



Universiteit  
Leiden  
The Netherlands

## **EU privacy and data protection law applied to AI: unveiling the legal problems for individuals**

Häuselmann, A.N.

### **Citation**

Häuselmann, A. N. (2024, April 23). *EU privacy and data protection law applied to AI: unveiling the legal problems for individuals*. Retrieved from <https://hdl.handle.net/1887/3747996>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/3747996>

**Note:** To cite this publication please use the final published version (if applicable).

## 2 Artificial Intelligence (AI)

This chapter aims to answer Subquestion 1, namely, what AI is and what AI disciplines exist.<sup>50</sup> It starts with existing definitions of AI (Section 2.1) and then provides an overview of the AI disciplines that seem to be the most problematic ones from a privacy and data protection perspective (Section 2.2). These disciplines include machine learning (Section 2.2.1), natural language processing (Section 2.2.2), computer vision (Section 2.2.3), affective computing (Section 2.2.4) and automated reasoning (Section 2.2.5). Section 2.3 answers Subquestion 1.

### 2.1 Definitions of AI

There is no officially agreed definition of Artificial Intelligence (AI). AI covers a wide range of concepts and terms, making it difficult to define. Available definitions often involve ambiguous terms such as ‘thinking’, ‘learning’ and ‘intelligence’. In 1968, Minsky defined AI as ‘the science of making machine do things that would require intelligence if done by men’.<sup>51</sup> Bellman defined AI in 1978 as ‘the automation of activities that we associate with human thinking, activities such as decision-making, problem solving, learning, creating, game playing, and so on’.<sup>52</sup> Nilsson described AI as ‘activity devoted to making machines intelligent, and intelligence is that quality that enables an entity to function appropriately and with foresight in its environment’.<sup>53</sup> Russell and Norvig organised definitions of AI into four categories: a) thinking humanly, b) acting humanly, c) thinking rationally and d) acting rationally.<sup>54</sup> According to Munakata, AI involves abilities such as ‘inference based on knowledge, reasoning with uncertain or incomplete information, various forms of perception and learning, and applications to problems such as control, prediction, classification, and optimization’.<sup>55</sup> More recent definitions are the ones adopted by the Organisation for Economic Co-operation and Development (OECD) and the National Institute of Standards and Technology (NIST). The OECD defines an AI system as a ‘machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments. Different AI systems vary in their levels of autonomy and adaptiveness after deployment’.<sup>56</sup> NIST defines AI as a ‘branch of computer science devoted to developing data processing systems that performs functions normally associated with human intelligence, such as reasoning, learning, and self-improvement’.<sup>57</sup>

<sup>50</sup> A modified version of this chapter was published in Bart Custers, Eduard-Fosch Villaronga (eds) *Law and Artificial Intelligence* (Asser Press 2022). See Andreas Häuselmann, ‘Disciplines of AI: An Overview of Approaches and Techniques’ 43-70.

<sup>51</sup> Marvin Minsky, *Semantic Information Processing* (MIT Press 1968).

<sup>52</sup> Richard Bellman, *An Introduction to Artificial Intelligence: Can computers think?* (Boyd & Faser 1978) 3

<sup>53</sup> Nils J Nilsson, *The Quest for Artificial Intelligence: A History of Ideas and Achievements* (Cambridge University Press 2010).

<sup>54</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 2.

<sup>55</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) xx.

<sup>56</sup> See < <https://oecd.ai/en/wonk/ai-system-definition-update> > accessed 8 February 2024.

<sup>57</sup> See < <https://csrc.nist.gov/topics/technologies/artificial-intelligence> > accessed 8 February 2024.

The field of AI may be divided into narrow and general AI. Narrow AI refers to systems that are able to solve a specific problem or performing a specific task. For an example on a narrow AI system, one can refer to IBM's 'Deep Blue' chess-playing computer. Deep Blue defeated the reigning world champion in chess, Garry Kasparov, in 1997.<sup>58</sup> This example indicates that computers can perform better than humans. However, this holds only true for a narrow domain, such as playing chess. General AI aims to build machines that generally perform on a human level and have a 'human-level' skillset. To achieve this goal, such a system must be able to mimic the functioning of the human brain in the most important aspects.<sup>59</sup> Unlike with narrow AI, general AI arguably has not been achieved yet despite rapid developments, for instance, ChatGPT. Although AI found its 'birth' at the Dartmouth Summer Research Project on AI in the summer of 1956 in New Hampshire,<sup>60</sup> there are many open challenges. According to Shi, AI research is still in its first stage since no breakthrough progress has been achieved for some key challenges such as common sense knowledge representation and uncertain reasoning.<sup>61</sup> Therefore, current AI systems must be considered examples of 'narrow' AI. However, computing power has become more affordable; the computers have become faster and contain larger memories. This led to the 'summer of AI' and it seems reasonable to expect major developments in the field of AI.

In his famous paper, called 'Computing Machinery and Intelligence',<sup>62</sup> Turing proposed the 'Imitation Game', which has later become known as the 'Turing test'.<sup>63</sup> Turing offered his test as a sufficient condition for the existence of AI.<sup>64</sup> This test involves three actors: (A) a machine, (B) a human and (C) another human called the interrogator (see Figure 1.1). In the Turing test, the human interrogator (C) stays in a room apart from the other two actors (A) and (B). The human interrogator knows the machine (A) and human (B) by labels (X) and (Y)<sup>65</sup> and therefore does not know which label is (A) or (B).<sup>66</sup> The object of the test is for the interrogator (C) to determine which of the other two actors is the human and which is the machine<sup>67</sup> by asking (X) and (Y) questions which they must answer.<sup>68</sup> In other words, the human interrogator engages in conversation with either a human or an AI natural language program which are both hidden from view. If the human interrogator cannot reliably

<sup>58</sup> <https://www.livescience.com/59065-deep-blue-garry-kasparov-chess-match-anniversary.html>, accessed 8 February 2024.

<sup>59</sup> Kevin Warwick, *Artificial Intelligence: The basics* (Routledge 2012) 65.

<sup>60</sup> Ronald R Kline, 'Cybernetics, Automata Studies, and the Dartmouth Conference on Artificial Intelligence' (2011) 4, EEE Computer Society, 5.

<sup>61</sup> Zhongzhi Shi, *Advanced Artificial Intelligence* (World Scientific 2011) 18.

<sup>62</sup> Alan Mathison Turing, 'Computing Machinery and Intelligence' (1950) Vol LIX Iss 236 Mind 433-460.

<sup>63</sup> Chris Bernhardt, *Turing's Vision: The Birth of Computer Science* (MIT Press 2016) 157.

<sup>64</sup> Stan Franklin, 'History, motivations, and core themes' in Frankish Keith and Ramsey William (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 17.

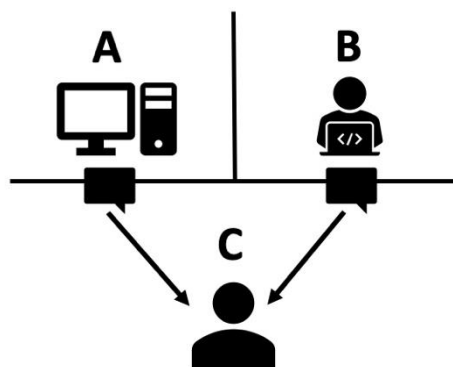
<sup>65</sup> Alan Mathison Turing, 'Computing Machinery and Intelligence' (1950) Vol LIX Iss 236 Mind 433-460.

<sup>66</sup> Chris Bernhardt, *Turing's Vision: The Birth of Computer Science* (MIT Press 2016) 157.

<sup>67</sup> Alan Mathison Turing, 'Computing Machinery and Intelligence' (1950) Vol LIX Iss 236 Mind 433-460.

<sup>68</sup> Chris Bernhardt, *Turing's Vision: The Birth of Computer Science* (MIT Press 2016) 157.

distinguish between the human and the program/machine, (artificial) intelligence is ascribed to the program.<sup>69</sup>



**Figure 1.1** Illustration of the Turing test created by the author.

There are plenty of definitions for AI, which involve ambiguous terms such as those already mentioned. In this thesis, AI refers to adaptive machines that can autonomously execute activities and tasks that require capabilities usually associated with humans. ‘Autonomously’ in this sense means that the machine has the ability to make its *own* decisions and perform tasks on the designer’s behalf.<sup>70</sup> ‘Adaptive’ refers to the machine’s ability to learn from, and adapt to its environment in order to preserve its autonomy in dynamic environments.<sup>71</sup> Adaptivity is very important, since only a machine that *learns* will succeed in a vast variety of environments.<sup>72</sup> Learning in this context corresponds to ‘adapt’ the performance according to previously made experiences based on statistics and probability calculations.<sup>73</sup> This definition aligns well with the ones adopted by the OECD and NIST.<sup>74</sup>

## 2.2 AI disciplines

Since AI covers a broad range of concepts, this research will pay particular attention to AI disciplines which could be problematic in the light of the fundamental rights to privacy and data protection. These AI disciplines are coloured blue in Figure 1.2.<sup>75</sup> The remaining disciplines (white) will not be discussed in this thesis.

<sup>69</sup> Stan Franklin, ‘History, motivations, and core themes’ in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 17, 18.

<sup>70</sup> Eduardo Alonso, ‘Actions and agents’ in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 235, 236.

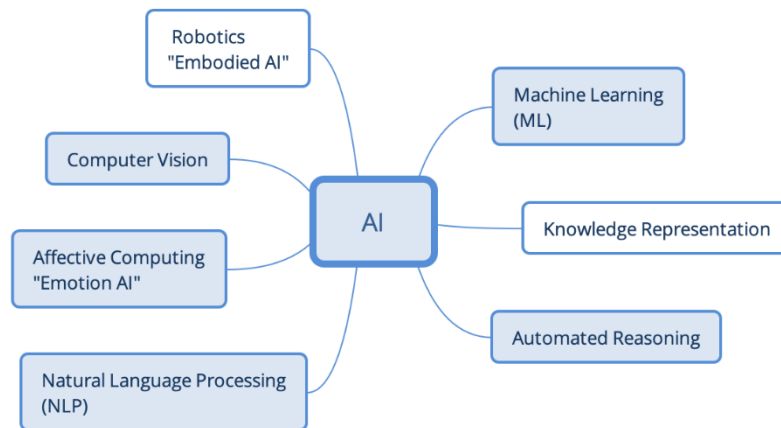
<sup>71</sup> *Ibid* 235.

<sup>72</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 39.

<sup>73</sup> Stefan Strauß, ‘From Big Data to Deep Learning: A Leap Towards Strong AI or Intelligentia Obscura’ (2018) 2 (3), *Big Data and Cognitive Computing* <<https://www.mdpi.com/2504-2289/2/3/16>> accessed 14 January 2019, 7.

<sup>74</sup> See <<https://oecd.ai/en/wonk/ai-system-definition-update>> and <<https://csrc.nist.gov/topics/technologies/artificial-intelligence>> respectively, accessed 8 February 2024.

<sup>75</sup> This figure shall not be considered as a complete overview of all AI disciplines, but serves as an illustrative overview for this thesis.



**Figure 1.2** Graph created by the author outlining the AI disciplines inspired by Russel/Norvig<sup>76</sup> and slightly adjusted by adding the field of affective computing.

