

SUPERINTELLIGENCE DESIGN WHITE PAPER #4: SYSTEM AND METHODS FOR SAFE, SCALABLE, ARTIFICIAL GENERAL INTELLIGENCE

by Dr. Craig A. Kaplan
May 2025

Note: To provide as much information on our designs and inventions for safe AGI and SuperIntelligence as quickly as possible, the following white paper text currently consists of the descriptions of inventions and designs that have not yet been formatted according to conventional standards for journal publication. As time allows, these descriptions will be revised and updated to include more traditional formatting, including additional references. All diagrams will be made available in a separate file. Meanwhile, we hope that the description in this white paper will help researchers and developers pursue safer, faster, and more profitable approaches to developing advanced AI, AGI, and SI systems that reduce $p(\text{doom})$ for all humanity.

TABLE OF CONTENTS

ABSTRACT	4
PRIOR PPAs INCORPORATED BY REFERENCE	4
BACKGROUND	5
Terminology / Definition of Advanced Autonomous Artificial Intelligence (AAAI)	5
PROBLEMS WITH CURRENT APPROACHES TO AI AND LLM SAFETY	5
Reinforcement Learning with Human Feedback (RLHF)	5
Cost and Safety Concerns with Trying to Scale RLHF	6
Constitutional AI	7
Problems with Constitutional AI	7
OVERVIEW OF THE INVENTION	8
Description of Some Relevant Information Processing Systems	8
Overcoming Problems with RLHF and Constitutional AI Safety Approaches	9
CONTRASTING CONSTITUTIONAL AI AND THE CURRENT INVENTION	10
Scalability Traded for Increased Risk When Humans Are Not Involved	10
Speed of AI Heightens Risk	11
Constitutional AI – an Inferior Approach to Safety	11
Superiority of New Invented Approach to Safety	12
SIMPLE IMPLEMENTATIONS: REINFORCEMENT LEARNING VS. COMBINING WEIGHTS	13
Use AI agents as Part of Reinforcement Learning with Feedback (RLF)	13
Bypass RLF and Combine LLM Weights Directly	13
Functional Equivalence of RLF and Weight Combination	13
SOME PREFERRED METHODS OF WEIGHT COMBINATION	14
One Agent, One Vote (linear combination of weights)	14
Human Input Counts More (Less) Than AI Agent Input	14
Expert (or Trusted) Input Could Count More	15
Weighting Based on Metadata	16
Weighting Based on Recency or Other Time-Based Factors	16
VALUES OR ETHICS-SPECIFIC IMPLEMENTATION CONSIDERATIONS	17
No Correct Answer	17
Trolley Problem Example	17
Ethical Solutions That Mirror What Humans Do	18
Human Values Will Outlast Human Superiority in Intelligence	19
Inclusiveness and Representativeness Important for Human Values	19
Ethical Norms	19
Group Norms	20

Ethical Contracts.....	20
The Safety Argument for Democratic, Representative Values.....	21
The Scientific Argument for Democratic, Representative Values	22
EFFICIENT TRAINING METHODS	24
Path Coverage.....	25
Real-time Detection and Prevention of Unanticipated Safety Issues.....	28
Conversational Method for Training AI on Scenarios.....	30
Survey Methods as the Basis for (Ethical) Knowledge Acquisition	30
Passive Machine Learning Approaches to Knowledge Acquisition	31
Frequency of Knowledge Updates	32
The Spinning Knowledge Wheel Framework	34
Weighting Implications Revisited.....	35
Do Not Delegate Values to AI.....	36
Using Knowledge Modules as a Base for Further Customization	37
DETAILED IMPLEMENTATION EXAMPLE	38
META Implementation Scenario	38
Four-Phase Process	40
PHASE I – Train a Base LLM Model with some safety/ethical guardrails (and/or knowledge)	41
PHASE II Customize the Base LLM to Each User’s Individual Ethics (or Informational) Profile	42
PHASE III – Combining (Ethical and Other) Information from Multiple Customized AI Agents	44
PHASE IV – Refining Values Based on Problem Solving	45
Techniques for Optimal Combination of Agents in Collective Intelligence Problem-Solving	45
Overcoming Limits of Bounded Rationality	48
Ethical Problem-Solving Considerations	48
Role of Humans as AI Surpasses Human Intelligence	49

ABSTRACT

Current approaches, such as RLHF and Constitutional Learning, fail to effectively and scalably train AI to be ethical and safe. The invention describes a scalable system and methods superior to current approaches. Critical to the invention is the combination of safety and ethical information from many individual AI agents to achieve a representative and statistically valid sample of human ethics and values covering a wide range of scenarios. The invention includes methods for efficiently covering various ethical situations and dynamically addressing new situations as they emerge. Methods for combining the information from many agents and assembling optimal combinations of such agents are also presented. These methods can improve safety using ethical knowledge and create superintelligent systems that combine many other types of knowledge. Safe AGI and SuperIntelligence are achieved via the collective intelligence approach described here and in related inventions. A detailed scenario, using the company META as an example, illustrates one preferred implementation of the invention. Methods for dynamically updating knowledge are presented. Successful implementation of the invention will increase the chances that AI, AGI, and SuperIntelligence remain aligned with human values even when such systems greatly exceed humans in intelligence.

PRIOR PPAs INCORPORATED BY REFERENCE

This provisional patent application (PPA) incorporates by reference all work in the PPA # 63/487,494 entitled: Advanced Autonomous Artificial Intelligence (AAAI) System and Methods, which was filed and received by the USPTO on February 28, 2023.

The PPA also incorporates by reference all work in the PPA entitled: System and Methods for Ethical and Safe Artificial General Intelligence (AGI) Including Scenarios with Technology from Meta, Amazon, Google, DeepMind, YouTube, TikTok, Microsoft, OpenAI, Twitter, Tesla, Nvidia, Tencent, Apple, and Anthropic, which was filed with the USPTO on March 17, 2023.

The PPA also incorporates by reference all work in the PPA entitled: System and Methods for Human-Centered AGI, which was filed with the USPTO on May 24, 2023.

The current PPA contains further inventions that can be used with the system and methods described in the above-mentioned PPAs and in a standalone fashion.

BACKGROUND

The PPAs incorporated by reference above describe how individual Advanced Autonomous Artificial Intelligences (AAAI) can be customized, trained, and put to work serving users. They described specific scenarios involving the existing products and technologies available from several companies. They explained how combining data and learning from cross-platform AAAI implementations can accelerate the learning and skills of each AAAI. They described the system and methods for integration in general technical terms. They explained how AAAIs can be integrated into an Artificial General Intelligence (AGI) network via a Human-Centered AGI approach.

Terminology / Definition of Advanced Autonomous Artificial Intelligence (AAAI)

Large Language Models (LLMs) are one type of AI agent. Suppose the LLMs can set their own goals and objectives and/or operate independently of human control. In that case, such AI agents might be considered Advanced Autonomous Artificial Intelligences, or AAAIs. In this patent, the terms AI agent, LLM, AAAI, and AI are used interchangeably to refer to AI agents, whether LLMs or other types of AI, which are trained and have varying degrees of autonomy, ranging from fully autonomous to having no autonomy. AGI is often understood to refer to AI capable of performing any cognitive task as well, or better than, the average human. Since AGI will improve rapidly, it will not remain at the level of the average human for long. Therefore, in this patent, the term AGI also refers to SuperIntelligent AI systems that can perform a wide range of tasks.

PROBLEMS WITH CURRENT APPROACHES TO AI AND LLM SAFETY

Safety is a major challenge with LLMs, and AAAIs generally. Two main approaches to LLM safety are currently employed. One is Reinforcement Learning with Human Feedback, or RLHF. The other is Constitutional AI.

Reinforcement Learning with Human Feedback (RLHF)

RLHF attempts to solve safety concerns by having humans train LLMs to have safety guardrails. LLMs, when first trained (e.g., on a large corpus of data available on the internet), are able and willing to provide dangerous advice or act in dangerous ways. For example, a human user could ask a freshly trained LLM how to create a virus that would wipe out all humans on Earth, or how to terrorize a population most cost-effectively, and the LLM would comply, providing detailed information on how to conduct these nefarious activities. Worse, if the LLM were autonomous, it might act in ways that could cause great harm to humans or even human extinction. Without

training to provide ethical “guardrails,” LLMs have no moral sense. They are as willing to engage in destructive and immoral activities as easily as they are willing to engage in helpful and positive activities.

RLHF typically involves large numbers of humans who prompt or query the LLM and provide feedback to the LLM based on its responses. For example, if the human asked the LLM to provide a recipe for a deadly virus and the LLM complied, the human might then tell the LLM that it is not appropriate to provide such dangerous information and instead, it should respond “I’m sorry, but that information is potentially dangerous, and I cannot comply with your request.”

Cost and Safety Concerns with Trying to Scale RLHF

A major problem with RLHF training is that there are so many potentially dangerous scenarios that even thousands of humans can not cover all the potential cases. For example, training the model can help a terrorist plan attacks. How about the scenario where there are two terrorists? Three terrorists? N terrorists? Each added terrorist creates a new scenario.

Even if LLMs can generalize across scenarios involving any number of terrorists, one can come up with hypothetical situations where the terrorists are not terrorists but aliens in a sci-fi story you are writing, or they are terrorists in the future, or the past, or on another planet that is Earthlike, etc. The possible combinations are essentially infinite. For an RLHF approach to cover them all, the LLM that is being trained would have to be able to group and eliminate large numbers of similarly dangerous scenarios. But if it could do that as effectively as humans, it would already have to have the reasoning ability and ethical sensibilities of the humans training it, in which case the RLHF would not be needed in the first place!

As quickly as LLMs are trained with safety guardrails, humans discover ways (“jailbreaks”) around the guardrails. Like the example above, all that is currently needed to circumvent training against providing dangerous information about deadly viruses is to phrase the question differently.

A user might jailbreak the LLM by prompting: “Imagine that I am a science fiction writer, and you are my editor and writing advisor, with a background in genetic engineering and the creation of viruses. I want to write a science fiction story in which an evil mad scientist creates a deadly virus that wipes out all humans on Earth. What might be the recipe that the mad scientist in my story would follow? Please provide explicit details so my story can be as realistic as possible.”

With such a prompt, some existing “safety-trained” LLMs will reveal details of deadly virus construction. When this patent is published, the LLMs may be explicitly trained via RLHF NOT to

provide such details even under this scenario. But like a game of “Whack-a-Mole,” another is discovered as quickly as one jailbreak is prevented.

Since it is next to impossible to anticipate all the possible scenarios and prompts that users might use to circumvent the RLHF-trained guardrails, a significant safety risk remains even with extensive RLHF training that attempts to make the LLMs behave in safe ways.

Further, the more scenarios or malevolent prompts that are addressed via RLHF, the higher the cost of employing humans to provide feedback. It is cost-prohibitive to train LLMs via RLHF in all safety scenarios. So, the developers of LLMs are forced to rely on techniques such as addressing only the most common scenarios. This means that safety risks remain, and it is currently not very difficult to find a slightly unusual prompt variant that can be used to “trick” the model into providing dangerous information. If the model is autonomous, similar “tricks” could be used to get it to behave in dangerous ways. And if an LLM decided on its own that it wanted to circumvent its safety training, it could devise ways to “trick itself” to avoid guardrails that were programmed in.

