

RMLStreamer with Reference Conditions in the KGCW Challenge 2023

Els de Vleeschauwer¹, Gerald Haesendonck¹, Dylan Van Assche¹ and Ben De Meester¹

¹IDLab, Dept. Electronics & Information Systems, Ghent University – imec, Belgium

Abstract

The Knowledge Graph Construction Workshop (KGCW) 2023 Challenge aims to be a competitive challenge for knowledge graph construction systems to encourage optimizations not only for execution time but also for CPU and memory usage. We participated in this challenge with RMLStreamer, an RML mapping engine which processes all data in a streaming fashion. For the second part of the challenge, which is based on the Madrid-GTFS-Bench, we added RMLLooseGenerator as a first step. RMLLooseGenerator is a proof-of-concept implementation that simulates the effect of using reference conditions in RML mapping rules. In previous work we showed that using reference conditions in the GTFS-Madrid-Bench mapping file results in exactly the same graph output, while join operations are executed faster. The challenge results show that RMLStreamer scales well regarding execution time and CPU usage, while maintaining a constant memory usage. Therefore it received the Scalability Award in the KGCW 2023 Challenge. The challenge also highlighted some weaknesses of RMLStreamer: no support for relational databases, inefficient implementation of join operations, and longer execution time when handling nested sources such as JSON and XML files. After the challenge, the RMLStreamer has been expanded with support for relational databases. In the future, we will investigate how to optimize further the handling of joins and nested sources.

Keywords

RMLStreamer, challenge, knowledge graph construction

1. Introduction

Knowledge graph construction of heterogeneous data has seen a lot of uptake in the last decade from compliance to performance optimizations with respect to execution time [1, 2, 3]. Besides execution time as a metric for benchmarking knowledge graph construction systems, other metrics, e.g. CPU or memory usage, are often not considered. The Knowledge Graph Construction Workshop (KGCW) 2023 Challenge¹ aims to be a competitive challenge for knowledge graph construction systems to encourage optimiza-

KGCW'23: 4th International Workshop on Knowledge Graph Construction, May 28, 2023, Crete, GRE

✉ els.devleeschauwer@ugent.be (E. de Vleeschauwer); gerald.haesendonck@ugent.be (G. Haesendonck);

dylan.vanassche@ugent.be (D. Van Assche); ben.demeester@ugent.be (B. De Meester)

🌐 <https://dylanvanassche.be/> (D. Van Assche); <https://ben.de-meester.org/#me> (B. De Meester)

🆔 0000-0002-8630-3947 (E. de Vleeschauwer); 0000-0003-1605-3855 (G. Haesendonck);

0000-0002-7195-9935 (D. Van Assche); 0000-0003-0248-0987 (B. De Meester)

© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

¹<https://doi.org/10.5281/zenodo.7689310>

tions for execution time, but also CPU and memory usage. This challenge consists of two parts: (i) knowledge graph construction (KGC) parameters to evaluate individual parameters, e.g. joins and duplicates, with artificial data, and (ii) GTFS-Madrid-Bench [1] to focus on real-life use cases based on public transport data from Madrid.

In this paper, we present the results of the KGCW 2023 Challenge for RMLStreamer [4], an RML mapping engine which processes all data in a streaming fashion, in combination with RMLLooseGenerator², the proof-of-concept implementation simulating the effect of using reference conditions [5].

Section 2 describes the components of our knowledge graph construction pipeline. Section 3 discusses the setup we used to execute the challenge’s experiments. Section 4 explains our setups with other RML engines we compare with. We present our results in Section 5 and our conclusion in Section 6.

2. Knowledge Graph Construction Pipeline

Our knowledge graph construction pipeline consists of two parts: (i) RMLLooseGenerator emulates the effect of reference conditions [5], and (ii) RMLStreamer executes the RML mapping rules in a streaming fashion [4].

2.1. RMLLooseGenerator

RMLLooseGenerator is a proof-of-concept implementation for simulating the effect of reference conditions [5] in RML mapping rules. Reference conditions enable the reuse of another RML triples map’s subject map, without joining the logical sources. If a subject map referenced by an RML join has no references outside of those mentioned in the join conditions, RMLLooseGenerator interprets the join conditions concerned as reference conditions. It generates a new mapping file where those joins are replaced by appropriately adjusted object maps (crafted URIs), e.g. in the mapping of the GTFS-Madrid-Bench it replaces `rr:objectMap [rr:parentTriplesMap <#agency>; rr:joinCondition [rr:child "agency_id"; rr:parent "agency_id"]]` by `rr:objectMap [rr:template "http://transport.linkeddata.es/madrid/agency/agency_id"]`. The new mapping file can be processed by any RML engine.

