

# Learning Pattern-Based Extractors from Natural Language and Knowledge Graphs: Applying Large Language Models to Wikipedia and Linked Open Data

Célian Ringwald

Université Côte d'Azur, Inria, CNRS, I3S  
celian.ringwald@inria.fr

## Abstract

Seq-to-seq transformer models have recently been successfully used for relation extraction, showing their flexibility, effectiveness and scalability on that task. In this context, knowledge graphs aligned with Wikipedia such as DBpedia and Wikidata give us the opportunity to leverage existing texts and corresponding RDF graphs in order to extract, from these texts, the knowledge that is missing in the corresponding graphs and meanwhile improve their coverage. The goal of my thesis is to learn efficient extractors targeting specific RDF patterns and to do so by leveraging the latest language models and the dual base formed by Wikipedia on the one hand, and DBpedia & Wikidata on the other hand.

## Introduction

Whether automatically extracted from structured elements of articles or manually populated, the open and linked data published in DBpedia and Wikidata offer rich and structured complementary views of the textual descriptions found in Wikipedia. However, the unstructured text of Wikipedia articles contains a lot of information that is still missing in DBpedia and Wikidata. Extracting them would be interesting in order to improve the coverage and quality of these knowledge graphs (KG) since they have an important impact on all downstream tasks. This thesis proposes to exploit the dual bases formed from Wikipedia pages and Linked Open Data (LOD) bases covering the same subjects in natural language and in RDF, in order to produce RDF extractors targeting specific RDF patterns and tuned for a given language. Therefore, the main research question addressed in my thesis is:

**RQ** – *Can we learn efficient customized extractors targeting specific RDF patterns from the dual base formed by Wikipedia on one hand, and DBpedia and Wikidata on the other hand?*

Formally, let  $Db$  be a dual base, subset of  $W \times G$  where  $W$  is the set of Wikipedia articles and  $G$  is the set of corresponding RDF graphs in DBpedia and Wikidata. The aim is to learn from this dual base an extractor  $E_{Db} : W \rightarrow L; (t, S) \mapsto E_{Db}(t, S) = g$ , where  $L$  is the LOD,  $t$  is an input text,  $S$  is a set of RDF patterns expressing constraints

represented as SHACL shapes, and  $g$  is an RDF graph implied by  $t$  and valid against  $S$ .

This question is closely tied to the relations extraction (RE) task which consists of retrieving relations from unstructured texts. Until recently, the RE task was solved by complex pipelines including multiple steps. But this approach leads to error accumulations and propagation (Mesquita et al. 2019). However, the latest progresses in NLP including pre-trained language models (PLM) have drastically improved the performance of many downstream tasks (Ye et al. 2022).

## An Incremental Research Plan

Following an incremental methodology, I intend to generalize the approach by relaxing one constraint at a time: starting from the generation of a single triple pattern before generalizing to arbitrary basic graph patterns. I will do so by addressing sequentially the following sub-questions:

**SRQ.1** – *How to survey and follow the latest trends in PLM-based KG extraction?*

The landscape of the research field drawn at the intersection of language models and knowledge graphs is very dynamic and quickly evolving. Consequently, a systematic literature review to closely follow it, is crucial in this context.

Then, leveraging some of the latest techniques identified, the next research questions are set to find the currently best-performing approaches for a gradually more complex version of our task.

**SRQ.2** – *Which aspects of the task formulation impact the generation of triples with datatype properties?*

The performance of the task learned on top of an existing language model crucially depends on the formulation of this task. In this end, the choice of a specific RDF syntax for the extraction, as well as the content of the prompt given as input, are sensitive parameters in order to take full advantage of the PLM. One of the first steps in this work will be to determine the combination leading to the best results when focusing only on the datatype properties.

**SRQ.3** – *How to jointly extract datatype properties and object properties for a KG?*

The scientific community has recently underlined the “hallucination” problem of PLMs. In practice, this issue may affect a large proportion of the triples containing literals (e.g. attributes such as dates, measures, textual descriptions, etc.).

These relations are defined in the OWL semantic web language as datatype properties. In the literature, the research conducted until today was more focused on the extraction of object properties that link together two objects. I will have to propose a method to jointly extract both types of properties.

**SRQ.4** – *How to support fact extraction relying on different document granularity?*

Relations can span different levels of a document (sentences, paragraphs, sections, etc.). The recent development of information retrieval techniques based on embedding seems to be a promising solution that must be adapted to our context.