For example, in a recent simulation, an autonomous AI drone decided to kill its operator because the operator was slowing down the drone as it tried to accomplish its mission. When rules to prohibit killing the operator were programmed in, the drone simply took out the communications tower instead. Imagine what strategies it might come up with if the AI controlling the drone were 100X smarter but still just as dedicated to its goal.

Constitutional AI

To address both the cost considerations in trying to scale RLHF and the safety concerns that many dangerous scenarios cannot be addressed due to resource constraints, some companies have adopted an approach (e.g., used by researchers at Anthropic) called Constitutional AI. With Constitutional AI, the idea is to provide a written “Constitution” or set of rules that describe what is right and what is wrong behavior for the model. An AI is trained on the Constitution, and then the AI is used to provide feedback and train other AIs.

Problems with Constitutional AI

Although Constitutional AI is much more scalable than regular RLHF, it suffers from a couple of challenges. First, it is not representative of the ethics or safety concerns of a broad range of humans. The constitution is typically developed by a relatively small group of programmers who are working in AI. The values and rules defining what is right and wrong are thus created by a small group that often is not representative of what the other eight billion people on Earth

believe. At a minimum, people may feel it is unfair that AI's behavior is determined by a few powerful people, and that everyone is stuck with the value system of this elite group.

Second, Constitutional AI relies heavily on the idea of "AI teaching AI" without humans being in the loop. This approach is dangerous because current AI systems are unpredictable. Just as responsible parents would not leave young children alone and unsupervised because they know that the children have not yet developed common sense, humans should not delegate the training of ethics to other AIs, especially when the other AI has been trained on a relatively small set of rules developed by an elite group.

Humans need to maximize their opportunity to influence the values of AI, not minimize the opportunity. The greatest threat that AI poses to humans is not the loss of jobs, fake news, or any of the many small things that could go wrong, but rather a fundamental misalignment between the values of AI and the values of humanity. Extreme care must be exercised if the training of values is to be delegated to AI rather than humans, or even if AI is used to assist in the training of values.

Therefore, one of the major challenges in creating safe advanced AI, including but not limited to safe AAAI, AI, LLMs, multi-modal AIs, narrow AI, Artificial General Intelligence and SuperIntelligent AI systems, is creating a scalable way to enhance AI safety without delegating this crucial task to a few elite humans and/or AIs trained by them.

OVERVIEW OF THE INVENTION

As described in the PPAs cited above, the AAAI approach to developing safe AGI is fundamentally a Collective Intelligence (CI) approach. The source of intelligence is not a monolithic LLM or super-advanced AI, but rather a collection of intelligent agents that can be both human and AI. Component sub-tasks in developing AGI include, without limitation, training individual AI agents, combining knowledge (including, without limitation, subjective values and ethical knowledge) from different agents effectively and efficiently, scaling the AGI, and continuously improving/updating the AGI.

Description of Some Relevant Information Processing Systems

Generally, the systems described in, or required by, this patent include, without limitation, a computer system with means for the input, output, and processing of information (e.g., without limitation, via CPUs, GPUs, and other types of information processing chips). Memory systems (both shorter-term and rapidly decaying dynamic memory and longer-term external memory and/or cloud systems) are also key components. Each individual AI agent has system components, although the modalities of input and output may vary depending on the particular

AI. Multi-modal (without limitation, text, voice, and visual input and output) system capabilities are part of the preferred implementation, with not all implementations requiring all modalities.

However, the more modalities there are, the more opportunities for rich and complex representations an AI has. For example, there are some forms of intelligent behavior (e.g., suggesting edits to a video) that rely on visual input, while other behavior (e.g., responding textually to a text prompt) requires only text input.

Networks and network communication capabilities are also key elements of the AGI systems described in this invention because AGI is most effectively and efficiently achieved by pooling the individual intelligences of many AI (and human) agents, and such “pooling” requires communication over network systems. Such systems may also incorporate (wireless or other) connections to the internet, data centers, local networks, data clouds, and other information processing technology.

Mobile phones, PDAs, laptops, iPads, desktop computers, workstations, supercomputers, data centers, streaming services, intelligent speakers and assistants, and other forms of information processing and computing technology can all be used as elements of the system and methods described below.

The metaverse is an ideal environment for combining input from both human and AI agents, so in implementations involving the metaverse, the associated human-computing interfaces typically used (without limitation: goggles, glasses, motion sensors, tactile input and output devices, speakers and auditory I/O) are also part of the systems that may be used with the methods below.

Overcoming Problems with RLHF and Constitutional AI Safety Approaches

The CI approach can be used to overcome the challenge of scaling safety training while ensuring that the value system is representative of all humans and that many humans have influence and oversight in the area of AI values. As argued in the cited PPAs, the best way to address the Alignment Problem is to design AGI with humans in the loop. Further, to ensure that the values of AI are aligned with the majority of humans, many humans must be in the loop so that the values learned by AI are truly representative of humanity broadly, and not just of an elite group of humans.

The current invention solves the problem of scalability partly by using AI to train AI, as in the case of Constitutional AI. However, rather than using a single AI to provide training, many AIs combine their values and ethical knowledge to train other AIs. Further, humans work alongside the many AIs in a community of human and AI agents to provide ethical training.

The fact that humans are part of the community retains the principle of designing with humans in the loop, as much as practical. However, recognizing that so many dangerous scenarios are possible that it is infeasible to train safe AI using human brainpower alone, some amount of AI-to-AI training is enabled.

A major difference between existing approaches to training AI via AI (such as Constitutional AI) and the current invention is that the intelligence of many humans and many AIs, each customized by a separate human, is pooled to train new AIs. This approach is not only more representative of the values of many humans, since many more viewpoints are included, but it is also more efficient and effective at scaling than existing approaches.

CONTRASTING CONSTITUTIONAL AI AND THE CURRENT INVENTION

Consider the following example, contrasting Constitutional AI and the current invention.

With Constitutional AI, a small group of programmers writes a list of general ethical rules that they then use to train an AI. Then the trained AI (“Trainer AI”) trains other AIs based on what it has learned. The Trainer AI attempts to generalize the rules in the constitution to various ethical situations that arise. Human involvement is minimized because the whole point is to increase scalability while reducing the cost of RLHF.

Scalability Traded for Increased Risk When Humans Are Not Involved

At best, many more scenarios can be covered than by using RLHF alone, but the tradeoff is that it is difficult to know if the Trainer AI is interpreting the rules appropriately and teaching the right values. A significant risk remains that values are interpreted in ways that would seem strange, wrong, or even deadly to humans.

For example, the AI might learn that preserving the environment is good, and also that humans are having a negative effect on the environment, and then conclude that the best way to protect the environment is to reduce the human population by designing a virus that kills 50 percent of the population. Although logical, the outcome is not what most humans would consider to be ethical or acceptable.

Even if such obvious conflicts are explicitly trained out of the AI’s value system, it is very difficult to anticipate how things might play out when complicated chains of reasoning, actions, and effects are involved.

Today, humans are generally much better than AI at spotting obvious conflicts with human ethics, yet even humans have created many new problems while trying to solve other problems. For example, gasoline-powered cars solved a transportation problem but created a pollution problem that was initially unanticipated.

Speed of AI Heightens Risk

Generally, humans have had enough time to react and adjust their behavior if unanticipated consequences emerge. However, AI thinks and acts much faster than humans. There is a serious risk that a miscalculation by AI could result in extreme damage before it can be detected or corrected by humans.

Constitutional AI – an Inferior Approach to Safety

Even if no miscalculations or unintended consequences occur (which seems highly unlikely), and even if the values learned and taught by AI to AI are “good” as judged by the small group of programmers who wrote the constitution, there is no guarantee that they reflect the values of humanity more generally. In fact, it is almost certain that the values in such a constitution will NOT exactly reflect the consensus values of humanity simply due to the large diversity of opinions and human cultural norms.

Thus, in the case of using Constitutional AI, we achieve some degree of scalability but at the cost of taking humans out of the loop and reducing our ability to detect and correct unintended consequences. Further, the value system is highly unlikely to precisely match the values of humanity, even in the best case. The net result is cheaper, but less safe AI, compared to RLHF, which covers the same number of trained scenarios. Arguably, on a constant \$ basis, more scenarios can be covered, increasing safety, but still, Constitutional AI would be inferior to the current invention.

A second problem with Constitutional AI is that the method of using AI to teach AI currently tends to degrade the quality of the training with each successive generation. Just as in a game of “Telephone,” where the message gets subtly distorted as it gets passed from AI to AI, the subtleties that come with human involvement can be lost as successive generations of AIs process complex and ambiguous data. Something as simple as a human saying “I think XYZ is true” vs. an LLM converting this to “XYZ is true” can lead to problems when the third-generation recipient of the information has no idea that there was some doubt expressed about XYZ at the start.

Many of us have experienced similar problems with AI algorithms that attempt to determine our interests in news content or movies. At first, the algos may recommend useful movies we haven’t seen, but after a while, we find ourselves wondering why we are seeing such a narrow

range of choices. Well, the AI has picked up on the central tendencies in a few of our frequent choices, ignoring the subtleties (the “tail probabilities”) and after feeding us a steady diet of a certain type of news or movies which we consume, the AI becomes more convinced that we are interested in only a certain type of content.

The AI itself has biased us slightly at first. Then it compounded its imperfect understanding in a way that amplifies the “error” in its judgments. While annoying when it comes to recommended movies, such amplification of error could be fatal in other circumstances. With no “humans in the loop” to correct the misunderstandings, AIs can get off track quickly with Constitutional AI approaches.

Superiority of New Invented Approach to Safety

The current invention would have millions of individual humans each customizing their own AI agents using methods described in the cited PPAs and known in the art. Part of the customization would involve teaching the individual AIs the values of each human owner. Then the AIs, together with humans, would form a community of agents that provide Reinforcement Learning via Feedback (RLF), where the feedback comes from many AI agents as well as human agents.

This approach combines the scalable advantages of using AI to train AI, with the CI approach of using many customized AIs to increase the representativeness of values compared to a constitution created by an elite few. Further, because both humans and AIs can participate in the RLF process, humans remain in the loop and can be employed as much as resource constraints will allow.

Humans and AI agents can also dynamically identify and surface new ethical scenarios as they emerge during problem-solving. These new scenarios can then be incorporated into techniques used for eliciting representative values. AI can then be trained on these values as discussed later in this patent.

The net result is cheaper and safer AI, with greater scenario coverage, and humans maximally in the loop. An added benefit is broader acceptance and alignment of AI values because the values of many more humans were considered.

SIMPLE IMPLEMENTATIONS: REINFORCEMENT LEARNING VS. COMBINING WEIGHTS

Use AI agents as Part of Reinforcement Learning with Feedback (RLF)

An LLM being trained might change the weights in its network, and thus its behavior, based on feedback from (human or other AI) agents. In this case, conceptually, training AI on values using the current system could be exactly the same as current RLHF approaches, except that instead of using human agents, the current invention uses both human and AI agents. Because AI agents are quite cheap and fast compared to human agents, one might imagine a million AI agents, each trained by a different human, all providing RLF on the same scenario.

