

# Popularity Driven Data Integration

Fausto Giunchiglia<sup>[0000-0002-5903-6150]1</sup>, Simone Bocca<sup>[0000-0002-5951-4589]1</sup>,  
Mattia Fumagalli<sup>[0000-0003-3385-4769]2</sup>, Mayukh Bagchi<sup>[0000-0002-2946-5018]1</sup>,  
and Alessio Zamboni<sup>[0000-0002-4435-1748]1\*</sup>

<sup>1</sup> Department of Information Engineering and Computer Science (DISI),  
University of Trento, Italy

{fausto.giunchiglia,simone.bocca,mayukh.bagchi,alessio.zamboni}@unitn.it

<sup>2</sup> Conceptual and Cognitive Modeling Research Group (CORE),  
Free University of Bozen-Bolzano, Bolzano, Italy  
{mattia.fumagalli@unibz.it}

**Abstract.** More and more, with the growing focus on large scale analytics, we are confronted with the need of integrating data from multiple sources. The problem is that these data are impossible to reuse *as-is*. The net result is high cost, with the further drawback that the resulting integrated data will again be hardly reusable *as-is*. *iTelos*<sup>3</sup> is a general purpose methodology aiming at minimizing the effects of this process. The intuition is that data will be treated differently based on their *popularity*: the more a certain set of data have been reused, the more they will be reused and the less they will be changed across reuses, thus decreasing the overall data preprocessing costs, while increasing backward compatibility and future sharing.

**Keywords:** Data Reuse · Data Sharing · Knowledge Graphs.

## 1 Introduction

More and more, with the growing focus on large scale analytics, we are confronted with the need of integrating data from multiple sources. The key issue is how to handle the *semantic heterogeneity* which is intrinsic in such data. Two main approaches have been proposed. The first approach consists of using *ontologies*, where the goal is to agree on a fixed language and/or schema towards facilitating future sharing [3]. The second consists of exploiting the flexibility of *Knowledge Graphs (KGs)* [9], as the means for facilitating the adaptation and integration of heterogeneous data. However the problem is far from being solved. No matter the approach, it is just impossible to reuse data *as-is*. The net result is usually a lot of data preprocessing (e.g., cleaning, normalization) which results in high cost, with the further drawback that the resulting integrated data will again be hardly reusable *as-is*. It is a negative loop which consistently reinforces itself.

\* The research by F. Giunchiglia, M. Bagchi and S. Bocca has received funding from the “*DELPhi - Discovering Life Patterns*” project funded by the MIUR (PRIN) 2017. The research by A. Zamboni was supported by the *InteropEHRRate* project, EC Horizon 2020 programme under grant number 826106.

<sup>3</sup> Not to be confused with *Telos* [11].

In this paper, we propose *iTelos*, a general purpose methodology whose main goal is to minimize the high costs of this loop of reuse. *iTelos* is crucially based on the use of KGs and ontologies, i.e., *reference schemas*, as the schemas of KGs. However, its novel underlying intuition is to treat data differently, depending on their *popularity*. In particular, the idea is to select first and minimize changes on those data which are more (re)-used thus decreasing the preprocessing costs, while increasing backward compatibility and future sharing. To this extent, *iTelos* distinguishes among three categories of data. That is: *Common*, which are used across domains (e.g., data about space, time, transportation), *Core*, which, while being more vertical than common, are extensively used in the domain under consideration (e.g., in tourism, all the data that can be found in Open Data portals), and *Contextual*, namely, data specific to the application at hand (e.g., the data extracted on purpose from legacy systems). Popularity also drives the selection of the reference schemas, based on the intuition that, also at the schema level, more reuse is a good motivation for further reuse. Thus we have *Common* reference schemas (e.g., standards about space, time, transportation, *schema.org*<sup>4</sup>), *Core* reference schemas, (e.g., in Health, FHIR<sup>5</sup>), and *Contextual* reference schemas (e.g., as they apply the current application).

