

RESEARCH CENTRE

**Inria Center
at the University of Lille**

IN PARTNERSHIP WITH:
CNRS, Université de Lille

2022

ACTIVITY REPORT

Project-Team
MAGNET

Machine Learning in Information Networks

IN COLLABORATION WITH: Centre de Recherche en Informatique,
Signal et Automatique de Lille

DOMAIN

Perception, Cognition and Interaction

THEME

**Data and Knowledge Representation and
Processing**

Inria

Contents

Project-Team MAGNET	1
1 Team members, visitors, external collaborators	2
2 Overall objectives	3
3 Research program	3
4 Application domains	6
5 Social and environmental responsibility	6
5.1 Footprint of research activities	6
5.2 Impact of research results	6
6 Highlights of the year	7
7 New software and platforms	7
7.1 New software	7
7.1.1 CoRTeX	7
7.1.2 Mangoes	7
7.1.3 metric-learn	8
7.1.4 MyLocalInfo	8
7.1.5 Voice Transformer	8
7.1.6 declearn	9
7.1.7 fairgrad	10
8 New results	10
8.1 Natural Language Processing	10
8.2 Decentralized Learning	11
8.3 Privacy and Machine Learning	12
8.4 Federated Learning in Medicine	14
8.5 Learning and Speech Recognition	15
8.6 Fairness and Transparency	17
8.7 Graph-Based Learning	17
9 Partnerships and cooperations	18
9.1 International initiatives	18
9.1.1 Associate Teams in the framework of an Inria International Lab or in the framework of an Inria International Program	18
9.1.2 Participation in other International Programs	19
9.2 European initiatives	20
9.2.1 Horizon Europe	20
9.3 National initiatives	21
9.3.1 ANR Pamela (2016–2022)	21
9.3.2 ANR DEEP-Privacy (2019–2023)	21
9.3.3 ANR-JCJC PRIDE (2020–2025)	21
9.3.4 ANR PMR (2020-2024)	22
9.3.5 FedMalin. INRIA Defi (2021-2024)	22
9.3.6 FLAMED: Federated Learning and Analytics on Medical Data. INRIA Action Exploratoire (2020-2024)	22
9.3.7 COMANCHE: Computational Models of Lexical Meaning and Change. INRIA Action Exploratoire (2022-2026)	23
9.3.8 IPoP, Projet interdisciplinaire sur la protection des données personnelles, PEPR Cybersécurité (2022-2028).	23
9.3.9 HyAIAI. INRIA Defi (2019-2022)	24

9.4	Regional initiatives	24
9.4.1	Chaire TIP: Transparent artificial Intelligence preserving Privacy	24
9.4.2	Fairness in Decentralized and Privacy Preserving Machine Learning STaRS (2021-2023)	24
9.4.3	CAPS'UL: CAmпус Participatif en Santé numérique du site Universitaire de Lille. Santé numérique, PIA4. (2022-2027)	24
10	Dissemination	25
10.1	Promoting scientific activities	25
10.1.1	Scientific events: organisation	25
10.1.2	Scientific events: selection	25
10.1.3	Journal	25
10.1.4	Invited talks	26
10.1.5	Scientific expertise	26
10.1.6	Research administration	26
10.2	Teaching - Supervision - Juries	26
10.2.1	Teaching	26
10.2.2	Supervision	27
10.2.3	Juries	28
10.3	Popularization	29
10.3.1	Articles and contents	29
11	Scientific production	29
11.1	Major publications	29
11.2	Publications of the year	30

Project-Team MAGNET

Creation of the Project-Team: 2016 May 01

Keywords

Computer sciences and digital sciences

- A3.1.3. – Distributed data
- A3.1.4. – Uncertain data
- A3.4. – Machine learning and statistics
 - A3.4.1. – Supervised learning
 - A3.4.2. – Unsupervised learning
 - A3.4.4. – Optimization and learning
 - A3.4.6. – Neural networks
 - A3.4.8. – Deep learning
- A4.8. – Privacy-enhancing technologies
- A9.4. – Natural language processing

Other research topics and application domains

- B2. – Health
- B9.5.1. – Computer science
- B9.5.6. – Data science
- B9.6.8. – Linguistics
- B9.6.10. – Digital humanities
- B9.9. – Ethics
- B9.10. – Privacy

1 Team members, visitors, external collaborators

Research Scientists

- Aurélien Bellet [INRIA, Researcher, HDR]
- Pascal Denis [INRIA, Researcher]
- Michael Perrot [INRIA, ISFP]
- Jan Ramon [INRIA, Senior Researcher, HDR]
- Damien Sileo [INRIA, ISFP, from Oct 2022]

Faculty Members

- Marc Tommasi [Team leader, UNIV LILLE, Professor, HDR]
- Mikaela Keller [UNIV LILLE, Associate Professor]

Post-Doctoral Fellows

- Carlos Mauricio Cotrini Jimenez [INRIA, from Mar 2022 until Aug 2022]
- Vitalii Emelianov [INRIA, from Nov 2022]
- Batiste Le Bars [INRIA]
- Cesar Sabater [INRIA, from Jul 2022]

PhD Students

- Mahsa Asadi [INRIA, until Nov 2022]
- Antoine Barczewski [INRIA, from May 2022]
- Moitree Basu [INRIA]
- Ioan Tudor Cebere [INRIA, from Nov 2022]
- Edwige Cyffers [UNIV LILLE]
- Marc Damie [INRIA, from May 2022]
- Le Dinh-Viet-Toan [UNIV LILLE, from Oct 2022]
- Aleksei Korneev [UNIV LILLE, from Dec 2022]
- Bastien Liétard [INRIA and UNIV LILLE, from Nov 2022]
- Gaurav Maheshwari [INRIA]
- Paul Mangold [INRIA]
- Amal Mawass [UNIV LILLE, until May 2022]
- Arijus Pleska [INRIA and UNIV LILLE]
- Cesar Sabater [INRIA, until Jun 2022]

Technical Staff

- Paul Andrey [INRIA, Engineer, from Jun 2022]
- Antoine Barczewski [UNIV LILLE, Engineer]
- Nathan Bigaud [INRIA, Engineer, from Apr 2022]
- Marc Damie [INRIA, Engineer, until Apr 2022]
- Antonin Duez [INRIA, Engineer, from Mar 2022]
- Rishabh Gupta [INRIA, Engineer, from Apr 2022]
- Jean-Paul Lam [INRIA, Engineer, from Mar 2022]
- Joseph Renner [INRIA, Engineer]
- Sophie Villerot [INRIA, Engineer]

Interns and Apprentices

- Maximiliano Vargas Vargas [INRIA]

Administrative Assistant

- Aurore Dalle [INRIA]

External Collaborator

- Remi Gilleron [UNIV LILLE, HDR]

2 Overall objectives

The main objective of MAGNET is to develop original machine learning methods for networked data. We consider information networks in which the data consist of feature vectors or texts. We model such networks as graphs wherein nodes correspond to entities (documents, spans of text, users, datasets, learners etc.) and edges correspond to relations between entities (similarity, answer, co-authoring, friendship etc.). In *Mining and Learning in Graphs*, our main research goal is to efficiently search for the best hidden graph structure to be generated for solving a given learning task which exploits the relationships between entities. In *Machine Learning for Natural Language Processing* the objective is to go beyond vectorial classification to solve tasks like coreference resolution and entity linking, temporal structure prediction, and discourse parsing. In *Decentralized Machine Learning* we address the problem of learning in a private, fair and energy efficient way when data are naturally distributed in a network.

The challenges are the dimensionality of the input space, possibly the dimensionality of the output space, the high level of dependencies between the data, the inherent ambiguity of textual data and the limited amount of human labeling. We are interested in making machine learning approaches more acceptable to society. Privacy, sobriety and fairness are important issues that pertain to this research line, and we are interested in the empowerment of end users in the machine learning processes.

3 Research program

The research program of MAGNET is structured along three main axes.

Axis 1: Mining and Learning in Graphs This axis is the backbone of the team. Most of the techniques and algorithms developed in this axis are known by the team members and have impact on the two other axes. We address the following questions and objectives:

How to adaptively build graphs with respect to the given tasks? We study adaptive graph construction along several directions. The first one is to learn the best similarity measure for the graph construction. The second one is to combine different views over the data in the graph construction and learn good representations. We also study weak forms of supervision like comparisons.

How to design methods able to achieve a good trade-off between predictive accuracy and computational complexity? We develop new algorithms for efficient graph-based learning (for instance node prediction or link prediction). In order to deal with scalability issues, our approach is based on optimization, graph sparsification techniques and graph sampling methods.

How to find patterns in graphs based on efficient computations of some statistics? We develop graph mining algorithms and statistics in the context of correlated data.

Axis 2: Machine Learning for Natural Language Processing In this axis, we address the general question that relates graph-based learning and Natural Language Processing (NLP): *How to go beyond vectorial classification models in NLP tasks?* We study the combination of learning representation, structured prediction and graph-based learning methods. Data sobriety and fairness are major constraints we want to deal with. The targeted NLP tasks are coreference resolution and entity linking, temporal structure prediction, and discourse parsing.

Axis 3: Decentralized Machine Learning and Privacy In this axis, we study *How to design private by design machine learning algorithms?* Taking as an opportunity the fact that data collection is now decentralized on smart devices, we propose alternatives to large data centers where data are gathered by developing collaborative and personalized learning.

Contrary to many machine learning approaches where data points and tasks are considered in isolation, we think that a key point of this research is to be able to leverage the relationships between data and learning objectives. Therefore, using graphs as an abstraction of information networks is a major playground for MAGNET. Research related to graph data is a transversal axis, describing a layer of work supporting two other axes on Natural Language Processing and decentralized learning. The machine learning and mining in graphs communities have evolved, for instance taking into account data streams, dynamics but maybe more importantly, focusing on deep learning. Deep neural nets are here to stay, and they are useful tools to tackle difficult problems so we embrace them at different places in the three axes.

MAGNET conducts research along the three axis described above but will put more emphasis on social issues of machine learning. In the context of the recent deployment of artificial intelligence into our daily lives, we are interested in making machine learning approaches more acceptable to society. Privacy, sobriety and fairness are important issues that pertain to this research line, but more generally we are interested in the empowerment of end users in the machine learning processes. Reducing the need of one central authority and pushing more the data processing on the user side, that is decentralization, also participates to this effort. Reducing resources means reducing costs and energy and contributes to building more accessible technologies for companies and users. By considering learning tasks in a more personalized way, but increasing collaboration, we think that we can design solutions that work in low resources regime, with less data or supervision.

In MAGNET we emphasize a different approach than blindly brute-forcing tasks with loads of data. Applications to social sciences for instance have different needs and constraints that motivate data sobriety, fairness and privacy. We are interested in weaker supervision, by leveraging structural properties described in graphs of data, relying on transfer and multi-task learning when faced with graphs of tasks and users. Algorithmic and statistical challenges related to the graph structure of the data still contain open questions. On the statistical side, examples are to take dependencies into account, for instance to compute a mean, to reduce the need of sampling by exploiting known correlations. For the algorithmic point of view, going beyond unlabeled undirected graphs, in particular considering attributed graphs containing text or other information and addressing the case of distributed graphs while maintaining formal guarantees are getting more attention.

In the second axis devoted to NLP, we focus our research on graph-based and representation learning into several directions, all aiming at learning *richer, more robust, and more transferable linguistic representations*. This research program will attempt to bring about strong cross-fertilizations with the other axes, addressing problems in graph, privacy and fairness and making links with decentralized learning. At the intersection between graph-based and representation learning, we will first develop graph embedding algorithms for deriving linguistic representations which are able to capture higher-level semantic and world-knowledge information which eludes strictly distributional models. As an initial step, we envision leveraging pre-existing ontologies (e.g., WordNet, DBpedia), from which one can easily derive interesting similarity graphs between words or noun phrases. We also plan to investigate innovative ways of articulating graph-based semi-supervised learning algorithms and word embedding techniques. A second direction involves learning representations that are more robust to bias, privacy attacks and adversarial examples. Thus, we intend to leverage recent adversarial training strategies, in which an adversary attempts to recover sensitive attributes (e.g., gender, race) from the learned representations, to be able to neutralize bias or to remove sensitive features. An application domain for this line of research is for instance speech data. The study of learning private representation with its link to fairness in the decentralized setting is another important research topic for the team. In this context of fairness, we also intend to develop similar algorithms for detecting slants, and ultimately for generating de-biased or “re-biased” versions of text embeddings. An illustration is on political slant in written texts (e.g., political speeches and manifestos). Thirdly, we intend to learn linguistic representations that can transfer more easily across languages and domains, in particular in the context of structured prediction problems for low-resource languages. For instance, we first propose to jointly learn model parameters for each language (and/or domains) in a multi-task setting, and leverage a (pre-existing or learned) graph encoding structural similarities between languages (and/or domains). This type of approach would nicely tie in with our previous work on multilingual dependency parsing and on learning personalized models. Furthermore, we will also study how to combine and adapt some neural architectures recently introduced for sequence-to-sequence problems in order to enable transfer of language representations.

