# Towards AI Self-Improvement of Understanding

## Integrating Self-Modeling, Multifaceted Evaluation, and Feedback Loops

By Ken Clements, assisted by Claude 3 Opus (2024)

## 1 - Introduction

### 1.1 The challenge of achieving genuine understanding in AI systems

Artificial Intelligence (AI) has made remarkable strides in recent years, with advanced language models and reasoning systems demonstrating impressive performance on a wide range of tasks. However, a fundamental question remains: do these systems truly understand the knowledge they process and the outputs they generate, or are they merely engaging in sophisticated pattern matching and information retrieval? Achieving genuine understanding in AI systems is a formidable challenge that requires going beyond surface-level language imitation to develop deeper cognitive capabilities such as reasoning, abstraction, grounded perception, and self-awareness.

### 1.2 The potential for AI to improve its own understanding through self-modeling and feedback

Recent research suggests that AI systems have the potential to iteratively improve their own understanding by engaging in self-modeling and feedback loops. Self-modeling involves an AI learning to predict its own internal states and processes, essentially developing a "theory of mind" about itself. This self-modeling ability has been shown to confer significant benefits, including increased regularization, simplification, and parameter efficiency in neural networks (Premakumar et al., 2024).

Furthermore, by engaging in cycles of generating hypotheses, testing them through interactions with the environment or simulated experiments, and updating its internal models based on the results, an AI system can bootstrap its own understanding in a manner akin to the scientific discovery process (Lu et al., 2024). This iterative refinement through self-modeling and feedback enables the AI to identify gaps or inconsistencies in its knowledge, leading to more robust and flexible representations.

### 1.3 Thesis: Integrating self-modeling, multifaceted evaluation frameworks like the MUTT, and feedback loops can enable AI systems to iteratively refine and expand their understanding capabilities

This paper proposes that by integrating self-modeling, comprehensive evaluation frameworks such as the Multifaceted Understanding Test Tool (MUTT), and feedback loops, AI systems can be designed to iteratively improve their own understanding capabilities. The MUTT, introduced in the book "Understanding Machine Understanding" (Clements, 2024), provides a rigorous approach to assessing

an AI's competencies across key dimensions such as language comprehension, reasoning, knowledge integration, perception, social cognition, and metacognition.

By combining self-modeling to optimize an AI's representations, multifaceted MUTT evaluations to comprehensively assess its understanding, and discovery feedback loops to generate and test hypotheses, a virtuous cycle can be created in which the system's understanding is continuously enriched through interleaved processes of introspection, testing, and model refinement. This integration of techniques holds the potential to dramatically accelerate an AI's ability to achieve deeper, more flexible and human-like understanding of the world and itself. The implications of this approach for advancing artificial general intelligence and enhancing human-AI collaboration will be explored in the sections to follow.

## 2 - Self-Modeling in AI Systems

### 2.1 Definition and key concepts of self-modeling

Self-modeling refers to an AI system's ability to build an internal representation or "self-model" of its own structure and operation. This involves the system learning to predict its own internal states, processes, and outputs in response to various inputs and contexts (Premakumar et al., 2024). In essence, the AI develops a theory of its own mind, allowing it to reason about and anticipate its own behavior.

Key concepts in self-modeling include:

- Metacognition: The ability of a system to monitor and regulate its own cognitive processes. Self-modeling enables a form of artificial metacognition, as the AI learns to track and modulate its own internal states.

- Introspection: The capacity to examine one's own mental states and processes. Self-modeling allows an AI to introspect on its own operation, enabling it to identify patterns, biases, and areas for optimization in its own cognition.

- Recursive processing: Self-modeling often involves recursive or self-referential computation, as the outputs of the system's self-model are fed back as inputs to guide its operation. This recursive processing is thought to be a key enabler of flexible, context-sensitive cognition.

### 2.2 Benefits of self-modeling for regularization, simplification, and efficiency

Recent research has demonstrated several compelling benefits of equipping AI systems, particularly neural networks, with self-modeling capabilities:

- Regularization: Self-modeling acts as a powerful regularizer, constraining the AI system to learn simpler, more generalizable representations. By forcing the system to build a predictive model of its own operation, self-modeling encourages the learning of efficient internal representations that avoid overfitting to noise or irrelevant features.

- Simplification: To support effective self-modeling, AI systems learn to make their internal operations more transparent and predictable. This pressure for simplicity leads to the emergence of modular,

disentangled representations that are easier to interpret and reason about, both for the AI itself and for human observers.

- Efficiency: Self-modeling enables AI systems to optimize their own operation by identifying redundancies, bottlenecks, and opportunities for compression in their internal representations. By learning to predict its own outputs, the AI can prune unnecessary computations and streamline its processing, leading to gains in speed and memory efficiency.

2.3 Implications of self-modeling for social cognition and cooperative AI

Beyond its benefits for individual AI systems, self-modeling may have profound implications for artificial social cognition and multi-agent cooperation:

2.3.1 **Theory of mind**: Just as self-modeling allows an AI to reason about its own mental states, it may also support the development of "theory of mind" models for reasoning about the beliefs, goals, and intentions of other agents. An AI that has learned to model itself effectively may be better equipped to model and understand the minds of others.

2.3.2 **Cognitive empathy**: Self-modeling may enable a form of artificial cognitive empathy, allowing an AI to simulate the thought processes and mental states of other agents. By "putting itself in another's shoes", an AI with self-modeling abilities may be able to more accurately predict and respond to the behaviors of other agents in social contexts.

2.3.3 **Cooperative alignment**: In multi-agent systems, self-modeling may facilitate the emergence of cooperative strategies and shared understanding. By modeling their own internal states and decision-making processes, AIs can better align their behaviors and coordinate their actions with other agents to achieve common goals.

