

# Notions of syntax and semantics for voice assistants in autonomous vehicles

Warrick Macmillan

Developed with the support of Ekaterina Komendantskaya

## 1 Abstract

We introduce a grammar for a controlled natural language (CNL) to give imperative commands for an envisioned voice assistant route-planner for a self-driving vehicle. The utility of the CNL is that it is inductively defined by a grammar : thereby, the sentences it admits, parsed as Abstract Syntax Trees (ASTs), can be manipulated as mathematical objects amenable to verification techniques. Using the TOUCHDOWN data set to empirically motivate common idioms and phrases our grammar should be capable of parsing, we give a denotational semantics from our ASTs to a Linear Temporal Logic (LTL) formulas, essentially expressing sequences of states which are amenable as specifications to downstream applications whose goal is the verification of various aspects of a vehicles behavior. This work contributes to a large existing literature, connecting the somewhat disparate research spaces including CNLs, verification for natural language-controlled robots, and semantic parsing.

## 2 Contributions

Initially motivated to give the system assurances against so-called substitution-based attacks, whereby impose “meaning equivalence” for synonymous expressions by imposing posterior conditions on the parse trees. Clustering via the tree structure to provide some the equivalence to meaningfully similar sentences was an initially enticing direction, as had been done with Komendantskaya and Heras’ work on Machine Learning for Proof General (ML4PG) [18]. However, this work had the advantage that there are multiple large, well-maintained Coq libraries which were amenable to clustering. Successful clustering results could give rise to proof developers seeking suggestions in their developments.

Our use case, the design of a non-existent language, left us with the conundrum of an impoverished data set to train over. There was no empirical data source from which to observe “natural ASTs”, and generating trees in an ad-hoc random basis would not likely provide real-world applicability. What has followed should be seen as a response to these constraints. Our intention was to find a data set suitable to give examples of non-trivial natural language utterances, in addition to finding a suitable semantic language with utility and applicability which is amenable to translation from sentences parsed by our grammar. Working from both the empirical and semantic directions seemed the most reasonable way to build a robust least prototype of a CNL .

The primary contributions of this undertaking so far are as follows :

- A Grammatical Framework (GF) grammar providing the definition of a CNL for suitable directives from a passenger to a driving agent
- A Haskell library mapping trees generated by our grammar a particularly well-behaved subset of LTL
- An Agda implementation of LTL with a standard semantic interpretation
- A refinement of the TOUCHDOWN dataset [8] suited to our needs of designing a better grammar

We suggest that while each of these components are still relatively primitive, they define a pipeline which has potential to provide both theoretical insights to researchers and suggest possible practical steps that can be taken to constructing robust voice applications in industrial settings. In addition to discussing our own contributions, we give a relevant and comprehensive literature review that embeds our work in the context of ongoing studies about these topics.

## 3 Overview

Perhaps the most pervasive question in the use and application of natural language technologies can be stated as follows : How does one optimize the system to provide for wide coverage of the domain while ensuring that system is robust? This question exemplifies the boundary of the “formalist”, verification-minded and “empiricist”, data-oriented camps in designing such technologies,

The statistical and machine learning methods applied to Natural Language Processing (NLP) tasks have made enormous gains over the past three decades. They take a more pragmatic approach : compromise robustness for wide coverage, as this means the tools will be usable and by non-experts. The belief espoused is that the machines should “learn” from us. Somewhat orthogonal techniques prioritize the formal approaches of the computational linguistics communities. These methodologies are often more concerned with theoretical justification and explainability.

While practical tools are a goal, building practical applications often isn’t linguistically informative and therefore the empiricists’ goals shouldn’t override work on building the theoretical models which enable our understanding of the machines. Those in the formalist camp, prioritizing theoretically informed systems, seek predictable and well-defined behavior for specific problem domains. Yet, these systems fail to generalize without an explosion in complexity when presented with data outside their domain.

Natural language is difficult because it is both structured with respect to “rules”, perhaps more descriptively titled *logical behaviors* which admit lots of predictability. Yet natural language continuously breaks or introduces exceptions to these rules, necessitating empirical and observational understanding. This makes it exceedingly hard to penetrate from exclusively the empirical or formalist approach. Many are led to wonder about the degree to which large amounts of linguistic data can be augmented with theoretical linguistic knowledge to create optimal and practical systems with respect to both breadth and depth of coverage of language phenomena. The ultimate question seeking compromise from both camps asks : how can we build machines which “understand” us (or at least our data), and which are comprehensible by us.

This problem acutely arises in when trying to design a voice assistant in the domain of commanding controllable robots, specifically, autonomous vehicles. For the actions a vehicle takes, mostly the motion and path decisions, must be formally specified or at the very least controlled via some computer system subject to mathematical formalism. Assuming the user directing the vehicle isn’t aware of these formalisms, it is incredibly difficult to design a verifiable controller capable of dealing with the breadth of language one may encounter in the wild.

The instructions an arbitrary user gives are not subject to the same formalities the system requires. For her commands may leave out necessary detail (“go into the other lane” with multiple lanes on either side), say something wrong with respect to reality (“go into the other lane” on a single lane road), or give a command the controller should recognize as possible but bad (“drive directly into the car ahead”). Additionally, the controller may need to recognize many ways many users may say “the same thing”, that is the same with respect to some semantic formalism. In a dual sense, the same utterance may admit two perfectly meaningful interpretations in two situations or contexts. The phrase “drive to the store with the dog” should account for whether the dog is inside of the car. It is obviously worrisome that a nefarious actor may somehow interfere with the controls at any stage by exploiting the manifold issues arising as stated above. A failure to adequately deal with these issues in the vast majority of cases is not reassuring if one believes many verifiability criteria are critical for such technologies to see adoption.

We therefore analyze our “big-picture question” above in the following “sub-question” : how can one map the manifold ways of presenting information to an autonomous robot into a rigorous and formally verifiable kernel which the controller can understand? Our proposed solution is to build a semantic parser from natural language commands to Linear Temporal Logic, whereby we can filter the many empirical natural language commands into a “canonical subset” defined by our CNL which are equivalent to (sets of) temporal logic formulas. We detail here both the progress to these ends, as well as the challenges.

## 4 Previous Work

This research broaches many different fields, many of which were unknown to the author prior to this work. Indeed, voice assistants may encompass almost any natural language processing task, and autonomous vehicles are one of the premier emerging robotics technologies (and certainly the most talked about in the popular zeitgeist). The cyberphysical systems at their intersection is likewise incredibly broad.

Limiting the scope of work in this context can be challenging, as so many different tools and ideas can be seen as relevant. We therefore try to very explicitly narrow our focus to investigate how feasible it is to build a language for an autonomous vehicle that exhibits predictable behavior and also satisfies verification properties - this includes a determination to what extent the properties can even be stated.

