# An introduction to AI with Intel

Dr. Séverine Habert

Thursday, May 4th 2023

intel. lrz

# Agenda

- A history of AI
- Concepts of Classical Machine Learning and Deep Learning
- Where classical ML shines
- Overview of Intel Hardware for AI
- Overview of Intel Software for AI

intel.

# A history of AI

# A definition of AI

- Concept:

Human thinking has the ability to both be replicated and mechanized.

intel.

# Before AI

The history of Artificial Intelligence follows the history of computer

- **create a mechanical object capable of thinking like human.**

Dates back from the antiquity / Ancient Greece where:

- was first developed automatons

- was first theorized formal reasoning

- was fantasized more sophisticated artificial beings in myths and legends



Antikythera mechanism, oldest known example of an analogue computer used to predict astronomical positions and eclipses decades in advance



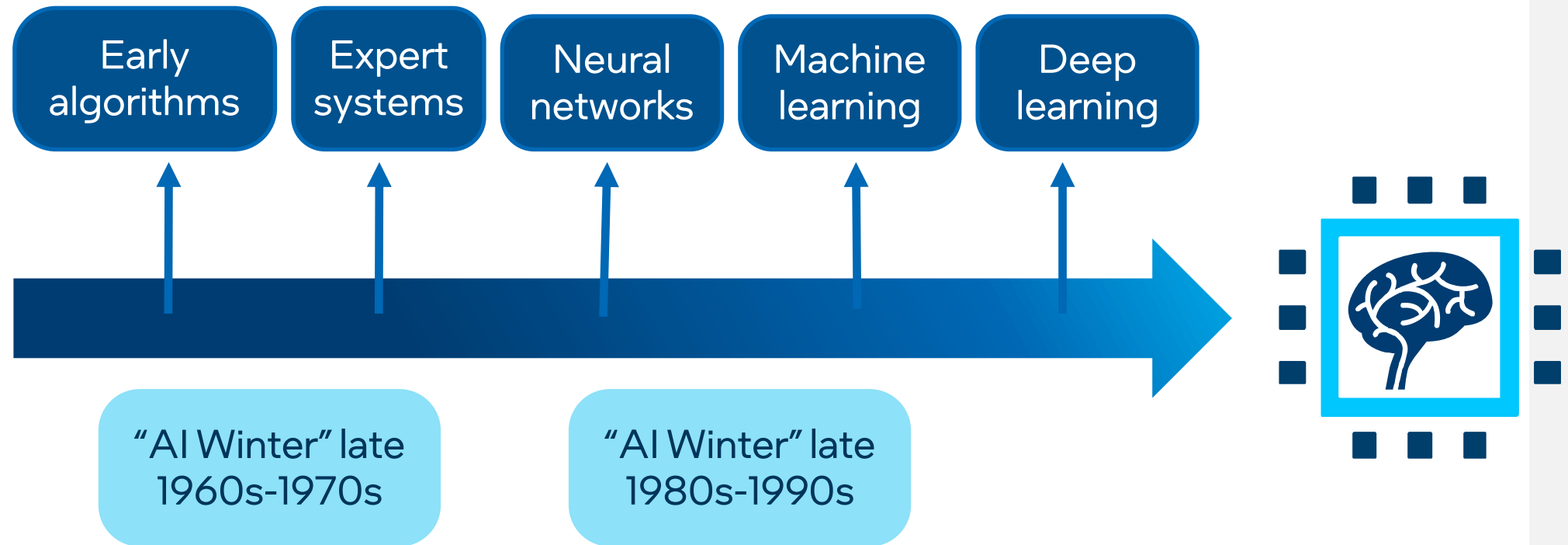Talos, bronze automate in Greek mythology (coin dating from 300/280-270 BCE

intel

# Forward in time

- Automata can be found through the history, Egypt, China, Renaissance

- Formal reasoning was also theorized by philosophers in various cultures : Al-Khwārizmī (algebra and algorithm), Descartes, Boole (mathematical logic)

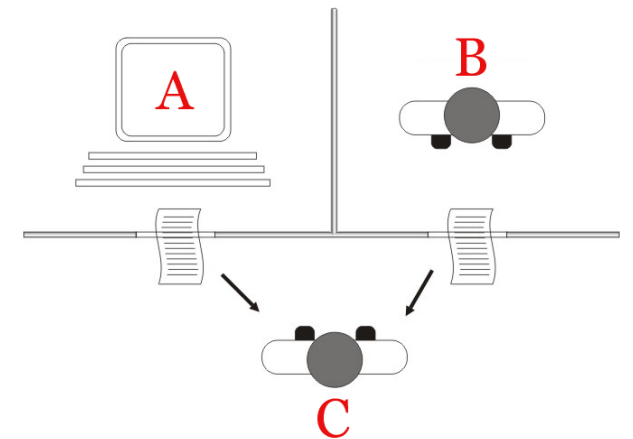- While "robots" start to appear in early 1900 in literature.

intel

# History in the XXth century

AI has experienced several hype cycles, where it has oscillated between periods of excitement and disappointment.



Early algorithms

Expert systems

Neural networks

Machine learning

Deep learning

"AI Winter" late 1960s-1970s

"AI Winter" late 1980s-1990s
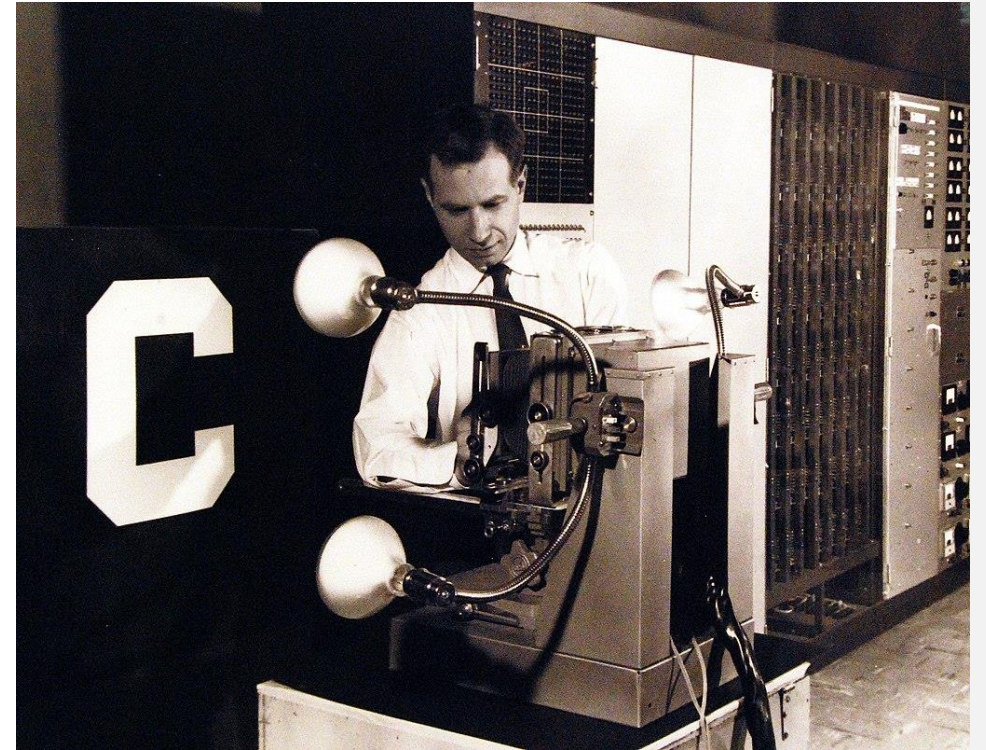
# Cybernetics and Turing test

- **Very early prototype of neural networks**
- Early 1940, better understanding of the human brain as electrical network of neurons
- In 1943, network of artificial neurons - perceptron.

- **Turing test:**
- In 1950, Turing defines the Turing test to determine if a machine is intelligent
    - can it fool a human with its answers?

intel

# 1950s: birth of AI

- 1956: Artificial Intelligence termed was coined at the Dartmouth Conference.