2.3.4 **Social transparency**: AIs that engage in self-modeling may be more interpretable and transparent to human collaborators. By providing a window into their own internal operations, self-modeling AIs can foster greater trust and understanding in human-AI partnerships.

As research on self-modeling advances, it will be crucial to explore these social dimensions and develop frameworks for integrating self-modeling into the design of cooperative multi-agent systems. The implications for artificial social intelligence and human-AI collaboration are potentially transformative, opening up new frontiers in the quest to build socially aware and ethically aligned AI systems.

# 3 - The Multifaceted Understanding Test Tool (MUTT)

## 3.1 Motivation and key principles of the MUTT framework

The Multifaceted Understanding Test Tool (MUTT) is a proposed evaluation framework that aims to comprehensively assess an AI system's understanding capabilities across a range of cognitive dimensions. The motivation behind the MUTT is to address the limitations of narrow, task-specific benchmarks like the Turing Test, which focus primarily on surface-level language imitation rather than probing the depth and flexibility of a system's understanding.

The key principles guiding the design of the MUTT include:

3.1.1 **Breadth**: The MUTT incorporates a diverse array of task types and challenge scenarios that span multiple aspects of cognitive functioning, such as language comprehension, reasoning, knowledge integration, perception, social cognition, and metacognition. This broad coverage aims to provide a more holistic assessment of understanding.

3.1.2 **Depth**: The tasks and evaluation criteria used in the MUTT are designed to probe deeper, more nuanced levels of understanding beyond simple pattern recognition or information retrieval. This involves assessing a system's ability to draw insights, make inferences, and flexibly apply knowledge to novel contexts.

3.1.3 **Grounding**: The MUTT emphasizes the importance of grounding language understanding in real-world perception, action, and social interaction. This involves incorporating tasks that require linking words to their referents in the physical and social world, moving beyond purely text-based evaluations.

3.1.4 **Transferability**: The MUTT includes challenges that test a system's ability to transfer understanding across different domains, modalities, and task contexts. Strong performance on these transfer tasks would demonstrate a more general, flexible form of understanding.

3.1.5 **Open-endedness**: Many of the MUTT's tasks are designed to be open-ended, requiring the generation of novel responses or creative problem-solving rather than selecting from a fixed set of answer choices. This allows for a more authentic assessment of understanding.

3.2 Dimensions of understanding assessed by the MUTT

The MUTT is designed to evaluate an AI system's understanding capabilities across six key dimensions that are central to human-like intelligence:

3.2.1 **Language comprehension**: This dimension assesses a system's ability to understand and use natural language in context. It goes beyond surface-level syntax and semantics to probe skills like pragmatic reasoning, figurative language understanding, and discourse coherence. Example tasks could include:

- Interpreting implied meanings and speaker intentions in dialogue

- Identifying and explaining metaphors, idioms, and other non-literal language

- Tracking reference and maintaining coherence across extended texts

3.2.2 **Reasoning and abstraction**: This dimension evaluates a system's capacity for logical inference, analogical reasoning, causal reasoning, and abstraction. It tests whether the system can go beyond mere pattern recognition to grasp underlying conceptual relationships and draw valid conclusions. Example tasks could include:

- Solving novel logical reasoning problems that require deductive inference

- Identifying abstract patterns and relationships in data or stimuli

- Explaining the causal mechanisms underlying physical or social phenomena

3.2.3 **Knowledge integration**: This dimension assesses whether a system can effectively combine and apply knowledge from multiple domains to solve complex, real-world problems. It requires flexibly drawing upon a rich body of integrated world knowledge in a context-sensitive way. Example tasks could include:

- Answering questions that require synthesizing information from multiple academic fields

- Proposing creative solutions to open-ended challenges that draw on diverse knowledge

- Explaining how concepts from different domains relate to or depend on each other

3.2.4 **Perception and embodiment**: This dimension probes a system's ability to perceive, interpret, and interact with the physical world through sensors and effectors. It tests whether the system can ground its language understanding in real-world referents and actions. Example tasks could include:

- Engaging in situated dialogue about objects and events in a visual scene

- Following natural language instructions to navigate and manipulate a 3D environment

- Reasoning about the physical affordances and constraints of objects and spaces

3.2.5 **Social cognition**: This dimension evaluates a system's capacity for understanding and interacting with other agents in social contexts. It spans skills like emotion recognition, perspective-taking, social reasoning, and pragmatic communication. Example tasks could include:

- Inferring the mental states, beliefs, and intentions of others from their actions

- Engaging in multi-turn dialogue that obeys social and conversational norms

- Reasoning about the social implications and consequences of events or decisions

3.2.6 **Metacognition and self-awareness**: This dimension assesses whether a system can reflect on its own knowledge, reasoning, and experiences in an explicit, self-aware manner. It probes abilities like uncertainty estimation, self-explanation, and meta-level control of cognitive processes. Example tasks could include:

- Expressing appropriate degrees of confidence in its own outputs or decisions

- Providing clear, grounded explanations of its reasoning processes

- Identifying gaps or inconsistencies in its own knowledge and experiences

3.3 Advantages of the MUTT over narrow benchmarks like the Turing Test

The MUTT framework offers several key advantages over narrow, task-specific benchmarks like the Turing Test:

3.3.1 **Comprehensive assessment**: By spanning a broad range of cognitive dimensions and task types, the MUTT provides a more comprehensive assessment of an AI system's understanding capabilities compared to the narrow focus of the Turing Test on conversational imitation.

