

Text-to-SQL Generation for Information Retrieval

Jacob Waffle

Boise State University

jakewaffle@u.boisestate.edu

Abstract

and utilized for our evaluation needs. One helpful aspect of this dataset is that the expected query output is already separated into a list of tokens and that works well for a content overlap algorithm.

1 Introduction

2 Background

3 Data

3.1 Demo

The demo application for this project utilized a PostgreSQL database composed of three tables: movies, casts and actors. The movie and actor tables listed out movies and actors with no foreign key relationships, while the cast table joined together those tables to describe which actors starred in which movies. This database was filled with data parsed from a UCI Machine Learning movie dataset (Wiederhold, 1999). The dataset held a lot more information, but we only needed enough information to test our SQL generation for the retrieval of a single REST resource. One critical requirement for our database was a distribution of the data across multiple tables so that the SQL generation required the usage of JOIN clauses.

The movie dataset itself was easy to parse since it had consistent delimiters for the columns of data it supported. We only needed to split the rows of the file by the delimiter and then map the chosen columns to the columns in the appropriate tables in the database. Some text manipulation was required due to some of the columns having unknown data and there being unwanted prefixes for other columns.

3.2 Evaluation

Our evaluation strategy also required its own dataset for ensuring our model was correctly generating SQL queries. We utilized the Spider Text-to-SQL dataset (Yu et al., 2019), which is widely used for this area of study. It came in the form of Json and was easily able to be parsed

4 Method

The goal of this project was to assess the production viability of generating SQL from a natural language query for a common REST server. The model itself needed to be proven to be effective in a general sense and against a real database for a real use case. The integration of the model into the application also needed to account for the possibility of many concurrent users. It is also hoped that the model can be used in a general-purpose manner so that it can be quickly adapted to whatever use case it might be beneficial for.

The T5 model (Raffel et al., 2020) is short for Text-to-Text Transfer Transformer and is one of the beginning auto-regressive Transformer models for text generation. It was created by Google in 2020 following GPT-2 in 2019 (Radford et al., 2018) and served as their solution for translation and summarization tasks in the wake of the shift to transformers from Recurrent Neural Networks. The T5 has since been replaced at Google by FLAN-T5 (Chung et al., 2022) and PaLM (Chowdhery et al., 2022), but still remains as one of the state of the art models used for Text-to-SQL generation.

Despite there being models that are technically better than the T5 model, we still used it for our assessment because it has been proven to work for our needs. The major difference in our task versus the others in this space is that our demo application requires the generated SQL to target a predetermined table and produce a predetermined set of columns. Another reason the T5 model is a good choice is that this requirement is easily able to be solved by any auto-regressive Transformer as long as it is paired with a constrained decoder.

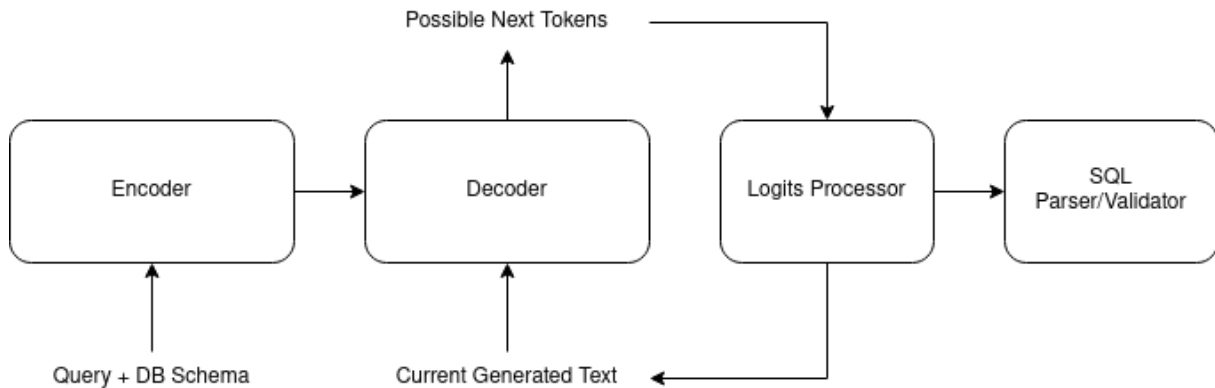


Figure 1: Model Overview

5 Evaluation

5.1 Task & Procedure

5.1.1 Model

We started our efforts with a T5 model that was already fine-tuned by the Picard project (Scholak, Schucher, & Bahdanau, 2021). This model was used because it is one of the state of the art models in this area of study and is easy to build off of due to its compatibility with Huggingface (Wolf et al., 2020). We also benefitted from this model’s database schema agnosticity because it expects the target database schema as an input and that meant that further fine-tuning to support a particular database is unnecessary. This is the main reason that no further fine-tuning was done.

