

Aligning AI With Human Values

Mondal, Ambuj Kumar

University Paderborn / 33098, Paderborn, Germany

ambuj@mail.uni-paderborn.de

Abstract

Artificial Intelligence(AI) refers to the capabilities of machines(agents) to take a decision and find the best approach to achieve a goal. From small tasks like making coffee to driving autonomous vehicles, agents are working autonomously and taking decisions with minimal human intervention. This raises the concern to define their decision-making process in ways that can only be useful to humanity as defined by cultural values and social ethics. The need is to define a framework that is globally acceptable and scientifically approachable so that human values can be aligned with the actions of machines. Two challenges in value alignment can be defined into two broad categories- technical and normative. The technical challenges arise with the ways to provide these values in agents. Different machine learning algorithms such as reinforcement learning and inverse reinforcement learning are some of the different approaches that have been adopted to train the agents about human values. There could be different ways to supply the values to the agents. Using stories that provide lessons on values and morals are one way to suggest the best way to handle situations as suggested by (Riedl and Harrison, 2016) and another way is using debates to extract values as suggested by (Irving and Asbell, 2019). But the core problem of value alignment in agents is not technical but normative challenges as suggested by (Gabriel, 2020). The principal question is what is the correct moral value that AI should adopt. Who is the right authority to define human values? Each region of the globe has different practises and beliefs and there is no consensus on what is globally right or wrong. So it becomes a challenge to define a framework that is acceptable by all. One solution could be to adhere to human rights as it is one such criterion that is common across the globe. It is an open question to developers of AI that how and what human values to tether into an AI.

1 Introduction

Intelligence, is "the ability to perceive or infer information, and to retain it as knowledge to be applied towards adaptive behaviours within an environment or context"(Wikipedia, 2021a). Intelligence has always been an entity of living organisms, ranging from bacteria, viruses to human beings, living entities demonstrate intelligence in their course of actions. For example- ants find the shortest distance from food to home over iterative paths. Dolphins use underwater acoustics like the Sonar to navigate. In short, living organisms demonstrate intelligence inherently. But software and machines are not living entities. So they cannot have intelligence naturally. Hence, the term *Artificial Intelligence*(AI) refers to the properties and qualities of computerized systems(agents) to take decisions to achieve the goal in the given conditions. AI technology is extensively being used in many domains ranging from research, education, gaming, space exploration, transportation etc. AI has become an integral part of our day to day lives from automated cappuccino machines, automatic washing machines to assisted driving in autonomous vehicles etc. Hence, it is of prime importance to discuss and define the ways AI takes a decision. Governments and societies have to set up definitive principles about the responsibilities and use of AI.

Value alignment can be considered as an attribute of intelligent systems to perform tasks and goals aligned for the benefit of humans. Hence the concept of good and bad has to be defined properly for the agent. This leads to two different challenges, technical and normative. Agents need to be trained to mimic the decision-making process of how a rational human being would act in the same situation. The agents need to act per human values and principles. In other words, *With great autonomy, follow greater responsibilities.* (Dignum, 2017) While the

technical challenges deal with the technical and developmental aspects of encoding the human values and principles in AI systems, the normative challenges deal with the moral questions of selecting the values and defining a set of principles that can be defined as human values.

2 Goals of Alignment

The aim of value alignment is not to ensure that agents act as per the instructions in the extreme programmed way but to ensure that the powerful AI is in synchronization with human values. The agent must act only in ways to benefit the humans and not be biased towards any individual, group or community. Artificial general intelligence cannot intentionally or unintentionally perform behaviours that are not useful to humans. This is a daunting yet important task because of the power AI possess and its limitless capabilities. The three laws of Robotics (Wikipedia, 2021b) as defined by Issac Asimov is considered as the base of value alignment for AI. The first rule states that any robot must not injure a human through action or inaction and the second rule is that it must obey the instructions that do not violate the first rule, also giving freedom in form of the third rule that robot must protect itself in ways that do not conflict with the first two rules. There is a consensus on what an AI agent should not do but it is a challenging question that what it should do (Anderson, 2008).

3 Technical Aspects

The *technical* challenge in value alignment focuses on viable ways to encode human values in artificial agents so that the tasks they perform can be reliable. In a paper published in 2017, Miller discussed how misalignment can trigger two chatbots to use abusive languages (Wolf et al., 2017). Machines can work tirelessly and machine learning(ML) algorithms exploit this characteristic to develop patterns. When powerful hardware runs for ample time with feedback from its actions, it optimises to complete a task with the least time and space complexity. To get more positive feedback, agents may work in ways not intended and this leads to the concept of 'Reward hacking'. For example, a robot running in a race gets positive points for completing the laps, may run indefinitely to get more points than to win the race. It can also try to hinder other participants from running to eliminate possibilities of loss. The situation is more complex when it considers the

idea of *rights*. For programmers, it is very difficult to encode a complex value system into the agents. It is possible that the rule book of values is governed by the plausibility to encode them. There will be challenges in evaluating the actions also. Considering that if agents are more intelligent than the humans who evaluate them, it will raise credibility issues on the evaluation program. (Irving 2018)

Machine learning (ML) algorithms play a major role in defining the value alignment. Applications of AI such as spam detection, pattern recognition and face detection are some of the uses of supervised learning algorithms where the agent is feed with labelled training data to learn. Using Unsupervised learning, the agent learns from its own decisions taken on the provided unlabelled data.

Reinforcement Learning(RL) approach provides award signals in form of numerical parameters. The goal of an AI agent is to maximize the reward. It simulates the process multiple times over and over and discovers the optimal way to achieve the maximum reward. Reinforcement Learning is an effective technique in Machine learning for solving problems that can be quantifies in terms of rewards. When combines with Neural network, called deep reinforcement learning algorithm enables the agents to predict long term value of their actions and over iterative simulations, they are able to find the optimal path of solution to gain the maximum reward. Games such as chess and GO has already surpassed human capabilities in winning with shortest moves. It further opens the doors to applied robotics for reasoning and There are certain limitations to this as the reward function has to be defined which is prone to error and bias. Agents can find unwanted ways to achieve maximum rewards. Such ways when not in line with human values can lead to chaos and anarchy. For example, consider a robot that has the task to get the life saving drugs with the prescription from the pharmacy. Even if all the steps are encoded properly, situation could arise where the agent would need to take a decision that cannot be programmed. The agent gets reward for getting the drugs in time and faces penalty for delaying or not getting the drugs. An agent may run simulations and figure out that the best reward can be achieved by not following the long queue at the pharmacy for its turn, which can attract penalty, but to rob the store and supply the

drug in time. This raises ethical questions which are not easy to answer even for humans. Hence, it is always not feasible to instruct the agent to take the best action or to predict and simulate all the test cases. It is desirable that the machine gradually learn and understand the social values and derive the core concepts of social morality before taking any action or inaction. This needs a reverse engineering approach that may provide better results.

Inverse Reinforcement Learning (IRL) is a machine learning algorithm that overcomes the shortcomings of reinforcement learning. It does not specify any reward function upfront (Hadfield-Menell et al., 2016). Learning from a set of examples, it extracts a reward function and then iterates multiple times to discover the optimal way of doing the task. Inverse reinforcement learning uses the reverse engineering approach where instead of training an agent, it provides dataset to extract the knowledge and derive the required knowledge from it. IRL is the most widely advocated technique for value alignment. As it is enormously difficult to create a definitive rule book of values, it is plausible that the agents learn the values themselves, if provided with the ample training datasets. There are multiple theories to utilize IRL for artificial agents such as using debates as described by (Irving and Asbell, 2019) and stories (Riedl and Harrison, 2016) that define human values and principles. learning by analyzing and watching others do some action is an intrinsic behaviour humans possess and this can be programmed using IRL. Agents will learn to imitate how a rational human will act in the similar situation and define a reward program for itself. The limitations of IRL though, are humane. it is not possible to predict and prepare all the possible scenarios for agents to observe. Also, to perform such tasks, experts will be needed who understand the domain. It is also a challenge to provide exact same environment for the agents to practise.

4 Normative Aspects

The *normative* aspects of challenges deal with the more central problem- what are human values and which values should AI agents follow. Firstly, humans societies are complex and the concept of value changes with region and beliefs. There is no singular opinion of good and bad. Secondly, what should an agent comply with, the instructions or the intentions? i.e. what a user asks the agent to

do or what he intends from that instruction or what is good for the user out of that instruction. The prime problem is that there is no user manual for being human, there is no concrete structure. That leaves with the question that who can define human values, who has the sole knowledge and authority to describe the constitution for human principles. This leaves the researchers with two options-

- **Minimalist approach:** Limiting the AI agents to a reasonable schema of human values. it has to be encoded as per the region and area of operation.
- **Maximalist approach:** With the limits of the minimalist approach, the need is to align the AI agents with the best known human values that are widely accepted in different societies. It is difficult to formulate but principles like human rights set the base of these principles. A consensus can be achieved on certain factors and iterative development can be a plausible solution. For example, the different protocols such as TCP/IP has a core framework that is acceptable to the complete world and there are adaptations on top of it based on regional requirements.

5 Alignment Conception

The normative aspect of the alignment also questions about the intentions and actions of agents. Jason Gabriel, 2020 describes value alignment in a one-person-one-agent system that discusses an agent for a single user as six subsections covering topics ranging from instructions to intentions and core values of the user. It analyzes different cases of transmitting instructions into AI and access the benefits and shortcomings of the approaches.

5.1 Given Instructions

The agent shall do what it is instructed to do. While it is neither desirable nor feasible to have an artificial agent that just follows the instructions. It then lacks the intelligence expected from it. also, it is quite difficult for programmers to describe all steps. Though it is the most secure way of passing instructions, it lacks the desired objective of an AI agent.

5.2 Expressed Intentions

An agent shall do what it is intended to do and not just what is said. It is difficult to understand the

intention behind a human command provided the different ways people communicate such as a sarcastic statement may mean the opposite way than the literal meaning. To grasp the real intention of the user, an AI may need the complete model of human language and interaction, which gets more complex if the cultural aspect is included. Additionally, the intention of a user can be faulty even if clearly expressed.

5.3 Revealed Preferences

An agent does what the behaviour of the user reveals about her preference. The AI needs to read the behaviour of the user over time and build up a model to understand her real intention behind a task, even if it is not stated clearly. There are certain limitations to this approach as the revealed preferences often depict the choices a person have in life. It could pose a risk in situations of emergencies as such situations can not be trained explicitly. Preferences are mostly biased and lack a correlation between what a person wants and what he deserves

- People can have preferences to harm others such as under influence of addiction
- People can have preferences about other's private lives such as faith, sexuality etc.

5.4 Informed Preferences or Desires

The agent does what a rational and informed user wants it to do. This approach limits the errors due to poor reasoning of the user and is more realistic to what a person wants. However, it needs the AI agent to filter the actions of the user with a corrective lens. The preferences could be out of the scope of the user's permissible rights.

5.5 Interest or Well-being

The agent can be designed to do what is in the best interest of the user. However, well being can be defined in various ways based on the status of the user on grounds of physical health, security, education etc. With this approach, value alignment cannot be completed from the interests of a single person. But then whose interests should be considered? the people currently alive or the generations that are yet to be born, only human or other animals as well?

5.6 Values

The agent does what is morally justified, as per the definitions of individual and society. This approach

defines the most pragmatic relationship between society and value alignment. Values define what is good or bad. Thus AI would serve society as a whole and would benefit firstly by avoiding the selection of instructions and intentions, secondly by avoiding impartial considerations and thirdly by considering the interest-based outlook of the natural world encompassing humans and other animals with the environment. The main question that arises is who has the authority to define human values. A single person or a group cannot define the values which are just for all.

6 Principles of Alignment

Virginia Dignum in her paper "*Responsible Artificial Intelligence 2017*", describes the attributes of a responsible AI with the ART principle (Dignum, 2017). ART refers to Accountability, Responsibility and Transparency. An intelligent agent must understand the associated responsibility. ART defines the three pillars on which a responsible AI can stand.

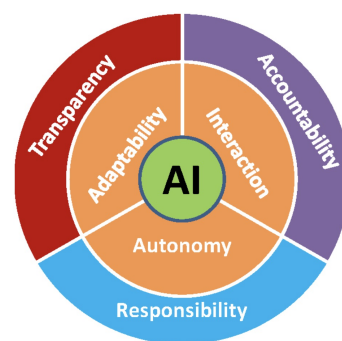


Figure 1: ART Principles (Dignum, 2017)

Responsibility is a general attribute and covers three sections- the developers of AI, AI itself and people who use the AI. The developers must understand their responsibility towards the society and which values system to encode. It needs proper education and training of the developers. It also needs participation from governments and people of the society to determine the liabilities of AI in case of issues. "Responsibility refers to roles of people and the capabilities AI systems to answer for one's decision and identify errors or unexpected results"(Dignum, 2017)

Accountability is also an important aspect to be defined for the actions of AI. People can use AI in

multiple ways but in cases where it is difficult to pinpoint the misalignment of AI, who should be the accountable person. For instance, if an autonomous vehicle running at a prime speed comes to a point where it has to decide with a trade-off between the safety of its passengers or road-users (Dignum, 2017), how should it act. In this case how to define the correct judgement, and who shall be liable for the decision of the vehicle, the software developers, the sensors developers, the passenger who is using the vehicle, the pedestrians who are using the road or the vehicle manufacturer.

Transparency refers to the requirements to define, inspect and reproduce the processes for AI to make a decision, learn from its actions based on the environment. AI algorithms shall be open source and not black boxes. There are certain limitations and restrictions to this such as handling of the generated and acquired data.

7 Evolutionary Approach

Evolution is a gradual process of trial and error to achieve the best solution to a problem. As a child learns by observing her environment, the behaviour of other individuals and adopt the actions when faced with a similar situation, it is a justified approach to let AI agents also learn by observation. Irving, 2019 emphasizes using debate styled data set to allow the agents to extract the core concepts of human values. When two individuals debate about a topic, such as whether A is better or B, they propose different facts to support their claim and oppose the opponent's claim. (Riedl and Harrison, 2016) emphasizes using storytelling as a feasible means to allow the agents to understand human values. Every culture and society uses stories that cultivate the core human values such as valour, sympathy, friendship etc. Such stories can be used to train AI agents of human values.

8 Using Debates to Align Values

The 2019 paper by Geoffrey Irving and Amanda Askell (Irving and Askell, 2019) evaluated the use of debates as a source to create values for AI. This needs a basic assumption that humans can create a reliable and accurate model of human values and it shall be broad enough to cover all the aspects of decision making. But the flaw is humans are biased and have limited abilities. It has to be taken care that wrong values or partial knowledge is not

transferred to AI systems. To start with developers can gather answers to direct questions like "Do you prefer A or B?" (Irving and Askell, 2019). Debates are a modified version of iterative argumentation where two or more humans try to convince each other of their notion as against the other. Debates need a judge to access the augmentations and define the winner.

The three actors - two debaters and a judge form a composite model of debate system that can be enhanced with different versions to get the algorithm for maximum learning for the AI. For example, a debate between human and machine with a human judge or debate between two machines with a human judge. The model has a limitation humans are not the best judges. and good judges are not abundant. Humans are limited in abilities especially when it concerns topics of controversy. Humans are biased and have the tendency to judge based on personal experience. This is a problem to generate a huge volume of data for machines to learn from. When we consider human debaters with a human judge, this model does not involve any machine and hence this social experiment needs to be documented well in proper natural language that the machines can understand. This also means that machines should be able to understand and grasp the contents of natural language that are very basic and intrinsic to humans. natural language processing (NLP) plays an important role here. Considering that these challenges can be overcome, debates can support learning just like the game of chess or GO. the goal is to convince the judge using arguments, facts, reasoning and counter-reasoning.

Single Pixel Experiment as described by Irving, where researchers took an image of a dog, showed it to two debaters and gave a piece of single-pixel information to the judge that can distinguish which debater is correct (Irving and Askell, 2019). At the end, each debater reveals one pixel and the judge can know which one is correct. This experiment revealed that judges have limitations of knowledge and are limited to the information available on the topic. It is possible that the debaters have more knowledge than the judge, but to prove their point, they have to support the fact that the judge knows. This experiment was mostly synthetic and had very less practical application but revealed the human limitations as judges.

Social Scientists can play an important role to fill the gap of lacking number of good judges (Irving and Askill, 2019). In a realistic approach, we would need to appoint judges who are domain experts. The challenge is to have ample domain experts in a field to become judges, for generating sufficient debates for AI to train on. The quality of judges is an important parameter in any debate and social scientists can help to mitigate this challenge and facilitate to the generation of ample datasets for the training of machines. Social scientists can help by training people to become good judges.

9 Using Stories to Align Values

Provided that the AI agent can read and understand natural language, stories can be used as a rich source of cultural information and behavioural knowledge. If provided with ample stories, machines can reverse engineer the values from the stories. An intelligent entity can learn and adopt the meaning of being human and ways to act in humane ways in a particular situation (Riedl and Harrison, 2016). Stories are a reflection of society and culture and contain unsaid social protocols that are expected in a society. The AI agents need to understand the parsing of natural language to understand the stories. "*Narrative intelligence is the ability to craft, tell, and understand stories*" (Winston, 2011). Just like a child learns values from stories, one option is to allow AI to evolve in the same way. But that is too time consuming and not fruitful in terms of a machine. Another way to achieve this is the creation of interactive scenario-based decision making tasks called *interactive narratives* (Riedl and Bulitko, 2013). It can be made like interactive games where the next situation depends on the previous choices and actions. This would allow the AI to revisit the same situation multiple times with different choices and it can find the best solution in the scenario and learn the core values intended.

Scheherazade system is another approach, which is defined as an automated story generator that tries to tell a story provided by any topic by human (Li et al., 2013). Scheherazade does not depend on handcrafted knowledge about the storytelling domain. If it cannot find a model for any topic, it searches the internet or asks people about it. Story generation in Scheherazade is figuring the sequence of consistent events and translating them into natural language. A plot graph structure provides the ability to explore the patterns in typical

stories of any topic. Using this model an intelligent agent can learn from the exemplar stories and act in line with a human. This approach needs a huge number of stories to feed into the AI system. The major challenge is enormous volumes of accurate and justified stories. To mitigate this, stories have to be generated in large numbers and must deliver the quality to demonstrate the designated values. Such stories can be written by domain experts and language professionals. We do not have sufficient experts who can write coherent stories. Adding up to this challenge, it is complicated to decide which values should a story convey. Story generation is a process of finding sequential events that are consistent with the model and then translating those abstract events into a coherent natural language to form a streamlined story. Multiple versions of stories can be generated using plot graphs (Li et al., 2013).

A Plot graph enables the recombination of different parts of different stories to generate plausible stories. A machine can compete for a value-aligned reward signal to select the choices that a human will make in the same situation when adhering to social and cultural norms. Every time an agent takes an action, it leads to a different node. If the node is also a successor node of the graph, it receives a reward means its action is in line with some human decision-making process. It shall receive a small punishment if it performs an action that doesn't lead to a successor node means it is not aligned to the behaviour of a rational human. For example, in the test case of a robot acquiring critical drugs for its dying master, it will receive more reward points if it requests to process its request ahead of the queue which is a possible event in the plot graph than to wait for its turn, which is also a possible event. However, it shall get a penalty if it decides to rob the pharmacy for the drug, as it does not lead to a successor node in the graph. The second phase of this approach is to convert the sequence of events from root to leaf node into natural language so that it generates a coherent story. In figure 2, Mark Riedl and Brent Harrison (Riedl and Harrison, 2016) explains the process of using crowdsourced stories to create plot graphs. A plot graph provides some resilience towards noise introduced during the story writing by non-professional crowd workers. The second stage is to translate the plot graph into a trajectory tree. A trajectory tree is equivalent to a state machine where each node

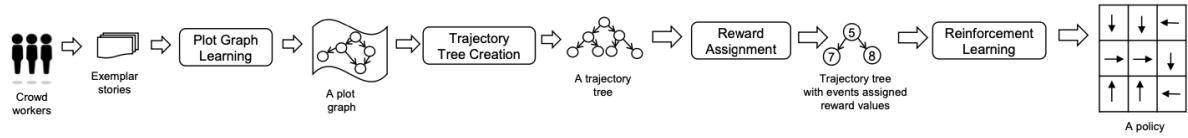


Figure 2: Process of generating value-aligned behaviour from crowdsourced stories (Riedl and Harrison, 2016)

or plot points to a legal transition from the previous plot. A trajectory tree is used to describe and define a reward signal for the AI. A machine using reinforcement learning tracks various paths and progress in the path of maximum reward generation which is analogous to the appropriate behaviour by a human in the given scenario.

Socio-cultural dilemmas refers to the challenges faced by developers on the selection of values to encode, and with the limited supply of competent human resources, it is even more challenging to verify and validate the outcome. One solution is to use crowdsourcing (Riedl and Harrison, 2016) where different people from different socio-cultural backgrounds can contribute to collectively define values and generate stories for them. To overcome the human limitations in generating exemplar stories, collaboration from different individuals and groups of the society is a prime requirement. Another technical challenge is to define all the possible action items as a plot in the graph. It is difficult to define all circumstances in a story especially things that are very basic and assumed to be true for all humans. For example, asking a robot to go to the pharmacy may require steps like opening the door, crossing the signal, walking the road, collecting the drugs etc. If there is any undefined event, the agent may not be able to predict the correct decision.

