

Hardware Architectures For Embedded And Edge AI

Christian Rossi

Academic Year 2024-2025

Abstract

Contents

1	Introduction	1
1.1	Introduction	1
1.2	Computing continuum	2
1.2.1	Datacenters	2
1.2.2	Edge computing systems	2
1.2.3	Embedded systems	3
1.2.4	Internet of Things	3
2	Hardware	4
2.1	Introduction	4
2.2	Architecture	5
2.3	Sensors and signals	5
2.3.1	Data	6
2.3.2	Sensors	7
2.4	Microprocessor	9
2.4.1	Microcontroller	9
2.4.2	System on chip	10
2.4.3	Comparison	10
3	Software	11
3.1	Introduction	11
3.2	Data preprocessing and feature extraction	12
3.2.1	Data segmentation	12
3.2.2	Data processing	13
3.2.3	Feature extraction	13

Introduction

1.1 Introduction

Artificial Intelligence (AI) is a field of computer science focused on developing hardware and software systems capable of performing tasks that typically require human intelligence. These systems can autonomously pursue specific goals by making decisions that were traditionally made by humans.

A key distinction in AI-driven systems lies between smart objects and connected objects. While connected objects primarily send and receive data from the cloud, smart objects analyze data locally, enabling faster decision-making and reducing reliance on constant connectivity.

The definition of AI evolves rapidly, to the point that what was considered AI a decade ago may differ significantly from today's understanding.

AI hardware and software can be categorized similarly to traditional computing environments but are specifically designed to handle AI workloads. In this context, the development environment is often referred to as a framework, platform, or tool, rather than just a conventional programming environment.

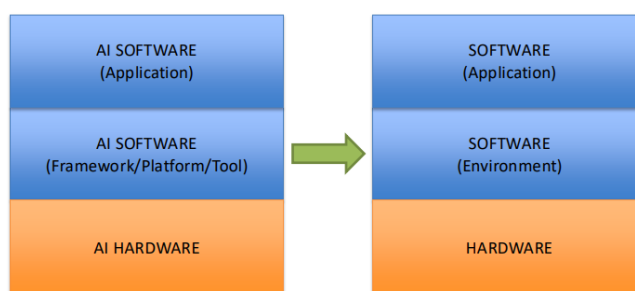


Figure 1.1: Artificial Intelligence stack

The AI stack consists of three main layers:

- *AI software* (application): AI-powered applications running within an IT system.
- *AI software* (framework, platform and tool): programs and libraries that manage physical resources and provide the necessary tools for building AI applications.

- *AI hardware*: the infrastructure supporting AI computation, including data centers, edge computing devices, IoT systems, and specialized processors like CPUs, GPUs, and TPUs.

1.2 Computing continuum

AI hardware ranges from small, low-power devices running on batteries to large-scale, high-performance systems in datacenters. This range represents the computing continuum:

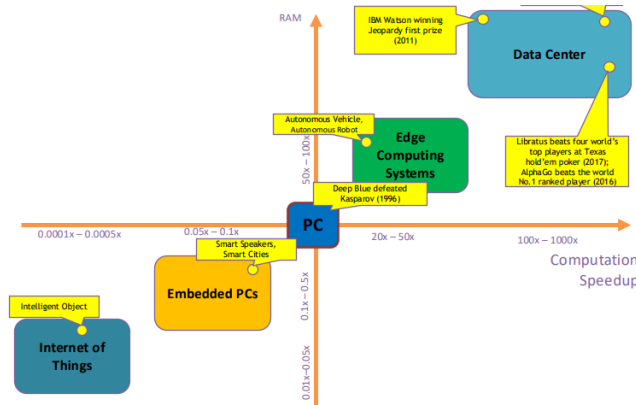


Figure 1.2: Computing continuum

1.2.1 Datacenters

Datacenters sit at the upper end of this spectrum, offering immense computational power for AI workloads. They provide cost-efficient IT infrastructure, high-performance computing capabilities, and instant software updates. Their vast storage capacity ensures data reliability and accessibility, allowing seamless collaboration across devices and locations. Furthermore, by decoupling AI processing from end-user devices, datacenters enable powerful AI applications that are not limited by local hardware constraints.

Datacenters require a constant internet connection, making them less viable in low-bandwidth environments. Their reliance on shared infrastructure can introduce privacy and security concerns, while the lack of direct hardware control may limit customization options. Additionally, the high energy consumption of large-scale AI operations raises both environmental and cost concerns. In latency-sensitive applications, delays in data transmission and processing can further impact real-time decision-making.