In terms of technological transfer, we maintain collaborations with researchers in the humanities and the social sciences, helping them to leverage state-of-the-art NLP techniques to develop new insights to their research by extracting relevant information from large amounts of texts.

The third axis is on distributed and decentralized learning and privacy preserving machine learning. Recent years have seen the evolution of information systems towards ubiquitous computing, smart objects and applications fueled by artificial intelligence. Data are collected on smart devices like smartphones, watches, home devices etc. They include texts, locations, social relationships. Many sensitive data—race, gender, health conditions, tastes etc— can be inferred. Others are just recorded like activities, social relationships but also biometric data like voice and measurements from sensor data. The main tendency is to transfer data into central servers mostly owned by a few tier parties. The situation generates high privacy risks for the users for many reasons: loss of data control, unique entry point for data access, unsolicited data usage etc. But it also increases monopolistic situations and tends to develop oversized infrastructures. The centralized paradigm also has limits when data are too huge such as in the case of multiple videos and sensor data collected for autonomous driving. Partially or fully decentralized systems provide an alternative, to emphasis data exploitation rather than data sharing. For MAGNET, they are source of many new research directions in machine learning at two scales: at the algorithmic level and at a systemic level.

At the algorithmic level the question is to develop new privacy preserving algorithms in the context of decentralized systems. In this context, data remains where it has been collected and learning or statistical queries are processed at the local level. An important question we study is to take into account and measure the impact of collaboration. We also aim at developing methods in the online setting where data arrives continuously or participants join and leave the collaboration network. The granularity of exchanges, the communication cost and the dynamic scenarios, are also studied. On the privacy side, decentralization is not sufficient to establish privacy guarantees because learned models together with the dynamics of collaborative learning may reveal private training data if the models are published or if the communications are observed. But, although it has not been yet well established, decentralization can naturally increase privacy-utility ratio. A direction of research is to formally prove the privacy gain when randomized decentralized protocols are used during learning. In some situations, for instance

when part of the data is not sensitive or when trusted servers can be used, a combination between a fully decentralized and a centralized approach is very relevant. In this setting, the question is to find a good trade-off between local versus global computations.

At the systemic layer, in MAGNET we feel that there is a need for research on a global and holistic level, that is to consider full processes involving learning, interacting, predicting, reasoning, repeating etc. rather than studying the privacy of isolated learning algorithms. Our objective is to design languages for describing processes (workflows), data (database schema, background knowledge), population statistics, privacy properties of algorithms, privacy requirements and other relevant information. This is fully aligned with recent trends that aim at giving to statistical learning a more higher level of formal specifications and illustrates our objective for more acceptable and transparent machine learning. We also work towards more robust privacy-friendly systems, being able to handle a wider range of malicious behavior such as collusion to obtain information or inputting incorrect data to obtain information or to influence the result of collaborative computations. From the transfer point of view, we plan to apply transparent, privacy-friendly in significant application domains, such as medicine, surveying, demand prediction and recommendation. In this context, we are interested to understand the appreciation of humans of transparency, verifiability, fairness, privacy-preserving and other trust-increasing aspects of our technologies.

4 Application domains

Our application domains cover health, mobility, social sciences and voice technologies.

Health Privacy is of major importance in the health domain. We contribute to develop methods to give access to the use of data in a private way rather than to the data itself centralized in vulnerable single locations. As an example, we are working with hospitals to develop the means of multicentric studies with privacy guarantees. A second example is personalized medicine where personal devices collect private and highly sensitive data. Potential applications of our research allow to keep data on device and to privately compute statistics.

Social sciences Our NLP research activities are rooted in linguistics, but learning unbiased representations of texts for instance or simply identifying unfair representations also have impacts in political sciences and history.

Music information retrieval By using analogies between language and music (symbolic notation) we tackle music information retrieval tasks such as style classification and structure detection.

Voice technologies We develop methods for privacy in speech that can be embedded in software suites dedicated to voice-based interaction systems.

5 Social and environmental responsibility

5.1 Footprint of research activities

Some of our research activities are energy intensive and we will work to reduce this carbon footprint in the future. Parts of the new research project FedMalin (See Section 9.3.5) is dedicated to this objective for the Federated Learning setting.

5.2 Impact of research results

The main research topics of the team contribute to improve transparency, fairness and privacy in machine learning and reduce bias in natural language processing.

6 Highlights of the year

Two European projects have been accepted. JAN RAMON will coordinate the FLUTE project (starting in 2023) and is an active member of the TRUMPET project. Both projects are related to privacy and security in the health domain.

Magnet is involved in the PEPR on CyberSecurity. The IPoP project has been selected. Magnet will conduct research on privacy preserving machine learning and fairness in machine learning.

PASCAL DENIS is the principal investigator of the new INRIA Exploratory project (Action Exploratoire) called COMMANCHE on Computational Models of Lexical Meaning and Change. The project opens new collaborations for Magnet in the domain of Social Sciences (linguistics) with labs in Lille, Paris and Nancy.

AURÉLIEN BELLET and GIOVANNI NEGLIA have launched the FedMalin INRIA challenge on federated learning.

BRIJ MOHAN LAL SRIVASTAVA, NATHALIE VAUQUIER and EMMANUEL VINCENT (from Multispeech) have created [Nijta](#), a startup on speech anonymization.

Two PhD defenses: Cesar Sabater ([41]) and Mahsa Asadi ([40]).

DINH-VIET-TOAN LE succeeded in the competitive selection and got a funding from the AIPhD program. His PhD is supervised by MIKAELA KELLER and LOUIS BIGO from the AlgoMus team in CRIStAL. The PhD opens a new application of Magnet's research on NLP to the musical domain.

7 New software and platforms

7.1 New software

7.1.1 CoRTeX

Name: Python library for noun phrase COreference Resolution in natural language TEXTs

Keyword: Natural language processing

Functional Description: CoRTeX is a LGPL-licensed Python library for Noun Phrase coreference resolution in natural language texts. This library contains implementations of various state-of-the-art coreference resolution algorithms, including those developed in our research. In addition, it provides a set of APIs and utilities for text pre-processing, reading the CONLL2012 and CONLLU annotation formats, and performing evaluation, notably based on the main evaluation metrics (MUC, B-CUBED, and CEAF). As such, CoRTeX provides benchmarks for researchers working on coreference resolution, but it is also of interest for developers who want to integrate a coreference resolution within a larger platform. It currently supports use of the English or French language.

Contact: Pascal Denis

Participant: Pascal Denis

Partner: Orange Labs

7.1.2 Mangoes

Name: MAGnet liNGuistic wOrd vEctorS

Keywords: Word embeddings, NLP

Functional Description: Mangoes is a toolbox for constructing and evaluating static and contextual token vector representations (aka embeddings). The main functionalities are:

- Contextual embeddings: Access a large collection of pretrained transformer-based language models, Pre-train a BERT language model on a corpus, Fine-tune a BERT language model for a number of extrinsic tasks, Extract features/predictions from pretrained language models.
- Static embeddings: Process textual data and compute vocabularies and co-occurrence matrices. Input data should be raw text or annotated text, Compute static word embeddings with different

state-of-the-art unsupervised methods, Propose statistical and intrinsic evaluation methods, as well as some visualization tools, Generate context dependent embeddings from a pretrained language model.

Future releases will include methods for injecting lexical and semantic knowledge into token and multi-model embeddings, and interfaces into common external knowledge resources.

URL: <https://gitlab.inria.fr/magnet/mangoes>

Contact: Nathalie Vauquier

7.1.3 metric-learn

Keywords: Machine learning, Python, Metric learning

Functional Description: Distance metrics are widely used in the machine learning literature. Traditionally, practitioners would choose a standard distance metric (Euclidean, City-Block, Cosine, etc.) using a priori knowledge of the domain. Distance metric learning (or simply, metric learning) is the sub-field of machine learning dedicated to automatically constructing optimal distance metrics.

This package contains efficient Python implementations of several popular metric learning algorithms.

URL: <https://github.com/scikit-learn-contrib/metric-learn>

Contact: Aurélien Bellet

Partner: Parietal

7.1.4 MyLocalInfo

Keywords: Privacy, Machine learning, Statistics

Functional Description: Decentralized algorithms for machine learning and inference tasks which (1) perform as much computation as possible locally and (2) ensure privacy and security by avoiding that personal data leaves devices.

Contact: Nathalie Vauquier

7.1.5 Voice Transformer

Name: Voice Transformer

Keywords: Speech, Privacy

Scientific Description: Three privacy-driven voice transformation methods have been implemented, improved and compared in this library: VoiceMask: This voice transformation method described in [Qia+18, Qia+17] performs frequency warping based on the composition of a quadratic function and a bilinear function using two different parameters. VTLN-based conversion: This voice conversion method described in [SN03] represents each speaker by a set of spectra for k phonetic classes, and maps the original speech to the target speaker's voice by finding the transformation parameters that minimise the distance between target class spectra and transformed original class spectra. The VPC method is inspired from the speaker anonymisation method proposed in [Fan+19], which performs voice conversion based on x-vectors [Sny+18], a fixed-length representation of speech signals that form the basis of state-of-the-art speaker verification systems. We have brought several improvements to this method such as pitch transformation, and new design choices for x-vectors selection

[Fan+19] F. Fang, X. Wang, J. Yamagishi, I. Echizen, M. Todisco, N. Evans, and J.-F. Bonastre. "Speaker Anonymization Using x-vector and Neural Waveform Models". In: Proceedings of the 10th ISCA Speech Synthesis Workshop. 2019, pp. 155–160. [Qia+17] J. Qian, H. Du, J. Hou, L. Chen, T. Jung, X.

Li, Y. Wang, and Y. Deng. “Voice- Mask: Anonymize and Sanitize Voice Input on Mobile Devices”. In: arXiv preprint arXiv:1711.11460, 2017 abs/1711.11460 (2017). [Qia+18] J. Qian, H. Du, J. Hou, L. Chen, T. Jung, and X.-Y. Li. “Hidebehind: Enjoy Voice Input with Voiceprint Unclonability and Anonymity”. In: Proceedings of the 16th ACM Conference on Embedded Networked Sensor Systems. SenSys ’18. Shenzhen, China: Association for Computing Machinery, 2018, pp. 82–94. [SN03] D. Sundermann and H. Ney. “VTLN-based voice conversion”. In: Proceedings of the 3rd IEEE International Symposium on Signal Processing and Information Technology (IEEE Cat. No.03EX795). 2003, pp. 556–559.

[Sny+18] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur. “X-vectors: Robust DNN embeddings for speaker recognition”. In: Proceedings of ICASSP 2018 - 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2018, pp. 5329–5333.

Functional Description: COMPRISE Voice Transformer is an open source tool that increases the privacy of users of voice interfaces by converting their voice into another person’s voice without modifying the spoken message. It ensures that any information extracted from the transformed voice can hardly be traced back to the original speaker, as validated through state-of-the-art biometric protocols, and it preserves the phonetic information required for human labelling and training of speech-to-text models.

Release Contributions: This version gives access to the 2 generations of tools that have been used to transform the voice, as part of the COMPRISE project (<https://www.compriseh2020.eu/>). The first one is a python library that implements 2 basic voice conversion methods, both using VLTN. The second one implements an anonymization method using x-vectors and neural waveform models.

News of the Year: We modified the x-vector based transformer by fixing the percentile-based pitch conversion method, using conda in Docker to fix issues with the Python version, and adding data from the speaker pool to simplify quick start.

URL: https://gitlab.inria.fr/comprise/voice_transformation

Contact: Marc Tommasi

Participants: Nathalie Vauquier, Brij Mohan Lal Srivastava, Marc Tommasi, Emmanuel Vincent, Md Sahidullah

7.1.6 declearn

Keyword: Federated learning

Scientific Description: declearn is a python package providing with a framework to perform federated learning, i.e. to train machine learning models by distributing computations across a set of data owners that, consequently, only have to share aggregated information (rather than individual data samples) with an orchestrating server (and, by extension, with each other).

