

**European Joint Doctorate in
Data Engineering for Data Science (DEDS)**

Doctoral Project Plan¹

Thesis Title

First Name Last Name

1 Project Summary

The term Big Data refers to data sets that are too large or complex to be dealt with by traditional data processing software. In particular, Big Data captures to 4 dimensions: Velocity, Variety, Veracity, and Volume. This large quantity of data available nowadays has a lot of statistical and business value, therefore their analysis is of core importance for business decision-making.

Data Integration is the set of processes to gather and bridge data from heterogeneous sources together in order to have a unified view. The premise of data integration is to make data more freely available and easier to consume and process by systems and users. Research on Data Integration started more than 50 years ago [10, 18] but as we entered the Big Data era, new challenges arose, which include scaling [7] data integration while guaranteeing the privacy [27] of the individuals involved in the datasets. Over the last decade new techniques have been applied for improving computational performance, consisting of the use of parallelizing the computation by using big data processing platforms [7], or algorithmically, namely using summarization techniques [4], used for approximate fast and approximate querying, improving the performance of Machine Learning processes [12, 1, 15] as well as in some data integration scenarios.

Analyzing Big Data may involve accessing sensitive information, especially in the health sector. For these reasons, in Data Integration, techniques like Secure Multi-Party Computation and Differential Privacy [8] have been used. One big concern in private Data Integration is that it is difficult to guarantee privacy, accuracy, and good performance at the same time.

Machine Learning and Data Integration have really close relationship [6]. In particular, it is possible to leverage the first to improve the performance of the second and vice-versa. An evolving branch of Machine Learning is Federated Learning, which consists in building a Machine Learning model in a federated setting (when the data is distributed across edge devices) and model in a collaborative way, without moving the data to a central server.

This Ph.D. aims to explore summarization techniques and Federated Learning for scaling Data Integration tasks while guaranteeing privacy and achieving good performance.

2 Scientific Content of the Doctorate Project

2.1 Background

2.1.1 Data Integration

Data Integration (DI) is the practice of consolidating data from disparate sources into a unified view. It has been studied since the birth relational databases. It is characterized by three main

¹Choose the appropriate heading. Based on the PhD Study Plan of Aalborg University, available at <http://www.phd.teknat.aau.dk/intranet/phd-study-plan>.

steps:

1. **Schema alignment**, a process that takes as input a set of different schemas on the same domain and outputs a *mediated schema*, an *attribute matching* and a *schema mapping*.
2. **Record linkage**, also referred as entity resolution, computes a partitioning of the set of records from different datasets, such that each partition identifies the records that refer to a distinct entity.
3. **Data fusion** aim is to identify which are the best records to represent a specific entity, when a source provide conflicting values.

2.1.2 Federated Learning

Federated Learning is a machine learning approach where a model is trained across multiple decentralized edge-devices, in an edge computing fashion. Each device trains a local model and then, either in a centralized, decentralized or heterogeneous approach, build a global model. It differs from a distributed machine learning as the data is not expected to be identically distributed. Besides the advantage of having a distributed computation, guaranteeing more efficiency, it gained a lot of popularity due to the fact that data is not exchanged between the parts involved, thus guaranteeing privacy.

It has gained a lot of popularity both in research and industry, in particular in transportation [9], Industry 4.0 [24] and digital health [20].

2.1.3 Differential Privacy & Synopses for Big Data

Differential Privacy (DP) is a technique for sharing datasets' information without compromising the privacy of the individuals. The idea is to add noise to the data such that the new distribution is close to the real one, but not equal.

Synopses or summaries are a set of technique and probabilistic data structures to compute compact description of big datasets. They gained a lot of popularity over the last decade due to the rise of Big Data.

2.2 State of the Art

2.2.1 Privacy-aware Data Integration in the Big Data Era

Privacy in the context of data management has gained a lot of popularity over the last decade, as public awareness about issues in management sensitive data increased. Due to this, privacy became of central importance in the field of Big Data Management, analytics and processing.

