



TAILOR

Foundations of Trustworthy AI – Integrating Reasoning, Learning and Optimization

TAILOR

Grant Agreement Number 952215

Integrated Learning, Reasoning and Optimisation in Practice v.1 Report

Document type (nature)	Report
Deliverable No	4.3
Work package number(s)	4
Date	23 June 2022
Responsible Beneficiary	P 5 - KU Leuven
Editor(s)	Luc De Raedt
Publicity level	Public
Short description (Please insert the text in the Description of Deliverables in the Appendix 1.)	Integrated Learning, Reasoning and Optimisation in Practice v.1

History			
Revision	Date	Modification	Author
1.0	30 June 2022	Minor	Luc de Raedt

Document Review		
Reviewer	Partner ID / Acronym	Date of report approval
Philipp Slusalek	DFKI	25 June, 2022
Michela Milano	UNIBO	25 June, 2022

Table of Contents

Summary of the report	2
Introduction to the Deliverable	2
Organisation	3
Motivation and Aim	4
Description of the Tables	5
Analysis of the Tables	8
Towards Creating the Next Generation of Challenges	9

Summary of the report

Objectives: The question addressed in this WP is how to integrate learning, reasoning and optimisation, that is, how to computationally and mathematically integrate different AI paradigms. The most apparent difference between paradigms lies in the representations that are used and so an operational way to answer the question is to tightly integrate different representations as to offer both learning, reasoning and/or optimisation in common frameworks. This theme will therefore design representational systems with accompanying inference, learning and optimisation algorithms that can support trustworthy artificial intelligence. It will also study applications in two different domains. The WP is divided into four main Tasks, and is connected to other WPs by two tasks.

Introduction to the Deliverable

There are two deliverables for WP 4, that are both divided into an intermediate report (v1 M22) and a final report (v2 due at the end of the project).

Deliverables

D4.1: **Foundations**, techniques, algorithms and tools for integrating learning, reasoning and optimisation. (report) Report on the scientific challenges tasks T4.1 & T4.2.

D4.3: Integrated learning, reasoning and optimisation **in practice** (report). Report on the scientific challenges tasks T4.3 & T4.4.

This TAILOR WP has largely focused on two types of meetings and workshops. In the first type, there has been an emphasis on foundations, techniques, and tools for integrating learning, reasoning and optimization. In this type of workshop, the four scientific topics that characterize the first four tasks of WP4 within TAILOR have been covered. This has not only provided us with insight into the foundations and challenges connected to this WP, it has also delivered a number of interesting tutorials and survey papers, that have partly or fully been inspired by TAILOR and that led to novel insights and often also collaborations.

Deliverable 4.1 starts with these results, and then outlines the other results obtained within the WP. The second type of meeting was connected to the important taskforce of WP4 around benchmarks, datasets and systems. Given the plethora of different systems, representations and datasets, it is not easy to see the general picture in this diverse landscape. Therefore, we decided to start up a taskforce that would collect existing data, systems and study and compare them in order to get insight into the current and future abilities of integrated learning, reasoning and optimization approaches. This is the topic of Deliverable 4.3 and promises to result in publications summarizing useful observations and insights about the **practice** of *integrated learning, reasoning and optimisation* approaches.

Thus rather than dividing the deliverables along the task dimensions T4.1 / 2 vs T4.3 / 4 we found it more appropriate to report on the foundational issues in D4.1 and focus on the results of the taskforce in D4.3 as this is related to the potential and practice of WP 4 techniques.

Organisation

Main Contributors

Marco Lippi (UNIMORE)
Francesco Giannini (CINI)
Andrea Passerini (UNITN)
Emanuele Sansone (KUL)
Luc De Raedt (KUL)

Other People Involved

Neil Yorke-Smith (TU Delft), Sebastijan Dumancic (TU Delft), Tias Guns (KUL), Michele Lombardi (UNIBO), Debjit Paul (EPFL), Boi Faltings (EPFL), Kristian Kersting (TUDA), Devendra Dhami (TUDA), Mehdi Ali (FhG), Jens Lehmann (FhG), Michele Lombardi (UNIBO), Andrea Borghesi (UNIBO)