The approach taken therefore sets out to build a semantic parser, even if primitive, that serves as a Petri dish through which many of the deeper questions in this space may be viewed.

### 4.1 GF, Parsers, and Personal Work

The questions of designing an expressive formal language, with roots in Frege [12], manifested more recently in the natural language semantic tradition of Montague [28], who proposed an interpretation of English in a typed higher order logic with a focus on quantifiers. Aarne Ranta, a student of Martin-Löf, attempted to reformulate Montague’s work in an intuitionistic setting [31], thereby amenable to a natural treatment via computer programs [ml79]. In implementing a parser from natural language to a dependent type theory, Ranta discovered that dual sugaring (pretty printing) transformation of a tree to a string could allow for a general mechanism of purely syntax-based translation. This work culminated in Grammatical Framework (GF) [32].

Grammatical Framework became a full research project, allowing for the simple specification of a parser using a statically typed programming language whereby the grammar rules could be seen as types. Separate concrete syntaxes cohering with a given abstract syntax allowed for language-specific parsing, sugaring, and translation. The GF “standard library”, the Resource Grammar Library (RGL) [33], allows one to get off-the-shelf grammatical constructions for more than 30 languages, with English being the most comprehensive. The RGL therefore allows the grammar writer to focus on the semantic domain of the application the grammar is being developed for. In addition to this, one can embed a grammar as a Generalized Algebraic Datatypes (GADTs) in Haskell via the Portable Grammar Format (PGF) [1]. One can get run-time support for parsing and linearization directly in Haskell, in addition to manipulating the trees by pattern matching over them as Haskell programs.

A reflection on these historical developments reveals that GF is intimately tied to both the formal/informal distinction in addition to the syntactical and semantical approaches present in computational linguistics. These dual characteristics very much inform our problem as well. In the context of designing a voice assistant for, whereby one can give commands like “turn right after the woman with the big dog”, we desire that the intensional belief a user has about her utterance is consistent with the extensional behavior of the vehicle. This can be done through an intermediary mapping to a formal semantic representation. Ensuring that the syntactic content of a voice director’s (well-formed) utterance maps predictably to the logical form is important from the verificationist perspective : one wants to maximize the *syntactic completeness* of the system [23].

In a dual situation we briefly mention, one can imagine our voice assistant as giving the user feedback, responding with clarifications (“we will turn after the big cafe even though the other route may have less traffic”), questions (“do you mean this or that person?”), or even possible illocutionary directives (“we won’t drive over the speed limit in a school zone”), requiring the computer to generate an utterance after it has made some internal determination. This internal deliberation must be a program, possibly expressed inside or outside our semantical space,

identifying multiple routes in the clarification, multiple possible states in the question regarding two people, or constraints based off external circumstances like speed limits. In each case, the formation of a natural language utterance requires the computer to generate natural language which must conform to both a program's structure and behavior, but which also may be clear and recognizable to the user.

Independently of *how* the robot determines a program whose meaning it needs to convey to a user, the property of providing a natural language utterance which fluently conveys meaning to a native speaker is called *semantic adequacy* [23]. Determining a reasonable syntax and semantics for a controlled natural language should most certainly conform to the dual standards of syntactic completeness and semantic adequacy, if the voice assistant is to be held to any kind of regulatable standard.

## 4.2 Semantical Representations

We choose LTL as our semantic form in large part due to its relative expressivity for the kinds of verification conditions one might anticipate an autonomous vehicle needing to carry out, in addition to its ubiquitous appearance in the existing literature. Nonetheless, it is obvious their are many types of logical conditions LTL doesn't immediately support, and other logics, particularly ones which allow one to reason about space in its relation to time, would be an ideal direction to look. This line of research is probably more suited to people developing systems at later stages of development, where empirical observations may be collected in the wild. The nuances of where an autonomous navigator responding to a human agent can go wrong, and the most amenable set of verification conditions to prevent this, will ultimately have to be gained through trial and error.

### 4.2.1 Notions of Semantics

We also note that a notion semantics, having many connotations and interpretations in different fields, is subject to many interpretations. When we say semantics

- In linguistics, semantics, in its most natural understanding, should be interpreted as intended meaning. Different theoretical branches of meaning may include a logical meaning, as in the case of Montague semantics, or a meaning as it arises in the use and context of culture, as is the case of Cognitive Semantics.
- In programming languages, the semantics of a syntactic entity most commonly means the mathematical behavior (denotational semantics) and behavior during execution (operational semantics)
- In statistical notions of semantics, one often seeks the ability of one to capture meaning via language use, most common in contemporary contexts, its practical uses. Frequently Word2Vec [26] is referenced in this context, although the advent of transformers in recent years has largely usurped this technology

The problem presented in our work, of speaking to a machine, presents challenges in that it requires notions of semantics from disparate disciplines, which themselves have little overlap (at least as treated in the existing literature). This is because we are attempting to witness an utterance as a natural, native linguistic phenomena with an indented speaker meaning, a program whose syntax is defined via the CNL, and a statistical observation defined over some probability distribution. More concretely,

- How is the speakers meaning interpreted as if intended to be understood by other native speakers?
- How does the speakers meaning manifest as a formal program a computer can evaluate?
- How can we identify a speakers meaning in a possibly infinite space of utterances and contexts in which those utterances arise, neither of which can be formally defined *a priori*?

Although the inter-relatedness of various semantic theories is a much bigger project than we can give space to here, it should be granted that problem we address forces one, both implicitly and explicitly, to try to grapple with them.

We chose *the syntax of LTL* as the *semantics of our CNL* which is defined by filtering a “naturally observed” corpus to a primitive grammar and then be used to fit unseen utterances by fine-tuning a transformer-based language model to the corpus and grammar.

- Meaning for a computer can be analyzed in a variety of ways
  - The meaning of an utterance is a logical formula following Montague’s lead, with temporal operators taking precedence over quantifiers.
  - That the passenger’s utterance should be determined as a speech act which carries illocutionary force and intention. The computer’s response can be seen as conforming to or negotiating with the desires of the user, subject to the computers internal constraints and possible contextual information about which user may be unaware. Applying speech act theory in the context of human computer interaction has a long history [43]
- The meaning of the syntactic formula, can be interpreted in many possible ways
  - A possible program to be evaluated. In a case where temporal logic formulas are interpreted as types, Functional Reactive Programming [41]. One could then see the evaluation of a program conforming to some operational semantics coherent with the type system.
  - A (possibly verifiable) motion planner [35] [5] [19]
  - A dialogue state, in the envisioned Question/Answer context, whereby the computer must provide feedback to the user based of contextual information
- The meaning from the mostly unseen utterances is given a canonical form, and the canonicalization process is a transformation via the vector-space and distributional notions of meaning implicit in the transformer neural networks structure

We don’t intend to exhaust the list of possibilities here, neither in our description of the many meanings of “semantics”, nor in how our taxonomy of semantics can be understood in the context of our solution to the problem of giving navigation commands to a autonomous driver. We intend to clarify some of the many subtleties and terminological confusions arising from many communities of researchers. We suggest that working towards a unified view of what kinds of semantic notions we want to deal in this particular domain may inform better solutions to the problem at hand.