Methods of AI that combine AI with Robotics, i.e. ‘Embodied Artificial Intelligence’, are out of scope of this thesis. Applications of Embodied AI such as driverless vehicles, surgical robots and companions pose different questions such as liability issues or ethical issues in the context of robot-human interactions.<sup>77</sup> However, these questions are not in the scope of this research.

AI systems need to translate input into information or knowledge so that it can be processed to select output (action).<sup>78</sup> The discipline of AI research commonly referred to as knowledge representation focusses on the computers capabilities to store what it knows and hears.<sup>79</sup> Since research in this discipline of AI focusses on conceptual issues<sup>80</sup> not related to privacy and data protection, it will not be discussed here. However, the subfield of automated reasoning, which is a fundamental part of knowledge representation, will be discussed due to its implications on automated decision-making.

### 2.2.1 Machine learning (ML)

ML may be considered a discipline or one of the tools of AI.<sup>81</sup> I follow the former approach in this thesis and acknowledge that ML is often combined with other AI disciplines. Computer science has traditionally aimed to manually program computers. ML however aims to have computers program themselves based on experience.<sup>82</sup> In other words, the goal of ML is to adapt to new circumstances and to detect and extrapolate patterns.<sup>83</sup> Murphy defines ML as ‘a set of methods that can

<sup>76</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 2, 3.

<sup>77</sup> Cándido García Molyneux, Rosa Oyarzabal, ‘What Is a Robot (Under EU Law)?’ (2018) Vol 1 RAIL: The Journal of Robotics, AI & Law 11, 12.

<sup>78</sup> Stan Franklin, ‘History, motivations, and core themes’ in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 24.

<sup>79</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 2.

<sup>80</sup> E.g. the issue of whether or not to represent knowledge, Franklin Stan, ‘History, motivations, and core themes’ in Frankish Keith and Ramsey William M (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 24, 25.

<sup>81</sup> Vijay Kotu, Bala Deshpande, *Data Science* (2<sup>nd</sup> edn Elsevier 2019) 2.

<sup>82</sup> Tom M. Mitchell, ‘The discipline of Machine Learning’ (2006) 1 <<http://www.cs.cmu.edu/~tom/pubs/MachineLearning.pdf>> accessed 8 February 2024.

<sup>83</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 2.

automatically detect patterns in data, and then use the uncovered patterns to predict future data or to perform other kinds of decision making under uncertainty'.<sup>84</sup> ML can simply be described as the set of computational methods that use experience to improve its performance or to make accurate predictions.<sup>85</sup> This is achieved by using ML algorithms, algorithms that learn from experience.<sup>86</sup> Put simply, an algorithm is typically a numerical process that consists of a sequence of well-defined steps leading to the solution of a particular type of problem.<sup>87</sup> Experience refers to the data from the past available to the algorithm for analysis.<sup>88</sup> Learning in this context is about making computers modify or adapt their performance (actions) so that these actions become more *accurate*.<sup>89</sup> ML uses data-driven methods, combining fundamental concepts in computer science with approaches from statistics, probability and optimisation.<sup>90</sup> In fact, the probabilistic approach in ML is closely related to the field of statistics, but differs slightly in terms of its emphasis and terminology. The probabilistic approach is particularly helpful for handling ambiguous cases.<sup>91</sup> The main goal of ML is to generate accurate predictions for unseen data and to design efficient algorithms to produce these predictions.<sup>92</sup>

Before the specific *kind* of ML called deep learning (DL) will be discussed in Section 2.2.1.4, some of the most widely used ML *methods* will be elaborated on first in Sections 2.2.1.1 and 2.2.1.3. These methods are called supervised, unsupervised and reinforcement learning. In practice, the distinction between supervised and unsupervised learning is not always clear-cut. Therefore, semi-supervised learning creates a continuum between supervised and unsupervised learning: The algorithm is provided with a few labelled examples (supervised learning) but also has the task to uncover hidden patterns and structures in the data (unsupervised learning).<sup>93</sup> Another method deployed in ML is reinforcement learning (RL). RL is becoming increasingly relevant, in particular in natural language processing, a discipline of AI which aims to enable computers to process human language (see Section 2.2.2).

### 2.2.1.1 Supervised machine learning

Supervised ML aims to learn a mapping from input  $x$  to output  $y$ , given a labelled set of input-output pairs called the *training set* or training data. It can be used to make predictions on *new* input through generalisation.<sup>94</sup> Generalisation refers to the ability of the algorithm to categorise new examples that

<sup>84</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 1.

<sup>85</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 1.

<sup>86</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 97 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>87</sup> Yadolah Dodge, 'Algorithm' in: *The Concise Encyclopedia of Statistics* (Springer New York 2006) 1-2.

<sup>88</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 1.

<sup>89</sup> Steven Marsland, *Machine Learning: An Algorithmic Perspective* (2<sup>nd</sup> edn Chapman & Hall 2015) ch 1.2.1.

<sup>90</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 1.

<sup>91</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 1, 4.

<sup>92</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 2.

<sup>93</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 695.

<sup>94</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 3.

differ from the ones used during the training phase.<sup>95</sup> In the supervised ML approach, the learning algorithm receives several examples, each *labelled* with the correct label (training data). Consider, for example, several labelled pictures with different animals (lions, horses, and cows). The goal is that the algorithm automatically recognises the correct label for the training data and *predicts* the value of unseen (unlabelled) inputs.<sup>96</sup> In other words, the aim is that the algorithm *generalises* accurately by producing a model that can classify input *not seen* during training.<sup>97</sup> The user who provides the correct labels to the algorithm is the teacher, knowing for each input the correct output. Therefore, this is called ‘supervised’ learning: the algorithm learns under the supervision and guidance of the teacher.<sup>98</sup> To measure the accuracy of the model generated by the algorithm, the teacher provides the algorithm with a set of examples that are *different* from the set of training.<sup>99</sup> Hence, the teacher feeds the algorithm with new pictures containing lions, horses and cows and evaluates the accuracy of the model, namely, whether the algorithm recognised the animals correctly. The algorithm learns by adjusting the relevant parameters so that the model makes the most accurate predictions on the data.<sup>100</sup>

There are basically two techniques used for supervised machine learning: classification and regression.<sup>101</sup> As indicated by its name, classification refers to situations where the predicted attribute is categorical, and regression applies to situations where the predicted attribute is numeric.<sup>102</sup> Classification orders data into exhaustive and exclusive groups or classes on the basis of their similarity. Consequently, all data can only be assigned to one class.<sup>103</sup> The example with the animal referred to the classification technique. Regression is suitable when the prediction to be made by the algorithm should be a numerical value. Regression could be described as a statistical approach that is used to identify the relationship between variables.<sup>104</sup> Therefore, the regression technique could be used to predict the number of people likely to click on an online advertisement based on the ad content and the user’s previous surfing history. Other real-world examples using regression are predicting stock market prices given current market conditions or predicting the age of a viewer watching a given video on YouTube.<sup>105</sup>

<sup>95</sup> Christopher M Bishop, *Pattern Recognition and Machine Learning* (Springer 2006) 2.

<sup>96</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 7.

<sup>97</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 39.

<sup>98</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn Springer 2008) 38.

<sup>99</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3<sup>rd</sup> edn, Pearson Education 2016) 695.

<sup>100</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 39.

<sup>101</sup> Michele Uselli, *R machine learning essentials* (Packt Publishing 2014) 155.

<sup>102</sup> *Ibid* 154.

<sup>103</sup> Toon Calders, Bart Custers, ‘What is Data Mining and How Does it Work?’ in Bart Custers et al. (eds) *Discrimination and Privacy in the Information Society* (Springer 2013) 32.

<sup>104</sup> However, note that decision tree regression would not be considered as traditional statistics.

<sup>105</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 9.

### 2.2.1.2 Unsupervised machine learning

Unlike supervised ML, the algorithm only receives *unlabelled* training data.<sup>106</sup> That means that the algorithm is not told what the desired output is for each form of input and unsupervised ML does not require a human expert to manually label the data.<sup>107</sup> Due to the fact that there is no external comparison between actual and ideal output by the teacher, this approach is called unsupervised: There are no correct answers available.<sup>108</sup> Therefore, the algorithm tries to discover patterns in the input even though no explicit feedback is supplied.<sup>109</sup> The goal of unsupervised ML is to identify associations and patterns among a set of input data and categorise them accordingly.<sup>110</sup> It can be difficult to quantitatively evaluate the performance of the model, since there are no labelled examples available.<sup>111</sup> Two branches of techniques used for unsupervised learning are clustering and dimensionality reduction.<sup>112</sup>

Clustering in this context means dividing detected patterns into groups or clusters. Similar patterns are placed in the same group, while all others are put in different groups.<sup>113</sup> Simply put, clustering refers to the partition of unlabelled items into homogeneous regions.<sup>114</sup> Clusters may overlap, while classifications do not (see Section 2.2.1.1). Clustering is particularly performed to analyse very large data sets. A common example is to use clustering in the context of social network analysis, where the clustering algorithm tries to identify ‘communities’ within large groups of people.<sup>115</sup> The same applies to e-commerce, where users are clustered into groups based on their purchasing or online behaviour, which enables online shops to send customised targeted ads to each group.<sup>116</sup>

Dimensionality reduction aims to represent data with fewer dimensions<sup>117</sup> and is useful to project high-dimensional data to a lower dimensional subspace to capture the ‘essence’ of the data.<sup>118</sup> By reducing the dimensions, *hidden patterns* and *structures* in the data may be observed, and non-informative features are discarded. Dimensional representations often produce *better predictive accuracy* because they focus on the essence of the object and filter out non-essential features.<sup>119</sup> Dimensionality reduction is commonly used to pre-process digital images, in computer vision tasks<sup>120</sup> (see

<sup>106</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 7.

<sup>107</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 9.

<sup>108</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn Springer 2008) 38.

<sup>109</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 694.

<sup>110</sup> Hastie Trevor, Tibshirani Robert, Friedman Jerome, *The Elements of Statistical Learning* (2<sup>nd</sup> edn 2008) xi; Steven Marsland, *Machine Learning: An Algorithmic Perspective* (2<sup>nd</sup> edn Chapman & Hall 2015) ch 1.3.

<sup>111</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 7.

<sup>112</sup> Michele Usuelli, *R machine learning essentials* (Packt Publishing 2014) 164.

<sup>113</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn Springer 2008) 72.

<sup>114</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 2.

<sup>115</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 2.

<sup>116</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 11.

<sup>117</sup> Ethem Alpaydin, *Introduction to Machine Learning* (4th edn MIT Press 2020) 137, 138.

<sup>118</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 11.

<sup>119</sup> *Ibid*, 12.

<sup>120</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 2.