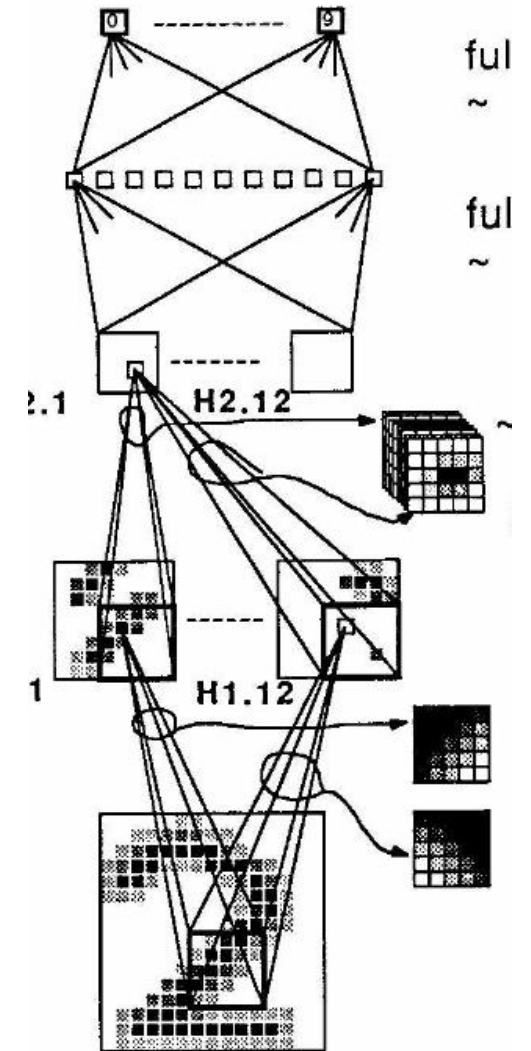- 1957: Frank Rosenblatt implemented the perceptron algorithm.



Camera system of the Mark 1 Perceptron

intel

# The First "AI Winter"

- 1969: limitations of the Perceptron algorithm
  - Can not represent non-linear functions

- 1966 & 1973: Two US government backed committee highlights AI's failure to live up to promises.

- No more funding -> the first "AI Winter."

# 1980's AI Boom

- Expert Systems become the first AI algorithm to be widely used in enterprise

- 1986: "Backpropogation" algorithm is able to train multi-layer perceptrons

- 1989: first theorical demonstration of backpropagation on Convolutional Neural Networks by Yann LeCun
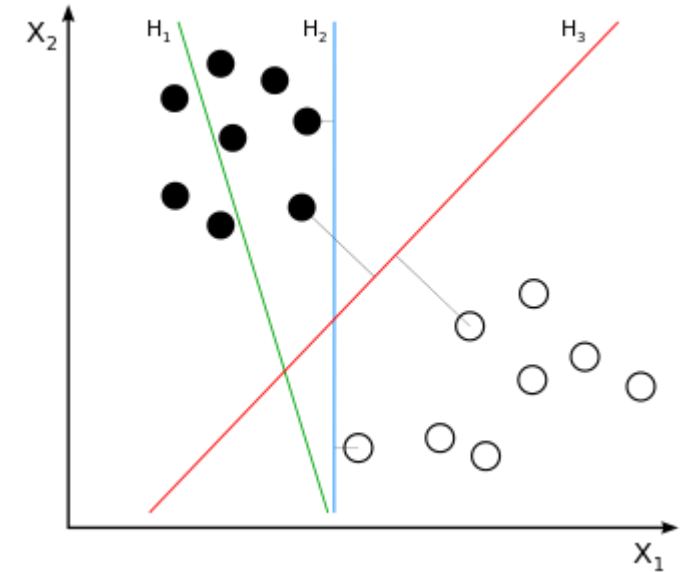
# Another AI Winter (late 1980's – early 1990s)

- Expert systems' progress on solving business problems slowed.

- Neural networks didn't scale to large problems.

- Interest in AI in business declined.

intel.

# Late 1990's to early 2000's: Classical Machine Learning



- Advancements in the Support Vector Machine algorithm led to it becoming the machine learning method of choice.



- The Deep Blue chess expert system beat world chess champion Garry Kasparov.
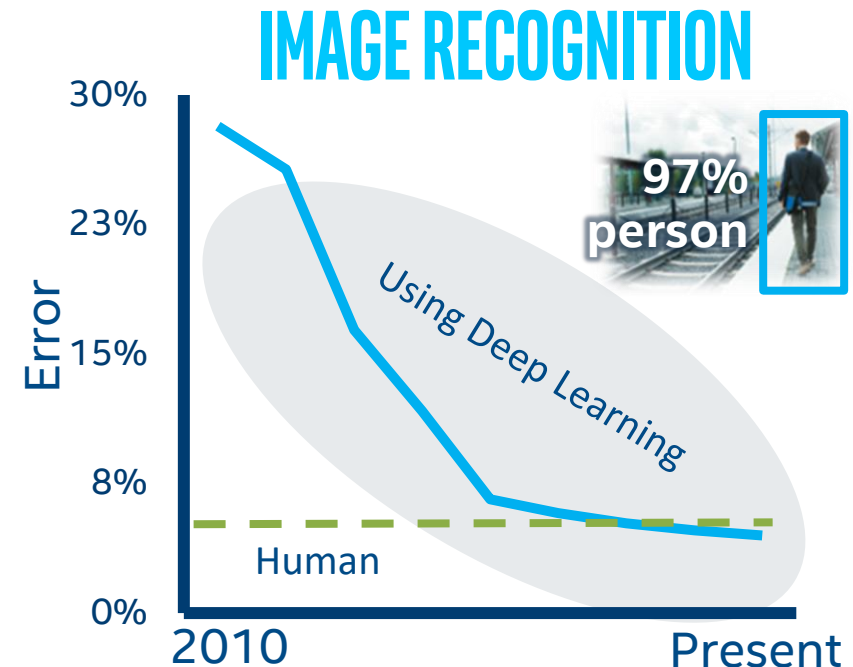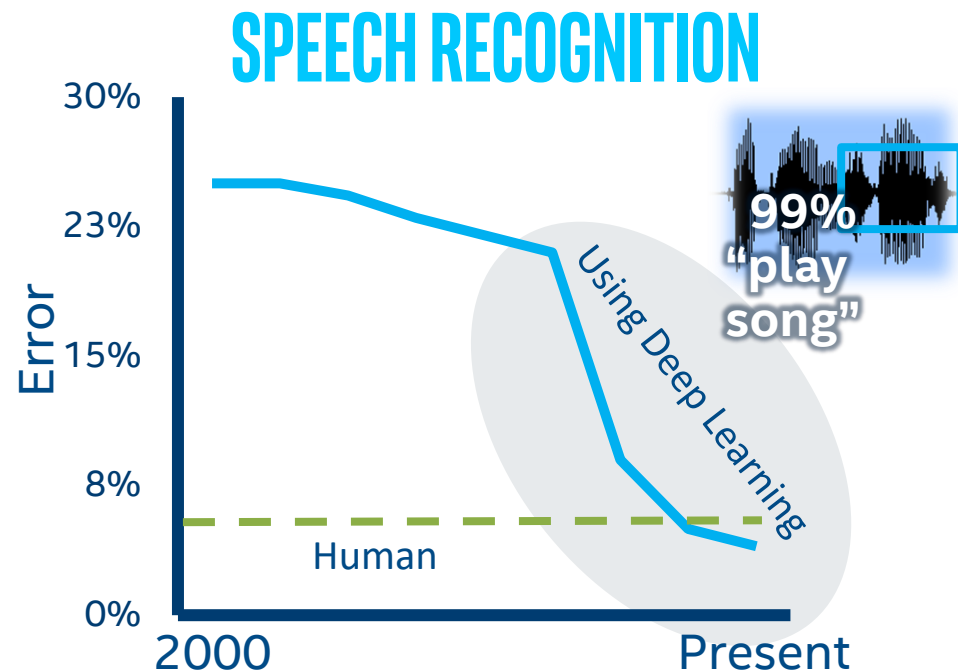
# 2000's: Rise of Deep Learning

- Model architectures already exist since the 80-90's such as CNN and RNN

- New architecture appears such as LSTM

- Bigger dataset and more powerful machines provide the essential tools for Deep Learning

- 2004: GPU implementation of Neural Networks

- 2009: The ImageNet database of human-tagged images is presented at the CVPR conference.

# Deep Learning – Breakthroughs (2012 – Present)

- 2011-2012: DanNet and AlexNet starts to win challenges in CV
- 2015: Deep learning platform TensorFlow is developed.



