstrations. While collecting and anno-088 tating symbolic logical forms is an expensive 089 process, we can more easily collect execu-090 tions of these logical forms as trajectories in 091 our environment. Our system takes in natural 092 language as input, grounds to symbolic tem-093 poral logical form, and checks this against a 094 human-annotated output trajectory as a form 095 of weak supervision. Our work is novel in 096 that we ground to LTL logical form (that can 097 handle temporal constraints) and moreover al-098 lows using methods like LTL progression, to 099 check satisfiability against trajectories in order 100 to reward the learned latent LTL forms. In 101 comparison to previous approaches that build 102 execution models or planners in order to re-103 ward logical forms, our approach is an order 104 of magnitude more efficient in terms of com-105 plexity, given that we only need to progress 106 an LTL expression against a fixed-length tra-107 jectory, rather than execute an expression and 108 solve a planning problem each time. We eval-109 uate our method on a benchmark dataset and 110 highlight the benefits of using a more expressive language, such as LTL, that can handle 111 temporal order. 112

## 2 Related Work

113

114

115 Previous work has used LTL to formulate tasks 116 in RL, either by creating reward functions 117 that maximise the probability of satisfying the 118 LTL formula (Wen et al., 2017; Littman et al., 119 2017) or by guiding policy search with a mea-120 sure of distance to satisfaction of the task (Li 121 et al., 2017). Other work exploits the struc-122 ture of LTL to decompose a task into subtasks 123 (Toro Icarte et al., 2018) to deal with tem-124 poral abstraction. To the best of our knowl-125 edge, there is currently no work that attempts to ground tasks given in natural language to temporal logic without annotation of logical forms, to exploit LTL structure for further planning. 126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

Language  $\rightarrow$  Logical Form There is ample prior work on supervised semantic parsing to LTL (Gopalan et al., 2018) and to logical forms other than LTL (Zettlemoyer and Collins, 2012; Tang and Mooney, 2000; Berant et al., 2013; Yu et al., 2018). Relevant to our work is past work which learns semantic parsers without explicit annotation of logical forms, by allowing the execution of the learned logical form to act as supervision e.g., conversational logs (Artzi and Zettlemover, 2011), system demonstrations (Chen et al., 2015; Goldwasser and Roth, 2014; Artzi and Zettlemoyer, 2013; Williams et al., 2018) and question-answer pairs (Clarke et al., 2010; Liang et al., 2013).

Language  $\rightarrow$  Plan Also relevant is the body of recent work which seeks to map natural language directly to action sequences, e.g. Mei et al. (2016); Misra and Artzi (2015); Branavan et al. (2009); Misra et al. (2017); Fried et al. (2018). Such methods are typically trained end-to-end, and do not pass through an explicit intermediate logical form, as we do in this work.

Logical Form (LTL)  $\rightarrow$  Plan (Dzifcak et al., 2009; Gopalan et al., 2018) explore grounding language to LTL and then planning with an existing planner, given LTL task specifications. However (Gopalan et al., 2018) use a standard sequence-to-sequence model trained on annotated utterances, and they note that this method fails to exhibit compositionality or generalise to unseen logical forms. Several recent approaches use Deep Q-Networks with LTL specifications (Toro Icarte et al., 2018), by making using of LTL based rewards (Littman et al., 2017) where the input to the Q-value function is both the state and the progressed LTL

168task. Other work uses hierarchical RL methods169(Sutton et al., 1999; Kulkarni et al., 2016) in170the options framework, by creating one option171per proposition with terminal states defined by172states in which the proposition is true; giving173the state and progressed LTL task as input to a174meta-controller.

## **3** Grounding Language to Linear Temporal Logic with Trajectories

We propose to ground natural language to tem-poral logical form, where grounded LTL ex-pressions are only supervised by trajectories in the environment. Our model therefore receives natural language instructions (as a sequence of words) as input and is supervised by an exe-cuted path (as a sequence of points specifying (x, y) locations) in the environment. Once the natural language instruction has been grounded to an LTL expression, this can be used to plan in a given environment using existing methods such as Q-learning with LTL-based rewards. This entire pipeline therefore ensures that the intended meaning of the instruction is first rep-resented in logical form (that satisfies temporal constraints) and then used to navigate in the en-vironment (to reach the final goal location). In this section we give an overview of the differ-ent components that our approach relies on. 