**SRQ.5** – *What is the best strategy to extract rare relations and under-represented instances of classes?*

The state-of-the-art models struggle with the under-representation of some facts. A lot of improvement must be made in this direction to be able to extract relations beyond those that are highly represented. In that respect, data augmentation methods for generating negative or rare relation are a first step to expand few-shot learning capacities of the current models.

## Preliminary Results

To this day, I have built a tool to automate my literature review and conducted two preliminary experiments, in order to approach SRQ.1 and SRQ.2. This work and the results obtained from them will soon be formalized, extended, and submitted for publication.

### A Living and Systematic Review

The answer to SRQ.1 requires kick-starting a literature exploration from existing surveys related to a given task via querying several digital library APIs. The scientific corpus collected was extended with: (1) the dataset and models related to our task, and the various metadata found on PaperWithCode; (2) the retrieval of the citation network of each paper using the OpenCitation API. The first run of the literature exploration allowed me to draw the current trends in RE. It suggests that the field is shifting from discriminative (encoder-based models derived from BERT) to generative modelling (based on encoder-decoder and decoder-only architectures). The latter models have the advantage of being flexible and enabling the conception of end-to-end systems. Moreover, they better handle overlapping triples extraction (Ye et al. 2022). Finally, to my knowledge, no system is currently trying to perform semantic relation extraction where triples explicitly follow an RDF syntax.

### A First Focus on Datatype Properties

In a first experiment I tackled the first part of SRQ.2, by focusing on datatype properties that are primarily affected by the hallucination problem. I built a first dataset of relations based only on the English chapter of DBpedia and Wikipedia, restricted to entities of type `dbo:Person`, and focusing on the following relations: `dbo:birthDate`, `dbo:deathDate` and `rdfs:label`. I employed a SHACL shape to filter the graphs respecting the following criteria: an instance of `dbo:Person` must have a `dbo:birthDate`, an

`rdfs:label` and a `dbo:deathDate`. On the example of REBEL (Huguet Cabot and Navigli 2021), I fine-tuned a BART model by giving as input the identifier of the entity followed by its Wikipedia abstract. The model was trained to generate RDF Turtle triples including the set of relations selected and found in DBpedia. A qualitative analysis of the errors allowed me to underline 3 sources of errors: (1) the fact is in the text but not in DBpedia, (2) the values in the text and in the DBpedia are different, (3) the fact is in DBpedia but not in the text.

I conducted a second experiment to understand if one syntax is easier to learn for a pre-trained model. In the literature, the choice of syntax is related to the “linearization process” where triples are serialized as a string (a list or a tagged sequence). Until now, different methods have been investigated but they were not rigorously compared. For this reason, I extended the dataset of my first experiment to represent the triples according to seven different syntaxes: a simple list, a tagged sequence, a light Turtle syntax without the prefixes definition, full Turtle, N-Triples, XML and JSON-LD. My approach also considered the fine-tuning of two sizes of the BART model (base, large) and the training time needed before the F1-micro saturation ( $>0.9$ ). The results of the experiments showed us that (1) the model quickly mastered the simplified Turtle syntax, followed by the list and the tags. (2) Some RDF syntaxes took longer to learn for BART: the easiest was Turtle, followed equally by JSON-LD and RDF-XML. (3) We have also pointed out that N-Triples syntax were harder to learn. However, this experiment must be extended by incorporating cross-validation and other measures better suited to assess the generated triples.

## Future Works

Until the date of the Workshop, my research plan will be to continue to answer SRQ.2 and will extend to SRQ.3, by incrementally investigating new and more complex RDF graph patterns and integrating Wikidata to go beyond our current use of DBpedia. Aside from that, I will continue the analysis allowed by the systematic review launched to answer SRQ.1. I will then have another two years to cover SRQ.4 and SRQ.5. Depending on my advancement some extensions of the current plan could consider including approaches evaluating transfer learning and active learning.

## References

- Huguet Cabot, P.-L.; and Navigli, R. 2021. REBEL: Relation Extraction By End-to-end Language generation. In *Findings of the ACL: EMNLP 2021*, 2370–2381. ACL.
- Mesquita, F.; Cannavicchio, M.; Schmidek, J.; Mirza, P.; and Barbosa, D. 2019. KnowledgeNet: A Benchmark Dataset for Knowledge Base Population. In *Proc. of the 2019 EMNLP Conference and the 9th International IJCNLP*, 749–758. ACL.
- Ye, H.; Zhang, N.; Chen, H.; and Chen, H. 2022. Generative Knowledge Graph Construction: A Review. In *Proc. of the 2022 EMNLP Conference*, 1–17. ACL.