In the particular case of Data Integration, privacy-preserving techniques has been used extensively in literature, especially for record linkage. In particular, differentially-private record linkage and cryptography has been used extensively [16, 26, 3, 17, 11, 19, 2, 17]. As regards schema matching and data fusion, there are fewer works on guaranteeing privacy, most of the work is based on guaranteeing efficiency, by using both rule-based and learning-based approaches [22, 21, 23].

Schema matching. Original schema matching techniques can be categorized as it follows:

- **Schema-level matchers**, where only the metadata is considered (e.g. column labels, data-types).
- **Content-level matchers**, where the content of the columns is used for matching by using probabilistic approaches [25, 5].
- **Hybrid matchers** that combine the two matchers just described.

In schema matching the problem of *volume* does not take place, namely there is no need to increase computational performance, rather exploring new techniques for having more precise matchers could be necessary. Nevertheless, the availability of big data, brought recent research towards **learning-based** matchers, where pre-existing mappings are used to train machine learning algorithms that

Record linkage. Up to 2017, a lot of Privacy-Preserving Record Linkage (PPRL) protocols proposed for secure two-party private record linkage were not able to meet the following three requirements altogether: (1) **full end-to-end privacy**, (2) **perfect precision and recall** for the matching records and (3) **sub-quadratic computational complexity** [14]. After [14] pointed out this problem, new approaches were proposed and resulted in meeting all the three requirements described above [13], but still limited to certain adversarial models, that do not match well with real world scenarios (e.g. *honest-but-curious*), and does not take into account the multi-party scenario, where the problem of *volume*, *veracity* and *velocity* needs to be addressed [26]. Although there have been many prior works that looked into different security and privacy aspects of PPRL, there still remain major challenges that need to be addressed. As a starting point, more research is needed to understand membership inference and attribute inference attacks on PPRL techniques. For example, a malicious party may create many random records with the intent of matching with the other parties' actual records [11]. To tackle this problem, blockchain technologies are recently being proposed in literature [19].

Section 3. Co-supervisors/Candidate Co-operation Agreements

The project will be carried out in three years during which the PhD student will stay in one research institution and one university. During the first and the third year, the candidate will work in Athena Research Center (ARC) under the supervision of Prof. Minos Garofalakis (ARC). During the second year, the program will take place in Universitat Politecnica de Catalunya (UPC) under the supervision of Prof. Oscar Romero (UPC). The project will be a joint work of all parties, hence close co-operation is expected in the following way.

The progress of the project will be validated through frequent meetings between the candidate and his supervisors. The candidate will meet on a weekly basis with his home supervisor and one or two times per month with his host supervisor (the opposite when he will be hosted at UPC). Following typical business practice, the expectations and tasks planned for each meeting will be clearly communicated in advance, with a reasonable notice, both from the supervisors to the candidate and vice-versa. Standard tools of the trade will be used to boost collaboration, such as a shared repository for documents and code artifacts (e.g., Mendeley Library, Github), communication platforms (e.g. Skype, Teams).

Section 4. Proposed Education and Training Programme

The DEDS education and training programme is composed of several activities.

- **Research**, where doctoral candidates work on a novel research problem guided by two supervisors who will advise them to gradually become independent researchers.
- **Research-specific courses**, aimed at providing doctoral candidates with focused state-of-the-art technical skills pertaining to their research topic.
- **Innovation and entrepreneurship courses**, aimed at complementing the scientific training of doctoral candidates with business-related aspects such as entrepreneurship, intellectual property rights, etc.
- **Methodological and communication courses**, aimed at introducing the necessary research methods and communication skills.

- **Language courses**, aimed at introducing the local language at each partner university.
- **Summer and winter schools**, wherein candidates will obtain feedback about their research from invited researchers and practitioners, as well as get international contacts in both academia and industry.
- **Tutoring**, whereby candidates will be involved in teaching activities (e.g., supervising student projects and delivering exercises) while being coached by their supervisors or other experienced staff.
- **Knowledge dissemination and participation to scientific events**, aimed at allowing doctoral candidates to present and confront their findings, thereby familiarising themselves with essential practices such as peer-review and public debating.
- **External cooperation and secondments**, aimed at ensuring that the candidate participate actively in another research environment outside his/her home and host universities. These activities are realised typically with DEDS partner organisations.

