



For each paper we reviewed, we identify (1) what is the target text task (e.g., machine translation, question answering), (2) what is the explanation technique it uses (e.g., human-readable rules, highlighting important words), and (3) who are the appropriate target users (AI experts or lay users).

Valuable insights can be obtained by exploring our findings. For example, AI engineer can use XAIT to quickly explore different possible options and identify the most suitable technique(s) for their use cases. However, it is not easy to present our findings in an intuitive and accessible way so that users can quickly find the most relevant information. Clearly, listing our results in a spreadsheet will not serve the purpose.

XAIT was designed to address these issues. The core of XAIT is a tree-like structure that organizes research publications, explanation techniques, and target users in a taxonomy that is based on the one originally presented in [1, 3]. The taxonomy presented in [1, 3] was for general XAI tasks, the main goal of which is to suggest suitable toolkits for different XAI tasks. We reused the components (both the taxonomy structure and the available toolkits) that are text related. More importantly, we extend their taxonomy by adding more nodes, including relevant research publications, and suggesting target users according to our findings.

In this demo, we illustrate work in progress towards our overall goal to build a system that can recommend a list of possible explanation techniques as well as provide a list of existing approaches that adopted these explanation techniques.

## 2 OVERVIEW OF XAIT

As shown in Figure 1, our taxonomy is visually presented to the user. This visualization provides a global and structural view of different explanation techniques included in our taxonomy. The exploration starts from the root node, and the user can navigate down to the leaf nodes (according to the choices s/he made). Each inner node is a burst node that splits the exploration into different paths. The leaf nodes of the taxonomy contain the actual explanation techniques, toolkits (if available), target users, and existing publications that adopted these techniques. Each node also has a clickable link that can pop up a dialog that provides some high-level description of the purpose of this node. Below we introduce several key concepts used in XAIT.

**Static or Interactive Explanation.** Whether the explanation is provided in a static way or an interactive way.

**Data Explanation or Model Explanation.** Whether the explanations are generated for data or an AI model.

**Local or Global Explanation.** Whether the goal is to provide only local explanation (i.e., explain the prediction of a particular input instance) or global explanations (i.e., explain a model’s decision making process in general).

**Surrogate model.** In some machine learning scenarios, when we want to explain the predictions of a non-interpretable

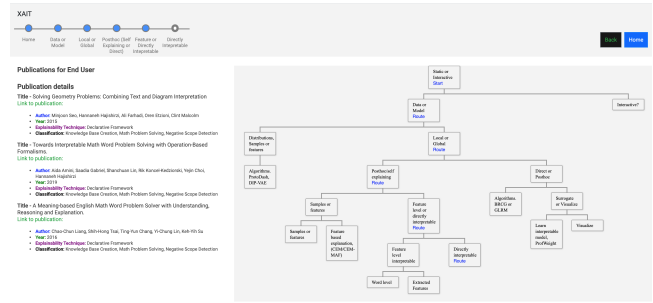


Figure 2: Publications and toolkits for directly and globally interpretable models.

model, we learn a second interpretable surrogate model that approximates the predictions of the first model.

**Posthoc or self-explaining.** Whether the explanation is generated via a self-explaining interpretable model (the learned model itself also generates some sort of explanations for its prediction) or post-hoc process (i.e., the learned model is not interpretable, and we need some post-hoc processes to explain its prediction, for example, learn an interpretable surrogate model).

**Feature-level or directly interpretable.** The explanations are based on feature-level elements (e.g., sparse word embeddings, or highlighted words) or the explanations are directly interpretable by human users (e.g., human-readable rules).

## 3 A CONCRETE DEMO EXPERIENCE

Assuming there is an AI engineer who needs to design an entity resolution (the task of identifying and linking different representations of the same real-world objects) model for end users who want to have a globally transparent AI model. The AI engineer can start with XAIT by selecting *static* mode at the root node, then selects *global* to indicate that s/he wants to provide global explanation. Then, s/he can select *Direct* to indicate that s/he wants to provide human-comprehensible explanations. After that, a leaf node with a list of relevant publications with detailed information and toolkits will be presented to the user (see Figure 2).

## REFERENCES

- [1] Vijay Arya, Rachel KE Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C Hoffman, Stephanie Houde, Q Vera Liao, Ronny Luss, Aleksandra Mojsilović, et al. 2019. One explanation does not fit all: A toolkit and taxonomy of ai explainability techniques. *arXiv preprint arXiv:1909.03012* (2019).
- [2] Shipi Dhanorkar, Christine Wolf, Kun Qian, Anbang Xu, Lucian Popa, and Yunyao Li. 2020. Tutorial: Explainability for Natural Language Processing. In *(AAACL 2020)*, to appear.
- [3] Michael Hind. 2019. Explaining explainable AI. *XRDS: Crossroads, The ACM Magazine for Students* 25, 3 (2019), 16–19.