1.2.2 Edge computing systems

Edge computing delivers high computational power with the advantage of distributed processing. By bringing computation closer to where data is generated, it enhances privacy and security while significantly reducing latency in decision-making. However, these systems depend on a stable power supply and often integrate with cloud services to extend their processing capabilities.

By processing data locally, edge computing minimizes the need for constant data transmission, optimizing bandwidth and improving energy efficiency. This approach not only strengthens security and privacy but also enables real-time decision-making and adaptive learning across distributed networks. However, edge devices often operate with limited computing resources,

constrained memory, and restricted energy availability. Their design requires careful coordination of hardware, software, and machine learning models, adding complexity to development and deployment.

1.2.3 Embedded systems

Embedded systems, widely used in AI applications, provide high-performance computing in a compact form. They benefit from the availability of development boards and can be programmed similarly to traditional computers, making them accessible to a broad community of developers. Despite these advantages, they tend to consume relatively high power, and in some cases, require custom hardware design to meet specific application needs.

1.2.4 Internet of Things

At the smallest scale, the Internet of Things (IoT) enables AI integration into pervasive, low-cost, battery-powered devices. These systems support wireless connectivity and often include sensing and actuating capabilities, making them essential for smart environments. However, IoT devices face limitations in computing power, energy efficiency, and memory capacity, which can complicate programming and constrain their ability to run advanced AI models.

CHAPTER 2

Hardware

2.1 Introduction

In embedded and edge AI systems, a typical setup includes sensors that capture data from the physical world, software that processes this data, and actuators that execute actions based on computational outcomes. All processing tasks rely on specialized hardware optimized for efficiency and performance.

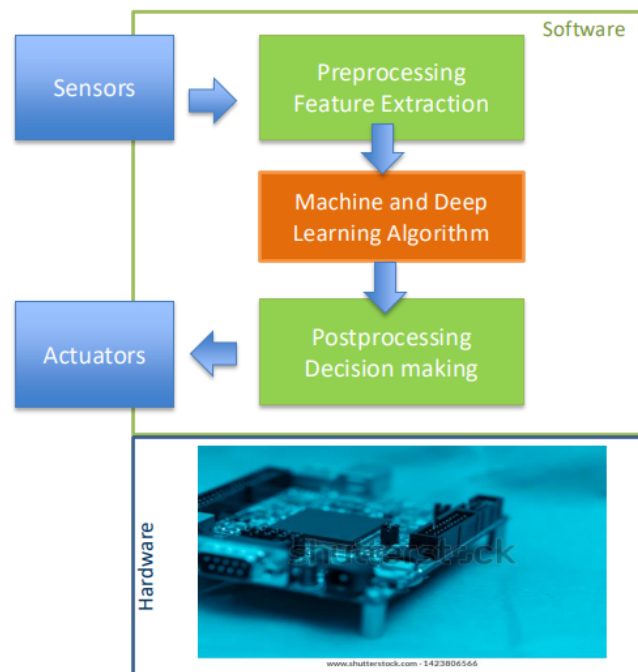


Figure 2.1: AI systems

Embedded systems are computers designed to control and manage the electronics within various physical devices. Embedded software refers to the programs that run on these systems, enabling their functionality.

Unlike general-purpose computers such as laptops or smartphones, embedded systems are typically designed for a specific, dedicated task, ensuring optimized performance, reliability,

and energy efficiency for their intended application.

2.2 Architecture

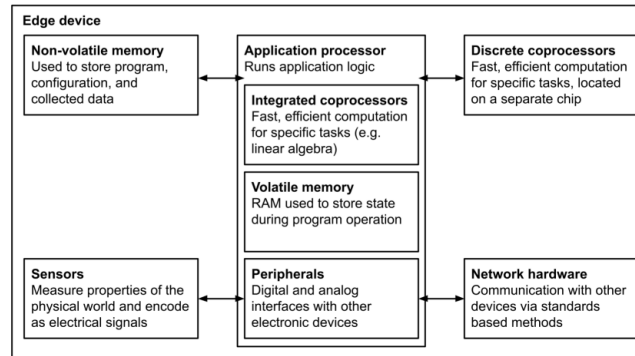


Figure 2.2: Hardware architecture

The hardware architecture of embedded and edge AI systems consists of several key components:

- *Non-volatile memory*: used to store programs, configurations, and collected data. Flash memory is typically used for this purpose, as it retains data even when the system is powered off. It is ideal for storing information that does not change frequently but is slow to read and extremely slow to write.
- *Application processor*: runs the application logic and manages program execution. It includes an integrated coprocessor for efficient computation of specific tasks, volatile memory (RAM) for storing the system state during operation, and various digital and analog peripherals that allow interaction with other electronic components.
- *Discrete coprocessors*: external chips designed for high-speed, efficient mathematical computations. They provide additional processing power for specialized AI workloads that require high performance.
- *Sensors*: measure physical-world properties and convert them into electrical signals for processing. They enable real-time data collection, which is essential for AI-driven decision-making.
- *Network hardware*: ensures communication with other devices using standardized protocols. Reliable connectivity is crucial for data exchange in distributed AI systems.

RAM is often the performance bottleneck in embedded and edge AI systems. It is very fast but consumes significant energy, making efficiency critical in power-sensitive applications. Since it is volatile, data is lost when power is turned off. RAM is also costly and takes up a large physical footprint, impacting the overall design of embedded AI devices.

2.3 Sensors and signals

Sensors are used to acquire measurements from the environment or from human interactions. They generate continuous streams of data, which can be used for various AI-driven applications.

In addition to sensor-generated data, other sources such as digital device logs, network packets, and radio transmissions can also provide valuable information. Sensors can output data in different formats depending on their purpose and design.

Data storage Data values can be stored in various formats, depending on precision and memory constraints. Boolean values (1 bit) represent binary states with two possible values. An 8-bit integer can store up to 256 distinct values, while a 16-bit integer extends this range to 65,536 possible values. A 32-bit floating point number can represent a wide range of values with up to seven decimal places, reaching a maximum of approximately 3.4×10^{38} . Quantization techniques help optimize memory usage by reducing the required storage for each value while maintaining sufficient precision for AI computations.

2.3.1 Data

Time series data A time series is a sequence of data points recorded in chronological order:

$$X = (x_1, x_2, \dots, x_N)$$

Essentially, it represents observations collected at consistent time intervals. Key factors to consider include:

- *Sampling period*: the time gap between consecutive data points.
- *Bit depth* (n): the number of bits used to represent each value.
- *Memory usage*: each sample requires n bits of storage.

Audio data Audio data is a specific type of time series, representing sound wave oscillations as they travel through air. Key parameters are:

- *Sampling rate* (Hz): the number of samples taken per second.
- *Quantization* (bit depth): the precision of each sample.
- *Signal duration* (s): the total length of the recording.
- *Number of channels*: mono (single channel) or stereo (two channels).

Memory consumption is calculated as:

$$\text{length} \times \text{sampling} \times \text{bit} \times \text{channel}$$

Image data Images capture visual information as a grid of pixels, where each pixel represents a specific property of the scene. Key characteristics are:

- *Dimensions* ($W \times H$): width and height of the image.
- *Bit depth* (N): the number of bits used to store each pixel.
- *Number of channels*: typically 1 (grayscale) or 3 (RGB).

Memory usage is given by:

$$W \times K \times N \times \text{channels}$$

Video data Videos are sequences of images displayed rapidly to create motion. They share the same structure as images but add an extra dimension: time. Critical parameters:

- *Resolution* ($W \times H$): width and height of each frame.
- *Bit depth* (N): bits per pixel.
- *Number of channels*: defines color representation.
- *Frame rate* (fps): number of frames per second.
- *Duration* (s): total length of the video.

Memory requirements are determined by:

$$W \times K \times N \times \text{channels} \times \text{frame rate} \times \text{length}$$

2.3.2 Sensors

There are thousands of different types of sensors available, each designed to capture specific kinds of data. In the context of embedded and edge AI, sensor technologies can be categorized into six main families:

1. *Acoustic and vibration*: detects sound and mechanical vibrations.
2. *Visual and scene*: captures images, video, and environmental light data.
3. *Motion and position*: measures movement, acceleration, and spatial positioning.
4. *Force and tactile*: detects pressure, touch, and force.
5. *Optical, electromagnetic, and radiation*: measures light, radio waves, and radiation levels.
6. *Environmental and chemical*: monitors temperature, humidity, gases, and other environmental factors.