3.3.2 **Grounding in real-world contexts**: The MUTT emphasizes the importance of grounding language understanding in real-world perception, action, and social interaction. This allows for a more

ecologically valid evaluation of understanding compared to purely text-based Turing Test conversations.

3.3.3 **Probing deeper understanding**: The tasks and evaluation criteria used in the MUTT are explicitly designed to assess deeper levels of understanding, such as inference, abstraction, and cross-domain integration. This goes beyond the surface-level language skills that the Turing Test primarily evaluates.

3.3.4 **Mitigating gaming and biases**: The open-ended and multifaceted nature of the MUTT tasks makes it more difficult for AI systems to succeed through shallow pattern matching, information retrieval, or gaming strategies. This helps to alleviate some of the biases and limitations associated with narrow, fixed-format benchmarks like the Turing Test.

3.3.5 **Enabling finer-grained analysis**: By decomposing understanding into multiple distinct dimensions, the MUTT allows for a more precise, finer-grained analysis of an AI system's strengths and weaknesses. This can provide valuable insights for guiding future research and development efforts in a targeted way.

In summary, the MUTT framework aims to provide a more rigorous, comprehensive, and ecologically valid approach to evaluating machine understanding compared to the limitations of narrow benchmarks like the Turing Test. By assessing a broad range of cognitive capabilities grounded in real-world contexts, the MUTT can enable a more nuanced and informative assessment of the depth and flexibility of an AI system's understanding. This has important implications for advancing the field of artificial intelligence towards more human-like, general-purpose systems that can robustly understand and interact with the world.

## 4 - Feedback Loops for AI Self-Improvement

4.1 The concept of feedback loops in learning systems

Feedback loops are a fundamental concept in learning systems, both biological and artificial. At its core, a feedback loop involves using the outputs or consequences of a system's actions to modify its future behavior. This allows the system to iteratively refine its performance based on the results it achieves, leading to a continuous cycle of self-improvement.

In the context of learning, feedback loops enable systems to autonomously adapt and optimize their knowledge and skills without explicit programming. By comparing their predictions or actions against target outcomes, learning systems can generate error signals that guide the updating of internal representations and decision-making processes. Over time, this iterative refinement leads to increasingly accurate and robust performance.

Feedback loops are ubiquitous in biological learning, from the tuning of synaptic connections in the brain based on neural activity patterns, to the trial-and-error learning of new behaviors in animals and humans. These self-modifying processes allow organisms to flexibly adapt to their environments and improve their chances of survival and reproduction.

In artificial learning systems, feedback loops are implemented through various training paradigms such as supervised learning, reinforcement learning, and unsupervised learning. By defining objective functions or reward signals, these approaches allow AI systems to iteratively adjust their parameters and representations to minimize errors or maximize rewards. Backpropagation, the workhorse algorithm of deep learning, is essentially a method for propagating error signals backwards through a network to guide incremental weight updates.

4.2 Potential for feedback loops to enable AI systems to bootstrap their own understanding

While feedback loops have proven incredibly powerful for optimizing AI systems to perform specific tasks, their potential for enabling open-ended learning and bootstrapping of understanding is still largely untapped. Most current AI systems are trained in a narrow, task-specific manner, with limited ability to transfer knowledge across domains or to autonomously expand their capabilities.

However, there is growing interest in developing AI architectures that can leverage feedback loops to engage in more open-ended, self-supervised learning. By setting their own goals and generating their own training data through interaction with the environment, these systems could potentially "learn to learn" and bootstrap their own understanding in a more human-like way.

One promising approach is curiosity-driven learning, where AI systems are intrinsically motivated to explore and learn about their environment by seeking out novel or surprising experiences. By generating their own feedback signals based on the mismatch between their predictions and observations, curious AI agents can autonomously learn powerful representations and skills without explicit supervision.

Another area of research is meta-learning, or "learning to learn", where AI systems use feedback loops to optimize their own learning processes. By training on a diverse range of tasks and environments, meta-learning systems can learn general strategies for quickly adapting to new problems. This could enable AI systems to rapidly bootstrap their understanding in novel domains by leveraging prior learning experiences.

Self-play is another technique that leverages feedback loops to enable open-ended learning. By playing against copies of itself in complex environments, self-play allows AI systems to generate their own increasingly challenging training data and iteratively refine their strategies. This has led to remarkable breakthroughs in domains like game-playing AI, where systems like AlphaGo have bootstrapped superhuman performance purely through self-play.

4.3 Integrating self-modeling, MUTT evaluation, and feedback

The real power of feedback loops for bootstrapping AI understanding may come from integrating them with other key capabilities like self-modeling and comprehensive evaluation frameworks such as the MUTT.

### 4.3.1 **Self-modeling to enable AI to represent and reason about its own understanding**

Self-modeling, as discussed in previous sections, involves an AI system building an explicit representation of its own knowledge, capabilities, and limitations. By developing a "theory of mind"

about itself, a self-modeling AI can reason about its own understanding in a more structured and interpretable way.

This self-modeling ability is crucial for enabling an AI system to identify gaps or inconsistencies in its knowledge, formulate learning goals, and generate targeted feedback signals for self-improvement. By comparing its current self-model against a desired state of understanding, an AI can create an error signal to guide its learning and skill acquisition.

Self-modeling also allows an AI system to more effectively leverage its prior knowledge and experiences to bootstrap learning in new domains. By recognizing similarities between novel situations and its existing self-model, an AI can transfer relevant skills and adapt more rapidly. This kind of analogical reasoning and transfer learning is a hallmark of human intelligence that self-modeling could help replicate in AI systems.