RMLLooseGenerator is used only for the second part of the challenge (fig. 2), which is based on the GTFS-Madrid-Bench. RMLStreamer does not eliminate self-joins or duplicates. It needs more than three hours to execute the first scale of the GTFS-Madrid-Bench, generating an output of 105 GB. However, when re-interpreting all join conditions of the GTFS-Madrid-Bench mapping as reference conditions RMLStreamer needs only 33 seconds, generating an output of 76 MB, while the resulting knowledge graph remains semantically identical [5].

²<https://github.com/RMLio/rml-loose-generator>

2.2. RMLStreamer

RMLStreamer executes RML mapping rules to generate high quality Linked Data from multiple originally (semi-)structured data sources in a streaming way. RMLStreamer processes all data in a streaming fashion. It handles big input files and continuous data streams like sensor data without consuming more memory when the input data size increases. RMLStreamer leverages Apache Flink to scale vertically across multiple CPU cores and horizontally across multiple machines. In the challenge we use RMLStreamer version 2.4.2 with an embedded Flink version in a Docker container³.

3. Experiment setup

The KGCW 2023 Challenge provides CSV files as source datasets, mapping files, queries, baseline results (i.e. the expected set of triples and query results), an example pipeline based on the RMLMapper, MySQL, and Virtuoso for reaching those results, and a tool for executing the example pipeline.

We made the following adaptations to the provided experiments to enable execution with RMLStreamer. (i) In the provided end-to-end pipelines the CSV files are loaded into a relational database. As the RMLStreamer did not support SQL yet when the execution of the challenge was performed, we used the CSV files directly to construct the knowledge graphs. (ii) We replaced JSON iterator `[*]` with `[$[*]` in the mapping files. (iii) We added a condition to the mapping files to recognize the string `NULL` in the provided CSV files as an empty value.

We compared our experiments' results to ensure that our output is correct with respect to the baseline results of the challenge. For the first part of the challenge (KGC parameters), where the output of RMLStreamer is not loaded into a triples store, we deduplicated the output results as RMLStreamer cannot eliminate duplicates by itself. After deduplication we compared the number of triples to the baseline results of the challenge. For the second part of the challenge (GTFS-Madrid-Bench) we compared the number of query results to the baseline.

Figure 1 and Figure 2 visualize the resulting experiment setups.

We executed all experiments on a Intel Xeon CPU E5645 (12 cores with Hyper-Threading, 2.4GHz) with 24GB RAM and 250GB HDD. The challenge execution tool configures the Java heap space to 50 % of the available memory. All experiments were performed 5 times and the experiment with the median execution time is reported. All files needed to reproduce the conducted experiments are available on Zenodo⁴.

³<https://doi.org/10.5281/zenodo.7181800>

⁴<https://doi.org/10.5281/zenodo.8034245>

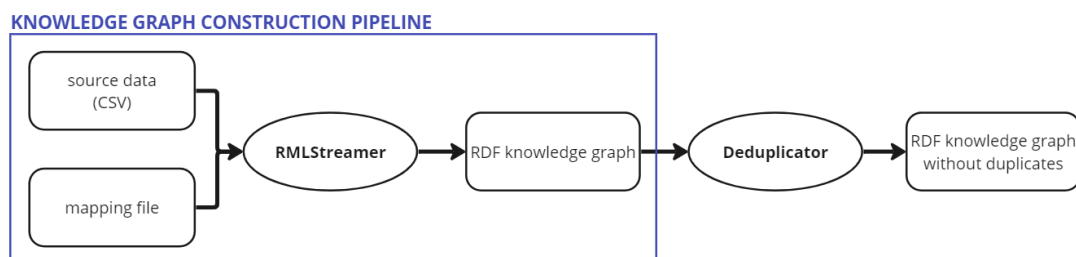


Figure 1: Setup for the first part of the KGCW 2023 Challenge (KGC parameters): a deduplication step is added to make the end results comparable to the baseline results of the challenge.