Section 2.3) and applied in natural language processing (see Section 2.2.2), e.g., for acoustic signals.<sup>121</sup>

### 2.2.1.3 Reinforcement learning (RL)

Reinforcement learning (RL) is a distinct method in ML that differs from supervised and unsupervised ML approaches. In RL, the algorithm interacts with its environment and the method is inspired by behavioural psychology.<sup>122</sup> RL algorithms modify or acquire new behaviours incrementally and use *trial-and-error* experience without requiring complete knowledge or control of the environment.<sup>123</sup> Unlike supervised learning, RL learns with a ‘critic’ who does not instruct the algorithm what to do, but rather provides it with feedback in the form of a reward or punishment.<sup>124</sup> The reward depends on the correctness of the decision (the action by the agent).<sup>125</sup> In RL, the decision-maker is called the agent which interacts with everything outside the agent, called the environment. The agent and environment interact continuously: the agent selects actions, and the environment responds to these actions and presents new situations to the agent.<sup>126</sup> The agent has no prior knowledge of what action to take; it learns from interaction with the environment.<sup>127</sup> The object of the agent is to maximise its reward over a course of interactions with the environment.<sup>128</sup> Therefore, the agent uses the received feedback to update its knowledge so that it learns to perform actions that return the highest reward.<sup>129</sup>

An illustrative example is a machine (agent) that learns to play chess. The chessboard is the environment of the agent that must decide over a sequence of actions, namely, ‘moves’ on the chessboard (environment) to achieve a certain goal, namely, winning the game. In RL, the agent evolves and learns while analysing the consequences of its actions with the feedback received from the environment.<sup>130</sup> This is different from the unsupervised ML approach, where no feedback is distributed. RL also differs from supervised ML because the agent does not learn from the initially labelled training data, but from the interaction with the environment based on feedback in the form of a punishment or reward.<sup>131</sup> Combining it with deep learning techniques has made ‘deep RL’ increasingly successful in addressing challenging sequential decision-making problems such as mastering the game ‘Go’<sup>132</sup> or

<sup>121</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 11.

<sup>122</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 104 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>123</sup> Vincent François-Lavet et al, ‘An Introduction to Deep Reinforcement Learning’ (2018), Vol. 11, No. 3-4 Foundations and Trends in Machine Learning, 2, 15.

<sup>124</sup> Ethem Alpaydin, *Introduction to Machine Learning* (4th edn MIT Press 2020) 570.

<sup>125</sup> Andries P Engelbrecht, *Computational Intelligence – An Introduction* (2nd edn John Wiley & Sons 2007) 83.

<sup>126</sup> Zhongzhi Shi, *Advanced Artificial Intelligence* (World Scientific 2011) 365.

<sup>127</sup> Andries P. Engelbrecht, *Computational Intelligence – An Introduction* (2 edn John Wiley & Sons 2007) 83.

<sup>128</sup> Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, *Foundations of Machine Learning* (MIT Press 2012) 8.

<sup>129</sup> Ethem Alpaydin, *Introduction to Machine Learning* (4th edn MIT Press 2020) 570.

<sup>130</sup> Zhongzhi Shi, *Advanced Artificial Intelligence* (World Scientific 2011) 362.

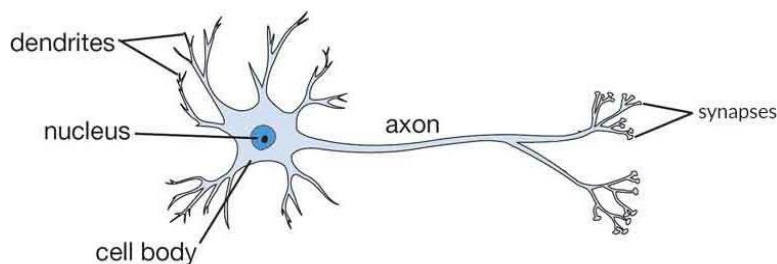
<sup>131</sup> Sumit Das et al., ‘Applications of Artificial Intelligence in Machine Learning: Review and Prospect’ (2015), Vol. 115, No. 9 International Journal of Computer Applications 31, 32.

<sup>132</sup> See <<https://www.deepmind.com/research/highlighted-research/alphago>> accessed 8 February 2024.

beating the world's top professionals in poker.<sup>133</sup> Its adaptive capabilities make RL very suitable for interactive applications. For example, deep RL is applied for dialogue systems and conversational agents, in particular for digital assistants and chatbots.<sup>134</sup> The most impressive and current example is ChatGPT provided by OpenAI. ChatGPT is a large language model trained to produce text. It was optimised by using reinforcement learning with human feedback.<sup>135</sup> Deep RL seems to possess promising potential for real-world applications such as robotics, self-driving cars, finance and smart grids.<sup>136</sup> Current ML applications based on the supervised method for natural language processing and speech recognition require vast amounts of labelled training data. This issue could be eliminated by applying deep RL methods.<sup>137</sup>

#### 2.2.1.4 Artificial Neural Networks and deep learning

The human brain consists of a very large number of processing units called neurons.<sup>138</sup> These neurons have an output fibre called an axon and a terminal fibre called a synapse. The axons split up and connect to several dendrites, which are the input pathways of other neurons through the junction terminal synapse.<sup>139</sup> Because the neurons of the human brain are connected, it is called a neural network. Figure 1.3 shows a typical biological neuron.



**Figure 1.3** Biological neuron illustrated by Navdeep Singh.<sup>140</sup> Used with permission.

Although it is not entirely clear how the neural network of human brains actually works, it is considered to be the fundamental functional source of intelligence, which includes perception, learning and cognition.<sup>141</sup> The characteristic of a neural network is that the neurons operate in parallel and transfer

<sup>133</sup> See <<https://www.nature.com/articles/d41586-019-02156-9>> accessed 8 February 2024.

<sup>134</sup> Iulian Serban et al. 'A Deep Reinforcement Learning Chatbot' (2017) 1 <<https://arxiv.org/pdf/1709.02349.pdf>> accessed 8 February 2024.

<sup>135</sup> See FAQs about ChatGPT provided by OpenAI: < <https://help.openai.com/en/articles/6783457-what-is-chatgpt> > accessed 8 February 2024.

<sup>136</sup> Vincent François-Lavet et al., 'An Introduction to Deep Reinforcement Learning' (2018) Vol. 11 No. 3-4 Foundations and Trends in Machine Learning 3.

<sup>137</sup> Deng Li and Liu Yang, 'Epilogue: Frontiers of NLP in the Deep Learning Era' in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 316.

<sup>138</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 86.

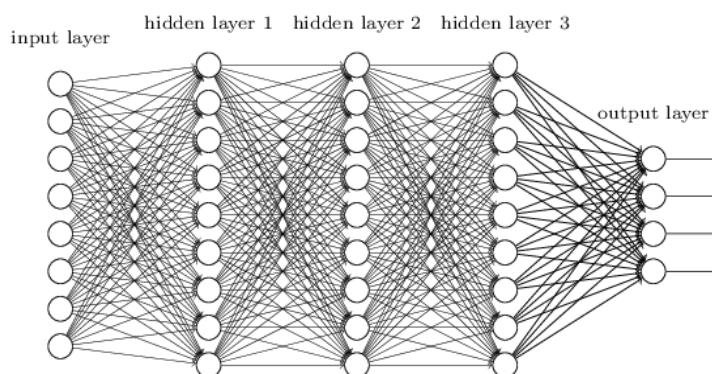
<sup>139</sup> Tommy Chow, Siu-Yeung Cho, *Neural Networks and Computing: Learning Algorithms and Applications* (Imperial College Press 2007) 2.

<sup>140</sup> Navdeep Singh Gill, 'Overview of Artificial Neural Networks and its application' <<https://www.xenonstack.com/blog/artificial-neural-network-applications/>> accessed 8 February 2024.

<sup>141</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 7.

information among themselves over the synapses so that the neurons are connected and influence each other.<sup>142</sup> The brain is believed to learn by examples and experience and to be highly capable of adapting to external changes.<sup>143</sup>

A single biological neuron would be too simple to make decisions like humans do. Similarly, a single artificial neuron would not be able to cope with challenging decision-making and prediction processes. Hence, to unleash the full potential of artificial neurons, they must operate in parallel and transfer information among themselves. That is why researchers such as Rumelhart and others in 1986 attempted to design artificial neural networks (ANN) with the aim to allow an arbitrarily connected neural network to develop an internal structure that is appropriate for a particular task.<sup>144</sup> ANNs can be simply described as an abstract model that is inspired by knowledge of the inner workings of the human brain that can be programmed on a computer. ANNs consist of artificial neurons and interconnections similar to the human brain. The network receives input, performs internal processes such as the activation of the neurons and finally yields output.<sup>145</sup> However, ANNs are generally not designed to be realistic models of the human brain. The neural perspective on deep learning is motivated by two main ideas: first, that the brain provides an example that intelligent behaviour is possible; and second, that it is possible to create machine learning models that shed light on the principles of the brain and human intelligence.<sup>146</sup> The pattern of connections between the artificial neurons is called the architecture or topology of the ANN and consists of distinct layers of neurons. The layers depend on the model used.<sup>147</sup> Each of the layers has a certain number of neurons which is usually determined by a specific application problem the model aims to solve. An example of a deep ANN is given in Figure 1.4.



**Figure 1.4** Example of a deep artificial neural network illustrated by Michael Nielsen.<sup>148</sup> Used with permission.

<sup>142</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 86.

<sup>143</sup> Tommy Chow, Siu-Yeung Cho, *Neural Networks and Computing: Learning Algorithms and Applications* (Imperial College Press 2007) 2.

<sup>144</sup> David Rumelhart, Geoffrey Hinton, Ronald Williams 'Learning representations by backpropagating errors' (1986) Vol. 323 Nature 533.

<sup>145</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 3, 7.

<sup>146</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 13 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>147</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 9.

<sup>148</sup> Michael Nielsen, 'Why are deep neural networks hard to train' in: *Neural Networks and Deep Learning* (Determination Press 2015) <<http://neuralnetworksanddeeplearning.com/chap5.html>> accessed 8 February 2024.

Generally, there are one input layer, one output layer and *any number* of hidden layers. Neurons of the input layer are connected to the neurons of the hidden layer through edges, and the neurons of the hidden layer(s) are connected to the output layer. A weight is associated to each edge. The input layer (see on the left side of Figure 1.4) consists of neurons that receive their input directly from the data, and its function is to merely send out input signals to the hidden layer neurons; it does not compute anything.<sup>149</sup> The hidden layer then applies computation methods to the inputs that depend on the model used for the neural network, transforming the received inputs to something the output layer can use. Hidden means that the values in these layers are not given in the data, but the model has the task of determining which concepts are useful for explaining the relationships in the observed data.<sup>150</sup> It then sends its output to the next layer, in the present case, to hidden layer 2, which sends it to hidden layer 3 and subsequently to the output layer (see the right side of Figure 1.4). Subsequently, the role of the output layer is to produce the output of the entire network. The output of ANNs can then be used to extract a prediction or a decision.

Deep learning (DL) is a particular kind of ML that represents the world as a nested hierarchy of concepts.<sup>151</sup> The human brain seems to execute many levels of processing with increasing levels of abstraction.<sup>152</sup> DL seems to resemble this by computing more abstract concepts in terms of less abstract ones.<sup>153</sup> Most of the models used for supervised and unsupervised ML have a simple two-layer architecture.<sup>154</sup> This is different with DL models, which use many different layers. Approaches in DL feed a large set of input data into the ANN that produces successive transformations of the input data, where each hidden layer combines the values in its preceding layer and learns more complicated functions of the input.<sup>155</sup> Then, the final transformation predicts the output.<sup>156</sup> The deep learning approach avoids the requirement that the human operator must specify all the knowledge which the computer requires. Deep learning solves this by enabling the computer to build complex concepts out of simpler concepts. When illustrating the approach in a graph by building the concepts on top of each other, that graph is deep, with many layers. Therefore, the approach is called deep learning (see Figure 1.4).<sup>157</sup> DL draws inspiration from many fields, especially from linear algebra and probabilistic statistics. Foundation models, namely models that are trained on broad data using self-supervision

<sup>149</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 10.