Motivation and Aim

Building systems that can integrate learning, reasoning and optimization has long been a dream for artificial intelligence. One of the major challenges, within this context, is certainly to evaluate novel ideas and frameworks on appropriate benchmarks. Too often, in fact, the tasks and the datasets that are considered and proposed for experimental evaluation are tailored to some algorithms or methodologies, and limited to ad-hoc scenarios and application domains. More in general, they lack an open and wider perspective to test the considered approaches across a variety of different tasks and under different conditions, making experimental comparisons hard to obtain. In addition, too often novel systems that aim to integrate learning, reasoning and optimization still rely on old-fashioned data and tasks: while a comparison with standard benchmarks is always useful to have an idea of the performance of an approach with respect to some reference point, we argue that the time is ripe for considering new challenges, which can drive the development of new integrated systems. To make an example, several classic datasets in image classification, such as MNIST or CIFAR, have been used for a wide variety of artificial tasks, each time with a specific goal: to propose a setting for few-shot learning, to introduce explicit knowledge for reasoning, to integrate rules and constraints for collective classification. In this sense, they have nowadays become real benchmarking frameworks. However, these datasets offer a limited playground for the development of systems integrating different paradigms (e.g. MNIST is limited both from the learning/perceptual perspective, as it is mainly devised to solve simple digit recognition tasks, and also from the reasoning perspective, as enabling forms of reasoning restricted to operations on natural numbers).

Consequently, can we define a set of requirements for a challenge that goes beyond those currently available?

Can we do this with the goal of obtaining a benchmarking framework that meets these requirements and that can still be implemented in a reasonable time? Possibly building on top of existing ones?

To address these questions, the TAILOR project has established a taskforce working across the different tasks of WP 4, identifying the following phases: (i) to analyze the current state-of-the-art for what concerns the existing datasets and corpora at the intersection of learning, reasoning and optimization; (ii) to study their limitations; (iii) to analyze the existing systems that have been applied to such data; (iv) to provide a list of the desiderata that new benchmarks should include; (v) to propose novel ideas for the evaluation and comparison of different approaches. This is all intended to provide insight into the abilities and limitations of current and future learning and reasoning systems.

It is worth mentioning that the goal is not just to list data collections, but especially to highlight which tasks can be applied to such data (i.e., in the form of benchmarks), and how a more extensive benchmarking framework could be designed, by unifying and composing a variety of heterogeneous tasks, working on the same original data collection. **As a consequence, the ultimate goal of the taskforce is to provide a suite of benchmarks which enable the creation of new tasks at a minimal cost and also provide a methodological evaluation to assess the performance of hybrid systems, which integrate the paradigms of learning, reasoning and optimization, thus providing insight into the practice and driving also future research.**

The **expected outcomes** at the end of the project are:

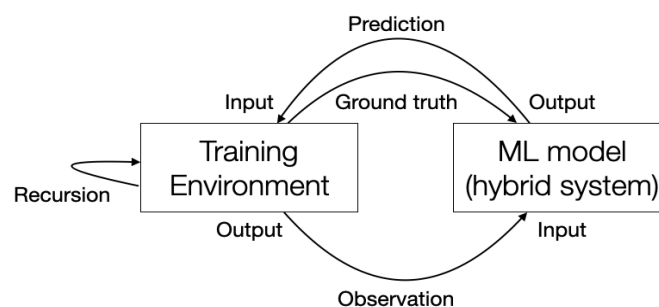
- insights into the abilities and limitations of the hybrid systems we study
- a number of publications comparing such systems, both theoretically and empirically
- a number of new challenges for hybrid systems

In the final deliverable, we also intend to go one step further and highlight some actual applications developed with integrated learning, reasoning and optimisation techniques.

Description of Tables - Datasets and Systems

The taskforce focused on contexts where low-level data is combined with knowledge, so that perception and reasoning skills need to be integrated in order to solve the tasks at hand. Knowledge could be either implicit (i.e., specific of the domain, derived from commonsense, encoded into data structures) or explicit (i.e., made available in the form of logic predicates and rules, or as constraints); either exact or uncertain; features could be just numeric or symbolic, or a combination of the two; examples could be either independent or involved in a variety of relations. All these characteristics have a strong impact on the categories of systems that can handle the corresponding datasets and which benchmarks can be defined upon.

To handle the complexity of the problem, the taskforce has produced two tables, one for datasets/benchmarks and one for systems. The general structure of the tables is based on a reinforcement learning setting, where the machine learning model (namely our hypothetical hybrid system) and the training environment interact with each other. The machine learning model provides a prediction for each observation. The training environment provides ground truth feedback to the machine learning model based on the received prediction and generates new observations optionally using historical information (i.e. in the form of a recursion). The overall setting is depicted in the following figure.