### 3.1 Linear Temporal Logic

**LTL Syntax:** LTL has the following grammatical syntax:

$$\phi ::= \pi |\neg \phi| \phi \land \varphi | \phi \lor \varphi | \diamond \phi | \Box \phi | \bigcirc \phi | \phi \cup \varphi$$

where the operators  $\neg$ ,  $\land$ ,  $\lor$  are the logical connectives for *negation*, *and*, *or* and the temporal operators are  $\diamond$  for *eventually*,  $\Box$  for *globally*,  $\bigcup$  for *until* and  $\bigcirc$  for *next*. We can also define the symbols true and false through the following equivalences: true  $\equiv \phi \lor \neg \phi$  and false  $\equiv \neg$  true. Our set of propositions

*P* consists of observable elements in the environment that trajectories can pass through e.g (at\_object, is\_intersection, is\_corridor). All LTL expressions are constructed from the set of propositional symbols *P* and the extended set of operators defined above i.e., the Boolean operators  $\land, \lor, \neg$  and the temporal operators  $\bigcirc$ , U. From these we can define  $\Box$  (*always*) and  $\diamond$  (*eventually*) for e.g.,  $\diamond \phi$  = true U  $\phi$ .

**LTL Semantics:** Given the observable elements in the environment that form atomic propositions, the truth value of an LTL formula is determined relative to a sequence of truth assignments  $\sigma = \langle \sigma_1, \sigma_2, \sigma_3, ... \rangle$  where each state  $\sigma_i$  is an assignment of true or false values to propositions. A proposition  $\rho \in \sigma_i$  indicates that the proposition  $\rho$  is true in the state  $\sigma_i$ .

**LTL Progression:** Given a sequence of truth assignments and an LTL task specification, an LTL formula can be *progressed* along the sequence. For example, the task  $\diamond(p \land \bigcirc \diamond q)$  (i.e., eventually p and eventually q) can be progressed to  $\diamond q$  (i.e., eventually q) once the agent reaches a state where p is true. In an RL setting, the LTL formula can be updated to reflect the agent's actions and the parts of the formula that have been satisfied so far.

$prog(\sigma_i, p)$	true if $p \in \sigma_i$ , where $p \in P$
$prog(\sigma_i, p)$	false if $p \notin \sigma_i$ , where $p \in P$
$prog(\sigma_i, \neg \phi)$	$ eg prog(\sigma_i, \phi)$
$prog(\sigma_i, \phi_1 \land \phi_2)$	$prog(\sigma_i,\phi_1)\wedgeprog(\sigma_i,\phi_2)$
$prog(\sigma_i, \phi_1 \lor \phi_2)$	$prog(\sigma_i,\phi_1) \lor prog(\sigma_i,\phi_2)$
$prog(\sigma_i, \bigcirc \phi)$	$\phi$
$prog(\sigma_i, \phi_1  U  \phi_2)$	$prog(\sigma_i, \phi_2) \lor (prog(\sigma_i, \phi_1))$
	$\wedge \phi_1 \cup \phi_2)$

Table 1: Semantics of progression functions for different logical and temporal operators. The progression function takes in the current state and LTL formula that is updated after application of the function.

