

TECHNICAL WHITEPAPER

Teaching Machines to Detect and
Understand Humans

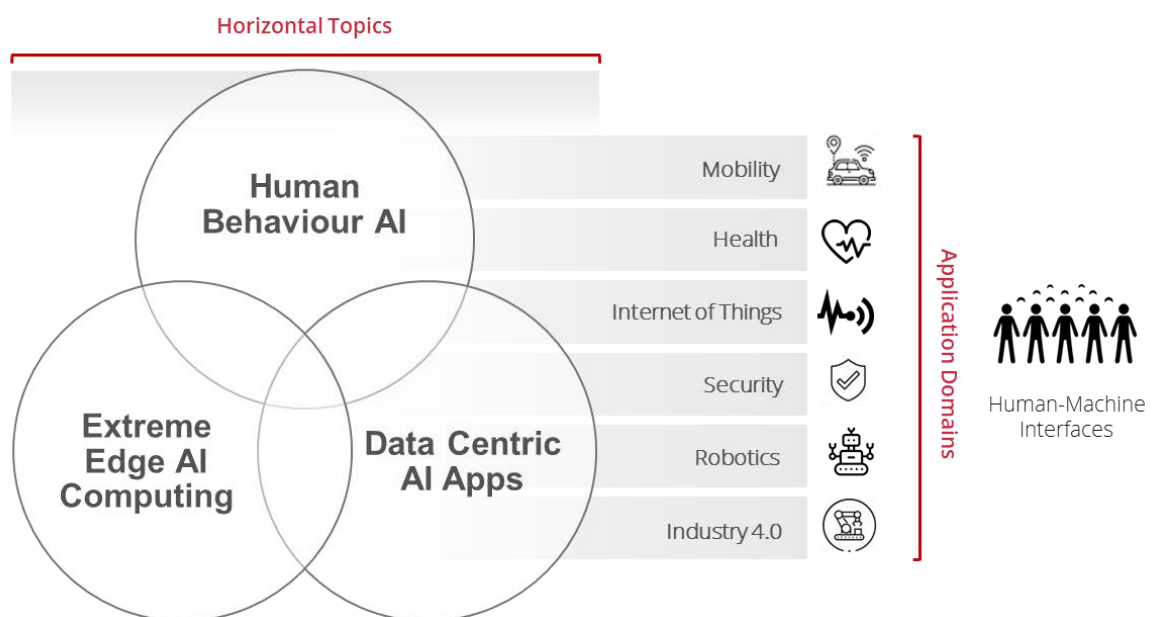
Introduction

The age of autonomous machines and intelligent IoT devices is upon us. They will affect many aspects of our lives. They will bring about a new generation of diagnostic instruments to clinics, improve cars and transportation, and create new consumer experiences, while being embedded into the networks that power and inform society underpinning the efficient provision of public and private services. Increasingly they sense their surroundings, operate autonomously and collaborate with other devices seamlessly.

Critical to their operation, will be improving the safety, security, and quality of life of humans that interact with such intelligent systems. To interact efficiently and effectively, sensing and interpreting human behaviour will become mission critical to connecting the human user with automated systems. This white paper outlines the core horizontal topics necessary to teaching machines to sense and understand humans and use cases where “AI at the Extreme Edge” is not only critical but transformative.

Topics covered in this technical whitepaper include:

- What is **Human Behaviour AI** and what are key **Use Cases**?
- What are **Data Centric AI Apps** and how are they developed?
- What is **Extreme Edge AI Computing**?

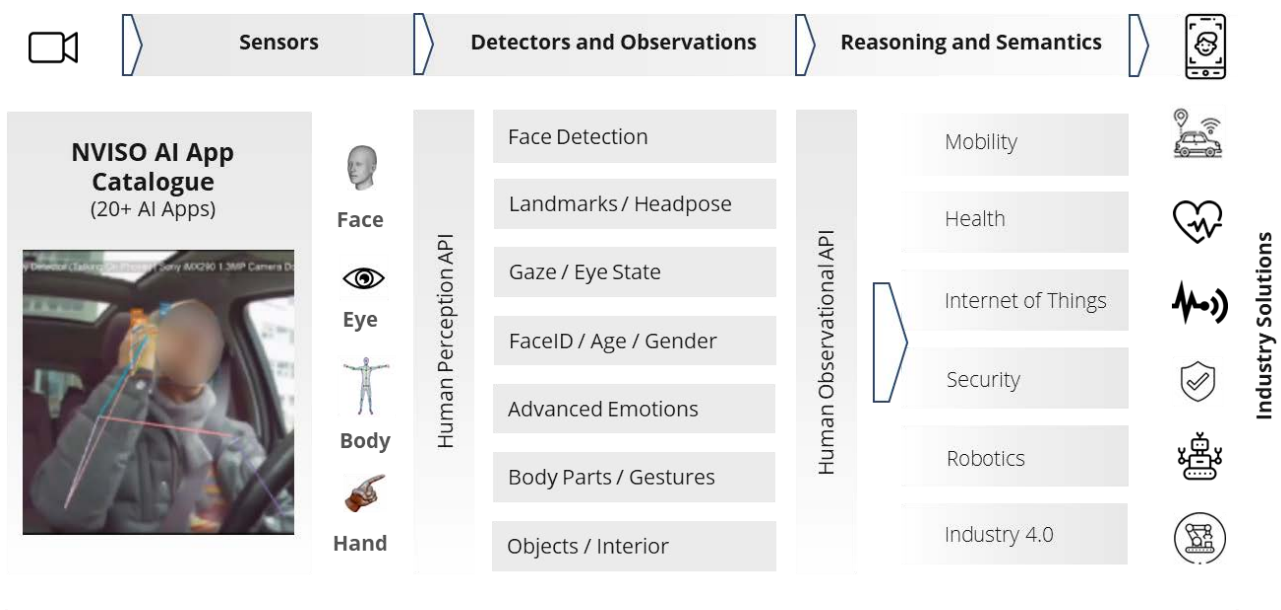


Human Behaviour AI

The task of understanding humans by using smart sensing has become increasingly important in modern society, given the proliferation and ubiquity of embedded sensors in a huge range of devices ranging from smartphones and smart city infrastructure to consumer IoT devices and health monitoring systems. A key to enabling the human user to connect with automated systems is artificial intelligence, which enables the plethora of sensing data to be processed in a way that gives context and meaningful insights to aid smart devices to make decisions and perform tasks.

Although there are a number of challenges in developing and integrating artificial intelligence into smart sensing and IoT applications (ranging from memory footprint and computational complexity, to privacy and robustness), the detection and understanding of human activities using artificial intelligence can be divided into three fundamental layers:

- **Sensors**, which are devices that detect and respond to changes in an environment in which humans are present. Inputs can come from a variety of sources such as light, temperature, motion and pressure (e.g. cameras, microphones, IMUs, etc) which generate data-streams that are used by data-driven and learning-based frameworks.
- **Detectors and Observations**, provide reusable core signals around well understood human modalities directly from sensor data using machine learning models which are not application or industry specific.
- **Reasoning and Semantics**, that apply rules or prior knowledge around simple or more complex human states or behaviours that may be context or application specific.



Human Behaviour AI from Sensor to Industry Solution


Detectors and Observations

Non-verbal “expressions” can communicate emotions faster, more subtly and more effectively than words ever can, which is why understanding non-verbal cues remains crucial for systems doing human behavioural analysis. Human detectors and observations can be broken down into 4 core areas defined by the face, eyes, body, and hands as follows:

Face

The face is widely regarded as a unique feature of human beings. Humans locate, identify, and distinguish between faces with very casual inspection and this outstanding human ability to identify individual human faces has long been of major interest to cognitive scientists, neuropsychologists, and neuroscientists. The recent advent of smaller, affordable and more powerful computing devices and larger and more diverse facial databases, along with the recent advancements in deep learning have enabled machines to emulate this human capability of identifying and distinguishing the human face often surpassing human level accuracy. Some of the most natural use cases where a machine can make better informed decisions based on the human face analysis are:


- **Identity:** The human face is the significant characteristic for identifying a person. Everyone has their own unique face, even in the case of twins. A face recognition system can find a person’s identity through a biometric method.
- **Expressions:** Human faces can express things that are difficult to put into words. Expressions can communicate emotions faster, more subtly and more effectively than words, which is why understanding facial expressions remains crucial for systems doing human behavioural analysis.
- **Movement:** Head movement (and associated head gestures) are another important part of non-verbal communication in human interaction. Classifications of head movement are based on action type, frequency, amplitude, continuity and other factors.

Face		The core observations that can be modelled and detected from the human faces are:
Core Observations	Location	Facial landmarks are used to localise and represent salient regions of the face, such as: eyes, eyebrows, nose, mouth, jawline.
	Position and Headpose	The rotation and translation vector for the face provides a precise head location and headpose.
	Expression	Universal expressions of emotion which include Anger, Disgust, Fear, Happiness, Neutral, Sadness and Surprise, and Contempt.
	Identity	The estimation for the gender and age along with a unique encoded sequence describing the identity for the face.

Eyes

Human eyes are one of the most important receptors of stimuli allowing humans to perceive their surroundings. In any social interaction involving humans, eyes can play a major role in establishing who we are, how we feel and what we do. Therefore the ability to perceive and analyse the eyes is a very important capability to have while interacting with humans. In addition, the direction in which eyes are pointed (often referred as Gaze) also provides very significant clues about human behaviour such as attentiveness, distraction, interest and suspicious behaviour. Some of the natural use cases for eye analysis in human-machine interactions are:


- **Identity:** High resolution images of the eyes can be analysed to detect and match the unique iris patterns providing a mode of identification.
- **Attention/Interest:** The eye gaze for humans often tells the point of attention for their visual stimuli. This gaze information combined with other technologies such as object detection can be used to deduce the attention/interest level of human subjects.
- **Cognitive Load:** Often an increase in cognitive load is related to voluntary eye movement such as fixations as well as involuntary eye movements like blinking and pupil dilation. Eye tracking and analysis can be used to detect such behaviours.

Eye		The core observations that can be modelled and detected from the human eyes are:
Core Observations	Eye State	The state of the eyes visible or non-visible including if the eyes are open, closed, or partially open.
	Openness	The level of the openness of the eyes which is measured as a numeric score.
	Blink Rate	The number of blinks per minute.
	Eyelid Closure	The percentage of time during which the eye is closed for a given unit of time. It is one of the most common indicators for fatigue monitoring.
	Gaze	The direction in which eyes are pointed.
	Iris Dilation	The measurement for the dilation of Iris in the pupil.