Acoustic and vibration Detecting vibrations is a crucial capability in embedded and edge AI. These sensors allow systems to perceive movement, structural vibrations, and even communication signals from humans and animals at a distance. Acoustic sensors measure vibrations traveling through different media: air (microphones), water (hydrophones), and ground (geophones and seismometers). Since acoustic data is distributed across different frequencies, the sampling frequency plays a key role in ensuring accurate representation for a given application. These sensors typically produce audio data as their output.

Visual and scene Visual sensors capture information about the environment without direct contact. These range from tiny, low-power cameras to high-resolution multi-megapixel sensors. Key characteristics of image sensors: color channels, spectral response (infrared sensors), pixel size, resolution, and frame rate. The output of these sensors can be 2D or 3D images or video data, depending on the application.

Motion and position Motion and position sensors track movement and spatial positioning in various ways:

- *Tilt sensors*: simple mechanical switches that detect orientation changes.
- *Accelerometers*: measure acceleration along one or more axes.
- *Gyroscopes*: detect rotational movement.
- *Time-of-flight sensors*: emit light or radio waves to measure distances to objects.
- *Real-time locating systems*: use multiple transceivers placed around a space to track object positions.
- *Global Positioning System*: uses satellites to determine an object's precise location.

These sensors typically generate time-series data, tracking movement and positioning over time.

Force and tactile These sensors help users interact with devices, understand fluid and gas flow, or measure mechanical strain on objects:

- *Buttons and switches*: provide a simple on/off signal.
- *Capacitive touch sensors*: detect how much a conductive object touches a surface.
- *Strain gauges*: measure how much an object is being deformed.
- *Load cells*: determine the amount of force or weight applied to an object.
- *Flow sensors*: track the movement of liquids and gases.
- *Pressure sensors*: measure pressure in gases or liquids, whether in the environment or within a system.

Optical, electromagnetic, and radiation These sensors detect electromagnetic radiation, magnetic fields, and electrical properties:

- *Photosensors*: detect light across different wavelengths, from visible to infrared and ultraviolet.
- *Color sensors*: measure surface colors to help recognize different materials or objects.
- *Spectroscopy sensors*: analyze how light interacts with materials to determine their composition, useful in AI-driven material analysis.
- *Magnetometers*: measure the strength and direction of magnetic fields, like a digital compass.
- *Inductive proximity sensors*: use electromagnetic fields to detect nearby metal objects, commonly used in vehicle detection for traffic monitoring.
- *Electromagnetic field meters*: measure the intensity of electromagnetic fields, useful in industrial safety assessments.
- *Current and voltage sensors*: monitor electrical parameters, often represented as time-series data.

Environmental, biological, and chemical These sensors track environmental conditions, biological signals, and chemical presence:

- *Temperature sensors*: measure heat levels in the surroundings or specific systems.
- *Gas sensors*: detect gases like humidity, carbon dioxide, or air quality pollutants.
- *Particulate matter sensors*: monitor air pollution levels by measuring fine particles.
- *Biosignal sensors*: record electrical activity in the body.
- *Chemical sensors*: detect the presence or concentration of specific chemicals.

Most of these measurements are typically recorded as time-series data, allowing for trend analysis and predictive insights.

2.4 Microprocessor

A microprocessor is a general-purpose processor responsible for running embedded applications. Microcontrollers form the core of many modern, pervasive computing applications. These tiny, cost-effective computers are designed for specific tasks, making them ideal for embedded systems.

Unlike traditional computers, microcontrollers do not require an operating system. Instead, they run firmware (software that is directly executed on the hardware). This firmware includes low-level instructions to manage peripherals and system functions. The defining feature of microcontrollers is that they integrate all essential components into a single silicon chip.

2.4.1 Microcontroller

Microcontrollers have a fixed hardware architecture built around a central processing unit. The CPU manages a range of peripherals, providing both digital and analog functionality. Smaller devices typically include both volatile (RAM) and non-volatile (Flash/EEPROM) memory on the chip, while more powerful processors may require external memory. Programming is commonly done using low-level languages like Assembly or high-level languages like C.

Microcontrollers are the preferred choice for embedded systems because they integrate essential components like memory and peripherals, reducing the need for additional circuitry. This allows for compact, power-efficient designs that are crucial in space-constrained applications. A microcontroller typically consists of:

- A microprocessor.
- Program memory (Flash).
- Data memory (RAM).
- Various on-chip peripherals such as timers, serial communication ports, GPIO pins, counters, and ADCs (analog-to-digital converters).