For an LTL expression  $\phi$  and a state  $\sigma_i$ , we define the semantics of the progression func-

252 tion  $prog(\sigma_i, \phi)$  as in Table 1. At each point 253 in time i, we can therefore update the LTL ex-254 pression to reflect which parts of the original formula have so far been satisfied or unsatis-255 fied. We can do this because if a sequence of 256 truth assignments (i.e., a trajectory) satisfies an 257 LTL formula at time *i* if the formula progressed 258 through  $\phi_i$  is then satisfied at time i + 1. 259

260 Program Representation: Formally, we 261 represent LTL expressions (i.e., executable 262 programs) as a sequence of tokens that 263 describe a possibly recursive sequence of 264 functions in postfix notation. Every token 265 in an expression is either a function of fixed 266 arity (i.e., one or two arguments), constant, 267 variable or  $\lambda$  term that is used to define 268 Boolean functions. Previous work that parses to logical representations uses atomic types 269 such as Boolean, Integer, Size, 270 Shape, Colour, Side or composite 271 types Set(?), Func(?, ?). For our 272 purpose, since our LTL task specifications 273 specify goal locations in the environment, 274 our LTL vocabulary consists of observable 275 elements in the environment of atomic types 276 (Object, Floor, Wall) and either unary 277 or binary logical operators (Func). All valid 278 and syntactically correct programs have a 279 return type of Boolean. For example, the 280 LTL expression  $(a \land (\diamond b))$  converted to 281 postfix notation (a  $b \diamond \wedge$ ) is linearised to 282 allow easier execution with a stack based on 283 the semantics of operators and propositions. 284

### 3.2 Planning in an Environment

285

286

287

288

289

290

291

292

293

**Markov Decision Processes:** A Markov Decision Process (MDP) is a tuple  $M = (S, A, T, R, \gamma)$  where S and A are a finite sets of states and actions respectively,  $T : S \times$  $A \times S \rightarrow [0, 1]$  is the transition function,  $R : S \times A \times S \rightarrow \Pr(R)$  is the reward function and  $\gamma$  is the discount factor. An agent can take an action (from the set of all actions) and

World	Operation	Complexity		
Database	execution	$O( W  \times  s )$		
SQL:SELECT*MAX(state.area) FROM state				
CCG parse: $\lambda x. \texttt{flight}(x) \land \texttt{from}(x, bos)$				
Gridworld	progression execution	$\begin{array}{c} O( S \times  s ) \\ O( S ^{ s }) \end{array}$		
LTL: at_lamp U ◊ (at_chair)				
CCG: $\lambda$ a $(\lambda x.chai)$	a.move(a) $\land$ post(a, r(x),you)) $\land$ pre(a,f $\lambda$ x.lamp(x)))	intersect Front (you,		

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

334

335

Table 2: Comparing complexities of different feedback systems for different world representations and logical forms. The first rows show example execution models in other domains i.e., SQL query or a CCG parse (lambda calculus expression) executed against a database, while the latter section compares feedback methods in navigation domains i.e., progression vs. execution or planning for LTL forms or a CCG parse.

change its state, to move within states according to the transition function, hence accumulating reward according to the reward function. An agent learns a *policy*  $\pi$  i.e., a probability distribution over state and action pairs, that allows it to determine the actions it should take in each state with probability  $\pi(a|s)$ . The stateaction value or Q-value denoted  $Q^{\pi}(s, a)$  is the expected discounted return of selecting action a in state s and then selecting actions according to  $\pi$ , and a policy is *optimal* if the expected discounted reward gained by following the pilicy is maximal for every state  $s \in S$ . Given the Q-value function  $Q^*$ , the optimal policy is to then select the action a in every state s with the highest value of  $Q^*(s, a)$ .

**Q-Learning:** The Q-learning algorithm is a well-known off-policy algorithm, that learns a *target* policy while using some other *behaviour* policy for action selection. We begin the Q-learning process by initialising the Q-values of all state-action pairs to zero, and at every step

336 use some behaviour policy to pick an action *a* 337 for the current state *s*, therefore returning a new 338 state *s*' an reward *r* from the environment. The 339 estimation of Q(s, a) at the current timestep is 340 then updated as shown in the equation below 341 where  $\alpha$  is the *learning rate*.

342

343 344

345

346

347

348

349

350

351

352

353

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

The Q-learning algorithm is guaranteed to converge to optimal Q-values if the behaviour policy visits every state-action pair an infinite number of times. In practice, this is done by setting the behaviour policy to be  $\epsilon$ -greedy on the target policy, to ensure that on each step, the behaviour policy selects random actions with probability  $\epsilon$  and the action with the highest Qvalue with probability  $1 - \epsilon$ .