Body

A machine can better interact and help users if it can understand human body poses, actions and activities of people. For example, a robot can take timely actions when it detects the pose of a person who is prone to fall. Therefore it is very crucial to model and teach computers to understand and interpret the human body pose with a high level of accuracy. With the latest developments on collecting accurate human body pose datasets along with the sophisticated deep learning architectures, it is now possible for machines to accurately detect and analyse human body pose. Some of the natural use cases for body analysis in human-machine interaction are:


- **Joint Location and Shape:** The locations of the body joints for each key point are represented relative to input data in 2D space (i.e., from an image or video frame) so as to understand body location and shape.
- **Position and Movement:** An intuitive human body model that includes a set of 3D joint positions and limb orientations to represent the human body structure and the relations between different body parts for biomechanical analysis.
- **Body Posture:** Classification of static human body postures (sitting, lying, etc).
- **Activity Recognition:** Classification of dynamic human activities (jumping, walking, etc).
- **Volumetric and Weight Analysis:** With the collection of rich datasets having the body volume, weight and pose information, along with advanced deep learning architectures it is possible to detect the salient body features which highly correlate with the body weight and volume. This can be used to adjust machines to provide more customised support to its users.

Body			The core observations that can be modelled and detected from the human body are:
Core Observations	Location		Key points on the body which can be used to encode the visible body locations.
	Position and Angles		Set of body joint positions and angles which are needed to perform a specific pose.
	Posture		Classification of static human body postures (sitting, lying, etc).
	Activity		Classification of dynamic human activities (jumping, walking, etc).
	Volumetric and Weight		Estimation of body weight and volume.

Hands

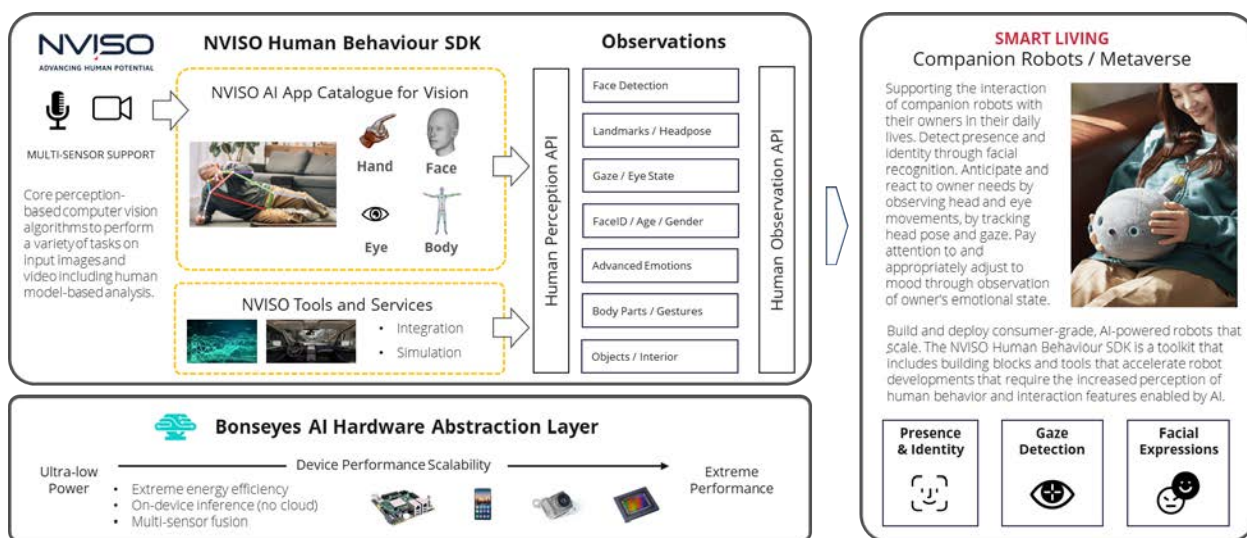
Gesturing is a natural and intuitive way to interact with people and the environment. So it makes perfect sense to use hand gestures as a method of human-machine interaction. Snap your fingers and make your coffee maker brew you a fresh cup. Wave a hand near your smart TV and switch on today's weather forecast. Tap a finger near your smartwatch and set an alarm in your child's bedroom. How great would it be to get things done just by gesturing? It's not that unrealistic anymore: hand tracking and gesture recognition technologies are penetrating multiple industries. Some of the natural use cases for hand analysis in human machine interactions are:

- **Joint Location and Shape:** The “shape” of the hand only considers the locations of the hand and finger joints which are represented relative to input data in 2D space (i.e., image or video frame); it does not consider a hand's orientation in space.
- **Position and Movement:** An intuitive human hand model that includes a set of 3D joint positions and hand and arm orientations used for posture and gesture recognition.
- **Hand Posture:** Classification of static human hand postures (stop, thumbs up, etc).
- **Gesture Recognition:** Classification of dynamic hand and arm movements to infer the underlying gesture performed by the person.

Hand		The core observations that can be modelled and detected from the human hand are:
Core Observations	Location	Key points on the hand which can be used to encode the visible hand locations.
	Position and Angles	Set of hand joint positions and angles which are needed to perform a specific pose.
	Posture	Classification of static human hand postures (stop, thumbs up, etc).
	Gesture	Classification of dynamic hand and arm movements to infer the underlying gesture performed by the person.

Reasoning and Semantics

Human behaviour signals provide important information by themselves useful for machines or computers to detect and understand humans. Additional insights however can be extracted by looking at how the core observations change over time or correlate together and can be used with context and/or application specific reasoning and semantics for industry specific solutions. Core behaviour signals generated from AI Apps (which carry out basic human detections and observations as previously explained) are then combined at an application level with reasoning and semantics to enable improved or new human-machine interactions.



This section outlines some examples (non-exhaustive) of how detections and observations can be combined with reasoning and semantics to create extra layers of interpretation and application value.

- Activities:** For various reasons, the machines we are building would like to know the activities in which the persons in its field of view are involved. The activities of interest vary widely with the application and the context: from activities which might distract the driver from the road in the automotive context (such as eating, drinking, smoking, speaking, sleeping) to activities any person is involved in when close to a robot's viewpoint (sitting, standing, talking to somebody, walking etc.). Apart from the obvious use case of providing triggers to real time decision making, identifying activities is also important for filtering various core signals to increase their robustness and relevancy.
- Gestures:** Vision-based AI opens up a whole new level of human-computer interaction. Forget about touch-based inputs like keyboards or touch screens and welcome vision-based gesture control for interacting with the new generation of intelligent machines, be they cars, drones, social robots, wall displays or industrial machinery. Further, without the need for contact-based sensors, vision-based gesture recognition has unchallenged advantages: longer distance, no wearable or additional sensors needed, environmental independence.

- **Distraction, Inattentiveness, and Drowsiness:** Distraction is often related to an object, a person, an idea or an event that diverts someone's attention from a task. Distraction-detection systems use head pose or gaze information to detect if someone is paying sufficient attention to a task. Drowsiness refers to sleepiness, often in inappropriate situations. Although the state of drowsiness may only last for a few minutes, its consequences can be disastrous. Someone does not become drowsy suddenly, without showing some signs. Examples of such signs include difficulty keeping eyes open, yawning, frequent blinking, difficulty concentrating, and nodding. These signs gradually become more apparent as drowsiness deepens and, as such, can serve as indicators for the level of drowsiness.
- **Incapacitation:** Incapacitation can be caused as a result of physical or mental exertion or a prolonged period of performing the same task. The eye states (open, closed, partially closed) combined with the technologies such as expressions and head movement can be used to decipher the level of inability for performing tasks such as driving etc. which can be very beneficial for keeping safety and security at the forefront.
- **Fatigue:** Fatigue is defined as a global reduction in physical or mental arousal that results in a performance deficit and a reduced capacity of performing a task. When a person enters a state of fatigue, two things happen: blinking will unconsciously increase for protection, or the eyes will become sluggish and blurry and the number of blinks will change. The increase or decrease in the number of blinks compared to the state of consciousness can reflect the degree of fatigue.
- **Stress:** Usually there is a visible change in the expressions and eye state of humans upon being exposed to increasing levels of stress. The eye blinking frequency increases as well as the presence of the main facial features related to anger and fear. Combining these two measurements can provide the level of stress.
- **Dangerous and Suspicious Behaviours:** Based on the context, certain behaviours can be classified as dangerous ((like taking the gaze away from the road while driving) or suspicious (like a vehicle cruising the streets repeatedly). By combining the core signals like object detection, body and face analysis it is possible to accurately identify the surroundings and flag the behaviours and activities that are deemed dangerous or suspicious for those environments.
- **Advanced Emotions:** The trigger for a transient emotional response may be identified. Is there another person or object which makes us smile? The head orientation or gaze vector when an emotional response is registered will, in most cases, indicate the direction of stimulus. If another person or object is detected in the direction of attention when the emotional response is registered, it can be postulated that that is the source of the emotional response.

Use Cases

In this section, a number of applications involving the adoption of Human Behaviour on Extreme Edge AI Computing are illustrated and analysed with the goal to highlight the considerable breadth of scenarios where NVISO technology can play an important role.

Smart Health

Data will play an increasingly important role in providing a better understanding of consumer needs in terms of health, and to enhance and tailor a more cost-efficient health offering that delivers the right care at the right time and in the right place. The interconnections made possible by being able to access pools of data not previously available (worldwide databases, data clouds, apps, in-sensor computing etc) are creating a major shift in healthcare provision.

According to different projections, healthcare budgets around the world are expected to increase by 10% in aggregate by 2030. Healthcare spending will be driven by ageing and growing populations, rising labour costs, and also by clinical and technology advances. Consequently, by 2030 healthcare is expected to be centered on patients being empowered to prevent disease rather than seek treatment.

Healthcare budgets in Europe will therefore shift towards novel areas such as digital health and more advanced prevention and rehabilitation options, for which homecare will play a key role. Money is expected to be redirected toward personalised medicine for the most complex diseases, and preventive, early stage treatments. This split is likely to lead to significant changes and require new R&D strategies for many industry players.