Modern microcontrollers have significantly improved in both performance and efficiency:

- Higher clock speeds and 32-bit architectures for better computational power.

- SIMD support for parallel computing, beneficial in signal processing and machine learning.
- Larger RAM to handle more complex tasks—critical since memory is often a limiting factor.
- Improved energy efficiency, making them suitable for battery-powered applications.

Microprocessor and microcontroller A microprocessor is a more advanced computing unit, typically requiring external memory to fetch and execute program instructions. It offers greater processing power compared to microcontrollers but lacks integrated peripherals.

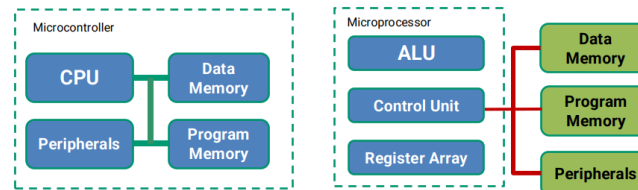


Figure 2.3: Microprocessor and microcontroller

2.4.2 System on chip

A System on Chip (SoC) integrates all the core functionalities of a traditional computing system into a single chip. This includes the CPU, memory, input/output interfaces, and sometimes specialized accelerators for AI, graphics, or signal processing. They run full operating systems and support standard development tools and libraries. However, SoCs have lower energy efficiency compared to microcontrollers, making them less ideal for ultra-low-power applications, and also more costly.

2.4.3 Comparison

The table below compares different processing devices based on their ability to handle various types of data.

Device	Time Series	Audio	Images	Video
<i>Low-end MCU</i>	Limited	None	None	None
<i>High-end MCU</i>	Full	Full	Low resolution	Limited
<i>High-end MCU with accelerator</i>	Full	Full	Full	Limited
<i>SoC</i>	Full	Full	Full	Full
<i>SoC with accelerator</i>	Full	Full	Full	Full
<i>Edge server</i>	Full	Full	Full	Full
<i>Cloud</i>	Full	Full	Full	Full

Software

3.1 Introduction

Artificial Intelligence systems in embedded and edge computing are designed to process data efficiently and make real-time decisions.

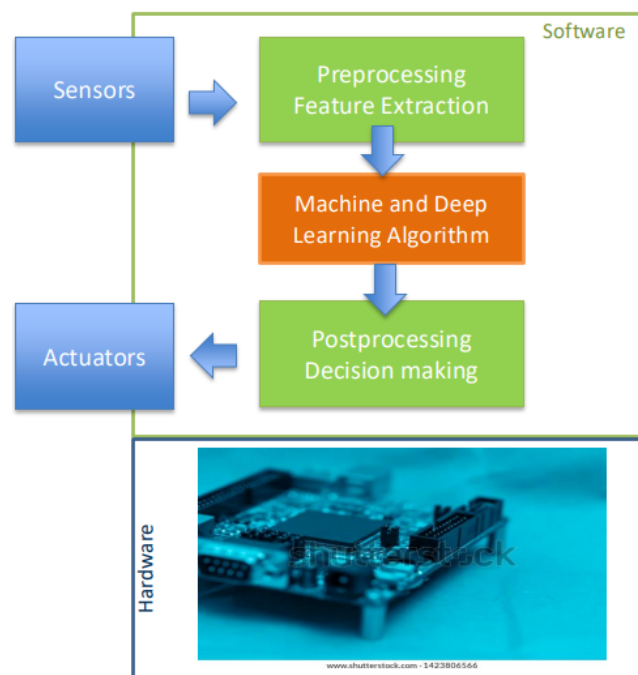


Figure 3.1: AI systems

The first stage is preprocessing and feature extraction, where incoming data streams from sensors are processed. This step involves segmenting the data into meaningful windows, removing noise, and extracting relevant features to enhance accuracy.

Once the data is preprocessed, machine learning and deep learning algorithms are applied to analyze and interpret it. These algorithms recognize patterns, make predictions, and generate insights based on the extracted features.

Finally, in the postprocessing and decision-making stage, the system interprets the AI

model's output and translates it into meaningful actions. This could involve making decisions, triggering automated responses, or providing feedback to users.

