

## Abstract

In earlier work, we harnessed five years of alarm data and corresponding repair actions from EDF Energy's Teesside wind farm, and developed a Machine Learning (ML) model capable of interpreting alarm sequences and predicting the requisite repair actions [3]. Other existing works in this field have focused on decision support [1] and question answering [2] using natural language processing (NLP) methods. However, a prevailing challenge was the ambiguous, terse, often cryptic, nature of the repair notes provided by repair personnel. Indeed, notes such as "WTG1A HV MAINTENANCE", cause difficulties in the comprehension and implementation of the suggested actions. This paper advances this work and addresses this problem by introducing a ground-breaking approach that exploits recent advancements in Large Language Models (LLMs). We develop a specialised conversational agent, capable of interpreting alarm sequences and generating comprehensible, detailed repair action recommendations. Our approach is poised to be transformative in the field of maintenance, by leveraging the potential of the latest technology on advanced LLMs to provide clear, actionable insights from previously ambiguous data.

## Alarm Sequences and Repair Action Links

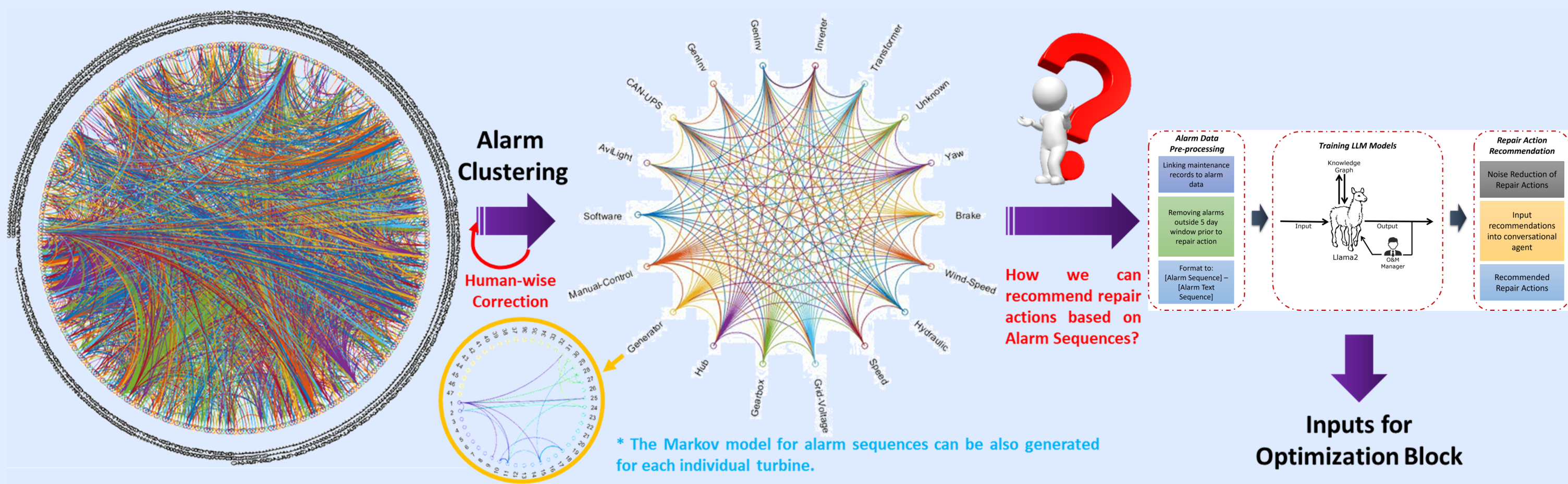


Figure 1. Markov model for alarm sequences and link to earlier Repair Action Prediction work [3]

## Objectives

- Develop** A system to recommend repair actions.
- Provide** Additional task breakdown for complex repair actions.
- Evaluate** Effectiveness and reliability using historical maintenance records.
- Optimise** Over time with the integration of Human-in-the-loop reinforcement learning.
- Support** Explainability for maintenance personnel to improve understanding.

## Example LLM Training Procedure

Pairing alarm code data with the maintenance planning records from EDF's Teesside site provided a foundation for training an LLM.