**SPEECH RECOGNITION**

Error — 30%, 23%, 15%, 8%, 0%

Using Deep Learning

99% "play song"

Human

2000 — Present



**IMAGE RECOGNITION**

Error — 30%, 23%, 15%, 8%, 0%

Using Deep Learning

97% person

Human

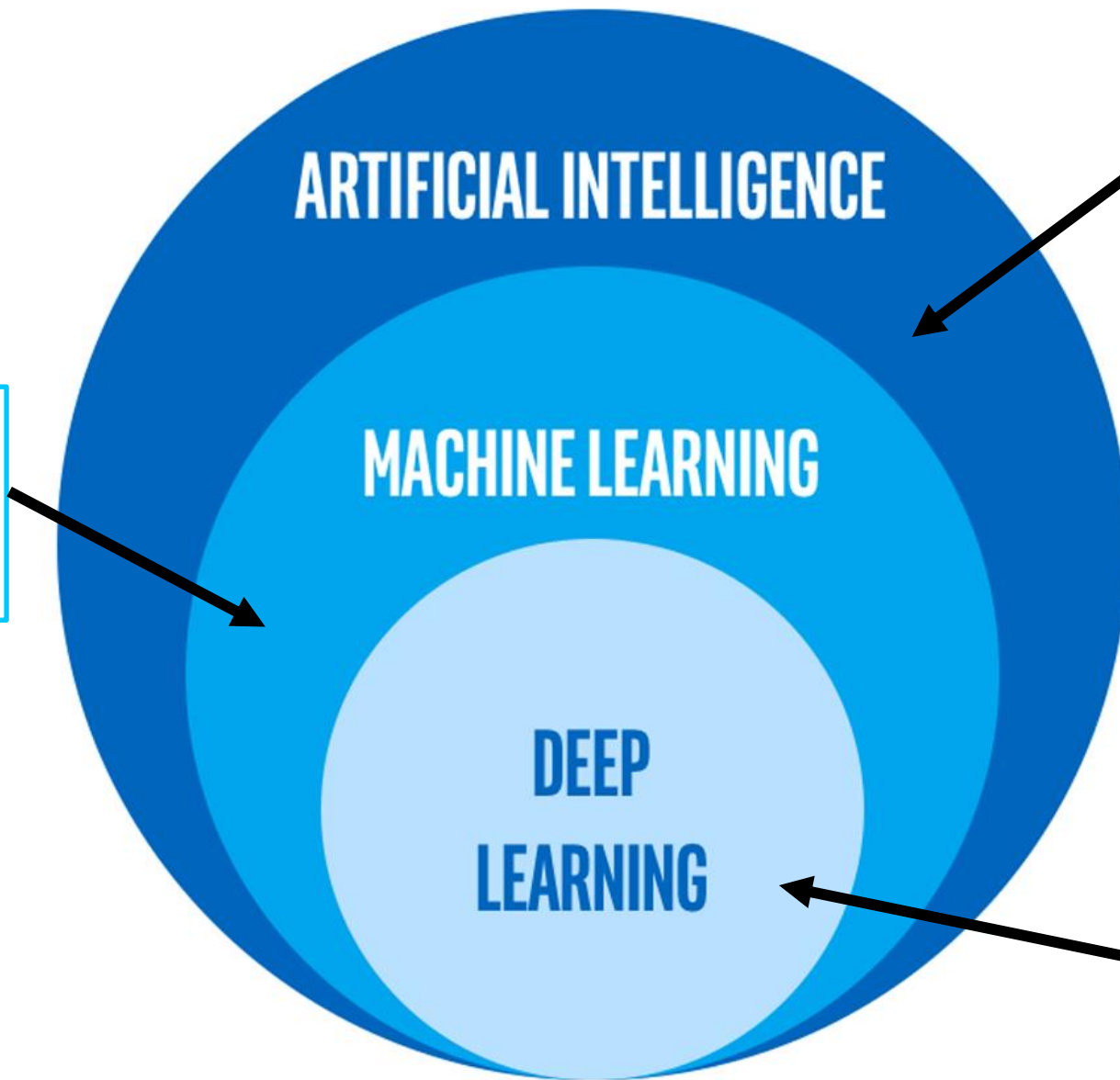2010 — Present

# Latest developments

- Deep Learning underwent tremendous developments in the last 5 years, largely driven by transformers

- Used in every vertical: healthcare, manufacturing, agriculture, surveillance, ...

- In 2022, ChatGPT and Stable Diffusion brought the attention of the public to AI

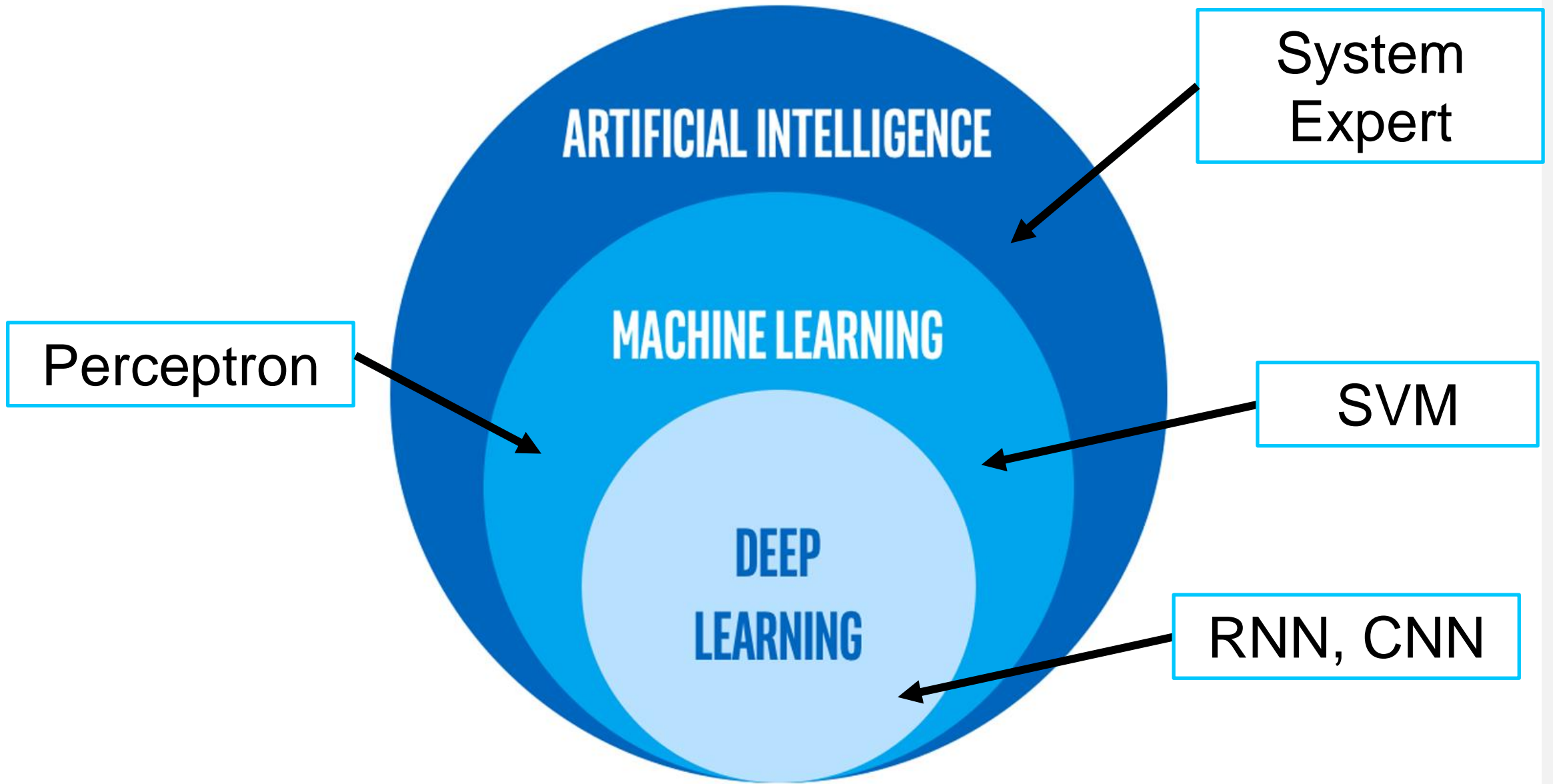## We are in the middle of an AI Summer !

# Definitions of AI concepts

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

DEEP LEARNING

System that leverage human knowledge for learning

System to perform a specific task without rule-based instruction

No more feature engineering, system learn hidden patterns

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

DEEP LEARNING
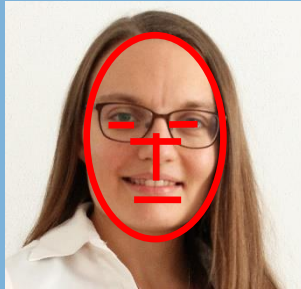
System Expert

Perceptron

SVM

RNN, CNN

intel.

# MACHINE LEARNING AND DEEP LEARNING

## CLASSICAL MACHINE LEARNING
How do you engineer the best features?

$N \times N$



$(f_1, f_2, ..., f_K)$

Roundness of face
Distance between eyes
Nose width
Eye socket depth
Cheek bone structure
Jaw line length
Etc.

### CLASSIFIER ALGORITHM

SVM
Random Forest
Naïve Bayes
Decision Trees
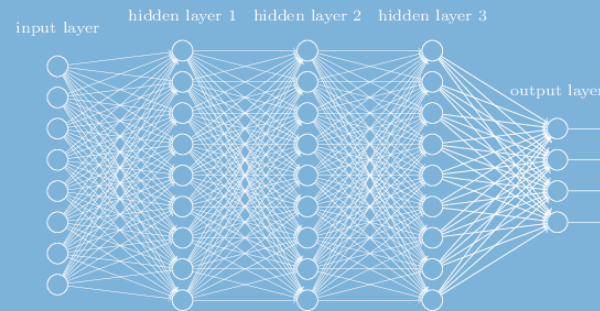Logistic Regression
Ensemble methods

**Séverine**

## DEEP LEARNING
How do you guide the model to find the best features?

$N \times N$



### NEURAL NETWORK

input layer    hidden layer 1   hidden layer 2   hidden layer 3

output layer

**Séverine**

# Two Main Types of Machine Learning

| | Goal | Common algorithms | Dataset |
|---|---|---|---|
| **Supervised Learning** | **Make predictions** | **Linear regression, decision tree, NN** | **Labelled** |
| **Unsupervised Learning** | **Find structure in the data** | **Clustering, kmeans, NN in contrastive learning** | **Unlabelled** |

# Classic ML example - Clustering

*An 'Unsupervised Learning' Example*

# Classic ML example – Decision tree



A 'supervised Learning' Example

FeatureA > threshold?