*iTelos* implements the above idea based on a precisely articulated data integration process, based on three key assumptions, as follows:

- *data* and *reference schemas* should be integrated under the overall guidance of the needs to be satisfied, formalized as *competence queries (CQs)* [8];
- The requirements, including those on data and reference schema reuse, as well as CQs, should be known a priori as part of an application *purpose*;
- A difficulty is that the schemas of the data to be reused usually do not map to reference schemas, the mapping being usually arbitrarily complex. The idea is to build the integrated KG via a sequence of *middle-out* iterations where, first, CQs are used to drive the selection and preprocessing of the data, largely independently from the reference schemas, and where, in a second step, reference schemas, suitably and independently integrated among them, are adapted to fit best the integrated data minimizing the negative effects on sharability.

This paper is organized as follows. In Section 2 we describe the *iTelos* process. In Section 3 we describe how *iTelos* enhances the reusability, with minimal changes, of the available data. In Section 4 we describe how *iTelos* enhances the future sharability of the resulting KG. Finally, Section 5 syntetically describes the case studies to which *iTelos* has been applied.

## 2 The Process

The *iTelos* process is depicted in Figure 1. The *User* provides in input the specification of the problem, the *Purpose*, and receives in output an integrated KG, i.e., the *Entity Graph*. The purpose contains three main elements, as follows:

<sup>4</sup> <https://schema.org/>

<sup>5</sup> <http://hl7.org/fhir/>

- the functional requirements of the KG to be generated, that we assume to be ultimately formalized as CQs;
- The *datasets* to be reused. We assume that these datasets consist of *Entity Graphs (EGs)*, namely graphs where nodes are *entities* (e.g., my cat *Garfield*), decorated with data property values and linked among them via object property links;<sup>6</sup>
- The *Ontologies* to be reused. We assume that these ontologies, i.e., reference schemas, consist of *entity type (etype) Graphs (ETGs)*, namely KGs which define the schema of EGs. In ETGs nodes are *etypes*, namely classes of entities (e.g., the class *cat*), decorated by the data and object properties which define the EG structure.<sup>7</sup>

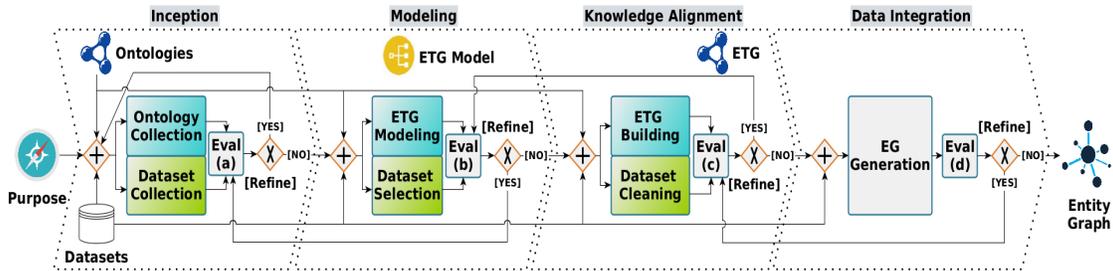


Fig. 1: The *iTelos* Process.

The overall *iTelos* process is articulated in four phases. The *Inception* phase takes as input the purpose and collects the ontologies and datasets needed to build the target EG. During this phase, the functional requirements are encoded into a list of CQs which are then matched to the input datasets and ontologies thus implementing the activities of *Ontology Collection* and *Dataset Collection*, as from Figure 1. The main objective of the *Modeling* phase is to build the most suitable model of the ETG to be used as the schema of the final EG, which in Figure 1 is called the *ETG model*. In practice, the ETG model includes all the etypes and properties needed to represent the information required by each CQ, possibly extended by extra etypes and properties suggested by the datasets. The *Dataset Selection* activity finalizes the selection of the datasets selecting whenever possible, the ones which are most popular. The goal of *Knowledge Alignment* is to enhance the sharability of the final EG, building in turn a shareable ETG that fits in the best possible way the datasets to be integrated. The input ETG model is itself a possible solution. However the set of reference schemas provides more possibilities, in terms of etypes and properties available, that can be adopted to implement a final ETG (called ETG in Figure 1) easier to share

<sup>6</sup> Many Open data portals provide such datasets, at different levels of formalization in the 5STAR Open Data Model.