### 4.3 Robot Motion Planning from Natural Language

The challenge of designing a system which generates robot control strategies from human language has to balance the expressiveness of task specification, complexity of environment, and provable correctness [2]. In this context, we assume that expressivity of the language itself should reflect the complexity of the environments, thereby being adequately descriptive. The criteria of correctness : that the language itself is well-represented in the LTL semantics - the system being is syntactically complete - is the focus of these investigations. Our work additionally, is the only work we know of which actually seeks autonomous vehicles as the central motivation, rather than more general robotics applications.

Just as important as producing a well-formed and meaningful LTL formula, but not explored is here, is that the evaluation gives a faithful controller to navigating a complex environment. For instance, in [30] the authors indicate how to actually ground basic propositions from language to paths in a space, while our model, outputting formulas with non-grounded base predicates, is merely concerned with logical structure.

Similarly, in [5], the authors develop a Verifiable Distributed Correspondence Graph (V-DCG) model whereby LTL formulas are used to ensure grounded instruction sequences are

consistent. This work builds on other work of Kress-Gazit et al. [21], whereby the The Situated Language Understanding Robot Platform (SLURP) allows translation of arbitrary natural language into LTL - with less reliance on grammar formalisms - and also supports a notion of feedback. Our work is directly focused on the *centrality of the grammar* (which we will emphasize continuously) as the key intermediary phase when balancing formal and empirical interests, whereas this work is more concerned with the controllers generated as the end result of a pipeline where the intermediary grammatical structure may not be so relevant.

Our GF implementation, seeing the grammar as a necessary part of the verifiability (in that we can systematically map our sequences of commands to logical formulas representing sequences of states), also makes the possibility of supporting multi-lingual verifiability more immediate. Our system does not support this currently, but can easily be adjusted to so with the help of GF’s functors (roughly adapted from Standard ML’s functors) and the RGL. The lack of wide-coverage support of our grammar is possible to remediate through possible fine-tuning of a large language model to a data-set which coheres to the language our GF grammar generates, and we detail this in our discussion below [TODO : link].

Another approach seeks to train a natural language to LTL planner using both NNs and reinforcement learning [42]. Their work also uses a simultaneous CFG to generate *semantically inadequate* sentences with corresponding LTL formulas from which they can direct machines to follow the instructions, and then have users describe the robot behavior in a more natural form. Despite this, their approach uses the machine to generate sentences and corresponding situations, most of which are “nonsense” and need to be filtered out, thereby leaving the narrations upon which their system leaves devoid of a genuine empirical data source. In addition, their corpus only contains 266 words, still not the size one would need for our system. Finally, our suggested use of a pre-trained language model fine-tuned to the semantic parsing task gives us more flexibility in that the neural network and the verifiable grammar and semantics in the kernel could allow us to focus on the problems of breadth and depth somewhat independently.

The same group, in [19], explore the most general possible end-to-end utterance to planner pipeline without intermediary states, namely, a symbolic representation. While this “cutting out the middle woman” mentality may be an idealistic long-term vision, it makes the system much too much of a black box - even though they are able to reason about their system’s behavior through the use of attention maps. For the fine-tuned verification conditions about the linguistic utterances our work explores, the intermediate symbolic representations give a more explainable, predictable, and regulatable system.

- 
- More similar to Tellex et al’s approach [29],
- One of the

#### 4.3.1 Temporal Logic for NL verification

Modal logics, specifically those dealing with stateful staging of events like Linear Temporal Logic (LTL), Computation Tree Logic (CTL), Signal Temporal Logic (STL), have been used extensively in the specification and verification of properties of robotics systems, including autonomous vehicles [].

With verification being a core motivation of our work, we take for granted that these different logics have differing interpretations, utilities, or manifestations in their different applied settings. As LTL is often seen as one of the “primitive” temporal logic, we chose it as our semantic space despite its limitations. We appreciate that future work will need to expand the scope of which logic (or possibly *logics*) the machine may use to verify behavior, in addition to the mathematical models most amenable to verification of a logical formula.

We recognize that there are many degrees of freedom in the both the syntactic and semantic formalisms chosen, as well as their evaluation or grounding within physical environments.

- Parsing formalisms - phrase-based grammars (GF), categorial grammars, or even perhaps dependency grammars for wide-coverage usage
- Semantic formalisms including spatial logics and vector space semantics
- Mechanisms for evaluating and syntactic and semantic choices
- Robotics domains outside of the autonomous vehicle space, for which natural language is still an ideal interface. Ideally we'd like one system to be generally applicable to many areas, independent of certain nuances in a given space

#### 4.3.2 NL to TL

Here we show mainly relevant research for NL to LTL.

The applications of LTL in machine learning are vast, and the scope of our specific application is still unclear, but nevertheless, we give a literature review of methods and applications relevant for our work.

- This paper [11] from 2009 uses a categorial grammar approach, but more or less can serve as an idea template for us, also nice pictures with grammar rules and formulas
- Also, a highly relevant template combines Natural Language, LTL, with the idea of having a verifiable pipeline [21]

The production of an ontology of common actions and the type of formulas that they produce—for example, safety conditions, adding goals, constraining the initial state—in their negated and positive forms would be a step toward a more general solution to the problem of mapping natural language to LTL.  
Previous work has relied heavily on grammar formalisms to ease... [21]

- LTL formulas can be transformed into automata which can then be used as reward functions for reinforcement learners, as in [7]
- The following is one of the more relevant quotes from a paper reviewing the whole space of English to LTL translations

Overall, the typical approach followed by these studies can be summarized as follows: given an input English utterance, preprocess it to extract syntactical information, which may include part of speech tagging, dependency parsing, semantic role labelling, and so on. Then, enrich the input with these pieces of information. Finally, run an attribute grammar-based parser, or rely on some hand-made rules, to derive a translation into a target logical format. A notable exception is the work of [89], where a fully-supervised learning setting is considered. [6]

- Translating between English and STL can be done via a large language model [16] [under review], but the domain specificity of the problems are still significant enough to suggest that it will be years before an automated semantic parser is available, if it is even possible.
- Could ask Lapata in Edinburgh, whose work [10] is relevant and well-cited (although they use an encoder-decoder method)
- 
- 
- 

#### 4.3.3 Tellex

Stephanie Tellex has written extensively about natural language inputs and interfaces with robots. Although she has not specifically written about autonomous vehicles, the domains have enough intersection to warrant careful consideration of much of her work, especially the recent stuff.