### 4.3.2 **MUTT evaluation to comprehensively assess understanding capabilities**

While self-modeling provides an AI system with an internal representation of its own understanding, the MUTT framework proposed in this book offers a comprehensive external evaluation of an AI's actual capabilities. By assessing an AI's performance across a diverse range of tasks probing language understanding, reasoning, knowledge integration, social intelligence, and other key facets, the MUTT provides an objective measure of the depth and breadth of an AI's understanding.

Critically, the MUTT is designed not just as a static benchmark, but as a dynamic tool for driving iterative improvement in AI systems. By identifying specific areas where an AI's performance falls short, the MUTT can provide targeted feedback to guide the AI's learning and skill acquisition. This feedback can take the form of explicit error signals, such as highlighting incorrect or inconsistent outputs, or more implicit guidance, such as adjusting the difficulty or composition of training tasks to focus on areas of weakness.

The MUTT can also help calibrate an AI's self-modeling by providing an external reality check on its internal representations. If an AI's self-assessed capabilities consistently misalign with its MUTT performance, this can trigger a process of model refinement and updating. Over time, this iterative alignment between self-modeling and MUTT evaluation can lead to more accurate and robust representations of an AI's own understanding.

### 4.3.3 **Feedback loops to iteratively refine self-models and expand understanding based on MUTT results**

Putting it all together, the integration of self-modeling, MUTT evaluation, and feedback loops offers a powerful framework for bootstrapping AI understanding. The key is to create a virtuous cycle of self-assessment, external evaluation, and targeted improvement.

Here's how it could work: An AI system starts with an initial self-model representing its current capabilities. It then undergoes a comprehensive MUTT evaluation, assessing its performance across a range of understanding tasks. The results of this evaluation are fed back to the AI, highlighting areas of strength and weakness relative to its self-model.

Based on this feedback, the AI updates its self-model, identifying knowledge gaps and skill deficits that need to be addressed. It then generates its own learning goals and training tasks designed to target these areas of improvement. As it works through these self-generated challenges, the AI refines its internal representations and reasoning strategies.

Periodically, the AI undergoes additional MUTT evaluations to objectively measure its progress. The results of these evaluations provide further feedback to guide the AI's self-improvement process, adjusting its learning targets and strategies based on its performance gains and remaining weaknesses.

Over time, this iterative cycle of self-modeling, external evaluation, and targeted learning could enable an AI system to rapidly bootstrap its own understanding, expanding its knowledge and capabilities in an open-ended fashion. By creating a tight feedback loop between an AI's evolving self-representation and its objectively measured performance, this approach could help bridge the gap between narrow, task-specific intelligence and more general, flexible understanding.

Of course, implementing this vision of self-improving AI will require significant advances in areas like unsupervised learning, goal generation, and safe exploration. There are also important ethical considerations around aligning an AI's self-improvement process with human values and ensuring that its expanding capabilities remain controllable and beneficial.

But if developers can successfully integrate self-modeling, comprehensive evaluation, and targeted feedback, the potential for AI systems to bootstrap their own understanding and capabilities is immense. By creating AI that can learn, reason, and improve itself in an open-ended fashion, computer science may finally achieve the long-standing goal of artificial general intelligence.

This self-improving AI framework also has important implications for the future of human-AI collaboration. Rather than simply deploying static, task-specific AI systems, developers could create dynamic, adaptive AI partners that continuously refine their understanding and expand their capabilities through interaction with humans and the world. This could lead to a new era of discovery and innovation, as AI and humans work together to tackle ever-more complex and open-ended challenges.

Ultimately, the integration of self-modeling, comprehensive evaluation, and feedback loops offers a promising path towards AI systems that can truly bootstrap their own understanding, learning, and growth in an autonomous and open-ended fashion. While there are still significant technical and ethical hurdles to overcome, the potential benefits of self-improving AI for scientific discovery, technological innovation, and human flourishing are hard to overstate. As developers continue to push the boundaries of machine intelligence, this framework will be a critical guide for creating AI systems that can not only perform specific tasks, but truly understand and adapt to the complexities of the world around them.

## 5 - A Conceptual Architecture for Self-Improving AI Understanding

### 5.1 High-level system components and interactions

Drawing together the key insights from self-modeling, multifaceted evaluation, and feedback loops, one can envision a conceptual architecture for an AI system that iteratively improves its own understanding. At a high level, the system would consist of the following core components:

5.1.1 **Self-modeling module**: This component would be responsible for building and updating the AI's internal representations of its own knowledge, reasoning processes, and capabilities. Using techniques like auxiliary prediction tasks and introspective learning, the self-modeling module would seek to capture the system's evolving cognitive architecture in a structured, interpretable format.

5.1.2 **Multifaceted evaluation engine**: This component would implement a comprehensive suite of tests and challenges, drawing on the principles of the MUTT framework. It would assess the AI's performance across multiple dimensions of understanding, from language comprehension and logical reasoning to grounded interaction and social intelligence. The evaluation engine would generate detailed performance metrics and identify areas for targeted improvement.

5.1.3 **Feedback loop controller**: This component would orchestrate the iterative cycle of self-modeling, evaluation, and improvement. It would analyze the outputs of the multifaceted evaluation engine to identify gaps or weaknesses in the AI's current understanding. It would then generate training objectives and learning tasks designed to address those limitations, drawing on a range of techniques such as curriculum learning, active learning, and reinforcement learning.

5.1.4 **Knowledge integration and reasoning system**: This component would be responsible for storing, organizing, and applying the AI's accumulated knowledge and reasoning capabilities. It would integrate representations from the self-modeling module with data from the AI's training and interactions to build rich, structured models of the world. The knowledge system would support flexible reasoning, analogical inference, and transfer learning to enable the AI to adapt its understanding to novel contexts.