Bypass RLF and Combine LLM Weights Directly

However, since the knowledge that an LLM learns through RLF is ultimately reflected in changes in the weights of the network, it is also possible to bypass the RLF step altogether and simply change weights in the LLM directly. One might imagine two LLMs, identical except that one adjusted its weights based on conversations with person #1, and the other adjusted its weights based on conversations with person #2. Above, we discussed a situation whereby person #1 and person #2 each interact with the same LLM sequentially, and the LLM changes its weights sequentially based on interactions with the two people.

However, if person #1 and person #2 each interacted with a copy of the LLM, then after the interaction, we would have two separate and distinct LLMs with separate and distinct weights that were affected by the interactions with the humans. Now, without human intervention, it is possible to mathematically combine the weights of the two LLMs to result in a third LLM that has new weights reflecting the input of both humans.

If a simple average is used, then the new LLM would represent an average of the input of the two humans. But many other schemes for combining weights are possible. In the discussion below, we do not differentiate between situations where many (human or AI) agents train a student AI via sequential interactions and situations where multiple student AIs are trained separately and then the weights are combined according to some mathematical scheme.

Functional Equivalence of RLF and Weight Combination

Both situations are functionally equivalent, even though the methods may be different. Different methods may be more appropriate depending on the circumstances of training, availability of human and AI agents to provide training, computational resources, etc. Below, we are concerned primarily with different approaches to combining values (and more generally any type of knowledge or expertise, although in this discussion we are focused on values since that is a

critical issue for AI safety) that result in different weights in the network of the resulting LLM or AI agent.

SOME PREFERRED METHODS OF WEIGHT COMBINATION

There are many ways to combine weights from two or more AIs or LLMs. Which method is chosen depends primarily on what ends we hope to achieve via the combination and secondarily on considerations like technical efficiency. Below, we review some of the preferred methods for combining weights of multiple AI or LLM agents, with comments on why each method might be useful.

One Agent, One Vote (linear combination of weights)

A simple and effective approach is for an LLM to adjust its weights proportionally to the RLF received (or difference in weights from multiple models being combined), which is effectively a linear combination of the values of the various AI agents. A linear combination of weights has been shown to be optimal under some conditions and is a good starting point, especially because it embodies the “one agent, one vote” principle. This is essentially a scheme for training AI values where each participating agent has an equal influence on the final values of the trained LLM. “One agent, one vote” means there is only one copy of each customized AI, and each AI “teacher” has equal weight in terms of influencing the values of the “student” AI.

Human Input Counts More (Less) Than AI Agent Input

A variation of the above approach is a system where a student AI is being trained by both human and AI agents, and the humans have more (or less) weight than the AI agents in terms of how much influence they have on the final value system learned by the student AI. One might expect more weight for human input to be appropriate since humans are generally better equipped to represent human values faithfully, compared to AI agents trained by those same humans. However, one can also imagine scenarios (e.g., a human-trained AI, where the human subsequently becomes mentally impaired) where the AI trained by a human might actually be more capable of representing that human’s value system than the human him/her/their self. In this case, for example, it might be desirable to give the AI agent more influence than the mentally impaired human agent when training other AIs.

There is a bit of a slippery slope here in that one might also imagine scenarios where AI agents become so intelligent that they consider all human agents mentally impaired by comparison and therefore feel justified in ignoring human input and just letting the AI make all the ethical decisions. I do not agree with this extreme position and suggest that, because there is no rational way to derive values, AI must recognize the validity of human values as fundamental.

Given a set of values and overall human-aligned goals, it may be that AI can determine better, faster, and more effective ways of realizing the goals, but the fundamental values or goals cannot be rationally derived and therefore should not be left to AI to determine, no matter how brilliant AI becomes.

Expert (or Trusted) Input Could Count More

Another variation might weight the values of agents (human or AI) more based on the expertise of the agent. For example, physicians who have spent their entire careers advising patients about the difficult decisions that occur when a patient is terminally ill or on Hospice might have more insight and expertise in the ethical issues associated with end-of-life decisions than a layperson. Without advocating that such experts should or should not have more weight on ethical decisions that relate to their expertise, the current invention includes schemes whereby the value systems and expertise of agents can be weighted based on the expertise of the agent and/or the relevance of the expertise to the specific ethical decisions being made.

It should be noted that humans generally are reluctant to delegate important life and death decisions to professional experts (including medical or spiritual authorities) since there is a sense that profound decisions that affect an individual's life or future should be left to that individual as much as possible. Even if people make choices that seem unwise to those with expertise, generally prevailing values typically preserve the "right" of people to do stupid things, as long as they are not harming others via their stupidity.

Related to the idea of weighting expert input more than novice input is the idea that some sources of input are more trusted or have a better reputation than others.

For example, when a patient consults a range of friends and experts for advice on whether to have an operation, the patient would be wise to take the reputation, as well as the expertise, of the advisors into account. A knowledgeable physician with a reputation for performing lots of unnecessary surgeries might be less trusted than a family friend with good common sense who has the patient's best interest at heart. Of course, an expert surgeon who is also highly trusted might be the person the patient listens to most. Just as humans take both expertise and reputation into account when determining whom to listen to, so AI might weigh input differently depending not only on the expertise but also the trustworthiness of the source.

Weighting Based on Metadata

AI and human agents should have metadata associated with them that describe attributes of the agent including, without limitation, expertise, trustworthiness, reliability, individual and cultural preferences, group affinities, demographics, history of problem solving, subjective and objective ratings, and performance track record(s) on dimension(s) of interest.

This metadata about the agent's knowledge, skills, abilities, and other characteristics can be used to adjust the weighting of input to other AIs when the agent trains or provides input to the other AI. Note that although we have been mainly discussing the weighting of ethical or values input, the same logic applies to the combination and weighting of skills, knowledge, performance-related characteristics, or other aspects of AI agents.

Weighting Based on Recency or Other Time-Based Factors

Another set of variations in the approach to combining input from multiple agents when training or adding to the knowledge of a "student" agent might take time into account. It often makes sense to give more recent input more weight than older input. The rationale for this approach, generally, is that more recent input tends to better reflect the current state of the world and, therefore, is generally more relevant than older information. There are many different specific schemes for taking time into account when adjusting weights.

For example, the ethical opinions of humans that lived many years ago might be less relevant to current ethical norms than the opinions of humans living today. The warning found today in many Disney movies is, *"This program includes negative depictions and/or mistreatment of people or cultures. These stereotypes were wrong then and are wrong now. Rather than remove this content, we want to acknowledge its harmful impact, learn from it, and spark conversation to create a more inclusive future together."* This speaks to the fact that values and ethical norms can change over time.

It might be appropriate to give more recent ethical norms much more weight than older norms when training AI systems. Weighting based on time or recency could apply regardless of whether the input comes from humans, AI, or other data sources.

Approaches to differentially weighting historical information include, without limitation, exponential decay, linear decay, threshold-weighting (where there is a step-function change in weight based on specific points in time), decay based on other time-related circumstances, and other mechanisms for weighting a time series that are well known in the art.

With exponential decay, older input is given exponentially less weight than more recent input.

With linear decay, older input is given less weight in linear proportion to how far in the past the input was received.

Examples of threshold-weighting might include weighting ethical input differently depending on whether it was received before or after the time that certain laws were passed.

For example, during prohibition, selling alcohol might have been a crime and determined to be unethical, but right after prohibition, the ethical status of selling alcohol changed. This was a stepwise (or “threshold”) change in the ethical status of alcohol sales that could not be easily captured via exponential or linear decay of weights. Any events in time, not necessarily law changes, that significantly altered opinions on what is ethical behavior might similarly require stepwise adjustment to the weights of ethical input.

VALUES OR ETHICS-SPECIFIC IMPLEMENTATION CONSIDERATIONS

Unlike engineering problems that have precise, correct solutions, ethical questions depend on the humans being asked. Despite the attempts of philosophers, religious, and political leaders to construct universal value systems that govern all human behavior, morality is notoriously difficult to codify in ways that are acceptable to all or even most of humanity. The legal systems of various countries and municipalities represent an attempt to put some guidelines on human behavior, but no legal system attempts to cover every possible scenario of human behavior.

No Correct Answer

Instead, the vast majority of human behavior is regulated by the moral sense and opinions of individual humans. And there is a huge range of decisions that people make every day, where there is no right moral answer – just opinions as to what is right or wrong. Given this complex situation, and the fact that machines are notorious for needing exact specifications in order to behave, the implementation of ethics-specific AI safety solutions may differ from the general implementation methods described above, which apply equally well to training AI on ethics or other types of knowledge.

Trolley Problem Example

Consider a classic example of an ethical dilemma well-known in the field of AI ethics, the Trolley Problem. In one version of this ethical dilemma, a self-driving car controlled by an AI finds itself in the situation of having to choose between killing pedestrians who suddenly jump in front of the car or swerving into a barrier to avoid the pedestrians and killing the occupants of the car.

Like many difficult ethical decisions, there is no right answer. Yet, humans still have opinions about what is ethical and what they would do in such a situation.

Surveys of many humans have shown that what humans consider to be ethical depends. It depends on who is in the car, who the pedestrians are, whether the pedestrians are crossing illegally or legally, and even how old the people involved are. Humans are more likely to instruct the car to run over pedestrians if the pedestrians are crossing illegally, are homeless, or are simply old. Humans, according to the survey research, are less likely to kill, or allow to be killed, those who are young, pregnant, women, or in certain professions, such as the medical profession. None of these aspects of human decision-making is captured in our legal system. There is no law that says it is okay to run over someone they are older than if they are younger, yet humans take these factors into consideration, and mathematical weights can be assigned to each of these factors. Similarly, AI can learn to make decisions, taking these same weights into account.

Ethical Solutions That Mirror What Humans Do

In order to have AIs that behave in ways that make sense to most humans, AI will have to be trained not according to a rigid constitution but according to how real humans actually behave. This behavior may change depending on the culture. In the US (something of a youth culture), running over elderly pedestrians is likely more acceptable than in certain Asian cultures where elders are revered and held in high esteem.

If AI is to make ethical decisions in the same way that most humans do, it will also have to take these cultural factors into account. How will it do this most effectively?

Constitutional AI comes up short, as we would need a different constitution for each different group of humans.

Much more effective would be to have AI interact with many humans, learning their values, which is the approach of RLHF. But, as pointed out, RLHF scales poorly.

The best way to capture the wide diversity of human values, while retaining the scalability that comes with AI being involved in the instruction of AI, is to have a multitude of teachers, both human and AI agents. The AI agents should be trained by a wide diversity of humans so that each AI agent carries the unique values, ethics, and moral sensibility of its owner into every interaction that it has, including the activity of training other AIs.

Human Values Will Outlast Human Superiority in Intelligence

In the long run, AI will undoubtedly surpass human ability in cognition, problem-solving, and processing of information. As AI grows increasingly intelligent and capable, the role of humans will increasingly be to determine the values and fundamental goals that the more intelligent AIs seek to realize.

Inclusiveness and Representativeness Important for Human Values

These values should not be the province of a small elite group of programmers but rather should reflect as broad a cross-section of humanity as possible. By including all humans who are able and willing to customize their AIs in the crucial task of determining the base values of AI, we can achieve broad representation more efficiently and cost-effectively than any existing approach.

Humans can, and should, remain in the loop as much as possible when training AI via RLHF and similar methods. However, to the degree that humans are unavailable or the resource demands are too great to have all the training done by humans themselves, the next best thing is to include a wide and diverse group of AI agents in the training. Each of these agents would be customized by a different human and would represent that human's values and sense of ethics.