The aim of declearn is to provide both real-world end-users and algorithm researchers with a modular and extensible framework that:

- builds on abstractions general enough to write backbone algorithmic code agnostic to the actual computation framework, statistical model details or network communications setup - designs modular and combinable objects, so that algorithmic features, and more generally any specific implementation of a component (the model, network protocol, client or server optimizer...) may easily be plugged into the main federated learning process - enabling users to experiment with configurations that intersect unitary features - provides with functioning tools that may be used out-of-the-box to set up federated learning tasks using some popular computation frameworks (scikit-learn, tensorflow, pytorch...) and federated learning algorithms (FedAvg, Scaffold, FedYogi...)
- provides with tools that enable extending the support of existing tools and APIs to custom functions and classes without having to hack into the source code, merely adding new features (tensor libraries, model classes, optimization plug-ins, orchestration algorithms, communication protocols...) to the party

At the moment, declearn has been focused on so-called "centralized" federated learning that implies a central server orchestrating computations, but it might become more oriented towards decentralized processes in the future, that remove the use of a central agent.

Functional Description: This library provides the two main components to perform federated learning: - the client, to be run by each participant, performs the learning on local data et releases only the result of the computation - the server orchestrates the process and aggregates the local models in a global model

URL: <https://gitlab.inria.fr/magnet/declearn/declearn2>

Contact: Aurélien Bellet

Participants: Paul Andrey, Aurélien Bellet, Nathan Bigaud, Marc Tommasi, Nathalie Vauquier

Partner: CHRU Lille

7.1.7 fairgrad

Name: FairGrad: Fairness Aware Gradient Descent

Keywords: Fairness, Fair and ethical machine learning, Machine learning, Classification

Functional Description: FairGrad is an easy to use general purpose approach in Machine Learning to enforce fairness in gradient descent based methods

URL: <https://github.com/saist1993/fairgrad>

Authors: Gaurav Maheshwari, Michael Perrot

Contact: Michael Perrot

8 New results

8.1 Natural Language Processing

Chop and change: Anaphora resolution in instructional cooking videos, [33]

Linguistic ambiguities arising from changes in entities in action flows are a key challenge in instructional cooking videos. In particular, temporally evolving entities present rich and to date understudied challenges for anaphora resolution. For example "oil" mixed with "salt" is later referred to as a "mixture". In this paper we propose novel annotation guidelines to annotate recipes for the anaphora resolution task, reflecting change in entities. Moreover, we present experimental results for end-to-end multimodal anaphora resolution with the new annotation scheme and propose the use of temporal features for performance improvement.

Improving Tokenization Expressiveness With Pitch Intervals, [48]

Training sequence models such as transformers with symbolic music requires a representation of music as sequences of atomic elements called tokens. State-of-the-art music tokenizations encode pitch values explicitly, which complicates the ability of a machine learning model to generalize musical knowledge at different keys. We propose tracks for a tokenization encoding pitch intervals rather than pitch values, resulting in transposition invariant representations. The musical expressiveness of this new tokenization is evaluated through two MIR classification tasks: composer classification and end of phrase detection. We release publicly the code produced in this research.

Fair NLP Models with Differentially Private Text Encoders, [28]

Encoded text representations often capture sensitive attributes about individuals (e.g., race or gender), which raise privacy concerns and can make downstream models unfair to certain groups. In this work, we propose FEDERATE, an approach that combines ideas from differential privacy and adversarial training to learn private text representations which also induces fairer models. We empirically evaluate the trade-off between the privacy of the representations and the fairness and accuracy of the downstream model on four NLP datasets. Our results show that FEDERATE consistently improves upon previous methods, and thus suggest that privacy and fairness can positively reinforce each other.

Anaphora Resolution in Dialogue: System Description (CODI-CRAC 2022 Shared Task), [21]

We describe three models submitted for the CODI-CRAC 2022 shared task. To perform identity anaphora resolution, we test several combinations of the incremental clustering approach based on the Workspace Coreference System (WCS) with other coreference models. The best result is achieved by adding the "cluster merging" version of the coref-hoi model, which brings up to 10.33% improvement over vanilla WCS clustering. Discourse deixis resolution is implemented as multi-task learning: we combine the learning objective of corefhoi with anaphor type classification. We adapt the higher-order resolution model introduced in Joshi et al. (2019) for bridging resolution given gold mentions and anaphors.

8.2 Decentralized Learning

Refined Convergence and Topology Learning for Decentralized SGD with Heterogeneous Data, [42]

One of the key challenges in decentralized and federated learning is to design algorithms that efficiently deal with highly heterogeneous data distributions across agents. In this paper, we revisit the analysis of the popular Decentralized Stochastic Gradient Descent algorithm (D-SGD) under data heterogeneity. We exhibit the key role played by a new quantity, called neighborhood heterogeneity, on the convergence rate of D-SGD. By coupling the communication topology and the heterogeneity, our analysis sheds light on the poorly understood interplay between these two concepts. We then argue that neighborhood heterogeneity provides a natural criterion to learn data-dependent topologies that reduce (and can even eliminate) the otherwise detrimental effect of data heterogeneity on the convergence time of D-SGD. For the important case of classification with label skew, we formulate the problem of learning such a good topology as a tractable optimization problem that we solve with a Frank-Wolfe algorithm. As illustrated over a set of simulated and real-world experiments, our approach provides a principled way to design a sparse topology that balances the convergence speed and the per-iteration communication costs of D-SGD under data heterogeneity.

D-Cliques: Compensating for Data Heterogeneity with Topology in Decentralized Federated Learning, [22]

The convergence speed of machine learning models trained with Federated Learning is significantly affected by heterogeneous data partitions, even more so in a fully decentralized setting without a central server. In this paper, we show that the impact of label distribution skew, an important type of data heterogeneity, can be significantly reduced by carefully designing the underlying communication topology. We present D-Cliques, a novel topology that reduces gradient bias by grouping nodes in sparsely interconnected cliques such that the label distribution in a clique is representative of the global label distribution. We also show how to adapt the updates of decentralized SGD to obtain unbiased gradients and implement an effective momentum with D-Cliques. Our extensive empirical evaluation on MNIST and CIFAR10 demonstrates that our approach provides similar convergence speed as a fully-connected topology, which provides the best convergence in a data heterogeneous setting, with a significant reduction in the number of edges and messages. In a 1000-node topology, D-Cliques require 98% less edges and 96% less total messages, with further possible gains using a small-world topology across cliques.

PEPPER: Empowering User-Centric Recommender Systems over Gossip Learning, [15]

Recommender systems are proving to be an invaluable tool for extracting user-relevant content helping users in their daily activities (e.g., finding relevant places to visit, content to consume, items to purchase). However, to be effective, these systems need to collect and analyze large volumes of personal data (e.g., location check-ins, movie ratings, click rates .. etc.), which exposes users to numerous privacy threats. In this context, recommender systems based on Federated Learning (FL) appear to be a promising solution for enforcing privacy as they compute accurate recommendations while keeping personal data on the users' devices. However, FL, and therefore FL-based recommender systems, rely on a central server that can experience scalability issues besides being vulnerable to attacks. To remedy this, we propose PEPPER, a decentralized recommender system based on gossip learning principles. In PEPPER, users gossip model updates and aggregate them asynchronously. At the heart of PEPPER reside two key components: a personalized peer-sampling protocol that keeps in the neighborhood of each node, a proportion of nodes that have similar interests to the former and a simple yet effective model aggregation function that builds a model that is better suited to each user. Through experiments on three real datasets implementing two use cases: a location check-in recommendation and a movie recommendation, we demonstrate that our solution converges up to 42% faster than with other decentralized solutions providing up to 9% improvement on average performance metric such as hit ratio and up to 21% improvement on long tail performance compared to decentralized competitors.

Collaborative Algorithms for Online Personalized Mean Estimation, [14]

We consider an online estimation problem involving a set of agents. Each agent has access to a (personal) process that generates samples from a real-valued distribution and seeks to estimate its mean. We study the case where some of the distributions have the same mean, and the agents are allowed to actively query information from other agents. The goal is to design an algorithm that enables each agent to improve its mean estimate thanks to communication with other agents. The means as well as the number of distributions with same mean are unknown, which makes the task nontrivial. We introduce a novel collaborative strategy to solve this online personalized mean estimation problem. We analyze its time complexity and introduce variants that enjoy good performance in numerical experiments. We also extend our approach to the setting where clusters of agents with similar means seek to estimate the mean of their cluster.

8.3 Privacy and Machine Learning

An Accurate, Scalable and Verifiable Protocol for Federated Differentially Private Averaging, [16]

Learning from data owned by several parties, as in federated learning, raises challenges regarding the privacy guarantees provided to participants and the correctness of the computation in the presence of malicious parties. We tackle these challenges in the context of distributed averaging, an essential building block of federated learning algorithms. Our first contribution is a scalable protocol in which participants exchange correlated Gaussian noise along the edges of a graph, complemented by independent noise added by each party. We analyze the differential privacy guarantees of our protocol and the impact of the graph topology under colluding malicious parties, showing that we can nearly match the utility of the trusted curator model even when each honest party communicates with only a logarithmic number of other parties chosen at random. This is in contrast with protocols in the local model of privacy (with lower utility) or based on secure aggregation (where all pairs of users need to exchange messages). Our second contribution enables users to prove the correctness of their computations without compromising the efficiency and privacy guarantees of the protocol. Our construction relies on standard cryptographic primitives like commitment schemes and zero knowledge proofs.

Differentially Private Federated Learning on Heterogeneous Data, [32]

Federated Learning (FL) is a paradigm for large-scale distributed learning which faces two key challenges: (i) training efficiently from highly heterogeneous user data, and (ii) protecting the privacy of participating users. In this work, we propose a novel FL approach (DP-SCAFFOLD) to tackle these two challenges

together by incorporating Differential Privacy (DP) constraints into the popular SCAFFOLD algorithm. We focus on the challenging setting where users communicate with a "honest-but-curious" server without any trusted intermediary, which requires to ensure privacy not only towards a third party observing the final model but also towards the server itself. Using advanced results from DP theory and optimization, we establish the convergence of our algorithm for convex and non-convex objectives. Our paper clearly highlights the trade-off between utility and privacy and demonstrates the superiority of DP-SCAFFOLD over the state-of-the-art algorithm DP-FedAvg when the number of local updates and the level of heterogeneity grows. Our numerical results confirm our analysis and show that DP-SCAFFOLD provides significant gains in practice.

GAP: Differentially Private Graph Neural Networks with Aggregation Perturbation, [34]

In this paper, we study the problem of learning Graph Neural Networks (GNNs) with Differential Privacy (DP). We propose a novel differentially private GNN based on Aggregation Perturbation (GAP), which adds stochastic noise to the GNN's aggregation function to statistically obfuscate the presence of a single edge (edge-level privacy) or a single node and all its adjacent edges (node-level privacy). Tailored to the specifics of private learning, GAP's new architecture is composed of three separate modules: (i) the encoder module, where we learn private node embeddings without relying on the edge information; (ii) the aggregation module, where we compute noisy aggregated node embeddings based on the graph structure; and (iii) the classification module, where we train a neural network on the private aggregations for node classification without further querying the graph edges. GAP's major advantage over previous approaches is that it can benefit from multi-hop neighborhood aggregations, and guarantees both edge-level and node-level DP not only for training, but also at inference with no additional costs beyond the training's privacy budget. We analyze GAP's formal privacy guarantees using Rényi DP and conduct empirical experiments over three real-world graph datasets. We demonstrate that GAP offers significantly better accuracy-privacy trade-offs than state-of-the-art DP-GNN approaches and naive MLP-based baselines. Our code is publicly available at github.com/sisaman/GAP.

Differentially Private Coordinate Descent for Composite Empirical Risk Minimization, [29]

Machine learning models can leak information about the data used to train them. To mitigate this issue, Differentially Private (DP) variants of optimization algorithms like Stochastic Gradient Descent (DP-SGD) have been designed to trade-off utility for privacy in Empirical Risk Minimization (ERM) problems. In this paper, we propose Differentially Private proximal Coordinate Descent (DP-CD), a new method to solve composite DP-ERM problems. We derive utility guarantees through a novel theoretical analysis of inexact coordinate descent. Our results show that, thanks to larger step sizes, DP-CD can exploit imbalance in gradient coordinates to outperform DP-SGD. We also prove new lower bounds for composite DP-ERM under coordinate-wise regularity assumptions, that are nearly matched by DP-CD. For practical implementations, we propose to clip gradients using coordinate-wise thresholds that emerge from our theory, avoiding costly hyperparameter tuning. Experiments on real and synthetic data support our results, and show that DP-CD compares favorably with DP-SGD.