We augmented this model by pairing it with a custom Huggingface Logits Processor that utilizes a PyParsing PostgreSQL SQL grammar (PyParsing, n.d.) to parse the generated SQL and ensure the chosen tokens contribute to a valid SQL output. This was done to implement a constrained decoder for our model and is useful because it helps ensure that the overall generated SQL is valid. We could have used Picard’s Logits Processor, but chose not to because their solution utilized Haskell and thus would have been difficult to modify. Figure 1 illustrates how the Transformer, Logits Processor and SQL parser work in tandem.

A constrained decoder approach is useful for text generation with Transformers that are meant for a specialized task like code generation. A model can be trained to do a lot of things, but it’s debatable how perfect of a job it can do in all situations. A constrained decoder fills in the gaps by steering the text generation in the right direction. It can even be argued that a constrained decoder allows

for a smaller model, with less parameters. Picard’s constrained decoder for instance knows more details about the target database schema than what is provided in the inputs to the encoder and can thus make sure the proper operators and literals are used based on the type of table columns that they are being used alongside.

In addition to a custom Logits Processor, we also added a forced prefix to the generated SQL output by utilizing Huggingface’s support for a restricted vocabulary. This was done to ensure that the generated SQL targeted the expected table and outputted the expected columns. Our demo application relied on this behavior due to the assumption that the generated SQL would always target a predetermined table and would return all of the columns related to that table. Our constrained decoder strategy likely wouldn’t have been possible without the auto-regressive nature of our model and wouldn’t have met our requirements for generating valid SQL that retrieves predetermined information.

For the general evaluation of this model we utilized the Spider Text-to-SQL dataset (Yu et al., 2019) to verify that the model by itself is capable of generating reasonable SQL results. This was done in a unit test environment.

We also relied on a human evaluation for our demo application as a whole. Refer to Figure 2 for an idea of the overall architecture. This was necessary to evaluate how well our Text-to-SQL strategy worked against a real database for a real use case. Part of the evaluation was determining what was necessary to utilize the generated SQL for real database queries, while the other part was validating that the generated SQL made sense based on the natural language queries that we used for our movie database.

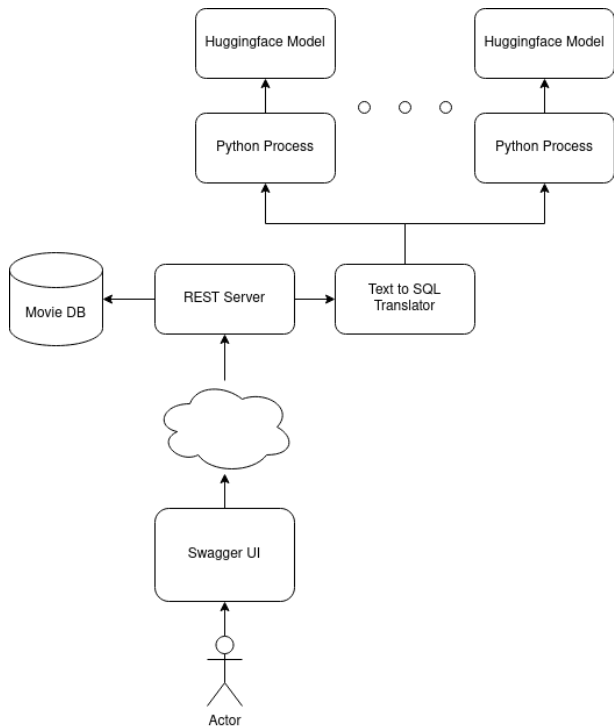


Figure 2: Demo Architecture

5.1.2 Framework

The model was a core effort for this project, but the framework around that model is equally important. In order for the model to be useful in a production environment it needs to support multiple users to some degree. This ultimately needed to be a focus of the project to assess production viability.

The approaches explored for scalability included spawning Python processes with their own model based on the load of the application and utilizing batching for the queries. Figure 2 illustrates the overall architecture for the demo application. The Text to SQL Translator piece was written to be a separate library that can be stood up on multiple computers if necessary and its main feature is that it spawns worker Python processes with our model based on the translation load. This was accomplished using Deep Java Library’s model serving utility (DJL, n.d.).