<sup>150</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 6 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>151</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 8 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>152</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 95.

<sup>153</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 8 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>154</sup> Kevin P Murphy, *Machine Learning: A Probabilistic Perspective* (MIT Press 2012) 995.

<sup>155</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 104.

<sup>156</sup> Yoav Goldberg, *Neural Network Methods in Natural Language Processing* (Morgan & Claypool Publishers 2017) 2.

<sup>157</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 1, 5 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

and that can be adapted to a wide range of tasks, are based on deep neural networks.<sup>158</sup> Typical examples of foundation models include large language models (LLMs) as introduced in the AI discipline natural language processing (Section 2.2.2).

Interestingly, achievements in modern DL have been made with an astonishingly small number of neurons contained in the ANNs when compared with neural networks of the human brain. Although today's ANNs are considered quite large from a computational perspective, they are smaller than the neural networks of relatively primitive animals such as frogs. Goodfellow, Bengio and Courville, leading scholars in the field, predict that ANNs will not reach the same number of neurons as the human brain possesses before the 2050s unless new technologies enable faster scaling.<sup>159</sup>

However, most current DL models lack reasoning and explanatory capabilities, making them vulnerable to produce unexplainable outcomes. Despite the recent success of DL, DL methods based on ANN generally lack interpretability.<sup>160</sup> Foundation models and LLMs are no exception.<sup>161</sup> Interpretability remains a challenge due to the hierarchical and nonlinear structure of ANNs and the central concept in DL called connectionism. With deep learning models, each artificial neuron works *independently* by computing a relatively simple task, and therefore *partially* contributes to the output produced by the ANNs.<sup>162</sup> ANNs produce output based on the central concept in DL called *connectionism*, where the idea is that a large number of simple computational units (artificial neurons) achieve intelligent behaviour when networked together.<sup>163</sup> Consequently, combining the characteristic of artificial neurons to work independently with the concept of connectionism leads to a situation where thousands or hundreds of thousands of artificial neurons work in parallel in an ANN with hidden layers to jointly calculate certain output.<sup>164</sup> Hence, it seems neither possible to understand which artificial neuron contributed to a distinct part of the output nor to understand what happened in the intermediate (hidden) layers of the ANN.<sup>165</sup> In other words, it is not possible to extract any underlying rules that may be implied by the DL model.<sup>166</sup> This holds even true for DL algorithms using the supervised learning method, where the algorithm cannot learn without being given correct sample patterns. Therefore, even if an ANN has successfully been trained to achieve its goal, the many

<sup>158</sup> Rishi Bommasani et al, 'On the Opportunities and Risks of Foundation Models' (2022) Center for Research on Foundation Models Stanford University 1, 3 <<https://arxiv.org/pdf/2108.07258.pdf>> accessed 8 February 2024.

<sup>159</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 21 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>160</sup> Deng Li and Liu Yang, 'A Joint Introduction to Natural Language Processing and Deep Learning' in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 11, 12.

<sup>161</sup> Melanie Mitchell, David C Krakauer, 'The debate over understanding in AI's large language models' (2023) Vol 120 Iss 3 PNAS 1-5.

<sup>162</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 44.

<sup>163</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 16 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>164</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 155.

<sup>165</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 155.

<sup>166</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 44.

numeric values of the weights produced by the model do not have a meaning to the supervisor.<sup>167</sup> Clearly, the model is parameterised by all these weights, but it remains unclear how these weights have been calculated and to what extent the various input variables contributed to the outcome. ANNs in use can be *updated dynamically* as new data are fed into the network.<sup>168</sup> Subsequently, this updates the weights produced by the model because they are learnt from experience. These updates contribute to further challenges regarding the interpretability of DL approaches.<sup>169</sup>

DL is well suited to deal with complex sensor data such as input from cameras and microphones that proved to be difficult to process when using conventional computational methods.<sup>170</sup> This applies in particular to cognitive tasks which include natural language processing and speech recognition or face recognition, which are discussed below.<sup>171</sup> Current research in DL attempts to decode speech directly from the human brain. Such approaches record the activity in the cortex to decode the characteristics of the produced speech.<sup>172</sup> State-of-the-art deep neural network models arguably contribute to an improved overall accuracy in speech reconstruction from neural recordings in the human auditory cortex.<sup>173</sup> The short-term goal of these research projects is to help individuals that are unable to communicate due to injuries or neurodegenerative disorders by creating a synthesised version of their voice that can be controlled by the activity of their brain speech centres.<sup>174</sup> However, the long-term goal of this could be much broader and very different. Facebook announced that it wants to ‘build a non-invasive, wearable device that lets people type simply by imagining themselves talking.’<sup>175</sup>

### 2.2.2 Natural language processing (NLP)

Natural language processing (NLP), a subfield of AI, aims to give computers the ability to process human language. This interdisciplinary field comprises many concepts and methods such as speech and language processing, human language technology, natural language processing, computational

<sup>167</sup> Toshinori Munakata, *Fundamentals of the New Artificial Intelligence* (2<sup>nd</sup> edn, Springer 2008) 12, 25, 35.

<sup>168</sup> A production model has fixed weights after training. To continuously update weights is possible, but by no means necessary.

<sup>169</sup> Paul De Laat, ‘Algorithmic Decision-Making based on Machine Learning from Big Data: Can Transparency restore Accountability’ (2017) Vol. 31 Issue 4 *Philosophy & Technology* 14 <<https://link.springer.com/article/10.1007%2Fs13347-017-0293-z>> accessed 8 February 2024.

<sup>170</sup> Tommy Chow, Siu-Yeung Cho, *Neural Networks and Computing: Learning Algorithms and Applications* (Imperial College Press 2007), 1/2.; Mitchell Tom T., *Machine Learning* (Mc-Graw-Hill 1997) 95.

<sup>171</sup> Tommy Chow, Siu-Yeung Cho, *Neural Networks and Computing: Learning Algorithms and Applications* (Imperial College Press 2007) 2.

<sup>172</sup> David A. Moses et al, ‘Real-time decoding of question-and-answer speech dialogue using human cortical activity’ (2019) 10 *Nature Communication* <<https://www.nature.com/articles/s41467-019-10994-4.pdf>> accessed 8 February 2024.

<sup>173</sup> Minda Yang et al, ‘Speech Reconstruction from Human Auditory Cortex with Deep Neural Networks’ (Interspeech Conference, Dresden, September 2015) 1124 <<https://dblp.org/db/conf/interspeech/interspeech2015.html>> accessed 8 February 2024.

<sup>174</sup> Nicholas Weiler, ‘Breakthrough device translates brain activity into speech’ (University of California, 25 April 2019) <<https://www.universityofcalifornia.edu/news/synthetic-speech-generated-brain-recordings>> accessed 8 February 2024.

<sup>175</sup> ‘Imagining a new interface: Hands-free communication without saying a word’ (Tech@Facebook, 30 March 2020) <<https://tech.fb.com/imagining-a-new-interface-hands-free-communication-without-saying-a-word/>> accessed 8 February 2024.

linguistics, and speech recognition and synthesis.<sup>176</sup> NLP includes both the generation and understanding of natural language.<sup>177</sup> The advances in NLP have led to the development of large language models (LLMs). LLMs are advanced language models with massive parameter sizes (billions to trillions)<sup>178</sup> and strong learning capabilities.<sup>179</sup> These models can perform various NLP tasks, such as translation, text summarisation, and question-answering.<sup>180</sup> ChatGPT is the current prime example.

From an engineering perspective, NLP intends to develop novel practical applications to facilitate interactions between computers and human languages.<sup>181</sup> Current NLP systems require large amounts of labelled data.<sup>182</sup> Speech recognition is a typical application of NLP, and its aim is to *automatically transcribe* the sequence of spoken words. It may be defined as the process of converting a speech signal to a sequence of words by means of an algorithm implemented by a computer program.<sup>183</sup> In particular, speech recognition does not concern *understanding* but is simply responsible to *convert* language from spoken words to text form.<sup>184</sup> The observable ‘physical’ signal of natural language is called text in symbolic form, and its counterpart is the speech signal, that is, the continuous correspondence of spoken texts.<sup>185</sup> Speech recognition is based on the acoustic signal captured by a microphone as input. The classes are the words that can be uttered. A word is a sequence of phonemes that are the basic speech sounds.<sup>186</sup> Therefore, speech recognition converts phonemes (speech signal) into text. A specific challenge in speech recognition is that different people pronounce the same word differently due to factors related to age, gender or accent, which makes it more difficult to recognise the words.<sup>187</sup> Another challenge is that a common conversational utterance involves multiple queries with disfluencies such as pauses and hesitations. However, current NLP systems embedded in virtual assistants typically focus on ‘unnatural’ and one-sided interactions without hesitation or disfluency. For this reason, speech recognition involving conversational speech is a challenging task.<sup>188</sup>

<sup>176</sup> Daniel Jurafsky, James H Martin, *Speech and Language Processing* (2 edn, Pearson Education Limited 2014) 1.

<sup>177</sup> Stan Franklin, ‘History, motivations, and core themes’ in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 26.

<sup>178</sup> Melanie Mitchell, David C Krakauer, ‘The debate over understanding in AI’s large language models’ Vol 120 Iss 3 PNAS 1-5; Rishi Bommasani et al, ‘On the Opportunities and Risks of Foundation Models’ (2022) Center for Research on Foundation Models Stanford University 1, 3 <<https://arxiv.org/pdf/2108.07258.pdf>> accessed 8 February 2024.

<sup>179</sup> Yupeng Chang et al, ‘A Survey on Evaluation of Large Language Models’ (2023) 1, 4 <<https://arxiv.org/pdf/2307.03109.pdf>> accessed 8 February 2024.

<sup>180</sup> Yiheng Liu et al, ‘Summary of ChatGPT-Related research and perspective towards the future of large language models’ (2023) Vol 1 Meta-Radiology 1 – 14.

<sup>181</sup> Deng Li and Liu Yang, ‘A Joint Introduction to Natural Language Processing and Deep Learning’ in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 1.

<sup>182</sup> Deng Li and Liu Yang, ‘Epilogue: Frontiers of NLP in the Deep Learning Era’ in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 316.

<sup>183</sup> Abhang Priyanka, Gawali Bharti, Mehrotra Suresh, *Introduction to EEG- and speech-based emotion recognition* (Elsevier Inc 2016) 13.

<sup>184</sup> Gokhan Tur et al, ‘Deep Learning in Conversational Language Understanding’ in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 24.

<sup>185</sup> Ibid 24.

<sup>186</sup> Ethem Alpaydin, *Machine Learning: The New AI* (3<sup>rd</sup> edn MIT Press 2016) 67.

<sup>187</sup> Ibid.

<sup>188</sup> Shuo-zhiin Chang et al, ‘Turn-Taking Prediction for Natural Conversational Speech’ (Interspeech Conference Incheon, September 2022) <<https://arxiv.org/pdf/2208.13321.pdf>> accessed 8 February 2024.