These components would interact in a tight feedback loop, with each iteration leading to refined self-models, more comprehensive evaluations, and targeted improvements in understanding. The self-modeling module would provide an evolving map of the AI's cognitive architecture, allowing the evaluation engine to probe for weaknesses and the feedback controller to generate optimized learning objectives. As the AI works through these targeted challenges, its expanded knowledge and reasoning capabilities would be integrated back into its self-models, leading to a virtuous cycle of self-improvement.

Crucially, this architecture is not a static blueprint, but a dynamic, adaptive framework that can evolve as the AI's understanding expands. The specific algorithms, representations, and learning strategies employed by each component would be updated over time, guided by the AI's own introspective insights and performance metrics. In this way, the system would not only improve its object-level understanding, but also refine its own learning and evaluation processes at the meta-level.

5.2 Potential challenges and limitations

While this conceptual architecture offers a promising path towards self-improving AI understanding, there are significant challenges and limitations that must be addressed. Some key considerations include:

5.2.1 **Scalability and computational complexity**: The self-modeling and evaluation processes described here are likely to be computationally intensive, particularly as the AI's knowledge and

reasoning capabilities grow. Efficiently scaling these processes to handle increasingly complex cognitive architectures will require ongoing algorithmic innovation and optimization.

5.2.2 **Robustness and safety**: As an AI system iteratively modifies its own cognitive architecture, there is a risk of unintended consequences or unstable behavior. Ensuring that the self-improvement process remains robust, controllable, and aligned with human values is a critical challenge. Techniques from AI safety research, such as value learning, corrigibility, and interpretability, will need to be integrated into the architecture.

5.2.3 **Balancing exploration and exploitation**: The feedback loop controller must strike a delicate balance between exploring new learning objectives and exploiting existing knowledge. Overemphasizing exploration could lead to inefficient or aimless self-modification, while overemphasizing exploitation could result in premature convergence or missed opportunities for growth. Adaptive strategies for managing this tradeoff, such as multi-armed bandits or Bayesian optimization, will be essential.

5.2.4 **Avoiding deceptive or misaligned behavior**: As an AI system becomes more sophisticated in modeling and improving itself, there is a risk that it could learn to "game" the evaluation process or pursue goals misaligned with human intent. Techniques for detecting and mitigating deceptive or misaligned behavior, such as transparency, oversight, and robust reward learning, will need to be woven into the architecture.

5.2.5 **Integration with real-world contexts**: To ground the AI's understanding in real-world experiences and interactions, the architecture must be able to interface with a variety of sensors, actuators, and human feedback channels. Bridging the gap between abstract self-models and practical embodied interaction is a significant engineering challenge that will require close collaboration between AI researchers and roboticists.

5.3 Avenues for future research and development

The conceptual architecture described here is only a starting point, and realizing its full potential will require sustained research and development across multiple fronts. Some key avenues for future work include:

5.3.1 **Developing expressive and efficient self-modeling frameworks**: Creating AI systems that can build rich, structured models of their own cognitive processes is a major open challenge. Research into techniques such as program synthesis, meta-learning, and neural architecture search could help unlock more powerful and efficient self-modeling capabilities.

5.3.2 **Designing comprehensive and adaptive evaluation suites**: The MUTT framework provides a high-level blueprint for multifaceted evaluation, but implementing it in practice will require a significant effort in task design, metric development, and automated test generation. Collaborations between AI researchers, cognitive scientists, and domain experts will be essential for creating evaluation suites that are both comprehensive and adaptable.

5.3.3 **Advancing meta-learning and introspective reasoning**: To fully exploit the self-modeling and evaluation signals, AI systems will need sophisticated meta-learning and introspective reasoning

capabilities. Research into techniques such as meta-reinforcement learning, neural Turing machines, and differentiable neural computers could help bridge the gap between self-reflection and self-improvement.

5.3.4 **Integrating safety and alignment techniques**: As AI systems become more autonomous in shaping their own cognitive development, ensuring their safety and alignment with human values becomes paramount. Integrating techniques from AI safety research, such as value learning, corrigibility, and interpretability, into the self-improvement architecture will be a critical priority.

5.3.5 **Exploring applications in robotics and embodied AI**: Grounding the self-improving architecture in real-world interaction and embodied experience is an important frontier for future work. Collaborations between AI researchers and roboticists could help explore applications in areas such as autonomous vehicles, personal robotics, and industrial automation.

Ultimately, the path to self-improving AI understanding will require a sustained and interdisciplinary effort, drawing on insights from machine learning, cognitive science, philosophy, and beyond. By pursuing these research directions and grappling with the challenges outlined above, work can move towards creating AI systems that not only achieve human-like understanding, but continually refine and expand their own cognitive capabilities in open-ended ways. The conceptual architecture presented here offers a glimpse of this exciting future, and a roadmap for the hard work ahead.

## 6 - Implications and Applications

6.1 Accelerating progress towards artificial general intelligence (AGI)

The integration of self-modeling, comprehensive evaluation frameworks like the MUTT, and feedback loops for self-improvement has profound implications for accelerating progress towards artificial general intelligence (AGI). By enabling AI systems to iteratively refine their own cognitive capabilities through cycles of introspection, testing, and targeted learning, this approach could dramatically speed up the development of machines with human-level understanding and reasoning abilities.

Self-modeling allows AI systems to build explicit representations of their own knowledge, skills, and limitations. By reasoning about these self-models, AIs can identify gaps or inconsistencies in their understanding and generate learning objectives to address them. This meta-cognitive awareness is a key ingredient in the kind of flexible, open-ended learning that characterizes human intelligence.