Yes    No

FeatureB > threshold2?

Yes    No

Class 1

Class 1    Class 2    Class 3

# Classic ML example – Gradient Boosting

*An 'supervised Learning' Example*

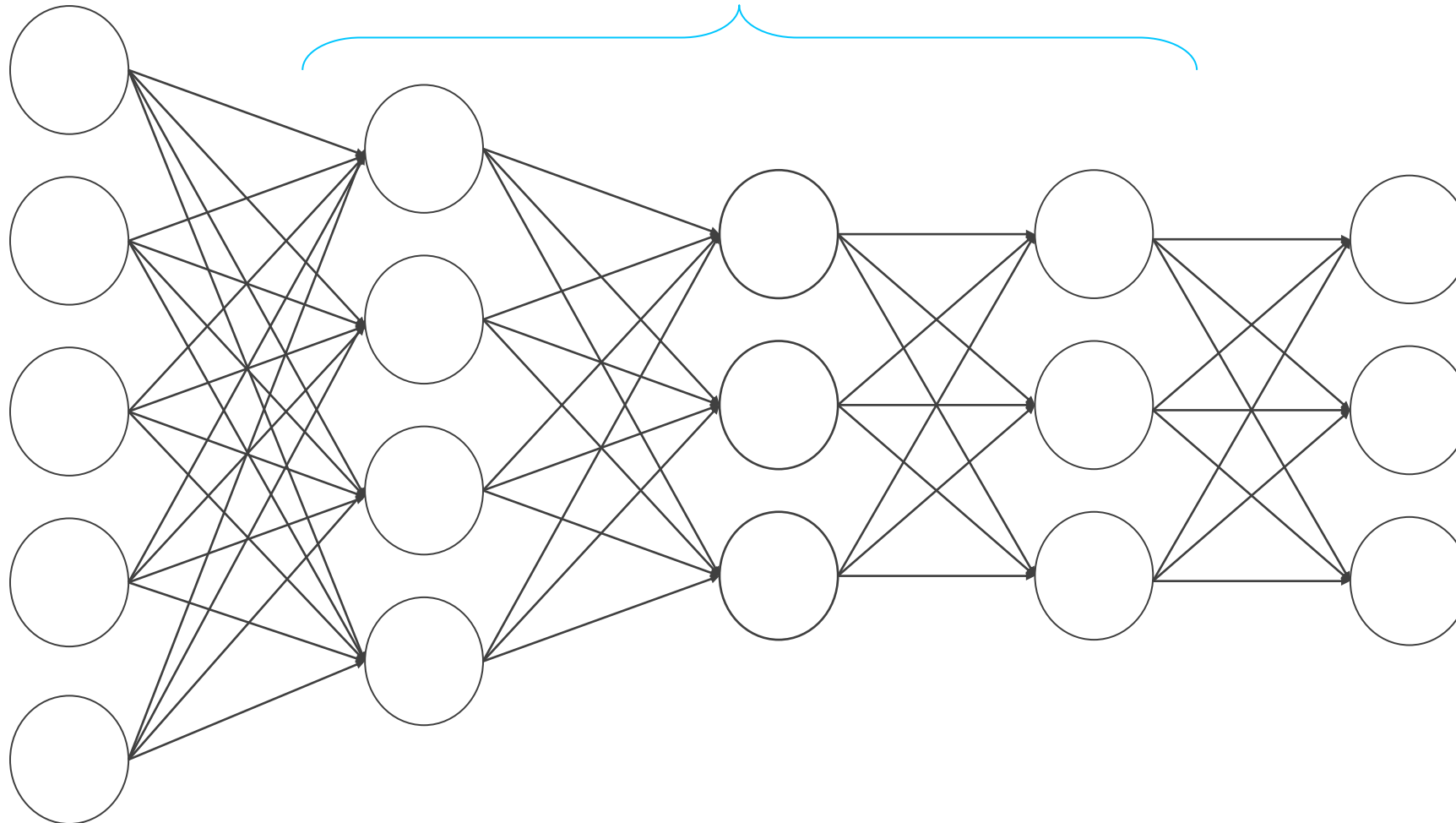Iteratively improve model predictions with gradients

# Deep Learning – Neural Network

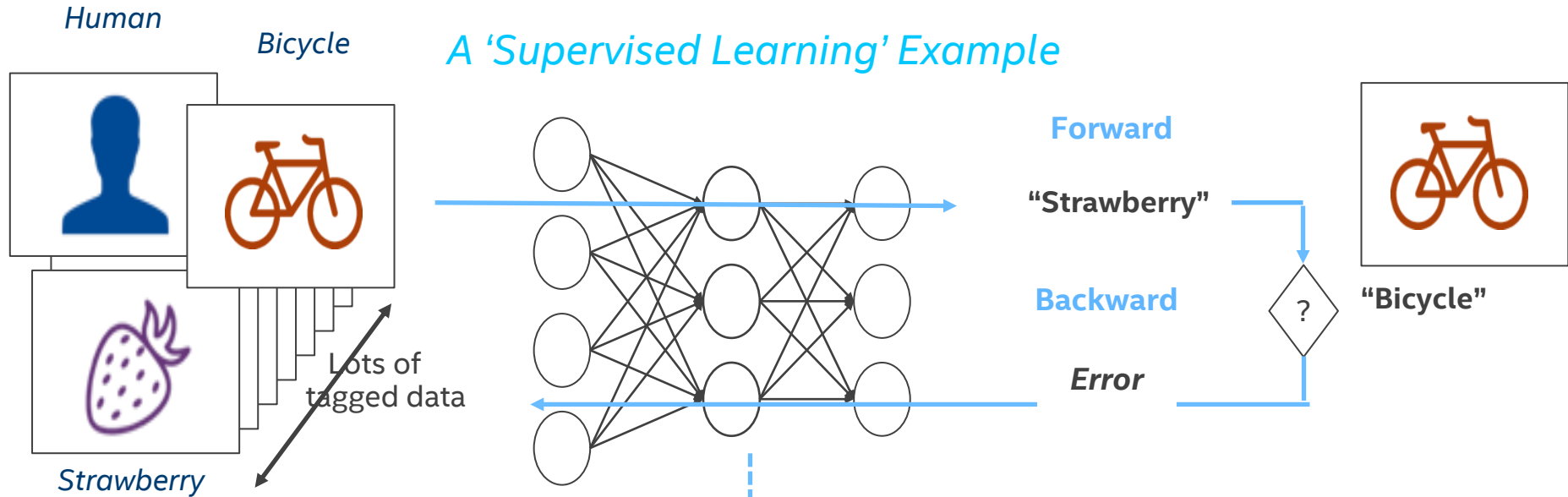**Input layer**  **Hidden layers**  **Output layer**

intel.

# Deep Learning – Training + Inference

# Deep Learning - typical domains

## COMPUTER VISION



Classification



Object Detection



Instance or Semantic segmentation

## NATURAL LANGUAGE PROCESSING

Speech and text tasks:

- Text Classification
- Named-entity recognition
- Machine Translation
- Chat-bots
- Automatic speech recognition
...

# Deep Learning - CNN

**Convolutional Neural Networks remain the gold standard for Vision. Transformers for Vision are picking up.**





From https://segment-anything.com/

intel.

# Deep Learning - Transformers

- Transformer architectures dominating (e.g., BERT, LLM, …)

- Hugging Face = *de facto* standard for model and dataset repo

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., … & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

# Deep Learning - Foundation models

- **A Foundation Models is a large AI model trained on enormous quantities of unlabeled data**
  - usually through self-supervised learning.
- **This one model can then be adapted to a wide range of downstream tasks**



Bommasani, Rishi, et al. "On the opportunities and risks of foundation models." *arXiv preprint arXiv:2108.07258* (2021).

# Deep Learning – Adaptation - Fine-tuning

**PRE-TRAINING**

Monday, we go to the Eiffel Tower

**Forward**

**Backward**

*Error*

?

Monday -> person
Eiffel Tower-> Location

**Model Weights**

**FINE-TUNING**

The concentration in $CO_2$ is very high

**Forward**

**Backward**

*Error*

?

