Explanation Ontology in Action: A Clinical Use-Case *

Abstract. We addressed the problem of a lack of semantic representation for user-centric explanations and different explanation types in our Explanation Ontology (https://purl.org/heals/eo). Such a representation is increasingly necessary as explainability has become an important problem in Artificial Intelligence with the emergence of complex methods and an uptake in high-precision and user-facing settings. In this submission, we provide step-by-step guidance for system designers to utilize our ontology, introduced in our resource track paper, to plan and model for explanations during the design of their Artificial Intelligence systems. We also provide a detailed example with our utilization of this guidance in a clinical setting.

Resource: https://tetherless-world.github.io/explanation-ontology

Keywords: Modeling of Explanations and Explanation Types \cdot Supporting Explanation Types in Clinical Reasoning \cdot Tutorial for Explanation Ontology Usage

1 Introduction

Explainable Artificial Intelligence (AI) has been gaining traction due to increasing adoption of AI techniques in high-precision settings. Consensus is lacking amongst AI developers on the type of explainability approaches and we observe a lack of infrastructure for user-centric explanations. User-centric explanations address a range of users' questions, have different foci, and such variety provides an opportunity for end-users to interact with AI systems beyond just understanding why system decisions were made. In our resource paper [2], we describe an Explanation Ontology (EO), which we believe is a step towards semantic encoding of the components necessary to support user-centric explanations. This companion poster focuses on describing the usage steps (Section 2) that would

^{*} Copyright ©2020 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

serve as a guide for system developers hoping to use our ontology. We demonstrate the usage of the protocol (Section 3) as a means to support and encode explanations in a guideline-based clinical decision support setting.

2 Usage Directions for the Explanation Ontology

Protocol 1 Usage of Explanation Ontology at System Design Time

Inputs: A list of user questions, knowledge sources and AI methods Goal: Model explanations that need to be supported by a system based on inputs from user studies

The protocol:

1. Gathering requirements

- (a) Conduct a user study to gather the user's requirements of the system
- (b) Identify and list **user questions** to be addressed

2. Modeling

- (a) Align **user questions** to explanation types
- (b) Finalize **explanations** to be included in the system
- (c) Identify **components** to be filled in for each explanation type
- (d) Plan to populate slots for each explanation type desired
- (e) Use the **structure of sufficiency conditions** to encode the desired set of explanations

System designers can follow the usage directions for EO at design time when planning for the capabilities of an AI system. The guidance aims to ensure that end-user requirements are translated into user-centric explanations. This protocol guidance is supported by resources made available on our website. These resources include queries to competency questions to retrieve sample user questions addressed by each explanation type³ and the components to be filled for each explanation type. Additionally, the sufficiency conditions that serve as a means for structuring content to fit the desired explanation type can be browsed via our explanation type details page.

3 Clinical Use Case

We applied the protocol to understand the need for explanations in a guidelinebased care setting and identify which explanation types would be most relevant

 $^{^3}$ https://tetherless-world.github.io/explanation-ontology/competency questions/ $\# {\rm question2}$

 $^{^4}$ https://tetherless-world.github.io/explanation-ontology/competency questions/ # question3

 $^{^{5}}$ https://tetherless-world.github.io/explanation-ontology/modeling/ # modeling explanations

for this clinical use case. [2] Further, we utilized the EO to model some of the explanations we pre-populated into the system prototype. Hereafter, we describe our application of the EO usage guidelines and highlight how we used a user study to guide system design. As a part of the user study, we first held an expert panel interview to understand the clinicians' needs when working with guideline-based care. We utilized the expert panel input to design a cognitive walkthrough of a clinical decision support system (CDSS) with some explanations pre-populated at design time to address questions that clinicians would want answered.

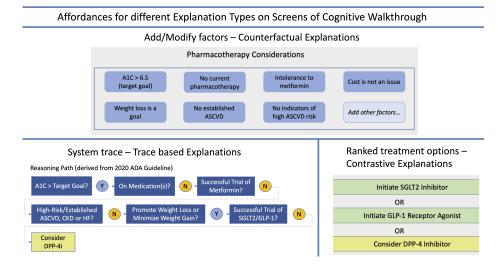


Fig. 1. An overview of the user interface affordances we designed into the prototype to accommodate the various explanation types.

The prototype of a CDSS we designed for the walkthrough, included allowances for explanations on different screens of the system (Fig. 1), in all, allowing clinicians to inspect a complicated type-2 diabetes patient case. Some examples of the pre-populated explanations are: contrastive explanations to help clinicians decide between drugs, trace based explanations to expose the guideline reasoning behind why a drug was suggested, and counterfactual explanations that were generated based on a combination of patient factors. Other examples were provided to us by clinicians during the walkthrough.