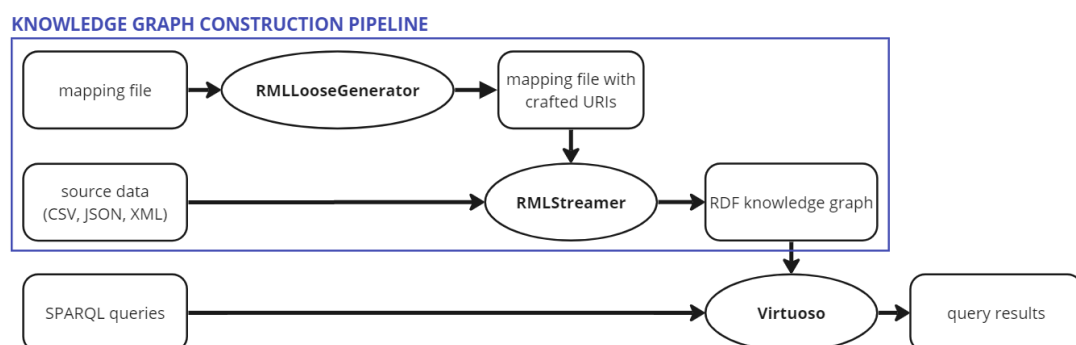


Figure 2: Setup for second part of the KGCW 2023 Challenge (GTFS-Madrid-Bench): joins are replaced by crafted URIs to reduce the knowledge graph construction execution time and the size of the resulting RDF knowledge graph.

4. Comparison with other RML engines

We also executed the experiments in our setup (section 3) with RMLMapper⁵ [6] for both parts of the challenge, and with Morph-KGC⁶ [2] for the GTFS-Madrid-Benchmark part. This way, we can compare the results of RMLStreamer on the same setup with other implementations.

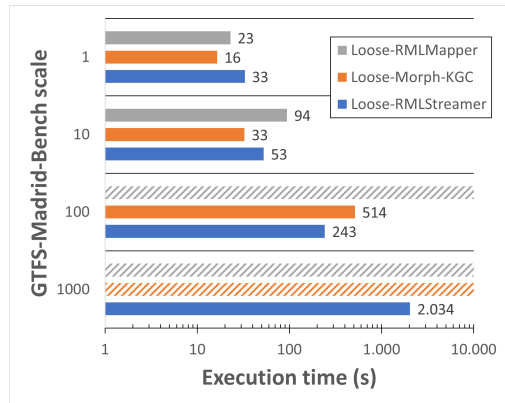
5. Results

In Figure 3 and Table 1 we included the measured execution time, CPU and memory usage of the knowledge graph construction pipeline (RMLLooseGenerator and RMLStreamer) for selected experiments. The complete overview of the results, as it was submitted to the KGCW 2023 Challenge, is available on Zenodo⁷.

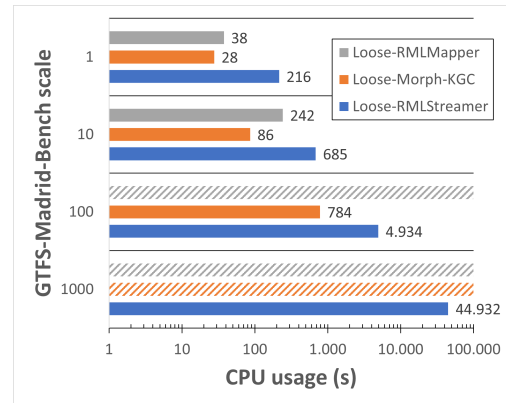
⁵<https://github.com/RMLio/rmlmapper-java>

⁶<https://doi.org/10.5281/zenodo.5543552>

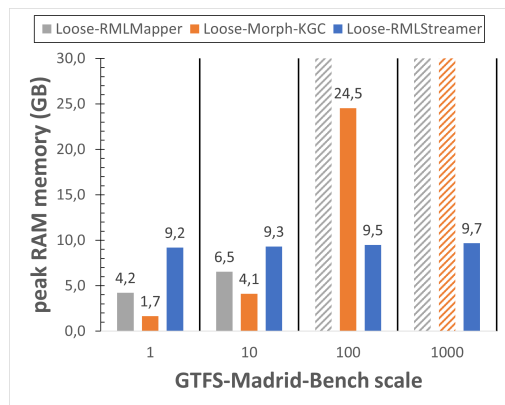
⁷<https://doi.org/10.5281/zenodo.8034245>