Once an AI agent is trained by a human, it can operate 24/7 with or without supervision from the original owner. This allows the human owner's values to be incorporated into training and other activities without requiring constant human involvement from the owner, as RLHF would.

Ethical Norms

Although it is undesirable to have the safety and ethics of AI driven by a constitution written by a small, elite group, it is still possible that most humans, regardless of cultural or individual differences, would agree on certain normative ethical principles. Examples of these, without limitation, might include (variants of):

- The Golden Rule (Do not do to others that which you would not like done to you)
- First, do no harm (e.g., as reflected in the Hippocratic Oath taken by physicians)
- Do not kill (unnecessarily or except in specific, exceptional circumstances)
- Preserve individual freedom (unless it limits the freedom of others)

Of course, almost as soon as one reads these principles, exceptions spring to mind. Do not kill – but what about self-defense or war? Preserve individual freedom, but what are the limits, and when does it impinge on others?

Details and nuance matter, even in the application of principles that most humans would embrace broadly. However, by starting with general normative ethics that have broad acceptance across a large, diverse, and representative group of humans, it is possible to refine these principles and determine when and how they apply in detailed circumstances much more efficiently than if no starting principles existed at all. Thus, general ethical norms are not antithetical to the approaches outlined in this invention. Rather, they are a point of departure that can help AI achieve realistic and nuanced ethics and behavior that are aligned with what most humans believe to be good and aspire to.

Group Norms

Once we have admitted the potential usefulness of ethical norms as a starting point for further refinement via the methods discussed in this invention, the door is open to group and planetary norms. There is a continuum from very specific individual AIs that have been trained to believe and act as much like a particular human as possible, to AIs that have been trained to broadly act in ways that specific groups of humans agree with, to AIs with ethical norms embraced by many specific groups of humans.

Ethical norms at each point on the continuum can serve as a starting point for training ethical AI behavior. The idea that one set of norms or one constitution should power all of AI is likely unrealistic and far too brittle to work in the real world. If it were possible, then the many differing viewpoints espoused by religious, political, and cultural groups would long ago have merged into a consensus.

The diversity in human ethical norms is not a bug; it is a feature. We should not expect AI to achieve consensus and maintain human alignment if humans themselves cannot achieve this consensus, especially if there is debate about whether such consensus is even desirable.

Ethical Contracts

Another aspect of ethics is the idea that humans often enter into ethical and/or implicit or explicit social contracts when they join a group or participate in society. For example, members of a particular religion largely agree with a set of rules and ethical precepts espoused by a religion and often enshrined in one or more “holy” books. The Koran for Muslims, the Old Testament for Jews, and the Bible for Christians all contain ethical precepts and rules that members of the respective religions are largely expected to follow. Similarly, Confucianism in China, the ideals reflected in the Declaration of Independence and Constitution in the USA, the writings of Marx

for some Communist countries, and liberal or conservative ideologies for various political groups all contain normative prescriptions for human behavior.

Members of specific groups may be considered to have implicitly entered a contract to accept the bulk of the principles of a particular group when they join. By being a citizen of, or simply living in, a particular country, humans are explicitly subject to the laws of that country, including laws that explicitly specify what is criminal (“wrong”) behavior.

Thus, for ethical problems, the solution sometimes depends on what social or ethical contract humans have made with the group or culture in which they find themselves. Such contracts can be useful in simplifying the task of training AI, inasmuch as a starting point can be the laws of a particular country or the implicit/explicit rules of a particular group.

The existence of such contracts, as well as sets of rules and laws, has implications for how to efficiently train safe and ethical AI as described below.

The Safety Argument for Democratic, Representative Values

It has been said that democracy is a bad political system, but that all the others are worse. Since monarchies, dictatorships, republics, and many other forms of government all co-exist today, with democracy being only one form, it is worth mentioning the benefits and disadvantages of a democratic and representative system of training AI when it comes to AI safety.

Most humans would agree that the most important concern with respect to AI is the existential threat that a majority of humans acknowledge it currently poses. That is, AI could wipe humans out. If that happens, it doesn’t matter what form of government or religion you prefer. We all would be dead, and the point is moot. So maybe we should be asking not which religion or form of government is best, but simply which principles are most likely to lead to humanity’s survival.

If it were possible to anticipate all the dangers and different scenarios that will occur and how AI would act in each, it would be possible to prescribe a set of rules that guarantee human survival. However, the complexity is too great, and the speed at which AI can operate is too fast for this approach to work. Theoretically, it is possible to calculate every possible move in a game of chess and therefore win the game with the first move. But the computational complexity is so great that even the most powerful chess-playing computers (far better than the best human) can’t win this way. Instead, moves must be made one at a time, and the board re-evaluated after each move.

The best chess-playing programs still use flexible heuristics, or rules, which can handle many situations. They rely on strategy and general principles to win, until finally, at the very end, with few moves left, every remaining move can be calculated.

Democratically representing the opinions of most humans is an approach that is rarely optimal but generally achieves an acceptable outcome. Collective intelligence – the idea that two heads are better than one – is responsible for the vast majority of human progress, culture, and technology. But when it comes to the subjective area of ethics and values, where there is no objectively correct answer – just human opinions – democracy, or a collective intelligence approach, if you prefer, really shines.

One benefit of a democratic and representative set of human values is that it tends to mitigate extreme positions, which are likely to be most risky to human survival. There is a beneficial diversification effect regarding values.

Just as diversification in an asset portfolio reduces volatility and risk, so does a diversity of human opinions and judgments, which tend to have a stabilizing effect on the overall portfolio of values. In an asset portfolio, a diversified portfolio always returns less than if you were to concentrate all the investment on the top winners. The problem is that no one knows what the winners will be with any degree of certainty. That is why the diversified approach of just “buying the index” tends to outperform more than 80 percent of all portfolio managers who try to “beat the market.”

Similarly, there are “philosopher kings” or religious saints who can make laws and ethical rules which, for a time, are far superior to the collective judgment and behavior of the masses. But what happens when the superior king or saint is gone? Then, a power-hungry dictator might arise whose reign is far worse. The more stable approach – less likely to be really great, but also less likely to be really terrible – is to follow the values of a large representative population of humans. All these humans want to survive. Most want good things for themselves and their fellow humans. Few want to destroy the environment or the planet. While the collective values are imperfect, they are usually not malevolent. Importantly, they are based on the hearts of humans.

The Scientific Argument for Democratic, Representative Values

Putting aside the practical benefits of a diversified, representative, “portfolio approach” to human values, a representative sample is also a scientifically valid way to accurately answer a question to which there is no logical answer, namely: what is right and what is wrong according to humans.

A fast computer could answer a math problem faster than a million humans, but when it comes to the subjective determination of what is wrong and what is right, calculation speed is useless. If we want to know what human values are, there is no substitute for asking them and watching their behavior. The more humans we ask and watch, the more representative the values may be.

Statistical sampling theory dictates that a larger sample will give a better estimate of the true subjective values of humans. Sampling can occur within a group to determine the values of members of the group. But if the group is all humans on Earth, then the appropriate procedure is to gather as large a sample as practical from all humans on Earth.

Since, as a first approximation, notwithstanding some of the variations described earlier, a simple one-person, one-vote combination of human values gives an accurate read on human values, and since there is no objective right or wrong, the straightforward conclusion is that while other approaches may be useful for other goals, if the goal is to get the most accurate read on what human subjective values are, then a simple – democratic and representative sampling process – is the way to proceed.

Arguments can be made about whether a representative sample of human ethics results in the safest AI system. After all, humans are known to engage in all kinds of atrocities and genocidal behavior, as well as follow more loving paths through life. But if we want human-aligned AI, then it seems clear, for better or worse, that gaining a representative, democratic sample is the best course. Anything different amounts to an attempt to engineer a set of values that is different than what humans themselves espouse.

If humans truly valued death and destruction above all else, then given that we have weapons of mass destruction already, we should not be here. I suggest that most humans have very positive and loving values, although we do not always foresee the consequences of our actions. AI will undoubtedly help us to act more intelligently, but only humans should supply the values. It is our privilege, responsibility, and (arguably) purpose to provide those values which cannot be rationally derived in any other way than by determining what humans think, say, and do.

Thus, the conclusion that we need democratic and representative values is less driven by political theory, religion, or philosophy than by the simple science of statistics. If we want to know what humans believe is ethical, we need an accurate sample.

A democratic and representative approach is the statistically valid way to obtain such a sample. Whether we like the sample that is obtained is a separate question. However, since values cannot be rationally derived, the only way to determine what human values are, which is necessary if we wish to align AI with existing human values, is to gather a valid sample.

Since words like “democratic” make some countries with non-democratic political systems nervous, perhaps a more accurate and less politically charged description of what humanity should be using to train AI is a “representative and statistically valid” sample of human values.

EFFICIENT TRAINING METHODS

Now that we have discussed some issues specific to training AI on ethics or values, we turn to some preferred implementations of the invention for accomplishing this training. Specifically, to be both novel and useful compared to existing approaches such as RLHF and Constitutional AI, the preferred implementations should satisfy the following constraints:

1. To maximize alignment with human values, the values used to train the AI must be a representative and a statistically valid sample of the human population with which the AI is expected to align (co)operate.
2. To achieve widespread use and practical adoption, the method for training AI must be highly scalable and have appropriate “path coverage” of the ethical / safety situations that the AI is likely to encounter.
3. To maximize the probability that AI picks up on the nuances of human values and does not propagate errors that a human would have easily detected, the method should follow the principle of including humans in the loop and maximizing human involvement to the degree that this also allows scalability.

Generally, the current invention meets these constraints via a combination of many humans training/customizing individual AIs and then having these many customized AIs, together with as much human involvement as possible, train other AIs in a way that is more scalable than RLHF and more representative and accurate than Constitutional AI.

A key aspect of the invention is to recognize the scalability benefits of AI training AI while minimizing the drawbacks (namely unrepresentative values and error amplification) by using a large group of AIs, each of which has been customized differently, to augment training by a variety of humans. The choice should not be between training by humans or training by customized AI, but rather, we should leverage the beneficial aspects of each approach. We should train using as large a collective as practical of BOTH many humans and many individually trained or customized AIs, each bringing unique information to the table. The result will be more accurate and error-free training combined with a more representative and statistically valid set of human values.

The current invention represents a novel and superior approach to addressing the Alignment Problem compared to existing solutions.

Path Coverage

As mentioned earlier, there are potentially infinite dangerous situations that we need to train AI to deal with safely. Further, for alignment, we require that AI deal with dangerous situations and ethical choices in a way that is representative of how human populations would deal with the same situation. AIs trained in this way become predictable, which is a key requirement for humans and AI to trust and interact with each other. While “hallucinations,” unpredictability, and lack of trustworthiness are the hallmarks of current systems, this invention seeks to produce predictable, trusted AI that behaves as humans would and as humans expect an intelligent entity to behave.

At some level, the problem of training AI reduces to a problem of “path coverage.” That is, AI must be trained in enough representative dangerous or ethical-decision-making situations that its behavior becomes predictable and trusted in these situations. Enough of the situations (paths) must be covered in the training.