- Grounding with an intermediate symbolic state, no LTL, but possibly relevant for paper generally. She also cites [22], a seminal paper in this area

Instruction following is a supervised learning problem where the agent must predict a trajectory that would satisfy an input natural language command. [13]

- The review paper [25] making recommendations has a section on robustness, but this is mostly for the sake of allowing sharing of interfaces and efficacy, no mention of verification (which is what we're primarily after)
- They design a NL -> LTL for drones that are grounded to actual landmarks [4]
- The group builds a trained pipeline that uses an object oriented template-instance methodology to generalize to different ontological categories in [17] [under review]
- In [29] build learn a semantic parser from NL to LTL (so that the language is grounded) where they collect executions of the LTL formulas in different environments using a weakly-supervised training method with reinforcement learning Part if the paper has to do with the execution of the command being dependent on the path taken by the robot executing the command, not just meeting the goal requirements, thereby giving a complexity bonus in comparison to previous work. She also evaluates the model on the [22] data set

## 5 Work

### 5.1 GF Grammar

## 6 TODO

### 6.1 Grammar modulo wordnet

### 6.2 LTL in Agda

Along with colleagues from Singapore Management University, we have begun an Agda implementation [24] of LTL which will serve as the semantic space for our parsed utterances. Our method, uses a deep embedding, as opposed to the shallow embedding in [9], although the temporal encoding of paths as streams was directly adapted from this paper.

This implementation will hopefully allow us to prove decidability of LTL in a relatively straightforward manner. Other than the assurance that our implementation is correct, we hope this will allow us to feed the formula into some SAT or SMT solver so-as to actually allow verification of the behavior of a vehicle with respect to an utterance.

[TODO : Help from Matthew?]

### 6.3 AST -> Agda

### 6.4 ML training/verification stuff

Help from Marco, Nathalia if interested?

## 7 Publications Description

Realizing that the structure of the paper is amenable to large changes, I'm posting a summary of relevant publications here.

### 7.1 Statistical (pre-trained) Language Models

The first set of publ

- In [38] [under review], the authors show how, using a *synchronous context-free grammar* (SCFG) to define a minified CNL with a parallel and dually parsable semantic form, that one can use a large pre-trained language model as a front-end to filter a much wider

syntax into the CNL. I postulate GF’s expressivity is more expressive than the SCFG, at least based off a tertiary reading in the index, and therefore if we carved out a subset of commands to cohere with our LTL (and maybe some other temporal or even spatial-temporal logics in the future), our model would be amenable to a similar “out-of-the box” semantic parser that could actually be used for verification. This paper borrows the idea of “semantic parsing as paraphrasing” from [3]

- In [36], the authors advocate for getting rid of parsers altogether, although this naively takes for granted large public data-sets, none of which exist for an autonomous vehicle and temporal logic formalism
- [15] [under review] claims that Bert is robust, analyzing claims of four papers, including the one which uses a wordnet attack

## 7.2 Voice assistants for autonomous vehicles

The public company Cerence [] is already designing voice assistants for autonomous vehicles, for which it has a large software stack between the voice processing to actual control of current automative components. In addition to its technologies, many of which aren’t accessible to external researchers due to intellectual property restrictions, Cerence has contracts with large automakers [...]. It is therefore natural to inquire, what a small team with varied backgrounds and not nearly the same expertise nor experience within the technological team at Cerence can provide.

First, we believe that the focus on verification, insofar as we envision it, is unlikely to be of current concern at Cerence due to the fact that their products are still being developed, and the primary goal of producing a working product is likely to precedence over preventing non-existent hostile actors.

Additionally, it is going to have to be determined by [verification of self-driving cars generally : software, hardware, behavior in a real environment, etc]

## 8 Future Work

The “sets of” clause references the inevitable ambiguity of parses even from a big enough parser, even if the size of the canonical expressions is vastly smaller than the domain of expressions mapping to them.