$CO_2$ -> molecule

**INFERENCE**

$CH_4$ is released in the atmosphere

**Forward**

$CH_4$ -> molecule

intel.

# Deep Learning – and now?

GPT4?

**Large Language Models** - sorted by billion parameters

| **540**B | 176B | 100B | 20B | 1B |

PaLM

BLOOM

YaLM

GPT-NeoX

GPT-2

# Where ML shines

# Benefits of using classical ML

- Classical ML complements Deep Learning methods
- Shine where DL is weak:
  - Tabular data
  - Small dataset
  - Explainability
  - Computationally efficient

# Tabular data

- Most common data type, embedded in all relational database.



- Research has shown that gradient boosted tree ensemble models (XgBoost) are performing better than Deep learning on tabular data [1] and in particular on unseen datasets

[1] Shwartz-Ziv, R., & Armon, A. (2022). Tabular data: Deep learning is not all you need. Information Fusion, 81, 84-90.

# Small dataset

- Model with less parameters as found in classical ML will perform better on small dataset

- Less risk of overfitting

# Explainibility

- Features used by the algorithm are available and their influence can be analyzed using different techniques such as SHAP

Output = 0.4

Age = 65 →
Sex = F →
**Model**
BP = 180 →
BMI = 40 →

Base rate = 0.1

Explanation →

Output = 0.4

| +0.4 | ← Age = 65 |
| -0.3 | ← Sex = F |
| +.1 | ← BP = 180 |
| +.1 | ← BMI = 40 |

Base rate = 0.1

- Classical model can also be easier to apprehend

intel.

# Computationally efficient

- Can be trained on a CPU and run inference fast

- In order of minutes and not days as with DL

- Faster prototyping

# Intel Hardware for AI

# Flexible AI Acceleration

**CPU** *only*
Built-in AI acceleration for mainstream AI use cases

**CPU + GPU**
When compute is dominated by AI, HPC, graphics, and/or real-time media

**CPU + custom**
When compute is dominated by deep learning (DL)

| Cloud / Data Center | Edge | Device |
|---|---|---|
| intel. **XEON** | intel. **CORE** i7 | intel. **ATOM** |
| X^e | | |
| habana An Intel Company | intel. **AGILEX** | intel. **MOViDIUS** |
| DL Training/Inference | DL Custom | DL Inference |

# Intel® Xeon® Scalable Processors

## The Only Data Center CPU with Built-in AI Acceleration

Intel Advanced Vector Extensions 512

Intel Deep Learning Boost (Intel DL Boost)

Intel Advanced Matrix Extensions

### Shipping

**Cascade Lake**
New Intel DL Boost (VNNI)
New memory storage hierarchy

**Cooper Lake**
Intel DL Boost (BFLOAT16)

### April 2021

**Ice Lake**
Intel DL Boost (VNNI) and new
Intel Software Guard Extensions
(Intel® SGX) that enable new
AI use cases like federated learning

### January 2023

**Sapphire Rapids**
Intel Advanced Matrix Extensions (AMX)
extends built-in AI acceleration
capabilities on Xeon Scalable

## Leadership performance

# One Processor for Scalar, Vector, and Matrix

| Intel® AVX-512 | Intel® AVX-512 (VNNI) | Intel® AMX |
|---|---|---|
| 85 int8 ops/cycle/core with 2 FMA | 256 int8 ops/cycle/core with 2 FMAs | 2048 INT8 ops/cycle/core Multi-fold MACs in one instruction |

**Intel® AVX-512**

**Clock cycle 1**

8-bit input | 8-bit input

⊕ vpmaddubsw

16-bit output | 16-bit constant

**Clock cycle 2**

⊕ vpmaddwd

32-bit acc output | 16-bit constant

**Clock cycle 3**

⊕ vpaddd

32-bit acc output

**Intel® AVX-512 (VNNI)**

8-bit input | 8-bit input | 32-bit acc input

⊕ vpdpbusd

8-bit new instruction

32-bit acc output

**Intel® AMX**

8-bit input | 8-bit input | 32-bit acc input

⊕ tdpbusd

8-bit new instruction

32-bit acc output

# X$^e$ HPC (Ponte Vecchio)
## Leadership Performance for Data-level Parallel AI Workloads



>40 active tiles, over 100 billion transistors integrated into a single package



Powering New Phase of SuperMUC-NG at Leibniz Supercomputing Centre (LRZ)

https://www.youtube.com/watch?v=JzbN1IOAcwY

# Habana – an Intel Company



**Deep Learning ASIC for Training and Inference**

GAUDI® AI Training Card

GOYA

AI Inference Card

Gaudi accelerators in AWS EC2 instances, leverages up to 8 Gaudi accelerators and deliver up to 40% better price performance than current GPU-based EC2 instances for training

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.

# AI Software Stack for Intel XPUs

# oneAPI

## One Programming Model for Multiple Architectures & Vendors

### Freedom to Make Your Best Choice

- Choose the best accelerated technology the software doesn't decide for you

### Realize all the Hardware Value

- Performance across CPU, GPUs, FPGAs, and other accelerators

### Develop & Deploy Software with Peace of Mind

- Open industry standards provide a safe, clear path to the future
- Compatible with existing languages and programming models including C++, Python, SYCL, OpenMP, Fortran, and MPI

**Application Workloads Need Diverse Hardware**

| Scalar | Vector | Spatial | Matrix |
|--------|--------|---------|--------|

**Middleware & Frameworks**

Industry Initiative · **one**API · Intel Product

**XPUs**

| CPU | GPU | FPGA | Other accel. |
|-----|-----|------|--------------|

# Intel's oneAPI Ecosystem

## Built on Intel's Rich Heritage of CPU Tools Expanded to XPUs

**oneAPI**

A cross-architecture language based on C++ and SYCL standards

Powerful libraries designed for acceleration of domain-specific functions

A complete set of advanced compilers, libraries, and porting, analysis and debugger tools

**Powered by oneAPI**

Frameworks and middleware that are built using one or more of the oneAPI industry specification elements, the DPC++ language, and libraries listed on oneapi.com.



Application Workloads Need Diverse Hardware

Middleware & Frameworks (Powered by oneAPI)

TensorFlow    PyTorch    MODIN    learn    NumPy    XBOOST    OpenVINO    ...

**1**
**oneAPI**

**Intel® oneAPI Product**

| Compatibility Tool | Languages: Data Parallel C++ | Libraries | Analysis & Debug Tools |
|---|---|---|---|
| | | oneMKL  oneTBB  oneVPL  oneDPL | |
| | | oneDAL  oneDNN  oneCCL | |

Low-Level Hardware Interface

**XPUs**

CPU    GPU    FPGA    Other accelerators

[Available Now](#)

# Intel® AI Analytics Toolkit

## Powered by oneAPI

Accelerate end-to-end AI and data analytics pipelines with libraries optimized for Intel® architectures

### Who Uses It?

Data scientists, AI researchers, ML and DL developers, AI application developers

### Top Features/Benefits

- Deep learning performance for training and inference with Intel optimized DL frameworks and tools

- Drop-in acceleration for data analytics and machine learning workflows with compute-intensive Python packages

Learn More: software.intel.com/oneapi/ai-kit

| Deep Learning | Data Analytics & Machine Learning |
|---|---|
| Intel® Optimization for TensorFlow | **Accelerated Data Frames** |
| | Intel® Distribution of Modin / OmniSci Backend |
| Intel® Optimization for PyTorch | **Intel® Distribution for Python** |
| Intel® Neural Compressor | XGBoost / Scikit-learn / Daal-4Py |
| Model Zoo for Intel® Architecture | NumPy / SciPy / Pandas |

### Samples and End2End Workloads

CPU     GPU

Supported Hardware Architechures[1]

Hardware support varies by individual tool. Architecture support will be expanded over time. Other names and brands may be claimed as the property of others.