Speech signals cannot only reveal the intended message, but also the *identity of the speaker* because the ways in which prosodic characteristics are manifested in speech disclose important information regarding the identity of the speaker.<sup>189</sup> Prosody refers to the study of the intonational and rhythmic aspects of language.<sup>190</sup> Systems in the domain of speaker verification are capable of using the voice of an individual in order to identify an unknown person (speaker identification), verify the identity of a person (speaker verification) and classify specific characteristics like age or gender (speaker classification).<sup>191</sup> Text-based verification of an individual through voice analysis is technically possible with a very short text such as ‘Ok Google’, which takes approximately 0.6 seconds if uttered by an individual.<sup>192</sup> Hence, speaker identity is embedded in the speaker’s voice and can be recognised using automatic speaker recognition systems, which apply DL approaches.<sup>193</sup>

Current research in speech recognition focusses on emotion recognition from speech signals, a major subject in human-computer interaction. This research focusses on how speech is modulated when a speaker’s emotion changes from neutral to another emotional state. For example, it has been observed that speech in anger or happiness shows longer utterance duration and higher pitch and energy value with deep length.<sup>194</sup> Speech emotion recognition may be used for various areas, such as call centres, smart devices or self-driving cars.<sup>195</sup> A real-world application of affective computing (AC) that aims to derive emotional states from speech is Amazon’s ‘Halo’ wearable, which analyses voice tones to detect user emotions.<sup>196</sup> The recent success in NLP and speech recognition has been powered by using the DL approach in ML, currently with supervised ML methods such as classification as described in Section 2.2.1.1. Therefore, the current bottleneck of these approaches is that they require large amounts of labelled data and lack reasoning abilities. However, it is tried to overcome this bottleneck by applying the unsupervised learning paradigm and particularly deep RL methods in NLP and speech recognition.<sup>197</sup> Deep learning has been successfully applied to real-world tasks in AI, in particular in

<sup>189</sup> Leena Mary, *Extraction of Prosody for Automatic Speaker, Language, Emotion and Speech Recognition* (2<sup>nd</sup> edn Springer 2019) 1, 8.

<sup>190</sup> Daniel Jurafsky, James H Martin, *Speech and Language Processing* (2 edn, Pearson Education Limited 2014) 238.

<sup>191</sup> Soufiane Hourri, Jamal Kharroubi, ‘A deep learning approach for speaker recognition’ (2020) Vol. 23 Iss. 1 International Journal of Speech and Technology 123.

<sup>192</sup> Gregor Heigold et al, ‘End-to-End Text-Dependent Speaker Verification’ (2015) <<https://arxiv.org/pdf/1509.08062.pdf>> accessed 8 February 2024.

<sup>193</sup> Leena Mary, *Extraction of Prosody for Automatic Speaker, Language, Emotion and Speech Recognition* (2<sup>nd</sup> edn Springer 2019) 7. See precedent references regarding DL approaches.

<sup>194</sup> Abhang Priyanka, Gawali Bharti, Mehrotra Suresh, *Introduction to EEG- and speech-based emotion recognition* (Elsevier Inc 2016) 14, 105.

<sup>195</sup> See services of the company audeering <<https://www.audeering.com/>> accessed 8 February 2024.

<sup>196</sup> Alex Hern, ‘Amazon’s Halo wristband: the fitness tracker that listens to your mood’ *The Guardian* (London, 28 August 2020) <<https://www.theguardian.com/technology/2020/aug/28/amazons-halo-wristband-the-fitness-tracker-that-listens-to-your-mood>> accessed 8 February 2024; Austin Carr, ‘Amazon’s New Wearable Will Know If I’m Angry. Is That Weird?’ *Bloomberg* (New York, 31 August 2020) <<https://www.bloomberg.com/news/newsletters/2020-08-31/amazon-s-halo-wearable-can-read-emotions-is-that-too-weird>> accessed 8 February 2024.

<sup>197</sup> Deng Li and Liu Yang, ‘Epilogue: Frontiers of NLP in the Deep Learning Era’ in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 316.



speech recognition as a part of the virtual personal assistants such as Google Assistant, Amazon Alexa, Microsoft Cortana or Apple Siri.<sup>198</sup>

### 2.2.3 Computer vision (CV)

Computer vision (CV) is a subfield of AI devoted to perceive objects, i.e. the automated understanding of visual images and comprises many fields of applications.<sup>199</sup> The goal of object detection is to detect all instances of objects from a known class, such as people, cars or faces in an image.<sup>200</sup> CV can also be described as the science and technology of machines that ‘see’, which refers to the ability of the machine to extract information from an image necessary to solve a task.<sup>201</sup> CV aims to infer properties from the observed visual data, which originate from a variety of sensors such as cameras, laser scans, etc.<sup>202</sup> CV algorithms reconstruct the properties of one or more images, such as shape, illumination and colour distributions. Researchers in computer vision develop mathematical techniques to recover the three-dimensional shape and appearance of objects in imagery. Real-world applications include optical character recognition (OCR) for automatic number plate recognitions (of vehicles), medical imaging for preoperative and intra-operative imagery, automotive safety to detect unexpected obstacles such as pedestrians on the street, surveillance to monitor intruders and fingerprint recognition for automatic access authentication.<sup>203</sup>

CV techniques are also currently used to identify individuals based on their gait. Biometric research implies that gait, i.e. the manner in which individuals walk, constitutes a unique identifier like a fingerprint or iris.<sup>204</sup> The biometrics necessary for gait identification may be captured in public places and from a distance in a rather ubiquitous manner. Methods used for identification are model-based approaches which consider the human body or its movements to acquire gait parameters (e.g., step dimensions, cadence, human skeleton, body dimensions) as well as model-free approaches that acquire gait parameters by that rely on gait dynamics and the measurement of geometric representations such as silhouettes.<sup>205</sup>

<sup>198</sup> Gokhan Tur et al, ‘Deep Learning in Conversational Language Understanding’ in Deng Li and Liu Yang (eds) *Deep learning in natural language processing* (Springer 2018) 23.

<sup>199</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 3, Stan Franklin, ‘History, motivations, and core themes’ in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 26.

<sup>200</sup> Yali Amit, Pedro Felzenszwalb, ‘Object Detection’ in Katsushi Ikeuchi (ed) *Computer Vision – A Reference Guide* (Springer 2014) 537.

<sup>201</sup> Sota R. Yoshida, *Computer Vision* (Nova Science Publisher 2011) vii.

<sup>202</sup> Varun Jampani, ‘Learning Inference Models for Computer Vision’ (Dissertation, Universität Tübingen 2016) 1.

<sup>203</sup> Richard Szeliski, *Computer Vision: Algorithms and Applications* (Griets David, Schneider Fred Springer eds 2011) 3, 5.

<sup>204</sup> Ale Sokolova, Anton Konushin ‘Methods of Gait Recognition in Video’ (2019) Vol 45 No 4 Programming and Computer Software 213.

<sup>205</sup> Jure Kovač, Vitomir Štruc, Peter Peer ‘Frame-based classification for cross-speed gait recognition’ (2019) Vol 78 Multimedia Tools and Applications 5621, 5622.

Another real-world example is Amazon Go. Amazon Go is a checkout-free grocery store which is equipped with state-of-the-art cameras and sensors. Amazon Go is powered by computer vision, DL and sensor fusion<sup>206</sup> in order to track shoppers and their purchases. Sensor fusion exploits the best features of sensors (for example, cameras and small Bluetooth radio transmitters called ‘beacons’) installed in a given environment. It is particularly helpful in situations where the sensors themselves are not self-sufficient to achieve a certain goal, for example, comprehensive and precise tracking of shoppers.<sup>207</sup> In Amazon Go stores, shoppers enter by scanning an Amazon Go smartphone app and sensors track items that the shoppers take from the shelves. Once picked up, the items are automatically charged to the Amazon accounts of the shoppers when they leave the store. Where Amazon Go’s inventory system cannot detect the object the user removed from the shelf, the system ‘may consider past purchase history’ of the user.<sup>208</sup>

Face recognition is one of the CV applications of particular relevance for this thesis. Section 2.2.3.1 introduces face recognition and Section 2.2.3.2 explains face recognition applications applying deep learning.

### 2.2.3.1 Face recognition

Face recognition refers to the technology capable of identifying or verifying the identity of subjects in images or videos based on biometric data.<sup>209</sup> It is one of the major biometric technologies and has become increasingly relevant due to the rapid advances in image capture devices and the availability of huge amounts of face images on the web.<sup>210</sup> Unlike other biometric identification methods, such as iris recognition (which requires individuals to get significantly close to a camera), face recognition can be used from a distance and in a covert manner.<sup>211</sup> Therefore, the range of potential applications for face recognition is wide because it can be easily deployed.<sup>212</sup>

<sup>206</sup> See <https://www.amazon.com/b?ie=UTF8&node=16008589011> and Vasilios Mavroudis, Michael Veale ‘Eavesdropping Whilst You’re Shopping: Balancing Personalisation and Privacy in Connected Retail Spaces’ (Living in the Internet of Things Conference, London, March 2018) 6 <<https://ieeexplore.ieee.org/document/8379705>> accessed 8 February 2024.

<sup>207</sup> For example, cameras offer a high level of precision, but might be too expensive to cover the whole shop. Beacons are not self-sufficient to provide tracking data for customer analysis, but can cover a wider operational range. Combined by means of sensor fusion, the sensors allow precise consumer path tracking. See Mirco Sturari et al, ‘Robust and affordable retail customer profiling by vision and radio beacon sensor fusion’ (2016) Vol. 81 Pattern Recognition Letters 30, 31, 40.

<sup>208</sup> Dilip Kumar et al. ‘Detecting item interaction and movement’ US Patent Number US 10268983 (Assignee: Amazon Technologies, Inc.) April 2019 at 9 <<https://patentimages.storage.googleapis.com/01/0b/6e/de57009f5670ae/US20150019391A1.pdf>> accessed 8 February 2024.

<sup>209</sup> Daniel Trigueros, Li Meng, Margaret Hartnett, ‘Face recognition: From Traditional to Deep Learning Methods’ (2018) 1 <<https://arxiv.org/pdf/1811.00116.pdf>> accessed 8 February 2024.

<sup>210</sup> Stan Li, Anil Jain, ‘Introduction’ in Li Stan, Jain Anil (eds) *Handbook of Face Recognition* (2<sup>nd</sup> edn, Springer 2011) 1.

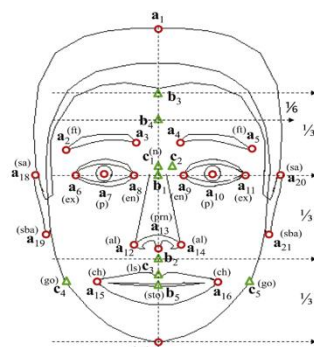
<sup>211</sup> Stan Li, Anil Jain, ‘Introduction’ in Li Stan, Jain Anil (eds) *Handbook of Face Recognition* (2<sup>nd</sup> edn, Springer 2011) 1; Daniel Trigueros, Li Meng, Margaret Hartnett, ‘Face recognition: From Traditional to Deep Learning Methods’ (2018) 1 <<https://arxiv.org/pdf/1811.00116.pdf>> accessed 8 February 2024.

<sup>212</sup> Daniel Trigueros, Li Meng, Margaret Hartnett, ‘Face recognition: From Traditional to Deep Learning Methods’ (2018) 1 <<https://arxiv.org/pdf/1811.00116.pdf>> accessed 8 February 2024.