Importantly, the two tables take two opposite perspectives of the same setting. Indeed, the table about datasets/benchmarks focuses on the training environment, hence the input and the output are the prediction of the ML model and the new observation, respectively. The table about the ML model swaps the input and the output. In other words, the input consists of the observation provided by the training environment, whereas the output consists of the prediction for the observation.

The ground truth feedback from the training environment to the ML model represents the supervisory information which can be used to drive the learning, reasoning and optimization of the ML model.

Table about Datasets/Benchmarks

The table contains a list of datasets for each task of TAILOR WP 4:

1. **Task 4.1, Learning and Reasoning:** List of all datasets including **explicit knowledge**.
2. **Task 4.2, Learning and Optimisation:** List of all datasets related to **constraints and optimization problems**.
3. **Task 4.3, Knowledge Graphs, Embeddings, Ontologies:** List of all datasets relying on **relational knowledge** and **embedding representation**, such as **KG**.
4. **Task 4.4, Perception, Spatial Reasoning, and Vision:** List of all datasets with connection to learning and reasoning with **implicit knowledge**.

Information is structured according to general content (grey columns, such as URL of dataset, license, brief description etc.), the training environment (blue columns, defining the input, the output, the ground truth feedback and whether recursion is used or not) and the evaluation procedure (green columns, such task, metrics and baselines).

The table about datasets/benchmarks is shown in the next page.

Table about Systems

The table contains a list of systems for each task of TAILOR WP 4. As mentioned earlier, the table takes the perspective of the ML model.

Information is structured according to learning systems (blue columns, related to learning component dealing with perception tasks), reasoning systems (green columns, related to the reasoning/optimization component dealing with high-level tasks) and their integration (orange columns)

The table about datasets/benchmarks is shown in the next page.

Analysis of the Tables

From the analysis of the tables, the taskforce collected a list of desiderata describing some properties that potential new datasets and benchmarking frameworks should include. We hereby describe such properties, trying to focus on different aspects: the nature of data; paradigms and tasks for learning, optimization, and reasoning; novel metrics to measure performance; practical issues dealing with software platforms, tools, and implementation; novel domains of interest that have been seldom investigated within this context.

Concerning data. By considering the data level only, a clear starting point is to combine low-level data (images, videos, text, signals) with knowledge of some kind. This knowledge could be implicit or explicit, exact or uncertain. One desirable feature would be to enable the possibility to ask different questions within the same dataset, thus by exploiting different sets or types of knowledge across different tasks. For this reason, considering multiple data sources (e.g., multimodal data) could be an interesting additional feature, as well as to include a dynamic dimension to tackle evolving data. That of temporal data is indeed a challenging domain that has seldom been considered, and which would need to rethink paradigms and tasks for experimental evaluation.

Concerning paradigms and tasks. Regarding the tasks and the paradigms for learning, reasoning or optimization, the taskforce identified a crucial element of novelty in interactive learning, where humans could interact with the systems, by providing various forms of feedback, from simple labels to critiques, explanations and arguments. This will habilitate interactive debugging during learning and foster interpretability and trustworthiness, and it will be especially relevant in systems that consider lifelong or continual learning and again the temporal dimension, allowing them to adapt to distribution and knowledge drifts and to small-data regimes (i.e., few-shot learning). Additionally, it would also be interesting to jointly consider multiple learning tasks within a single benchmark, since this would allow testing multiple skills at once of the systems.

Concerning performance. Another point that was raised by the analysis of the tables is how performance should be measured. Besides considering classic metrics that essentially focus on accuracy, benchmarks that aim to include and exploit background knowledge should also measure the interpretability of the results (following the recent trends in eXplainable AI) and possibly the coherence of the predictions with the available knowledge. Energy efficiency to reduce the carbon footprint is yet another dimension to consider.

Concerning implementation. From a more practical perspective, it has been noted that the comparison of the same system across different benchmarks, or of different systems on the same benchmark, is made difficult by the heterogeneity in the formalisms used to represent data and to model background knowledge. A standardization of frameworks would represent a crucial step to improve such comparisons and to advance the state-of-the-art: this could be enabled by providing APIs to the systems, by providing knowledge in different formats, or by including benchmarks within existing platforms such as OpenML.

