2022 International Conference on Software Engineering Pittsburgh, PA, May 25th 2022



[Journal First - ACM TOSEM] In-IDE Code Generation from Natural Language: Promise and Challenges



Frank Xu @frankxu2004



Bogdan Vasilescu @b_vasilescu



Graham Neubig @gneubig







@gneubig



Carnegie Mellon University



Institute





"I suspect that machines to be programmed in our native tongues — be it Dutch, English, American, French, German, or Swahili are as damned difficult to make as they would be to use."

Edsger W. Dijkstra. 1979. On the foolishness of "natural language programming." In Program Construction. Springer, 51–53.

Carnegie Mellon University School of Computer Science









Dijkstra certainly got the "difficult to make" part right ...

It took ~40 years

[EMNLP 2018]

TRANX: A Transition-based Neural Abstract Syntax Parser for Semantic Parsing and Code Generation

Pengcheng Yin, Graham Neubig

Language Technologies Institute Carnegie Mellon University {pcyin,gneubig}@cs.cmu.edu

Abstract

We present TRANX, a transition-based neural semantic parser that maps natural language (NL) utterances into formal meaning representations (MRs). TRANX uses a transition system based on the *abstract syntax descrip*tion language for the target MR, which gives it two major advantages: (1) it is highly accurate, using information from the syntax of the target MR to constrain the output space and model the information flow, and (2) it is highly generalizable, and can easily be applied to new types of MR by just writing a new abstract syntax description corresponding to the allowable structures in the MR. Experiments on four different semantic parsing and code generation tasks show that our system is generalizable, extensible, and effective, registering strong results compared to existing neural semantic parsers.¹

1 Introduction

Carnegie Mellon U School of Com Semantic parsing is the task of transducing natural language (NL) utterances into formal meaning representations (MRs). The target MRs can be defined according to a wide variety of for- ing specifically designed neural networks for each

Because of these varying formalisms for MRs, the design of semantic parsers, particularly neural network-based ones has generally focused on a small subset of tasks — in order to ensure the syntactic well-formedness of generated MRs, a parser is usually specifically designed to reflect the domain-dependent grammar of MRs in the structure of the model (Zhong et al., 2017; Xu et al., 2017). To alleviate this issue, there have been recent efforts in neural semantic parsing with general-purpose grammar models (Xiao et al., 2016; Dong and Lapata, 2018). Yin and Neubig (2017) put forward a neural sequence-to-sequence model that generates tree-structured MRs using a series of tree-construction actions, guided by the task-specific context free grammar provided to the model a priori. Rabinovich et al. (2017) propose the abstract syntax networks (ASNs), where domain-specific MRs are represented by abstract syntax trees (ASTs, Fig. 2 Left) specified under the abstract syntax description language (ASDL) framework (Wang et al., 1997). An ASN employs a modular architecture, generating an AST us-

[ACL 2020]

Incorporating External Knowledge through Pre-training for Natural Language to Code Generation

Frank F. Xu^{*}, Zhengbao Jiang^{*}, Pengcheng Yin, Bogdan Vasilescu, Graham Neubig Carnegie Mellon University

{fangzhex, zhengbaj, pcyin, vasilescu, gneubig}@cs.cmu.edu

Abstract

Open-domain code generation aims to generate code in a general-purpose programming language (such as Python) from natural language (NL) intents. Motivated by the intuition that developers usually retrieve resources on the web when writing code, we explore the effectiveness of incorporating two varieties of external knowledge into NL-to-code generation: automatically mined NL-code pairs from the online programming QA forum StackOverflow and programming language API documentation. Our evaluations show that combining the two sources with data augmentation and retrieval-based data re-sampling improves the current state-of-the-art by up to 2.2% absolute BLEU score on the code generation testbed CoNaLa. The code and resources are available at https://github.com/ neulab/external-knowledge-codegen.

Introduction

Semantic parsing, the task of generating machine executable meaning representations from natural



Figure 1: Our approach: incorporating external knowledge by data re-sampling, pre-training and fine-tuning.

2018; Iyer et al., 2018; Yin and Neubig, 2019) used a variety of models, especially neural architectures, to achieve good performance.

