

AI Supported Knowledge Graph Design & Generation

Marco Ratta^[000–0003–3788–6442]

The Open University, Walton Hall, Milton Keynes (UK) marco.ratta@open.ac.uk

Abstract. Knowledge Graphs (KG) have risen to be a powerful mechanism to represent data. Despite this most data sources are generally still represented via heterogeneous non-graph data structures. Converting these into KGs necessitates considerable effort from experts, proving this to be a time consuming process. While tools have been developed to aid KG builders, a gap still exists in terms of technologies that support the automation of designing KG building pipelines. Addressing this gap motivates this research. The aim is to first understand the problem at the knowledge level and, inspired by the recent release of generative tools such as *GPT-Engineer*, to put forward a conversational agent aimed at assisting the user in building their pipelines. We report on the preliminary findings that we have so far reached during the first year of research in deriving the requirements for building KG generating pipelines from the literature.

Keywords: Knowledge Graph Generation · Knowledge Graph Design · Artificial Intelligence · Knowledge Graph Refinement

1 Introduction and Motivation

Knowledge Graphs [13] have risen to be a powerful mechanisms to represent data and have been shown to play an important role for many applications, including conversational agents, data integration, and recommender systems [28].

However, most data sources come in heterogeneous non-graph data structures, schemes, and formats [6]. Converting these sources into KG representations involves a range of specific tasks. These span from mapping tasks, such as mapping a specific data type to its specific URI format, through to ontology engineering tasks such as aligning data specific vocabulary with known vocabularies. A central problem is also the organisation of these tasks into a pipeline whereby one is able to go from the original data to the KG. This process necessitates a considerable effort from both KG engineers and domain experts, showing this to be an often time consuming and complicated process [2].

While advances have been made in terms of supporting KG engineers in their work, “existing techniques offer [...] solutions, each covering a specific aspect of KG construction, but automatically orchestrating [...] the whole construction process remains the challenge” [14]. Addressing this challenge poses questions such as, for example, given data sources and a target ontology can a pipeline from

the sources to the target KG be generated automatically? If human intervention is required, what is the nature and extent of this input? And what is the role and extent of background knowledge? Addressing this challenge and questions is the motivation behind this research.

To accomplish this, we first aim to understand the problem of KG design and generation at the knowledge level, beginning with a review of the relevant literature to look at how KGs have concretely been built. The purpose is to abstract the required KG tasks and operations that ought to be supported, and from these to build a systematic model of the same and to combine them into pipelines inspired by problem solving methods. Influenced by recent generative tools such as *GPT-Engineer*¹, where the user interacts conversationally with the agent to develop code, the final aim is to build a conversational agent that leads a user through a problem-solving tree on the basis of its own knowledge base and the information provided by the user at each step. This whilst employing "under the hood" either rule-based or generative AI based techniques and/or tools to solve the required tasks and at the end outputs a KG view of the input source data.

Thus, our objectives are: (i) enriching the theoretical understanding of KG design and generation by looking into the existence of abstractable KG building tasks (i.e. that go beyond a single use-case) and using these to formulate abstract high-level tasks, (ii) extending the application domain of AI to the design of KG generating pipelines through the orchestration of tasks and (iii) supporting the community of practitioners by releasing novel methods for designing and organising KG building.

The remainder of the paper is structured as follows. First, we delve into the state of the art, then we state our research questions and contribution. Subsequently we proceed to describe our intended methodology and the evaluation approach. Finally we report on the early findings that we have so far achieved in the abstraction of general KG construction tasks from the literature during this first year of research.

2 State of the Art

As it is not assumed that the required tasks and the methods for solving them are known beforehand, we take a knowledge engineering approach to tackling our problem. This entails that the practice of KG engineers and creators is taken as "a source of knowledge that cannot be obtained from anywhere else" [22]. Central to a knowledge engineering approach are also problem solving methods (PSMs) [9], which are domain-independent reasoning components specifying reusable patterns of behaviour. Furthermore, when we speak of automation we usually refer to the automation of a process at the knowledge level, where we consider a solution to the problem via the design of actions, goals and body of a hybrid symbolic-generative agent.