Privacy Amplification by Decentralization, [24]

Analyzing data owned by several parties while achieving a good trade-off between utility and privacy is a key challenge in federated learning and analytics. In this work, we introduce a novel relaxation of local differential privacy (LDP) that naturally arises in fully decentralized algorithms, i.e., when participants exchange information by communicating along the edges of a network graph without central coordinator. This relaxation, that we call network DP, captures the fact that users have only a local view of the system. To show the relevance of network DP, we study a decentralized model of computation where a token performs a walk on the network graph and is updated sequentially by the party who receives it. For tasks such as real summation, histogram computation and optimization with gradient descent, we propose simple algorithms on ring and complete topologies. We prove that the privacy-utility trade-offs of our algorithms under network DP significantly improve upon what is achievable under LDP, and often match the utility of the trusted curator model. Our results show for the first time that formal privacy gains can

be obtained from full decentralization. We also provide experiments to illustrate the improved utility of our approach for decentralized training with stochastic gradient descent.

High-Dimensional Private Empirical Risk Minimization by Greedy Coordinate Descent, [45]

In this paper, we study differentially private empirical risk minimization (DP-ERM). It has been shown that the worst-case utility of DP-ERM reduces polynomially as the dimension increases. This is a major obstacle to privately learning large machine learning models. In high dimension, it is common for some model's parameters to carry more information than others. To exploit this, we propose a differentially private greedy coordinate descent (DP-GCD) algorithm. At each iteration, DP-GCD privately performs a coordinate-wise gradient step along the gradients' (approximately) greatest entry. We show theoretically that DP-GCD can achieve a logarithmic dependence on the dimension for a wide range of problems by naturally exploiting their structural properties (such as quasi-sparse solutions). We illustrate this behavior numerically, both on synthetic and real datasets.

Private Sampling with Identifiable Cheaters, [17]

In this paper we study verifiable sampling from probability distributions in the context of multi-party computation. This has various applications in randomized algorithms performed collaboratively by parties not trusting each other. One example is differentially private machine learning where noise should be drawn, typically from a Laplace or Gaussian distribution, and it is desirable that no party can bias this process. In particular, we propose algorithms to draw random numbers from uniform, Laplace, Gaussian and arbitrary probability distributions, and to verify honest execution of the protocols through zero-knowledge proofs. We propose protocols that result in one party knowing the drawn number and protocols that deliver the drawn random number as a shared secret.

Muffliato: Peer-to-Peer Privacy Amplification for Decentralized Optimization and Averaging, [25]

Decentralized optimization is increasingly popular in machine learning for its scalability and efficiency. Intuitively, it should also provide better privacy guarantees, as nodes only observe the messages sent by their neighbors in the network graph. But formalizing and quantifying this gain is challenging: existing results are typically limited to Local Differential Privacy (LDP) guarantees that overlook the advantages of decentralization. In this work, we introduce pairwise network differential privacy, a relaxation of LDP that captures the fact that the privacy leakage from a node u to a node v may depend on their relative position in the graph. We then analyze the combination of local noise injection with (simple or randomized) gossip averaging protocols on fixed and random communication graphs. We also derive a differentially private decentralized optimization algorithm that alternates between local gradient descent steps and gossip averaging. Our results show that our algorithms amplify privacy guarantees as a function of the distance between nodes in the graph, matching the privacy-utility trade-off of the trusted curator, up to factors that explicitly depend on the graph topology. Finally, we illustrate our privacy gains with experiments on synthetic and real-world datasets.

8.4 Federated Learning in Medicine

FLamby: Datasets and Benchmarks for Cross-Silo Federated Learning in Realistic Healthcare Settings, [35]

Federated Learning (FL) is a novel approach enabling several clients holding sensitive data to collaboratively train machine learning models, without centralizing data. The cross-silo FL setting corresponds to the case of few (2–50) reliable clients, each holding medium to large datasets, and is typically found in applications such as healthcare, finance, or industry. While previous works have proposed representative datasets for cross-device FL, few realistic healthcare cross-silo FL datasets exist, thereby slowing algorithmic research in this critical application. In this work, we propose a novel cross-silo dataset suite focused on healthcare, FLamby (Federated Learning AMple Benchmark of Your cross-silo strategies), to bridge the gap between theory and practice of cross-silo FL. FLamby encompasses 7 healthcare datasets with natural splits, covering multiple tasks, modalities, and data volumes, each accompanied with baseline

training code. As an illustration, we additionally benchmark standard FL algorithms on all datasets. Our flexible and modular suite allows researchers to easily download datasets, reproduce results and re-use the different components for their research. FLamby is available at www.github.com/owkin/flamby.

Passive Query-Recovery Attack Against Secure Conjunctive Keyword Search Schemes, [39]

While storing documents on the cloud can be attractive, the question remains whether cloud providers can be trusted with storing private documents. Even if trusted, data breaches are ubiquitous. To prevent information leakage one can store documents encrypted. If encrypted under traditional schemes, one loses the ability to perform simple operations over the documents, such as searching through them. Searchable encryption schemes were proposed allowing some search functionality while documents remain encrypted. Orthogonally, research is done to find attacks that exploit search and access pattern leakage that most efficient schemes have. One type of such an attack is the ability to recover plaintext queries. Passive query-recovery attacks on single-keyword search schemes have been proposed in literature, however, conjunctive keyword search has not been considered, although keyword searches with two or three keywords appear more frequently in online searches. We introduce a generic extension strategy for existing passive query-recovery attacks against single-keyword search schemes and explore its applicability for the attack presented by Damie et al. (USENIX Security '21). While the original attack achieves up to a recovery rate of 85% against single-keyword search schemes for an attacker without exact background knowledge, our experiments show that the generic extension to conjunctive queries comes with a significant performance decrease achieving recovery rates of at most 32%. Assuming a stronger attacker with partial knowledge of the indexed document set boosts the recovery rate to 85% for conjunctive keyword queries with two keywords and achieves similar recovery rates as previous attacks by Cash et al. (CCS '15) and Islam et al. (NDSS '12) in the same setting for single-keyword search schemes.

8.5 Learning and Speech Recognition

Enhancing speech privacy with slicing, [30]

Privacy preservation calls for speech anonymization methods which hide the speaker's identity while minimizing the impact on downstream tasks such as automatic speech recognition (ASR) training or decoding. In the recent VoicePrivacy 2020 Challenge, several anonymization methods have been proposed to transform speech utterances in a way that preserves their verbal and prosodic contents while reducing the accuracy of a speaker verification system. In this paper, we propose to further increase the privacy achieved by such methods by segmenting the utterances into shorter slices. We show that our approach has two major impacts on privacy. First, it reduces the accuracy of speaker verification with respect to unsegmented utterances. Second, it also reduces the amount of personal information that can be extracted from the verbal content, in a way that cannot easily be reversed by an attacker. We also show that it is possible to train an ASR system from anonymized speech slices with negligible impact on the word error rate.

Privacy and utility of x-vector based speaker anonymization, [19]

We study the scenario where individuals (speakers) contribute to the publication of an anonymized speech corpus. Data users then leverage this public corpus to perform downstream tasks (such as training automatic speech recognition systems), while attackers may try to de-anonymize it based on auxiliary knowledge they collect. Motivated by this scenario, speaker anonymization aims to conceal the speaker identity while preserving the quality and usefulness of speech data. In this paper, we study x-vector based speaker anonymization, the leading approach in the recent Voice Privacy Challenge, which converts an input utterance into that of a random pseudo-speaker. We show that the strength of the anonymization varies significantly depending on how the pseudo-speaker is selected. In particular, we investigate four design choices: the distance measure between speakers, the region of x-vector space where the pseudo-speaker is mapped, the gender selection and whether to use speaker or utterance level assignment. We assess the quality of anonymization from the perspective of the three actors involved in our threat model, namely the speaker, the user and the attacker. To measure privacy and utility, we use respectively the linkability score achieved by the attackers and the decoding word error rate incurred by

an ASR model trained with the anonymized data. Experiments on LibriSpeech dataset confirm that the optimal combination of design choices yield state-of-the-art performance in terms of privacy protection as well as utility. Experiments on Mozilla Common Voice dataset show that the best design choices with 50 speakers guarantee the same anonymization level against re-identification attack as raw speech with 20,000 speakers.

Retrieving Speaker Information from Personalized Acoustic Models for Speech Recognition, [31, 37]

The widespread of powerful personal devices capable of collecting voice of their users has opened the opportunity to build speaker adapted speech recognition system (ASR) or to participate to collaborative learning of ASR. In both cases, personalized acoustic models (AM), i.e. fine-tuned AM with specific speaker data, can be built. A question that naturally arises is whether the dissemination of personalized acoustic models can leak personal information. In this paper, we show that it is possible to retrieve the gender of the speaker, but also his identity, by just exploiting the weight matrix changes of a neural acoustic model locally adapted to this speaker. Incidentally we observe phenomena that may be useful towards explainability of deep neural networks in the context of speech processing. Gender can be identified almost surely using only the first layers and speaker verification performs well when using middle-up layers. Our experimental study on the TED-LIUM 3 dataset with HMM/TDNN models shows an accuracy of 95% for gender detection, and an Equal Error Rate of 9.07% for a speaker verification task by only exploiting the weights from personalized models that could be exchanged instead of user data.

Differentially private speaker anonymization, [18]

Sharing real-world speech utterances is key to the training and deployment of voice-based services. However, it also raises privacy risks as speech contains a wealth of personal data. Speaker anonymization aims to remove speaker information from a speech utterance while leaving its linguistic and prosodic attributes intact. State-of-the-art techniques operate by disentangling the speaker information (represented via a speaker embedding) from these attributes and re-synthesizing speech based on the speaker embedding of another speaker. Prior research in the privacy community has shown that anonymization often provides brittle privacy protection, even less so any provable guarantee. In this work, we show that disentanglement is indeed not perfect: linguistic and prosodic attributes still contain speaker information. We remove speaker information from these attributes by introducing differentially private feature extractors based on an autoencoder and an automatic speech recognizer, respectively, trained using noise layers. We plug these extractors in the state-of-the-art anonymization pipeline and generate, for the first time, differentially private utterances with a provable upper bound on the speaker information they contain. We evaluate empirically the privacy and utility resulting from our differentially private speaker anonymization approach on the LibriSpeech data set. Experimental results show that the generated utterances retain very high utility for automatic speech recognition training and inference, while being much better protected against strong adversaries who leverage the full knowledge of the anonymization process to try to infer the speaker identity.

The VoicePrivacy 2020 Challenge: Results and findings, [20] and [47]

This paper presents the results and analyses stemming from the first VoicePrivacy 2020 Challenge which focuses on developing anonymization solutions for speech technology. We provide a systematic overview of the challenge design with an analysis of submitted systems and evaluation results. In particular, we describe the voice anonymization task and datasets used for system development and evaluation. Also, we present different attack models and the associated objective and subjective evaluation metrics. We introduce two anonymization baselines and provide a summary description of the anonymization systems developed by the challenge participants. We report objective and subjective evaluation results for baseline and submitted systems. In addition, we present experimental results for alternative privacy metrics and attack models developed as a part of the post-evaluation analysis. Finally, we summarise our insights and observations that will influence the design of the next VoicePrivacy challenge edition and some directions for future voice anonymization research.

Privacy attacks for automatic speech recognition acoustic models in a federated learning framework, [36, 38]

This paper investigates methods to effectively retrieve speaker information from the personalized speaker adapted neural network acoustic models (AMs) in automatic speech recognition (ASR). This problem is especially important in the context of federated learning of ASR acoustic models where a global model is learnt on the server based on the updates received from multiple clients. We propose an approach to analyze information in neural network AMs based on a neural network footprint on the so-called Indicator dataset. Using this method, we develop two attack models that aim to infer speaker identity from the updated personalized models without access to the actual users' speech data. Experiments on the TED-LIUM 3 corpus demonstrate that the proposed approaches are very effective and can provide equal error rate (EER) of 1-2%.

8.6 Fairness and Transparency

Fairness Certificates for Differentially Private Classification, [46]

In this work, we theoretically study the impact of differential privacy on fairness in binary classification. We prove that, given a class of models, popular group fairness measures are pointwise Lipschitz-continuous with respect to the parameters of the model. This result is a consequence of a more general statement on the probability that a decision function makes a negative prediction conditioned on an arbitrary event (such as membership to a sensitive group), which may be of independent interest. We use the aforementioned Lipschitz property to prove a high probability bound showing that, given enough examples, the fairness level of private models is close to the one of their non-private counterparts.