Please detail in the following subsections your personalised education and training programme, taking into account your previous background and future career prospects. This programme must be approved from both co-supervisors.

Activities adding at least 30 ECTS credits must be outlined. A tabular listing of all activities performed or to be performed during the doctorate project is to be included. Group the activities according to the categories specified above. For each activity, the title, time, location, organiser, and ECTS credits must be included together with an indication of whether the activity has been completed. Please use this table:

| Activity | Place/Organised by | ECTS | General/Project course | Status |
|----------|--------------------|------|------------------------|--------|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

TPR: The table of activities must be updated with planned and hitherto completed activities.
RPR: The contents and the extent of the completed activities must be reported. It is expected that all training activities have been finalised in order to devote the last year of the project for finalising the Doctoral Dissertation.

Section 4.1. Planned Courses

Courses adding at least 20 ECTS credits must be outlined.

Only courses at doctorate level are approved. If a course at master level is deemed to be highly relevant for the doctorate project, the co-supervisors can establish a study group on the topic, which includes the master course and additional reading/discussion to bring it up to doctorate level. A written report on participation in a study group must be completed to get course credit. To ensure the scientific level, the study circle must be headed by a member of the scientific staff, who is Professor or Associate Professor (senior scientist level). A 2-3 ECTS study circle organised by the co-supervisors on the state of the art in the research field of the doctorate study is recommended.

TPR: The course table should be updated with more specific information for the completed courses, as well as with the rest of planned courses for the rest of the doctorate project.

RPR: The table must be updated, reporting the complete set of courses that are completed in the doctorate project.

Section 4.2. Knowledge Dissemination and Participation to Scientific Events

Detail the plan for dissemination of knowledge and findings from the project.

This could for instance be:

- Poster presentations at conferences/seminars.
- Presentations at conferences/seminars
- Newspaper articles or other popular presentations
- Teaching (lecturing and project supervision)

Each participation to scientific events must be accompanied by a written report by the doctorate candidate that relates the specific activity to the doctorate project. This report must be of general value for the project. Activities that relate to workshop and conference participation must not exceed 6 ECTS credits.

TPR: Updated plan for dissemination of knowledge and findings from the doctorate project other than those listed in Section 2 must be specified.

RPR: The final realisation of the knowledge and findings dissemination from the doctorate project other than those listed in Section 2 must be reported.

Section 4.3. External Co-operation

The doctorate candidate will spend time studying both in Greece ‘ and Spain. Furthermore, a secondment of three months will take place, where the candidate will join Spring Techno, where he will work on of a complex Federated Learning scenario with real data. During the following three years, all ESRs will meet in four different winter/summer schools to present their work, receive feedback, exchange ideas, and get exposed to new challenges. During these schools, candidates will have the opportunity to get in touch with academic and non-academic partners, presenting them their findings, reflecting on new opportunities, and opening the way for further collaboration. Finally, the candidate may co-operate with external researchers or research teams, in case that his work can be combined or merged with similar works of others.

Section 5. Agreements on Immaterial Rights to Patents

Patents and immaterial rights will be handled according to general rules applied by Athena Research Center, National and Kapodistrian University of Athens, and Universitat Politecnica de Catalunya.

Section 6. Financing Budget

This project is one of the 15 ESRs of Data Engineering for Data Science PhD programme, which is funded by the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 955895. The funding covers expenses related to the successful completion of the project, such as work equipment, research experiments, training activities and others that are relevant to the programme.

Section 7. Career Development Plan

In this section, the candidate's career plan and development are described. When the thesis is handed in (M4), this section is revisited in a self-contained document called "Career Development Plan" (CDP) to be signed by the candidate and the supervisors.

Section 7.1. Long-Term Career Objectives

DPP: Describe long-term career goals (over 5 years) and how to become able to reach those goals.

TPR, RPR, CDP: Update as needed if the career plans have evolved.