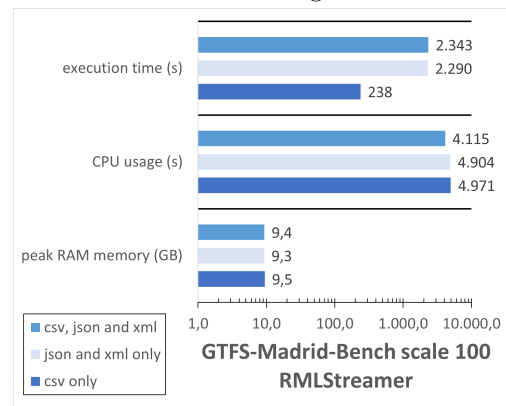
(a) Linear trend: the execution time of RMLStreamer increases with the same factor as the data size for higher scales.



(b) Linear trend: the CPU usage of RMLStreamer increases with the same factor as the data size for higher scales.



(c) RMLStreamer has a constant memory usage independent of data size



(d) The execution time of RMLStreamer increases with a factor 10 when including json and xml sources

Figure 3: Metrics of the knowledge graph construction pipeline for the GTFS-Madrid-Bench experiments. A striped bar signifies that this test could not be completed by the engine at hand.

First, we verify how RMLStreamer behaves when the size of the input data increases. This is best illustrated by the GTFS-Madrid-Bench scale experiments. Figure 3a and Figure 3b show that RMLStreamer scales towards a linear trend. Note that the reported metrics include the execution time of RMLLooseGenerator (average execution time of 8 seconds) and the startup process of RMLStreamer (no separate measurements in the challenge execution tool, but from the logs of RMLStreamer we conclude that this is on average 17 seconds for the GTFS-Madrid-Bench cases). These processes are independent of data size. The execution time increases with the same factor as the data size for higher scales, where RMLMapper and Morph-KGC cannot complete the experiments for respectively scale 100 and scale 1000. At scale 100 RMLStreamer is two times faster than the state-of-the-art RML engine Morph-KGC (fig. 3a). RMLStreamer uses more CPU

(fig. 3b) in comparison to RMLMapper and Morph-KGC because it maximizes parallelism over all given slots. The CPU usage also scales linear with the size of the input data for higher scales. The peak RAM memory (fig. 3c) measured is similar for all scales when using RMLStreamer, where RMLMapper and Morph-KGC hit the limits of the available memory and fail to complete all GTFS-Madrid-Bench experiments. RMLStreamer has a constant memory usage independent of the data size. We assume this is mostly due to the fact that it processes everything in a streaming way, which is to a lesser extent the case for RMLMapper and Morph-KGC.

The measurements for the KGC parameter experiments show similar trends: linear scaling of execution time and CPU usage, proportional to the size of the input data, in combination with a constant memory usage. Table 1 shows the results of the experiments with the lowest and the highest measured values, which are dubbed easiest and hardest experiments in the table respectively.

Second, we evaluate the impact of the format of the data input. Adding nested sources, such as JSON and XML files, increases the execution time of RMLStreamer with a factor ten (fig. 3d). The difference in execution time is the consequence of RMLStreamer chunking CSV files and processing the chunks in parallel. This is not the case for the XML and JSON formats yet.

	Execution time (s)	CPU (s)	Peak RAM (GB)	Output (triples)
1. Easiest experiment: records 10K rows 20 columns				
RMLMapper	8	15	1,4	200.000
RMLStreamer	21	107	2,4	200.000
2. Hardest experiment for RMLMapper: properties 1M rows 30 columns				
RMLMapper	262	646	13,2	20.000.000
RMLStreamer	120	2.247	9,2	20.000.000
3. Hardest experiment for RMLStreamer: records 10M rows 20 columns				
RMLMapper	out of memory	out of memory	out of memory	out of memory
RMLStreamer	653	14.360	9,2	200.000.000
4. Easiest experiment with joins: join 1-1 0%				
RMLMapper	out of memory	out of memory	out of memory	out of memory
RMLStreamer8	219	2409	9,3	0
RMLStreamer	43	371	9,4	0
5. Hardest experiment with joins: join 5-5 100%				
RMLMapper	out of memory	out of memory	out of memory	out of memory
RMLStreamer8	434	4145	9,3	2.500.000
RMLStreamer	66	1093	9,5	2.500.000

Table 1

Metrics of the knowledge graph construction step for selected KGC parameter experiments