FairGrad: Fairness Aware Gradient Descent, [43]

We tackle the problem of group fairness in classification, where the objective is to learn models that do not unjustly discriminate against subgroups of the population. Most existing approaches are limited to simple binary tasks or involve difficult to implement training mechanisms. This reduces their practical applicability. In this paper, we propose FairGrad, a method to enforce fairness based on a reweighting scheme that iteratively learns group specific weights based on whether they are advantaged or not. FairGrad is easy to implement and can accommodate various standard fairness definitions. Furthermore, we show that it is comparable to standard baselines over various datasets including ones used in natural language processing and computer vision.

8.7 Graph-Based Learning

Linear Programs with Conjunctive Queries, [23]

In this paper, we study the problem of optimizing a linear program whose variables are answers to a conjunctive query. For this we propose the language LP(CQ) for specifying linear programs whose constraints and objective functions depend on the answer sets of conjunctive queries. We contribute an efficient algorithm for solving programs in a fragment of LP(CQ). The naive approach constructs a linear program having as many variables as elements in the answer set of the queries. Our approach constructs a linear program having the same optimal value but fewer variables. This is done by exploiting the structure of the conjunctive queries using hypertree decompositions of small width to group elements of the answer set together. We illustrate the various applications of LP(CQ) programs on three examples: optimizing deliveries of resources, minimizing noise for differential privacy, and computing the s-measure of patterns in graphs as needed for datamining.

A Revenue Function for Comparison-Based Hierarchical Clustering, [44]

Comparison-based learning addresses the problem of learning when, instead of explicit features or pairwise similarities, one only has access to comparisons of the form: *Object A is more similar to B than to C*. Recently, it has been shown that, in Hierarchical Clustering, single and complete linkage can be directly implemented using only such comparisons while several algorithms have been proposed to emulate the

behaviour of average linkage. Hence, finding hierarchies (or dendrograms) using only comparisons is a well understood problem. However, evaluating their meaningfulness when no ground-truth nor explicit similarities are available remains an open question. In this paper, we bridge this gap by proposing a new revenue function that allows one to measure the goodness of dendrograms using only comparisons. We show that this function is closely related to Dasgupta's cost for hierarchical clustering that uses pairwise similarities. On the theoretical side, we use the proposed revenue function to resolve the open problem of whether one can approximately recover a latent hierarchy using few triplet comparisons. On the practical side, we present principled algorithms for comparison-based hierarchical clustering based on the maximisation of the revenue and we empirically compare them with existing methods.

Limits of Multi-relational Graphs, [13]

Graphons are limits of large graphs. Motivated by a theoretical problem from statistical relational learning, we develop a generalization of basic results from graphon theory into the "multi-relational" setting. We show that their multi-relational counterparts, which we call multi-relational graphons, are analogically limits of large multi-relational graphs. We extend the cutdistance topology for graphons to multi-relational graphons and prove its compactness and the density of multi-relational graphs in this topology. In turn, compactness enables to prove the large deviation principle for Multi-Relational Graphs (LDP) which enables to prove the most typical random graphs constrained by marginal statistics converge asymptotically to constrained multirelational graphons with maximum entropy. We show the equivalence between a restricted version of Markov Logic Network and Multi-Relational Graphons with maximum entropy.

Machine Learning

Optimal Tensor Transport, [27]

Optimal Transport (OT) has become a popular tool in machine learning to align finite datasets typically lying in the same vector space. To expand the range of possible applications, Co-Optimal Transport (Co-OT) jointly estimates two distinct transport plans, one for the rows (points) and one for the columns (features), to match two data matrices that might use different features. On the other hand, Gromov Wasserstein (GW) looks for a single transport plan from two pairwise intra-domain distance matrices. Both Co-OT and GW can be seen as specific extensions of OT to more complex data. In this paper, we propose a unified framework, called Optimal Tensor Transport (OTT), which takes the form of a generic formulation that encompasses OT, GW and CoOT and can handle tensors of any order by learning possibly multiple transport plans. We derive theoretical results for the resulting new distance and present an efficient way for computing it. We further illustrate the interest of such a formulation in Domain Adaptation and Comparison-based Clustering.

Robust Kernel Density Estimation with Median-of-Means principle, [26]

In this paper, we introduce a robust nonparametric density estimator combining the popular Kernel Density Estimation method and the Median-of-Means principle (MoM-KDE). This estimator is shown to achieve robustness to any kind of anomalous data, even in the case of adversarial contamination. In particular, while previous works only prove consistency results under known contamination model, this work provides finite-sample high-probability error-bounds without a priori knowledge on the outliers. Finally, when compared with other robust kernel estimators, we show that MoM-KDE achieves competitive results while having significant lower computational complexity.

9 Partnerships and cooperations

9.1 International initiatives

9.1.1 Associate Teams in the framework of an Inria International Lab or in the framework of an Inria International Program

LEGO

Participants: Aurélien Bellet (*contact person*), Pascal Denis.

Title: LEarning GOod representations for natural language processing

Duration: 2016 – 2022

Coordinator: Fei Sha (feisha@usc.edu)

Partners:

- University of Southern California

Summary: The proposed research lies in the intersection of Machine Learning and Natural Language Processing (NLP), addressing the following core problems in both communities: what are the right representations for structured data and how to learn them automatically? And how to apply such representations to complex and structured prediction tasks in NLP?

9.1.2 Participation in other International Programs**SLANT: Bilateral ANR project with Luxembourg**

Participants: Pascal Denis (*contact person*), Aurélien Bellet, Mikaela Keller, Gau-rav Maheshwari.

Acronym: SLANT

Title: Spin and bias in Language Analyzed in News and Texts

Duration: December 2019 – June 2023

Coordinator: Philippe Muller, IRIT, Toulouse

Partners: IRIT (Toulouse), SnT (Luxembourg)

Abstract: There is a growing concern about misinformation or biased information in public communication, whether in traditional media or social forums. While automating fact-checking has received a lot of attention, the problem of fair information is much larger and includes more insidious forms like biased presentation of events and discussion. The SLANT project aims at characterizing bias in textual data, either intended, in public reporting, or unintended in writing aiming at neutrality. An abstract model of biased interpretation using work on discourse structure, semantics and interpretation will be complemented and concretized by finding relevant lexical, syntactic, stylistic or rhetorical differences through an automated but explainable comparison of texts with different biases on the same subject, based on a dataset of news media coverage from a diverse set of sources. We will also explore how our results can help alter bias in texts or remove it from automated representations of texts.

IMPRESS: Bilateral Inria-DFKI project

Participants: Pascal Denis (*contact person*), Rémi Gilleron, Priyansh Trivedi.

Acronym: IMPRESS

Title: Improving Embeddings with Semantic Knowledge

Duration: Oct 2020-Sept 2023

Coordinator: Pascal Denis and Ivana Kruijff-Korabayova (DFKI)

Partners: Sémagramme (Inria Nancy), DFKI (Germany)

Abstract: Virtually all NLP systems nowadays use vector representations of words, a.k.a. word embeddings. Similarly, the processing of language combined with vision or other sensory modalities employs multimodal embeddings. While embeddings do embody some form of semantic relatedness, the exact nature of the latter remains unclear. This loss of precise semantic information can affect downstream tasks. Furthermore, while there is a growing body of NLP research on languages other than English, most research on multimodal embeddings is still done on English. The goals of IMPRESS are to investigate the integration of semantic knowledge into embeddings and its impact on selected downstream tasks, to extend this approach to multimodal and mildly multilingual settings, and to develop open source software and lexical resources, focusing on video activity recognition as a practical testbed.

9.2 European initiatives

9.2.1 Horizon Europe

TRUMPET [TRUMPET project on cordis.europa.eu](https://cordis.europa.eu/trumpet)

Participants: Jan Ramon (*contact person*), Aurélien Bellet, Marc Tommasi, Aleksei Korneev, Sophie Villerot.

Title: TRUstworthy Multi-site Privacy Enhancing Technologies

Duration: From October 1, 2022 to September 30, 2025

Partners:

- Institut National de Recherche en Informatique et Automatique (INRIA), France
- TIME.LEX (time.lex), Belgium
- Technovative Solutions LTD, United Kingdom
- Fundacion Centro Tecnoloxico de Telecomunicacions de Galicia (GRADIANT), Spain
- Commissariat à l'Énergie Atomique et aux Energies Alternatives (CEA), France
- Instituto Romagnolo per lo Studio dei Tumori Dino Amadori - IRST SRL (IRST), Italy
- Centre Hospitalier Universitaire de Liege (CHUL), Belgium
- Universidad de Vigo (UVIGO), Spain
- Arteevo Technologies LTD (ARTEEVO), Israel

Coordinator: Ruth Muleiro Alonso (GRAD)

Summary In recent years, Federated Learning (FL) has emerged as a revolutionary privacy-enhancing technology and, consequently, has quickly expanded to other applications.

However, further research has cast a shadow of doubt on the strength of privacy protection provided by FL. Potential vulnerabilities and threats pointed out by researchers included a curious aggregator threat; susceptibility to man-in-the-middle and insider attacks that disrupt the convergence of global and local models or cause convergence to fake minima; and, most importantly, inference attacks that aim to re-identify data subjects from FL's AI model parameter updates.

The goal of TRUMPET is to research and develop novel privacy enhancement methods for Federated Learning, and to deliver a highly scalable Federated AI service platform for researchers, that will enable AI-powered studies of siloed, multi-site, cross-domain, cross border European datasets with

privacy guarantees that exceed the requirements of GDPR. The generic TRUMPET platform will be piloted, demonstrated and validated in the specific use case of European cancer hospitals, allowing researchers and policymakers to extract AI-driven insights from previously inaccessible cross-border, cross-organization cancer data, while ensuring the patients' privacy. The strong privacy protection accorded by the platform will be verified through the engagement of external experts for independent privacy leakage and re-identification testing.

A secondary goal is to research, develop and promote with EU data protection authorities a novel metric and tool for the certification of GDPR compliance of FL implementations.

The consortium is composed of 9 interdisciplinary partners: 3 Research Organizations, 1 University, 3 SMEs and 2 Clinical partners with extensive experience and expertise to guarantee the correct performance of the activities and the achievement of the results.

9.3 National initiatives

9.3.1 ANR Pamela (2016–2022)

Participants: Marc Tommasi (*contact person*), Aurélien Bellet, Rémi Gilleron, Jan Ramon, Mahsa Asadi.

The Pamela project aims at developing machine learning theories and algorithms in order to learn local and personalized models from data distributed over networked infrastructures. Our project seeks to provide first answers to modern information systems built by interconnecting many personal devices holding private user data in the search of personalized suggestions and recommendations. More precisely, we will focus on learning in a collaborative way with the help of neighbors in a network. We aim to lay the first blocks of a scientific foundation for these new types of systems, in effect moving from graphs of data to graphs of data and learned models. We argue that this shift is necessary in order to address the new constraints arising from the decentralization of information that is inherent to the emergence of big data. We will in particular focus on the question of learning under communication and privacy constraints. A significant asset of the project is the quality of its industrial partners, Snips and Mediego, who bring in their expertise in privacy protection and distributed computing as well as use cases and datasets. They will contribute to translate this fundamental research effort into concrete outcomes by developing personalized and privacy-aware assistants able to provide contextualized recommendations on small devices and smartphones.

9.3.2 ANR DEEP-Privacy (2019–2023)

Participants: Marc Tommasi (*contact person*), Aurélien Bellet, Pascal Denis, Jan Ramon, Brij Mohan Lal Srivastava, Rishabh Gupta.

DEEP-PRIVACY proposes a new paradigm based on a distributed, personalized, and privacy-preserving approach for speech processing, with a focus on machine learning algorithms for speech recognition. To this end, we propose to rely on a hybrid approach: the device of each user does not share its raw speech data and runs some private computations locally, while some cross-user computations are done by communicating through a server (or a peer-to-peer network). To satisfy privacy requirements at the acoustic level, the information communicated to the server should not expose sensitive speaker information.

9.3.3 ANR-JCJC PRIDE (2020–2025)

Participants: Aurélien Bellet (*contact person*), Marc Tommasi, Jan Ramon, Edwige Cyffers, Batiste Le Bars, Paul Mangold, Tudor Cebere.