¹ <https://github.com/gpt-engineer-org/gpt-engineer>

The application domain of this research falls under the umbrella of data integration processes. Lenzerini [17] formalises such a process. This combines data from heterogeneous sources into a single unified view available to a user or client. Its basic components are the heterogeneous data sources, a global schema acting as a mediator and view of this underlying data and the mappings that connect the information contained in the sources to the schema of the global view. For us this global view is a KG.

By KG we understand a graph intended to accumulate and convey knowledge of an object whose nodes represent entities of interest and whose edges represent relations between these entities [13]. A number of approaches and tools have been developed to aid KG engineers in their work. Central to the task of KG construction has been the effort put into the development of W3C recommended R2RML/RML mapping languages for the creation of RDF mappings[7] and the tools that materialise them [6][1]. Further technologies include methods for resolving data heterogeneity such as Ontology Based Data Access (OBDA) [4][3], Extract Transform Load (ETL) based tools [25], and SPARQL based tools [16]. Further, there are also multi-agency [21] and deep-learning [27] approaches for schema matching and ontology mapping, and a variety of fuzzy logic, probabilistic soft logic and machine/deep learning methods [14].

One of our hypotheses is that we can frame the problem of designing KG generation pipelines as one of graph-to-graph manipulation. The *Facade-X* methodology [2] is therefore one of the starting points of the artifact development process as it can resolve from the start the heterogeneity of the data sources into graphs. This methodology assumes a two-phase iterative process: (a) re-engineering, where the task is to resolve the heterogeneous formats of the sources into a graph, and (b) re-modelling, where the main objective is to add semantics [5]. High-level tasks undertaken during these include: (i) specifying a representative ontology for the sources, (ii) extracting the information, (iii) generating the graph from the mappings and (iv) KG evaluation and quality assessment. Crucially, this hypothesis and starting point allows to connect our approach with possible implementations from the field of graph transformation, which focuses on geometrical and rule based methods for graph to graph manipulations [11], from ontology matching and alignment [24][15] from knowledge graph refinement [23] and from ontology design patterns [8] and engineering [10].

Recently a lot of the attention has focused on large language models (LLMs) as possibly universal means for the development of AI and real-world programming tasks [26], although not without debate in the community [20]. Other than the previously mentioned *GPT-Engineer*, *PandasAI*² provides another example of using a conversational agent for working with databases. The use of LLMs for the generation of KGs has also been the object of investigation by the academic community [19].

Finally, we aim to implement a methodology based on the principles of design science [12]. This is an iterative process that can be viewed as three closely related cycles of activities that are happening simultaneously: the relevance cycle,

² <https://github.com/Sinaptik-AI/pandas-ai>

the rigor cycle and the design cycle. The relevance cycle is concerned with the introduction of input requirements from the contextual environment into the development process and introducing the research artifacts into environmental field testing. The rigor cycle aims to ensure the proper theoretical grounding of our development work. The design cycle aims at generating design alternatives and evaluating them against the requirements provided by the other cycles.

3 Problem Statement and Contributions

The overarching research question of the project is the following:

(RQ1) How can the development of Knowledge Graph generating processes from structured and semi-structured data sources be automated?

The answer to this question can be broken into two interconnected parts, each with their own sub-questions. For the first part, concerned with requirements and theoretical rigor, we ask:

(RQ2) are there general features, specifically tasks, in Knowledge Graph generation pipelines that are shared and that can be abstracted?

(RQ3) how can we systematise our findings into a model of the process of designing Knowledge Graph generating pipelines at the knowledge level?

For the second part, concerned with implementation, our target is that given a set of data sources and a set of ontologies, the system aids in the designing of a KG generating pipeline derived from the given inputs. Our assumption is that this can be achieved through a hybrid symbolic and generative AI system, supported by breaking down the problem into tasks that can be addressed by PSMs that are to be applied to graph-to-graph transformations. But is this actually required and is it beneficial with respect to an approach based say solely on using an LLM and prompt engineering? Therefore,

(RQ4) Does the breaking down of the problem into tasks and methods help with the creation of KG generating pipelines?