354 Function Approximation and DON: For 355 the Q-learning algorithm to keep track of state-356 action pairs, in the simplest form of Q-learning, 357 a Q-table is used to store the required values. 358 This however, is impractical for large (or in-359 finite) state spaces. Therefore, we instead use 360 some form of function approximation on the Q-361 value function i.e., we define the Q-value func-362 tion as a function of state features, where Q-363 value updates involve updating the function, instead of updating entries in a Q-table. Deep Q-364 365 Networks (DQN) (Mnih et al., 2015) use neural networks for this function approximation of the 366 Q-value function. An experience replay buffer 367 and a target network are used to stabilise learn-368 ing and the agent's experiences (of the form 369 (s, a, r, s')) are stored in the bugger and sam-370 pled to train the network over time. The Q-371 learning updates are computed with respect to 372 a target network which is only updated period-373 ically to attempt to decrease the chance of the 374 policy diverging. 375

376Learning with Temporal non-Markovian377Rewards: In order to allow handling of tem-

poral constraints, we use the concept of a non-Markovian reward decision process that incorporates the LTL expression into the environment MDP. Instead of only considering the previous state, the reward function R is therefore defined over state histories  $R : S^* \rightarrow Pr(R)$ . Given this MDP, the Q-value function of a policy  $\pi$  is then defined over sequences of states:

$$\mathbf{Q}^{\pi}(\langle s_0, .., s_t \rangle, a) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R(\langle s_0, .., s_{t+k+1} \rangle) | A_t = a \right]$$

The reward functions for these MDPs are defined in terms of completing the task defined by the LTL formula. We incorporate the LTL task specification through a labelling function  $L: S \rightarrow 2^P$  where P is the set of propositional symbols. This is therefore a 7-tuple  $\tau = \langle S, A, T, P, L, \phi, \gamma \rangle$ , where S, A, T and  $\gamma$  are defined as in an MDP and refer to states, actions, transitions and the discount factor respectively. P is the set of propositional symbols that form LTL task specifications, L is the labelling function defined above and  $\phi$  is the set of tasks. As in general RL settings, the agent does not know the transition probability distribution of the domain, but has access to the labelling function and the task.

### **3.3** Environments and Data

SAIL from (MacMahon et al., 2006) is a navigation dataset containing route instructions annotated with trajectories for three different environments, each composed of connected hallways with different patterns (grass, brick, wood, gravel, blue, flower, or yellow octagons) on floors, paintings (butterfly, fish, or Eiffel Tower) on walls and objects (hat rack, lamp, chair, sofa, barstool, and easel) at intersections. Figure 1 shows an example environment annotated with object positions, wall paintings and floor patterns.

The challenge of learning to ground natural language to logic stems from the fact that instructions given by humans are complex, free-

### Confidential Review Copy. DO NOT DISTRIBUTE.

420	No. of instruction sequences	706
421	No. of sentences	3233
	Avg. sentences per sequence	4.61
422	Avg. tokens per sentence	7.94
423	Vocabulary size	522
104	Avg. trajectory length (per sentence)	3.05
424	Max. trajectory length (per sentence)	31
425		I

Table 3: Statistics from the SAIL dataset from (MacMahon et al., 2006) that contains natural language instructions annotated with trajectories.

430 form and of variable length (single sentences 431 in isolation or full paragraphs). While this task 432 and dataset bears superficial similarity to other 433 navigation benchmarks, the language and poli-434 cies required for this task are quite different — the proportion of instructions to actions is 435 much higher, the interpretation of language is 436 highly compositional and the length of the in-437 structions vary widely. This task has therefore 438 been the subject of focused attention in seman-439 tic parsing, resulting in a range of different ap-440 proaches that attempt to learn how to plan in 441 such settings. 442

#### 3.4 Model

426

427

428

429

443

444

447

Our model draws insights from previous work 445 (Guu et al., 2017; Goldman et al., 2017) that 446 train semantic parsers from denotations with an algorithm that searches through a space of 448 programs at training time, in order to find the 449 correct program. We set this up as a reinforce-450 ment learning problem, where our agent sam-451 ples a token from the vocabulary to generate 452 a sequence of tokens that in our case, form an 453 LTL expression. At the end of this sequential 454 prediction of tokens, each LTL expression is 455 given a binary reward of 1 or 0, by progress-456 ing the LTL expression along the trajectory in 457 the environment.