Section 7.2. Objectives Covered in Project

Describe which development objectives will be/have been achieved in the project with respect to

1. Research skills and techniques
2. Research management and co-operation
3. Communication skills
4. Other professional training
5. Networking activities and opportunities
6. Other activities with professional relevance

CDP: Include also published and accepted papers as well as completed course activities such that the document is self-contained.

Section 8. References

References

- [1] Jesus Antonanzas, Marta Arias, and Albert Bifet. "Sketches for Time-Dependent Machine Learning". In: 1 (2021), pp. 1–9. arXiv: 2108.11923. URL: <http://arxiv.org/abs/2108.11923>.
- [2] Luca Bonomi et al. "Frequent grams based embedding for privacy preserving record linkage". In: *ACM International Conference Proceeding Series* (2012), pp. 1597–1601. DOI: 10.1145/2396761.2398480.
- [3] Chris Clifton and Gunther Schadow. "Privacy Preserving Data Integration and Sharing". In: (), pp. 1–10.
- [4] Graham Cormode et al. "Synopsis for massive data: Samples, histograms, wavelets, sketches". In: *Foundations and Trends in Databases* 4.1-3 (2011), pp. 1–294. ISSN: 19317883. DOI: 10.1561/19000000004.
- [5] Tamraparni Dasu et al. "Mining database structure; or, how to build a data quality browser". In: (2002), p. 240. DOI: 10.1145/564691.564719.
- [6] Xin Luna Dong and Theodoros Rekatsinas. "Data integration and machine learning: A natural synergy". In: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2019), pp. 3193–3194. DOI: 10.1145/3292500.3332296.

- [7] Xin Luna Dong and Divesh Srivastava. “Big data integration”. In: *Proceedings - International Conference on Data Engineering* September (2013), pp. 1245–1248. ISSN: 10844627. DOI: 10.1109/ICDE.2013.6544914.
- [8] Cynthia Dwork and Aaron Roth. “The algorithmic foundations of differential privacy”. In: *Foundations and Trends in Theoretical Computer Science* 9.3-4 (2013), pp. 211–487. ISSN: 15513068. DOI: 10.1561/04000000042.
- [9] Ahmet M. Elbir, Burak Soner, and Sinem Coleri. *Federated Learning in Vehicular Networks*. 2020. DOI: 10.48550/ARXIV.2006.01412. URL: <https://arxiv.org/abs/2006.01412>.
- [10] Ivan P. Fellegi and Alan B. Sunter. “A Theory for Record Linkage”. In: *Journal of the American Statistical Association* 64.328 (1969), pp. 1183–1210. DOI: 10.1080/01621459.1969.10501049.
- [11] Aris Gkoulalas-Divanis et al. “Modern Privacy-Preserving Record Linkage Techniques: An Overview”. In: *IEEE Transactions on Information Forensics and Security* 16 (2021), pp. 4966–4987. ISSN: 15566021. DOI: 10.1109/TIFS.2021.3114026.
- [12] Rémi Gribonval et al. “Sketching Datasets for Large-Scale Learning (long version)”. In: (2020), pp. 1–35. arXiv: 2008.01839. URL: <http://arxiv.org/abs/2008.01839>.
- [13] Adam Groce, Peter Rindal, and Mike Rosulek. “Cheaper Private Set Intersection via Differentially Private Leakage”. In: *Proceedings on Privacy Enhancing Technologies* 2019.3 (2019), pp. 6–25. DOI: 10.2478/popets-2019-0034.
- [14] Xi He et al. “Composing Differential Privacy and Secure Computation: A case study on scaling private record linkage”. In: *Proceedings of the ACM Conference on Computer and Communications Security* (2017), pp. 1389–1406. ISSN: 15437221. DOI: 10.1145/3133956.3134030. arXiv: 1702.00535.
- [15] Jiawei Jiang et al. “SketchML: Accelerating distributed machine learning with data sketches”. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data* (2018), pp. 1269–1284. ISSN: 07308078. DOI: 10.1145/3183713.3196894.
- [16] Basit Khurram and Florian Kerschbaum. “SFour: A protocol for cryptographically secure record linkage at scale”. In: *Proceedings - International Conference on Data Engineering* 2020-April (2020), pp. 277–288. ISSN: 10844627. DOI: 10.1109/ICDE48307.2020.00031.
- [17] Mehmet Kuzu et al. “Efficient privacy-aware record integration”. In: *ACM International Conference Proceeding Series* (2013), pp. 167–178. DOI: 10.1145/2452376.2452398.
- [18] Ramez El-Masri and Gio Wiederhold. “Data Model Integration Using the Structural Model”. In: *Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’79. Boston, Massachusetts: Association for Computing Machinery, 1979, 191–202. ISBN: 089791001X. DOI: 10.1145/582095.582127. URL: <https://doi.org/10.1145/582095.582127>.
- [19] Thiago Nóbrega, Carlos Eduardo S. Pires, and Dimas Cassimiro Nascimento. “Blockchain-based Privacy-Preserving Record Linkage: enhancing data privacy in an untrusted environment”. In: *Information Systems* 102 (2021), p. 101826. ISSN: 03064379. DOI: 10.1016/j.is.2021.101826. URL: <https://doi.org/10.1016/j.is.2021.101826>.
- [20] Prayitno et al. “A systematic review of federated learning in the healthcare area: From the perspective of data properties and applications”. In: *Applied Sciences (Switzerland)* 11.23 (2021). ISSN: 20763417. DOI: 10.3390/app112311191.