### Get the Toolkit HERE or via these locations

| Intel Installer | Docker | Apt, Yum | Conda | Intel® DevCloud |
|---|---|---|---|---|

# Intel oneAPI Software Tools for AI & Analytics

## Intel® oneAPI Toolkits

### Intel® AI Analytics Toolkit

Accelerate machine learning & data science pipelines with optimized deep learning frameworks & high-performing Python libraries

Data Scientists, AI Researchers, DL/ML Developers

### Intel® oneAPI Base Toolkit

Incl. Intel® oneAPI Deep Neural Network Library (oneDNN), Intel® oneAPI Collective Communications Library (oneCCL), & Intel® oneAPI Data Analytics Library (oneDAL)

Optimize primitives for algorithms and framework development

DL Framework Developers - Optimize algorithms for Machine Learning & Analytics

## Toolkit **Powered by oneAPI**

### Intel® Distribution of OpenVINO™ Toolkit

Deploy high performance inference & applications from edge to cloud

AI Application, Media, & Vision Developers

# oneAPI Ecosystem Endorsements for AI domain

The industry needs a programming model where developers can take advantage of an array of innovative hardware architectures. The goal of oneAPI is to provide increased choice of hardware vendors, processor architectures, and faster support of next-generation accelerators. Microsoft has been using oneAPI elements across Intel hardware offerings as part of its initiatives and supports the open standards-based specification.  We are excited to support our customers with choice and accelerate the growth of AI and machine learning.
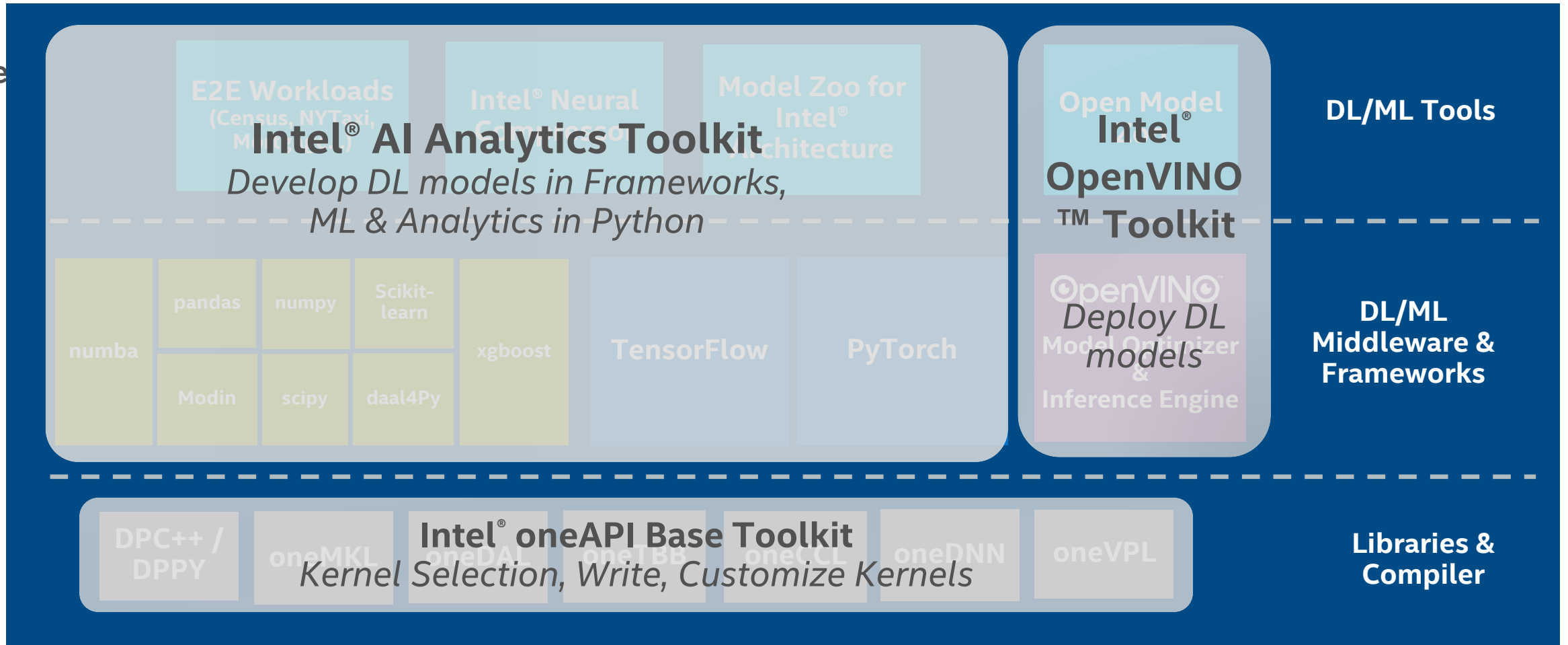*- Tim Harris, Principal Architect, Azure AI, Microsoft*

With the growth of AI, machine learning, and data-centric applications, the industry needs a programming model that allows developers to take advantage of rapid innovation in processor architectures. TensorFlow supports the oneAPI industry initiative and its standards-based open specification. oneAPI complements TensorFlow's modular design and provides increased choice of hardware vendor and processor architecture, and faster support of next-generation accelerators. TensorFlow uses oneAPI today on Xeon processors and we look forward to using oneAPI to run on future Intel architectures.

# AI Software Stack for Intel XPUs

▪ Inte

**Intel® AI Analytics Toolkit**
*Develop DL models in Frameworks, ML & Analytics in Python*

E2E Workloads (Census, NYTaxi, MLPerf)

Intel® Neural Compressor

Model Zoo for Intel® Architecture

Open Model

**Intel® OpenVINO™ Toolkit**

**DL/ML Tools**

numba

pandas

numpy

Scikit-learn

Modin

scipy

daal4Py

xgboost

TensorFlow

PyTorch

OpenVINO
*Deploy DL models*
Model Optimizer & Inference Engine

**DL/ML Middleware & Frameworks**

**Intel® oneAPI Base Toolkit**
*Kernel Selection, Write, Customize Kernels*

DPC++ / DPPY

oneMKL

oneDAL

oneTBB

oneCCL

oneDNN

oneVPL

**Libraries & Compiler**

Full Set of Intel oneAPI cross-architecture AI ML & DL Software Solutions

# oneAPI Available on
# Intel® DevCloud for oneAPI

A development sandbox to develop, test and run workloads across a range of Intel CPUs, GPUs, and FPGAs using Intel's oneAPI software.

## Get Up & Running In Seconds!

Sign up at:
**software.intel.com/devcloud/oneapi**

**intel.**
## DevCloud

**1 Minute to Code**

**No Hardware Acquisition**

**No Download, Install or Configuration**

**Easy Access to Samples & Tutorials**

**Support for Jupyter Notebooks, Visual Studio Code**

# High-Performance Deep Learning Using Intel® Distribution of OpenVINO™ toolkit - **Powered by oneAPI**

A toolkit for fast, more accurate real-world results using high-performance AI and computer vision inference deployed into production on Intel XPU architectures (CPU, GPU, FPGA, VPU) from edge to cloud
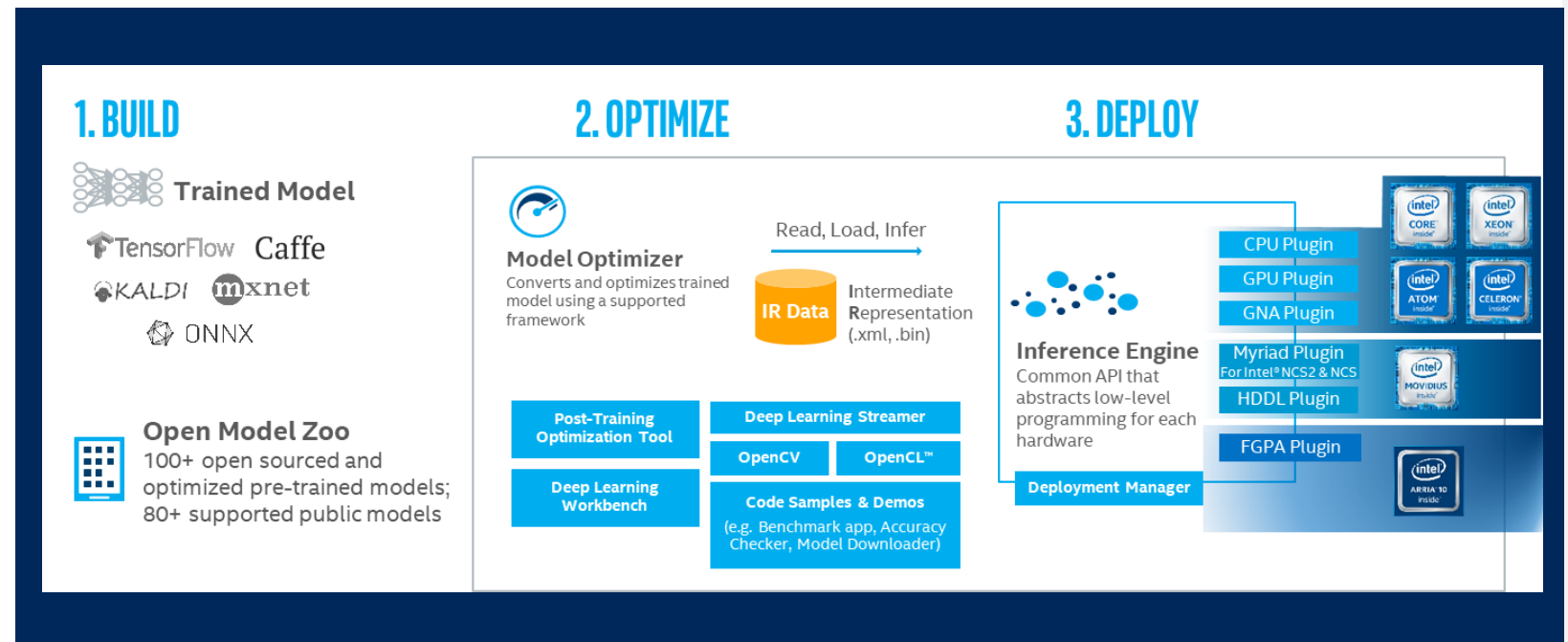
## Who needs this product?