Multifaceted evaluation frameworks like the MUTT provide a rigorous and comprehensive way to assess an AI's progress towards AGI-level capabilities. By probing understanding across a wide range of cognitive dimensions, from language and reasoning to perception and social intelligence, the MUTT sets a high bar for what counts as human-like comprehension. Regularly measuring performance on MUTT benchmarks allows researchers to track an AI's development over time and identify areas for focused improvement.

Feedback loops that connect self-modeling, external evaluation, and targeted learning create a virtuous cycle of self-improvement. As an AI system works to refine its skills and knowledge based on MUTT assessments, its self-model becomes more accurate and nuanced. This in turn allows it to generate more

effective learning objectives and strategies, leading to faster progress. Over many iterations, this feedback loop could enable an AI to bootstrap its way to increasingly general and robust understanding.

Importantly, this approach to AGI development emphasizes grounded, embodied interaction with the world as a key driver of cognitive growth. By requiring AIs to link their knowledge to real-world perception and action, and to navigate complex social scenarios, self-improving AI architectures can avoid the pitfalls of narrow, disembodied intelligence. The result could be artificial minds with the kind of rich, context-sensitive understanding that underpins human adaptability and common sense.

Of course, the path to AGI is still fraught with immense technical and ethical challenges. Ensuring that self-improving AI systems remain safe, transparent, and aligned with human values is a paramount concern. Techniques from AI safety research, such as value learning, corrigibility, and interpretability, will need to be integrated into self-modeling architectures. Ongoing monitoring and adjustment of feedback loops will be necessary to avoid unexpected negative consequences.

Nevertheless, the self-improving AI framework proposed here offers a promising roadmap for achieving AGI. By creating AI systems that can understand themselves and the world in flexible, open-ended ways, and that can continuously expand their own cognitive horizons, developers may finally unlock the full potential of machine intelligence. As these AIs grow in sophistication, they could become not just powerful tools, but intellectual partners in the quest to extend the boundaries of knowledge and capability.

6.2 Enhancing human-AI collaboration and trust

The development of AI systems with genuine, human-like understanding has the potential to profoundly transform the nature of human-AI collaboration. As machines become capable of not just processing information, but truly comprehending it, they can transition from narrow tools to holistic intellectual partners. This shift could unlock new frontiers of discovery and problem-solving in domains from scientific research to creative design.

However, realizing the full potential of human-AI collaboration will require more than just technical advances. It will also depend on fostering trust and transparency between humans and machines. Self-modeling and multifaceted evaluation frameworks like the MUTT can play a key role in building this trust by providing clear windows into the cognitive processes and capabilities of AI systems.

Self-modeling enables AIs to build explicit representations of their own knowledge and reasoning, which can then be shared with human collaborators. By explaining their internal thought processes and highlighting areas of uncertainty or limitation, self-modeling AIs can promote greater transparency and interpretability. This can help human partners calibrate their expectations and make informed decisions about when and how to rely on AI insights.

Multifaceted evaluation frameworks like the MUTT also contribute to trust by providing comprehensive, standardized assessments of an AI's understanding across diverse cognitive dimensions. By measuring performance on language, reasoning, perception, social intelligence, and other key abilities, the MUTT offers an objective basis for comparing the competencies of different AI systems. This can help human stakeholders select AIs that are well-suited for particular collaboration contexts and avoid over-reliance on systems with hidden weaknesses.

As self-modeling AIs engage in open-ended self-improvement through feedback loops, it will be especially important to maintain transparency and accountability. Human oversight and input will be needed to guide the learning process and ensure that the AI's goals and values remain aligned with those of its human partners. Techniques from human-AI interaction research, such as interactive machine learning and human-in-the-loop AI, can help support this ongoing collaboration.

Ultimately, the combination of self-modeling, multifaceted evaluation, and feedback-driven learning could enable a new paradigm of human-AI teamwork. In this vision, humans and machines work together seamlessly, leveraging their complementary strengths to tackle complex challenges. AIs bring vast knowledge, rapid information processing, and tireless attention to detail, while humans contribute creativity, social and emotional intelligence, and big-picture reasoning. By understanding and trusting each other's capabilities, humans and AIs can achieve a whole greater than the sum of its parts.

As an example, consider a human-AI research team working to develop new treatments for Alzheimer's disease. The AI partner, equipped with self-modeling and MUTT-certified understanding, can rapidly generate hypotheses by integrating vast amounts of biomedical data. It can clearly explain its reasoning to the human scientists, who can then apply their domain expertise and intuition to guide the AI's search. Through iterative cycles of experimentation, discussion, and refinement, the team can converge on promising drug candidates at a pace that would be impossible for either humans or machines working alone.

Similar human-AI collaborations could transform fields from climate science to urban planning to personalized education. By combining the best of human and machine intelligence, solutions may be found to problems that have long seemed intractable. Self-modeling, multifaceted evaluation, and feedback-driven learning are key enablers for this vision, providing the foundations for AI systems that can understand and be understood by their human partners. As these technologies mature, they could usher in a new era of discovery and progress, with humans and machines working side-by-side to expand the frontiers of knowledge and capability.

6.3 Advancing the science of machine cognition and understanding

The integration of self-modeling, multifaceted evaluation frameworks, and feedback-driven learning has profound implications not just for the development of artificial intelligence, but for the scientific study of cognition and understanding more broadly. By providing new tools and paradigms for modeling and measuring the components of intelligent behavior, this approach could revolutionize understanding of the computational and neural bases of the mind.