Face recognition systems operate in a face verification (authentication) and/or face identification (recognition) mode. The former involves a one-on-one match that compares a query face image of the person whose identity is claimed (e.g., for self-serviced immigration clearance using E-passports). The latter involves one-to-many matching, which compares a query face image against multiple face images in a database to associate the identity of the query face. Usually, finding the most similar face is not sufficient and a confidence threshold is specified. Therefore, only those faces whose similarity score is above the threshold are reported.<sup>213</sup> Face recognition systems are usually built on four building blocks:

1. Face detection, which finds the position of a face in an image
2. Face normalisation, which normalises the face geometrically and photometrically
3. Face feature extraction performed to extract salient information which is useful to distinguish faces such as reference points located at fixed locations in the face (e.g., position of eyes, nose, lips)
4. Face matching, where extracted features from the input face are matched against one or many of the enrolled faces in the database<sup>214</sup>

The facial features used for the third building block may be grouped into two classes of features: continuous and discrete. Continuous features are real valued numbers and are extracted using distances and angles between facial landmarks such as forehead height, eyebrow length, nose height, chin height, ears length, mouth length etc. Discrete features represent a finite number of categories, for example, the shape of the eyebrow or nose root width.<sup>215</sup> Figure 1.5 provides an example of such features.



**Figure 1.5** Face layout illustrated by Tome et al.<sup>216</sup> with examples of facial features extracted by using distances and angles between facial landmarks such as eyebrows, eyes, nose and lips. Used with permission.

<sup>213</sup> Stan Li, Anil Jain, 'Introduction' in Li Stan, Jain Anil (eds) *Handbook of Face Recognition* (2<sup>nd</sup> edn, Springer 2011) 3.

<sup>214</sup> Ibid, 4; Daniel Trigueros, Li Meng, Margaret Hartnett, 'Face recognition: From Traditional to Deep Learning Methods' (2018) 1 < <https://arxiv.org/abs/1811.00116> > accessed 8 February 2024.

<sup>215</sup> Pedro Tome et al., 'Facial soft biometric features for forensic face recognition' (2015) Vol 257 *Forensic Science International* 271, 273.

<sup>216</sup> Ibid.

### 2.2.3.2 DL and face recognition

Current face recognition applications use DL methods based on convolutional neural networks (CNN) which are trained with very large datasets.<sup>217</sup> A CNN is a specific kind of neural network for processing data that has a known grid-like typology. For example, image data can be thought of as a 2D grid of pixels. As the name indicates, a CNN employs a mathematical operation called convolution, which is a specialised kind of linear operation.<sup>218</sup> Notably, the performance of a face recognition system largely depends on a variety of factors such as illumination, facial pose, expression, age span, hair and motion.<sup>219</sup> Whereas the building blocks of face recognition systems and the general architecture of the ANN are predetermined by the developer of the system, the ANN itself decides how to create the optimal score for determining similarity in the face matching building block mentioned in Section 2.2.3.1. Therefore, it remains often unclear how the similarity score is calculated by the ANN, even to the developer of the system.<sup>220</sup> Another issue is that face recognition systems perform poorly in recognising individuals of different ethnicities. For example, Hewlett Packard face recognition software could not recognise dark-coloured faces as faces.<sup>221</sup> A ‘passport robot’ in New Zealand rejected the passport picture of an Asian man because the ‘subject’s eyes are closed’ although his eyes were open.<sup>222</sup>

However, face recognition systems are widely used in commercial applications and consumer products with built-in AI capabilities. Examples are cars with on-board cameras to deploy biometric identification and monitor driving behaviour<sup>223</sup> or connected retail spaces.<sup>224</sup> Furthermore, there is a trend to improve face recognition systems with the ability to monitor and analyse the emotions in real-time based on extracted biometric data and facial expressions. The gained knowledge is then used to build specific customer profiles.

<sup>217</sup> Daniel Trigueros, Li Meng, Margaret Hartnett, ‘Face recognition: From Traditional to Deep Learning Methods’ (2018) 1 <<https://arxiv.org/abs/1811.00116>> accessed 8 February 2024.

<sup>218</sup> Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning* (MIT Press 2016) 326 <[www.deeplearningbook.org](http://www.deeplearningbook.org)> accessed 8 February 2024.

<sup>219</sup> Stan Li, Anil Jain, ‘Introduction’ in Li Stan, Jain Anil (eds) *Handbook of Face Recognition* (2<sup>nd</sup> edn, Springer 2011) 3.

<sup>220</sup> Yana Welinder, Aeryn Palmer, ‘Face Recognition, Real-Time Identification, and Beyond’ in Selinger Evan, Polonetsky Jules, Tene Omer (eds) *The Cambridge Handbook of Consumer Privacy* (Cambridge University Press 2018) 104.

<sup>221</sup> Frederik Zuiderveen Borgesius, ‘Discrimination, artificial intelligence, and algorithmic decision-making’ (2019) Report for the Anti-discrimination department of the Council of Europe, 17 <<https://rm.coe.int/discrimination-artificial-intelligence-and-algorithmic-decision-making/1680925d73>> accessed 8 February 2024.

<sup>222</sup> Regan James, ‘New Zealand passport robot tells applicant of Asian descent to open eyes’ (2016) *Reuters* <<https://www.reuters.com/article/us-newzealand-passport-error/new-zealand-passport-robot-tells-applicant-of-asian-descent-to-open-eyes-idUSKBN13W0RL>> accessed 8 February 2024.

<sup>223</sup> See <<https://visagetechnologies.com/application-fields/driver-monitoring/>> accessed 8 February 2024

<sup>224</sup> See <<https://www.einfochips.com/blog/facial-recognition-in-retail-enhance-in-store-customer-experience-and-improve-retailer-operations/>> accessed 8 February 2024.

### 2.2.4 Affective computing (AC)

Affective computing (AC), sometimes called ‘emotion AI’, is computing that relates to, arises from or influences emotion.<sup>225</sup> AC is a scientific and engineering endeavour inspired by psychology, neuroscience, linguistics and related areas.<sup>226</sup> Affective states are considered to be experiential phenomena such as emotions and moods.<sup>227</sup> Emotions form an important part of human intelligence and daily life, be it for decision-making, social interaction, perception or learning. In other words, emotions play a pivotal role in functions considered essential to intelligence.<sup>228</sup> Picard, the pioneer in the field of AC, therefore, concludes that if computers are to be genuinely intelligent, they too should have emotional capabilities.<sup>229</sup> In this thesis, the focus lies on affect detection from facial expressions and speech, since they may be easily deployed compared to more invasive approaches that include measurement of physiological factors such as cardiac activity (heart rate) or skin conductance (sweat).

The following sections elaborate on affect detection from facial expressions (Section 2.2.4.1), speech (Section 2.2.4.2) and discuss multimodal approaches in which different methods of AC are combined to detect emotions (Section 2.2.4.3).

#### 2.2.4.1 Facial expressions

Facial expressions are probably the most natural way humans express their emotions.<sup>230</sup> According to Darwin’s evolutionary theory of emotions, emotion expressions help in regulating the social interaction and increase the likelihood of survival.<sup>231</sup> Due to the developments in technology, it is possible to detect facial information automatically in real-time, for example, with the use of a simple video camera. However, automatic detection of emotions derived from facial expressions and their interpretation is not simple and context-driven.<sup>232</sup> Physically, a facial expression is a change in the face due to movements of several muscles demonstrating an emotional state. An emotional state is an individual’s transient reaction to specific encounters with the environment, one that occurs and disappears depending on particular conditions. For example, someone is feeling or reacting with anger at a particular time and place.<sup>233</sup> A facial expression is communicated by a transient flexing of facial

<sup>225</sup> Rosalind W Picard, ‘Affective Computing’ (1995) MIT Media Laboratory Perceptual Computing Section Technical Report No 321 at 1 <<https://hd.media.mit.edu/tech-reports/TR-321.pdf>> accessed 8 February 2024.

<sup>226</sup> Rafael Calvo et al, ‘Introduction to Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 2.

<sup>227</sup> Steffen Steinert, Orsolya Friedrich, ‘Wired Emotions: Ethical Issues of Affective Brain–Computer Interfaces’ (2020) Vol 26 Science and Engineering Ethics 351, 352.

<sup>228</sup> Rosalind W Picard, *Affective Computing* (MIT Press 1997) 47.

<sup>229</sup> Ibid preface x.

<sup>230</sup> Rafael Calvo et al, ‘Introduction to Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 4.

<sup>231</sup> Avinash Awasthi, Manas K. Mandal, ‘Facial Expressions of Emotions: Research Perspectives’ in Manas K. Mandal, Avinash Awasthi (eds) *Understanding Facial Expressions in Communication* (Springer 2015) 3.

<sup>232</sup> Catherine Marechal et al, ‘Survey on AI-Based Multimodal Methods for Emotion Detection’ in Joanna Kołodziej, Horacio González-Vélez (eds) *High-Performance Modelling and Simulation for Big Data Applications* (Springer 2019) 314, 315.

<sup>233</sup> Richard S Lazarus, *Emotion and Adaption* (OUP 1991) 46, 47.

futures such as mouth, eyes and eyebrows due to the contraction of the muscles that make up the face.<sup>234</sup> These muscle contractions are controlled by two different areas of the brain, one controlling voluntary movements and the other involuntary reactions.<sup>235</sup> Facial expressions can easily be used for emotion detection because it only requires a simple video camera to register facial information automatically and in real-time.<sup>236</sup> Two approaches to measuring facial expressions will be discussed here: message-based and sign-based approaches.

Based on the assumption that the face provides a direct ‘readout’ of emotion, the message-based approach makes inferences about emotion or the affective state by assigning facial expression and movements to ‘basic emotions’ according to Ekman.<sup>237</sup> Facial movements and the ‘basic emotions’ hypothesised are illustrated in Figure 1.6.



**Figure 1.6** Facial movements and hypothesised ‘basic’ emotion categories illustrated by Barret et al.<sup>238</sup> Used with permission.

It should be noted that this approach is problematic since the meaning of an expression depends on the context. For example, smiles accompanied by cheek raising express enjoyment, the same smile combined with head lowering and turning to the side convey embarrassment. Additionally, facial expressions can be posed or faked.<sup>239</sup>

The sign-based approach measures anatomic facial signs and then uses experimental or observational methods to discover the relation between these signs and emotion.<sup>240</sup> In 1978, the psychologists Ekman and Friesen proposed a model for measuring facial muscle contractions involved in facial expression called ‘Facial Action Coding System’ (FACS).<sup>241</sup> FACS is now a common standard used to

<sup>234</sup> Alice Caplier, ‘Visual Emotion Recognition: Status and Key Issues’ in Catherine Pelachaud (ed) *Emotion-oriented Systems* (Wiley-ISTE 2012) 107, 109.

<sup>235</sup> Hyisung C. Hwand, David Matsumoto, ‘Emotional Expression’ in Catharine Abell, Joel Smith (eds) *The Expression of Emotion* (CUP 2016) 139, 140.

<sup>236</sup> Catherine Marechal et al, ‘Survey on AI-Based Multimodal Methods for Emotion Detection’ in Joanna Kolodziej, Horacio Gonzalez-Vélez (eds) *High-Performance Modelling and Simulation for Big Data Applications* (Springer 2019) 314.