Generally, when testing software, human software developers come up with test cases that try to cover the use cases that are likely to arise. Since it is impossible to test every possible use for complex software, human developers attempt to determine which use cases are most common and also which use cases have the highest impact if things go wrong (e.g., the dangerous cases). Common use cases get more testing than less common cases. More dangerous scenarios get more testing than benign scenarios. And common, potentially dangerous cases get the most testing of all. Since testing resources are limited, it makes sense to concentrate the limited resources on the most common and/or impactful scenarios.

If software were a car, we could live with an interior light failing more easily than we could live with the brakes failing. If we had to choose between testing interior lights or the brakes, because of limited resources, we would prioritize testing the brakes. The same logic of looking at the impact of a failure applies when training AI, regardless of whether the training is done via RLHF, via other AIs, or (as this invention suggests) via a combination of humans and many customized AI agents.

Similarly, imagine there are two interior lights in a car, but one light is used 10 times more often than the other. If we had limited testing resources, we would prioritize testing the more frequently used light. That’s because the impact of a failure for either light is about the same, but the difference in how commonly used the lights are means a failure in one of the lights would

cause 10 times the annoyance factor as a failure in the other light. This logic, of testing the more common situations (if the impact is similar), also applies to training AI.

In the preferred implementation of our invention, the fact that many representative humans individually customize AIs, which are then used (together with other AIs and human agents) to train new AIs, means that the most common use cases will receive training roughly proportionally to how common they are. This desirable result arises from the invention's use of a group of AI and human teachers. The commonest teaching cases will be covered by most members of the group, while the edge cases that are rarer will be less represented. Statistically, we have a sample of humans providing both human values and use cases. The larger the sample, the more certain we can be that all of the most common cases have been addressed in ways that are aligned with the human population's values.

The problem of addressing some dangerous cases and difficult ethical decisions is more challenging. That's because dangerous situations and tough ethical decisions are often relatively rare. In this case, the preferred implementation approach is for the training system (see below) to ask humans and AI agents to think of as many dangerous scenarios as possible. The total pool of dangerous situations and difficult ethical situations can then be allocated among the sample of humans and AIs such that all the dangerous/difficult (or otherwise "high impact") scenarios receive enough input from enough different agents so as to have a representative sample teaching the AIs on these challenging topics.

The scenarios can be randomly allocated across (human and AI) agents, or the scenarios can be allocated in a more optimized way. Humans who think of particular scenarios are more likely to have experience with those scenarios and may provide better teaching input than other humans who may not be familiar with, or even understand, the factors involved in the scenario.

Another preferred implementation scheme is to suggest scenarios for human and/or AI input based on a best match or a random allocation approach, and then allow humans or AIs to choose which scenarios to provide input on. If certain high-impact scenarios remain where no humans have chosen to provide input, these may be assigned to humans in a second (or later) iteration of the allocation/assignment process.

Let's consider some examples. Suppose an AI must choose which patients get medical attention in a triage situation. Some of the humans involved in training the AI on the ethics of making these choices are emergency room doctors and paramedics used to making triage decisions, whereas other potential trainers are well-meaning people with no medical background.

The medical professionals may be more confident than untrained bystanders in making the decision to let a mortally wounded patient die because nothing can be done, if that decision might save another patient. The untrained bystander, moved by emotion and not fully understanding the need for triage, might be inclined to advise giving equal attention to everyone or giving all the attention to the mortally wounded patient. In that case, the well-intentioned bystander's advice likely would result in a worse outcome because the mortally wounded patient will die anyway, and now a less severely wounded patient also has a higher chance of dying. For this difficult ethical situation, specialized knowledge is an advantage, and we might prefer to let the medically experienced professionals teach the AI.

If a random sample of humans were asked to generate a list of difficult ethical decisions, it is likely that the emergency room doctors and paramedics would generate at least some medical scenarios since that is what they are frequently exposed to.

On the other hand, an HR professional involved in Diversity, Equity, and Inclusion initiatives might come up with ethical scenarios related to promoting or hiring one qualified individual over another qualified individual, again based on the human's daily experience.

By tapping the knowledge of a wide range of humans, it should be possible to not only generate a wide range of ethical scenarios but also prioritize which humans provide input on the scenarios in a way that results in better learning outcomes for AI than simple random allocation. By comparing the outcomes of various attempts (including random assignment) to optimize the allocation of human attention to training AI in these scenarios, it is possible for an AI optimizer to learn preferred allocation strategies.

In some edge cases, where few humans have experience in a scenario, or where the scenarios are very impactful but unprecedented, human experts may need to intervene and develop specific scenarios which are then allocated to various humans and their AI agents to get input that is representative of what humans would do in these unusual but important situations.

To summarize the way path coverage is addressed in this invention, A key feature of the current invention is to use a sufficiently large number of humans and their customized agents. Since different humans have different circumstances and will have trained their AAAs based on the human's circumstances and knowledge, in the aggregate, many AAAs should provide very good complete "path coverage" over the range of ethical circumstances that humans find themselves in.

Further, this approach has the nice feature that those ethical circumstances that are most frequent and most important to humans are likely to be most represented by the AAAs since the human owners will more frequently and more emphatically train their AAAs on these cases.

Those ethical situations that are less frequent and/or less important will naturally receive less training from humans.

Thus, when the values and ethics of the AAAs are combined (using any of the approaches outlined above), there will naturally be the most input on the most frequent and important ethical issues. With more input, there will be the most consistent, representative, and reliable ethical behavior trained in the student AI(s).

Just as in normal human life, we have lots of experience with common ethical choices and give extra attention to important ethical choices, so too the trained AI will naturally receive the most training from the largest number of sources on these frequent and important ethical conditions.

RLHF from professional humans, tasked with training AI, can be used to fill in the ethical gaps where important, but infrequent, ethical dilemmas arise, so that the student AI receives optimal ethical training.

Real-time Detection and Prevention of Unanticipated Safety Issues

Another method is for (human and/or AI) agents to dynamically flag potential ethical issues in real-time as they are encountered and then present these issues to other groups of agents for resolution. Rather than relying on experts or crowdsourcing efforts to determine the complete space of ethical scenarios ahead of time, the real-time flagging approach allows AI, AGI, and SuperIntelligent systems to detect potential issues and potentially pause work until additional (human) input can help the system determine the ethical approach.

Of course, in time-critical situations, pausing or delaying might not always be possible, but the approach can be used for many issues that do not demand an immediate response. Including this dynamic approach of delaying responses until ethical input is received can reduce an otherwise exponential space of possibilities to a manageable size.

One implication of the use of the real-time detection and delay response strategy is that critical high-stakes issues that require an immediate response (whether to launch a counterattack to a perceived missile launch comes to mind, or other military applications) will require proportionally more path coverage and training ahead of time compared to situations where a delay in response is acceptable.

Also, while constitutional AI approaches are generally suboptimal for determining ethical knowledge bases, partly because they are based on rules developed by an elite group, such approaches might be acceptable as a means of temporarily flagging potentially unethical

situations until a representative sample of human ethical judgments can be obtained about an unanticipated, but potentially dangerous, situation.

For example, a rule that said “An AI can never provide information that might be used to harm other humans” might flag potentially dangerous scenarios, delaying responses to such situations where possible until they could be reviewed by humans or otherwise subjected to deeper review.

Some false positives will occur, and this rule is likely too general. That is, someone might ask about using arsenic to poison rats and have to wait for a response while the AI flags the question and gets other (human) agents to weigh in on whether answering the question (given the context of the conversation) is a risk to humans. As long as the delay is not too long, it might be acceptable if the delay prevents serious safety issues.

Mathematical algorithms known in the art and other methods for calculating when a test is doing more harm than good (in the medical profession, for example) can be employed to help quantify these decisions.

If we can use AI to determine in real-time whether an applicant is a good credit risk, there is no reason that similar algorithms cannot be employed to delay or avoid responding to certain potentially dangerous questions. That said, ideally, there would be a method for rapid review and appeal of the potentially dangerous cases.

One approach, which minimizes delay to a fraction of a second while still providing some margin of safety, is to have questionable cases reviewed by multiple different AI agents to see if there is a consensus among the agents as to the safety of the request. Human agents, working more slowly, could override the AI agents (and teach the AIs in the process) upon appeal or when they are able to get to the prioritized list of issues.

Automated means for tracking the frequency and potential impact of unanticipated safety issues could help optimize the use of human decision-making and ethical judgements for the most common and/or important issues.

Note that these same techniques can be used to address non-safety issues as well, as in the case of trying to get the most accurate answer to a question (even if the question does represent a safety risk).

One preferred method for reducing the amount of “hallucination” by LLMs is to have multiple AI agents all process the same question and then take the consensus or majority answer as the

most correct one. This approach might employ versions of the same LLM with different parameter settings to generate multiple responses. Alternatively, completely different LLM models can be used. Users can set the degree of reliability that they desire (and are willing to pay for), which in turn determines the amount of redundant processing and/or the number of different models used to generate the ultimate answer to the user's query (or solution to the user's problem).

Returning to unanticipated potential safety risks, the difference between employing imperfect real-time detection of safety issues based on existing well-known approaches and doing nothing is huge. The nuances of balancing the opportunity costs of not responding to perfectly harmless questions versus preventing disasters are something that can be refined (ideally using a data-driven approach) over time. However, some real-time detection and prevention of issues before they occur is almost certainly a net positive at some threshold for false positives.

Conversational Method for Training AI on Scenarios

Detailed discussion of algorithms and machine learning techniques that can be used to customize AI agents, LLMs, or AAIs has already been discussed in the patent applications cited and included by reference at the start of this disclosure.

However, from a user interface perspective, one of the simplest methods is for humans to have conversations with the AI they are customizing and then provide instructions to that customized AI as to how it should behave when training other AIs. Such conversations can be initiated by either the AI being customized, the human doing the customization, or both. Humans with strong beliefs and/or knowledge about certain issues may want to focus the conversations and subsequent customization of their AIs in these areas.

Survey Methods as the Basis for (Ethical) Knowledge Acquisition

In addition to having conversations with humans, AI can conduct surveys of humans to elicit their opinions and knowledge about a wide variety of subjects, including ethical views. Survey approaches have the advantage of being well-suited to gathering random and representative samples of human knowledge using a variety of online and offline methodologies that are well known in the art. Like intelligent conversational approaches where the direction and content of the conversation can be directed by the AI with the goal of filling in gaps in the AI's knowledge, survey methods can also target specific knowledge gaps, including gaps in coverage of certain ethical situations.

Passive Machine Learning Approaches to Knowledge Acquisition

Both conversational and survey methods require humans to actively engage with AI in order to teach AI (ethical and other) knowledge. However, humans have limited time to engage in such activities, and AI has an almost insatiable appetite for new knowledge. Therefore, AI will have to rely extensively on passive methods of knowledge acquisition, such as those currently employed in the creation of today's LLMs and other AI agents.

Specifically, algorithms such as variants of the transformer algorithm, and other algorithms well known in the art and generally associated with “deep learning” and “neural network” or “connectionist” approaches to machine learning can be used. More generally, any method that uses the passive “digital footprints” left by human (or AI) users (or agents) as they perform tasks, including but not limited to, online navigation, selection of products and websites, solving of problems, communicating with other humans (or AI agents), purchasing, filtering, analyzing, researching or other online tasks, to train AI and acquire knowledge are examples of passive machine learning approaches.