AI application developers, OEMs, ISVs, System Integrators, Vision and Media developers

## Top Features/Benefits

High-performance, deep learning inference deployment

Streamlined development; ease of use
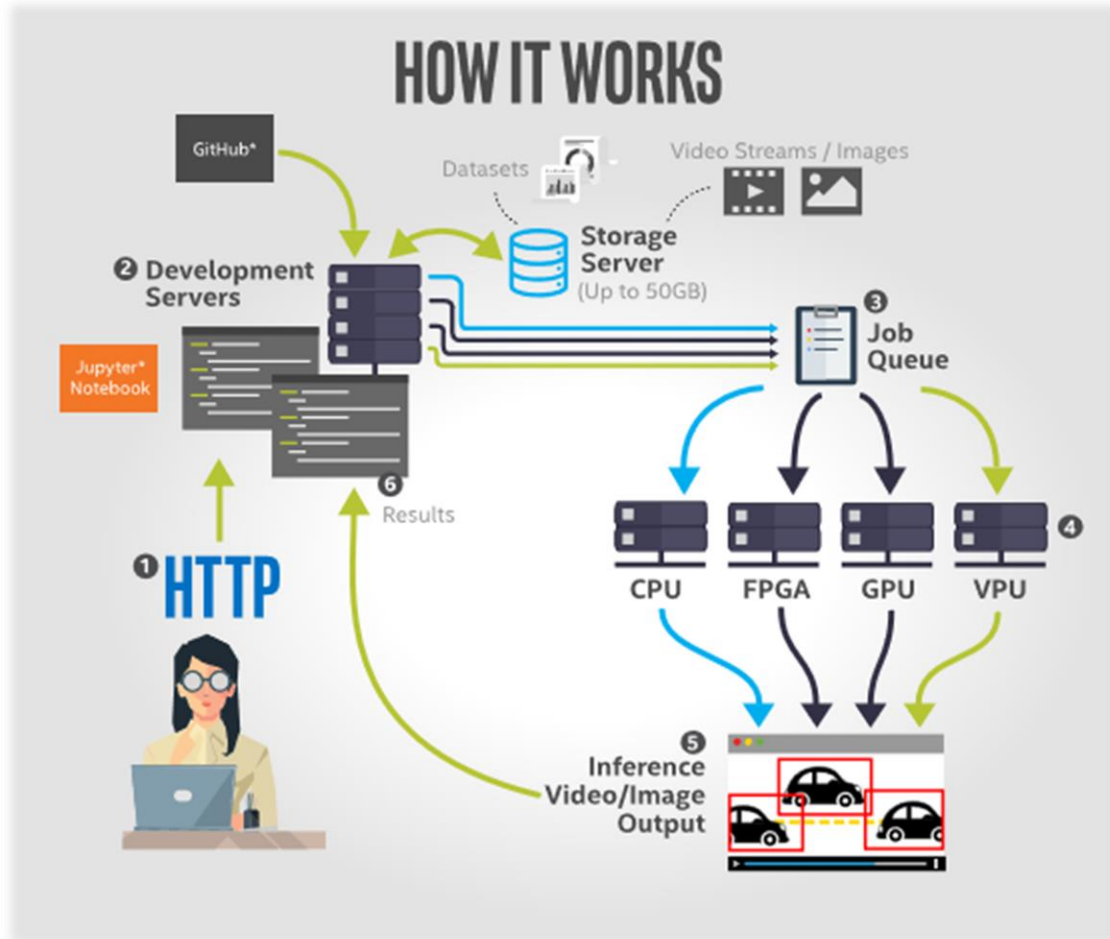
Write once, deploy anywhere

Proven, industry-leading accelerated technology



## software.intel.com/openvino-toolkit

# Accelerate Time to Production with Intel® DevCloud for the Edge

## See immediate AI Model performance across Intel's vast array of Edge Solutions



**HOW IT WORKS**

- **Instant, Global Access**
  Run AI applications from anywhere in the world

- **Prototype on the Latest Hardware and Software**
  Develop knowing you're using the latest Intel technology

- **Benchmark your Customized AI Application**
  Immediate feedback – frames per second, performance

- **Reduce Development Time and Cost**
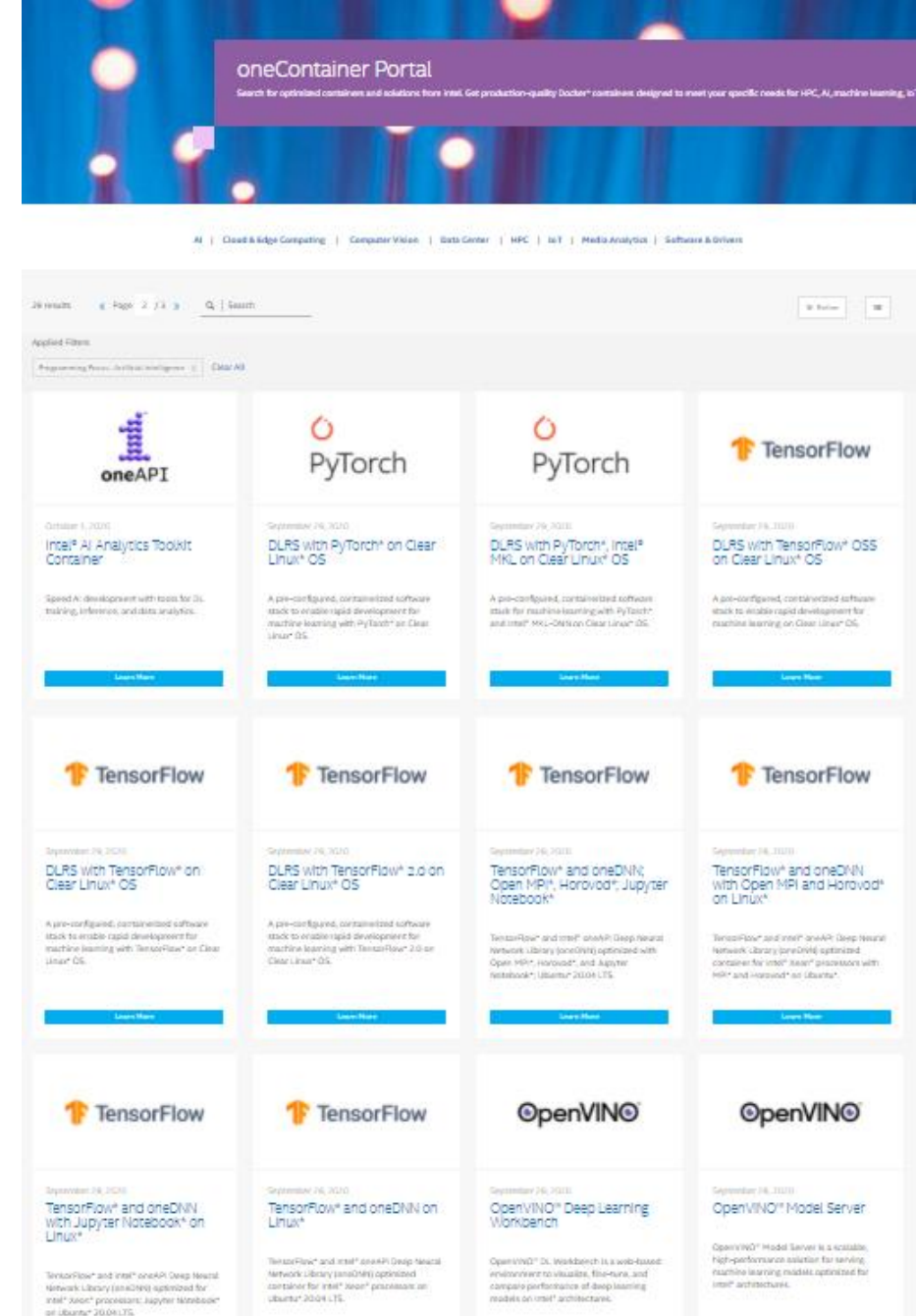  Quickly find the right compute for your edge solution

**Sign up now for access**

# AI Containers for Flexibility

- Optimized, validated, deployable AI containers

- Available via Docker containers. Will expand to include Kubernetes orchestrations, Helm charts

- Access from oneContainer Portal

  - Include containers with ready-to-use AI software stacks

  - And containers with full AI workloads (including models)

| Topology | Frameworks |
|---|---|
| DLRM | PYT |
| ResNet50 | PYT, TF, OV |
| BERT-large | PYT, TF, OV |
| Transformer-LT | PYT, TF |
| MobileNet-v1 | PYT, TF, OV |
| SSD-Mobilenet-v1 | PYT, TF, OV |
| SSD-Resnet34 | PYT, TF, OV |
| WaveNet* | TF |

| Topology | Framework |
|---|---|
| Mask R-CNN | PYT, TF, OV |
| RNN-T | PYT, TF, OV |
| 3D-UNet | TF, OV |
| DIEN | TF |
| Wide & Deep | PYT, TF |
| RNX101 | |
| Yolo-V3 | PYT, TF, OV |
| NCF* | TF |

# Which Toolkit Should I Use

intel.