458 Our training samples are of the form (x, t)459 where x is a natural language instruction and 460 t is a trajectory in the environment, and our 461 model generates program tokens  $z_1, z_2$ .. from left to right using a neural encoder-decoder model (Sutskever et al., 2014). We encode every utterance x with a bidirectional LSTM (Hochreiter and Schmidhuber, 1997) to create a contextualised representation  $h_i$  for ev-Our decoder is then ery input token  $x_i$ . a feed-forward network with attention (Bahdanau et al., 2014) over the output from the encoder, that takes as input, the last K tokens that were decoded. Formally, the probability of a decoded LTL expression is the product of the probability of its tokens conditioned on the history i.e.,  $p_{\theta}(z|x) = \prod_{t} p_{\theta}(z_t|x, z_{1:t-1})$  and the probability of a decoded token comes from the learned parameters and embedding matrices as shown in the equations below.

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

494

495

496

497

498

499

500

501

502

503

**Encoder** The utterance  $u_m$  is encoded with a bidirectional LSTM:

$$h_i^F = LSTM(h_{i-1}^F, \phi_u(u_m, i))$$
  

$$h_i^B = LSTM(h_{i-1}^B, \phi_u(u_m, i))$$
  

$$h_i = [h_i^F; h_i^B]$$

**Decoder** Let  $f(z_{1:t-1})$  refer to the execution history and  $e_m$  to the input embedding to the decoder. We can then compute an attention vector  $c_t$  as follows:

$$q_t = ReLU(W_q[e_m; f(z_{1:t-1})])$$
  

$$\alpha_i \propto \exp(q_t^T W_a h_i) i - 1, ..., |u_m|)$$
  

$$c_t = \sum_i \alpha_i h_i$$

After combining the above i.e., concatenating  $q_t$  with  $c_t$ , we then produce the program token by token. Formally, we produce a distribution over the set Z of program tokens, computed with a softmax, as shown in the equation below, where  $\phi_z(z_t)$  is the embedding for a token  $z_t$ .

$$p_{\theta}(z_t|x, z_{1:t-1}) \propto \exp(\phi_{z_t}^T W_s[q_t; c_t])$$

### Confidential Review Copy. DO NOT DISTRIBUTE.



Figure 2: Skeleton of our weakly supervised model that learns to compose a logical meaning representation by only giving binary feedback signal through LTL progression against a ground truth trajectory. This model can then produce intermediate logical forms for any instruction that can then be given to a planner that solves the task to find a path.

We also keep track of the *execution history* i.e., the k most recent tokens  $z_{t-k:t-1}$  and concatenate their embeddings. Our exploration strategy is the same as in (Guu et al., 2017) that uses an  $\epsilon$ -greedy randomised beam search to deal with the problem of biased exploration that can result in visiting only a few states during exploration. Like regular beam search, at each iteration, the set of all continuations for the next token are sorted by the model probability  $p_{\theta}(s|x)$ . However, we inject random noise into exploration. Instead of simply selecting the k highest scoring continuations (for a beam of size k), we choose continuations one by one without replacement, from the set of all possible continuations. At each selection step from the remaining pool, we either unformly sample a random continuation with probability  $\epsilon$  or pick the highest scoring continuation from the pool with probability  $1 - \epsilon$ .

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

541

544

538 Figure 2 shows a skeleton of our proposed 539 model that is trained to compose together an 540 LTL formula for the input natural language utterance, supervised by a trajectory, by pro-542 gressing the LTL expression along this trajec-543 tory in the environment. This gives us a score of the satisfiability of the produced LTL ex-545 pressions that allows ranking and learning of logical meaning representations.