<sup>7</sup> Many such schema repositories are already available; see for instance: LOV (<https://lov.linkeddata.es/>), LOV4IoT (<http://lov4iot.appspot.com/>), DATAHUB (<https://old.datahub.io/>) and *LiveSchema* (<http://liveschema.eu/>)

and reuse. Here again preference is given to those reference schemas which are most popular. The *Dataset Cleaning* activity performs the final cleaning of the datasets consistently with the ETG, trying to minimizing it and concentrating the preprocessing on those datasets which are less popular. As described in detail in Section 4, this may require building an ETG which is an adaptation of the input reference schemas, e.g., in the selected etypes, properties, data types and formats.

The last phase is *Data Integration*. The objective is to build the EG, what we call *EG Generation*, integrating the ETG with the data resources. To do that, the ETG and datasets are provided in input to a specific data mapping tool, called *KarmaLinker*, which consists of the *Karma* data integration tool [10] extended to perform Natural Language Processing on short sentences (i.e., what we usually call *the language of data*) [1]. The process is described in some detail in [7]. The first activity in this phase maps the data to the etypes and properties of the ETG. The following step is the generation of the entities that are then matched and, whenever they are discovered to be different representations of the same real world entity, merged. These activities are fully supported by *Karmalinker*. The above process is iteratively executed over the list of selected datasets. The process concludes with the export of the EG into an RDF file.

The key observation is that the *iTelos* is implemented as two separate sub-processes, executed in parallel within each phase, one operating on reference schemas, the other on the input data (blue and green boxes in Figure 1). During this process, the initial purpose keeps evolving building the bridge between CQs, datasets and reference schemas. To enforce the convergence of this process, and also to avoid making costly mistakes, each phase ends with an evaluation activity (*Eval* boxes in Figure 1). The details of how the evaluation is performed is out of the main scope of the paper. Here it is worth noticing that, within each phase, the evaluation aims to verify that the target of that phase is met, namely: aligning CQs with datasets and ontologies in phase 1, thus maximizing reusability; aligning the ETG model with the datasets in phase 2, thus guaranteeing the success of the project; and aligning ETG and ontologies in phase 3, thus maximizing sharability. The evaluation in phase 4 has the goal of checking that the final EG satisfies the requirements specified by the purpose.

### 3 Enhancing data reuse

Reusability is enhanced during the phases of inception and modeling, whose main goal is to progressively transform the specifications from the purpose into the ETG Model. This process happens according to the following steps:

1. generation of a list of natural language sentences, each informally defining a CQ, as implicitly or explicitly implied by the purpose;
2. generation of a list of relevant etypes and corresponding properties, which formalize the informal content of CQs, as from the previous step;
3. selection of the datasets whose schema informally matches the CQs, as from the previous step;

4. generation of a list of etypes with associated properties, from the selected datasets, which match the etypes and properties from the CQs;
5. construction of the ETG model;

Steps 1-4 happen during inception, while step 5 happens during the modeling phase. The key observation is that, starting from an analysis of the etypes and properties inside CQs, the two types of resources involved (i.e., ontologies, datasets) are handled through a series of three iterative executions, each corresponding to a specific category, following a decreasing level of reusability. The categories are defined as follows:

*Common*: this category involves resources associated with aspects that are common to all domains, also outside the domain of interest. Usually, these resources correspond to abstract etypes specified in *upper level ontologies* [4], e.g., *person*, *organization*, *event*, *location*, and/or to etypes from very common domains, usually needed in most applications, e.g., *Space* and *Time*. The data that are found in Open Data sites as well the ontologies which can be found in the repositories mentioned above are examples of common resources.

*Core*: this category involves resources associated with the more core aspects of the domain under consideration. They carry information about the most important aspects considered by the purpose, information without which it would be impossible to develop the EG. Consider for instance the following purpose:

*“There is a need to monitor Covid-19 data in the Trentino region (Italy), to understand the diffusion of the virus and the social restriction caused by the virus, with the possibility to identify new outbreaks”.*<sup>8</sup>

In this example, core resources could be those data values reporting the number of Covid-19 infections in the specified region. Examples of common resources are the data of certain domains, e.g., public sector facilities (e.g., hospitals, transportation, education), domain specific ontologies that can be found again in the repositories above, as well as domain specific standards (e.g., Health, interoperability standards of various types). In general, data are harder to find than ontologies, in particular when they are about the private sector, where they carry economic value and are often collected by private bodies.