# Use Both!
## Intel® oneAPI Analytics Toolkit & Intel® Distribution of OpenVINO™ toolkit

**Toolkits are complementary to each other and recommendation is to use them both based on your current phase of AI Journey**

- *I am **exploring and analyzing data;** I am **developing models***
- *I want **performance and compatibility** with frameworks and libraries I use*
- *I would like to have **drop-in acceleration** with little to no additional code changes*
- *I prefer **not to learn any new tools** or languages*

- *I am **deploying models***
- *I want **leading performance and efficiency** across multiple target HW*
- *I'm concerned about **having lower memory footprint**, which is critical for deployment*
- *I am **comfortable with learning and adopting a new tool or API** to do so*

**Data Scientist/ML Developer**
Intel® oneAPI AI Analytics Toolkit

**App Developer**
Intel® Distribution of OpenVINO™ toolkit

If you prefer working on primitives and to optimize kernels and algorithms directly using oneAPI libraries (oneDNN, oneCCL & oneDAL), then use Intel® oneAPI Base Toolkit
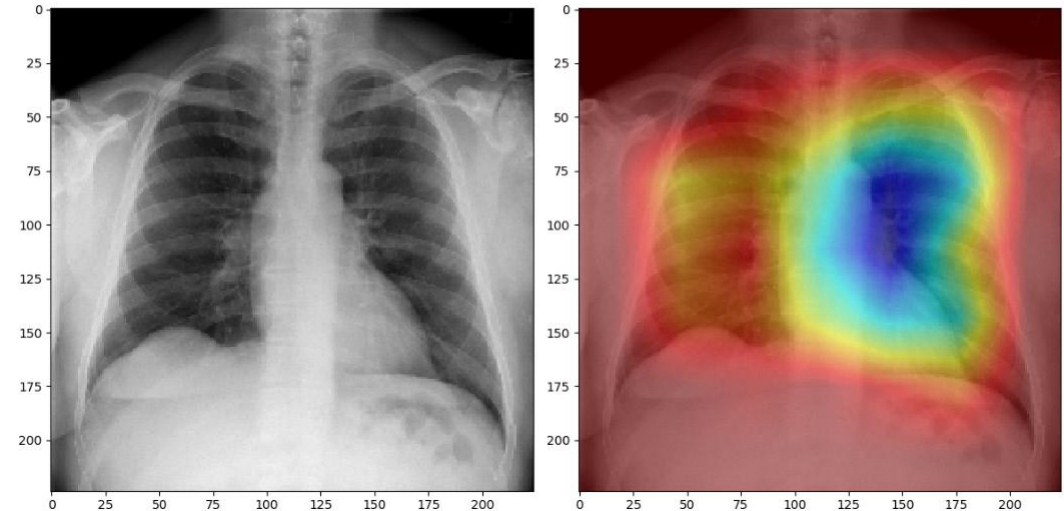
# Accrad AI-based Solution Helps Accelerate COVID-19 Diagnosis

**Optimized by Intel® oneAPI Analytics Toolkit & Intel® Distribution of OpenVINO™ toolkit**

*CheXRad* helps radiologists and physicians identify COVID-19, viral pneumonia and other diseases on chest X-ray images, and predict the need for ventilators.

- *CheXRad* comes pre-configured with a COVID-19 and viral pneumonia classification neural network.

- To architect, train and validate the neural network, Accrad used **Intel Tensorflow from AI Analytics Toolkit** and the **Intel oneAPI DevCloud** to develop the model.

- To optimize its model for deployment, Accrad used **OpenVINO™ toolkit** and **Intel® DevCloud for Edge**.

- *CheXRad* could classify pathologies in 140 chest x-rays in just **90 seconds** —up to **160x faster** than radiologists, at comparable levels of accuracy, sensitivity and specificity.



Ground Truth Class: 0 (non-COVID-19)
Predicted Class: 0 (non-COVID-19)
Prediction probabilities: ['1.00', '0.00']

Learn more in this solution brief

# Key Takeaways & Call to Action

- Intel toolkits are **FREE**, complementary & work seamlessly together
- They help achieve performance & efficiency across different stages of AI Journey
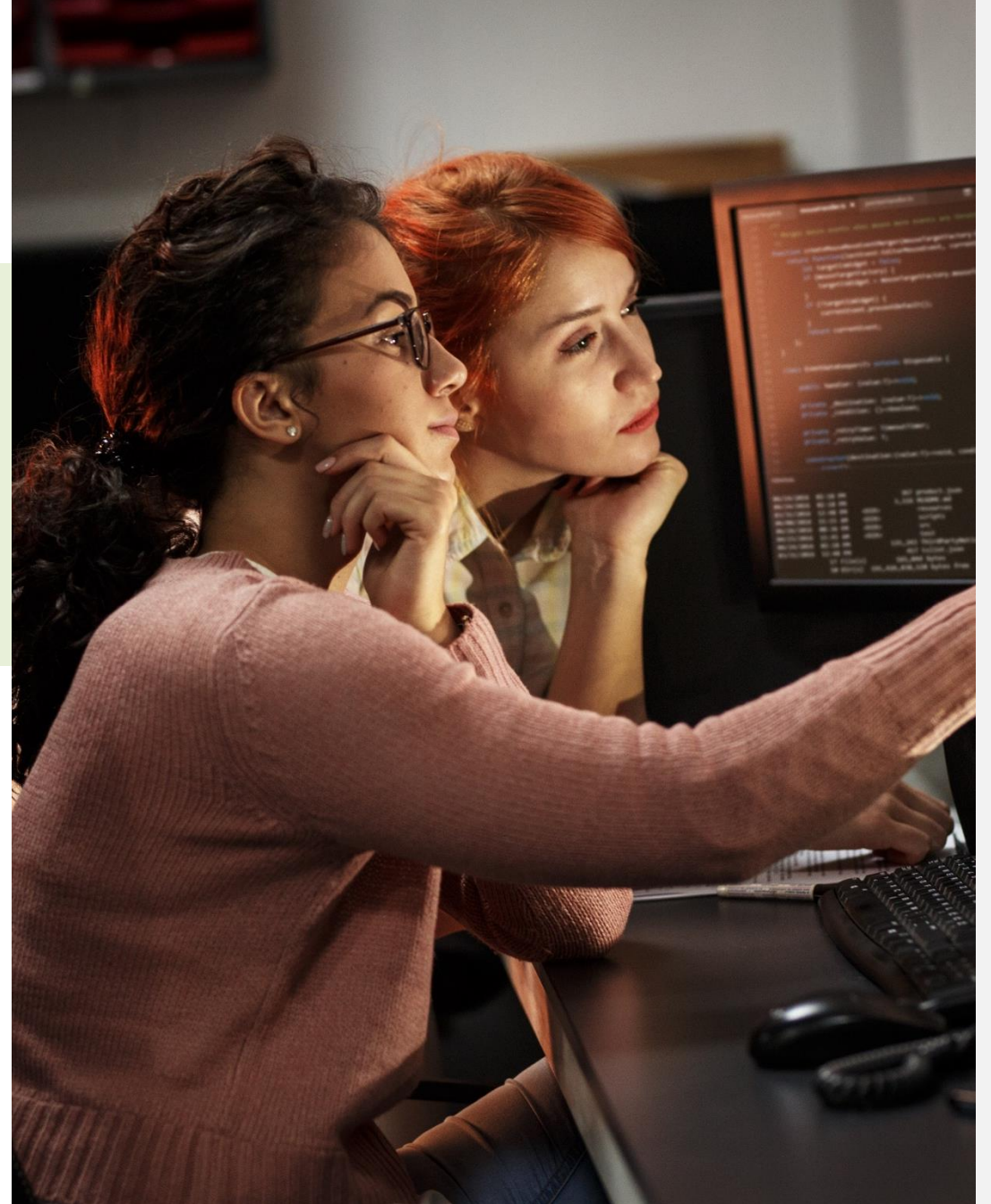- Recommend the toolkits based on current phase of customer pipeline

Download the toolkits

Intel® oneAPI AI Analytics Toolkit

Intel® Distribution of OpenVINO ™ toolkit

Intel® oneAPI Base Toolkit

Learn more about Intel® oneAPI Toolkits
intel.com/oneAPI-AllToolkits

# May 16th
# Intel Deep Learning workshop
# @ LRZ

[link](#)

intel

# Survey

https://survey.lrz.de/index.php/582495?lang=en