However, open-domain code generation for general-purpose languages like Python is challenging. For example, given the intent to *choose* a random file from the directory contents of the C *drive*, '*C*:\\', one would expect the Python code snippet random.choice(os.listdir(`C:\\')), that realizes the given intent. This would involve not just generating syntactically correct code, but also using (and potentially combining) calls to APIs and libraries that implement some of the desired functionality. As we show in \S 3, current code generation models still have difficulty generating the





Dijkstra certainly got the "difficult to make" part right ...

It took ~40 years



Carnegie Mellon University School of Computer Science





df.drop(df.columns[[0]])



Code retrieval

del df['column name']

(Training set example is somewhat different still) df = df.drop(df.columns[[0]], axis=1)



What about usability?

BLEU score X on some benchmark

Carnegie Mellon University School of Computer Science



Usefulness in real-world scenarios

VS.



Our human study





Step 1: Develop an instrumented IDE plugin

	樻 main	.py	×						
l		#	Example	code,	write	your	program	here	
	3								
1	3								
I									
I									
I									

Carnegie Mellon University School of Computer Science







Step 1: Develop an instrumented IDE plugin

orogra
CODE
jttb
and se
DDE
Э Э





I

Telemetry



Step 1: Develop an instrumented IDE plugin

🛃 main.py 🗵						
1	# Example code, write your program					
2						
3	x = [1, 2, 3]					
4						
5	# BEGIN AUTO-GENERATED CODE -					
6	# <u>bifob6hryl</u> 4hb8jc1 <u>smdajttb</u> -					
7	# query: reverse a list x					
8	$\stackrel{\frown}{}$ # to remove these comments and ser					
9	x = reversed(x)					
10	# END AUTO-GENERATED CODE					
11						





57

Telemetry



Step 2: Compile a set of realistic tasks

Popular coding education websites

Stack Overflow questions





7 representative task categories (with 2 tasks each)



Step 2: Compile a set of realistic tasks

Category	
Docio Duthon	T1-1 Randomly generate a
Dasic Python	T1-2 Date & time format p
	T2-1 Read, manipulate, and
гпе	T2-2 Text processing abou
\cap	T3-1 File and directory cop
03	T3-2 File system informati
Wah Saraning	T4-1 Parse URLs and speci
web Scraping	T4-2 Extract table data and
Web Server & Client	T5-1 Implement an HTTP
WED SELVEL & CHEIII	T5-2 Implement an HTTP
Data Analysis & MI	T6-1 Data analysis on auto
Data Milary 515 & Will	T6-2 Train and evaluate a
Data Vigualization	T7-1 Produce a scatter plo
	T7-2 Draw a figure with th

Carnegie Mellon University School of Computer Science STREDEL N# NEULAB

Tasks

and sort numbers and characters with dictionary parsing and calculation with timezone d output CSV files it encoding, newline styles, and whitespaces pying, name editing ion aggregation ific text chunks from web page d images from Wikipedia page server for querying and validating data client interacting with given blog post APIs omobile data of performance metrics and prices multi-class logistic regression model given dataset t given specification and dataset nree grouped bar chart subplots aggregated from dataset



11

Step 3: 31 participants completing tasks







Highlight: No significant improvement in quantitative measures of speed, code quality, or program correctness when using the plugin



STREDEL N# NEULAB

Carnegie Mellon University School of Computer Science



13

Highlight: Code generation and retrieval are useful in different settings.

Code generation methods work better for simpler queries on basic Python functionality

 140^{-1} 120 100 Tokens 80 -Code retrieval is preferred for more 60 complicated 40 functionality 20 0 Original Length

STREDEL N# NEULAB

Carnegie Mellon University School of Computer Science



Final Length Edit Distance



Highlight: In general, users liked using the plugin! But we have a loooong way to go to make code generation systems practical