Machine learning (ML) is ubiquitous in AI-based services and data-oriented scientific fields but raises serious privacy concerns when training on personal data. The starting point of PRIDE is that personal data should belong to the individual who produces it. This requires to revisit ML algorithms to learn from many decentralized personal datasets while preventing the reconstruction of raw data. Differential Privacy (DP) provides a strong notion of protection, but current decentralized ML algorithms are not able to learn useful models under DP. The goal of PRIDE is to develop theoretical and algorithmic tools that enable differentially-private ML methods operating on decentralized datasets, through two complementary objectives: (1) prove that gossip protocols naturally reinforce DP guarantees; (2) propose algorithms at the intersection of decentralized ML and secure multi-party computation.

9.3.4 ANR PMR (2020-2024)

Participants: Jan Ramon (*contact person*), Aurélien Bellet, Marc Tommasi, Cesar Sabater.

Given the growing awareness of privacy risks of data processing, there is an increasing interest in privacy-preserving learning. However, shortcomings in the state of the art limit the applicability of the privacy-preserving learning paradigm. First, most approaches assume too optimistically a honest-but-curious setting. Second, most approaches consider one learning task in isolation, not accounting for the context where querying is a recurring activity. We will investigate new algorithms and models that address these shortcomings. Among others, (i) our algorithms will combine privacy-preserving properties of differential privacy with security offered by cryptography and (ii) based on models of information flows in integrated data handling processes, we will build more refined models analyzing the implications of repeated querying. We will demonstrate the utility of our new theory and algorithms by proposing strategies to realistically apply them in significant real-world problems illustrated through use cases in the medical domain

9.3.5 FedMalin. INRIA Defi (2021-2024)

Participants: Aurélien Bellet (*contact person*), Jan Ramon, Marc Tommasi, Michaël Perrot, Batiste Le Bars, Edwige Cyffers, Paul Mangold, Tudor Cebere.

In many use-cases of Machine Learning (ML), data is naturally decentralized: medical data is collected and stored by different hospitals, crowdsensed data is generated by personal devices, etc. Federated Learning (FL) has recently emerged as a novel paradigm where a set of entities with local datasets collaboratively train ML models while keeping their data decentralized.

FedMalin is a research project that spans 10 Inria research teams and aims to push FL research and concrete use-cases through a multidisciplinary consortium involving expertise in ML, distributed systems, privacy and security, networks, and medicine. We propose to address a number of challenges that arise when FL is deployed over the Internet, including privacy and fairness, energy consumption, personalization, and location/time dependencies.

FedMalin will also contribute to the development of open-source tools for FL experimentation and real-world deployments, and use them for concrete applications in medicine and crowdsensing.

9.3.6 FLAMED: Federated Learning and Analytics on Medical Data. INRIA Action Exploratoire (2020-2024)

Participants: Aurélien Bellet (*contact person*), Marc Tommasi, Paul Mangold, Nathan Bigaud, Paul André.

Flamed is about decentralized approaches for AI in medicine. The main objective is to operate data analysis and machine learning tasks in a network of hospital units without any data exchange. This approach helps to solve data privacy and sovereignty issues while taking advantage of the statistical power of federation and collaboration. This research is done in collaboration with the Lille Hospital.

9.3.7 COMANCHE: Computational Models of Lexical Meaning and Change. INRIA Action Exploratoire (2022-2026)

Participants: Pascal Denis (*contact person*), Mikaela Keller, Bastien Liétard.

Comanche proposes to transfer and adapt recent Natural Language representation learning algorithms from deep learning to model the evolution of the meaning of words, and to confront these computational models to theories on language acquisition and the diachrony of languages. At the crossroads between machine learning, psycholinguistics and historical linguistics, this project will make it possible to validate or revise some of these theories, but also to bring out computational models that are more sober in terms of data and computations because they exploit new inductive biases inspired by these disciplines.

In collaboration with UMR SCALAB (CNRS, Université de Lille), l'Unité de Recherche STIH (Sorbonne Université), et l'UMR ATILF (CNRS, Université de Lorraine).

9.3.8 IPoP, Projet interdisciplinaire sur la protection des données personnelles, PEPR Cybersécurité (2022-2028).

Participants: Aurélien Bellet (*contact person*), Jan Ramon, Marc Tommasi, Michaël Perrot, Cesar Sabater, Edwige Cyffers, Paul Mangold, Tudor Cebere.

Digital technologies provide services which can greatly increase quality of life (e.g. connected e-health devices, location based services, or personal assistants). However, these services can also raise major privacy risks, as they involve personal data, or even sensitive data. Indeed, this notion of personal data is the cornerstone of French and European regulations, since processing such data triggers a series of obligations that the data controller must abide by. This raises many multidisciplinary issues, as the challenges are not only technological, but also societal, judiciary, economic, political and ethical.

The objectives of this project are thus to study the threats on privacy that have been introduced by these new services, and to conceive theoretical and technical privacy-preserving solutions that are compatible with French and European regulations, that preserve the quality of experience of the users. These solutions will be deployed and assessed, both on the technological and legal sides, and on their societal acceptability. In order to achieve these objectives, we adopt an interdisciplinary approach, bringing together many diverse fields: computer science, technology, engineering, social sciences, economy and law.

The project's scientific program focuses on new forms of personal information collection, on Artificial Intelligence (AI) and its governance, data anonymization techniques, personal data management and distributed calculation protocol privacy preserving infrastructures, differential privacy, personal data legal protection and compliance, and all the associated societal and ethical considerations. This unifying interdisciplinary research program brings together internationally recognized research teams (from universities, engineering schools and institutions) working on privacy, and the French Data Protection Authority (CNIL).

This holistic vision of the issues linked to personal data protection will on the one hand let us propose solutions to the scientific and technological challenges and on the other help us confront these solutions in many different ways, in the context of interdisciplinary collaborations, thus leading to recommendations and proposals in the field of regulations or legal frameworks. This comprehensive consideration of all the issues aims at encouraging the adoption and acceptability of the solutions

proposed by all stakeholders, legislators, data controllers, data processors, solution designers, developers all the way to end-users.

9.3.9 HyAIAI. INRIA Defi (2019-2022)

Participants: Jan Ramon (*contact person*), Marc Tommasi.

HyAIAI is an Inria Defi about the design of novel, interpretable approaches for Artificial Intelligence.

Recent progress in Machine Learning (ML) and especially Deep Learning has made ML pervasive in a wide range of applications. However, current approaches rely on complex numerical models: their decisions, as accurate as they may be, cannot be easily explained to the layman that may depend on these decisions (ex: get a loan or not). In the HyAIAI IPL, we tackle the problem of making “Interpretable ML” through the study and design of hybrid approaches that combine state of the art numeric models with explainable symbolic models. More precisely, our goal is to be able to integrate high level (domain) constraints in ML models, to give model designers information on ill-performing parts of the model, and to give the layman/practitioner understandable explanations on the results of the ML model.

9.4 Regional initiatives

9.4.1 Chaire TIP: Transparent artificial Intelligence preserving Privacy

Participants: Jan Ramon (*contact person*), Cesar Sabater, Jean-Paul Lam, Antonin Duez, Sophie Villerot.

While AI techniques are becoming ever more powerful, there is a growing concern about potential risks and abuses. As a result, there has been an increasing interest in research directions such as privacy-preserving machine learning, explainable machine learning, fairness and data protection legislation.

The overall goal of the TIP project is to develop, exploit and explain a sound understanding of privacy-preserving strategies in larger AI-based processes involving massive numbers of agents among whom part may be malicious.

9.4.2 Fairness in Decentralized and Privacy Preserving Machine Learning STaRS (2021-2023)

Participants: Michaël Perrot (*contact person*).

Machine Learning is becoming ubiquitous in our everyday lives. It is now used in digital assistants, for medical diagnosis, for autonomous vehicles, Its success can be explained by the good performances of learned models, sometimes reaching human-level capabilities. However, simply being accurate is not sufficient to ensure that the learning approaches are socially acceptable, in particular if the models are to be largely deployed. Hence, Fairness and Privacy have been extensively studied as standalone trustworthiness notions. However, in practice, it is often mandatory that a model has both properties and thus, jointly studying the two notions is important. This is particularly relevant in decentralized settings where the data is owned by multiple entities that would like to collaborate to learn efficient and fair models but wish to keep their own data private. The goal of this project is twofold: (i) propose new approaches to learn fair and privacy preserving models in a decentralized setting and (ii) provide theoretical guarantees on the trustworthiness level of the learned models that may serve as certificates for the stakeholders.

9.4.3 CAPS'UL: CAmпус Participatif en Santé numérique du site Universitaire de Lille. Santé numérique, PIA4. (2022-2027)

Participant: Marc Tommasi (*contact person*).

Participation to this regional initiative about training in the health domain. Magnet contributes by the design of algorithms for data anonymization and data generation for specific software training.

10 Dissemination

10.1 Promoting scientific activities

10.1.1 Scientific events: organisation

- AURÉLIEN BELLET co-organized the [Privacy Preserving Machine Learning \(PPML'22\) workshop](#) at FOCS 2022.
- Since May 2020, AURÉLIEN BELLET is co-organizing the [Federated Learning One World webinar](#) (900+ registered attendees).

10.1.2 Scientific events: selection

- AURÉLIEN BELLET served as Area Chair for ICML'22, NeurIPS'22, AISTATS'23 and AFCP@NeurIPS'22, and as PC member for SaTML, PPAI@AAAI'23, FL@NeurIPS'22, CAP'22, APVP'22.
- JAN RAMON was PC member of AAAI'22, AFCE@Neurips'22, AISTATS'22, BigData'22, BNAIC'22, DS'22, ECMLPKDD'22, PhDforum@ECML'22, ICDM'22, ICLR'22, ICML'22, IDA'22, IJCAI'22, doctoral program @IJCAI'22, MLG@KDD'22, MLG@ECMLPKDD'22, NeurIPS'22, PPAI@AAAI'22, SDM'22, UAI'22 and XKDD'22.
- PASCAL DENIS was Action Editor for ACL Rolling Review, as well as Area Chair for ACL'22 and NAACL'22. He was also PC member of COLING'22, EMNLP'22, and CODI@COLING'22.
- MARC TOMMASI served as Area Chair for ECML'22 and as PC member of NeurIPS'22, CAP'22.
- MICHAËL PERROT served as PC member of ICML'22, ECML'22, NeurIPS'22, AISTATS'23.
- MIKAELA KELLER was PC member of TALN'22
- RÉMI GILLERON served as PC member of ICML'22, ICLR'22, CAP'22, NeurIPS'22, AISTATS'22.
- MARC DAMIE served as reviewer for PETS'22.
- PAUL MANGOLD served as reviewer for ECML'22 PhD forum, NeurIPS'22, AISTATS'22.

10.1.3 Journal

Member of the editorial boards

- AURÉLIEN BELLET is Action Editor for Transactions of Machine Learning Research (TMLR).
- JAN RAMON is member of the editorial boards of Machine Learning Journal (MLJ), Data Mining and Knowledge Discovery (DMKD), Journal of Machine Learning Research (JMLR), ECML-PKDD Journal track. JAN RAMON is action editor of Data Mining and Knowledge Discovery (DMKD).
- PASCAL DENIS is standing reviewer for Transactions of the Association for Computational Linguistics (TACL).

Reviewer - reviewing activities

- MIKAELA KELLER reviewed for Scientific Data

10.1.4 Invited talks

- AURÉLIEN BELLET gave invited talks at [Google Federated Learning and Analytics Workshop 2022](#), [CAP'22](#) (keynote speaker), [JDS'22](#) (keynote speaker), DeepMind Paris, [Learning and Optimization in Luminy \(LOL2022\)](#), [OpenMined Medical Federated Learning Program](#), [Privaski'22](#), [Journée Private Math-Stat'22](#), [Comète Workshop on Ethical AI](#), [Qarma seminar](#), [Séminaire Scientifique STIC AmSud](#), [GDR IA](#), [Qwant seminar](#), [AXA seminar](#), [Idemia seminar](#), [Seminar "Machine Learning in Montpellier"](#), [Fairness in Machine Learning Day at Lille](#)
- MICHAËL PERROT presented his work of Fairness in Machine Learning during the [Comète Workshop on Ethical AI](#) in September, 2022.
- BATISTE LE BARS presented his work [42] at the [EPFL-INRIA workshop](#), in the MILES team of the Lamsade lab in PSL and in the [Federated Learning Workshop at NeurIPS'22](#).
- MARC DAMIE gave a talk for the Dutch Cyber Security Best Research Paper award ceremony (paper accepted at USENIX Security 21).