For the weakly supervised model, we perform a grid search over hyperparameters to maximise accuracy over the validation set. These hyperparameters include the learning rate, the value of  $\epsilon$  and the embedding size.

Planning Mechanism: To plan in a given environment with the grounded LTL task specification, we use Deep Q-Networks with LTL specifications, following the methodology from (Littman et al., 2017) to learn LTL-based rewards. We use standard RL to solve the cross-product MDP where the input to the Qvalue function is both the state and the progressed LTL task at every step. This is the state-of-the-art approach to learn with reward functions specified in LTL to allow us to deal with temporal non-Markovian constraints.

The DON implementation is baed on the OpenAIBaselines (Dhariwal et al., 2017). To train the network, we use experience replay minibatches of experiences are sampled from an experience replay buffer of a fixed size. The feature vector at each state is computed based on the distance of every object from the agent. We use feedforward networks with 2 hidden layers and 64 ReLu units, trained us-

588 ing an Adam optimiser with a learning rate of 589 0.0001. At every step, the DQN network learns 590 by randomly sampling 32 transitions from the replay buffer. We set the size of the buffer to 591 25,000 and update the target network at every 592  $100^{th}$  step and use a discount factor of 0.9. 593

#### 4 **Experimental Evaluation**

594

595

601

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

596 To evaluate models on the SAIL dataset, previ-597 ous works compare the agent's end state to a la-598 belled state s' i.e., the end point of the groundtruth trajectory, for all instructions in the test 599 600 set. Models are tested on both the single and multi-sentence instructions to assess their ability to correctly interpret the intent of the in-602 struction and navigate to the goal state in the 603 environment. 604

**Comparison to prior work:** We compare our model to existing work that report results on the SAIL dataset, by cross-validating over the three maps and reporting a final accuracy. Most similar to our setup is the approach from (Artzi and Zettlemoyer, 2013) that train a CCG semantic parser with supervision from trajectories, however they do this by manually building execution models for the feedback mechanism which is more computationally intractable ---the worst case complexity is exponential in the number of variables.

System	Single	Multi
Chen and Mooney (2011)	54.4	16.18
Chen (2012)	57.28	19.18
+ additional data	57.62	20.64
Kim and Mooney (2012)	57.22	20.17
Artzi and Zettlemoyer (2013)	65.28	31.93
Andreas and Klein (2015)	59.60	-
Mei et. al. (2017)	71.05	30.34
Ours	66.92	20.17

Table 4: Evaluation of systems on the SAIL dataset. The second and third columns shows accuracies of the reaching the goal state for commands that take the form of single sentences or entire paragraphs.

Goal-state accuracy To directly compare our approach with existing work, measure the goal-state accuracy i.e., the ability of the model to reach the correct final location in the environment after correctly interpreting and grounding the natural language instruction. Specifically, for each instruction in the test set, our models produces an LTL logical form, given to the DQN-LTL model that performs Q-learning to solve the navigation task to reach the final goal state. As shown in table 4 we compare to other work that is different in the form of supervision provided and model architecture used, but evaluates on the same end-task i.e., final goal-state accuracy over the SAIL dataset. These include algorithms supervised with logical forms (Chen and Mooney, 2011; Chen et al., 2015) that learn semantic parsers for natural language instructions as well as ones that involve strategies for online learning of lexicons for the semantic parsers (Chen, 2012) and ones that use contextual information (Chen et al., 2010) for better language understanding. We also compare models supervised with paths — the supervised alignment-based models (Andreas and Klein, 2015) that build grounding graphs representations to execute instructions and neural sequence-to-sequence models (Mei et al., 2016) that translate natural language instructions to actions that an agent can execute in the environment.