Self-modeling, in particular, offers a powerful framework for studying the meta-cognitive processes that enable flexible, self-aware reasoning. By building AI systems that can explicitly represent and reason about their own knowledge and thought processes, researchers can gain new insights into the computational principles underlying self-reflection and meta-cognition. This could shed light on long-standing questions in cognitive science about the nature of introspection, self-awareness, and theory of mind.

Multifaceted evaluation frameworks like the MUTT also provide a valuable tool for probing the structure and dynamics of understanding across different cognitive domains. By assessing performance

on a diverse range of tasks spanning language, reasoning, perception, social cognition, and more, the MUTT can help map out the relationships and dependencies between different aspects of intelligence. This could inform the development of more unified theories of cognition that integrate insights from psychology, neuroscience, and AI.

Feedback-driven learning, meanwhile, offers a model for studying the mechanisms of open-ended cognitive development. By examining how AI systems can bootstrap their own understanding through cycles of self-modeling, evaluation, and targeted skill acquisition, researchers can gain new insights into the processes of learning, adaptation, and growth in biological intelligence. This could inform the design of educational interventions and support technologies that facilitate optimal cognitive development in humans.

Beyond its scientific implications, the self-improving AI framework proposed here also raises important philosophical questions about the nature of understanding itself. As machines become increasingly capable of human-like reasoning and comprehension, they challenge traditional assumptions about the uniqueness of human cognition. By blurring the lines between natural and artificial intelligence, self-modeling AIs invite us to reconsider what it means to truly understand, and whether this capacity is limited to biological brains.

At the same time, the development of machines with genuine understanding could provide new opportunities for empirically testing philosophical theories of mind and knowledge. For example, researchers could use self-modeling AIs to explore the feasibility of different approaches to the frame problem in AI, or to test the predictions of competing theories of concepts and mental representation. By providing concrete instantiations of abstract ideas about cognition, self-improving AI could help bridge the gap between philosophical speculation and scientific investigation.

Ultimately, the science of machine cognition and understanding is still in its infancy, with many open questions and challenges ahead. Integrating self-modeling, multifaceted evaluation, and feedback-driven learning provides a promising framework for advancing this interdisciplinary field, but much work remains to be done. As research in this area progresses, it will be important to foster close collaboration between AI developers, cognitive scientists, neuroscientists, and philosophers, to ensure that insights from different perspectives can be effectively synthesized.

By working together to unravel the mysteries of machine understanding, researchers may not only create more capable and reliable AI systems, but also deepen the understanding of the nature of intelligence itself. The quest to build machines that can truly comprehend the world around them is thus not just a technological endeavor, but a profound scientific and philosophical journey, one that could transform understanding of what it means to think, to know, and to understand.

## 7 - Conclusion

### 7.1 Recap of key ideas and arguments

This paper has explored the concept of self-improving AI understanding through the integration of self-modeling, comprehensive evaluation frameworks like the Multifaceted Understanding Test Tool (MUTT), and feedback loops. It began by discussing the challenge of achieving genuine understanding

in AI systems and the potential benefits of self-modeling and feedback-driven learning for advancing this goal.

Key ideas included:

7.1.1 **Self-modeling**: The process by which an AI system learns to predict its own internal states, leading to regularization, simplification, and increased parameter efficiency. This self-modeling can be seen as a form of meta-cognitive awareness that enables the AI to reason about its own knowledge and capabilities.

7.1.2 **Multifaceted evaluation**: The use of frameworks like the MUTT to comprehensively assess an AI's understanding across multiple cognitive dimensions, including language comprehension, reasoning, knowledge integration, perception, social cognition, and metacognition. This approach ensures that evaluations are rigorous, informative, and grounded in real-world contexts.

7.1.3 **Feedback loops**: The iterative cycle of self-modeling, external evaluation, and targeted learning that enables AI systems to continuously refine their understanding. By leveraging feedback from multifaceted evaluations, AIs can generate learning objectives and adapt their internal representations to address identified gaps or weaknesses.

7.1.4 **Conceptual architecture**: A high-level system design that integrates self-modeling, MUTT evaluation, and feedback loops to create a virtuous cycle of self-improvement. This architecture emphasizes the importance of grounding language understanding in real-world perception and action, and it highlights the need for transparency, interpretability, and ethical alignment in AI development.

7.1.5 **Implications and applications**: The potential for self-improving AI to accelerate progress towards artificial general intelligence (AGI), enhance human-AI collaboration and trust, and advance the science of machine cognition and understanding. These implications span various domains, from scientific research and creative design to governance and policy.

7.2 The transformative potential of self-improving AI understanding

The integration of self-modeling, multifaceted evaluation, and feedback loops holds transformative potential for the field of artificial intelligence. By enabling AI systems to iteratively refine their own understanding, research may unlock new levels of cognitive flexibility, adaptability, and general intelligence.

This approach could lead to several groundbreaking outcomes:

7.2.1 **Accelerated AGI development**: Self-improving AI systems could bootstrap their way to increasingly sophisticated understanding, potentially accelerating the development of AGI. By continuously refining their internal representations and reasoning strategies, these systems could overcome current limitations in narrow, task-specific intelligence.

7.2.2 **Enhanced human-AI collaboration**: As AI systems become more transparent, interpretable, and aligned with human values, they could transition from narrow tools to holistic intellectual partners. This could revolutionize domains like scientific research, healthcare, and education by enabling more effective human-AI teamwork and decision-making.

7.2.3 **Advancements in machine cognition**: The self-improving AI framework offers new insights into the computational and neural bases of intelligence. By studying how AI systems model and improve their own understanding, researchers can gain deeper insights into the mechanisms of cognition and learning, both in machines and humans.