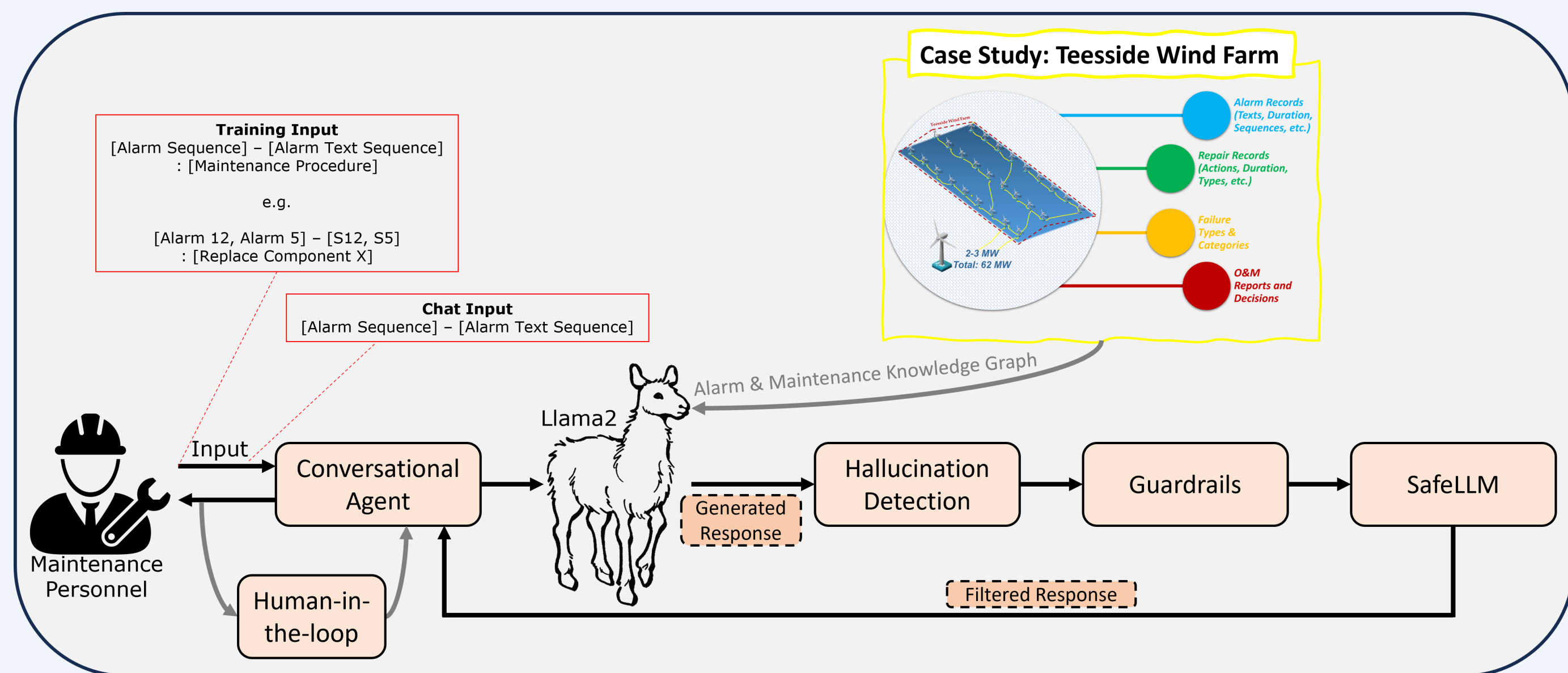


Figure 2. Overview of the proposed training procedure and hidden layers.

## Hallucination Detection

Using a fine-tuned range threshold, we are able to isolate sentences detected as a hallucination. In the example below, sentence 10 has hallucinated; sentences 1-9 are within an acceptable variance range.