(RQ5) Can the process of converting and integrating structured and semi-structured sources be adequately modelled as a problem of graph-to-graph transformations?

It is also the case that, for example, in a task of entity alignment where the same entity is referenced by name in one source and by an identification number in a second sources requires a background knowledge that only a user acquainted with the sources can provide. Thus, it can be asked

(RQ6) how much of the process of KG design and generation can be automated?

(RQ7) Which parts of the process require necessary human intervention (e.g. providing background knowledge)? What is the nature and extent of this intervention?

(RQ8) What is the impact of this intervention on the generated pipelines? Can these be linear from source to target or must they include iterative problem solving structures?

4 Research Methodology and Approach

To carry out the research programme we intend to implement a methodology based on the principles of design science as follows.

The first step is to conduct a literature review of relevant sources. As parts of the relevance cycle and the rigor cycle are concerned with introducing requirements into the development process, the plan is to collect these empirically by looking at how KG builders have operated concretely. The target is then twofold. We look first for the KG task required to answer the research questions of the first part (RQ1-3) and we collect a variety of other data related to concrete KG pipelines. This includes for example the provenance of the sources (i.e. numerical, textual, etc.), their formats (csv, xml, etc.), the producing stakeholders of the graphs (semantic web researchers, government institutions, etc.) or the tools that have been used to generate the graphs. As this PhD is in its first year this step is being carried out at the moment.

The second step is to classify and systematise the tasks abstracted in the first to build the knowledge level model that will contain the answers to the research questions of the first part of the PhD (RQ1-3). Following this is the development on the basis of this model of a series of problem solving methods. These initial three steps are required in order to complete the theoretical background of the rigor cycle.

Once the requirements have been formalised, the next step is to begin to implement the design cycle and the application of the problem solving methods that have been developed to graph-to-graph transformations. Within this a number of crucial design choices will also require implementation. These include, for example, determining the distribution of the tasks between the artifact and the user (RQ6), the specification of the role of the user interacting with the system and the necessary background knowledge (RQ7-8) and the task distribution model between symbolic and generative AI methods.

Finally, to conclude the relevance and design cycle, we plan to implement an evaluation of the artifact produced through a compare and contrast approach, using a dataset of KGs that have been built with already tried and tested methods and thus also evaluating some of our initial assumptions (RQ4-5).

5 Evaluation Plan

We provide below a high-level description of how we plan to evaluate the work.

The plan is to perform a final-state based evaluation of the outcomes (*EV1*) and a fine-grained evaluation of the process leading to those outcomes (*EV2*).

In terms of outcomes, we consider the final KG output (*EV1.1*) and the KG building pipeline that the system plans (*EV1.2*). For (*EV1.1*) we plan on using a selection of (i) quantitative metrics [18] and (ii) competency questions. For (*EV1.2*) we plan on evaluating (i) functional correctness (ii) pipeline execution feasibility and (iii) pipeline generation efficiency, i.e. the time taken for this generation to occur and whether it is sustainable for a user workflow.

For the fine-grained evaluation of the process (*EV2*), this means evaluating the recommended transformation tasks to perform, the recommendations (vocabulary terms, types, relations, extension actions, etc.) that the system provides for these actions, the ability to break each of them down into cognitively meaningful steps (explainability), the extent to which the user is capable of tailoring the behaviour of the system and the ability to generate correct and executable code. To measure this we plan on relying on precision and recall metrics and an expert based evaluation looking for (i) accuracy, (ii) quality, (iii) effort, (iv) explainability and (v) user agency, understood as the ability of the user to be in control of the process.

We plan on implementing our evaluation work by instantiating the following KG building scenarios:

(S1) *Given the data sources and multiple ontologies, select the subset of representative ontologies and generate the KG building pipeline.*

(S2) *Given the data sources and the target ontology, generate the KG building pipeline.*

(S3) *Given applicable tasks and data sources, build a KG graph generating the relevant semantics.*

To do so we plan on making use of benchmarks already published by the Semantic Web community, such as *GTFS-Madrid*³, *LUBM*⁴, *BSBM*⁵ and *NPD*⁶ and a selection of the KGs surveyed for the literature review for which a published expert manual/semi-automatic/automatic solution for the KG exists.