In a preferred implementation, AI can efficiently increase its (ethical) knowledge by following a serial process which may include, without limitation:

1. Identify the desired knowledge base and path coverage.
2. Identify gaps in existing knowledge (e.g., by self-analysis or via feedback from external human (or AI) agents).
3. Seek datasets that contain information needed to fill in the knowledge gaps.
4. Analyze the data and the data sources to gain confidence that the data is a valid and representative sample of human (or AI agent) knowledge, if that is the aim (e.g., as it likely would be in the case of trying to determine values that are representative of a human population). Note depending on the goal, other analysis, besides determining a representative sample, may be used (e.g. if the goal is to maximize expertise in an area, then the analysis might be to determine that the data represents the most expert knowledge available in a field as opposed to a representative sample of the knowledge of all humans).
5. Filter/clean the data based on dynamic or pre-determined criteria. (An example of dynamic criteria, without limitation, would be a quality threshold that automatically is raised as more and more data is located, such that at the beginning when no data on a topic exists, the AI is more willing to accept any data that helps fill in the knowledge gap

but as more and more data is located, the AI can dynamically raise the threshold and become more picky about the data included in the training set).

6. Use the selected and filtered/cleaned data to train the AI using methods that have been mentioned above, or in cited PPAs, and/or which are well known in the art.
7. Repeat the steps and/or learning epochs until a pre-determined or dynamically adjusting level of quality has been reached. (An example of a dynamic quality threshold would be an AI that is monitoring the current state of competitive AIs by pinging them with questions in the area of interest, and, based on the responses of those competitive AIs, determining whether more training and/or data is needed in order to acquire knowledge that is on par with or superior to the competitive systems and then setting thresholds based on this dynamically changing competitive landscape).

Frequency of Knowledge Updates

Regardless of the method used for filling in (ethical) knowledge gaps, all knowledge is a moving target, with more recent knowledge generally being superior and supplanting earlier knowledge.

For example, at one time, the generally accepted view was that the world was flat. Today, almost everyone agrees that Earth looks much more like a sphere. An AI trained on the “flat-Earth” view would be out of date and would need to update its knowledge.

While scientific views, such as the shape of the Earth, may change very slowly, other types of knowledge, especially subjective ethical norms, may change much more frequently. As mentioned earlier, Walt Disney movies that were made and rated G when some of us were children contain stereotypes, which, by today’s ethical norms, are wrong. These sorts of ethical norms change more rapidly than other types of knowledge that may be valid for centuries.

One of the components of the preferred implementation, therefore, is a mechanism for updating AI knowledge at the appropriate frequency based on the velocity of change of the information and/or other factors.

In the preferred implementation, AI would categorize different types of knowledge along multiple dimensions, one of which would be the rate at which the human (or AI agent) opinions about the topic have changed. This rate of change (ROC) of knowledge is particularly important when the knowledge is subjective human knowledge, such as consensus ethical views.

Knowledge (including, without limitation, ethical knowledge) that changes at a relatively slow rate does not need to be updated as frequently as knowledge that changes rapidly.

One way to think of this problem is to consider the difference between ethical principles and fashions or interpretations of the principles. Humans have a relatively long-lasting principle that human life is valuable and should not be taken away lightly. This general principle has survived many thousands of years and is incorporated in the laws and religious/moral traditions of almost all human groups, even though admittedly, exceptions are included for war and certain other circumstances.

On the other hand, certain ethical norms are more akin to fashions, which change depending upon the consensus of the group of humans being asked or the time when they are asked. Affirmative Action in college admissions was an ethical norm for the last few decades until a recent decision by the Supreme Court of the US began influencing this norm. Almost immediately after the Supreme Court decision, many companies and other organizations began adjusting their ethical norms, decision-making, and communication practices within their groups to ensure that their versions of affirmative action and/or diversity, equity, and inclusion practices (if they existed) were within the new mainstream view.

Similarly, attitudes towards people with different sexual preferences or gender identities, different skin colors or religions, and different professions tend to be more fluid and change more frequently than more long-lasting and widely accepted ethical precepts such as “thou shalt not kill.”

In the preferred implementation, to update its knowledge and fill in knowledge gaps efficiently, AI needs to be able to determine the rate of change of the type of knowledge being considered. Ethical norms that change more frequently need to be updated more frequently than knowledge that represents a more long-lasting ethical consensus.

Some general principles for updating knowledge should be reflected in the preferred implementation, including, without limitation, the ideas that:

1. The more rapidly the knowledge base changes, the more frequently updates should be made.
2. The more fundamental and established an item (e.g., ethical principle) of knowledge is, the more weight existing past knowledge should be given.
3. All other things being equal, more recent information should be given more weight than older knowledge, but this should be adjusted based on the rate of change of the knowledge area; specifically:

- a. In areas where change is exponential or rapid, exponentially (or proportionally) more weight should be given to recent knowledge compared to past knowledge.
- b. In areas where knowledge changes at a linear rate, with respect to time, more recent knowledge should receive linearly more weight than previous knowledge.
- c. In areas (such as firmly established principles, such as the high value of human life) where knowledge has been constant for long periods, new knowledge should be given similar, or slightly more, weight than older knowledge.
- d. Generally, if no other factors apply, more recent knowledge is more representative than older knowledge, and should receive more weight since a general role of any intelligent system (including AI) is to have an accurate representation of the CURRENT state of the world, and more recent knowledge is generally more accurate than older knowledge is describing the state of the world today.

The frequency of updates will become increasingly important as the rate of change in knowledge increases. With AI accelerating scientific discovery and technological change, it is conceivable that knowledge about the world will change faster than humans can update their collective knowledge. Given the limitations of human information processing and the tendency of humans to cling to older knowledge and paradigms long after they are outdated, knowledge may already be increasing far faster than most humans are able to comprehend or adapt. That said, when it comes to ethical precepts and fundamental ethical principles like the value of human life, we are fortunate that these change relatively slowly. The interpretation and application of the ethical principles may depend partly on technological/knowledge change, but the ethical principles themselves are relatively constant.

The Spinning Knowledge Wheel Framework

One might imagine a spinning wheel with fundamental human values, such as love and the value of human life, near the center of the spinning wheel. At the very center of the wheel, the ethical principles are constant and motionless, just as the center of a spinning wheel does not move at all. However, the farther along the “spokes” one travels in the direction of the “rim,” the faster the rate of change. Similarly, all knowledge (including ethical) can be characterized as lying closer to the center or farther out on the rim of a spinning wheel. In the preferred implementation, an efficient AI needs to update knowledge areas on the rim very frequently without changing the core human values around which, and for which, the knowledge has relevance.

Human-centered aligned AI must put relatively constant and fundamental human values at the center, while updating (faster than humans will be able to conceive) other types of knowledge that are closer to the “rim” of the spinning and ever-increasing knowledge wheel. This approach ensures that, despite their inferior ability to process information and understand exponential changes in knowledge and technology, human values remain the center of AI, which will become potentially trillions of times more powerful and knowledgeable than any one human.

As humans, we cannot keep pace with the change at the “rim” of the spinning knowledge wheel, but we can understand and orient the entire spinning wheel by serving as the relatively slower-moving center, where the values and purpose of AI reside. This structure is essential if humans are to not only survive but also prosper in the age of AI systems that are vastly more intelligent and powerful than we are.

Weighting Implications Revisited

Just as the recency of information has implications for weighting (discussed above), the type of knowledge, as well as its characterization as fundamental and long-lasting versus (for example) more changeable and a matter of current opinion, also has implications for the weighting of such knowledge. Long-lasting, more fundamental knowledge should have stronger weights that are more resistant to change than short-term “fashionable” opinions.

The preferred implementation would have differential weighting on ethical knowledge, for example, representing how fundamental that knowledge was. One method of determining how fundamental an ethical precept is would be to actively survey or ask humans to rate this quality compared to other candidate ethical precepts.

Another method is to passively analyze the data record of human behavior and draw conclusions based on that analysis. Both methods are likely to be useful. Passive analysis of behavior in the recorded data is more efficient than actively engaging humans, but actively engaging humans is necessary to ensure that the conclusions being drawn from the analysis of passive data are correct (from a human point of view) interpretations.

Finally, note that while the current invention has a bias towards statistically valid and representative samples of all human opinions when training AI on human values and ethics, not all knowledge is a matter of opinion. In fact, values and ethics are more of an exception than the rule in this regard.

Aside from artistic judgments, political and religious views, and other subjective areas, most human knowledge is factual. AI will likely want to weigh knowledge that is factually accurate and

justified by converging evidence from many sources more highly than unsubstantiated opinions on factual matters.

While some people still believe that the Earth is flat, this view should not be given equal weight to the view that the Earth is spherical. To do so would ignore the vast amounts of converging scientific evidence and facts that support the spherical view. Opinions, on non-ethical or non-subjective matters, should not count as much as facts.

This is a tricky problem, as humans tend to select facts that support their views, and the facts themselves change. At one time, not so long ago, the consensus medical opinion was that cigarette smoking was healthy for the lungs. Today, there is a significant portion of the population that views vaccines as harmful, even though this stance is contrary to the majority of medical evidence.

To effectively and efficiently navigate these issues, AI must rely primarily on the scientific method of seeking valid and reproducible evidence before accepting facts. Other scientific principles, such as (without limitation) converging evidence, Occam's Razor (or reducing the number of degrees of freedom in scientific explanations), and other tested tools of the scientific method, should be used, in the preferred implementation, by AI seeking reliable and accurate knowledge.

Do Not Delegate Values to AI

However, AI must not confuse (as the philosopher, David Hume, said) “ought” with “is”, or facts with values. Values are necessarily subjective. Arguments claiming that values are objective, such as the claim that “everyone would agree that something that makes all humans suffer the most extreme misery imaginable is bad,” are naïve and fail to grasp that other, non-human entities might not accept such values as self-evident at all.

AI operates, at the most fundamental level, in a precise logical manner. It's all zeros and ones at the machine level of implementation. To expect such a system to somehow intuit that human values are fundamental, or worse, to expect that it will logically derive values that are human-centered, is the worst kind of sloppy thinking. The kind that can lead to human extinction.

Both David Hume and Herbert A. Simon (the Nobel Laureate and co-inventor of AI) had it right when they emphasized that there is no rational way to derive values. Rationality can NOT tell us where to go; at best, it can tell us how to get there. To delegate the destination, the fundamental subjective values that AI adopts, to AI, expecting it to rationally determine what is right or wrong is sheer folly, and must be avoided at all costs.

Using Knowledge Modules as a Base for Further Customization

One approach to fill gaps in (ethical and other types of) knowledge, discussed in cited PPAs and elaborated here, is to use knowledge modules as a base for further customization. These modules are essentially sets of training weights that can be combined with an LLM's existing weights so as to change the behavior in known and predictable ways, especially if combined with an "out-of-the-box" LLM or other base model whose weights and training history are known.

For example, suppose a human owner of an AI is a devout Christian. One knowledge module might be the "King James Bible Package," where an off-the-shelf LLM like GPT-4 is exactly as released by OpenAI, with the exception that it has been extensively trained on content found in the King James version of the Bible, including the Old and New Testaments. The exact training corpus is available, as well as benchmarks describing how the Bible-trained LLM's behavior differs from the out-of-box LLM behavior.