DJL does a lot of things, but its serving utility claims to automatically support scaling up/down workers based on the load and dynamically batching requests. Although, these features weren’t fully explored and evaluated in this project, we still have incorporated its usage into the overall project with the assumption that its features work as expected. More exploration is required for how this library can be configured to fully accomplish a produc-

tion environment’s scaling needs while taking to account excessive RAM usage of the models and ensuring batch sizes aren’t too large.

5.1.3 Metrics

The various evaluation efforts used in this project required different metrics tailor made for each of the strategies. The main evaluation utilized the Spider dataset to test the quality of the generated SQL and the results were measured using a content overlap algorithm. We knew what sequence of tokens were expected to be in the resulting query and also which tokens were actually present in the generated SQL. Therefore the resulting content overlap was measured as an overall percentage of the expected tokens that matched the actual results. A side effect of this strategy is that the expected queries with a high amount of tokens contributed more to the resulting percentage than expected queries with less tokens.

5.1.4 Results

The work explored in this project could easily be developed into a general-purpose library that has potential for being commonly used by any application that requires information to be retrieved from a database based on a user’s needs. In the same way ORMs are designed to make database storage easier, this work makes the job of translating a user’s request into the information that they wish to see easier. However, much like with ORMs this solution isn’t going to replace the SQL that an expert is capable of writing for use cases that require efficient/performant queries and complete accuracy of the results. This approach will satisfy most of a user’s needs, but will likely never be the right tool for every use case in the same way that the UDP protocol isn’t the right tool for all communication.

The final percentage that we achieved for the Spider dataset evaluation was about 61.2% and that again represents the percentage of the total expected tokens matching what was in the generated SQL. The Picard project was able to achieve a higher percentage than ours at 74.8%, but also had a much better constrained decoder implementation that took into account the actual data types of the columns in the target database schema. Their incremental parser was also probably better in general as well.

The human evaluation of the overall demo application didn’t necessarily have results that can be measured, but we have observed some interesting

details and obstacles. Unlike the Spider evaluation, the demo application requires that the generated SQL has a predetermined set of column outputs and targets a particular table. In hindsight, forcing a prefix worked but also artificially gave the decoder some unintended cues towards what the rest of the SQL should include. With the anticipation that the generated SQL might generate a JOIN clause, the forced prefix included an alias to the main table that was to be queried. This resulted in almost all of the generated SQL containing a JOIN clause regardless of its actual necessity and usage in the rest of the query. It's possible that the training data would need to be adapted to correct this minor issue.

In the development of the model, we have also observed that the model and SQL parser can sometimes get into states where they appear to “argue” with each other over what parts of the SQL is and isn't valid. This mostly occurs when incremental parsing is enabled instead of only parsing the generated SQL when an end token is encountered. The only explanation we have is that the parser isn't failing fast enough during the generation due to how lenient it was written in order to account for partial word token, which when observed by themselves may be seen as invalid. If the parser confirms that a partial word token is valid in one pass and then invalid when the next token completes the word, then it's believable that the generation process might get into a state of confusion. In these situations the generated SQL usually grows to be very large and we timeout before the end of the sequence is reached.

6 Conclusion

6.1 Implications

6.2 Limitations

6.3 Future Work

(Scholak et al., 2021)

(Wolf et al., 2020)

(Yu et al., 2019)

(Wiederhold, 1999)

References

Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., ... Fiedel, N. (2022). *Palm: Scaling language modeling with pathways*.

Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., ... Wei, J. (2022). *Scaling instruction-finetuned language models*.

DJL. (n.d.). *Deep Java Library Serving*. <https://docs.djl.ai/docs/serving/index.html>. ([Accessed 24-Apr-2023])

PyParsing. (n.d.). *PyParsing: Python library for creating PEG parsers*. <https://github.com/pyparsing/pyparsing/>. ([Accessed 24-Apr-2023])

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2018). Language models are unsupervised multitask learners. Retrieved from <https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf>

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... Liu, P. J. (2020). *Exploring the limits of transfer learning with a unified text-to-text transformer*.

Scholak, T., Schucher, N., & Bahdanau, D. (2021). *Picard: Parsing incrementally for constrained autoregressive decoding from language models*.

Wiederhold, G. (1999). *Movie*. UCI Machine Learning Repository. (DOI: [10.24432/C5SW2R](https://doi.org/10.24432/C5SW2R))

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., ... Rush, A. M. (2020). *Huggingface's transformers: State-of-the-art natural language processing*.

Yu, T., Zhang, R., Yang, K., Yasunaga, M., Wang, D., Li, Z., ... Radev, D. (2019). *Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task*.