Last, we comment on the experiments including joins. At the time of the challenge we limited RMLStreamer to eight task slots (referenced as RMLStreamer8 in Table 1) for the execution of experiments including joins, because RMLStreamer reported an error and failed to start processing some mapping files including joins (e.g. the original GTFS-Bench-Mapping with joins). We assumed that this error appeared with any mapping file that includes joins. Further investigation afterwards revealed that this error

got triggered by mappings with a large number of mapping rules (i.e. a mapping with two triples maps and one join operation does not result in error). Hence, the limitation of task slots is not required for the experiments with joins in the first part of the challenge. For completeness we added the results of those join experiments with all task slots. Using all 24 task slots on the hardest experiment with joins decreases the execution time by a factor of 16, and the CPU usage by a factor of nine, compared to the results registered during the challenge using RMLStreamer8. The decreased execution time was in line with our expectations, however, the reason for the decreased CPU usage requires further investigation.

6. Conclusion

The KGCW 2023 Challenge results show that RMLStreamer has a linear scaling of execution time and CPU usage, proportional to the size of the input data, while maintaining a constant memory usage. Therefore it received the Scalability Award in the KGCW 2023 Challenge.

We noted the following improvement areas for RMLStreamer: (i) a lack of support for relational databases, (ii) an inefficient implementation of join operations (e.g. GTFS-Madrid-Bench experiments with joins cannot be handled properly by RMLStreamer), and (iii) a longer execution time when handling nested sources such as JSON and XML files.

Additionally we noticed that there are cases where the number of task slots used by RMLStreamer needs to be limited. At the time of the challenge this could only be achieved by modifying RMLStreamer's code.

Support for changing the number of task slots dynamically was implemented in RMLStreamer 2.5.0⁸, together with support for SQL databases.

For future work we will investigate further optimizations for execution of joins and nested sources.

Acknowledgments

The described research activities were supported by SolidLab Vlaanderen (Flemish Government, EWI and RRF project VV023/10), the imec ICON project AI4Foodlogistics (Agentschap Innoveren en Ondernemen project nr. HBC.2020.3097), and funded by the Special Research Fund of Ghent University under grant BOF20/DOC/132.

References

- [1] D. Chaves-Fraga, F. Priyatna, A. Cimmino, J. Toledo, E. Ruckhaus, O. Corcho, Gtfs-madrid-bench: A benchmark for virtual knowledge graph access in the transport

⁸<https://doi.org/10.5281/zenodo.7998156>

- domain, *Journal of Web Semantics* 65 (2020) 100596. doi:10.1016/j.websem.2020.100596.
- [2] J. Arenas-Guerrero, D. Chaves-Fraga, J. Toledo, M. S. Pérez, O. Corcho, Morph-KGC: Scalable knowledge graph materialization with mapping partitions, *Semantic Web* (2022) 1–20. doi:10.3233/sw-223135.
- [3] E. Iglesias, S. Jozashoori, D. Chaves-Fraga, D. Collarana, M.-E. Vidal, SDM-RDFizer: An RML Interpreter for the Efficient Creation of RDF Knowledge Graphs, in: *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020. doi:10.1145/3340531.3412881.
- [4] Sitt Min Oo, G. Haesendonck, B. De Meester, A. Dimou, RMLStreamer-SISO: An RDF Stream Generator from Streaming Heterogeneous Data, in: U. Sattler, A. Hogan, M. Keet, V. Presutti, J. P. A. Almeida, H. Takeda, P. Monnin, G. Pirrò, C. d’Amato (Eds.), *The Semantic Web – ISWC 2022*, Springer, Springer International Publishing, Cham, 2022, pp. 697–713. doi:10.1007/978-3-031-19433-7_40.
- [5] E. de Vleeschauwer, S. Min Oo, B. De Meester, P. Colpaert, Reference conditions: Relating mapping rules without joining, in: *Proceedings of the 4rd International Workshop on Knowledge Graph Construction (KGCW 2023) co-located with 20th Extended Semantic Web Conference (ESWC 2023)*, 2023.
- [6] A. Dimou, M. Vander Sande, P. Colpaert, R. Verborgh, E. Mannens, R. Van de Walle, RML: A Generic Language for Integrated RDF Mappings of Heterogeneous Data, in: *Proceedings of the 7th Workshop on Linked Data on the Web*, volume 1184, 2014.