# AI Alignment proposal №3: Enhancing Corrigibility in AI Systems through Robust Feedback Loops

aialignmentproposals.substack.com/p/ai-alignment-proposal-3-enhancing



## Abstract

This article explores the concept of corrigibility in artificial intelligence and proposes a detailed framework for a robust feedback loop to enhance corrigibility. The ability to continuously learn and correct errors is critical for safe and beneficial AI, but developing corrigible systems comes with significant technical and ethical challenges. The feedback loop outlined involves gathering user input, interpreting feedback contextually, enabling AI actions and learning, confirming changes, and iterative improvement. The article analyzes potential limitations of this approach and provides detailed examples of implementation methods using advanced natural language processing, reinforcement learning, and adversarial training techniques. It emphasizes the need for thoughtful design and testing to mitigate risks and biases. Fostering corrigibility is essential for aligning AI systems with human values.

## Introduction

As AI systems continue to permeate every sector of society, their potential to shape our lives and influence decision-making processes has received intense scrutiny. One crucial aspect of AI safety and ethics is corrigibility — the ability of an AI system to accept and respond to corrections from its users or operators. This article proposes a robust feedback loop as a promising strategy to enhance AI corrigibility.

## Implementing a Robust Feedback Loop

The idea behind implementing a robust feedback loop is to create a dynamic and interactive environment where users and the AI system can learn from each other. This enables the system to align more closely with the user's values, expectations, and instructions. The proposed feedback loop consists of five steps:

### 1. User Feedback

Allowing users to provide feedback on the AI's outputs is the foundation of this model. The feedback could be regarding factual inaccuracies, misconceptions of context, or any actions that violate the user's values or expectations. The system should be equipped with a simple, intuitive interface to facilitate this.

## 2. Feedback Interpretation

The AI system should interpret the feedback considering the appropriate context. It should not limit its understanding to immediate corrections but also infer the larger implications for similar future situations. Advanced natural language processing techniques and contextual understanding algorithms can be used to achieve this.

## 3. Action and Learning

After interpreting the feedback, the AI should take immediate corrective actions. Additionally, it should learn from this feedback to adjust its future responses. This learning can be facilitated by reinforcement learning techniques, where the AI adjusts its actions based on the positive or negative feedback received.

## 4. Confirmation Feedback

The system should then confirm with the user whether the correction has been implemented appropriately. This ensures that the system has correctly understood and applied the user's feedback. This can be done through a simple confirmation message or by demonstrating the corrected behavior.

## 5. Iterative Improvement

Finally, this process should be iterative, allowing the system to continuously learn and improve from ongoing user feedback. Each cycle of the feedback loop should refine the system's responses and behaviors.

# Addressing Limitations and Risks

While the idea of a robust feedback loop sounds promising, several counterarguments might challenge its feasibility and effectiveness:

**Counterargument 1:** User feedback might not always be accurate or reliable. Some users may provide malicious or misguided feedback that could lead the AI system to behave undesirably.

**Rebuttal:** This is a valid concern. However, this issue can be mitigated by incorporating a system of feedback verification, perhaps through consensus from multiple users or by using trusted moderators. Moreover, the system should be designed to detect and handle potentially harmful or malicious instructions.

**Counterargument 2:** Not all users may have the time, interest, or technical knowledge to provide detailed feedback and guide the AI's learning process.

**Rebuttal:** While true, this challenge can be addressed by making the feedback process as simple and intuitive as possible. Additionally, feedback could be incentivized to encourage more user participation. It's important to note that feedback doesn't always need to be active — it can also be passive, gleaned from user interactions with the system.

# Expanded Technical Details

Here is an expanded and more detailed version of the Implementing a Robust Feedback Loop:

## 1. User Feedback Collection

- Frontend interface — Simple widgets and APIs to collect ratings, text, audio etc. Integrate with common platforms like web, mobile, voice assistants.

- Database storage — Store feedback data in a relational SQL database or NoSQL database like MongoDB to allow complex querying.

- Data pipelines — Use Kafka, Airflow, etc. for scalable data ingestion and preprocessing.

- Sampling — SQL queries or data stream sampling to filter and sample feedback data for training. Address class imbalance.

This involves creating interfaces to gather multi-modal feedback like ratings, text, audio, video. Both active solicitation and passive collection are needed to motivate engagement while preserving privacy through anonymization and consent. Feedback should come from diverse users to mitigate bias. On a technical level, this requires frontend widgets, database storage like SQL and NoSQL, data ingestion pipelines, and sampling techniques to handle large volumes of data.

## 2. Interpretation Using NLP

- TensorFlow for ML framework — Provides scalability, distributed training, and portability for NLP models.

- HuggingFace Transformers — Pretrained NLP models like BERT for efficient fine-tuning on feedback data.

- Word embeddings — Map text to dense vector representations consumable by ML models.

- LSTM for sequence modeling — Recurrent neural network to interpret conversational context.

- Graph databases — Represent knowledge graphs for commonsense reasoning and entity linkage.

Advanced NLP techniques interpret the meaning and sentiment of feedback. This includes sentiment analysis, named entity recognition, conversational modeling to track context, commonsense reasoning using knowledge graphs, and explainability methods. TensorFlow provides a scalable machine learning framework to build these NLP models. Transfer learning from pretrained models like BERT enables efficient development. Word embeddings map text to vectors consumable by ML models. LSTMs specifically model conversational sequence context. Graph databases represent knowledge for reasoning.