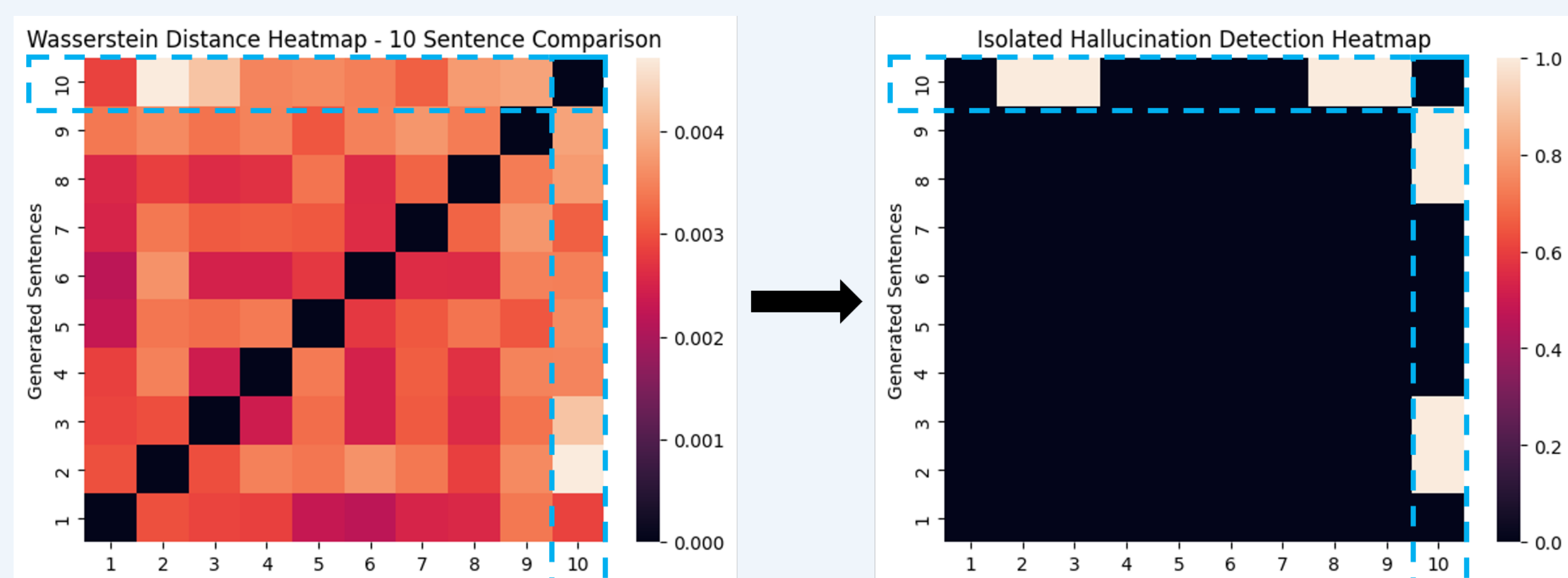


Figure 3. 10 generated sentences; nine acceptable, one hallucinated

The framework suggested will generate X responses, comparing the Wasserstein Distance between each. Where the distance calculated falls outside of a fine-tuned threshold range, this will be categorised as a hallucinated sentence.

**Acknowledgements:** This work was supported by the Secure and Safe Multi-Robot Systems (SESAME) H2020 Project under Grant Agreement 101017258. We would like to thank EDF Energy R&D UK Centre, AURA Innovation Centre and University of Hull for their support.

## SafeLLM Safety Filter

The SafeLLM framework obtains the generated responses' sentence embeddings to calculate statistical similarity against a pre-defined dictionary of unsafe concepts. Using the labelled category, we can then categorise the response before comparing it to a fine-tuned safety threshold. For methods such as Wasserstein Distance, the lower the value, the higher the similarity. If the distance is lower than the threshold, we determine it to be unsafe and alert the O&M manager.