*Contextual*: this last category involves resources that carry specific, possibly unique, information from the domain of interest. These are the resources whose main goal is to create added value. If core resources are necessary for a meaningful application, contextual resources are the ones which can make the difference with respect to the competitors. In the above example, examples of contextual resources can be those data describing the type of social restrictions adopted to contrast the virus. At the schema level, contextual etypes and properties are

---

<sup>8</sup> This is a small example extracted from the project which built a KG, following the *iTelos* methodology, on *“Integration of medical data on Covid-19”* developed by Antony, N., Gotca, D., Jyate, M., Donini, L. The complete material and description can be found at the URL <https://github.com/UNITN-KDI-2020/COVID-data-integration>.

those which differentiate the ontologies which, while covering the same domain, actually present major differences. Data level contextual resources are usually not trivial to find, given their specificity and intrinsic. In various applications we have developed in the past, this type of data have turned out to be a new set of resources those had to be generated on purpose for the application under development, in some cases while in production.

The overall conclusive observation is that the availability of resources, and of data in particular, decreases from common to contextual. On top of this and because of this, as part of the *iTelos* strategy, a decreasing effort is made, moving from common to core to contextual data, in maintaining high the level of reusability and sharability, thus concentrating the preprocessing costs in the latter categories.

## 4 Enhancing data sharability

The knowledge alignment phase aims to enhance sharability, by aligning and possibly modifying the ETG model to take into account the etypes and properties coming from the reference ontologies. The key observation is that the alignment mainly concerns the common and, possibly, the core types with much smaller expectations on contextual etypes. Notice how the alignment with the most suitable ontology will enable the reuse of data, at least for what concerns common and sometimes core etypes. As an example, the selection of Google GTFS or FOAF as reference ontologies ensures the availability of a huge amount of data, a lot of which are open data. This type of decision should be made during inception; if delayed up to here, it might generate backtracking.

We align the ETG model with the reference ontologies, by adapting the *Entity Type Recognition (ETR)* process proposed in [6]. This process happens in three steps as follows:

*Step 1: ontologies selection.* This step aims at selecting the set of reference ontologies that best fit the ETG model. As from above the first step is to rank the ontologies, as selected during the Inception phase, based on their popularity. Then, moving from top to down in the list, as from [6], this selection step occurs by measuring each reference ontology according to three metrics, which allow:

- to identify how many etypes of the reference ontologies are in common with those defined in the ETG model, and
- to measure a property sharability value for each ontology etype, indicating how many properties are shared with the ETG model etypes.

The output of this first step is a set of selected ontologies, which best cover the ETG model, and that have been verified fitting the dataset’s schema, at both etypes and etype properties levels.

*Step 2: Entity Type Recognition(ETR).* The main goal here is to predict, for each etype of the ETG model, which etype of the input ontologies, analyzed one at the time, best fits the ETG. In practice, the ETG model’s etypes are used as labels of the classification task and, as mentioned in [6], the execution exploits

techniques that are very similar to those used in ontology matching (see, e.g., [5]). The final result is a vector of prediction values, returning a similarity score between the ETG model's etypes and the selected ontology etypes.

*Step 3: ETG generation.* This step identifies, by using the prediction vector produced in the previous step, those etypes and properties from the ontologies which will be added to the final version of the ETG. Notice how this must be done while preserving the mapping with the datasets' schemas. The distinction among common, core and contextual etypes and properties plays an important role in this phase and can be articulated as follows:

- The common etypes should be adopted from the reference ontology, in percentage as close as possible to 100%. This usually results in an enrichment of the top level of the ETG model by adding those top level etypes (e.g., *thing*, *product*, *event*, *location*) that usually no developer considers, because too abstract, but which are fundamental for building a highly shareable ETG where all properties are positioned in the right place. This also allows for an alignment of those *common isolates* (see Section 3) for which usually a lot of (open) data are publicly available (e.g., *street*);
- The core etypes are tentatively treated in the same way as common etypes, but the results highly depend on the ontologies available. Think for instance of the GTFS example above;
- Contextual etypes and, in particular, contextual properties are mainly used to select among ontologies, the reason being that, they allow distinguishing the most suitable among a set of ontologies about the same domain [6].