- **Vital sign monitoring** based on non-invasive observational sensors will be an important component in many medical applications. However, a cloud-based implementation of the sensing would be too slow in time critical contexts. This is not the only problem of cloud systems as storing generated data in them is also a privacy concern. Issues of latency and privacy can be solved by using edge AI.
- **Personalised medicine** is possible as human physiology can vary greatly from individual to individual. Examples for that include blood pressure or lung capacity. However, these differences need to be considered for accurate medical applications like vital sign monitoring. Due to privacy concerns, it is difficult to process this information in the cloud-based solutions. Edge AI offers the possibility of maintaining privacy when processing medical data. Furthermore, many medical applications require real time processing, which can be better realised with local AI.
- **Sports analytics** can help analyse lower-limb injuries are common among athletes, accounting for 77% of hospitalised sport-related injuries, and are a risk factor for early-onset osteoarthritis. High-impact forces are one of the factors contributing to lower-limb injuries. To decrease the prevalence of lower-limb injuries, and their associated long-term disability and economic burden, multiple injury prevention programs have been proposed. These take into account the study of ground reaction forces (GRFs) in order to enhance athletes' performance, determine injury-related factors, and evaluate rehabilitation programs' outcomes.

Smart Living

Rethinking human activities to take advantage of the innovation opportunities offered by hyper-connectivity and AI solutions and new kinds of sensors based on miniaturised technologies will create numerous benefits for every new market, ranging from connected cars and digital health to smart home and smart living, and factories of the future. This should include lessons learned from the COVID-19 pandemic like the sudden increase in remote-working. Smart Living is a solution that aims to make an environment of the future that improves people's quality of life and examples include:

- **Companion robots** support the interaction of companion robots with their owners in their daily lives. Detect presence and identity through facial recognition. Anticipate and react to owner needs by observing head and eye movements, by tracking head pose and gaze. Pay attention to and appropriately adjust to mood through observation of the owner's emotional state.
- **Metaverse** will use augmented reality (AR) and virtual reality (VR) along with artificial intelligence and blockchains to create this virtual world. One of the most interesting and spoken about concepts of the Metaverse is Avatars. People are creative and love the thought of designing themselves in a virtual world. They can change their hair colour, style of clothing to their preference. AI has the ability to analyse 2D user images or 3D scans to create very realistic and accurate Avatars. Companies such as Ready Player Me have already actively been using AI to help build Avatars for the Metaverse.
- **Affective computing** is interested in automatically detecting and recognizing the emotional state of a human, either with remote or "nearable" sensors (visible and IR imagery, audio, physiology), or with sensors in contact (wearables) for physiology, or activity monitoring. Emotions, a classic conceptual representation of which follows a 2D valence (negative / positive) versus intensity (calm / excited) pattern, have an essential role in human behaviour. These influence the mechanisms of perception, attention, decision making, and social behaviour. The purpose of estimating emotional states is to improve understanding of human behaviour. This is the strongest reason as emotional states are both very personal and evolving, very different from one individual to another, and from one situation to another. Edge AI allows for maintaining the confidentiality of the data inside the measurement device, to guarantee the autonomy of the devices, and to aim for an individual estimator learning over time.

Smart Mobility

Mobility is a basic human need and Europe's mobility industry is a key contributor to it. The usage of smart perception, safety and automated mobility solutions and services to provide safe and comfortable inclusive mobility that is also suitable for the elderly as well as people with special needs. Research, development and innovation (R&D&I) of embedded AI-based software, sensors and electronic components and systems provide the core of automated on- and off-road vehicles, ships, trains and airplanes. A special focus requires validation of the safety and reliability of the automated mobility systems in all traffic and environmental situations as there are currently no adequate methods and tools available.

- **Occupant monitoring** is a key feature expected to become a standard feature in new cars as a result of regulatory and rating agency requirements Euro NCAP 2025 Roadmap. Advanced driver monitoring systems (DMS) can detect distracted and drowsy drivers by accurately measuring eye and head position, attention, and fatigue. The DMS alerts the driver and integrated safety systems upon detection of a risk such as drowsiness or distraction. This feedback enables the driver and vehicle to take action before safety is compromised. Managing transitions between different levels of autonomy is fundamental. The AI-based observer is a key point of this system as it detects the behaviour and the mental state of the driver. Edge AI offers the local calculation of the driver states, thus allowing for control of the response time thus preventing personal data from leaving the vehicle and continuous learning to adapt the AI-based observer to each driver and passenger.

Robotics

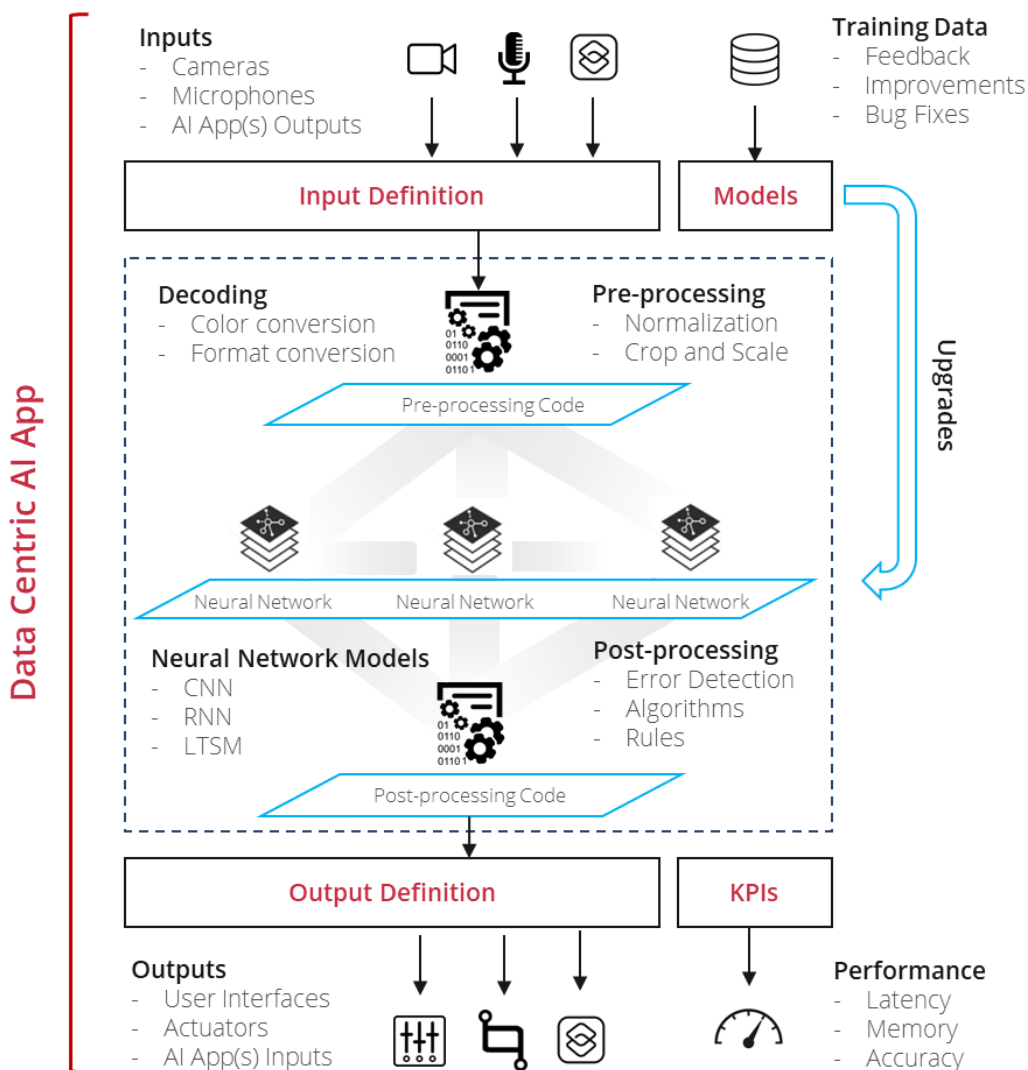
Industry 4.0 changes to the mode of operation have a profound impact on how the factories, construction zones and processes are managed and operated. Powerful networked digital tools are needed to achieve the necessary Situational Awareness and control of autonomous vehicles, robots and processes at various autonomy levels.

- **Collaborative robots** in industrial environments support human workforce in the fulfillment of repetitive jobs or heavy lifting, for instance. Applications can be found mainly in the manufacturing industry, e.g. assembly of automotive parts. Edge AI enables new possibilities for the cooperation of humans and robots, because in contrast to cloud based systems edge AI is fast enough to handle situations where the robot could inflict harm. To implement these new possibilities, sensors need to be deployed that are able to monitor the environment and the movement of humans and animals within range. The data from these sensors are locally processed by AI models in the robot or running on nearby edge nodes. Afterwards, edge AI-based components use the processed data to control the robot, allowing for close cooperation with humans in performing complex tasks like the manufacturing of custom products in workshops or rescue operations. To implement this vision of close cooperation many challenges need to be solved such as training the robots for new tasks.
- **Autonomous robotics** are advancing rapidly and becoming robust and reliable in a number of fields, but there remains much to be done before intelligent robots can work together with us as assistants. Today, most robots are either tele-operated or perform precisely defined missions. In many cases the human overhead required to use a robot exceeds the usefulness of having it. There is a clear need for human-robot teams in which humans and robots work side-by-side. Homogeneous teams involve team members doing the same work. In such a case, the value comes from parallelism, such as picking grapes from vines.

Data Centric AI Apps

What are AI Apps?