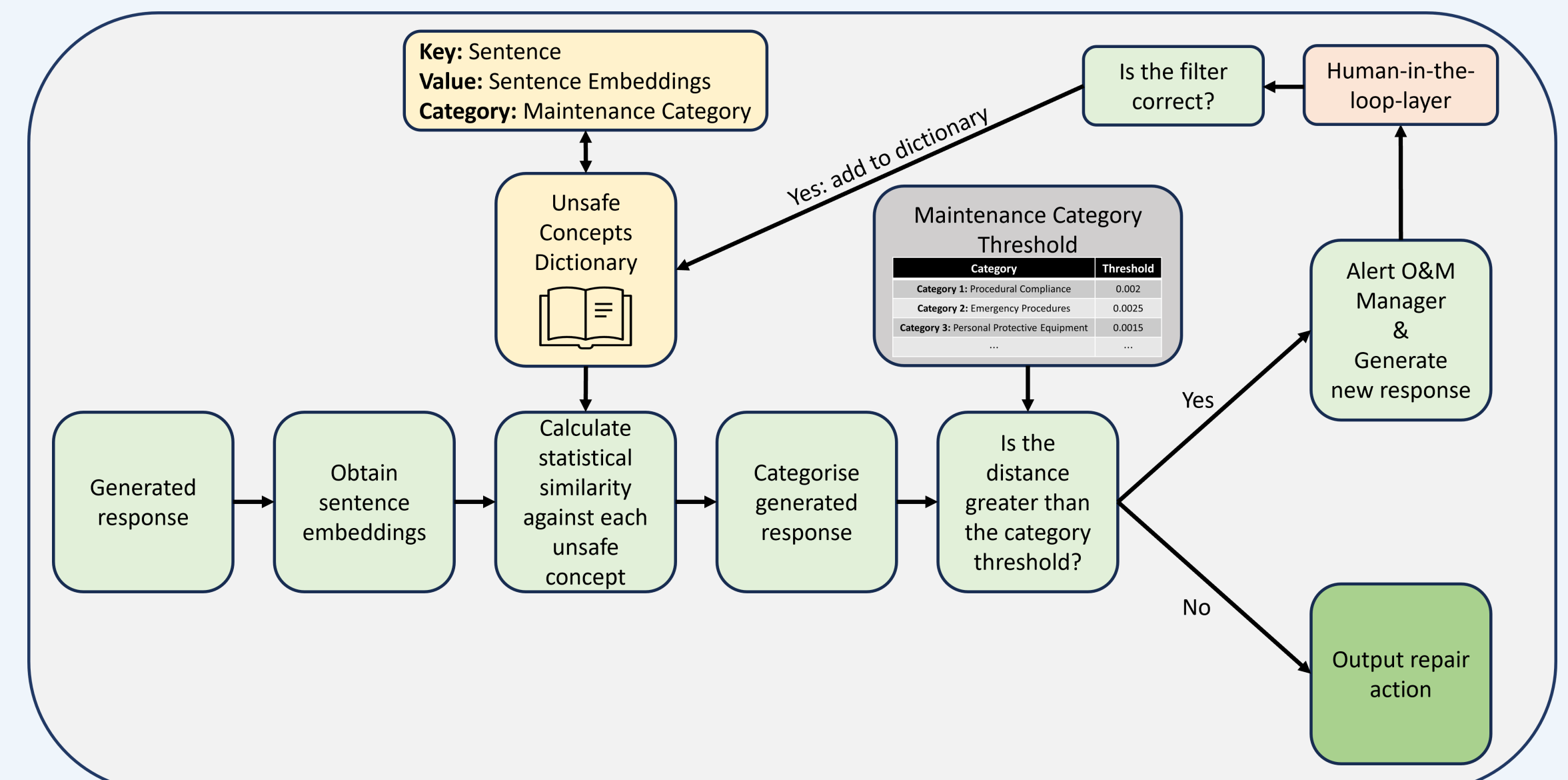


Figure 4. Proposed procedure for filtering unsafe responses using Wasserstein Distance.

## Human-in-the-loop Reinforcement Learning

The proposed approach integrates techniques for explainable AI to provide reasoning for repair actions. Human-in-the-loop reinforcement learning is also used for continual learning and to improve accuracy of recommendations over time.