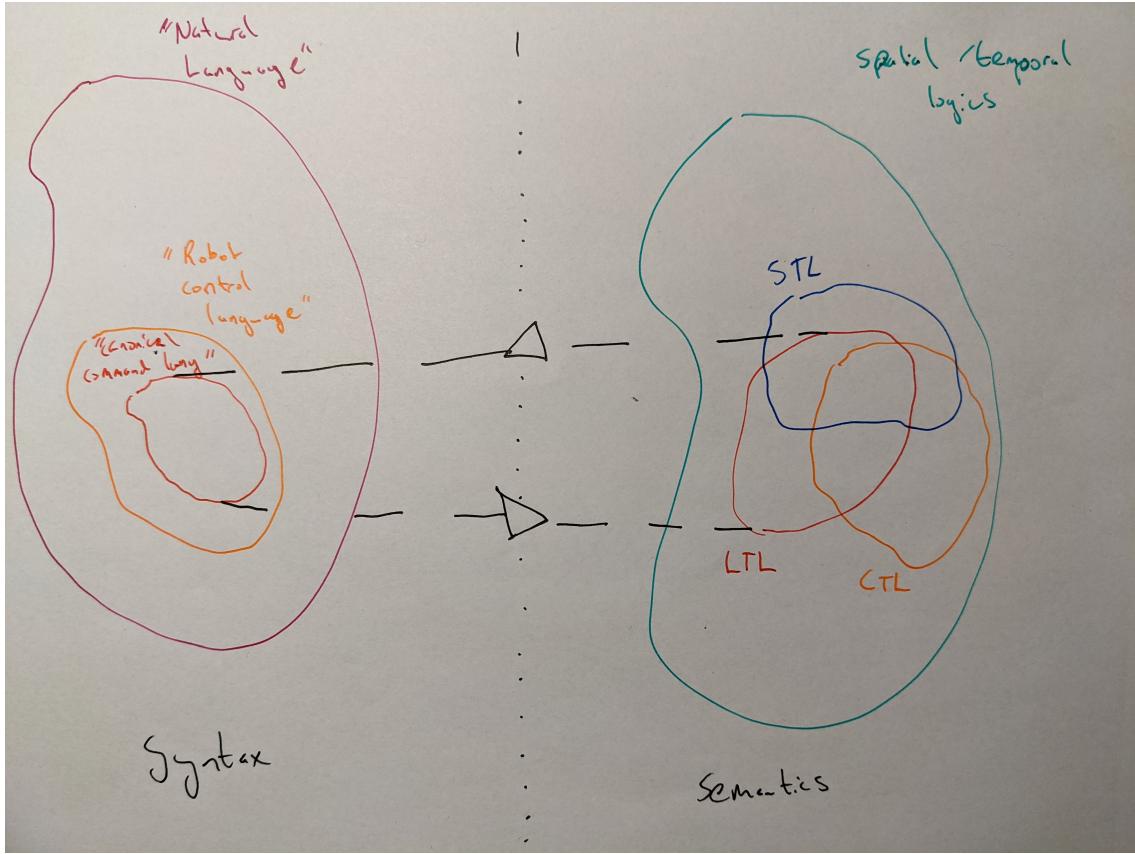


Figure 1: Language and Logical Spaces of Concern

We begin in [Figure 1](#) with a high level overview of this semantic parsing system, whereby the space of natural language syntax can be mapped to some formal language semantic space (and possibly have some kind of inverse mapping). We note that “Natural Language”, while an idealized notion, can be thought of the space of interpretable utterances. The relatively small subset of these utterances which one might give to a robot, labeled “Robot Control Language”, is the ideal breadth our system would support, is still actually very large. We therefore applying another filter, to the “Canonical Command Language” which is inductively defined via some relatively thin set of grammar rules, which simultaneously generate and parse expressions in some logic. Although we target LTL because of its prominence in the literature and relatively straightforward implementation and interpretation, it should be noted that there are other temporal logics which may well be more expressive and better suited to the actual problem of synthesizing controllers.

Due to the recent influx of transformer based language models like Bert and GPT-3, we take for granted that the easiest way to target our “Robot Control Language” will be through fine-tuning one of these models, as shown in [Figure 2](#). These transformers, trained on a separate corpus like Wikipedia, can be mapped to some suitable set of robot commands, even though these types of expressions will have a sparse presence in the corpus the model was initially trained on (presumably Marco will know more about this than me).

In this context, we can then further refine the language to something less natural, but more well-behaved. The whole proposed pipeline in [Figure 3](#), indicates using the methodology as used in [38], whereby the semantic parser should ideally be able to take any command from the Robot Control Language and turn it into a set of temporal logic formulas, distributed according to most likely interpretation.

Ideally, the downstream dialogue system should either be able to ask for clarification if two formulas are determined to be of some relative likelihood, reject a formula that is not determined

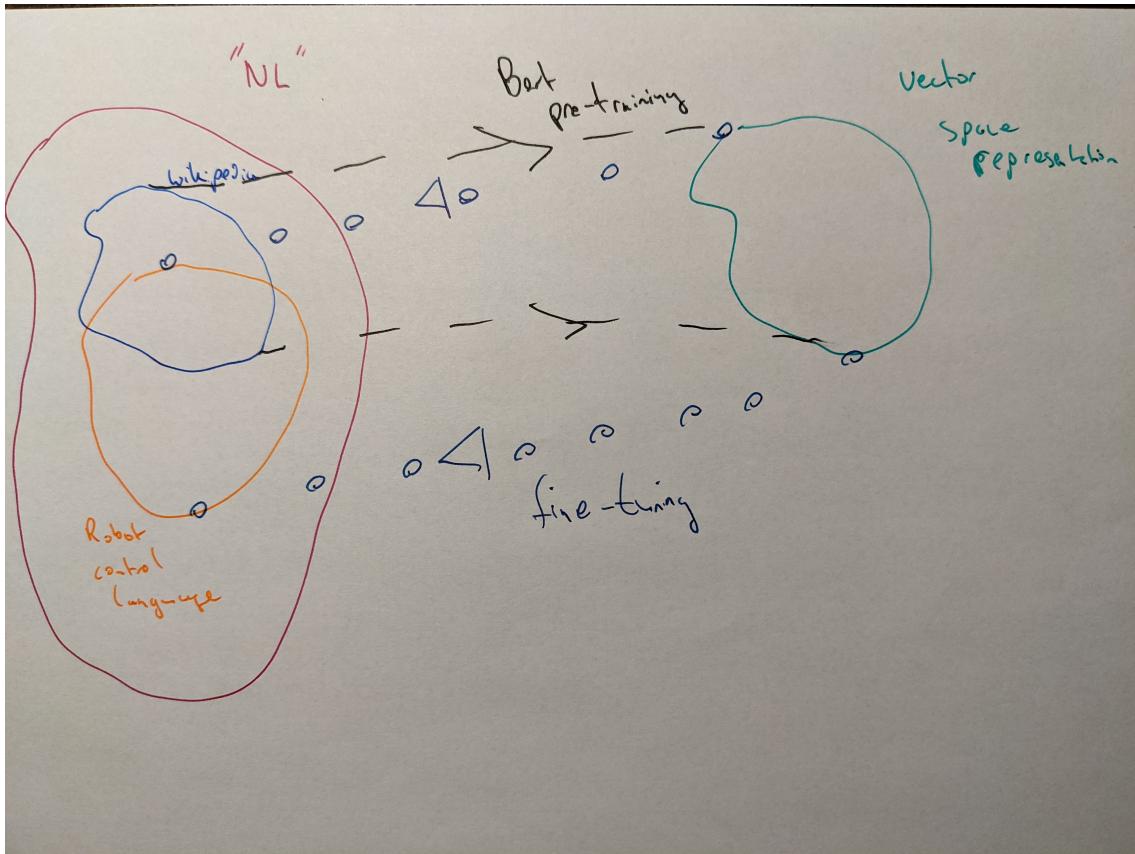


Figure 2: Transformer to Robot Control Language

to be achievable (for whatever reason), or synthesize a sequence of actions (and express those in the CNL) according to the possibly modified current path.

In theory, we can embed clauses which in turn reflect all of natural language : “Stop at the man who is watching the tv show on his phone about time traveler who goes back to the 12th century Mongolia, whereby the man, not speaking Mongolian ...” This is clearly outside the boundary of what the robot control language should support, and ideally would be accepted or rejected by the computer prior to the commands completion depending if there was a man looking at a phone. Our parser currently accepts strings in our primitive canonical language, designed in Grammatical Framework (GF), such as :

```
p "drive to the store , turn right and stop at the dog"
```

```
MultipleRoutes And (ConsPosCommand (SimpleCom (ModAction Drive (MkAdvPh To
(WhichObject The Store)))) (BasePosCommand (SimpleCom (ModAction Turn
(WherePhrase Right)))) (SimpleCom (ModAction Stop (MkAdvPh At (WhichObject The
Dog))))))
```

However, we may envision our system being able to accept an expression in the Robot Control Language like “hit the petal till we reach the store, hang a right, and halt when you see a cute little puppy”. We could certainly adjust our parser to accomodate this, but it would be one of many possible edge cases unlikely to be uttered. To accomodate many more such edge cases would cause an exponential blowup in the parser size (thereby slowing down parsing), but more importantly, cause the programmer a headache in building the parser, and then mapping the NL ASTs to a LTL form. If we treat  $F$  as the operating expressing the existence a future state,  $X$  as the next state, and  $G$  meaning the universal future, our desired LTL formula would