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

662

663

664

665

666

667

668

669

670

671

**Path accuracy** Unlike prior work, we also propose a more fine-grained analysis to evaluate *path accuracy*. Often times, the path taken to reach the final goal location is crucial - especially when the instruction specifies constraints on how to reach the goal, in complicated environments with several possible paths. Figure

#### 5 Conclusion

We propose a weakly supervised semantic parsing model that requires no supervision

### Confidential Review Copy. DO NOT DISTRIBUTE.

672 673 0 0 0 00 674 0 0 • • • 0 675 • 676 677 678 679 walk past both the lamp and the easel until you get to the red brick corridor • = lamp • = easel • = brick corridor 680 681 Figure 3: Example possible paths that require more fine-grained evaluation than just final goal location. 682 While both paths in the figure reach the correct goal 683 location (and get perfect accuracy under that met-684 ric) only the one on the left follows the constraints 685 specified by the instruction. 686 687 System Prec. Recl. Acc. 688 Mei et. al. (2017) 31.2 30.19 91.2 689 Ours 33.5 35.2 94.3 690 691 Table 5: Evaluating paths produced by systems 692 against ground-truth paths. Each path is represented 693 as a masked vector of locations in the environment, 694 thus allowing comparing of precision, recall and ac-695 curacy against the gold path. 696 697 of ground-truth logical forms during training. 698 Moreover, we propose to ground natural lan-699 guage to a formal language that allows han-700 dling of complex, temporal events that are typ-701 ically unable to be handled by most (tradi-702 tionally Markovian) methods. Our framework 703 exploits the structure of LTL by progressing 704 grounded LTL expressions along paths in the 705 environment, thus giving the required super-706 vision signal to parse language into logical 707 form. As opposed to other methods that em-708 ploy sequence-to-sequence models to map lan-709 guage directly to actions in the environment, 710 our method first formulates the meaning of 711 the natural language instruction in logical form 712 that is interpretable, and then uses this meaning

representation to formulate a plan to navigate

713

in an environment.

## 6 Future Work

Our concern in this project is to ground natural language instructions to logical representations that preserve the underlying meaning of the instruction that agents can then later use to correctly navigate in an environment to solve the task. Here, we primarily focus on the semantic parsing component, and show that grounding to LTL 714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

While this work does not plan or execute a policy in an end-to-end fashion, the use of LTL for task specifications can be beneficial for learning to plan, especially in hierarchical reinforcement learning settings that attempt to decompose tasks into subtasks. Future work will explore learning to ground language to collections of tasks specified in LTL, that can be composed together and learned in RL settings, by allowing agents to make use of LTL progression to extract shared subtasks and generalise this across different instances.

### References

- Jacob Andreas and Dan Klein. 2015. Alignmentbased compositional semantics for instruction following. *arXiv preprint arXiv:1508.06491*.
- Yoav Artzi and Luke Zettlemoyer. 2011. Bootstrapping semantic parsers from conversations. In *Proceedings of the conference on empirical methods in natural language processing*, pages 421–432. Association for Computational Linguistics.
- Yoav Artzi and Luke Zettlemoyer. 2013. Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, 1:49–62.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv* preprint arXiv:1409.0473.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on free-

base from question-answer pairs. In *Proceedings*of the 2013 Conference on Empirical Methods
in Natural Language Processing, pages 1533–
1544.

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

- Satchuthananthavale RK Branavan, Harr Chen, Luke S Zettlemoyer, and Regina Barzilay. 2009. Reinforcement learning for mapping instructions to actions. In Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1, pages 82–90. Association for Computational Linguistics.
- David L Chen. 2012. Fast online lexicon learning for grounded language acquisition. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 430–439. Association for Computational Linguistics.
- David L Chen, Joohyun Kim, and Raymond J Mooney. 2010. Training a multilingual sportscaster: Using perceptual context to learn language. *Journal of Artificial Intelligence Research*, 37:397–435.
  - David L Chen and Raymond J Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*.
  - Henry Chen, Austin S Lee, Mark Swift, and John C Tang. 2015. 3d collaboration method over hololens and skype end points. In *Proceedings* of the 3rd International Workshop on Immersive Media Experiences, pages 27–30. ACM.
- James Clarke, Dan Goldwasser, Ming-Wei Chang, and Dan Roth. 2010. Driving semantic parsing from the world's response. In *Proceedings of the fourteenth conference on computational natural language learning*, pages 18–27. Association for Computational Linguistics.
  - Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. 2017. Openai baselines. *GitHub, GitHub repository*.
- Juraj Dzifcak, Matthias Scheutz, Chitta Baral, and Paul Schermerhorn. 2009. What to do and how to do it: Translating natural language directives into temporal and dynamic logic representation for goal management and action execution. In 2009 IEEE International Conference

*on Robotics and Automation*, pages 4163–4168. IEEE.

