# Navigating the Fragmented Machine Learning Ecosystem for Hardware Devices

**Vaidheeswaran Archana**
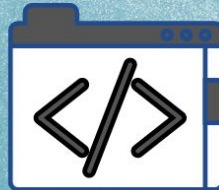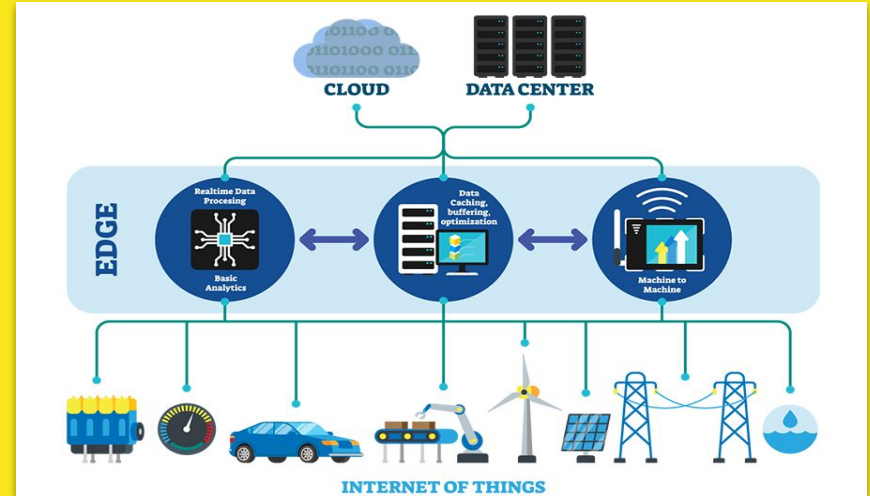
# Internet of Things

- IoT devices have a microcontroller and sensors

- They can collect and transfer data, as well as perform simple tasks

- IoT devices have become popular over the last few years

- They are used as smart home and fitness devices, but are also being used in industrial settings

- However **99% of IoT sensor data is discarded**

Number of global active Connections (installed base) in Bn

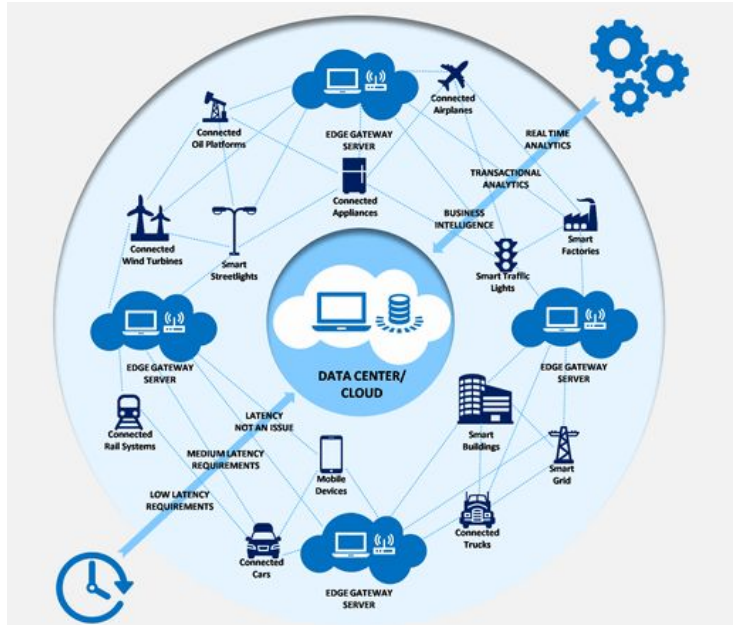| Year | Non-IoT | IoT | Total |
|------|---------|-----|-------|
| 2015 | 10.1 | 3.8 | 13.9 |
| 2016 | 10.3 | 4.7 | 15.0 |
| 2017 | 10.6 | 5.9 | 16.4 |
| 2018 | 10.8 | 7.0 | 17.8 |
| 2019 | 11.1 | 8.3 | 19.4 |
| 2020 | 11.3 | 9.9 | 21.2 |
| 2021 | 11.6 | 11.6 | 23.2 |
| 2022 | 11.9 | 13.5 | 25.4 |
| 2023 | 12.1 | 15.8 | 27.9 |
| 2024 | 12.4 | 18.5 | 30.9 |
| 2025 | 12.7 | 21.5 | 34.2 |

10%

Non-IoT
IoT

# Internet of Things and Edge Computing

- **Most DL models are deployed on cloud servers**
- **Running models closer to where the data is being generated is called edge computing**



Edge Computing Use Cases; innovationatwork.ieee.org

# Edge Computing



1– Increase in IoT devices causes an increase in cloud dependency

2-Need edge devices which have their own data centres

3-The computing that is performed in those data centers is Edge computing

4- Application: Security Cameras, Self Driving Cars

# Advantages of Edge Computing

- Reduced Latency

- Reduced Internet Bandwidth

- Increased Security

- Reduction in Dependence on Cloud Services

# Edge Computing Hardware

- **Microprocessors**
  - General Purpose
  - More Powerful
  - Consume More Power
- **Microcontrollers**
  - For Simpler Tasks
  - Less Computationally Powerful
  - Consumes Less Power
- **Accelerators**
  - Specialised for one task
  - Smaller and more power efficient
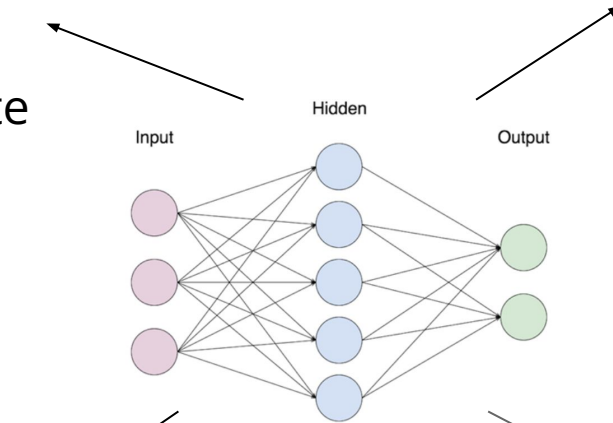  - Computationally powerful, but can only perform a specific task

# Need for Edge Computing Algorithms

**Power constraints:** Neural networks require massive amount of computational power and energy to execute on CPU and GPUs

**Composed of Floating Point Values:** Neural Networks are generally trained to preserve accuracy and not speed

**Memory Constraints:** Laptops and PC come with at least 4GB of RAM whereas Raspberry Pi 3 has 1GB of RAM

**Inference Efficiency:** We need model that takes less time for inference

# Edge Computing Frameworks

- TensorFlow Lite
  - Quantization
  - Pruning
  - Weight Clustering
  - Support for EdgeTPU
- OpenVINO
  - Quantization
  - Intermediate Representations
  - Support for NCS and other Intel Hardware
- Others
  - ONNX
  - PyTorch
  - DeepStream

# Knowledge Distillation

- Works by transferring the knowledge learned by a large teacher model to a smaller student model

- The student model is easier to execute

- The student model can be trained with unlabelled data

- Student models can often achieve similar or more accuracy than the teacher