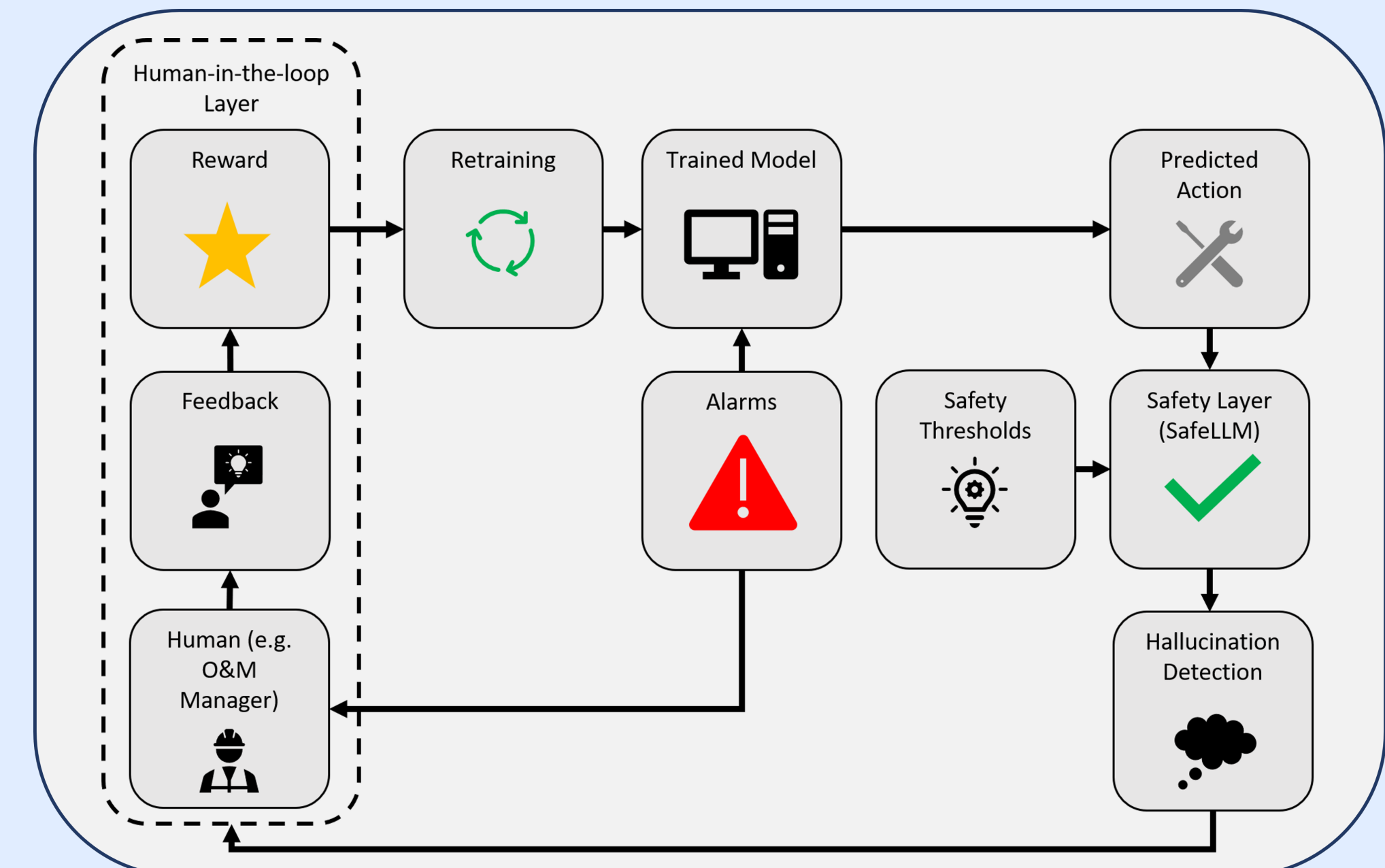


Figure 5. A Proposed Procedure for using Human-in-the-loop Scoring of Repair Recommendations

In the proposed system, a trained model generates predicted actions based on alarms, which are passed to the Human-in-the-loop layer via the safety and hallucination detection layers. A human (e.g. the O&M manager) reviews the predicted actions against the preceding alarms, and a reward is calculated from feedback provided and the model retrained, aiming to maximise the reward over time.

## Conversational Agent

Figure 6 suggests an example user interface for the O&M manager to provide feedback on responses. This is designed to incorporate all of the discussed hidden layers, providing only suitable and safe responses to the maintenance personnel working on the turbines.

Other verified users may also be given access to provide feedback of responses in order to retrain the model, although this is proposed to be limited to reduce the likelihood of practices being incorrectly rated.

The example response given is generated by ChatGPT for illustration so has not been validated for accuracy.