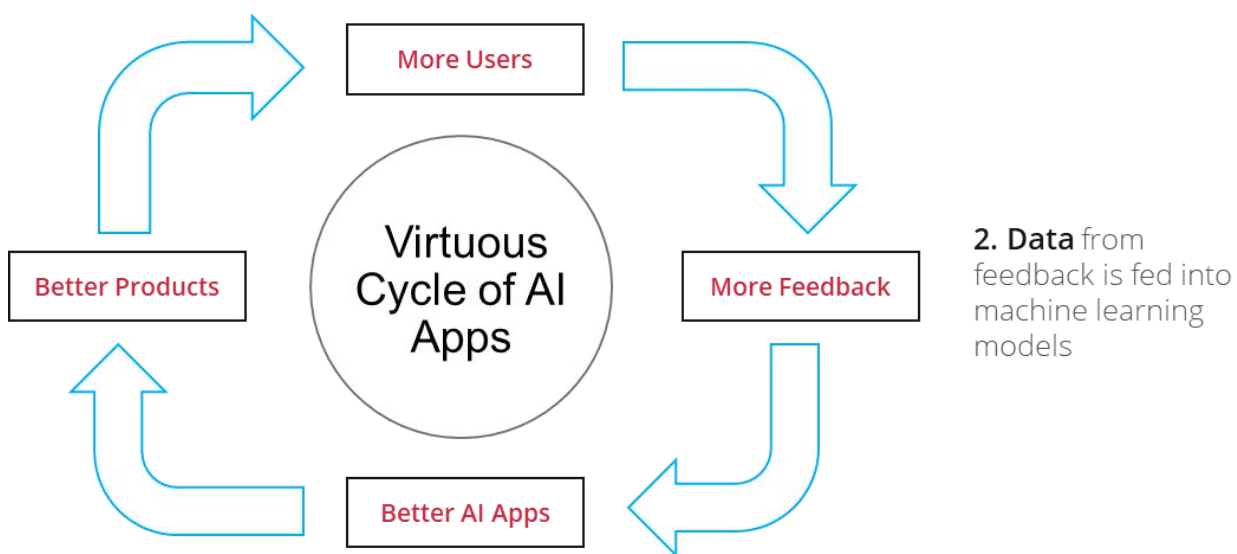
An Artificial Intelligence Application (AI App) is the entire software including any machine learning component(s) of a data-driven algorithm necessary to transform data input(s) into output(s). An AI App is represented by its architecture that shows how to transform its input(s) into its output(s) where the transformation is given in the form of learning algorithm(s) (normally specified in frameworks such as Pytorch, Tensorflow, etc.) including pre and post processing algorithms and code. Data Centric AI is the discipline of systematically engineering data used to build and continuously improve an AI system over time.



What is the virtuous cycle of AI Apps?

The virtuous cycle of AI Apps, also called the "AI Flywheel Effect," is one of the most exciting ideas in Artificial Intelligence and it's also incredibly simple. Essentially, when AI technologies are integrated with a product properly, they create a feedback loop where the product continuously improves with use, generating more usage and a better competitive position relative to other products.

1. **Product** get used, generates data and feedback



3. **AI Apps** improve the products, generating more usage

Although any product tends to improve with usage regardless of its underlying technology because a good team will use qualitative feedback and analytics data to bring it closer in line with user needs. This improvement, though, tends to reach an asymptote where additional usage and data no longer provide much marginal insight to the product.



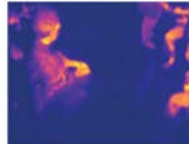
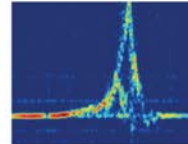

When a product's core technology is AI-driven, though, it adds another layer on top of the typical team-driven product improvement cycle. With today's AI technologies (for example Deep Learning), additional data continues to provide marginal improvement for a very long time, allowing the cycle to continue much further than it previously could and improving the core product for a long time.

What sensors can be used as input with AI Apps?

Humans observe the world in real-time through a combination of different modalities (vision, hearing, touch, taste, and smell) and can respond to a range of physical stimuli such as light, sound, gravity and acceleration, chemical substance, position, motion, and temperature. Machines, on the other hand, interpret the world through data that algorithms can be taught to process. To sense the presence of humans, a variety of sensors are used ranging from cameras to detect colour, depth, and temperature information. Radar and Wifi to detect distance and velocity. Microphones to detect sounds and speech. Each sensor type has specific characteristics (see below), and usage varies on use case and application requirements and constraints. However sensors for the detection of humans generally falls into three broad categories:

- Analysis of **visual information** (e.g. camera, depth, thermal). Over the past few years, quality mobile cameras have proliferated in devices ranging from smartphones, surveillance devices, and robotic vehicles, including autonomous and connected cars.
- Analysis of **audio information** (e.g. microphone). Conversation is natural and that's why it is a primary interface for human-machine interaction. Voice-based personal assistants (VPAs) are growing in popularity in smartphones, smart speakers, smart watches, wireless earbuds, cars, smart TVs and their remote controls.
- Analysis of **time-series information** (e.g. position and motion). Inertial Measurement Units (IMUs) within everyday consumer wearables (e.g. smartwatch or phone) offer a convenient and low-cost method to monitor the natural behaviour of people.
- Analysis of **meta-data information** (e.g. AI App outputs). AI Apps themselves can be used as "sensors" and can generate additional data that can be fed into AI Apps for further processing.

Human Behaviour AI can use one or all categories of data from such sensors.

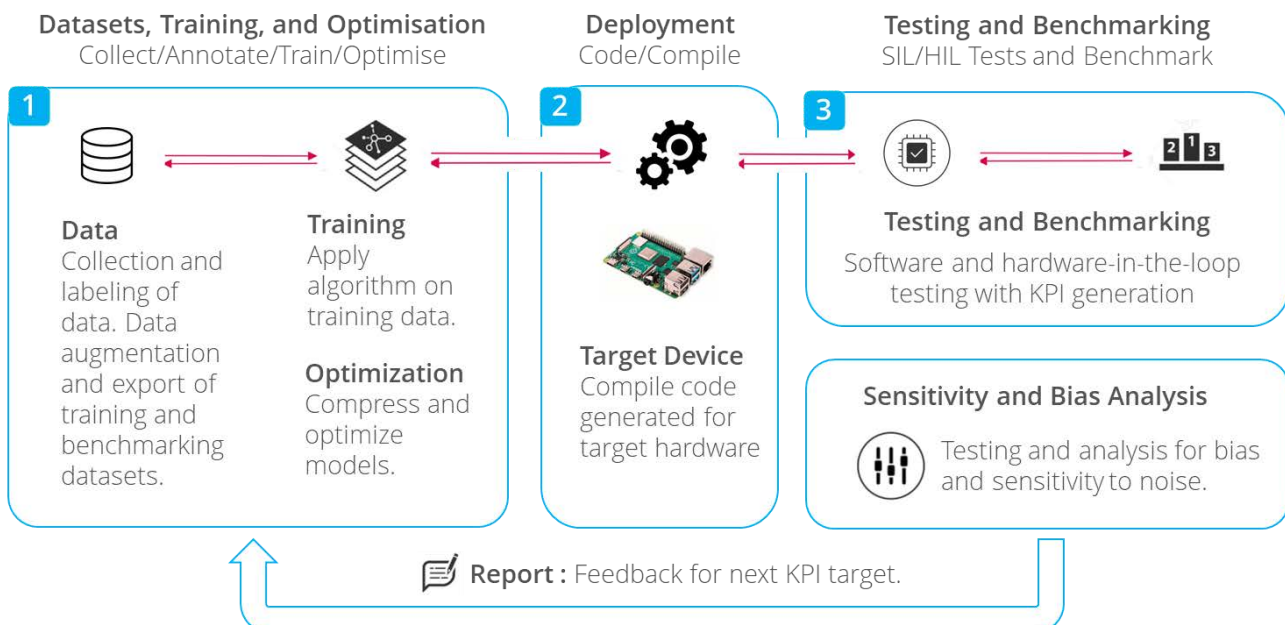
	Camera	Depth sensor	Thermal sensor	Radio sensor	Acoustic sensor
Sensory information	RGB, colour, video	Lidar	Infrared	Radar, Wi-Fi	Microphone
Function	Measures colour (visible light)	Measures distance to objects	Measures surface temperature	Estimates distance and velocity	Measures air pressure waves (sound)
Sampling rate	30 Hz (1,920 × 1,080)	30 Hz (1,280 × 720)	10 Hz (640 × 480)	800 Hz	44.1 kHz
Bit depth	24 bits	16 bits	16 bits	32 bits	16 bits
Uses	Object recognition, person detection	3D object detection, robotic navigation	Night vision, equipment safety	Motion detection, object detection	Speech recognition, event detection
Data visualization					

From: [Illuminating the dark spaces of healthcare with ambient intelligence](#)

How are AI Apps developed and tested?

To develop AI Apps for the Edge is far more demanding due to the limitations in terms of computation and storage abilities when compared to cloud-based systems. Additionally Edge systems have a wide range of hardware and software stacks which demands rigorous evaluation and testing to ensure correct performance with onboard sensors. While optimal model training and inference on resource-scarce devices are still a debated problem throughout academia and industry, there are four main steps in development and testing of AI Apps (although the actual process can be a more complicated depending a which product development life cycle model is employed during development e.g. for medical or automotive products):

- **Datasets, Training, and Optimisation:** Given the importance of data to data centric AI, it is crucial to ensure the quality and quantity of data that is used as the starting point of any AI App development. Once a training dataset(s) are ready, machine learning algorithm(s) are applied to make the neural network model understand the relationship between the data and the target objective. For edge-based systems the trained models are further compressed and optimised for edge constraints such as memory, latency, and power efficiency.
- **Deployment:** Deployment is the step of porting the AI App from the training environment to the target system, e.g., a low-power microcontroller, using code generation and cross-compilation.
- **Testing and Benchmarking:** The process of testing the trained and optimised model to make sure that it has been well trained and can generalise to unseen data on a target software and hardware platform with software-in-the-loop (SIL) and hardware-in-the-loop (HIL) testing. Verification on the target system of the expected performance is obtained with onboard sensors.



Datasets, Training, and Optimization

Data is the fundamental part of any AI App development workflow. Data refers to representative examples or cases from the target domain that characterises the problem to be solved. In AI modelling, a dataset is composed of many samples that are fed to an AI model, which will infer or predict a target element, e.g., classification, regression, etc. A dataset is an unstructured collection of data generally associated with a set of problems to solve. Such a dataset needs to be structured by means of a Datatool, which provides meaningful and reusable structured samples to train an AI model for such a problem.

Despite having such importance, machines do not understand what data represents. They do not understand why 'a' is 'a' and why it is written in this way or why 'this' means what it means. Most of us do not understand the food that we eat. The only thing that we know is that we have to eat, and we do so. Data for machine learning is food. The machine just consumes it and then learns the relations between different data to predict something rather than understanding the data itself.