10.1.5 Scientific expertise

- AURÉLIEN BELLET reviewed for the European Research Council (ERC - Starting Grants), and was a member of the recruitment committee of junior researchers at Inria Lyon and of associate professors at Saint-Etienne.
- JAN RAMON made reviews for the European Commission (review final report of 1 project), and was a member of the comité emploi recherche (CER) of INRIA-Lille.
- MARC TOMMASI was a member (scientific expert) of the recruitment committee of associate professors at Saint-Etienne and Lille.
- MIKAELA KELLER was a member (scientific expert) of the recruitment committee of associate professors at Saint-Etienne and Lille.

10.1.6 Research administration

- AURÉLIEN BELLET is member of the Operational Committee for the assessment of Legal and Ethical risks (COERLE).
- MARC TOMMASI is co-head of the DatInG group (4 teams, about 100 persons), member of the Conseil Scientifique du laboratoire CRISAL and member of the Commission mixte CRISAL/Faculty of Science, Lille University.
- PASCAL DENIS is a standing member of TALN-RECITAL permanent conference committee ("CPerm").
- PASCAL DENIS is a member of the CNRS GDR NLP Group.
- PASCAL DENIS is a member of the network "référénts données" at Inria and Université de Lille
- PASCAL DENIS is administrator of Inria membership to Linguistic Data Consortium (LDC).

10.2 Teaching - Supervision - Juries

10.2.1 Teaching

- Licence MIASHS: MARC TOMMASI, Data Science, 24h, L2, Université de Lille.
- Licence Informatique: MARC TOMMASI, Introduction to AI, 24h, L2, Université de Lille.
- Licence MIASHS: MIKAELA KELLER, Data Science, 12h, L2, Université de Lille.
- Licence MIASHS: MIKAELA KELLER, Data Science 2, 24h, L3, Université de Lille.

- Licence MIASSH: MIKAELA KELLER, Traitement de données, 24h, L2, Université de Lille.
- Master MIASSH: MIKAELA KELLER, Algorithmes fondamentaux de la fouille de données, 27h, M1, Université de Lille.
- Master Data Science: MIKAELA KELLER, Machine Learning 1, 24h, M1, Université de Lille.
- Master Computer Science: MIKAELA KELLER, Apprentissage profond, 24h, M1, Université de Lille.
- Master Computer Science: MIKAELA KELLER, Machine learning pour le traitement automatique du langage naturel, 24h, M2, Université de Lille.
- Master MIASSH: MICHAËL PERROT, Algorithmes fondamentaux de la fouille de données, 27h, M1, Université de Lille.
- Master Computer Science: MARC TOMMASI, Data Science, 48h, M1, Université de Lille.
- Master Computer Science: MARC TOMMASI, Semi-supervised learning and Graphs, 24h, M2, Université de Lille.
- Master Data Science: MARC TOMMASI Seminars 24h.
- Master Data Science: AURÉLIEN BELLET, Privacy Preserving Machine Learning, 24h, M2, Université de Lille and Ecole Centrale de Lille.
- Parcours Science des Données et Intelligence Artificielle: AURÉLIEN BELLET, Advanced Machine Learning, 6h, Ecole Centrale de Lille.
- Peyresq Summer School 2022: AURÉLIEN BELLET, Privacy-Preserving Machine Learning, 5h, PhD-level.
- Training in Privacy-Preserving and Federated Machine Learning, 3 days, researchers from L'Oréal R&D.
- Master Informatique: PASCAL DENIS, Foundations of Machine Learning, 46h, M1, Université de Lille.
- Master Sciences Cognitives: MICHAËL PERROT, Machine Learning for Cognitive Sciences, 8h, M2, Université de Lille.
- MARC TOMMASI is directeur des études for the Machine Learning master of Computer Science.

10.2.2 Supervision

- Postdoc: VITALII EMILIANOV. On the interactions between fairness and privacy in machine learning. November 2022. MICHAËL PERROT.
- Postdoc: BATISTE LE BARS. On collaboration graph design for decentralized learning. October 2021-. AURÉLIEN BELLET and MARC TOMMASI
- PhD defended in December 2022: MAHSA ASADI, On Decentralized Machine Learning, since Oct 2018. AURÉLIEN BELLET and MARC TOMMASI. (See [40])
- PhD in progress: NICOLAS CROSETTI, Privacy Risks of Aggregates in Data Centric-Workflows, since Oct 2018. FLORENT CAPELLI and SOPHIE TISON and JOACHIM NIEHREN and JAN RAMON.
- Phd in progress: ARIJUS PLESKA, Tractable Probabilistic Models for Large Scale Networks, since Oct 2018. JAN RAMON
- PhD in progress: MOITREE BASU, Integrated privacy-preserving AI, since 2019. JAN RAMON.
- PhD defended in June 2022: CÉSAR SABATER, Privacy Preserving Machine Learning, since 2019. JAN RAMON. (See [41])

- PhD in progress: PAUL MANGOLD. Decentralized Optimization and privacy. AURÉLIEN BELLET and MARC TOMMASI and JOSEPH SALMON, since October 2020.
- Phd in progress: GAURAV MAHESHWARI. Trustworthy Representations for Natural Language Processing, since Nov 2020. AURÉLIEN BELLET, MIKAELA KELLER, and PASCAL DENIS
- Phd in progress: PRIYANSH TRIVEDI. Enriching Linguistic Representations with External Knowledge, since Nov 2020. PHILIPPE DE GROOTE (Loria, Nancy) and PASCAL DENIS
- Phd in progress: EDWIGE CYFFERS. Decentralized learning and privacy amplification, since Oct. 2021. AURÉLIEN BELLET
- Phd in progress: MARC DAMIE. Secure protocols for verifiable decentralized machine learning, since May 2022. JAN RAMON with Andreas Peter (U. Twente, NL & U. Oldenburg, DE) Florian Hahn (University of Twente, NL).
- Phd in progress: TUDOR CEBERE. Privacy-Preserving Machine Learning, since Nov. 2022. AURÉLIEN BELLET
- Phd in progress: BASTIEN LIÉTARD. Computational Models of Lexical Semantic Change, since Nov. 2022. ANNE CARLIER (Université Paris Sorbonne), PASCAL DENIS and MIKAELA KELLER
- Phd in progress: DINH-VIET-TOAN LE. Natural Language Processing approaches in the musical domain : suitability, performance and limits, since Oct. 2022. MIKAELA KELLER and LOUIS BIGO
- Phd in progress: ALEKSEI KORNEEV. Trustworthy multi-site privacy-enhancing technologies, since Dec. 2022. JAN RAMON
- PhD in progress: ANTOINE BARCZEWSKI. Transparent privacy-preserving machine learning, since May 2022. JAN RAMON.
- Engineer: JEAN-PAUL LAM. Transparent privacy-preserving machine learning. JAN RAMON
- Engineer: ANTONIN DUEZ. Transparent privacy-preserving machine learning. JAN RAMON
- Engineer SOPHIE VILLEROT, ADT project Tailed: Trustworthy AI Library for Environments which are Decentralized, since Nov. 2020. JAN RAMON
- Engineer JOSEPH RENNER, Improving Word Representations with Semantic Knowledge, since Nov. 2020. PASCAL DENIS and RÉMI GILLERON
- Engineer RISHABH GUPTA, Disentanglement approaches for speech data. AURÉLIEN BELLET and MARC TOMMASI.
- Engineer PAUL ANDREY, Decentralized and Federated Learning with DecLearn. AURÉLIEN BELLET and MARC TOMMASI.
- Engineer NATHAN BIGAUD, Decentralized and Federated Learning with DecLearn. AURÉLIEN BELLET and MARC TOMMASI.

10.2.3 Juries

- AURÉLIEN BELLET was reviewer for the PhD thesis of Théo Ryffel (ENS/Inria) and Raphael Ettetdgui (PSL), and member of the PhD defense committee of Evrard Garcelon (FAIR Paris and CREST, ENSAE).
- MARC TOMMASI was member of the habilitation of Romain Azaïs (ENS Lyon).
- PASCAL DENIS was member of the PhD defense committee of Tom Bourgeade at IRIT, Université de Toulouse.
- PASCAL DENIS was member of the PhD award committee at TALN'22.

10.3 Popularization

10.3.1 Articles and contents

- AURÉLIEN BELLET participated to a round table organized by **Le Cycle de la Voix (Le Voice Lab)** and was **interviewed by CNIL**.
- MICHAËL PERROT was interviewed for the article L'Impasse de l'IA, Alexandra Pihen, Epsilon, June 2022.
- Marc Tommasi was **interviewed by CNIL**.

11 Scientific production

11.1 Major publications

- [1] A. Bellet, R. Guerraoui and H. Hendriks. 'Who started this rumor? Quantifying the natural differential privacy guarantees of gossip protocols'. In: *DISC 2020 - 34th International Symposium on Distributed Computing*. Freiburg / Virtual, Germany, Oct. 2020. URL: <https://hal.inria.fr/hal-02166432>.
- [2] A. Bellet, R. Guerraoui, M. Taziki and M. Tommasi. 'Personalized and Private Peer-to-Peer Machine Learning'. In: *AISTATS 2018 - 21st International Conference on Artificial Intelligence and Statistics*. Lanzarote, Spain, Apr. 2018, pp. 1–20. URL: <https://hal.inria.fr/hal-01745796>.
- [3] M. Dehouck and P. Denis. 'Phylogenetic Multi-Lingual Dependency Parsing'. In: *NAACL 2019 - Annual Conference of the North American Chapter of the Association for Computational Linguistics*. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Minneapolis, United States, June 2019. URL: <https://hal.archives-ouvertes.fr/hal-02143747>.
- [4] P. Kairouz, B. H. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings et al. 'Advances and Open Problems in Federated Learning'. In: *Foundations and Trends in Machine Learning* 14.1-2 (2021), pp. 1–210. URL: <https://hal.inria.fr/hal-02406503>.
- [5] E. Lassalle and P. Denis. 'Joint Anaphoricity Detection and Coreference Resolution with Constrained Latent Structures'. In: *AAAI Conference on Artificial Intelligence (AAAI 2015)*. Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI 2015). Austin, Texas, United States, Jan. 2015. URL: <https://hal.inria.fr/hal-01205189>.
- [6] G. Maheshwari, P. Denis, M. Keller and A. Bellet. 'Fair NLP Models with Differentially Private Text Encoders'. In: *Findings of the Association for Computational Linguistics: EMNLP 2022*. Abu Dhabi, United Arab Emirates, 2022. URL: <https://hal.inria.fr/hal-03905094>.
- [7] C. Pelekis, J. Ramon and Y. Wang. 'H¹-Holder-type inequalities and their applications to concentration and correlation bounds'. In: *Indagationes Mathematicae* 28.1 (2017), pp. 170–182. DOI: [10.1016/j.indag.2016.11.017](https://doi.org/10.1016/j.indag.2016.11.017). URL: <https://hal.archives-ouvertes.fr/hal-01421953>.
- [8] T. Ricatte, R. Gilleron and M. Tommasi. 'Skill Rating for Multiplayer Games Introducing Hypernode Graphs and their Spectral Theory'. In: *Journal of Machine Learning Research* 21 (2020), pp. 1–18. URL: <https://hal.inria.fr/hal-02566930>.
- [9] C. Sabater, A. Bellet and J. Ramon. 'An Accurate, Scalable and Verifiable Protocol for Federated Differentially Private Averaging'. In: *Machine Learning* (28th Oct. 2022). DOI: [10.1007/s10994-022-06267-9](https://doi.org/10.1007/s10994-022-06267-9). URL: <https://hal.inria.fr/hal-03820603>.
- [10] A. S. Shamsabadi, B. M. L. Srivastava, A. Bellet, N. Vauquier, E. Vincent, M. Maouche, M. Tommasi and N. Papernot. 'Differentially private speaker anonymization'. In: *Proceedings on Privacy Enhancing Technologies* 2023.1 (1st Jan. 2023). URL: <https://hal.inria.fr/hal-03588932>.

- [11] P. Vanhaesebrouck, A. Bellet and M. Tommasi. ‘Decentralized Collaborative Learning of Personalized Models over Networks’. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. Fort Lauderdale, Florida., United States, Apr. 2017. URL: <https://hal.inria.fr/hal-01533182>.
- [12] F. Vitale, N. Parotsidis and C. Gentile. ‘Online Reciprocal Recommendation with Theoretical Performance Guarantees’. In: *NIPS 2018 - 32nd Conference on Neural Information Processing Systems*. Montreal, Canada, Dec. 2018. URL: <https://hal.inria.fr/hal-01916979>.