<sup>237</sup> Jeffrey F. Cohn, Fernando De La Torre, ‘Automated Face Analysis for Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 132, 133.

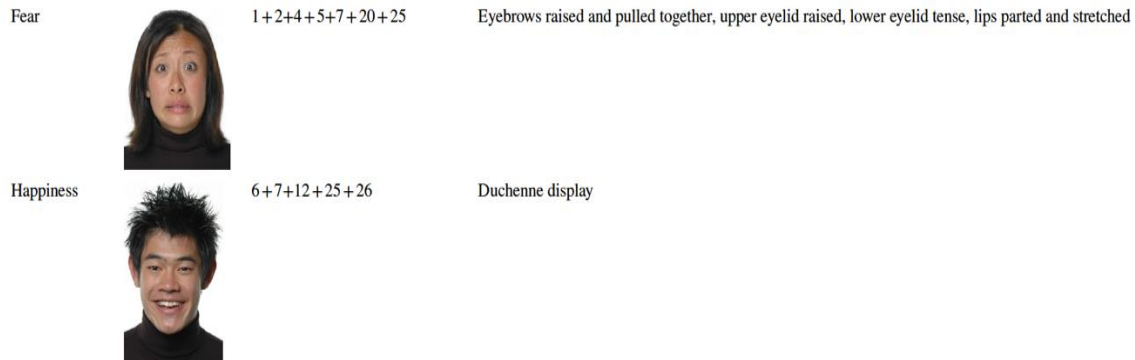
<sup>238</sup> Lisa Feldman Barrett et al. ‘Emotional Expressions Reconsidered’ (2019) Vol 20 (1) *Psychological Science in the Public Interest* 1, 19.

<sup>239</sup> Jeffrey F. Cohn, Fernando De La Torre, ‘Automated Face Analysis for Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 132.

<sup>240</sup> *Ibid* 133.

<sup>241</sup> Paul Ekman, Wallace V. Friesen, ‘Facial Action Coding System: A Technique for the Measurement of Facial Movements’ (1978) Consulting Psychologists Press.

systematically describe and quantify visible human facial movement<sup>242</sup> and describes facial activity in terms of anatomically-based action units (AU).<sup>243</sup> FACS defines 46 AUs to describe each independent movement of the face, including head and eye movements. FACS is used to verify the physiological presence of emotion. Due to its comprehensiveness, it also allows the discovery of new patterns related to emotional states.<sup>244</sup> FACS-coded facial events (AUs) such as ‘Inner Brow Raiser’, ‘Chin Raiser’, ‘Lip Corner Puller’ are classified into emotion categories by matching facial events with emotional events coded from previous empirical studies.<sup>245</sup> Figure 1.7 provides some examples of AUs.



**Figure 1.7** Facial expression examples for basic emotions ‘fear’ and ‘happiness’, the corresponding FACS action units and physical descriptions for each expression, illustrated by Keltner et al.<sup>246</sup> Used with permission.

Manual application of the FACS to videotaped behaviour is very time consuming. It takes approximately 100 hours to train a person to make judgements reliably and typically takes more than two hours to complete a one-minute video.<sup>247</sup>

Unsurprisingly, computer scientists started to use computer vision and graphics to automatically analyse and synthesise facial expression in automated face analysis (AFA) systems. Recently developed AFA systems claim to detect pain, frustration, emotion intensity, depression and psychological distress.<sup>248</sup> For example, a study aimed to predict depression, anxiety and stress levels from videos using the FACS approach built on ANN-based architecture.<sup>249</sup> Automated face analysis (AFA) systems seek to detect emotions using message-based and sign-based approaches. Such systems typically follow

<sup>242</sup> Lisa Feldman Barrett et al. ‘Emotional Expressions Reconsidered’ (2019) Vol 20 (1) *Psychological Science in the Public Interest* 1, 52.

<sup>243</sup> Jeffrey F. Cohn, Fernando De La Torre, ‘Automated Face Analysis for Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 132.

<sup>244</sup> Marian Stewart Bartlett et al, ‘Toward Automatic Recognition of Spontaneous Facial Actions’ in Paul Ekman, Erika L. Rosenberg, *What the Face Reveals* (2<sup>nd</sup> edn OUP 2005) 393, 394.

<sup>245</sup> Erika L. Rosenberg, ‘Introduction: The Study of Spontaneous Facial Expressions in Psychology’ in Paul Ekman, Erika L. Rosenberg, *What the Face Reveals* (2<sup>nd</sup> edn OUP 2005) 14, 16.

<sup>246</sup> Dacher Keltner et al. ‘Emotional Expression: Advances in Basic Emotion Theory’ (2019) Vol 43 Iss 2 *Journal of Non-verbal Behaviour* 133, 142.

<sup>247</sup> Marian Stewart Bartlett et al, ‘Toward Automatic Recognition of Spontaneous Facial Actions’ in Paul Ekman, Erika L. Rosenberg, *What the Face Reveals* (2<sup>nd</sup> edn OUP 2005) 394.

<sup>248</sup> Jeffrey F. Cohn, Fernando De La Torre, ‘Automated Face Analysis for Affective Computing’ in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 132, 133.

<sup>249</sup> Mihai Gavrilescu, Nicolae Vizireanu (2019) ‘Predicting Depression, Anxiety, and Stress Levels from Videos Using the Facial Action Coding System’ (2019) Vol 19 No 17, 1.

four steps: face detection, face registration, feature extraction and classification. In the first step, the face will be recognised using approaches from object detection in computer vision. During the face registration step, the face is rotated to an upright and frontal facing position to remove geometric differences.<sup>250</sup> In the feature extraction step, the algorithm extracts the main features of the face (e.g. mouth, eyebrows) and analyses movement, shape and texture composition of these regions to identify AUs.<sup>251</sup> After feature extraction, a machine learning (ML) component has the task to learn the relationship between the feature representation and the target facial expressions. Most of the current approaches use supervised learning<sup>252</sup>, with a tendency to also make use of deep learning and ANN methods.<sup>253</sup> Fully automatic FACS coding systems use state-of-the-art ML techniques that can recognise any facial action.<sup>254</sup>

#### 2.2.4.2 Speech in affective computing

Emotions of a person may be measured and quantified by observing speech signals from this person. This is exactly what speech-based emotion recognition systems aim at. Such systems are based on insight gained from research that investigates the mechanisms of emotional speech production.<sup>255</sup> Research in emotion recognition has shown that emotions in speech are related to prosody features such as pitch and energy.<sup>256</sup> Prosody refers to the study of the intonational and rhythmic aspects of language. Research has demonstrated specific associations between emotions such as fear, anger, sadness, joy and measures of pitch, voice level and speech rate.<sup>257</sup> Pitch is a perceptual property of a signal. The pitch of a sound is the mental sensation of fundamental frequency. In case a sound has a higher frequency, it is generally perceived as having a higher pitch.<sup>258</sup> The pitch of speech associated with emotions such as anger or happiness is higher than the pitch of speech associated with emotions such as sadness or disappointment.<sup>259</sup> In terms of speech rate, it has been shown that if the person who speaks is in an emotional state of anger or fear, the speech is usually faster. In case the person is bored or sad, then the speech is typically slower. Hence, effects of emotion tend to be present in features

<sup>250</sup> Michael Valstar, 'Automatic Facial Expression Analysis' in Manas K. Mandal, Avinash Awasthi (eds) *Understanding Facial Expressions in Communication* (Springer 2015) 144-150.

<sup>251</sup> Catherine Marechal et al, 'Survey on AI-Based Multimodal Methods for Emotion Detection' in Joanna Kołodziej, Horacio González-Vélez (eds) *High-Performance Modelling and Simulation for Big Data Applications* (Springer 2019) 315.

<sup>252</sup> Jeffrey F. Cohn, Fernando De La Torre, 'Automated Face Analysis for Affective Computing' in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 137.

<sup>253</sup> Panagiotis Tzirakis et al, 'End-to-End Multimodal Emotion Recognition using Deep Neural Networks' (2015) Vol. 14 No. 8 Journal of Latex Class Files, 1.

<sup>254</sup> Marian Stewart Bartlett et al, 'Toward Automatic Recognition of Spontaneous Facial Actions' in Paul Ekman, Erika L. Rosenberg, *What the Face Reveals* (2<sup>nd</sup> edn OUP 2005) 395.

<sup>255</sup> Chi-Chun Lee et al, 'Speech in Affective Computing' in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 171.

<sup>256</sup> Ricardo A. Calix, Leili Javadpour, Gerald M. Knapp, 'Detection of Affective States From Text and Speech For Real-Time Human-Computer Interaction' (2012) Vol 54 No 4 Human Factors and Ergonomics Society 530, 531.

<sup>257</sup> Christina Sobn and Murray Alpert, 'Emotion in Speech: The Acoustic Attributes of Fear, Anger, Sandess, and Joy' (1999) Vol 28 No 4 Journal of Psycholinguistic Research, 347.

<sup>258</sup> Daniel Jurafsky, James H Martin, *Speech and Language Processing* (2 edn, Pearson Education Limited 2014) 238.

<sup>259</sup> Ze-Jing Chuang, Chung-Hsien Wu, 'Multi-Modal Emotion Recognition from Speech and Text' (2004) Vol. 9 No. 2 Computational Linguistics and Chinese Language Processing, 45-62.



such as average pitch, pitch range and pitch changes, speech rate, voice quality and articulation.<sup>260</sup> Approaches in affective computing extract these acoustic signal features that characterise emotional speech. Machine learning algorithms map the automatically derived acoustic features described before to the desired emotion representations.<sup>261</sup> Research in the field aims to extract features from the voice to detect depressive people<sup>262</sup> or candidate stress levels during human resources interviews using ML and ANN.<sup>263</sup> Real-world applications of AC that aim to derive emotional states from speech are Amazon's 'Halo' wearable, which analyses voice tones to detect user emotions,<sup>264</sup> or Spotify's patented voice assistant,<sup>265</sup> which, based on commands or other utterances (e.g., 'ugh'), recognises when a user sounds sad and then offers encouragement by 'cheering' the user.<sup>266</sup> Methods applied to speech emotion recognition increasingly involve deep learning approaches.<sup>267</sup>

### 2.2.4.3 Multimodal approaches

Methods used in AC may be combined in multimodal approaches. For example, research in psychology aims to develop multimodal frameworks comprising audio-video fusion (facial expressions and emotions in speech) for the diagnosis, of depression to distinguish between people who suffer from depression and people who do not.<sup>268</sup>

Multimodal approaches may also include the detection from physiological factors such as cardiac activity (heart rate and heart rate variability). Research has shown that the variability of heart rate provides a novel marker to recognise emotions in humans.<sup>269</sup> Both heart rate and heart rate variability have been reported as indicators of fear, panic, anger and appreciation and are therefore used for affective computing.<sup>270</sup> Methods in AC can be integrated in commercial applications in order to track

<sup>260</sup> Rosalind W Picard, *Affective Computing* (MIT Press 1997) 179, 180.

<sup>261</sup> Chi-Chun Lee et al, 'Speech in Affective Computing' in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 173, 177.

<sup>262</sup> Marius Dan Zbancioc, Silvia Monica Feraru, 'A study about the automatic recognition of the anxiety emotional state using Emo-DB' (E-Health and Bioengineering Conference, Iasi, 2015) 1.