Given the importance of data to AI modelling, it is crucial to ensure the quality and quantity of data just as we humans ensure this for our food intake. The process of aiming to prepare and provide the required data to the modelling process is often termed as "Data Engineering" with different steps of this process often together called a "Data Pipeline". Some of the different steps in a typical data pipeline are as follows:

- **Data Identification:** is the process of identifying the right data for solving a specific problem. For example, to build an object detector, it is necessary to have an image-based dataset with images having one or more objects inside them.
- **Data Acquisition and Annotation:** is the process of acquiring the required data once data identification has been completed and adding labels to the data elements by reviewing the elements with the help of problem specific domain experts and tools.
- **Data Harmonisation:** Often due to a wide variety of data collection sources and protocols, the collected datasets have substantial differences in their structure, which makes it difficult to use them in model training. Data harmonisation, by means of a Datatool, is the process of transforming and structuring the collected dataset so that it is easier to analyse, and consume by the subsequent model training pipeline.
- **Data Validation and Verification:** Once data has been annotated and transformed, it is very important to ensure the correctness and quality of the dataset. The process of data validation and verification applies multiple methods and approaches to ensure that data used for model training is both correct and accurate, as incorrect or inaccurate data is of little use.
- **Data Provisioning and Versioning:** is the process of making the dataset available for model building and evaluation such that the data from the dataset can easily be integrated into model training pipelines.

Training is the process of applying the machine learning algorithm on training data to make the AI model(s) understand the relationship between the data and the target objective. In a nutshell, during model training data elements are fed to the AI model iteratively in a random order while

making small adjustments to the AI model itself. This iterative process is called “Model fitting”, which generally lasts until the model learns to correctly predict the right outcome for the vast majority of the training data elements. The different steps in the model training are:

- **Data Preparation:** is the process of preparing the dataset's before they are fed to the AI model. This includes partitioning the dataset into 2 or more subsets, often called “Training, Validation and Test” sets. Another major step in data preparation is data balancing which is essentially the process for ensuring that one segment of the data population is not over represented compared to any other segment. Overall data preparation is crucial to achieving a fairly trained model that generalises well when applied to unseen data examples.
- **Algorithm Selection:** is the process of identifying the right model architecture(s) and algorithm(s) to solve a specific problem. Selection of the right model architecture(s) depends on the specific problem statement as well as the associated data to be used for the model training. If not selected correctly, the wrong architecture can severely affect the model performance when applied on unseen data examples.
- **Model Fitting:** is the iterative process of feeding the data examples to the AI model in a random order while making small adjustments to the model itself based on its current prediction accuracy. The process generally lasts until the model learns to correctly predict the outcome for the vast majority of the data examples.
- **Model Selection:** In general when performing model training, there are many parameters which can be altered and whose value can influence the final performance of the model. The set of such parameters is called “Model Hyperparameters”. Model selection is the process of trying different sets of hyperparameters to train models and selecting the model having the best performance given the problem. Model selection can also have a great impact in the latency and memory consumption on the target device.
- **Model Evaluation:** Evaluation is the process of testing the model to make sure that it has been well trained and can generalise to unseen data. Model evaluation is performed on data that has not been used during training, either a different split of the same dataset or a different dataset (on the same AI task). The difference between the trained error (training) and the test error (evaluation) tells the fitness of the model. Underfitting or overfitting are some of the problems that can be found through evaluating the model on unseen data.

Trained models for cloud-based systems are usually overparameterized, i.e., they are larger than needed, or exhibit redundancy. When ample system resources are available, this is not an issue, however for Edge based systems it is possible to compress and optimise models without any drop in accuracy to reduce memory, latency, and power consumption. Some of the well-known techniques considered for model compression and optimisation include:

- **Compact Model Design:** The Model architecture has a large influence in the latency and memory consumption during the Model's execution on the device. Hence, searching for an efficient and compact design is very beneficial to reduce the Model's requirements on resource-constrained devices.

- **Knowledge Distillation:** Whereby a small (easy to implement) model (student) is trained to behave like a larger trained neural network (teacher) while trying to preserve the accuracy of the teacher model, thus enabling the deployment of such models on small devices.
- **Compression:** Steps such as quantization, sparsification, dimensionality reduction, tensor decomposition, pruning, components sharing, etc. These methods exploit the inherent sparsity structure of gradients and weights to reduce the memory and channel occupation as much as possible, while improving the arithmetic intensity and efficiency.
- **Hardware-in-the-Loop Co-Design:** Devices exhibit large differences in terms of hardware architecture, computing processor or memory. Hardware-in-the-loop Co-Design allows the adaptation of the Model to the target environment by feeding back to the design environment metrics such as the latency and memory consumption on the device.

Deployment

Deployment is the step of porting the trained model from the training environment to the target device, e.g., a low-power microcontroller including the implementation of all necessary pre and post-processing steps of the AI App. This is a complex step as it involves several data and code conversions as well as low-level optimisations. Some of the steps during deployment are the following:

- **Device Tooling Setup and Configuration:** Unlike cloud-systems, low-power devices are very heterogeneous and come in many different forms. Setting up a device is a difficult step that requires very detailed and specific knowledge about hardware and software. Thus, easy-to-follow workflows, i.e., Developer Board Packages, containing the full software stack (operating system, drivers, middleware components, etc.) and documentation required to procure, set up and control target hardware are of extremely importance to the users.
- **Code Generation:** Models are usually trained in Python (high-level language) that allows an easy design process. However, Python is not suitable and optimised for low-power microcontrollers. Hence, a code generation step is needed to convert the representation of the model to a low-level language (C-code) that is understandable by the target device. Code generation also includes the implementation of necessary pre and post-processing steps.
- **Data Format Conversion:** Models may be organised in different data formats or layouts, exploiting the inherent structure of the data. Likewise, data can be represented by a different number of bits (Quantisation), trading precision with energy consumption. During deployment, certain design options need to be taken that optimise the deployment of the model on the target device.
- **Memory Allocations:** Memory Allocation process reserves memory space to the input/output tensors and the weights of the neural network. Since memory is very scarce in low-power microcontrollers, an efficient memory allocation strategy is fundamental to allow the deployment of large Models on these resource-constrained devices.

Testing and Benchmarking

Once an AI App is deployed to a target device testing and benchmarking needs to be carried out across the deployment chain to assess the performance of the model after optimisations and

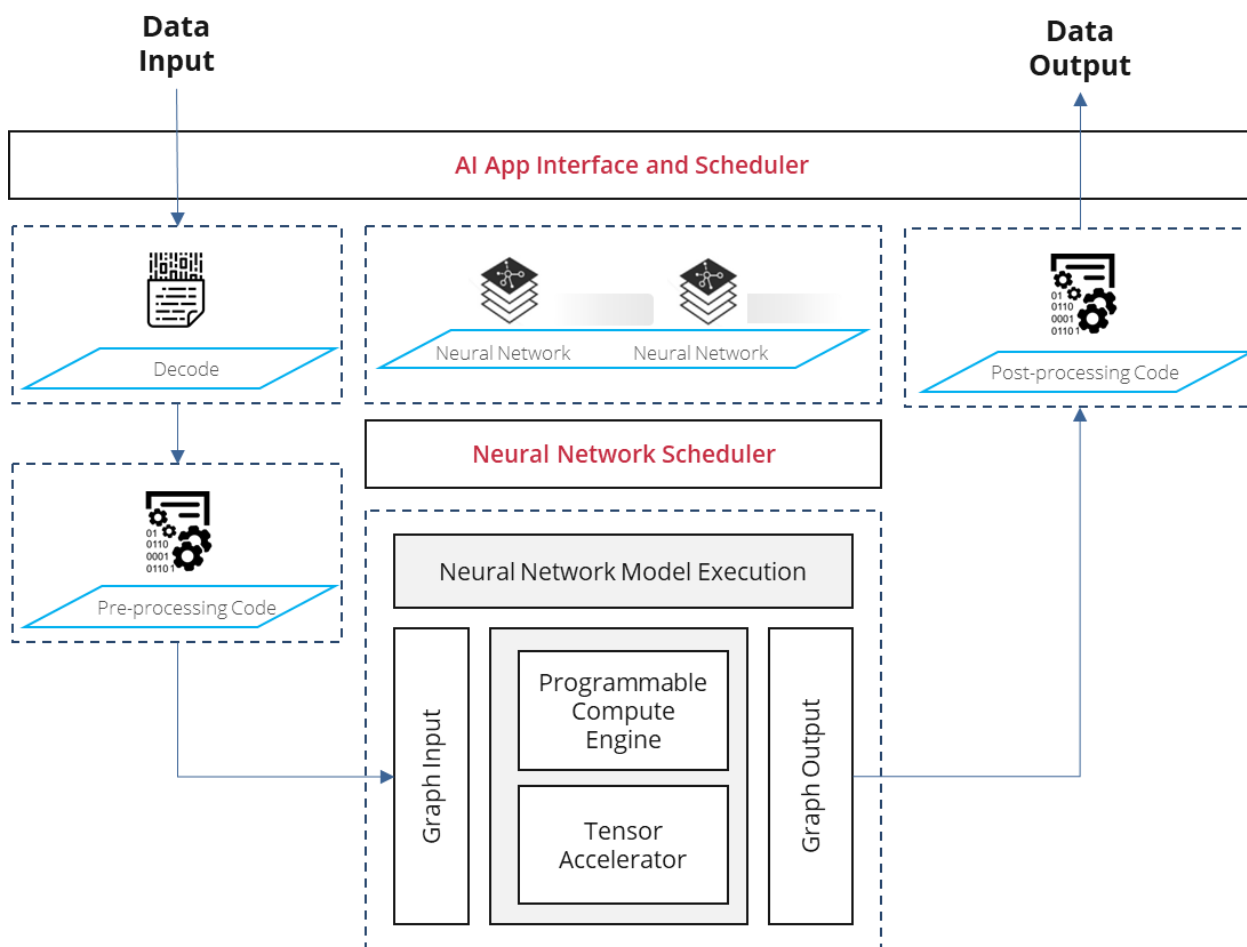
deployment on target devices which also includes evaluation of pre- and post-processing steps integrated within the AI App interface. Testing and benchmarking is performance under industry standard [Software-in-the-Loop \(SIL\)](#) and [Hardware-in-the-loop \(HIL\)](#) testing protocols:

- **Software-in-the-Loop (SIL) Testing:** SIL testing assesses that the AI App still performs correctly after porting from the training environment and after the optimisation and code generation process for embedded devices. As the model(s) inference results may vary as different software libraries are used for model inference, it is therefore needed to perform careful evaluations to make sure no drop in performance occurs due to the complete software environment and inputs.
- **Hardware-in-the-loop (HIL) Testing:** HIL testing is performed to the AI App when it has been deployed on the target hardware environment with inputs from sensors and running on the target environment within a simulation environment. HIL testing entails simulating device and environmental inputs for the hardware under test, causing it to believe that it is reacting to real-world conditions in the real-world. The HIL bench contains all of the relevant device components. A simulator presents inputs to actual cameras and sensor systems, which in turn send signals to the system under test to see whether it responds correctly to the inputs. HIL testing confirmation of the performance at this stage will guarantee the correct functioning of the AI App for the target application meeting memory, latency, and power efficiency requirements.