Following FAIR⁷ open science practices, we plan on releasing the whole source code on *GitHub* to make it available for reproducibility.

6 Preliminary Findings

We now report on the early findings from the literature review. This is being conducted on a selection of papers from four publications, the *Semantic Web Journal*⁸, the *Journal of Web Semantics*⁹ and the conference proceedings from the *Extended Semantic Web Conference* and the *International Semantic Web Conference*. These have been selected as the main venues of Semantic Web research. Isolating a small sample of relevant papers deemed relevant, we constructed a keyword search¹⁰ which we ran on the databases over the period 2013-2024. We ordered the results using the publications' search functions for relevance and have gone through the first 100 hits of each search. By reading the title and abstracts we targeted papers that describe an actual KG construction pipeline built from any source and excluded papers describing tools, systems or ontologies without any use in an actual KG construction process. 67 papers have been

³ <https://github.com/oeg-upm/gtfs-bench> ⁴ <https://swat.cse.lehigh.edu/projects/lubm/>

⁵ <http://wbsg.informatik.uni-mannheim.de/bizer/berlinsparqlbenchmark/>

⁶ <https://github.com/ontop/npd-benchmark> ⁷ <https://www.go-fair.org/fair-principles/>

⁸ <https://www.semantic-web-journal.net/> ⁹ <https://www.sciencedirect.com/journal/journal-of-web-semantics>

¹⁰ ("data transformation" OR "data transformations") AND ("knowledge graph" OR "knowledge graphs" OR "linked open data" OR "linked data")

selected to go to a further screening process where, applying the same criteria that a paper must describe an actual KG building process. 65 of them have been selected for reading.

The goal of the literature review is to provide an empirical ground, based on the published KGs in the scientific literature, to the tasks that a system built for producing a KG from given sources ought to be able to perform. Furthermore, we collect a number of other data related to the process of building KGs (e.g. ontologies re-used, tools used, custom solutions, etc.). After reviewing 25% of the papers, we have collected a number of statistics and abstracted a series of tasks (Table 1). An interesting datum is that so far in 62.5% of the papers reviewed the authors explicitly mention the use of custom scripts and solutions for either performing tasks or for orchestrating the whole pipeline, which we interpret as evidence for the existence of the problem that we wish to tackle.

Table 1: KG building tasks

Task ID: Name	Task Description
T1: URI design	Constructing an expression pattern for a URI.
T2: URI mapping	Mapping an element’s type to its specific URI pattern.
T3: Prefix mapping	Mapping a URI to its target prefix.
T4: Type mapping	Mapping a data type of the sources to the corresponding target type of the ontology.
T5: Entity linking	Mapping an entity/object in the source to the corresponding entity in an existing taxonomy/graph.
T6: Entity resolution	Establishing whether multiple objects refer to the same object.
T7: Domain ontology selection	Selecting ontologies that are relevant to the domain.
T8: Ontology feature selection	Selecting the parts of an ontology that are descriptively relevant.
T9: Ontology composition	Assembling a plurality ontology parts into a single domain model.
T10: Ontological requirements specification	Specifying the ontological requirements for representing the underlying sources.
T11: Database interlinking	Generating triples that link a RDF graph to another.
T12: Ontology extension	Constructing a custom extension of an ontology.
T13: Schema validation	Checking a graph conforms to the given schema or ontology.
T14: Content validation	Validating the semantic content of a graph.
T15: Syntax validation	Validating the syntactic content of a graph.