- [21] Sebastian Riedel et al. “Relation extraction with matrix factorization and universal schemas”. In: *NAACL HLT 2013 - 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Main Conference* June (2013), pp. 74–84.
- [22] Diego Rodrigues and Altigran da Silva. “A study on machine learning techniques for the schema matching network problem”. In: *Journal of the Brazilian Computer Society* 27.1 (2021). ISSN: 16784804. DOI: 10.1186/s13173-021-00119-5.
- [23] Khalid Saleem. “Schema Matching and Integration in Large Scale Scenarios”. In: (2009).
- [24] Bohdan Shubyn et al. “Federated Learning for Anomaly Detection in Industrial IoT-enabled Production Environment Supported by Autonomous Guided Vehicles”. In: *Computational Science – ICCS 2022*. Ed. by Derek Groen et al. Cham: Springer International Publishing, 2022, pp. 409–421. ISBN: 978-3-031-08760-8.
- [25] Marcin Szymczak et al. “Content Data Based Schema Matching”. In: *Challenging Problems and Solutions in Intelligent Systems*. Ed. by Guy de Trè et al. Cham: Springer International Publishing, 2016, pp. 281–322. ISBN: 978-3-319-30165-5. DOI: 10.1007/978-3-319-30165-5_14. URL: https://doi.org/10.1007/978-3-319-30165-5_14.
- [26] Dinusha Vatsalan et al. “Privacy-preserving record linkage for big data: Current approaches and research challenges”. In: *Handbook of Big Data Technologies* (2017), pp. 851–895. DOI: 10.1007/978-3-319-49340-4_25.
- [27] Shui Yu. “Big privacy: Challenges and opportunities of privacy study in the age of big data”. In: *IEEE access* 4 (2016), pp. 2751–2763.

**European Joint Doctorate in
Data Engineering for Data Science (DEDS)
Doctorate Project Plan²
Thesis Title
First Name Last Name**

This page must be completed and sent together with the project plan/report in a pdf file to the chair of the Candidate Progress Committee.

Project title:
 Name of doctorate candidate:
 Email:
 Supervisor:
 Home University:
 Co-supervisor:
 Host University:
 Secondment supervisor:
 Partner organisation:
 Date of enrolment:
 Expected date of completion:

Signatures

The Doctorate Candidate

.....
 Date:

The Supervisor from the Home University

Professor
 Date:

The Supervisor from the Host University

Professor
 Date:

The Secondment Supervisor

.....
 Date:

²Choose the appropriate heading among the three

The Chair of the Candidate Progress Committee

Professor

Date: