

# Individual assignments

Fall 2020 - IN5480

Mariann Gundegjerde margunde@uio.no

## Table of content

<b>Module 1 - First iteration</b>	<b>1</b>
1 Concepts, definition and history of AI and interaction with AI	1
2 Robots and AI systems	3
3 Universal Design and AI systems	5
4 Guideline for Human-AI interaction	6
5 References	7
<b>Module 2 - Second iteration</b>	<b>8</b>
1 Characteristics of AI-infused systems	8
2 Human-AI interaction design	9
3 Chatbots / conversational user interfaces	11
4 References	13
Appendix 1 - Feedback from iteration 1	13
<b>Module 3 - Third iteration</b>	<b>14</b>
1 Human AI collaboration	14
2 References	17
Appendix 2 - Feedback from iteration 2	17

# Module 1 – First iteration

## 1 Concepts, definition and history of AI and interaction with AI

The history of AI:

*When, and by whom, was the term first used?*

The term artificial intelligence was first used by John McCarthy in 1956. McCarthy was an American mathematician and logician and was known to be the first to use the word. Before the term was officially used, the potential of computing was discussed in the 50s by people like Alan Turing, Claude Shannon and Isaac Asimov. These three all wrote influential papers and did work that shaped the history of AI going further. (Grudin, 2009)

### Definitions

#### **Definition from the Oxford learning dictionary:**

“The study and development of computer systems that can copy intelligent human behaviour.”

This definition is somewhat of a basic understanding of what AI is, and the source for this is a dictionary, so it makes sense that the definition doesn't go into much detail. I would call this the basis of what AI could be described as.

#### **Definition from T. Bratteteig and G. Verne, 2018:**

“AI is a subfield of computer science aimed at specifying and making computer systems that mimic human intelligence or express rational behaviour, in the sense that the task would require intelligence if executed by a human.”

This definition is a recent one, since it's from an article published in 2018. This definition focuses on that an AI should be able to do and process tasks, and especially tasks that if they were done by a human would require some sort of intelligence. This allows us to imagine what kind of tasks an AI would be required to do according to this definition. Also the definition is written by two researchers working with participatory design, which the article where the definition is located is about. This could mean that the researchers don't look at the concept of AI from a very technical perspective, but more from a design-perspective.

**Definition from John McCarthy, 2007:**

“It is the science and engineering of making intelligent machines, especially intelligent computer programs. It is related to the similar task of using computers to understand human intelligence, but AI does not have to confine itself to methods that are biologically observable”

McCarthy describes AI as intelligent machines and computer programs. And it seems like he tries to describe that the AI doesn't always openly show human-like intelligence.

**My own definition:**

It seems like most definitions describe a somewhat intelligent computer system, that could in some way be compared to how humans are intelligent. It also seems like the AI needs to be able to accomplish tasks. My definition would therefore be something like this:

“Artificial intelligence (AI) is a technological system designed for performing tasks with a human-like intelligence.”

**A contemporary company that work with AI**

*Describe how this company present AI on their web pages. In what way does this company talk about AI, as a product, as a service, framework or “idea”?*

A company I found that uses AI is the norwegian bank Nordea. It seems like they use an AI-solution for connecting users with the right service personnel or unit within Nordea. On the front page of their website nothing is presented about AI or the use of an AI, but they have released a press release about their collaboration with an AI-specialized company called Feelingstrem, which delivered the AI-system they now use. The press release i found was published in 2017, and as I don't have Nordea as my personal bank, it was difficult to find out how this system actually works and if it's still in use.

## A fictional film series that is about the use and interaction with AI:

*Describe with your own word how human interaction with AI is portrayed in this work.*

Westworld is a show on HBO that describes a futuristic world, where robots have passed the Turing-test and can no longer be distinguished from humans. Westworld in itself is an amusement park, where humans can enjoy living a life without limits or consequences while interacting with the robots in any way they want. One could say that the villains in this show are the humans, because of the way they treat the robots and how they feel little to no remorse of their actions against them. In later seasons of the show, they are faced with a computer system that has predicted all foreseeable futures and in a way has chosen the best path that humanity can follow based on these predictions. This system is driven by an AI called Rehoboam, and Rehoboam gains control over the world by for example restricting work options for certain people, because it thinks that this is the only way to keep all of humanity on the right path.

## 2 Robots and AI systems

How the word Robot came about:

The word “Robot” was first introduced in a Czech play called R.U.R (Rossum’s Universal Robots), written by Karel Capek and performed in 1920. The word itself comes from a Slovenian or a Czech word “robota” which means forced labor.

Definitions:

**Definition from the Oxford learner dictionary:**

“A machine that can perform a complicated series of tasks by itself.”

Again, as the AI definition from a dictionary, this definition touches on the basis of what a robot is. It is an “independent” machine that can perform tasks.

**Definition by Trenton Schulz, 2020:**

“A robot ... refers to a physical object that interacts with the physical environment, either on its own or via a person, to accomplish a task.”

This definition defines more aspects of what a robot could be and how it can be described. It is something physical that accomplishes tasks by interacting with something. This definition gives you more of an actual picture of what a robot could be, in comparison with the definition from Oxford.

**My own definition:**

I found it difficult to come up with a definition of a robot. One could think of robots in movies and in books, and it's something totally different then when I imagine a functioning robot that exists today. So I decided to focus on what exists today, and try not to think of all possible future robots. My definition: “A robot is a physical object that can accomplish predetermined tasks by interacting with the physical environment or with a person.”

**The relation between AI and Robots:**

*Is “a robot” different from “an AI”? In what ways are they different and similar? Bring in the definitions that you described earlier about robots and AI for this discussion.*

I would argue that a robot is an AI-system in physical form. Many of the definitions describe a system that is able to perform tasks, and is somewhat “independent” when following the programmed pattern or task. Here I'm thinking “independent” in the way that it may not need a human to turn it on and off, or that the system could make decisions that would benefit the predetermined tasks. Both a robot and an AI are according to the definitions supposed to be able to perform tasks, and a robot is a physical object while an AI doesn't need to be.

Based on the definitions, more of the AI-definitions focus on intelligence and human-like intelligence, so it could seem that one doesn't have the same intelligence-standards for a robot such as for an AI-system.

## A contemporary physical robot:

*Either described in a research article - or a commercial robot, and describe how this robot moves and how a human user is interacting and using the robot in a specific situation.*

The iRobot Roomba is a robot vacuum cleaner that is supposed to do some of the cleaning jobs so that users don't need to do them. Users don't have to be present when the robot does its tasks, and you can program the robot to clean when it suits you the best. It moves around the room and changes direction every time it encounters a wall, a floor molding or an edge (like a staircase going downward). Once it is done it finds its way back to a charging station or if it's cleaning in a separate room from the cleaning station the user has to physically move it so it can manage to enter the charging station.

## 3 Universal Design and AI systems

### A definition of Universal Design:

#### **Definition from National Disability Authority:**

“Universal Design is the design and composition of an environment so that it can be accessed, understood and used to the greatest extent possible by all people regardless of their age, size, ability or disability.”

Universal design is in my opinion design that opens possibilities for everyone to use it. That there are always options to access it no matter what types of abilities you have as a user. The definition from NDA describes exactly this, that something should be designed so that it's accessible, understandable and can be used by everyone no matter their differences.

### The potential of AI:

*Describe the potential of AI with respect to human perception, human movement and human cognition/emotions.*

Systems using AI is as I've mentioned before often described as kind of independent, so this could maybe help people who were earlier dependent on the help of other people to be more independent in their everyday life. An example of this could possibly be a smart home where the user could use voice commands to control the system.

*Describe the potential of AI for including and excluding people.*

AI has great potential when it comes to including humans. For example when it comes to voice recognition or voice-over-systems. This helps people with vision related disabilities or maybe people have trouble using their hands to press buttons. It could also be a possibility that AI could exclude people if the systems weren't designed with universal design in mind. It could for example be made too complicated for elders to understand.

**Do machines understand?:**

The word understanding is in my opinion the way people perceive information and how they act on it. So in a way machines could do the exact same thing, they perceive the information they get and act on it according to the programmed rules of the system. I would say the biggest difference between a humans' understanding and an AIs understanding, is their ability to see the nuances in the information and interpret the information according to body language or tone.

## **4 Guideline for Human-AI interaction**

**Microsoft's 18 guidelines:**

*Please select one of the 18 guidelines from Microsoft, and describe this guideline with a different example than what is given by Microsoft.*

### **Microsoft guideline number 9: Support efficient correction**

The AI system needs to let the user correct the answers and decisions made by the AI, without disrupting the flow of the interaction. For example, the user needs to be able to let the AI know that it disagrees with a decision the AI makes.

**HCI guidelines:**

*Search, and find one set of HCI design guidelines. Discuss briefly similarities and differences between the HCI design guidelines and the Human-AI interaction guidelines.*

I chose to look at Norman's Seven Principles of Usability. The ability to see what is happening (visibility) and to get feedback is something that seems important in both sets of guidelines. I think feedback is the most relevant similarity between Norman's HCI principles and Microsoft's Human-AI interaction guidelines.

## 5 References

Oxford Learner's Dictionary, robot:

<https://www.oxfordlearnersdictionaries.com/definition/english/robot?q=robot> (Retrieved: 8/sept-2020)

Oxford Learner's Dictionary, artificial intelligence:

<https://www.oxfordlearnersdictionaries.com/definition/english/artificial-intelligence?q=artificial+intelligence> (Retrieved: 8/sept-2020)

Grudin, Jonathan. 2009. "AI and HCI: Two Fields Divided by a Common Focus". AI magazine 30, no 4

<https://aaai.org/ojs/index.php/aimagazine/article/view/2271> (7/sept-2020)

Verne, G, Bratteteig, 2018, Does AI make PD obsolete?; exploring challenges from Artificial Intelligence to Participatory design

<https://dl.acm.org/doi/10.1145/3210604.3210646> (Retrieved: 3/sept-2020)

McCarthy, John. 2007. "What is Artificial Intelligence?". Computer Science Department.

<http://jmc.stanford.edu/articles/whatisai/whatisai.pdf> (Retrieved: 9/sept-2020)

Schulz, Trenton. 2020. "Exploration of Moving Things in the Home.

<https://www.duo.uio.no/handle/10852/74061> (Retrieved: 8/sept-2020)

National Disability Authority. Centre for Excellence in Universal Design:

<http://universaldesign.ie/What-is-Universal-Design/> (Retrieved: 10/sept-2020)

Information about Noreas' use of an AI-system:

<https://www.nordea.com/no/presse-og-nyheter/nyheter-og-pressemeldinger/news-group/2017/ai-partnership-nordea.html> (Retrieved: 9/sept-2020)



## Module 2 – Second iteration

### 1 Characteristics of AI-infused systems

Key characteristics of AI-infused systems:

*AI-infused systems are 'systems that have features harnessing AI capabilities that are directly exposed to the end user' (Amershi et al., 2019). Identify and describe key characteristics of AI-infused systems.*

AI-infused systems are described as systems that may demonstrate unpredictable behaviour that can be disruptive, confusing, offensive, or even dangerous (Amershi et al., 2019). More easily they can be described as machine learning systems or just intelligent and smart systems. Examples of an AI-infused system could be an autocompletions system, or a search engine. A system that learns over time and reacts and behaves to tasks.

An AI-infused system focuses on the concepts of learning, improving, black box and that it's fuelled by large data sets. The AI-infused system needs the ability to learn over time and be able to improve in a way that helps the system work better. It might also be difficult to understand the system, as the black box characteristic says. Therefore, it's important for the system to open in the way that it shows why and what it's doing. Lastly, an AI-infused system is fuelled by large data sets, and the system is dependent on these data sets to work sufficiently. If the system is designed correctly, to capture data as it's used, this may help improve and benefit the system.

An example of an AI-infused system:

*Identify one AI-infused system which you know well, that exemplifies some of the above key characteristics. Discuss the implications of these characteristics for the example system, in particular how users are affected by these characteristics.*

A system I use quite regularly is Pinterest. This is an online community where you can share photos and videos, and put the content you like into folders called boards. The content presented to you on the main page is based on images you have pinned (saved) earlier and the other account that you follow. I have a Pinterest board dedicated to pictures of smart home organizing ideas, and my home page is therefore regularly filled up with pictures of home organizing equipment and ideas. If I for example spend some time pinning pictures of

Mariann Gundegjerde

halloween costumes, most of the home organizing pictures would disappear, because it seems like the homepage is filled with content most related to your recent activity on the app. Still, there will occasionally be home-organizing pictures presented since I regularly look this up and pin pictures in this category. The home page is mostly full of pictures related to your recent searches and pinnes, but with the occasional picture that relates to the board you have created.

It is obvious that the app learns and gets data from my activity, and presents new content based on this. While the app usually presents relevant content, there have been instances where some content doesn't fit my interests and the reason for it appearing is confusing. Other than continuing to pin content you like, you can also choose to hide content if you'd like. I can imagine this will help the system know if there is some content that you don't like, and therefore want to see less of. The same can be said if you are searching for content, instead of just scrolling through the home page content. Here you also have the ability to hide content, and again I can imagine this could help the algorithm decide what should be presented for the different searches.

## 2 Human-AI interaction design

Main takeaways from Amershi et al. (2019) and Kocielnik et al. (2019):

Amershi et al. (2019) presents 18 guidelines for designing Human-AI interaction. These guidelines have been tested and validated through four phases they describe thoroughly in the article. The goal of the guidelines are to give researchers a framework that can hopefully better their results and will facilitate future research into the refinement and development of principles for human-AI interaction. The 18 guidelines are separated into four categories, "Initially", "During interaction", "When wrong", and "Over time". The categories describe when they likely are to be applied during interaction with users.

Kocielnik et al. (2019) article also presents some strategies and guidelines for designing AI-infused systems. Though, Kocielnik's article "only" presents three techniques that are supposed to preserve user satisfaction and acceptance of an imperfect AI system. The AI-system they test in the article is an AI-powered Scheduling Assistant they implemented themselves.

## Design guidelines:

*Select two of the design guidelines in Amershi et al. (2019). Discuss how the AI-infused system you used as example in the previous task adheres to, or deviates from these two design guidelines.*

### **G4 - Show contextually relevant information** (Amershi et al., 2019)

This guideline describes how the system should display information relevant to the user's current task and environment. For my example of Pinterest, this principle seems to be what the system is built upon. Pinterest is filled with contextually relevant information, but as I explained in the previous task, there are errors some times and the system needs to have a clear way of letting the user decide and tell when information does not feel relevant. Now they have a button called "hide", but it might not be clear for all users what this button does and that it might be useful for them (and the system) if they used this button more frequently.

### **G7 - Support efficient invocation** (Amershi et al., 2019)

This guideline describes how the system should make it easy to invoke or request the AI system's services when needed. I found this guideline quite interesting when comparing it to Pinterest. When searching on the app, the system presents suggestions to what you might add to your search. For example if you type "Halloween", they suggest "Halloween decorations" and "Halloween costumes". In addition to this, related searches are always present when you enter a search category. If you search for "Organization" and enter this search, you will be presented with a lot of different organizing pictures. On the top of the app, the system presents related searches or words/phrases you can add on to your search to narrow it down. For example, the add-ons of organization could be "Kitchen", "Bathroom", or "Small closet". This gives the user a clear overview of what the system can do, and gives the user easy access to the services by showing them without the user needing to come up with search words themselves.

### 3 Chatbots / conversational user interfaces

Key challenges in the design of chatbots or conversational agents:

*Discuss key challenges in the design of chatbots / conversational user interfaces.*

AI-systems in general are described as challenging to design. In the Yang et al. (2020) article they describe how HCI-designers struggle to envision and prototype AI systems. While this article focuses on mainly human-AI design, many of the struggles described in the article can be connected to the design of chatbots. Chatbots are machine agents serving as natural-language user interfaces to data and service providers, typically in the context of messaging applications (Følstad & Brandtzæg, 2017).

To locate the key challenges for designing chatbots or conversational user interfaces, I used both the Følstad and Brandtzæg (2017) article and the Yang et al. (2020) article. Here are my findings:

When designing a chatbot, you're designing a conversation and this can be a challenge for a designer usually more used to designing visual layouts and interaction mechanisms. The designers challenge is to focus not on the explanatory tasks of the system, the task which explains the availability and goal of the system, but rather focus on the interpretational task, understanding what the user needs and requirements are (Følstad & Brandtzæg, 2017).

The key challenge at hand is understanding the chatbots capabilities and how to envision and craft thoughtful interactions (Yang et al., 2020). If a chatbot uses machine learning the interactions can change over time, and the designer needs to understand how a user would interact with the chatbot and what the expected answers could be and design thereafter. One need to understand the context and expected use cases to hopefully design a chatbot that could accomplish what it's set out to do. While it might be difficult to articulate what AI, or a chatbot, can and cannot do, it could be even more challenging to prototype and test early stage design ideas.

## Resolving challenges:

*Revisit Guidelines G1 and G2 in Amershi et al. (2019). Discuss how adherence to these could possibly resolve some of the challenges in current chatbots / conversational user interfaces.*

Guideline 1 and 2 from Amershi et al. (2019), describe how the system should make it clear to the user what the system can do and how well it can do it. If the challenge is as described above, that it might be difficult to design a chatbot because it is difficult to understand the scope of the design - then these two guidelines may help both the user and the designer.

If the goal of the chatbot-design, is for the chatbot to clearly state what it can and cannot do. Then before the chatbot is implemented, the designer needs to research what they want the system to do and not do. If the process before the actual implementation is good, and gives the designer the knowledge they need to communicate clearly to the users through the design, then hopefully this will have somewhat resolved the issue.

Liao et al. (2020) describes some goals and themes of motivation as to why explanations in AI are important for the user. These points may strengthen the guidelines and why they might help resolve the challenges in current chatbot-like systems. One of these goals describes how explainability in an AI-system may help enhance decision confidence from the users or it can help the designers to appropriately evaluate the capabilities of the system. The explainability Liao et al. (2020) describes is closely related, or could be seen as the same as the clarity that is described in Amershi et al. (2019)'s guideline 1 and 2. If the systems' explainability and clarity toward the user is good and understandable, then there is a high chance that many of the challenges could be resolved and dealt with, both during the design process and the implementation process.

## 4 References

Amershi, S., Weld, D., Vorvoreanu, M., Fournery, A., Nushi, B., Collisson, P., ... & Teevan, J. (2019). “*Guidelines for human-AI interaction*”. In Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (paper no. 3). ACM.  
(<https://www.microsoft.com/en-us/research/uploads/prod/2019/01/Guidelines-for-Human-AI-Interaction-camera-ready.pdf>)

Følstad, A., & Brandtzæg, P. B. (2017). “*Chatbots and the new world of HCP*”. interactions, 24(4), 38-42.  
(<https://dl.acm.org/citation.cfm?id=3085558>)

Kocielnik, R., Amershi, S., & Bennett, P. N. (2019). “*Will You Accept an Imperfect AI?: Exploring Designs for Adjusting End-user Expectations of AI Systems*”. In Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (paper no. 411). ACM.  
([https://www.microsoft.com/en-us/research/uploads/prod/2019/01/chi19\\_kocielnik\\_et\\_al.pdf](https://www.microsoft.com/en-us/research/uploads/prod/2019/01/chi19_kocielnik_et_al.pdf))

Liao, Q. V., Gruen, D., & Miller, S. (2020, April). “*Questioning the AI: Informing Design Practices for Explainable AI User Experiences*”. In Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (paper no. 463). ACM.  
(<https://dl.acm.org/doi/abs/10.1145/3313831.3376590>)

Yang, Q., Steinfeld, A., Rosé, C., & Zimmerman, J. (2020, April). “*Re-examining Whether, Why, and How Human-AI Interaction Is Uniquely Difficult to Design*”. In Proceedings of the 2020 chi conference on human factors in computing systems (Paper no. 164).  
(<https://dl.acm.org/doi/abs/10.1145/3313831.3376301>)

## Appendix 1 – Feedback from iteration 1

For my feedback on iteration 1 I was told to try to use more references, so that has been my main goal for iteration 2. I got good feedback on my own reflections, I only lacked in referencing the curriculum and other articles. Hopefully I managed to connect my own reflections to relevant articles this time.

## Module 3 – Third iteration

### 1 Human AI collaboration

*Phillips et al. (2016) give a taxonomy and examples of human-robots collaboration. Choose 2- 3 examples, describe their levels of autonomy as described in Shneiderman (2020) and reflect on advantages and disadvantages if we decrease/increase their current level of autonomy. Reflect on their current and needed explainability (Hagras, 2018; Smith-Renner et al. 2020).*

#### Taxonomy and examples of human-robots collaboration

Phillips et al. (2016) describe a way of designing robotic technologies to not only be used as tools, but to serve more as a teammate. They use the analogy of animal-teammate to give insight into how these technologies can and should be designed. We can use the relationship within a team as a design guidance, with for example looking at forms of communication and features that engender trust. This guidance and knowledge can help create characteristics for robotic technologies that influence a person's expectations and formation of mental models that sparks trust and better understanding. To describe the comparisons between human-animal teams and similar robotic technologies, Phillips et al. (2016) list some examples and I will describe two of them to further use in this assignment.

The first example I've chosen is the NAO robot. The article compares these robots to human-animal teams where the animal has the purpose of being emotional support. The NAO robot was used to encourage children in care-taking behaviour and is compared to animals used to teach social skills and help people who are dealing with anxiety-inducing or abusive situations.

The second example I've chosen is the Boston Dynamics Big Dog robot. The Big Dog is used to replace physical capabilities in the same way humans would use pack mules or horses to carry cargo or other humans. Big Dog is a military robot that is built to carry cargo in order to reduce soldier load.

#### Levels of autonomy

Shneiderman (2020) describes a two dimensional HCAI-system for explaining the level of autonomy and human control in a technology. This system is an extension of Sheridan and

Verplank's (1978) ten levels of autonomy which Shneiderman describes as one dimensional. Shneiderman's levels of autonomy is described as a square where there are either human-control or computer-control and high and low levels of automation. The goal is to create systems that are trusted, reliable and safe (TRS) and one can achieve this goal by creating a system that has a high level of human control and a high level of automation (Shneiderman, 2020, p.7).

The NAO robot has a high level in human control, due to the way it's made. NAO is a programmable robot, which I interpret as it can be made to do certain things and react to predefined parameters. It seems that the robot has a somewhat high level of automation as well. After being programmed to be and act a certain way, I would assume that the robot then acted on its own without much need of human influence.

If we were to decrease NAO's level of autonomy I would see this as a disadvantage. If the NAO robot is to work most effectively it should be able to do so with as little interference as possible, and if one for example had to press a lot of buttons to get NAO to act it might not feel as natural and effective to use it. And as described earlier, it was used to encourage children in having care-taking behaviour, and the reason for this maybe being an effective robot for this type of work could be its high level of automation. If one were to include machine learning in a robot like this, and therefore decrease its level of human-control, it could be both an advantage and a disadvantage. The risks could be higher, because of a chance of errors, though maybe not life threatening errors. The advantages could be that the robot could be even more effective, if it were able to adapt its behaviour to the users.

Similar to the NAO robot, the Big Dog has a high level of human control. Big Dog is usually controlled by a human operator that controls its behaviour and how it should traverse the terrain (Raibert et al., 2008, p.10823). This also shows that Big Dog has a somewhat low level of automation, because it is dependent on a human operator to operate. It would be an advantage if one could make the Big Dog more independent with a higher level of automation. Then the robot could move around the terrain on its own without needing the help of a human to make sure it doesn't fall.



## Current and needed explainability

The NAO robot and the Big Dog would need different levels of explainability due to the different types of users that are interacting with them. Hagraas (2018) discusses the need for user feedback to the system and if explainability would lead to frustration or clarity. It seems from the article that the users wished to give the system feedback after getting explanations from the system when encountering flaws. Giving explanations without means for feedback may reduce the users satisfaction with the system (Hagraas, 2018, p.10). This would probably align with the users of the Big Dog, where it would be very beneficial to be able to give the system feedback if the system described and explained an error that occurred. Though here it would be important that the explanations from the Big Dog would be accurate, since it is a military robot it might be critical to be able to locate and fix the errors at hand quickly. The NAO robot, which in the example described by Phillips et al. (2016), were assigned to work with kids, wouldn't need the same ability to receive feedback. Kids might not give reliable feedback that would be beneficial for the system. It would still be important for the system to show a high level of explainability in that it explained clearly if something was wrong and therefore the child would understand to ask assistance from an adult. If I were to assume that the NAO robot acted on voice commands, it would be important that the system could explain if it didn't understand the command and wanted the user to rephrase or talk more clearly.

Smith-Renner et al. (2020) describes some areas of XAI (Explainable AI) that are important to focus on. Transparency and safety is two of these areas. With the Big Dog this is an important area if the robot is to work with a high level of autonomy in the future. The system needs to be trusted in the choices it makes and show transparency in why those choices were made. Fairness and bias are two other areas mentioned by Smith-Renner et al. If the NAO robot were to reach an even higher level of autonomy with the help of machine learning, it would need to be trusted that the decisions it made were fair and reasonable, and not corrupted by bias from earlier knowledge and decisions.

## 2 References

E. K. Phillips, K. Schaefer, D. R. Billings, F. Jentsch, and P. A. Hancock, “*Human-Animal Teams as an Analog for Future Human-Robot Teams: Influencing Design and Fostering Trust*,” *J. Human-Robot Interact.*, vol. 5, no. 1, p. 100-125, Sep. 2015

DOI: [10.5898/JHRI.5.1.Phillips](https://doi.org/10.5898/JHRI.5.1.Phillips)

Hagras, H., “*Toward Human-Understandable, Explainable AI*”, *Computer*, 51, 9, 2018, 28-36 <https://ieeexplore.ieee.org/document/8481251>

Shneiderman, B., “*Human-Centered Artificial Intelligence: Reliable, Safe & Trustworthy*”, arXiv.org (February 23, 2020). <https://arxiv.org/abs/2002.04087v1> (Extract from forthcoming book by the same title)

Smith-Renner, A., Fan, R., Birchfield, M., Wu, T., Boyd-Graber, J., Weld, D.S., and Findlater, L. 2020. “*No Explainability without Accountability: An Empirical Study of Explanations and Feedback in Interactive ML*”. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–13. DOI: <https://doi.org/10.1145/3313831.3376624>

Raibert, M., Blankespoor, K., Nelson, G., Playter, R., the BigDog Team. “*BigDog, the Rough, Terrain Quadruped Robot*”. 2008. *IFAC Proceedings Volumes*, Volume 41, Issue 2, 2008, Pages 10822-10825, ISSN 1474-6670, ISBN 9783902661005, DOI: <https://doi.org/10.3182/20080706-5-KR-1001.01833>

## Appendix 2 – Feedback from iteration 2

I got good feedback from iteration 2, and didn't really get that many comments to things I should change. I got some feedback around my use of references and that I maybe didn't use them correctly. So I will look into that for this last iteration, but I am not sure what I should change. I will also continue to take advantage of the feedback I got on the first assignment, to use the syllabus articles carefully throughout my assignment and tie this up to my own reflections.