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

- Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018. Speakerfollower models for vision-and-language navigation. In *Advances in Neural Information Processing Systems*, pages 3318–3329.
- Omer Goldman, Veronica Latcinnik, Udi Naveh, Amir Globerson, and Jonathan Berant. 2017. Weakly-supervised semantic parsing with abstract examples. *CoRR*, abs/1711.05240.
- Dan Goldwasser and Dan Roth. 2014. Learning from natural instructions. *Machine learning*, 94(2):205–232.
- Nakul Gopalan, Dilip Arumugam, LL Wong, and Stefanie Tellex. 2018. Sequence-to-sequence language grounding of non-markovian task specifications. In *Robotics: Science and Systems*.
- Kelvin Guu, Panupong Pasupat, Evan Zheran Liu, and Percy Liang. 2017. From language to programs: Bridging reinforcement learning and maximum marginal likelihood. *CoRR*, abs/1704.07926.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. 2016. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in neural information processing systems*, pages 3675–3683.
- Xiao Li, Cristian-Ioan Vasile, and Calin Belta. 2017. Reinforcement learning with temporal logic rewards. In 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 3834–3839. IEEE.
- Percy Liang, Michael I Jordan, and Dan Klein. 2013. Learning dependency-based compositional semantics. *Computational Linguistics*, 39(2):389–446.
- Michael L Littman, Ufuk Topcu, Jie Fu, Charles Isbell, Min Wen, and James MacGlashan. 2017. Environment-independent task specifications via gltl. *arXiv preprint arXiv:1704.04341*.
- Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. 2006. Walk the talk: Connecting lan-

840guage, knowledge, and action in route instruc-841tions. Def, 2(6):4.

842 Hongyuan Mei, Mohit Bansal, and Matthew R Walter. 2016. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *Thirtieth AAAI Conference on Artificial Intelligence*.

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

- Dipendra Misra and Yoav Artzi. 2015. Reinforcement learning for mapping instructions to actions with reward learning.
- Dipendra Misra, John Langford, and Yoav Artzi. 2017. Mapping instructions and visual observations to actions with reinforcement learning. *arXiv preprint arXiv:1704.08795*.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529.
  - Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information* processing systems, pages 3104–3112.
  - Richard S Sutton, Doina Precup, and Satinder Singh. 1999. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211.
  - Lappoon R Tang and Raymond J Mooney. 2000. Automated construction of database interfaces: Integrating statistical and relational learning for semantic parsing. In Proceedings of the 2000 Joint SIGDAT conference on Empirical methods in natural language processing and very large corpora: held in conjunction with the 38th Annual Meeting of the Association for Computational Linguistics-Volume 13, pages 133–141. Association for Computational Linguistics.
- Rodrigo Toro Icarte, Toryn Q Klassen, Richard Valenzano, and Sheila A McIlraith. 2018.
  Teaching multiple tasks to an rl agent using ltl. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 452–461. International Foundation for Autonomous Agents and Multiagent Systems.
- Min Wen, Ivan Papusha, and Ufuk Topcu. 2017.Learning from demonstrations with high-level

side information. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence.*  882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

- Edward C Williams, Nakul Gopalan, Mine Rhee, and Stefanie Tellex. 2018. Learning to parse natural language to grounded reward functions with weak supervision. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pages 1–7. IEEE.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. arXiv preprint arXiv:1809.08887.
- Luke S Zettlemoyer and Michael Collins. 2012. Learning to map sentences to logical form: Structured classification with probabilistic categorial grammars. *arXiv preprint arXiv:1207.1420*.