T16: Graph syntax selection	Selecting the syntax in which to express a graph.
T17: Source to RDF mapping	Mapping source components to their respective RDF terms.
T18: Language tagging	Adding language tags to literals.
T19: Links explicitation	Rendering implicit links in the sources explicit via triples.
T20: Subclassing	Mapping a source class to an ontology class via <code>rdfs:subClassOf</code> .
T21: Subpropertying	Mapping a source property to an ontology property via <code>rdfs:subPropertyOf</code> .
T22: Ontology alignment	Linking a class/property of the ontology to a class and/or property of another ontology via a specific triple.
T23: Internal linking	Linking unconnected data in the sources with one another via a triple.
T24: External linking	Linking resources to their representations in external datasets.
T25: Data format modification	Modifying the format of a piece of data in the sources to another one in the graph.
T26: Data model transformation	Mapping RDF data in one data model to a different one.

As the review is still ongoing the above list does not intend to be exhaustive. Once finished, the next step is to analyse and systematise them into a knowledge level model. This includes categorising them in terms of thematic groups, identifying which tasks can be considered concrete or abstract, or which specialise or are parts of others. Furthermore, as possible future work stemming from this research, it would be of interest to compare the results obtained with the tasks highlighted by known published methodologies and to find possible overlaps and differences.

7 Conclusion

In this work I have proceeded to expose my early stage (first year) PhD research on the yet open problem of theorising and applying the methods and tools of artificial intelligence to the problem of automating the designing KG generating pipelines.

To conclude, a brief discussion of some of the challenges that may be faced is required. A risk with what has been proposed is the possibly vast scale of the endeavour, the mitigation of which could lead to require some scope narrowing as we will begin to attempt to move from the knowledge level representation of our object to the building of its implementation. The early stage nature of this project also still leaves a number of open questions. A central one that may

affect the final results and related to the assumptions that have been taken is the extent to which automation can be achieved. That is, how many of the overall tasks to be performed can be safely allocated to an agent and how many to the user, either by design or by necessity, cannot be answered at this stage. Providing an answer to these questions will indeed play a part in the future development of this PhD project.

Acknowledgements

This research was supported by the UK EPSRC Doctoral Training Partnership 2022-2024 Open University and partly by the EU-funded project Polifonia: a digital harmoniser of musical cultural heritage (Grant Agreement N. 101004746), <https://polifonia-project.eu>.

Supervisors: Dr. Enrico Daga and Dr. Paul Mulholland.

References

1. Arenas-Guerrero, J., Chaves-Fraga, D., Toledo, J., Pérez, M.S., Corcho, O.: Morph-KGC: Scalable knowledge graph materialization with mapping partitions. *Semantic Web* **15**(1), 1–20 (2024). <https://doi.org/10.3233/SW-223135>
2. Asprino, L., Daga, E., Gangemi, A., Mulholland, P.: Knowledge graph construction with a façade: a unified method to access heterogeneous data sources on the web. *ACM Transactions on Internet Technology* (2022)
3. Buron, M., Goasdoué, F., Manolescu, I., Mugnier, M.L.: Obi-wan: ontology-based rdf integration of heterogeneous data. *Proceedings of the VLDB Endowment* **13**(12), 2933–2936 (2020)
4. Corcho, O., Priyatna, F., Chaves-Fraga, D.: Towards a new generation of ontology based data access. *Semantic Web* **11**(1), 153–160 (2020)
5. Daga, E., Asprino, L., Mulholland, P., Gangemi, A.: Facade-x: an opinionated approach to sparql anything. *Studies on the Semantic Web* **53**, 58–73 (2021)
6. Dimou, A., Chaves-Fraga, D.: Declarative description of knowledge graphs construction automation: Status & challenges. In: *Proceedings of the 3rd International Workshop on Knowledge Graph Construction (KGCW 2022) co-located with 19th Extended Semantic Web Conference (ESWC 2022)*. vol. 3141 (2022)
7. Dimou, A., Vander Sande, M., Colpaert, P., Verborgh, R., Mannens, E., Van de Walle, R.: Rml: A generic language for integrated rdf mappings of heterogeneous data. *Ldow* **1184** (2014)
8. Falbo, R.d.A., Guizzardi, G., Gangemi, A., Presutti, V.: Ontology patterns: clarifying concepts and terminology. In: *Proceedings of the 4th Workshop on Ontology and Semantic Web Patterns*. vol. 1188 (2013)
9. Fensel, D., Motta, E.: Structured development of problem solving methods. *IEEE Transactions on Knowledge and Data Engineering* **13**(6), 913–932 (2001)
10. Fernández-López, M., Gómez-Pérez, A., Juristo, N.: Methontology: from ontological art towards ontological engineering (1997)
11. Heckel, R.: Graph transformation in a nutshell. *Electronic notes in theoretical computer science* **148**(1), 187–198 (2006)