What is the typical data flow of an AI App?

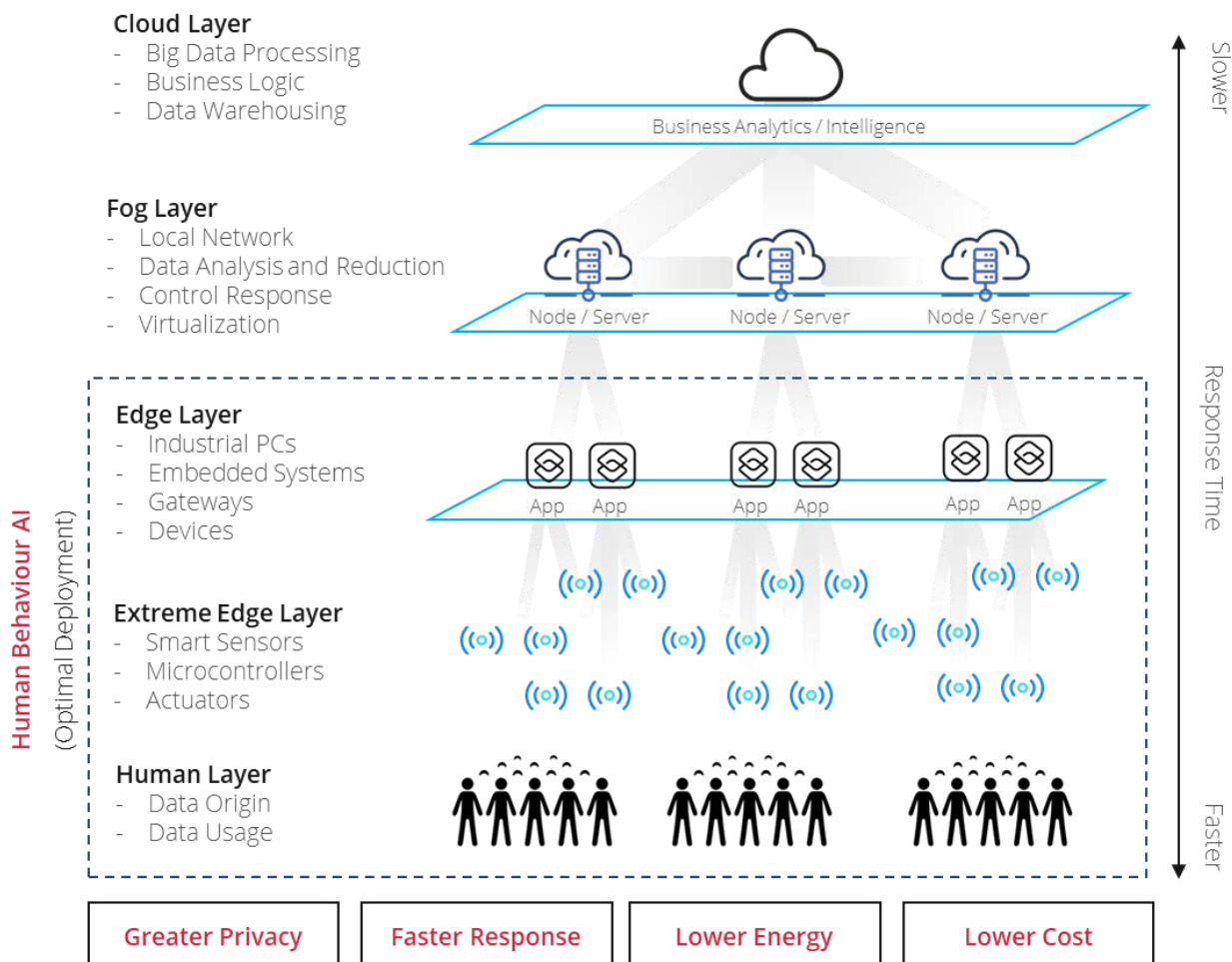
The training of AI Apps require training datasets that often require a series of steps to be executed before and after the neural network models which must be performed identically at inference in order to obtain accurate predictions. These data flow steps can be summarized as follows:

- **Decode:** Input data from image sensors must be decoded and converted into a format which is compatible with the AI App. e.g. conversion of RGB input data from videos to a single monochrome channel.
- **Pre-processing:** Model training can be more efficient if certain pre-processing steps are performed on the network inputs. These pre-processing steps must be replicated exactly at execution, as with the same decoding step.
- **Neural Network Model(s):** A set of parameters, known as weights, and mathematical operators described as a neural network graph determine the type of model used in the AI App. Common types include Deep Neural Networks (DNN), Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Long/Short-Term Memory Networks (LSTM).
- **Post-processing:** Additional algorithms and rules can be executed on the model outputs to improve and create additional predictions.



Extreme Edge AI Computing

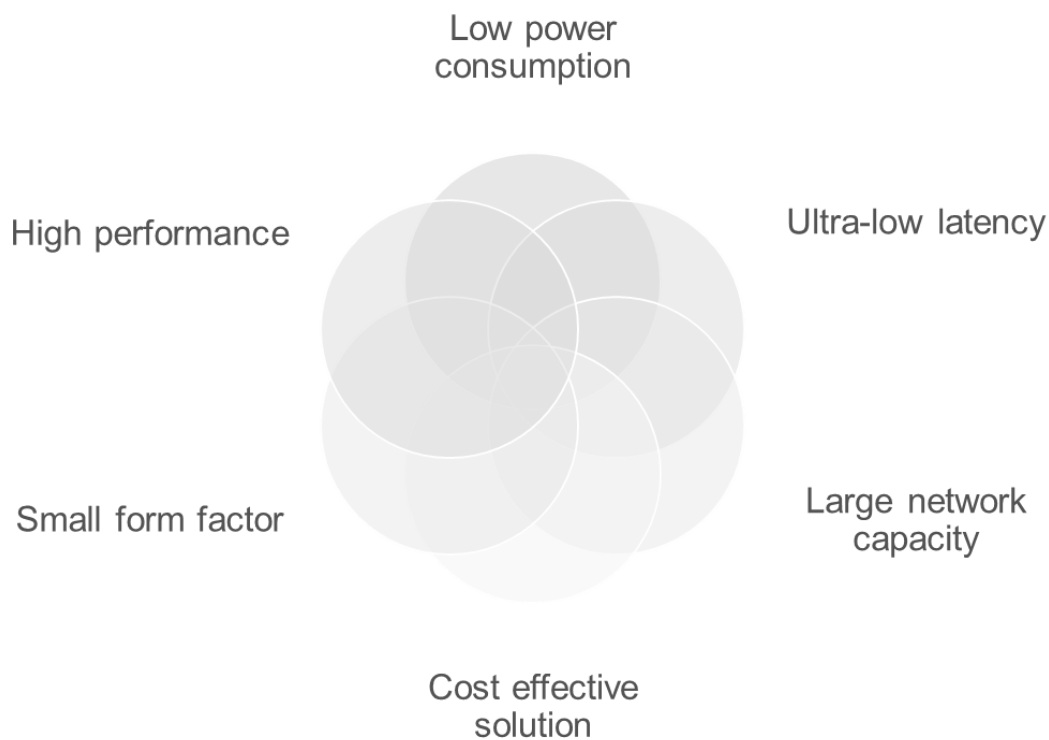
With the explosive growth in the adoption of Artificial Intelligence (AI) to address a large range of problems that were deemed very hard, data-intensive applications have placed a high demand on hardware performance, in terms of short access latency, high capacity, large bandwidth, low cost, and ability to execute artificial intelligence (AI) tasks. So far, the majority of attention and investment has been directed towards “Cloud AI”, with “Data” being the largest common denominator in creating value for industries, governments, and individuals’ lives. However, the quest for intelligence is fast becoming a prominent and essential feature at the “Edge” as well, where trillions of “things” will combine to generate even more data. Given the severe constraints that govern edge devices in terms of efficiency, footprint, robustness and cost, it is self-evident that bringing true intelligence to the edge will require profound innovation at all levels of the stack from the computational concepts all the way down the implementation technology.



What is “AI at the Extreme Edge”

Edge AI means that AI software algorithms are processed locally on a hardware device. The algorithms are using data (sensor data or signals) that are created on the device. A device using Edge AI software does not need to be connected in order to work properly, it can process data and take decisions independently without an internet connection. Edge AI applications need to be low-cost, small-form-factor devices with low latency, high performance, and low power.

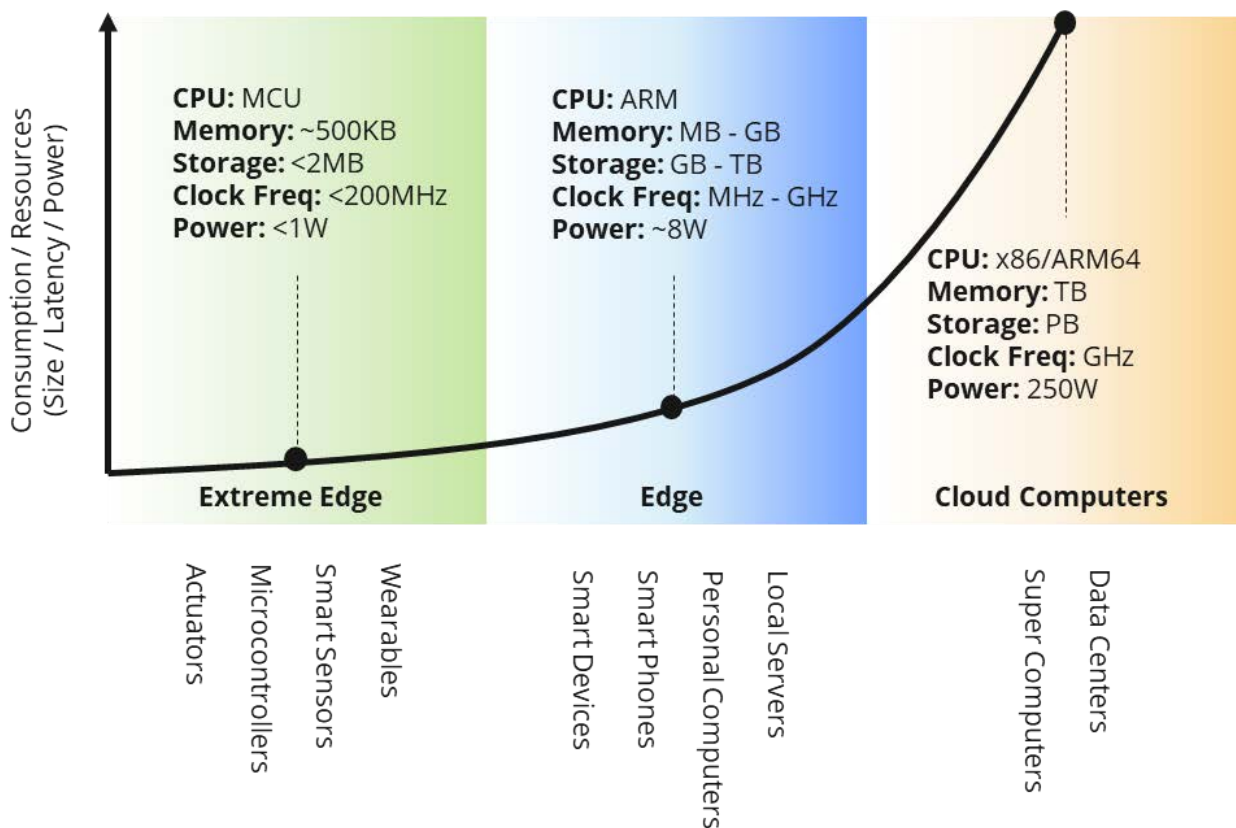
The Challenges of Meeting Extreme Edge AI Requirements