- *Combine code generation with code retrieval.* Our results suggest that some queries may be • Provide a unified and intuitive query syntax. We observed that users are not always formulatbetter answered through code retrieval techniques, and others through code generation. We ing queries in the way that we would expect, perhaps because they are used to traditional recommend that future research continue to explore these types of approaches jointly, e.g., search engines that are more robust to noisy inputs and designed for keyword-based search. using hybrid models [40, 41] that may be able to combine the best of both worlds. The NL2Code generation model we experimented with in this study was trained on natural language queries that are not only complete English sentences, but also include references to variables or literals involved with an intent, specially delimited by dedicated syntax (grave users' natural language queries can often be disambiguated by considering the local context accents). As our respondents commented in the post-test survey, getting used to formulatprovided by the source files they were working in at the time, which in turn could lead ing queries this way takes some practice. Future research should consider not only what to better performance of the code generation model. There is already convincing evidence is the most natural way for users to describe their intent using natural language, but also from prior work that considering a user's local context provides unique information about how to provide a unified query syntax for both code generation and code retrieval, to miniwhat code they might type next [111]. In addition, some work on code retrieval has also mize confusion. Robust semantic parsing techniques [8, 95] may also help with interpreting considered how to incorporate context to improve retrieval results [17]; this may be similarly ill-specified user queries.
- Consider the user's local context as part of the input. Our oracle comparison revealed that incorporated.
- Provide dialogue-based query capability. Dialogue-based querying could allow users to refine • Consider the user's local context as part of the output. Considering where in their local IDE their natural language intents until they are sufficiently precise for the underlying models users are when invoking an NL2Code assistant can also help with localizing the returned to confidently provide some results. Future systems may reference work on query reforcode snippets for that context. Some transformations are relatively simple, e.g., pretty printmulation in information retrieval, where the user queries are refined to improve retrieval ing and indentation. Other transformations may require more advanced program analysis results both for standard information retrieval [7] and code retrieval [39, 45]. In addition, in but are still well within reach of current technology, e.g., renaming variables used in the the NLP community there have been notable advancements recently in interactive semantic returned snippet to match the local context (the Bing Developer Assistant code retrieval parsing [51, 119], i.e., soliciting user input when dealing with missing information or ambiengine [115] already does this), or applying coding conventions [2]. guity while processing the initial natural language query, which could be of use as well.
- *Provide more context for each returned snippet.* Our study shows that NL2Code generation or • Consider new paradigms of evaluation for code generation and retrieval systems. Usage log data, retrieval systems can be useful when users already know what the right answer is, but they such as the ones we collected here, is arguably very informative and useful for researchers need help retrieving it. At the same time, many of our study participants reported lacking looking to evaluate NL2Code systems. However, compared to automated metrics such as sufficient background knowledge, be it domain-specific or API-specific, to recognize when a BLEU, such data is much less readily available. We argue that such data is worth collecting plugin-returned code snippet is the right one given their query, or what the snippet does in even if only in small quantities. For example, with little but high-quality data, one could still detail. Future research should consider incorporating more context and documentation totrain a reranker [125] to try to select the outputs that a human user selected; if the predictive gether with the plugin's results, which allows users to better understand the code, e.g., links power exceeds that of BLEU alone, then the trained reranker could be used to automatically to Stack Overflow, official documentation pages, explanations of domain-specific concepts, evaluate the quality of the generated or retrieved code more realistically than by using BLEU. other API usage examples. One example of this is the work of Moreno et al. [78], which retrieves usage examples that show how to use a specific method.

STREDEL N# NEULAB





Conclusion: We need more human-centered approaches to evaluate ML-based tools within the software development workflow

neulab/tranX-plugin Public									
<> Code 🕢 Issues 3 1. Pull requests 🕞 Actions 🗄 Projects 🖽 Wiki 🕩 Security 🗠 Insights									
ੇ ਸਿ master ਚ ਮੈਂ 1 bran	ch tags Go to file Add file	▼ Code ▼	About A plugin for code generation in PyCharm/IntelliJ using tranX						
frankxu2004 add	user id settin f8c1099 12 days ago	C 64 commits							
gradle/wrapper	update http client implementation a	2 years ago	C Readme						
imgs	update usage	2 years ago	 ☆ 33 stars ② 2 watching 						
src	also upload unselected queries to th	2 years ago	%0 forks						
 .gitignore	update http client implementation a	2 years ago							
README.md	add user id setting in readme	12 days ago	Releases 9						
build.gradle	bump ver	2 years ago	S Beta release 1.8 Latest						
🗋 gradlew	Initial commit of tranx plugin	3 years ago	on Sep 8, 2020						
🗋 gradlew.bat	Initial commit of tranx plugin	3 years ago	+ 8 releases						
Settings.gradle	Initial commit of tranx plugin	3 years ago	Packages						

https://github.com/neulab/tranX-plugin

Carnegie Mellon University School of Computer Science

STREDEL N# NEULAB



16