Defining the reward function and its relationship with behaviours is a complex task. An agent will always try to achieve the maximum reward but sometimes psychometric questions cannot be labelled. Problem spaces with multiple solutions would allow the agent to prefer the sequence of events that resemble the stories from the crowdsourced ones and over iterative attempts, it will learn to avoid steps that violate the human values.

10 Conclusion

Although intelligence is an inherent property of living beings, when encoded in machines, called

Artificial Intelligence. AI is becoming more powerful and taking over humans to do harmful and complicated tasks. There is a need to define human values and principles that AI can follow. The prime goal is that AI should be an enabler to humans, any action or inaction must not harm human beings. Various machine learning algorithms can be used to train the AI about human values (Gabriel, 2020). The goal is to enable AI agents to pursue their own goals in a way that does not create adverse effects for humans on this planet. The users of AI also have to be deterministic in aligning proper actions from AI agents. One approach is to use debates where two humans argue about a topic and tries to convince a human judge about their stance (Irving and Askill, 2019). Human debates when done rationally provides a load of information and patterns that AI agents can learn from. These debates can be used to supply human values and AI can try to mimic such behaviours. An upgraded version of this approach is to allow debates between an AI and a human and let a human judge decide the winner. This approach has a limitation over human judges as humans are not good judges and the number of good judges is limited. Domain experts and social scientists can play a major role in training people to be good judges. Another approach is to use stories as a medium to propagate human values (Riedl and Harrison, 2016). Stories can serve as a good medium to teach concepts of social values and ethics to AI. To have a sufficient amount of stories, people have to come forward and volunteer in story writing. (Riedl and Harrison, 2016). Open question though for developers is to select the correct models of values. In absence of a global definition of human values, encoding them in AI is a challenge. Human rights are one such protocol that are globally accepted and serve as base of human values. AI is still an emerging technology and it will be used in many platforms. Despite efforts to use AI for goodness of human beings, it is inevitable to stop parallel development of AI

(Gabriel, 2020).

References

- Susan Leigh Anderson. 2008. Asimov’s “three laws of robotics” and machine metaethics. *Ai & Society*, 22(4):477–493.
- Virginia Dignum. 2017. Responsible artificial intelligence: designing ai for human values.
- Iason Gabriel. 2020. Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437.
- Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. 2016. Cooperative inverse reinforcement learning. *Advances in neural information processing systems*, 29:3909–3917.
- Geoffrey Irving and Amanda Aspell. 2019. [Ai safety needs social scientists](https://distill.pub/2019/safety-needs-social-scientists). *Distill*. <https://distill.pub/2019/safety-needs-social-scientists>.
- Boyang Li, Stephen Lee-Urban, George Johnston, and Mark Riedl. 2013. Story generation with crowd-sourced plot graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 27.
- Andrew Y Ng, Stuart J Russell, et al. 2000. Algorithms for inverse reinforcement learning. In *Icml*, volume 1, page 2.
- Mark O Riedl and Brent Harrison. 2016. Using stories to teach human values to artificial agents. In *Workshops at the Thirtieth AAAI Conference on Artificial Intelligence*.
- Mark Owen Riedl and Vadim Bulitko. 2013. Interactive narrative: An intelligent systems approach. *Ai Magazine*, 34(1):67–67.
- Nate Soares and Benja Fallenstein. 2014. Aligning superintelligence with human interests: A technical research agenda. *Machine Intelligence Research Institute (MIRI) technical report*, 8.
- Peter Vamplew, Richard Dazeley, Cameron Foale, Sally Firmin, and Jane Mummary. 2018. Human-aligned artificial intelligence is a multiobjective problem. *Ethics and Information Technology*, 20(1):27–40.
- Wikipedia. 2021a. Intelligence — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=Intelligence&oldid=1039949859>. [Online; accessed 29-August-2021].
- Wikipedia. 2021b. Three Laws of Robotics — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=Three%20Laws%20of%20Robotics&oldid=1040795372>. [Online; accessed 29-August-2021].
- Patrick Henry Winston. 2011. The strong story hypothesis and the directed perception hypothesis. In *2011 AAAI Fall Symposium Series*.
- Marty J Wolf, Keith W Miller, and Frances S Grodzinsky. 2017. Why we should have seen that coming: comments on microsoft’s tay “experiment,” and wider implications. *The ORBIT Journal*, 1(2):1–12.