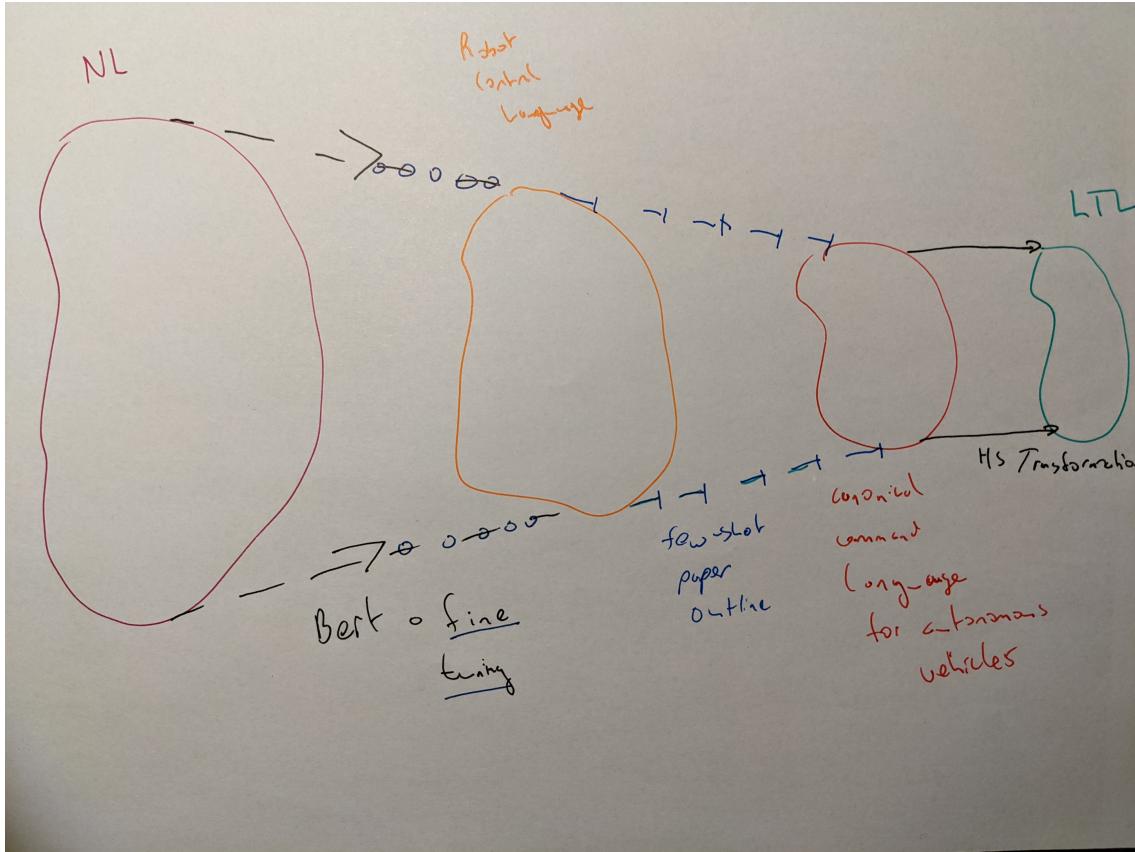


Figure 3: Pipeline from NL to LTL

most likely treat this as  $F(store \wedge (X \text{ turn\_right} \wedge (F(G \text{ dog}))))$ , although we propose that the actual grounding of these to images or controllable actions to some downstream system.

LTL has been a popular logic for specifying controllable robot behaviors, particularly with respect to verification of their behaviors. In [35], Rizaldi et al. prove logical correctness of a motion planner with respect to LTL formulas over maneuver automata formulas in the Isabelle/HOL theorem prover, a non-dependent cousin of Agda. We are choosing to deeply embed LTL in Agda for a few reasons, although the syntax of the embedding could easily be translated to any other dependently typed theorem prover, and with a little more effort probably any functional programming language. The composition of a “weakly verified” natural language front-end with a formally verified back-end such as in Rizaldi’s work would pave the way for a fully verified, utterance-to-vehicle-path pipeline for the autonomous vehicles.

The big question to address is what kind of verification conditions the natural language component should be subjected to, and what kind of attacks would be most important to preemptively anticipate. Substitution based attacks [37], for instance, have been consistently emphasized throughout our discussions so far. The question is, *where* in the pipeline it would be best to filter out the vulnerabilities, as well as *how*.

One possibility would be to define words modulo equivalent meanings using Wordnet [27] in the syntactic phase, either via training [34] (presumably during the fine-tuning to the Robot Command Language or our “canonicalization” from that). It has been suggested that Bert is already relatively robust against such attacks [15], but we nonetheless feel that even higher sensitivities of robustness may be better done at other phases in the pipeline.

Alternatively, one could just map these equivalent Wordnet forms to equivalent parse trees using the Portable Grammar Format (PGF) Haskell library, which essentially deeply embeds a GF grammar into a Generalized Algebraic Datatype (GADT).

For instance, if we abstract over all abstract syntax trees for our grammar using this library, we can define the following Haskell functions to equate a “female human” with a “woman”.

```
treeMapfemalePersonIsWoman :: forall a. Tree a -> Tree a
treeMapfemalePersonIsWoman (GModObj GFemale GPerson) = GWoman
treeMapfemalePersonIsWoman GWoman = (GModObj GFemale GPerson)
treeMapfemalePersonIsWoman gp = composOp treeMapfemalePersonIsWoman gp
```

There has been work integrating multiple language Wordnets with GF [40], so it would presumably be easy to integrate with our system, depending on how large we want the grammar to get.

As it is unclear what the best direction for this is, and how the attacker model in the context of an autonomous vehicle might work, all these decisions need to be made in the context of discussions within the group.

[ Addendum before meeting : ]

The TOUCHDOWN data set [8] seems like the most comprehensive and relevant data we'll find to fine-tune via one of these pre-trained models. Please see <https://github.com/lil-lab/touchdown>

The idea of domain specific pre-training can be traced to [14], where the authors introduce the concepts of *domain-adaptive pretraining* and *task-adaptive pretraining*, whereby this additional pre-training phase greatly improves efficacy of the LM on corpus and task data not well represented in the training data.