12. Hevner, A., Chatterjee, S., Hevner, A., Chatterjee, S.: Design science research in information systems. *Design research in information systems: theory and practice* pp. 9–22 (2010)
13. Hogan, A., Blomqvist, E., Cochez, M., d’Amato, C., Melo, G.d., Gutierrez, C., Kirrane, S., Gayo, J.E.L., Navigli, R., Neumaier, S., et al.: Knowledge graphs. *ACM Computing Surveys (CSUR)* **54**(4), 1–37 (2021)
14. Hur, A., Janjua, N., Ahmed, M.: A survey on state-of-the-art techniques for knowledge graphs construction and challenges ahead. In: *2021 IEEE Fourth International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*. pp. 99–103. IEEE (2021)
15. Jiménez-Ruiz, E., Grau, B.C., Zhou, Y., Horrocks, I.: Large-scale interactive ontology matching: Algorithms and implementation. In: *ECAI*. vol. 242, pp. 444–449 (2012)
16. Lefrançois, M., Zimmermann, A., Bakerally, N.: A sparql extension for generating rdf from heterogeneous formats. In: *The Semantic Web: 14th International Conference, ESWC 2017, Portorož, Slovenia, May 28–June 1, 2017, Proceedings, Part I 14*. pp. 35–50. Springer (2017)
17. Lenzerini, M.: Data integration: A theoretical perspective. In: *Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. pp. 233–246 (2002)
18. Lourdasamy, R., John, A.: A review on metrics for ontology evaluation. In: *2018 2nd International conference on inventive systems and control (ICISC)*. pp. 1415–1421. IEEE (2018)
19. Meyer, L.P., Stadler, C., Frey, J., Radtke, N., Junghanns, K., Meissner, R., Dziwis, G., Bulert, K., Martin, M.: Llm-assisted knowledge graph engineering: Experiments with chatgpt. *arXiv preprint arXiv:2307.06917* (2023)
20. Mitchell, M., Krakauer, D.C.: The debate over understanding in ai’s large language models. *Proceedings of the National Academy of Sciences* **120**(13), e2215907120 (2023)
21. Nagy, M., Vargas-Vera, M.: Towards an automatic semantic data integration: Multi-agent framework approach. *Semantic web* pp. 107–134 (2010)
22. Newell, A.: The knowledge level. *Artificial intelligence* **18**(1), 87–127 (1982)
23. Paulheim, H.: Knowledge graph refinement: A survey of approaches and evaluation methods. *Semantic web* **8**(3), 489–508 (2017)
24. Shvaiko, P., Euzenat, J.: Ontology matching: state of the art and future challenges. *IEEE Transactions on knowledge and data engineering* **25**(1), 158–176 (2011)
25. Sreemathy, J., Nisha, S., RM, G.P., et al.: Data integration in etl using talend. In: *2020 6th international conference on advanced computing and communication systems (ICACCS)*. pp. 1444–1448. IEEE (2020)
26. Vaithilingam, P., Zhang, T., Glassman, E.L.: Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In: *Chi conference on human factors in computing systems extended abstracts*. pp. 1–7 (2022)
27. Zhang, J., Shin, B., Choi, J.D., Ho, J.C.: Smat: An attention-based deep learning solution to the automation of schema matching. In: *Advances in Databases and Information Systems: 25th European Conference, ADBIS 2021, Tartu, Estonia, August 24–26, 2021, Proceedings 25*. pp. 260–274. Springer (2021)
28. Zou, X.: A survey on application of knowledge graph. In: *Journal of Physics: Conference Series*. vol. 1487, p. 012016. IOP Publishing (2020)