<sup>263</sup> Kevin Tomba et al, 'Stress Detection Through Speech Analysis' (2018) Vol 1 ICETE 2018, 560.

<sup>264</sup> Alex Hern, 'Amazon's Halo wristband: the fitness tracker that listens to your mood' *The Guardian* (London, 28 August 2020) <<https://www.theguardian.com/technology/2020/aug/28/amazons-halo-wristband-the-fitness-tracker-that-listens-to-your-mood>> accessed 8 February 2024; Austin Carr, 'Amazon's New Wearable Will Know If I'm Angry. Is That Weird?' *Bloomberg* (New York, 31 August 2020) <<https://www.bloomberg.com/news/newsletters/2020-08-31/amazon-s-halo-wearable-can-read-emotions-is-that-too-weird>> accessed 8 February 2024.

<sup>265</sup> Daniel Bromand et al, 'Systems and Methods for Enhancing Responsiveness to Utterances Having Detectable Emotion' US Patent Number US 10566010 B2 (Assignee: Spotify AB) February 2020 11 <<https://patentimages.storage.googleapis.com/2a/9d/2d/926b58a2bd956f/US10566010.pdf>>, accessed 8 February 2024.

<sup>266</sup> Josh Mandell, 'Spotify Patents A Voice Assistant That Can Read Your Emotions' *Forbes* (New York, 12 March 2020) <<https://www.forbes.com/sites/joshmandell/2020/03/12/spotify-patents-a-voice-assistant-that-can-read-your-emotions/>> accessed 8 February 2024.

<sup>267</sup> Haytham M Fayek, Margaret Lech, Lawrence Cavedon, 'Evaluating deep learning architectures for Speech Emotion Recognition' (2017) Vol 92 Neural Networks 60.

<sup>268</sup> Jyoti Joshi et al, 'Multimodal assistive technologies for depression diagnosis and monitoring' (2013) Vol 7 Journal on Multimodal User Interfaces, 217.

<sup>269</sup> Quintana Daniel et al. 'Heart rate variability is associated with emotion recognition: Direct evidence for a relationship between the automatic nervous system and social cognition' (2012) Vol 86 No 2 International Journal of Psychophysiology 168.

<sup>270</sup> Jennifer Healey, 'Physiological Sensing of Emotion' in Rafael Calvo et al (eds), *The Oxford Handbook of Affective Computing* (OUP 2015) 206.

and analyse customer behaviour in retail stores, so-called behaviour inference systems. Behavioural inference systems apply deep learning (DL) and affective computing in order to monitor and analyse the shopper's behaviour based on extracted physiological factors (heart rate) and facial expressions.<sup>271</sup> An example for such a system comprises six modules: a speech recognition module, a biofeedback model, a facial expression and emotion recognition module, a gaze detection module, an age and gender recognition module and an identification module.<sup>272</sup>

### 2.2.5 Automated reasoning (AR)

Automated reasoning (AR) aims to develop computers that can use stored information to answer questions and draw new conclusions.<sup>273</sup> It may be described as the science of developing methods that intend to replace human reasoning with procedures that perform individual reasoning automatically.<sup>274</sup> Automated reasoning is devoted to answering questions from diverse data without human intervention and includes decision-making. As a form of reasoning, decision-making focusses on an autonomous agent trying to perform a task for a human.<sup>275</sup> Reasoning problems are of practical significance, they arise naturally in many applications that interact with the world, for example, reasoning about knowledge in the sciences or natural language processing. Furthermore, reasoning algorithms form the foundation for theoretical investigations into general AI (human-level AI).<sup>276</sup> Reasoning is the process of obtaining new knowledge from a given knowledge, where certain transformation rules are applied that depend only on knowledge and can be done exclusively in the brain without involving senses.<sup>277</sup> Research in automated reasoning focusses on logical reasoning, probabilistic reasoning and common sense reasoning.<sup>278</sup> Logical reasoning attempts to avoid any unjustified assumptions and confines itself to inferences that are infallible and beyond reasonable dispute.<sup>279</sup> Probabilistic reasoning deals with uncertainty about knowledge and belief. Uncertainty may be approached by applying tools from probability theory and statistics. Research in probabilistic reasoning focusses on the representation of different types of uncertainty and uncertain knowledge, reasoning with these types of knowledge, and learning them. It facilitates the development of applied systems of practical importance, such as machine vision, medical diagnosis and natural language processing. Probabilistic reasoning models are close to ML and serve as a medium between ML and AR.<sup>280</sup>

<sup>271</sup> Andrea Generosi, Silvia Ceccacci, Maura Mengoni, 'A deep learning-based system to track and analyse customer behaviour in retail store' (IEEE 8<sup>th</sup> International Conference on Consumer Electronics, Berlin 2018) 36.

<sup>272</sup> Ibid 37.

<sup>273</sup> Stuart Russel, Peter Norvig, *Artificial Intelligence, A Modern Approach* (3rd edn, Pearson Education 2016) 2.

<sup>274</sup> Tudor Jebelean et al, 'Automated Reasoning' in Buchberger Bruno et al (eds) *Hagenberg Research* (Springer 2009) 63.

<sup>275</sup> Amir Eyal, 'Reasoning and decision making' in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 191.

<sup>276</sup> Ibid 191.

<sup>277</sup> Tudor Jebelean et al, 'Automated Reasoning' in Buchberger Bruno et al (eds) *Hagenberg Research* (Springer 2009) 63.

<sup>278</sup> Amir Eyal, 'Reasoning and decision making' in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 193.

<sup>279</sup> John Harrison, *Handbook of Practical Logic and Automated Reasoning* (Cambridge University Press 2009) 1.

<sup>280</sup> Amir Eyal, 'Reasoning and decision making' in Frankish Keith and Ramsey William M. (eds) *The Cambridge Handbook of Artificial Intelligence* (2014) 201.

For a very long time, scientists and philosophers have tried to understand and formalise how humans reason and whether reasoning methods may be automatised.<sup>281</sup> Achieving common sense reasoning capabilities in computational systems has been one of the goals of AI since its beginning in the 1960s.<sup>282</sup> Common sense reasoning constitutes a central part of human behaviour and is a precondition for human intelligence. Unsurprisingly, the creation of systems that exhibit common sense reasoning is a central goal towards achieving general AI. History in AI has proven that it is more difficult to develop systems with common sense reasoning capabilities compared to systems that solve explicit reasoning problems, such as chess-playing programs or expert systems that assist in clinical diagnosis. Part of this difficulty is due to the all-encompassing aspect of common sense reasoning: It requires many different kinds of knowledge. Furthermore, most common sense knowledge is implicit and therefore difficult to explain and compute, unlike expert-knowledge which is usually explicit. Therefore, implicit common sense knowledge must be made explicit in order to develop common sense reasoning systems.<sup>283</sup>

Other problems that impede the development of automated common sense reasoning are the lack of a precise meaning of ‘common sense reasoning’, how to take into account of polysemy, ambiguity and vagueness of natural language and the difficulty in modelling the role of various forms of implicit knowledge such as context, background knowledge and tacit knowledge.<sup>284</sup> Therefore, common sense reasoning capabilities are still a challenge in AI applications.<sup>285</sup> According to Oren Etzioni, who oversees the Allen Institute for Artificial Intelligence, AI ‘is devoid of common sense’.<sup>286</sup> Hence, to acquire common sense from massive amounts of data and implementing it in intelligent systems appears to be the next frontier in AI.<sup>287</sup> The lack of progress in providing general automated common sense reasoning capabilities underscores that this is a very difficult problem in the field of AI.<sup>288</sup> Common sense reasoning is not just the hardest problem for AI, it is also considered to be the most important problem.<sup>289</sup>

<sup>281</sup> Marco Gavanelli, Toni Mancini, ‘Automated Reasoning’ (2013) Vol. 7 No. 2 *Intelligenza Artificiale* 113.

<sup>282</sup> Brandon Bennet, Anthony G Cohn, ‘Automated Common-sense Spatial Reasoning: Still a Huge Challenge’ in Stephen Muggleton, Nicholas Chater (eds) *Human-Like Machine Intelligence* (Oxford University Press 2021) 405.

<sup>283</sup> Ernest Davis, Leora Morgenstern, ‘Introduction: Progress in formal common sense reasoning’ (2004) Vol 153 *Artificial Intelligence* 1.

<sup>284</sup> Brandon Bennet, Anthony G Cohn, ‘Automated Common-sense Spatial Reasoning: Still a Huge Challenge’ in Stephen Muggleton, Nicholas Chater (eds) *Human-Like Machine Intelligence* (Oxford University Press 2021) 406.

<sup>285</sup> Shoham Yoav et al, ‘The AI Index 2018 Annual Report’ (AI Index Steering Committee Stanford University 2018) 64 <[https://hai.stanford.edu/sites/default/files/2020-10/AI\\_Index\\_2018\\_Annual\\_Report.pdf](https://hai.stanford.edu/sites/default/files/2020-10/AI_Index_2018_Annual_Report.pdf)> accessed 8 February 2024.

<sup>286</sup> Cade Metz, ‘Paul Allen Wants to Teach Machines Common Sense’ *The New York Times* (New York, 28 February 2018) <<https://www.nytimes.com/2018/02/28/technology/paul-allen-ai-common-sense.html>> accessed 8 February 2024.

<sup>287</sup> Niket Tandon, Aparna S. Varde, Gerard de Melo, ‘Commonsense Knowledge in Machine Intelligence’ (2017) Vol 46 No 4 *SIGMOD Record* 49.

<sup>288</sup> Brandon Bennet, Anthony G Cohn, ‘Automated Common-sense Spatial Reasoning: Still a Huge Challenge’ in Stephen Muggleton, Nicholas Chater (eds) *Human-Like Machine Intelligence* (Oxford University Press 2021) 405.

<sup>289</sup> Gary Marcus, Ernest Davis, *Rebooting AI: Building Artificial Intelligence we can trust* (Pantheon Books 2019).

### 2.3 Conclusions

This chapter answered Subquestion 1, namely, what AI is and what disciplines exist therein. AI is an exciting, challenging and complex technology which accelerates at a tremendous pace. AI covers a broad range of approaches and techniques and at least five disciplines. These five disciplines are machine learning, natural language processing, computer vision, affective computing and automated reasoning.

As a major discipline of AI, *machine learning* (ML) is focussed on computers that program themselves based on experience. ML can be applied by means of several methods, ranging from supervised to unsupervised to reinforcement learning. *Deep learning* (DL) is a very powerful kind of machine learning considering that the achievements in the field have been reached with *artificial neural networks* (ANNs) comprising an astonishingly small number of neurons when compared with neural networks of the human brain. By means of *natural language processing* (NLP), machines can process human language. It includes both the generation and understanding of natural language. NLP significantly contributes to improved interactions between machines and humans. *Computer vision* (CV) facilitates the automated processing of visual images and thus enables machines to see. Face recognition, which is one of the applications of computer vision, empowers machines to identify or verify the identity of humans in images or videos based on biometric data. Because emotions form an important factor of human intelligence and daily life, *affective computing* (AC) aims to equip machines with emotional capabilities. Approaches in AC which derive emotions from facial expressions and speech may be easily deployed and widely used. Efforts in the discipline of *automated reasoning* (AR) seek to perform individual reasoning automatically.