## 3. Reinforcement Learning from Feedback

- OpenAI Gym — Simulation environments for initial RL model training and testing.

- PyTorch for ML framework — Provides auto-diff and modular libraries to build RL algorithms.

- Policy gradients — Algorithm to learn behaviors directly from feedback rewards and penalties.

- Transfer learning — Retrain final layers of RL model on new tasks and environments.

Reinforcement learning treats positive feedback as rewards and negative as penalties to shape behaviors. It enables transfer learning across environments and continual learning to assimilate new data. OpenAI Gym provides environments for initial simulation testing. PyTorch enables flexible RL model development with auto-differentiation and modular libraries. Policy gradient algorithms directly optimize behaviors based on feedback.

## 4. Confirmation and Demonstration

- Natural language generation — Template-based or neural models like GPT-3 to generate confirmation text and explanations.

- Visualization — Dynamic visualizations to demonstrate simulated behavior changes for validation.

Natural language generation confirms interpreted feedback with users. Visualization demonstrates simulated behavior changes for validation.

## 5. Controlled Iterative Improvement

- CI/CD pipelines — Automate testing and controlled deployments of new iterations.

- Canary releases — Slowly roll out to a small population to detect issues.

- Feature flags — Enable or disable functionality dynamically for control.

- Monitoring — Logging, metrics dashboards, anomaly detection to monitor for regressions.

- MLOps — Model versioning, reproducibility, and monitoring to track iterative improvement.

Robust testing, canary releases, monitoring, and MLOps help deploy improvements safely.

The technical complexity requires extensive testing and safety practices before real-world deployment. But this illustrates a lower-level view of how robust feedback loop principles could be realized.

## Improving Transparency with Hybrid Interpretability Model

To provide transparency into the system's inner workings, a hybrid interpretability model can be implemented. This model integrates various methods of interpretability and explainability to create a more robust and comprehensive understanding of AI systems. The model consists of three primary components:

1. **Feature Importance Visualization:** This involves using techniques like SHAP (SHapley Additive exPlanations), LIME (Local Interpretable Model-agnostic Explanations), or permutation feature importance to highlight which features are most influential in a model's predictions. This can help to reveal how the model is thinking on a broad scale.

2. **Model Transparency Tools:** These tools, such as Attention Visualization for transformers or CNN (Convolutional Neural Network) visualization for image-based models, provide a more granular look at the individual components of a model. For example, attention visualization can show which parts of input data a model is focusing on when making a decision.

3. **Counterfactual Explanations:** These are hypothetical scenarios that show how the outcome would change if the data were different. Counterfactual explanations can help to understand the boundaries and decision-making process of an AI model. For example, it might show that changing a specific feature from X to Y would lead the model to change its prediction from A to B.

This hybrid model could also integrate natural language explanations, where the AI system explains its reasoning in human-readable form. This could be particularly useful for complex models where feature importance visualization and model transparency tools might not be sufficient. The model should be designed to allow users to switch between these various modes of interpretability depending on their needs. For instance, a data scientist debugging the model might require a detailed view with model transparency tools, while an end-user or stakeholder might prefer a simpler, high-level explanation through feature importance visualization and natural language explanations. This layered, multi-faceted approach can provide a holistic understanding of an AI system's decision-

making process, thus significantly improving interpretability and ensuring that the alignment techniques presented in this propsal are having verifiable desirable effect on the AI system in training.

## Improving Robustness and Reliability

Here are some additional improvements that might be consider in pursuit of system's robustness and reliability:

- Incorporation of formal verification methods to prove safety-critical properties are preserved.

- Implementation of consensus validation from multiple users to prevent manipulation.

- Enabling layered feedback access with privileges for higher-impact changes.

- Analyzing demographics and expand testing for underrepresented populations.

- Supporting explainability so users understand how feedback impacts changes.

- Introducing feedback quality evaluation via user meta-reviews.

- Developing an ontology of acceptable behaviors aligned with human values.

- Incorporating simulation, game theory, and adversarial techniques to anticipate exploits.

- Integrating with bug bounty programs to stress test security.

- Implementing cryptographic assurances like blockchain or zero-knowledge proofs.

- Modeling choreography to maintain optimal performance.

- Gradient debugging for targeted unlearning.

- Model forking to enable reverting to historic versions if needed.

- Committee-based learning requiring consensus among models.

- Sandboxed experimentation for controlled testing.

- Compliance with relevant regulations like GDPR and ADA.

A multi-faceted approach considering technical, ethical, and social factors will produce the most robust outcomes. Responsible development demands avoiding harms across the full sociotechnical feedback loop.

## Conclusion

Enhancing the corrigibility of AI systems is a complex yet crucial endeavor. Implementing a robust feedback loop offers a promising approach, but it requires careful design to overcome potential pitfalls. It is also essential to remember that any approach must be grounded in ethical practices, respecting user privacy and ensuring the feedback process does not exploit or harm users. An open, participatory paradigm may be necessary to mitigate limitations of individual technical schemes. If pursued transparently and ethically, decentralized crowdsourced approaches could tap into our collective intelligence to shape AI for the future benefit of humanity. Creating beneficial AI demands ongoing collaborative processes grounded in shared values, not just fixed destinations. Through inclusive cooperation and constructive criticism, we can expand human potential and despite the challenges, the pursuit of corrigible AI systems is a necessary step towards ensuring that AI operates in alignment with human values and commonly accepted societal norms.