

Explanation: what does it mean for humans, for machines, for man-machine interactions?

Rémy Chaput¹, Amélie Cordier³, Alain Mille^{1,2}

April 8, 2021

¹Université de Lyon, Université Lyon1, LIRIS UMR CNRS 5205

²Coexistence, Lyon, France

³Lyon-iS-Ai, Lyon, France

**Adaptation journée
MADICS-Fender Rennes 8 Juillet
2021**

There is an important, growing, **impact** of Artificial Intelligence applications on society.

Even more so with Machine Learning and Deep Learning, but not limited to them.

⇒ Urgency of research on Explainable Artificial Intelligence

Although numerous works begin to take into account end users, they mostly do so by **pre-constructing** explanations for specific audience profiles.

We posit that explanation is a **complex process**, and must be **co-constructed** with the users, within their own context: task, responsibilities, knowledge, mental model of the system, etc.

**Explanations: what does it mean
for humans?**

Key observations

- Explaining is a **continuous** process
- Explaining is a **co-adaptive** process
- Explanation must be **triggered**
- We should facilitate **self-explanation**
- Explanation is an **exploration**
- **Contrast** situations can be explained **differently**

Explanations: what does it mean for AI?

Symbolic based expert systems

The symbolic AI theoretically allows to provide the reasoning steps pursued to propose an answer to a request.

The **origin** and the **justification** of the knowledge having been used to answer the request are **not generally available**.

The expert knowledge is not questionable and the explanation is reduced to the trace of the reasoning in a form often only useful for **debugging**.

Many **methods** and **techniques** to produce explanations; however, the explanation is often considered an **artifact** and not a **process**.

In particular, even when toolkits consider several audiences, the ability to construct the explanation by exploring is not focused.

In the field of social robotics and of collaborative robotics, robots are designed to **interact** with human beings. They have to be designed to ensure that humans **trust** robots they are working with, to make sure that they can operate **safely**.

In the field of Internet of Things (IoT), many people express concerns about **safety, security, data privacy and resilience**. But, it could be much more difficult to explain IoT systems because of the **complexity** emerging from a **vast network** of simple devices.

Bio-inspired AI methods do not provide symbolic explanations but offer simulators to become familiar with **emerging behaviors**.

The use of these methods is slowed down by the **difficulties** of explanation and research communities are active in associating information to **key moments** of the behavior as a possible support for explanations.

Towards an UXAI model

The question of explanation arises in **any decision support** agent. Artificial Intelligence techniques can help implementing a dynamic explanation process, constructed and conducted by the user of the system.

It is possible to design a **dynamic explanation process**, based on **explanatory agents** able to learn in **co-construction with users** how to explain the behavior of design support agents.

- Researchers build decision models, from expert knowledge and collected data. These are General Models (GMs).
- Designers use published GMs to build Applied Models (AMs), operationalized in the context of a task (a user task).
- Users use the AM in an application as a support for their activity.
[option] The interaction traces provide data that can be used to improve the AM. Researchers can also collect data to build the data corpuses that are fed into their own GM.

- Researchers publish a GM and an associate general model of explanation (UXAI-GM), which contains the necessary knowledge to explain it.
- Designers integrate into the explanations applied models (UXAI-AM) the possibility of co-constructing explanations with the users. A first UXAI-AM is deployed along with the AM, and both evolve synchronously.
- Users have not only an application but also an explanation agent with it. They can learn to understand the application's behavior, for example by self-explaining. The explanation agent stores in its memory the various explanations and their contexts, as situated explanations which can be reused.

Conclusion

Main issues for UXAI?

1. ethical problems of misuse;
2. respective responsibilities of AI agent / User;
3. economic consequences of a refusal of use as a precautionary principle.

What does UXAI can provide?

1. the user associates a self-explanation with the explanation provided by the agent;
2. the resulting explanation is contextualized with the help of the user for a similar usage situation;
3. contradictions between design logic and usage logic feed the agent design loop.

Thank You

Any questions?

References

-  Chakraborti, Tathagata, Sarath Sreedharan, and Subbarao Kambhampati (Feb. 2020). “The Emerging Landscape of Explainable AI Planning and Decision Making”. en. In: *arXiv:2002.11697 [cs]*. arXiv: 2002.11697.
-  Garcia-Magarino, Ivan, Rajarajan Muttukrishnan, and Jaime Lloret (2019). “Human-Centric AI for Trustworthy IoT Systems With Explainable Multilayer Perceptrons”. en. In: *IEEE Access* 7, pp. 125562–125574.
-  Hoffman, Robert R, Gary Klein, and Shane T Mueller (2018). “Explaining explanation for “explainable ai””. In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*. Vol. 62. 1. SAGE Publications Sage CA: Los Angeles, CA, pp. 197–201.



Swartout, William and Johanna D Moore (1985). "Explainable (and Maintainable) Expert Systems". In: *Proceedings of the 9th International Joint Conference on Artificial Intelligence*. Vol. 1. Los Angeles.