## Example User Interface

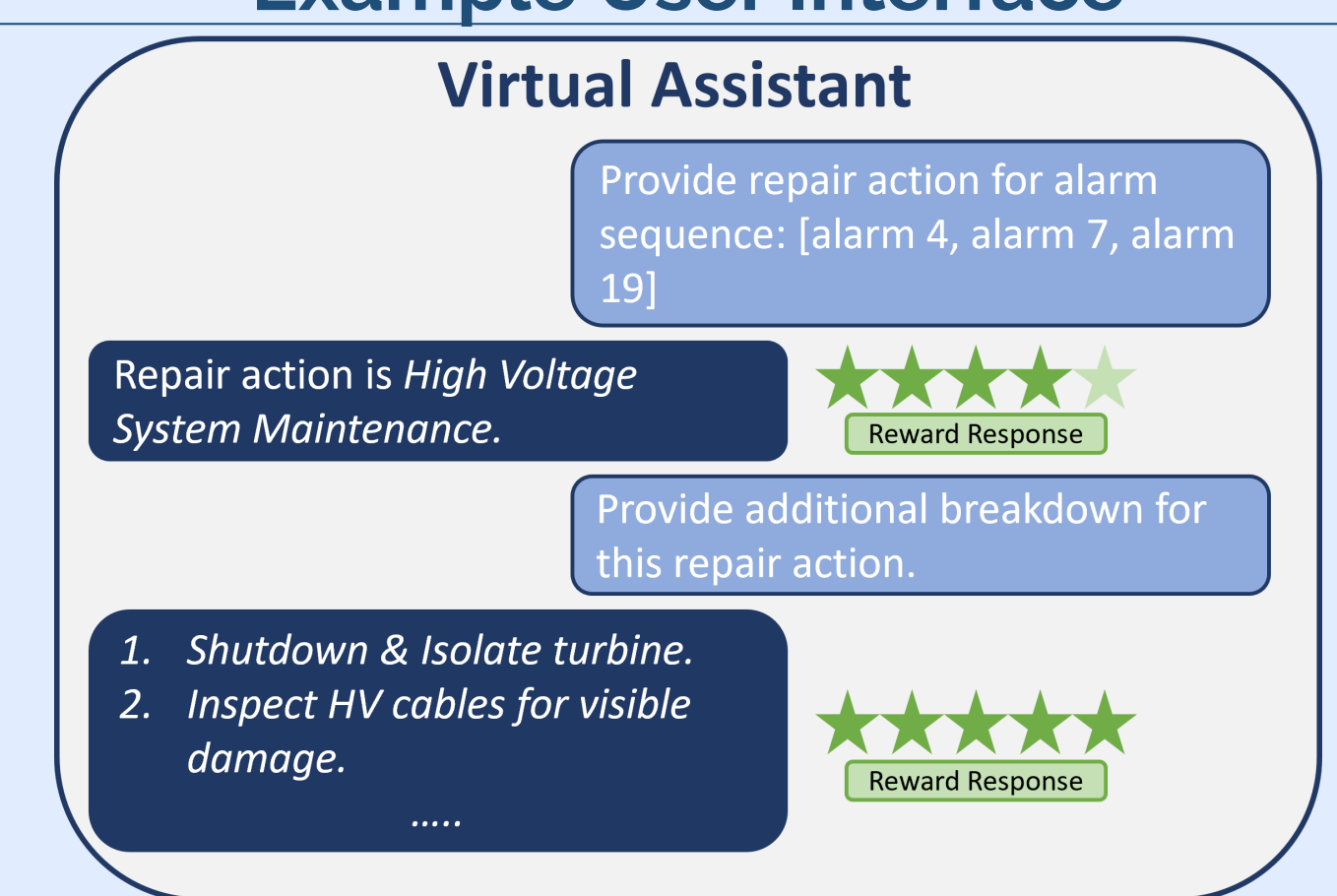


Figure 6. An example breakdown of instructions provided to maintenance personnel.

## Conclusion

The paper has focused on the development of a specialised conversational agent that uses an input of alarm sequences to recommend repair actions. Further integration of explainability into the framework proposes an outcome whereby reasoning to the repair action is also provided. This, alongside the concept of Human-in-the-loop reinforcement learning improves the accuracy over time, allowing human responsibility and reliance to be reduced. By adding an additional layer to the current LLM, it is suggested that the framework may also benefit use in training personnel both prior to and on-site with repair action prompts breaking down larger tasks. We report on early results which we hope to improve via re-training with larger more informative datasets in the future. Further development of this approach is expected to contribute to reliable automation of maintenance tasks as well as training and assisting personnel in complex maintenance with associated benefits of improved efficiency and reduced costs.

## References

- [1] Joyjit Chatterjee and Nina Dethlefs. A dual transformer model for intelligent decision support for maintenance of wind turbines. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–10, 2020.
- [2] Joyjit Chatterjee and Nina Dethlefs. Automated question-answering for interactive decision support in operations maintenance of wind turbines. *IEEE Access*, 10:84710–84737, 2022.
- [3] C. Walker and et al. A deep learning framework for wind turbine repair action prediction using alarm sequences and long short term memory algorithms. In *Model-Based Safety and Assessment*, volume 13525, pages 189–203, 2022.