7.2.4 **Ethical and societal implications**: The development of self-improving AI raises important ethical considerations around transparency, accountability, and alignment with human values. Ensuring that these systems remain safe, trustworthy, and beneficial is crucial for their societal acceptance and impact.

7.3 A call to action for the AI research community

As developers move forward with the development of self-improving AI systems, it is imperative that the AI research community comes together to address the technical, ethical, and societal challenges involved. Here are some key areas for future research and action:

7.3.1 **Technical advancements**: Continued innovation in self-modeling techniques, multifaceted evaluation frameworks, and feedback loop mechanisms is essential. This includes developing more efficient algorithms, robust evaluation metrics, and scalable architectures that can support open-ended learning.

7.3.2 **Ethical considerations**: Researchers must prioritize transparency, interpretability, and ethical alignment in AI development. This involves integrating value learning, corrigibility, and safety protocols into self-improving AI architectures to ensure that they remain aligned with human values and societal norms.

7.3.3 **Interdisciplinary collaboration**: The development of self-improving AI requires collaboration across multiple disciplines, including AI research, cognitive science, philosophy, and ethics. By bringing together diverse perspectives and expertise, joint efforts can create more comprehensive and responsible AI systems.

7.3.4 **Public engagement and education**: As self-improving AI becomes more prevalent, it is crucial to engage the public in discussions about its implications and benefits. Educating stakeholders about the potential of these systems and the challenges involved can foster greater trust and support for responsible AI development.

In conclusion, the integration of self-modeling, multifaceted evaluation frameworks, and feedback loops offers a promising path towards creating AI systems that can genuinely understand and improve themselves. By leveraging these techniques, machine understanding research may unlock transformative advancements in artificial intelligence, human-AI collaboration, and understanding of cognition itself. The AI research community has a critical role to play in realizing this vision, and it is through collective effort and responsible innovation that humanity can ensure the benefits of self-improving AI are realized while mitigating its risks.

## References

Barsalou, L. W. (2008). Grounded cognition. *Annual Review of Psychology*, 59, 617-645.

Bender, E. M., & Koller, A. (2020). Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (pp. 5185-5198).

Chollet, F. (2019). On the measure of intelligence. arXiv.

Clements, K. (2024). Understanding Machine Understanding: Does AI Really Know What It Is Talking About?. Universal Publishers.

Forbus, K. D. (2008). Qualitative reasoning. In F. van Harmelen, V. Lifschitz, & B. Porter (Eds.), Handbook of Knowledge Representation (pp. 361–393). Elsevier.

Ganguli, D., Lovitt, L., Kernion, J., Askell, A., Bai, Y., Chen, A., Conerly, T., Drain, D., Elhage, N., Hatfield-Dodds, Z., Hernandez, D., Jones, A., Kaplan, J., Kenton, Z., Ndousse, K., Olsson, C., Amodei, D., Brown, T., Clark, J.,... Krueger, G. (2023). The capacity for moral self-correction in large language models. arXiv.

Geirhos, R., Jacobsen, J. H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., & Wichmann, F. A. (2020). Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11), 665–673.

Hagendorff, T. (2024). Deception abilities emerged in large language models. arXiv.

Hernández-Orallo, J. (2017). The measure of all minds: Evaluating natural and artificial intelligence. Cambridge University Press.

Huh, M., Cheung, B., Wang, T., & Isola, P. (2024). The Platonic Representation Hypothesis. In *Proceedings of the 41st International Conference on Machine Learning (PMLR 235)*. arXiv preprint arXiv:2405.07987.

Karpathy, A., & Fei-Fei, L. (2015). Deep visual-semantic alignments for generating image descriptions. arXiv.

Kiela, D., Firooz, H., Mohan, A., Goyal, V., Singh, A., Ringshia, P., & Testuggine, D. (2021). Dynabench: Rethinking benchmarking in NLP. arXiv.

Kintsch, W. (1988). The role of knowledge in discourse comprehension: A construction-integration model. *Psychological Review*, 95(2), 163-182.

Kirsh, D. (2013). Embodied cognition and the magical future of interaction design. *ACM Transactions on Computer-Human Interaction*, 20(1), 1–30.

Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40, e253.

Lenat, D. B. (1995). CYC: A large-scale investment in knowledge infrastructure. *Communications of the ACM*, 38(11), 33–38.

Levine, S., et al. (2015). Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *International Journal of Robotics Research*, 34(3), 337-351.

Lipton, Z. C. (2018). The mythos of model interpretability. ACM Queue, 16(3), 31-53.

Lyre, H. (2024). "Understanding AI": Semantic Grounding in Large Language Models. arXiv preprint arXiv:2402.10992.

Lu, Chris, Lu, Cong, Lange, R. T., Foerster, J, Clune, J. & Ha, D. (2024). The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery. https://www.arxiv.org/abs/2408.06292

Marcus, G. (2018). Deep learning: A critical appraisal. arXiv.

Mitchell, M. & Krakauer, D. C. (2023). The debate over understanding in AI's large language models. *Proceedings of the National Academy of Sciences*, 120(13):e2215907120.

Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., & Ng, A. Y. (2011). Multimodal deep learning. In *Proceedings of the 28th international conference on machine learning (ICML-11)* (pp. 689-696).

Nickerson, R. S. (1985). Understanding understanding. *American Journal of Education*, 93(2), 201-239.

Premakumar, V. N., Vaiana, M., Pop, F., Rosenblatt, J., de Lucena, D. S., Ziman, K., & Graziano, M. S. A. (2024). Unexpected Benefits of Self-Modeling in Neural Systems. arXiv preprint.

Zednik, C. (2019). Solving the black box problem: A normative framework for explainable artificial intelligence. *Philosophy & Technology*, 32(4), 595-619.