The language models have been show [20]

## 8.1 Data Set

The most comprehensive data set known to us is the TOUCHDOWN data set [8], which can simultaneously serve as a data source to inform the actual grammar (ideally, we'd like to be able to generate a grammar from a data source)

authors use “interactive visual navigation environment based on Google Street View”

position of the agent relative to an object, and the position of two objects relative to one another “resolving the spatial descriptions”, but we focus only on the navigation part of the task

Positive

- Real-world observations
- Real descriptions of these observations
- Diverse, relatively large
- 

Follows “instruction writing, target propagation to panoramas, validation, and workers and qualification”, the validation here

That having the data grounded is both incredibly beneficial, but also makes designing the syntax and semantics tricky.

- Finding the touchdown only coarsely approximates general navigating in a city.
- The workers aren't in a place they know, so everything the reference is in their immediate visual environment
- and this is a “short-term” task, it requires no long-distance navigation and reasoning
- Limited to NYC, hectic urban environments (also, daytime)
- Working with panoramas is not necessarily a great simulation to a real environment
- “They are not permitted to write instructions that refer to text in the images, including street names, store names, or numbers”
- No temporal reasoning (as spatio-temporal is assumed)

Ideally our data set would then have the properties

- Long and short-term tasks
- Different cities, different languages (this will be dependent on the context of the data collected), for instance, where there are dirt roads
- Updates over time (the user can update a context locally or globally) on the road
- Users with various degrees of contextual information Contextual information - street place, names (named entity recognition) accessible to google - people's names (mom's house)

That the data collection task in an objective way is inherently tied to the way we approximate it in collecting data, thereby limited by our experimental apparatus and assumptions in designing the data set.

What is the actual feasibility of this stuff?

how well does a given logic allow us to reason about a space of instructions. What is the logic grounded to, how it is verified , all of these may effect the choice of formulas

we can't just design a perfect language to capture our meaning - goes back to the wishful thinking of Frege, but we can try to approximate it

next intersection (and next left?) versus next gas station versus “next to” will be a store on your left with stars next to the name.

there is could either be “all of the apartment buildings on your left”

if you are going the correct way. logical content You will know you're on the correct road if to your left there are planters in the middle of the street

## 8.2 Syntax and Semantics

When designing a grammar, we can pretend that initially

an abstract syntax, we have the following considerations

That when we are conditioning our syntactic model of empirical, noisy, and biased natural language data, so as to ideally generalize to unencountered phenomena.

A central insight ambiguity :

- Ontological semantic space. What are we trying to represent?
- Intended semantic space, the logical or formal system which our grammar will map to (via Haskell transformations)
- We want to account for some grammatical constructions via the abstract syntax, but outsource most of the grammaticality to the RGL
- Data source. How to conform to the data set in a way that's faithful but doesn't overfit (the overfitting can probably result in generating functions which are useless and either make our parser slow down or overgenerate)

While theres no clear way of relating the trade-offs, we can come up with some heuristics that shed light. Developing the “ontological design” allows one to capture the intuitive problem.

## 8.3 Ambiguity

What happens when we encounter ambiguity? For instance, in p ”go to the person with the dog .” The prepositional phrase ”with the dog” can either modify person (as an adjectival clause) or it can modify go (as an adverbial clause). Because the parser is designed to accommodate simple cases of both types of clauses, these ambiguities, even in simple sentences from our corpus, will grow quickly.

In the case of a vehicle, however, knowing the correct parse is dependent on the context in which the driver is going to the person : is the language grounded in the fact that there a dog in the car, or a person in the purview with a dog (or, most confusingly, perhaps both conditions are met, in which case more contextual information is required to disambiguate the correct parse).

For we can actually program the semantics to accommodate both scenarios, whereby  
 $F(\text{manwithdog} / \text{GFinish})$   $F(\text{man} / \text{GFinish}) / \text{Gwithdog}$

We can define our semantics to accommodate both interpretations, whereby the parses produce unique semantic conditions, and the LTL solver will have to see which condition is more easily satisfied. While this edge case may seem overly pedantic to consider, as one's intuition might suggest the first case to be overwhelmingly more natural, the

We should just make simplifying assumptions, though, at least for the purpose of a grammar like ours.

probably scrap

## 9 Alternative ideas

While the evaluation of machine learning systems provides assurances using different scores and metrics on different tasks assures one they may on average perform better than humans at certain tasks, the advent of adversarial attacks [39] with the intention of deceiving such a system by a hostile actor leads the system designer to desire, and possibly require additional verification about the system's behavior. In the context of natural language processing (NLP), where data sources rely on strings of text, these attacks can focus an array of features from spellings of individual words to rearranging entire sentences []. So-called synonym attacks, which adversarially target the system at the lexical level, can cause traditional NLP models to [...].

Aside from the user experience being compromised by a system which has been adversarially afflicted, there is also a possibility of physical danger for the passenger and other people in the vicinity. As voice directed robots have many possible points of failure, we focus on two types of verification for our system. Rather than focus on breadth of language coverage, which ML language models excel at due to their reliance on statistical modeling and tons of data, our system is narrowly focused as a proof-of-concept, from which it could either be extended by hand, or different components modified using other techniques and tools.

## References

- [1] Krasimir Angelov, Björn Bringert, and Aarne Ranta. “PGF: A Portable Run-time Format for Type-theoretical Grammars”. In: *Journal of Logic, Language and Information* 19.2 (Apr. 2010), pp. 201–228.
- [2] Calin Belta et al. “Symbolic planning and control of robot motion [Grand Challenges of Robotics]”. In: *IEEE Robotics Automation Magazine* 14.1 (2007), pp. 61–70.
- [3] Jonathan Berant and Percy Liang. “Semantic Parsing via Paraphrasing”. In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Baltimore, Maryland: Association for Computational Linguistics, June 2014, pp. 1415–1425.
- [4] Matthew Berg et al. “Grounding Language to Landmarks in Arbitrary Outdoor Environments”. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. 2020, pp. 208–215.
- [5] Adrian Boteanu et al. “A model for verifiable grounding and execution of complex natural language instructions”. In: *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2016, pp. 2649–2654.
- [6] Andrea Brunello, Angelo Montanari, and Mark Reynolds. “Synthesis of LTL Formulas from Natural Language Texts: State of the Art and Research Directions”. In: *26th International Symposium on Temporal Representation and Reasoning (TIME 2019)*. Ed. by Johann Gamper, Sophie Pinchinat, and Guido Sciavicco. Vol. 147. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2019, 17:1–17:19.