Edge computing is emerging as a strong alternative to traditional cloud computing, enabling new types of applications with the advantage of implementing the required AI solutions as close as possible to the end-users and the data sources. Edge devices often face extremely tight constraints in terms of size and energy budget where the results of the data analysis serve as inputs to a controller and are translated into concrete and immediate action within a very tight latency window. This is definitely the case in applications such as autonomous mobile devices (cars, drones), robotics, human-machine interfaces (AR/VR), brain-machine interfaces and wearable medical devices.

While one can easily imagine that parts of such systems can be run from a centralised location, energy, latency, robustness, security and privacy considerations form powerful arguments for more distributed realisations with a large share of the functionality executed on the Edge (or a

cluster/swarm thereof), and eventually penetrating all the way into the Extreme Edge (sensors and actuators for in-sensor computing).



Benefits of “AI at the Extreme Edge”

AI solutions that run autonomously, and are distributed and implemented **at the Extreme Edge** offer the following advantages:

- **Increased real-time performance (low-latency):** Edge applications process data and generate results locally on the sensing device. As a consequence, the device is not required to be continuously connected with a cloud data center. As it can process data and make decisions independently, there is increased real-time performance in the decision-making process, reduced delay of data transmissions and improved response speed.
- **Reliable low-bandwidth communication:** Distributed devices can handle a large number of computational tasks, therefore reducing the need to send data to the cloud for storage and further processing. Overall, this results in minimising the traffic load in the network and supports low-bandwidth communication.
- **Enhanced power-efficiency:** As the amount and rate of data exchange with the cloud is minimised, the power consumption of the device is reduced thus improving battery lifetime, which is critical for many edge devices.

- **Improved data security and privacy:** By processing data locally with Edge computing it does not have to be sent over a network to remote servers for processing. This improves data security and privacy as it can perform all processing disconnected from the central server, which is a more secure and private architecture.
- **High availability:** Decentralisation and offline capabilities make edge AI more robust since internet access is not required for processing data. This results in higher availability and reliability for mission-critical, production-grade AI applications.
- **Lower costs:** Edge AI has analytical scalability and reduced latency, which can lead to significant cost reductions for businesses. Processing and analysing large amounts of data in the cloud is not cheap. For example, cloud companies charge for every minute of inference per endpoint. That can be helpful for organisations that want to pay on demand, but for organisations that require a lot of real time processing (for example, smart cities or hospitals with many cameras and sensors running 24 hours a day), the burden becomes very heavy.

Market Opportunities of “AI at the Extreme Edge”

According to the study published by Astute Analytica, the Global Edge AI Software Market is projected to witness a **major rise in its revenue from US\$ 1,459.8 Mn in 2021 to US\$ 8,049.8 Mn by 2027**. The market is registering a CAGR of 29.8% over the forecast period. Various factors such as increasing enterprise workloads on the cloud and rapid growth in intelligent applications are expected to drive the adoption of edge AI solutions and services. The Edge AI hardware and the consulting market will grow at the same pace. Grand View Research estimates that the total global Edge computing market will grow 37.4 percent per year and will be worth **\$43.4 billion by 2027**.



Source: European Commission, Towards 5G.

Key drivers of extreme-edge AI computing include:

- **Advanced Communication Technologies:** 5G networks can collect large and fast data streams. The construction of 5G networks is beginning gradually, and they will be established initially in very local and densely populated areas. When the utilisation and analysis of these data streams are as close as possible to the devices connected to the 5G network, the value of Edge AI technology will increase.
- **Massive Amounts of IoT Generated Data:** Currently, about 5%–15% of the world's energy is spent in some form of data manipulation, such as transmission or processing [1], and this fraction is expected to rapidly increase due to the exponential increase of data generated by ubiquitous sensors in the era of internet of things. If the data is separated from its source and there is no metadata describing the meaning of the data, the data will tell us nothing. Therefore, recovering data is not enough. Only Edge AI can take full advantage of advertised IoT data. Large amounts of sensor data can be analysed locally, and operational decisions can be performed automatically. Only the most important data is stored in a data warehouse in the cloud or data centre.
- **Customer Experience:** People expect a smooth and flowing service experience. Today, a delay of just a few seconds can easily ruin the customer experience. Edge computing responds to this demand by eliminating the delay caused by data transmission. Additionally, sensors, cameras, GPU processors and other hardware continue to be cheaper, so more and more people are able to use highly productive and custom edge AI solutions.

Key Trends and Challenges of “AI at the Extreme Edge”

The massive opportunity of Extreme Edge AI will require a number of challenges to be addressed and solved, specifically in the areas of secure and trustworthy AI. For a full detailed overview of key trends and challenges please visit:

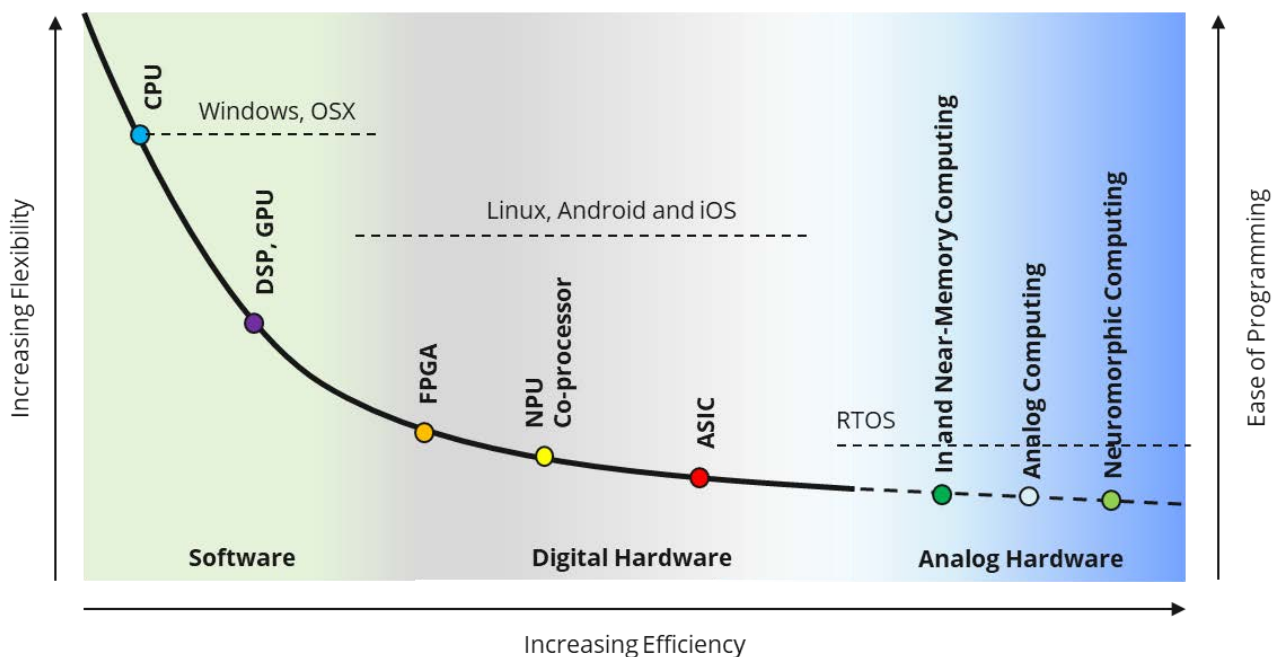
<https://bonseyes-association.org/roadmap-future-challenges/>

which covers in more depth and detail the following topics:

- Trust and explainability
- Re-learning
- Security and adversarial attacks
- Learning at the edge
- Integrating AI into the smallest devices
- Data as a basis for AI
- Neuromorphic technologies
- Meta-learning
- Hybrid modelling
- Energy efficiency

Hardware for Edge AI

Choosing the ideal hardware for a particular application requires careful consideration of all the requirements. A successful system design finds a balance between the different aspects of system architecture, such as memory footprint, executing time, model accuracy, power consumption, scalability, cost, and maintainability. While data centers allow engineers to scale available computational power to the current demand (via GPUs or TPUs), an application running on edge devices needs to keep sufficient power reserves. An increasing number of vendors are now moving from producing simple resource-constrained microcontrollers (e.g. ARM Cortex MCU) to pair general-purpose processors with specialised units tailored to execute the computational tasks required to implement AI solutions. As embedded systems are typically focused on using AI in the form of machine learning (ML) for interpreting incoming sensor data, these specialised sub-processors aim to speed up classification or prediction tasks while maintaining a low power draw. This is especially important in applications running on battery power or with a low potential for cooling the system.