Concerning domains. Finally, the analysis of the datasets table was very useful in highlighting how some domains are under-represented in the panorama of benchmarks that are usually considered. Planning is a clear example of an application domain that would be perfectly suitable for testing the integration of learning, reasoning and optimization, as it can easily provide both symbolic data, such as activity traces or maps, and numeric data, coming from perception. The medical and legal domains represent as well two scenarios where background knowledge provided by experts could be a crucial element to boost performance of purely data-driven systems: such knowledge could be provided in various formats, including knowledge graphs, ontologies, or even plain natural language. Visual question answering and conversational agents are instead two candidate applications to allow interaction with users and knowledge integration in the fields of computer vision and natural language processing: in the latter case, computational argumentation and argumentation mining could be an additional research field where symbolic knowledge is typically employed to encode argument models. Finally, safety-critical applications have also been identified as a domain where it is quite usual to have hard and soft constraints that intelligent agents have to satisfy when interacting with the environment.

Towards Creating the Next Generation of Challenges

While benchmarks are clearly extremely important in providing a common ground to quantitatively evaluate the performance of different solutions, in modern research on AI there is a concrete risk of benchmark *hyperspecialization* and *overfitting*, in which the goal of research becomes beating the state-of-the-art on a specific benchmark (or group of closely related benchmarks), and the longer-term objective of which the benchmark is an initial and very partial proxy is lost.

The taskforce organized a panel discussing these topics, and how to create novel challenges that allow to overcome the limitations of existing benchmarks and encourage the exploration of radically new ideas, in particular involving the combination of learning, reasoning and optimization. The panelists were Fosca Giannotti, Marco Gori, Kristian Kersting, Michèle Sebag and Joaquin Vanschoren, and the panel was moderated by Andrea Passerini.

A first critical aspect was identified in the obsolescence of benchmarks, which is especially important when talking about standard, static benchmarks, and calls for solutions involving evaluation of benchmark overfitting, benchmark evolution, dynamic benchmarking and the relation with lifelong and continual learning tasks.

A major requirement for long-term challenges was identified in the possibility of having a diverse set of tasks to be accomplished. This calls for solutions relying on interactive learning environments, most likely virtual ones, where a combination of broad perceptual and reasoning abilities are needed in order to successfully accomplish the tasks.

A second major requirement concerns the need to have the human in-the-loop of the process. This is in-line with the human-centric and trustworthy perspective on AI fostered by the EC, and poses a number of new challenges in how to make this interaction efficient and effective.

Finally, the evaluation metrics and process for these systems should be substantially revised. Standard measures like accuracy are clearly insufficient and need to be complemented with aspects involving energy efficiency, interpretability, reliability, but most importantly the utility of the *joint* system that combines machine(s) and human(s).

Appendix: Program of a WP 4 Workshop on This Deliverable

What Are the Next Measurable Challenges in AI? (March 3, 2022)

Building systems that can integrate learning, reasoning and optimization has long been a dream for artificial intelligence. One of the major challenges, within this context, is certainly to evaluate novel ideas and frameworks on appropriate benchmarks. Too often, in fact, the tasks and the datasets that are considered and proposed for experimental evaluation are tailored to some algorithms or methodologies, and limited to ad-hoc scenarios and application domains. More in general, they lack an open and wider perspective to test the considered approaches across a variety of different tasks and under different conditions, making experimental comparisons hard to obtain.

Can we define a set of requirements for a challenge/benchmark that goes beyond those currently available?

Can we do it with the goal of having a benchmark (or rather a benchmarking framework maybe) that meets these requirements and can still be implemented in a reasonable time? possibly building on top of existing ones?

Program

13:00-13:15 Doors open

Introduction

13:15-13:30 Introduction & Expectations - Luc De Raedt

13:30-14:00 Invited Talk: Lessons Learned at NeurIPS 2021 Datasets and Benchmarks - Joaquin Vanschoren

PART I (grounding the discussion in the literature)

14:00-14:30 Presentation Datasets/Systems Tables - Marco Lippi/Francesco Giannini

14:30-15:30 Discussion on Tables - Working groups

15:30-15:45 Break

PART II (widening the perspective)

15:45-16:45 Panel on Limitations of Existing Benchmarks and New Challenges - Andrea Passerini

- Fosca Giannotti
- Marco Gori
- Kristian Kersting
- Michele Sebag
- Joaquin Vanschoren

16:45-18:00 Discussion on Panel - Working groups

18:00-18:15 What's Next? - Luc De Raedt