11.2 Publications of the year

International journals

- [13] J. A. Alvarado, Y. Wang and J. Ramon. ‘Limits of Multi-relational Graphs’. In: *Machine Learning* (Nov. 2022). URL: <https://hal.inria.fr/hal-03881631>.
- [14] M. Asadi, A. Bellet, O.-A. Maillard and M. Tommasi. ‘Collaborative Algorithms for Online Personalized Mean Estimation’. In: *Transactions on Machine Learning Research* (2022). URL: <https://hal.inria.fr/hal-03905917>.
- [15] Y. Belal, A. Bellet, S. Ben Mokhtar and V. Nitu. ‘PEPPER: Empowering User-Centric Recommender Systems over Gossip Learning’. In: *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 6.3 (6th Sept. 2022), pp. 1–27. DOI: [10.1145/3550302](https://doi.org/10.1145/3550302). URL: <https://hal.archives-ouvertes.fr/hal-03867109>.
- [16] C. Sabater, A. Bellet and J. Ramon. ‘An Accurate, Scalable and Verifiable Protocol for Federated Differentially Private Averaging’. In: *Machine Learning* (28th Oct. 2022). DOI: [10.1007/s10994-022-06267-9](https://doi.org/10.1007/s10994-022-06267-9). URL: <https://hal.inria.fr/hal-03820603>.
- [17] C. Sabater, F. Hahn, A. Peter and J. Ramon. ‘Private Sampling with Identifiable Cheaters’. In: *Proceedings on Privacy Enhancing Technologies* 2023.2 (2022). URL: <https://hal.inria.fr/hal-03904200>.
- [18] A. S. Shamsabadi, B. M. L. Srivastava, A. Bellet, N. Vauquier, E. Vincent, M. Maouche, M. Tommasi and N. Papernot. ‘Differentially private speaker anonymization’. In: *Proceedings on Privacy Enhancing Technologies* 2023.1 (1st Jan. 2023). URL: <https://hal.inria.fr/hal-03588932>.
- [19] B. M. L. Srivastava, M. Maouche, M. Sahidullah, E. Vincent, A. Bellet, M. Tommasi, N. Tomashenko, X. Wang and J. Yamagishi. ‘Privacy and utility of x-vector based speaker anonymization’. In: *IEEE/ACM Transactions on Audio, Speech and Language Processing* (15th June 2022). URL: <https://hal.inria.fr/hal-03197376>.
- [20] N. Tomashenko, X. Wang, E. Vincent, J. Patino, B. M. L. Srivastava, P.-G. Noé, A. Nautsch, N. Evans, J. Yamagishi, B. O’Brien, A. Chanclu, J.-F. Bonastre, M. Todisco and M. Maouche. ‘The VoicePrivacy 2020 Challenge: Results and findings’. In: *Computer Speech and Language* 74 (July 2022), p. 101362. DOI: [10.1016/j.csl.2022.101362](https://doi.org/10.1016/j.csl.2022.101362). URL: <https://hal.archives-ouvertes.fr/hal-03332224>.

International peer-reviewed conferences

- [21] T. Anikina, N. Skachkova, J. Renner and P. Trivedi. ‘Anaphora Resolution in Dialogue: System Description (CODI-CRAC 2022 Shared Task)’. In: *CODI-CRAC 2022*. Gyeongju, South Korea, Oct. 2022. URL: <https://hal.inria.fr/hal-03925147>.
- [22] A. Bellet, A.-M. Kermarrec and E. Lavoie. ‘D-Cliques: Compensating for Data Heterogeneity with Topology in Decentralized Federated Learning’. In: *41st International Symposium on Reliable Distributed Systems (SRDS)*. Vienna, Austria, 2022. URL: <https://hal.inria.fr/hal-03905085>.
- [23] F. Capelli, N. Crosetti, J. Niehren and J. Ramon. ‘Linear Programs with Conjunctive Queries’. In: *25th International Conference on Database Theory (ICDT 2022)*. Edinburgh, United Kingdom, 29th Mar. 2022. URL: <https://hal.archives-ouvertes.fr/hal-01981553>.

- [24] E. Cyffers and A. Bellet. ‘Privacy Amplification by Decentralization’. In: Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, Virtual, Spain, 2022. URL: <https://hal-cnrs.archives-ouvertes.fr/hal-03906735>.
- [25] E. Cyffers, M. Even, A. Bellet and L. Massoulié. ‘Muffliato: Peer-to-Peer Privacy Amplification for Decentralized Optimization and Averaging’. In: Advances in Neural Information Processing Systems 35 (NeurIPS). New Orleans, United States, 2022. URL: <https://hal-cnrs.archives-ouvertes.fr/hal-03906768>.
- [26] P. Humbert, B. Le Bars and L. Minvielle. ‘Robust Kernel Density Estimation with Median-of-Means principle’. In: Proceedings of the 39th International Conference on Machine Learning (ICML). Vol. 162. Proceedings of Machine Learning Research. Baltimore, United States, 17th July 2022. URL: <https://hal.archives-ouvertes.fr/hal-02882092>.
- [27] T. Kerdoncuff, M. Perrot, R. Emonet and M. Sebban. ‘Optimal Tensor Transport’. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI. Vancouver, Canada, 22nd Feb. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03479241>.
- [28] G. Maheshwari, P. Denis, M. Keller and A. Bellet. ‘Fair NLP Models with Differentially Private Text Encoders’. In: Findings of the Association for Computational Linguistics: EMNLP 2022. Abu Dhabi, United Arab Emirates, 2022. URL: <https://hal.inria.fr/hal-03905094>.
- [29] P. Mangold, A. Bellet, J. Salmon and M. Tommasi. ‘Differentially Private Coordinate Descent for Composite Empirical Risk Minimization’. In: ICML 2022 - 39th International Conference on Machine Learning. Vol. 162. Baltimore, United States: PMLR, 2022, pp. 14948–14978. URL: <https://hal.inria.fr/hal-03424974>.
- [30] M. Maouche, B. M. L. Srivastava, N. Vauquier, A. Bellet, M. Tommasi and E. Vincent. ‘Enhancing speech privacy with slicing’. In: Interspeech 2022 - Human and Humanizing Speech Technology. Incheon, South Korea, 18th Sept. 2022. URL: <https://hal.inria.fr/hal-03369137>.
- [31] S. Mdhaffar, J.-F. Bonastre, M. Tommasi, N. Tomashenko and Y. Estève. ‘Retrieving Speaker Information from Personalized Acoustic Models for Speech Recognition’. In: IEEE ICASSP 2022. Singapour, Singapore, 2022. URL: <https://hal.archives-ouvertes.fr/hal-03539741>.
- [32] M. Noble, A. Bellet and A. Dieuleveut. ‘Differentially Private Federated Learning on Heterogeneous Data’. In: Proceedings of The 25th International Conference on Artificial Intelligence and Statistics (AISTATS). Virtual, Spain, 2022. URL: <https://hal.inria.fr/hal-03905078>.
- [33] C. Oguz, I. Kruijff-Korbayová, P. Denis, E. Vincent and J. van Genabith. ‘Chop and change: Anaphora resolution in instructional cooking videos’. In: Findings of ACL-IJCNLP 2022 - 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics - 12th International Joint Conference on Natural Language Processing. Taipei, Taiwan, 20th Nov. 2022. URL: <https://hal.inria.fr/hal-03807530>.
- [34] S. Sajadmanesh, A. S. Shamsabadi, A. Bellet and D. Gatica-Perez. ‘GAP: Differentially Private Graph Neural Networks with Aggregation Perturbation’. In: 32nd USENIX Security Symposium. Anaheim, United States, 2023. URL: <https://hal.inria.fr/hal-03905068>.
- [35] J. O. D. Terrail, S.-S. Ayed, E. Cyffers, F. Grimberg, C. He, R. Loeb, P. Mangold, T. Marchand, O. Marfoq, E. Mushtaq, B. Muzellec, C. Philippenko, S. Silva, M. Teleńczuk, S. Albarqouni, S. Avestimehr, A. Bellet, A. Dieuleveut, M. Jaggi, S. P. Karimireddy, M. Lorenzi, G. Neglia, M. Tommasi and M. Andreux. ‘FLamby: Datasets and Benchmarks for Cross-Silo Federated Learning in Realistic Healthcare Settings’. In: NeurIPS 2022 - Thirty-sixth Conference on Neural Information Processing Systems. Proceedings of NeurIPS. New Orleans, United States, 28th Nov. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03900026>.
- [36] N. Tomashenko, S. Mdhaffar, M. Tommasi, Y. Estève and J.-F. Bonastre. ‘Privacy attacks for automatic speech recognition acoustic models in a federated learning framework’. In: ICASSP 2022. Singapour, Singapore, 2022. URL: <https://hal.archives-ouvertes.fr/hal-03539742>.

National peer-reviewed Conferences

- [37] S. Mdhaffar, J.-F. A. Bonastre, M. Tommasi, N. Tomashenko and Y. Estève. ‘Extraction d’informations liées au locuteur depuis un modèle acoustique personnalisé’. In: *JEP 2022*. JEP 2022. île de Noirmoutier, France, 13th June 2022. URL: <https://hal.archives-ouvertes.fr/hal-03706944>.
- [38] N. Tomashenko, S. Mdhaffar, M. Tommasi, Y. Estève and J.-F. Bonastre. ‘On speaker verification from the neural network footprint of personalized acoustic models’. In: *Journées d’Études sur la Parole - JEP2022*. Île de Noirmoutier, France, June 2022. URL: <https://hal.inria.fr/hal-03626964>.

Scientific book chapters

- [39] M. Dijkslag, M. Damie, F. Hahn and A. Peter. ‘Passive Query-Recovery Attack Against Secure Conjunctive Keyword Search Schemes’. In: *Applied Cryptography and Network Security*. Vol. 13269. Lecture Notes in Computer Science. Springer International Publishing, 18th June 2022, pp. 126–146. DOI: [10.1007/978-3-031-09234-3_7](https://doi.org/10.1007/978-3-031-09234-3_7). URL: <https://hal.inria.fr/hal-03854573>.

Doctoral dissertations and habilitation theses

- [40] M. Asadi. ‘Identifying structure in online and collaborative learning problems’. Lille University, 25th Nov. 2022. URL: <https://hal.inria.fr/tel-03892355>.
- [41] C. Sabater. ‘Efficient and Robust Protocols for Privacy-Preserving Semi-Decentralized Machine Learning’. Université de Lille, 20th June 2022. URL: <https://hal.inria.fr/tel-03904039>.

Reports & preprints

- [42] B. L. Bars, A. Bellet, M. Tommasi, E. Lavoie and A.-M. Kermarrec. *Refined Convergence and Topology Learning for Decentralized SGD with Heterogeneous Data*. 17th Dec. 2022. URL: <https://hal.inria.fr/hal-03905091>.
- [43] G. Maheshwari and M. Perrot. *FairGrad: Fairness Aware Gradient Descent*. 22nd June 2022. URL: <https://hal.archives-ouvertes.fr/hal-03902196>.
- [44] A. Mandal, M. Perrot and D. Ghoshdastidar. *A Revenue Function for Comparison-Based Hierarchical Clustering*. 29th Nov. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03902209>.
- [45] P. Mangold, A. Bellet, J. Salmon and M. Tommasi. *High-Dimensional Private Empirical Risk Minimization by Greedy Coordinate Descent*. 2022. URL: <https://hal.inria.fr/hal-03714465>.
- [46] P. Mangold, M. Perrot, A. Bellet and M. Tommasi. *Fairness Certificates for Differentially Private Classification*. 28th Oct. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03902203>.
- [47] N. Tomashenko, X. Wang, E. Vincent, J. Patino, B. M. L. Srivastava, P.-G. Noé, A. Nautsch, N. Evans, J. Yamagishi, B. O’Brien, A. Chanclu, J.-F. Bonastre, M. Todisco and M. Maouche. *Supplementary material to the paper The VoicePrivacy 2020 Challenge: Results and findings*. 26th Sept. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03335126>.

Other scientific publications

- [48] M. Kermarec, L. Bigo and M. Keller. ‘Improving Tokenization Expressiveness With Pitch Intervals’. In: 23rd International Society for Music Information Retrieval Conference (ISMIR 2022), Late-Breaking Demo Session. Bangaluru, India, 4th Dec. 2022. URL: <https://hal.archives-ouvertes.fr/hal-03877642>.