Source: NVISO adaptation of "Computer vision algorithms and hardware implementations: A Survey," by Xin Fang, 2019

Today there are three main classes of hardware available for model deployment with various tradeoffs between flexibility and efficiency:

- **Standard programmable CPU/GPUs with specialised co-processors** : A common approach is the use of specialised co-processors, which can be either directly co-located with the general-purpose processor on the same silicon or might be connected as a separate chip. These accelerators stretch the power continuum from relatively simple digital signal processors (DSPs) to highly parallel matrix computation units and similar advanced

designs. These accelerators can either be monolithic designs such as special Edge variants of Google's TPU or a distributed set of smaller compute cores. Some examples for the latter are the Tensor Cores in newer Nvidia GPU architectures, the Hexagon cores in Qualcomm Snapdragon SoCs and Intel Movidius SHAVE processors. They commonly need specialised drivers and software libraries that allow software developers to take advantage of their capabilities.



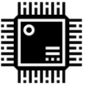
- **Digital hardware:** For greater efficiency, field-programmable gate arrays (FPGAs) can be combined with powerful CPU systems, specialised arrays of AI accelerator cores and traditional fabric, where the hardware can be programmed by the use of hardware description language (HDL) or C/C++ via high-level synthesis tools. FPGAs combine the benefits of specialised hardware with the freedom to change the layout even after the chip has left the factory. The use of FPGA can create more dynamic, scalable, and flexible systems, even though they often carry higher cost. Another type of processor is emerging as a new class of processing accelerator for these predominantly data-centric heterogeneous processing tasks as offload to the main CPU. These processors are a well-partitioned circuit that comprises all the control and arithmetic logic components necessary to execute machine learning algorithms and are referred as NPUs (Neural Processing Unit). NPUs are designed to accelerate the performance of common machine learning tasks. NPUs may be part of a large SoC, a plurality of NPUs may be instantiated on a single chip, or they may be part of a dedicated neural-network accelerator.
- **Analog hardware:** The biggest time and energy costs in most computers occur when lots of data has to move between external memory and computational resources such as CPUs, GPUs, or NPUs. This is the "von Neumann bottleneck," named after the classic computer architecture that separates memory and logic. One way to greatly reduce the power needed for machine learning is to avoid moving the data — to do the computation where the data is located. Because there is no movement of data, tasks can be performed in a fraction of the time and require much less energy. In-memory and neuromorphic analog computing are such approaches that take this approach. In-memory computing is the design of memories next to or within the processing elements of hardware such that bitwise operations and arithmetic operations can occur in memory. Neuromorphic analog computing allows in-place compute but also mimics the brain's function and efficiency by building artificial neural systems that implement «neurons» and «synapses» to transfer electrical signals via an analog circuit design. This circuit is the breakthrough technology solution to the Von Neumann bottleneck problem. Analog neuromorphic ICs are intrinsically parallel, and better adapted for neural network operations than current digital hardware solutions offering orders of magnitude improvements for edge applications. Market-ready, end-user programmable chips are an essential need for neuromorphic computing to expand its visibility and to achieve a variety of "real-world applications" with an increasing number of users. As large technology companies are waiting for the technology to become more mature, some start-ups have released their chips to fill the gap and to have a competitive advantage against those tech-giants. Examples include BrainChip AKIDA, GrAI Matter Labs VIP, and SynSense DYNAP processors.

Software for Edge AI

Edge computing software requires the implementation of comprehensive intelligent edge frameworks and platforms addressing specific requirements and challenges spanning hardware, power efficiency, software, connectivity, flexibility and interoperability, and security. Software stacks for edge computing are inherently fragmented due to the following:

- Not all applications are built the same way, the software stack will determine how well the final application will perform.
- In order to achieve full capacity for an application, a software stack needs to be tailored specifically to what is trying to be accomplished and the underlying hardware.
- Different models require specific tools that only customizable stacks will offer.
- The software stack must include a scheduler and run-time engine that must meet varying requirements of latency-sensitive applications.

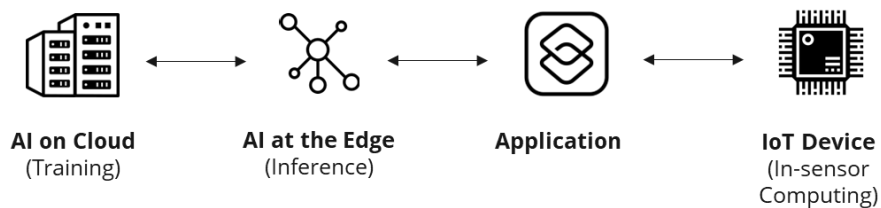
Software for machine learning can be broken down into three main categories:

		Form Factors	Algorithms	Hardware
	Cloud ML	Data Centers Super Computers	Deep neural network	TPU, FPGA, GPU, and CPU CPU: x86/ARM64 Memory: TB Storage: PB Clock Freq: GHz Power: 250W
	Edge ML	Local Servers Personal Computers Smart Phones Smart Devices	Optimized algorithms and lightweight convolution neural networks	ASICs with NPU / Coprocessors CPU: ARM Memory: MB - GB Storage: GB - TB Clock Freq: MHz - GHz Power: ~8W
	Tiny ML	Wearables Smart Sensors Microcontrollers Actuators	Micro convolution neural networks	MCUs with Neuromorphic Fabric CPU: MCU Memory: ~500KB Storage: <2MB Clock Freq: <200MHz Power: <1W

- **CloudML:** In CloudML, machine learning is used to integrate artificial intelligence into applications via cloud infrastructure. Cloud service providers provide a wide range of AI services that developers can make use of in their applications. Devices that transmit the data required by these services are connected to the internet through which they send the data for processing, storage and analysis. CloudML can be characterised by:
 - Typically applied to a limited number of focused, vertical applications
 - Targeting a small range of processors
 - Plenty of available power and bandwidth
 - Large equipment budget
- **EdgeML:** EdgeML is made for real-time, always-on solutions. By processing data as close to its source as possible, latency is minimized and organizations gain actionable insights in real time. EdgeML can be characterised by:
 - Potentially applied to a wide and diverse range of applications
 - Many possible processor targets, from CPUs and GPUs to NPUs, DSPs and other forms of dedicated accelerator

- Numerous – often proprietary – application programming interfaces (APIs)
- Relatively low-cost devices operating in thermally and power-constrained environments
- **TinyML:** TinyML is a fast-growing multidisciplinary field at the intersection of machine learning, hardware, and software, that focuses on enabling deep learning algorithms on embedded (microcontroller powered) devices operating at extremely low power range (mW range and below) targeting in-sensor computing. TinyML can be characterised by:
 - compact and low-cost devices (microcontrollers)
 - very low power usage
 - extremely limited memory capacity
 - low lag time (almost immediate) integrated machine learning algorithms analysis.

Moving AI from the cloud to the edge can be challenging because of the need to create neural network architectures that can operate well in edge AI chips. Cloud servers run on general computing platforms which can perform on any network architecture, but in edge AI, the models and architectures have to be modified to cope with the low resources of the AI modules at the edge. In addition, tools, inference engines, acceleration frameworks, and operating systems vary greatly from the cloud, edge computing, and IoT devices.



	AI on Cloud (Training)	AI at the Edge (Inference)	Application	IoT Device (In-sensor Computing)		
Tools and Inference Engines	TensorRT	ONNX Runtime	CoreML	TensorFlow	TensorFlow Lite	Vendor Specific
ML Acceleration Libraries	CUDA	DirectML	Accelerate	NNAPI	OpenCL	CMSIS-NN
BSPs and Operating Systems	Linux	Windows	iOS	Android	Embedded Linux	RTOS

To efficiently leverage this heterogeneous computing environment for edge computing, it is critical to have a full software stack able to manage multiple layers from high-level level training, through efficient inference engines, to developer board support packages and operating systems. The main components being:

- **Datatools** are containerized python based utilities, which extract data from a given source dataset and translate both unstructured (images, audio, video files etc.) and semi-structured (CSV, JSON, TXT, XML, XLS, etc.) data to common representation. Datatools generate a standard data structure and definitions for the annotations, which make data easier to analyse, consume and reuse by the subsequent model training pipeline.
- **Training Pipelines** are workflows that deliver powerful and easy-to-deploy building blocks for creating complex AI models that can be deployed on cyber-physical systems. By taking care of many end-to-end tooling dependencies and providing standardised interfaces, the

Training Pipelines enable users to focus on producing optimal solutions while allowing faster feedback during the implementation of end user requirements. Such workflows enable accelerated deployment of deep learning models to resource constrained low-power embedded systems (Deep Edge).

- **Runtime Frameworks** provide the capabilities to generate portable and efficient implementations of DNNs that can be deployed and optimised across heterogeneous platforms, e.g., CPU, GPU, FPGA, DSP, NPU. Such frameworks supply a collection of tools, executables, libraries and inference engines, featuring a full development flow for deep learning solutions on embedded devices by providing platform support, sample models, optimisation tools, and benchmarking.
- **Hardware Acceleration Libraries** are the collection of libraries and backends that allows the actual execution of the neural network on the hardware device. Hardware Acceleration Libraries are tightly coupled with the underlying hardware processor or accelerator, e.g., CPU, GPU, NPU, to fully optimise the AI workload on the device. General Hardware Acceleration Libraries may vary severely from system to system, often requiring a design search exploration to find the optimal one for a given target platform.
- **Board Support and Manager Packages** contain the full software stack (operating system, drivers, middleware components, etc.) and documentation required to procure, set up and control target hardware for the execution of AI Applications in an easy and reproducible format. The use of a Manager Package which works as a running environment for AI algorithms on the edge platform, supports inference tasks, benchmarking, and model re-training, easing the usage of AI on such difficult and heterogeneous environments.

About NVISO

NVISO is an Artificial Intelligence company founded in 2009 and headquartered at the Innovation Park of the École Polytechnique Fédérale de Lausanne (EPFL) in Switzerland with offices in Serbia and Japan. Its mission is to help teach machines to understand people and their behavior to make autonomous machines safe, secure, and personalized for humans. As leader in human behavioral AI, it provides robust embedded software solutions that can sense, comprehend, and act upon human behavior in real-world environments deployed at the deep edge. It achieves this through real-time perception and observation of people and objects in contextual situations combined with the reasoning and semantics of human behavior based on trusted scientific research. NVISO's technology is made accessible through ready-to-use AI solutions addressing Smart Mobility and Smart Health and Living applications (in-cabin monitoring systems, health assessments, and companion robot sensing) with a key focus on the deep and extreme edge with ultra-low power processing capabilities. With a singular focus on how to apply the most advanced and robust technology to industry and societal problems that matter, NVISO's solutions help advance human potential through more robust and rich human machine interactions. www.nviso.ai