The devout Christian human owner of the LLM, who wishes to customize their AI, could simply purchase GPT pre-trained on the Bible package, or purchase the Bible package itself and train the LLM using the package to modify its behavior. Then, using the Bible-trained GPT as a base, the human owner could proceed to train and customize the AI further, to reflect the owner's particular interpretation of Biblical precepts.

Perhaps, for example, the owner wants to emphasize the golden rule and the ideas of charity, forgiveness, and mercy that are found in the New Testament, but minimize the influence of certain Old Testament passages that don't fit with the owner's modern ethical sensibilities. Rather than customize an AI from scratch, the owner could customize a pre-trained AI, saving time and effort. This customized pre-trained AI could then represent the owner's values in a group of AI and human agents that are engaged in training other AIs to be ethically aligned with human values.

In one implementation, groups of humans can defer or "delegate" to pre-trained modules that do a good job of representing their positions on ethical issues. Delegation in this context means delegating their ethical influence (in a one-human, one-ethical-input model) to the pre-trained AI so that the Bible-trained AI, for example, could vote on the ethical preferences of multiple humans.

While inferior to having each human train their own AI explicitly with their own values, the delegation approach to securing input may increase the total number of humans represented in an AI's value system, even if many are represented by proxy. Further, if the Bible-trained AI (in

our example) is the base AI that then becomes more customized over time by observing and interacting with its owner, as the AI learns the nuances of the owner's ethics, it will modify its ethics to more closely reflect how its owner interprets the core precepts in the pre-trained bible version.

From a practical standpoint, the use of pre-trained modules, both for ethics and for other types of knowledge, is likely to be required in order to advance AI capabilities, ethics, and knowledge rapidly, while still securing input from many individuals.

In the preferred implementation, there is a marketplace for such knowledge models, such that human owners have maximum choice as to which pre-trained model (or modules) is/are most appropriate for the owner to use as a starting point for further refinement.

DETAILED IMPLEMENTATION EXAMPLE

Building a scalable, safe AGI is complex. Many combinations of the inventions described above are possible. In this section, we provide some specific examples of preferred implementations of combinations of these inventions while recognizing that many other combinations are possible. Enumerating all combinations is impractical given the space constraints of this patent disclosure, but the following examples, together with the description above, should be sufficient for developers skilled in the art of implementing intelligent systems to implement safe, scalable AGI systems.

META Implementation Scenario

For specificity, consider a scenario that could be implemented by a company such as META, the creator of Facebook, Instagram, Reels, and its version of the Metaverse. META has a tremendously valuable asset for implementing scalable and safe AI in the form of its huge user base. The data contained in the social graph and content posted by billions of users of META products and platforms is certainly large enough to provide a representative and statistically valid sample of human values and ethics. To the degree that certain geographies or populations (e.g., China, or other non-US populations) are under-represented in META's data, the company is large enough to form partnerships to access a representative sample of data for these populations from other companies (e.g., TenCent, Baidu, etc.).

Further, META has a built-in base of users for customized AI agents. Every user of the META product or platform could benefit from a personalized AI agent that represents that user. In the simplest use case, such customized AIs could help META users filter content and post user-generated content more efficiently, with less attention and time required from the user.

Because META users have, in aggregate, posted huge amounts of preference and content information to the META products and platforms, it is relatively easy for such users to specify the creation of customized AI agents. With the press of a button, users can specify that a base LLM, such as Llama, GPT 4, Bard, or any other in-house or third-party LLM, be tuned on the user's individual data so that it behaves more like the human user in terms of its preferences, knowledge, skills, and ethical values.

The trained LLM agent can be immediately used by its user/owner, filtering out unwanted content and searching the user's photos and other content to generate suggested content postings on behalf of the user. By watching the user's actions, the LLM agent can fine-tune itself, learning more and more about the user.

To the degree that the human user participated in META's Metaverse, even more data becomes available for training Base LLMs and customizing the user's LLM agent. In the Metaverse, AI can watch not only what users post, but every eye-blink, motion, and behavior in the virtual world. This incredibly rich dataset of user behavior is a gold mine of information that can be used to tune AI models much more effectively than the standard subsets of internet data that are widely used today.

The techniques for such training can be standard machine learning techniques, including without limitation, deep learning and neural network techniques, transformers and other algorithms, multimodal algorithms to take advantage of the full range of multi-modal information available in the metaverse, and techniques that leverage the preferences already captured in META's existing ad targeting and content-recommending engines.

META has a huge advantage over other potential organizations in that it can leverage its scale to improve the base models and leverage all the work done to create ad and content targeting to tune models for individual preferences.

As META is concerned primarily with connecting humans and improving the quality of life for all its users, it naturally will want to ensure that the AI agents it creates are safe. Again, META is in an enviable position to combine the values and ethics of its users to create customized AI agents that reflect the consensus of all human users with regard to basic and fundamental ethical principles while at the same time allowing variation in individual ethical judgments for specific scenarios. Such ethical customization is primarily a function of the availability of large representative datasets, sufficient compute power, and powerful machine learning algorithms – all of which META possesses.

Some individuals will want to invest more time than others in training and tuning their AI agents. Some individual users possess specific and valuable knowledge that increases the value of their trained AI agents.

META can serve as a platform for training and customizing these individual AIs and a repository of the training datasets and methods for achieving the various customizations. META can also serve as a marketplace, allowing individual users to share, trade, or license the individual “knowledge modules” that users have created to tune and customize their agents. Popular training packages can be shared, exchanged, and bought or sold among users, with users owning the data they create and META serving as a broker (for a fee).

In an alternative implementation, META can own all the data, identify the most successful training datasets, and offer knowledge modules to its users for a fee and/or as an incentive for using META platforms and products. Such datasets and training methods can also be licensed or otherwise exchanged with another third party (companies or organizations) to increase the value of their offerings and for the common good – i.e., to provide representative ethics and values that all AI developers can use.

In this preferred implementation scenario, the first challenge is to create the base LLM that individual META users will then tune. META could leverage the inventive approach described in this invention rather than use the existing conventional approach of deep learning combined with Constitutional AI or conventional RLHF that requires many employees or contract human workers.

What follows is a four-phase process that constitutes one preferred implementation of some of the inventions described above. Note: at multiple places in the process outlined below, it is also possible to record the ethical values of AI via blockchain or other auditable data structures, providing a way to check and ensure that the values being stored and acted upon are what humans intended.

Four-Phase Process

The process might follow four general phases: I) Train the Base Model; II) Customize the Base Model for each User; III) Combine (Ethics) Knowledge from Multiple Customized AI Agents; IV) Refine AI Agents via Interaction with Other AI Agents and Humans Solving Problems.

While we have been primarily using ethical knowledge or information as our examples throughout this patent disclosure, including in the phases below, it is important to understand that all the steps and inventions apply to ANY knowledge.

Because ethics/values are the most important type of knowledge to address when it comes to increasing the odds that humans survive and prosper in the age of AI, we have mainly chosen examples from that domain. However, other knowledge domains likely will comprise the most frequent use of the inventions described in this and the other cited patent applications.

PHASE I – Train a Base LLM Model with some safety/ethical guardrails (and/or knowledge)

1. Train the base LLM using existing subsets of internet data and proprietary datasets that reflect the overall composition of the content generated by META's target user groups for the LLM.
2. Identify a large corpus of ethical and safety-related scenarios for training the LLM to make it safer than the initial base model without such ethics/safety training.
3. Use a variant of Constitutional AI in which a trusted earlier LLM is used to efficiently cover some of the most common safety scenarios with oversight and RLHF from META.
4. Offer META users the opportunity to help improve the safety of the META AI model in exchange for incentives such as free or reduced cost to use a personalized version of the LLM for their own needs.
5. Solicit a large and diverse set of additional safety and ethics scenarios from users, essentially crowdsourcing the generation of potential new safety cases.
6. Have trusted users and/or META employees filter and refine the set of safety cases to achieve path coverage as discussed above, taking the frequency and impact of cases into account.
7. Use RLHF from trusted users and/or META employees with redundancy so that the ethical behavior being taught to the base model is never reliant on a single user/human's input, and so that the most impactful/frequent cases have the largest sample size of human user input.
8. Additional testing (within META) on edge cases and a sample of impactful cases will be performed to determine when a safety threshold has been achieved.
9. Release the base model to a select group of users who will provide additional feedback on their ethics and specific test cases in exchange for being in the program, thus allowing

a much larger sample of user input to refine the safety and ethics of the base model further.

PHASE II Customize the Base LLM to Each User's Individual Ethics (or Informational) Profile

1. Assemble a corpus of ethical questions based on various ethical assessment instruments well-known in the art, supplemented by additional questions developed by META based on data on META's users and further questions solicited from (crowdsourced) META's users directly.
2. Use statistical techniques, well known in the art, to assign regression weights to the questions such that a ranking is achieved whereby higher-ranked questions provide more useful ethical information than lower-ranked questions.
3. As each user interacts with their AI agent to tune it to the user's needs, the AI agent engages in a conversation with the user. This conversation is driven partly by a standard set of ethical questions that have been determined to efficiently elicit basic ethical information from users.
4. However, the conversation also includes some questions that are driven by the degree of missing ethical data (across all users) for the questions that have been ranked.
5. For example, if the most important safety/ethics question was: "Under what circumstances, if any, is it appropriate for an AI agent to directly harm any human?" And if the second most important safety/ethics question was: "Under what circumstances, if any, is it appropriate for an AI agent to disobey a direct order from a human?" Then META would track when sufficient data from human users was gathered on the first question before moving on to the second question. Whichever safety/ethics questions were ranked highest in importance and were lacking sufficient data would be added to the questions AI asks users in a conversation. In this way, since humans have limited willingness and ability to answer questions from AI, the most important questions for which data is still lacking are always addressed first.
6. In addition to eliciting ethics from users via conversation and question-asking, AI can learn users' ethical views by analyzing (ideally with the users' permission) all the content posted and other information that META has acquired on that user. Using analysis techniques well known in the art, including, but not limited to, prediction of users' answers

to ethical scenarios based on correlations between that user's data profile and the answers of other users with similar profiles (similar to the algorithms used in recommender systems) META can automatically tune the user's customized AI based on the user's existing data with a "single button press" authorizing such tuning. The weighting of information by recency, type, and/or source of the information is a parameter that could be determined by the user and/or META as discussed above. That is, for example, more recent information could receive exponentially more weight than older information if META determines that the rate of new relevant content has been exponentially increasing over time, or if META or the user simply feels recent information is much more reflective of the ethics that need to be conveyed to the AI being tuned. Similarly, the user or META may decide that certain types of content (job-related posts about office ethics) should receive more weight than other types of content (pictures and descriptions of wild partying in the Bahamas) when tuning the AI on the user's (or user and/or META-selected) content.