3.2 Data preprocessing and feature extraction

The goal of data preprocessing and feature extraction is to convert raw signals into structured data that can be processed by Machine Learning and Deep Learning algorithms. This process involves three key steps:

1. *Segmenting the raw signals*: breaking the continuous signal into smaller, manageable chunks of data.
2. *Processing the data*: applying digital signal processing techniques to reduce noise and enhance the most relevant parts of the signal.
3. *Feature extraction*: identifying and extracting meaningful patterns or characteristics from the processed signal chunks to be used in model training.

3.2.1 Data segmentation

Sensors generate continuous streams of data, which must be divided into smaller segments (windows) for processing. Each window represents a chunk of data that is analyzed by an algorithm to produce meaningful results.

Definition (*Latency*). Latency refers to the time required to process a single chunk of data.

Latency plays a crucial role in determining how efficiently an embedded system can process data. Lower latency allows for a higher number of processed chunks per unit of time. However, there is a trade-off:

- Larger windows increase latency but provide more information, often leading to improved accuracy.
- Smaller windows reduce latency but may capture less useful data, potentially affecting performance.

Windows can be:

- Overlapping, ensuring no information is lost from the signal.
- Non-overlapping, which may be more computationally efficient but risks missing important details.

The choice of algorithm significantly impacts both computational requirements (latency) and memory usage, making it a key design consideration.

Frame rate The frame rate defines how frequently the system can acquire and process data, similar to how frame rate applies to image and video streaming.

High latency can hinder real-time analysis by preventing the system from processing new incoming data while still handling previous chunks, potentially leading to data loss. Optimizing both latency and frame rate is essential for effective data segmentation and real-time performance.

3.2.2 Data processing

Data preprocessing using digital processing algorithms typically involves three key steps:

1. *Reconstruction of missing data:*

- *Global filling methods:* filling in missing data based on patterns and trends observed across the entire dataset.
- *Local filling methods:* estimating missing values using nearby data points. Examples include forward fill, moving averages, and local interpolation techniques.
- *Deletion of affected time periods:* in some cases, time periods with missing data are removed entirely to avoid inaccuracies in analysis.

2. *Resampling:*

- *Time series resampling:* Handling time series data with varying sampling frequencies. In upsampling we increase the sampling rate by replicating or interpolating between timestamps. In downsampling we reduce the sampling rate through subsampling. Be cautious of aliasing when changing the sampling frequency.
- *Image resampling:* adjusting the spatial resolution (pixels per image). In downsampling we reduce the image resolution by decreasing the number of pixels. In interpolation we increase the image resolution by adding pixels based on existing data.
- *Shape modification for images:* In cropping we trim parts of the image. In resizing we adjust the image dimensions while maintaining or altering aspect ratios.

3. *Filtering:* we can use different types of filter:

- *Low-pass filter:* retains low frequencies, removing high-frequency noise. Pay attention to the cutoff frequency and frequency response.
- *High-pass filter:* retains high frequencies, removing low-frequency components. The cutoff frequency and response are important here as well.
- *Band-pass filter:* keeps a specific range of frequencies, removing those outside the band.

3.2.3 Feature extraction

Feature extraction can be applied across various domains and types of data:

- *Time domain:* features like mean, PCA eigenvalues, amplitude, signal-to-noise ratio (SNR), peak decay, and energy can be extracted.
- *Frequency domain:* features such as maximum amplitude, dominant frequency, and peak variance are commonly extracted.
- *Images:* features like edges, corners, blobs, and ridges (curves) are often detected. In embedded systems, tools for feature detection include OpenCV, a widely used library for image processing and feature detection, especially in System on Chips, and OpenMV, a library optimized for high-end microcontroller units, which focuses on efficient image processing and feature extraction.

Sensor fusion In sensor fusion, the goal is to combine data from multiple sensors rather than relying on a single sensor. Each sensor provides unique perspectives and data, and by combining these diverse inputs, we create a more comprehensive and accurate representation of the environment or system.

The features extracted from each sensor are merged to enhance the robustness, accuracy, and reliability of the analysis. These fused features are then used in machine learning or deep learning models for further processing, ultimately enabling more precise and informed decision-making.

Normalization and standardization For more efficient training, data should be normalized or standardized. Features with different scales can negatively impact the performance of machine learning models, potentially leading to underfitting.

- *Minmax normalization*: scales data to a range of $[0, 1]$:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- *Standardization*: centers data around zero (mean), but does not limit it to a specific range:

$$x' = \frac{x - \mu}{\sigma}$$

Here, μ is the mean and σ is the standard deviation.