- [7] Alberto Camacho et al. “LTL and Beyond: Formal Languages for Reward Function Specification in Reinforcement Learning”. In: *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, July 2019, pp. 6065–6073.
- [8] Howard Chen et al. “TOUCHDOWN: Natural Language Navigation and Spatial Reasoning in Visual Street Environments”. In: *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [9] Solange Coupet-Grimal. “An Axiomatization of Linear Temporal Logic in the Calculus of Inductive Constructions”. In: *Journal of Logic and Computation* 13.6 (2003), pp. 801–813.
- [10] Li Dong and Mirella Lapata. “Language to Logical Form with Neural Attention”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 33–43.
- [11] Juraj Dzifcak et al. “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*. 2009, pp. 4163–4168.
- [12] Gottlob Frege. *Begriffsschrift*. Halle, 1879.
- [13] Nakul Gopalan et al. “Simultaneously Learning Transferable Symbols and Language Groundings from Perceptual Data for Instruction Following”. In: *Robotics: Science and Systems XVI, Virtual Event / Corvalis, Oregon, USA, July 12-16, 2020*. Ed. by Marc Toussaint, Antonio Bicchi, and Tucker Hermans. 2020.
- [14] Suchin Gururangan et al. “Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, July 2020, pp. 8342–8360.
- [15] Jens Hauser et al. *BERT is Robust! A Case Against Synonym-Based Adversarial Examples in Text Classification*. 2021. arXiv: [2109.07403 \[cs.CL\]](https://arxiv.org/abs/2109.07403).
- [16] Jie He et al. *From English to Signal Temporal Logic*. 2021. arXiv: [2109.10294 \[cs.CL\]](https://arxiv.org/abs/2109.10294).
- [17] Eric Hsiung et al. *Generalizing to New Domains by Mapping Natural Language to Lifted LTL*. 2021. arXiv: [2110.05603 \[cs.CL\]](https://arxiv.org/abs/2110.05603).
- [18] Ekaterina Komendantskaya, Jónathan Heras, and Gudmund Grov. “Machine Learning in Proof General: Interfacing Interfaces”. In: *Electronic Proceedings in Theoretical Computer Science* 118 (July 2013), pp. 15–41.
- [19] Yen-Ling Kuo, Boris Katz, and Andrei Barbu. “Deep compositional robotic planners that follow natural language commands”. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2020, pp. 4906–4912.
- [20] Jinhyuk Lee et al. “BioBERT: a pre-trained biomedical language representation model for biomedical text mining”. In: *Bioinformatics* 36.4 (Sept. 2019), pp. 1234–1240. eprint: <https://academic.oup.com/bioinformatics/article-pdf/36/4/1234/32527770/btz682.pdf>.
- [21] Constantine Lignos et al. “Provably Correct Reactive Control from Natural Language”. In: *Auton. Robots* 38.1 (Jan. 2015), pp. 89–105.
- [22] Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. “Walk the Talk: Connecting Language, Knowledge, and Action in Route Instructions”. In: *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 2*. AAAI’06. Boston, Massachusetts: AAAI Press, 2006, pp. 1475–1482.

- [23] Warrick Macmillan. “On the Grammar of Proof”. MA thesis. University of Gothenburg, 2021, p. 90.
- [24] Warrick Macmillan and Andreas Kallberg. *LTL in Agda*. <https://github.com/wmacmil/LTL-Agda>. 2021.
- [25] Matthew Marge et al. “Spoken language interaction with robots: Recommendations for future research”. In: *Computer Speech and Language* 71 (2022), p. 101255.
- [26] Tomas Mikolov et al. “Distributed Representations of Words and Phrases and Their Compositionality”. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*. NIPS’13. Lake Tahoe, Nevada: Curran Associates Inc., 2013, pp. 3111–3119.
- [27] George A. Miller. “WordNet: A Lexical Database for English”. In: *Commun. ACM* 38.11 (Nov. 1995), pp. 39–41.
- [28] Richard Montague. “The Proper Treatment of Quantification in Ordinary English”. In: *Approaches to Natural Language: Proceedings of the 1970 Stanford Workshop on Grammar and Semantics*. Ed. by K. J. J. Hintikka, J. M. E. Moravcsik, and P. Suppes. Dordrecht: Springer Netherlands, 1973, pp. 221–242.
- [29] Roma Patel, Stefanie Tellex, and Ellie Pavlick. “Learning to Ground Language to Temporal Logical Form”. In: (2019).
- [30] Erion Plaku and Sertac Karaman. “Motion planning with temporal-logic specifications: Progress and challenges”. In: *AI communications* 29.1 (2016), pp. 151–162.
- [31] A. Ranta. *Type Theoretical Grammar*. Oxford University Press, 1994.
- [32] Aarne Ranta. “Grammatical Framework”. In: *Journal of Functional Programming* 14.2 (2004), pp. 145–189.
- [33] Aarne Ranta. “The GF Resource Grammar Library”. In: *Linguistics in Language Technology* 2 (2009).
- [34] Shuhuai Ren et al. “Generating Natural Language Adversarial Examples through Probability Weighted Word Saliency”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, July 2019, pp. 1085–1097.
- [35] Albert Rizaldi et al. “A Formally Verified Motion Planner for Autonomous Vehicles”. In: *Automated Technology for Verification and Analysis*. Ed. by Shuvendu K. Lahiri and Chao Wang. Cham: Springer International Publishing, 2018, pp. 75–90.
- [36] Subendhu Rongali et al. “Don’t Parse, Generate! A Sequence to Sequence Architecture for Task-Oriented Semantic Parsing”. In: *Proceedings of The Web Conference 2020*. WWW ’20. Taipei, Taiwan: Association for Computing Machinery, 2020, pp. 2962–2968.
- [37] Suranjana Samanta and Sameep Mehta. “Towards Crafting Text Adversarial Samples”. In: *CoRR* abs/1707.02812 (2017). arXiv: [1707.02812](https://arxiv.org/abs/1707.02812).
- [38] Richard Shin et al. “Constrained Language Models Yield Few-Shot Semantic Parsers”. In: *CoRR* abs/2104.08768 (2021). arXiv: [2104.08768](https://arxiv.org/abs/2104.08768).
- [39] Christian Szegedy et al. “Intriguing properties of neural networks”. In: *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*. Ed. by Yoshua Bengio and Yann LeCun. 2014.
- [40] Shafqat Mumtaz Virk et al. “Developing an interlingual translation lexicon using Word-Nets and Grammatical Framework”. In: *Proceedings of the Fifth Workshop on South and Southeast Asian Natural Language Processing*. 2014, pp. 55–64.

- [41] Zhanyong Wan and Paul Hudak. “Functional Reactive Programming from First Principles”. In: *Proceedings of the ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation*. PLDI ’00. Vancouver, British Columbia, Canada: Association for Computing Machinery, 2000, pp. 242–252.
- [42] Christopher Wang et al. “Learning a natural-language to LTL executable semantic parser for grounded robotics”. In: *CoRR* abs/2008.03277 (2020). arXiv: [2008.03277](https://arxiv.org/abs/2008.03277).
- [43] Terry Winograd and Fernando Flores. *Understanding Computers and Cognition: A New Foundation for Design*. Addison-Wesley, 1987.