Below, we step-through how the protocol guidance (Section 2) can be used to model an example of a counterfactual explanation from this clinical use case. A counterfactual explanation can be generated by modifying a factor on the treatment planning screen of the CDSS prototype. For illustration's sake, let us suppose that this explanation is needed to address a question, "What if the patient had an ASCVD risk factor?" where the "ASCVD risk factor" is an alternate set of inputs the system had not previously considered. From our definition of a counterfactual explanation (https://tetherless-world.github.io/explanation-

4 S. Chari et al.

ontology/modeling/#counterfactual), in response to the modification a system would need to generate a system recommendation based on the consideration of the new input. More specifically, in our use case, a system would need to consider the alternate input of the ASCVD risk factor in conjunction with the patient context and consult evidence from the American Diabetes Association (ADA) guidelines [5] to arrive at a new suitable drug recommendation for the patient case. With the alternate set of inputs and the corresponding recommendation in place, the counterfactual explanation components can now be populated as slots based on the sufficiency condition for this explanation type. A Turtle snippet of this counterfactual explanation example can be viewed in Fig. 2. Additionally, there have been some promising machine learning (ML) model efforts in the explainability space [1,6] that could be used to generate system recommendations to populate specific explanation types. We are investigating how to integrate some of these AI methods with the semantic encoding of EO.

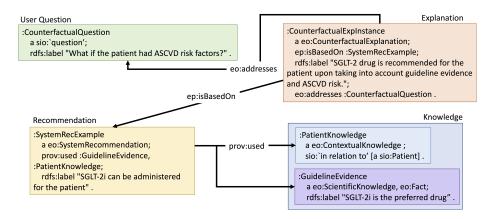


Fig. 2. An annotated turtle representation of a counterfactual explanation that the CDSS would output in response to an additional factor the system had not considered.

4 Related Work

In the explainable AI space, several research projects [4,7,3] have begun to explore how to support end-user specific needs such as ours. Wang et al. [7] present a framework for aligning explainable AI practices with human reasoning approaches, and they test this framework in a co-design exercise using a prototype CDSS. However, their framework isn't in a machine-readable format and hence is difficult to reuse for supporting explanations in a system. The findings from their co-design process that clinicians seek different explanations in different scenarios corroborate with those from our cognitive walkthrough, and EO can help support system designers in this endeavor. Similarly, Dragoni et al. [3] propose a rule-based explanation system in the behavior change space capable of providing trace based explanations to end-users to encourage better lifestyle and

nutrition habits. While this explanation system was tested with real target users and subject matter experts, it is limited in scope by the types of explanations it provides. Vera et al. [4] have released a question bank of explanation question types that can be used to drive implementations such as ours.

5 Conclusion

We have presented guidelines for using our Explanation Ontology in user-facing settings and have demonstrated the utility of this guidance in a clinical setting. This guidance describes an end-to-end process to support the translation of end-user requirements into explanations supported by AI systems. We are taking a two-pronged approach to pursue future work in operationalizing the use of our explanation ontology for clinical decision support systems. To this end, we are working towards building an explanation as a service that would leverage our explanation ontology and connect with AI methods to support the generation of components necessary to populate explanation types in different use cases. Further, we are implementing the clinical prototype as a functional UI with affordances for explanations. We expect the guidance along with open-sourced resources on our website will be useful to system designers looking to utilize our ontology to model and plan for explanations to include in their systems.

Acknowledgments

This work is done as part of the HEALS project, and is partially supported by IBM Research AI through the AI Horizons Network.

References

- Arya, V., Bellamy, R.K., Chen, P.Y., Dhurandhar, A., Hind, M., Hoffman, S.C., Houde, S., Liao, Q.V., Luss, R., Mojsilović, A., et al.: One Explanation Does Not Fit All: A Toolkit and Taxonomy of AI Explainability Techniques. arXiv preprint arXiv:1909.03012 (2019)
- 2. Chari, S., Seneviratne, O., Gruen, D.M., Foreman, M., Das, A.K., McGuinness, D.L.: Explanation Ontology: A Model of Explanations for User-Centered AI. In: Int. Semantic Web Conf. p. to appear. Springer (2020)
- 3. Dragoni, M., Donadello, I., Eccher, C.: Explainable ai meets persuasiveness: Translating reasoning results into behavioral change advice. Artificial Intelligence in Medicine p. 101840 (2020)
- 4. Liao, Q.V., Gruen, D., Miller, S.: Questioning the AI: Informing Design Practices for Explainable AI User Experiences. arXiv preprint arXiv:2001.02478 (2020)
- 5. American Diabetes Assoc.: 9. Pharmacologic Approaches to Glycemic Treatment: Standards of Medical Care in Diabetes-2020. Diabetes Care **43**(Suppl 1), S98 (2020)
- 6. Ribeiro, M.T., Singh, S., Guestrin, C.: Model-agnostic interpretability of machine learning. arXiv preprint arXiv:1606.05386 (2016)
- 7. Wang, D., Yang, Q., Abdul, A., Lim, B.Y.: Designing theory-driven user-centric explainable AI. In: Proceedings of the 2019 CHI Conf. on Human Factors in Computing Systems. pp. 1–15 (2019)