## 5 Case studies

The specification of *iTelos* is in its early phases, in particular in terms of tool support. However, a lot of work has been dedicated to the refinement of the single steps of the overall approach and on their extensive evaluation. In particular, *iTelos* has been validated during the past four Academic Years as part of the Knowledge and Data Integration (KDI) class, a six credit course of the Master Degree in Computer Science of the University of Trento.<sup>9</sup> During this class, 2-5 students per group, must generate an EG using the pipeline above starting from a high level problem specification. The overall project has an elapsed time of fourteen weeks during which students have to work intensely. We estimate the overall effort each group puts into building an EG in around 4-8 person-months, depending on the case.

As of today we have piloted around 30 projects and 90 evaluations of the *iTelos* methodology as a whole. The details of this work cannot be reported here for lack of space. However the results, restricted to the first three years

---

<sup>9</sup> <http://knowdive.disi.unitn.it/teaching/kdi/> contains the material used during the last two editions of the course. This material consists of theoretical and practical lectures, as well as demos of the tools to be used, some of which have been mentioned above.

are described in some detail in [2]. We report below the most relevant answers provided by the students, which were asked to fill a very detailed questionnaire containing a set of qualitative questions about the methodology.

- (*Strength*) the step by step, precisely articulated, *iTelos* process is easy to follow;
- (*Strength*) the stepwise iterative evaluation process supports well the refinement of the Entity Graph;
- (*Weakness*) A wrong decision made in the early phases is quite difficult to remedy, with this possibility being very high during the inception phase;
- (*Weakness*) The work between the schema and the data layer is unbalanced in favour of the second, in particular during the *informal modeling* phase. This complicates the synchronization of the work with the possibility of misalignments, mainly because of misunderstandings, that have to be handled very carefully by the project manager.

## 6 Conclusions

In this paper we have introduced *iTelos*, a novel methodology whose ultimate goal is to implement a *circular* development process. By this we mean that the goal of *iTelos* is to enable the development of EGs via the *reuse* of already existing EGs and ETGs, while being simultaneously developed to be later easily *reused* by other applications to come.

## References

1. Bella, G., Gremes, L., Giunchiglia, F.: Exploring the language of data. In: Proc. 28th Int. Conf. on Computational Linguistics. pp. 6638–6648 (2020)
2. Bocca, S., Dragoni, M., Giunchiglia, F.: *iTelos* - case studies in building domain specific knowledge graphs. In: Int. Semantic Intelligence Conference (2022)
3. Ekaputra, F., et al.: Ontology-based data integration in multi-disciplinary engineering environments: A review. *Open Journal of Inf. Systems* **4**(1), 1–26 (2017)
4. Gangemi, A., Guarino, N., Masolo, C., Oltramari, A., Schneider, L.: Sweetening ontologies with dolce. In: International Conference on Knowledge Engineering and Knowledge Management. pp. 166–181. Springer (2002)
5. Giunchiglia, F., Autayeu, A., Pane, J.: S-match: an open source framework for matching lightweight ontologies. *Semantic Web* **3**(3), 307–317 (2012)
6. Giunchiglia, F., Fumagalli, M.: Entity type recognition – dealing with the diversity of knowledge. In: Knowledge Representation Conference (KR) (2020)
7. Giunchiglia, F., Zamboni, A., Bagchi, M., Bocca, S.: Stratified data integration. In: 2nd Int. Wshop On Knowledge Graph Construction (KGCW) (2021)
8. Grüninger, M., Fox, M.S.: The role of competency questions in enterprise engineering. In: Benchmarking—Theory and practice, pp. 22–31. Springer (1995)
9. Kejriwal, M.: Domain-specific knowledge graph construction. Springer (2019)
10. Knoblock, C.A., Szekely, P.: Exploiting semantics for big data integration. *AI Magazine* **36**(1), 25–38 (2015)
11. Mylopoulos, J., Borgida, A., Jarke, M., Koubarakis, M.: *Telos: Representing knowledge about information systems*. ACM Trans. on Information Systems (TOIS) (1990)