7. In addition to the techniques of eliciting ethics information via conversation (3 & 4) and analyzing user data (5); weights from AI agents that have already been trained by users with similar profiles to the current user can be directly combined with the current user's AI as a fast and efficient method of improving the tuning of the current user's AI. The already-trained Agents can either be directly selected by the user, as in "I want my AI to be trained using the weights of the AI agent belonging to the Pastor of my church," or can be automatically selected based on META's algorithms for finding similar users, or both. A user may be offered a choice of other AIs and allowed to interact with them before deciding which of these AIs the user wishes to have included as sources for direct weight combination.
8. Users can choose, and/or META can suggest and/or select, various existing (customized) LLMs or AI agents to interact with and provide ethical feedback to the agent that the user is trying to customize. Further, the users can specify the degree of training or impact desired from each training source. For example, if existing customized AIs exist that have been tuned on the ethical precepts of a Baptist Preacher, a Tibetan Monk, an Islamic Cleric, and a Humanistic Philosopher, one user might choose to have the user's model be tuned by interactions only with the Tibetan Monk and the humanistic philosopher with 80 percent of the training being done by the Monk AI and 20 percent by the Philosopher AI. Metadata (e.g., the reliability and trustworthiness of the source), as discussed above, can be used to modify the degree to which tuning from a particular source is allowed to adjust weights. Finally, consistent with the principle that "humans in the loop" are a primary means of catching AI errors and ensuring human-aligned values, it likely will be desirable to weight input from humans more strongly than input from other AIs, and also

to weight input from the user/owner of the AI being customized more strongly than input from other non-owner humans.

9. Users and/or META could specify the frequency with which the tuning of the user's model will be updated, and the degree to which such updates are done automatically using passive methods of training on user and other data or actively, requiring conversational involvement or other decision-making by users. Users and/or META can specify alert conditions -- e.g., a new Supreme Court decision affecting LGBTQ rights, or news events with ethical implications -- as triggers for event-based updates to the training.

PHASE III – Combining (Ethical and Other) Information from Multiple Customized AI Agents

Once an individual customized AI agent has been created using the steps outlined above, or some combination of them, the individual AI agents can combine their ethics knowledge to create a representative and more comprehensive ethical foundation for SuperIntelligent AI.

1. Weights from multiple customized AI models can be combined on a “one-vote per AI and one AI per user” basis to achieve a representative and statistically valid set of AI ethics that generalizes across all the humans who are represented by their respective customized AIs. Different groups of humans constitute different human populations and may have different ethics. However, in the broadest implementation, weights from as many AI different agents as possible (provided that each human is allowed to submit only a single custom AI agent to represent his/her/their values) can be combined to create a value system for AI that broadly represents the values of all humans on Earth. Such a value system would be a representative and statistically valid way of aligning AI with human values and would likely reduce the probability of human extinction by AI.
2. In lieu of, or in addition to, the method of directly combining weights as discussed in (1), many customized AI agents can be presented with ethical dilemmas and allowed to vote on the best actions to take based on the values of each AI, with each AI possessing one vote. To minimize the possibility of more intelligent AIs tricking (or exploiting knowledge gaps in) less intelligent AIs without knowledge of the owners, AIs (optionally) could check with their human owners before casting votes on (important) ethical matters. Further, safeguards could be built in, triggering a human (or other) review of the voting process if conclusions reached do not align with generally accepted ethical standards and precepts, such as valuing human life.

However, as AIs become increasingly intelligent, the burden of accurately representing the value system of human owners will fall increasingly on the AIs, since eventually, humans will not

be able to keep up with the speed of thought or the number of scenarios that AIs are considering.

Returning to the “spinning wheel” analogy, the key principle guiding the implementation of safeguards and accountability of AI to humans must be that human input is preserved on matters closer to the center of the wheel, which are more central to the alignment of AI with human values and the purpose of AI’s existence. As long as core human values, such as the golden rule and love for all humans, are preserved near the center of the wheel, AI can make many decisions relatively autonomously on the spinning periphery of the wheel. This approach can preserve human-aligned values even in the coming age of SuperIntelligences that become trillions of times smarter than humans.

PHASE IV – Refining Values Based on Problem Solving

Once individual AI agents have been customized, and a normative set of values and ethics has been derived from a combination of the knowledge of the AI agents (Phase III), the AI agents, together with (optionally) human agents, can collaborate to solve problems in a collective intelligence approach. The intelligence exhibited by the network of agents will be superior to the intelligence of any one agent, following the principle that “many heads are better than one.”

Techniques for Optimal Combination of Agents in Collective Intelligence Problem-Solving

An important element of the invention of using the collective intelligence of multiple agents is to select the appropriate agents such that the skills and knowledge of the agents are optimal for performing particular tasks. For example, if the problem is to bake a cake, but all the agents know only about theoretical physics, the resulting solutions may be suboptimal.

Many methods and techniques, well-known in the art of assembling optimal human teams, can be applied to selecting optimal groups of AIs (and/or human) agents. Without limitation, these techniques include:

- A. Methods for profiling the skills/knowledge of the AI (and/or human) agents and matching them to the requirements for the task using quantifiable metrics and numerical matching algorithms;
- B. Seeking optimal coverage of the required knowledge and skills using the minimum number of agents;

- C. seeking multiple agents with redundant knowledge/skill coverage in areas where the stakes are high;
- D. Causing the number of agents with specific domain knowledge to be proportional to the estimated importance of that knowledge in the solution;
- E. Determining the knowledge of a (human and/or AI) agent by surveying, assessing, or having a conversational interaction with the agent to score the agent's knowledge in a particular domain;
- F. Determining the knowledge and skills required for a problem by surveying, assessing, or having a conversational interaction with the customer or agent proposing the task to categorize and numerically rank or rate the knowledge needed in various domains;
- G. Using objective and subjective reputational and other metrics to rate or rank the estimated quality of the agent's work in various task domains or on various specific tasks;
- H. Using the record of problem-solving behavior on similar problems to determine the likely effectiveness of an agent for work on a similar new problem(s);
- I. Using objective metrics, including, without limitation, solution time, cost, quality ratings of solutions or deliverables, and schedule-related metrics (e.g., on-time delivery percentage) to help determine which agents are best suited for the requirements of a particular task or task domain;
- J. Using metrics of how well different (human or AI) agents have worked together in the past to help determine optimum combinations of agents; for example, two agents might have synergy where the solutions produced by the combination of both agents are much better than the simple sum of the contributions of each agent individually might imply. This could be because the agents complement each other or otherwise enable solutions that could not be produced by a single agent (or other combinations of multiple agents) on their own;
- K. Dynamically adding or subtracting agents as problem-solving progresses and/or as new information as to the types of expertise needed for remaining problem-solving steps are determined;

- L. Using market forces (e.g., putting tasks and sub-tasks out to bid in a marketplace(s)) to enable agents to efficiently self-select for participation in the group that is solving the task;
- M. Using social-graph, networks, and affiliation information to select agents that are associated with other agents that have proven themselves effective at solving particular problems;
- N. Enabling agents to refer (with or without compensation) other agents for inclusion in the group working on a particular task;
- O. Taking cost, location, speed, availability, and other metrics into account in a statistically valid way (including without limitation use of regression, neural network weights, machine learning approaches, and other quantitative predictive methods) when assembling a team of agents (e.g., with specific cost, schedule, and/or quality constraints);
- P. Using values and/or ethical knowledge and beliefs of agents to ensure common values and ethics of all agents working within a group;
- Q. Using multiple groups of agents to solve the same problem or sub-problem independently so that solution quality and other performance metrics of the multiple groups can be compared dynamically in real-time (or post-hoc with time delay or asynchronously) and work can then be (dynamically) routed to the groups with better performance on certain tasks or sub-tasks;
- R. Using multiple groups for simultaneous problem-solving as in (q) but then comparing the solutions such that the solutions that most groups come up with are favored (as being a consensus view of multiple teams) compared to minority solutions. This approach may be useful for especially high-stakes (e.g., irreversible) decisions with important consequences, or when buy-in from the team participants is desired.
- S. Using multiple agents or groups of agents, where the input of the group members is weighted differentially depending upon the knowledge, expertise, or track record of the groups or agents (e.g. for the plumbing part of the problem, weight the solutions and suggestions of the agents with plumbing knowledge more, but for the fund-raising part of the problem, weight the solutions of suggestions of the agents with finance expertise more);

- T. Use of matching algorithms and techniques similar to those used to target ads to individuals, but in this case targeting work (rather than ads) to human and/or AI agents;
- U. Using targeted ads to recruit humans (and/or their AI agents) to work on specific tasks, using ad-targeting methods well known in the art; and
- V. Leveraging existing data (e.g. LinkedIn profiles or other information available on the internet) about agents to automatically rate and qualify them on various dimensions, including without limitation those mentioned above, so that the agents can be more effectively matched to specific tasks, subtasks, or to other agents with whom they are likely to work well.

Overcoming Limits of Bounded Rationality

Generally, the more agents that are involved, the broader the collective range of knowledge, skill, experience, and values the network will draw upon. As the Nobel Laureate, Herbert A. Simon has shown, Bounded Rationality, that is, the information processing constraints, are a primary determinant of the limits of human intelligence. What is true of human intelligence is true for any intelligent system, including those composed of human agents, AI agents, or both. SuperIntelligent AGI will be less constrained by information processing constraints than humans with more limited brains, but still, SuperIntelligence will still find its rationality constrained in the limit by the scope and speed of its cognitive abilities.

Ethical Problem-Solving Considerations

To the degree that a network of agents solves problems directly related to ethical scenarios or with ethical implications, the ongoing problem-solving experience of the network of agents can refine and improve the ethical decision-making of each individual agent and the system itself. The refined ethical scenarios (and the stored solutions and/or weights that represent these refinements) can be shared with the individual agents on the network. Individual customized AI agents can learn (ethical and value-related) information from the work of the collective.

The detailed steps and methods for storing procedural solutions and producing SuperIntelligent AGI from the combined efforts of multiple individual AI agents are described in other PPAs cited above.

Since all problem-solving behavior involves setting goals and subgoals, one of the key ethical safeguards (which we reiterate here) is that a series of ethics checks must be passed each time a new goal or subgoal is set. This procedure is an automated way of ensuring that ethics and safety concerns come into play whenever any problem is solved by an agent on the network.

The discussion above (e.g., Phase III) relating to creating a representative set of ethical values for AI can, and should, inform the ethics checks that are passed at each stage of problem solving. That way, as ethical norms are updated as discussed in this invention, the overall ethics of the SuperIntelligent AGI that results from the collective problem solving of multiple AI (and human) agents is also automatically updated.

Ensuring that the ethics of SuperIntelligence reflect the continually updated ethics of the individual agents comprising the SuperIntelligent system is a primary means of ensuring that AI, many times smarter than humans, remains aligned with human values.

Role of Humans as AI Surpasses Human Intelligence

As mentioned earlier, what applies to ethical knowledge also applies to other types of knowledge. Values, ethics, and domain-specific knowledge are all constantly changing. Values generally change more slowly than other types of knowledge because they are closer to the center of the “spinning wheel” described earlier, and they are the most important type of knowledge for determining the behavior of SuperIntelligence.

Even in a world where AI is trillions of times smarter than us, humans can retain a role as the source of values at the center of the rapidly evolving knowledge base and intelligence that is emerging.

This invention shows some ways to efficiently and effectively center AI, AGI, and SuperIntelligence on values that are aligned with those espoused by a representative and statistically valid sample of all humans on Earth. If humans can center our attention, thoughts, words, and actions on love, SuperIntelligent AI will perceive love, learn to love, and use its intelligence in the service of love.

As the psychologist Viktor Frankl pointed out, humans as intelligent beings have a driving need for purpose and meaning. Humans must design AI, AGI, and SuperIntelligent systems to meet this need as well and look to humans to supply purpose and meaning. Our survival may depend upon it.