# Agenda

- Edge AI applications

- The OpenVINO™ toolkit: An open-standard for building AI at the edge, in the cloud, or locally

- Deep-dive: Enterprise intelligence and Intel®'s portfolio

- Flexible edge and cloud computing paradigms

- Q&A

# AI is Everywhere – from Edge to Cloud

## Edge

- Real time data processing
- Wider reach
- Data sovereignty
- Cost efficiency

Frictionless Retail

Traffic Monitoring

Defect Detection

Depth Camera

IP Camera

NVR

Sensor

OpenVINO
oneAPI

OpenVINO  intel

# Enterprise Intelligence at the Edge with Intel®
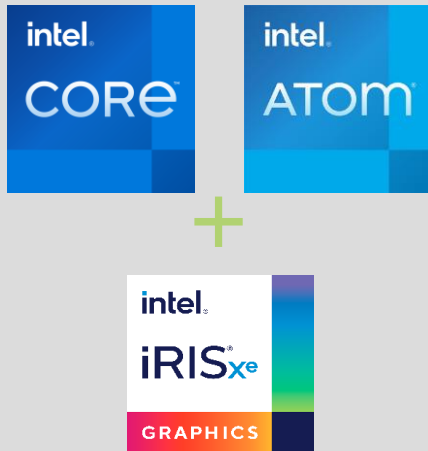
## Node
### Fine-tuning, Inference

## Cluster
### Light Training, Tuning, Peak Inf.

# Choose Your Compute

## Light AI

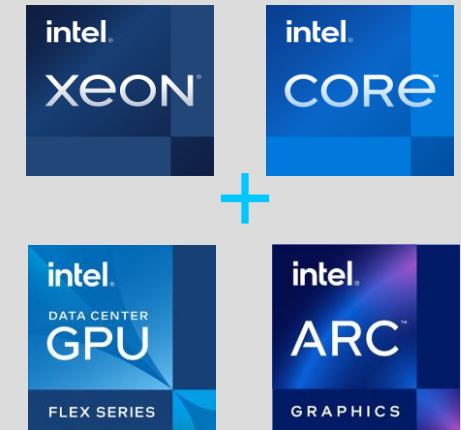Efficiency for sub-100 W designs

CPU AI, built-in GPU, built-in NPU

intel CORE    intel ATOM

+

intel iRISxe GRAPHICS


Anomaly

## Medium AI

Scale up perf/W for diverse system designs

Discrete GPU & built-in CPU AI

intel XEON    intel CORE

+

intel ARC GRAPHICS    intel iRISxe GRAPHICS


Small Gen AI

## Heavy AI

Optimize for peak perf and density

Optimize with high-end Discrete GPU

intel XEON    intel CORE

+

intel DATA CENTER GPU FLEX SERIES    intel ARC GRAPHICS

OpenVINO  intel

© 2024 Intel

# OpenVINO: An Open Standard for Building AI at the Edge

© 2024 Intel

# Unlocking Software Optimizations with Hardware Using the Intel® Edge AI Portfolio

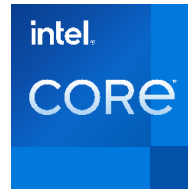| | | | |
|---|---|---|---|
| Supports myriad AI use cases | Wide range of AI performance | Real-time and offline execution | **Software Benefits with OpenVINO™** |
| Open-source for AI, DL, Inference | Performance-optimized | Cross-platform Support | **OpenVINO™ Toolkit Value Differentiators** |
| **Flexible video streams**<br><br>Ranges from few to many streams, and TOPS in low-end NVRs to TOPS for on-prem servers | **Temperature range**<br><br>Ranges from -5 to +105 degree C | **TDP/Power**<br><br>Ranges from sub 10 W to over 300 W | **Hardware Benefits with Intel®'s Edge AI Portfolio** |

OpenVINO™  intel®

# Compute for AI: Intel®'s Platforms



## Edge AI Platforms

Partner edge platforms using Intel® Arc™ GPU

## Intel® Edge AI Box

- Single NN Pipeline
- Multi NN Pipeline
- Data Fusion Pipeline

Seamlessly Integrated into Existing Camera and Video Deployments

© 2024 Intel

8

# Challenges and Opportunities
# with AI at the Edge

OpenVINO intel.

# Intelligent Queue Management at the Edge with OpenVINO™ and YOLOv8

**Objective:** Optimize the queuing process and reduce wait times via object detection.

**Challenges:**
- Real-time scalability
- Device setup and calibration
- Model performance
- Low-power

**Intel's solution:**
- Fast and efficient inference with optimized YOLOv8 models using OpenVINO

# Mobile Multi-modal Assistant with MobileVLM and OpenVINO™

**Objective:** Use a mobile chatbot to answer questions about images

**Challenges:**
- Fast token generation
- Memory-efficiency
- Model size

**Intel's solution:**
- Compress and quantize LLM models for faster, efficient local inference



**GPU**

# Snapshot of MobileVLM Output

# Document Visual Question Answering with Pix2Struct and OpenVINO™ on CPU

**Objective:** Use a low-power chatbot to answer questions about documents on the fly.

**Challenges:**
- Visual understanding
- Memory-efficiency
- Model size

**Intel's solution:**
- Compress and quantize mult modal models for faster, efficient local inference



**Intel® Core™ Ultra 9 processor CPU**

# Snapshot of Pix2Struct Output

© 2024 Intel
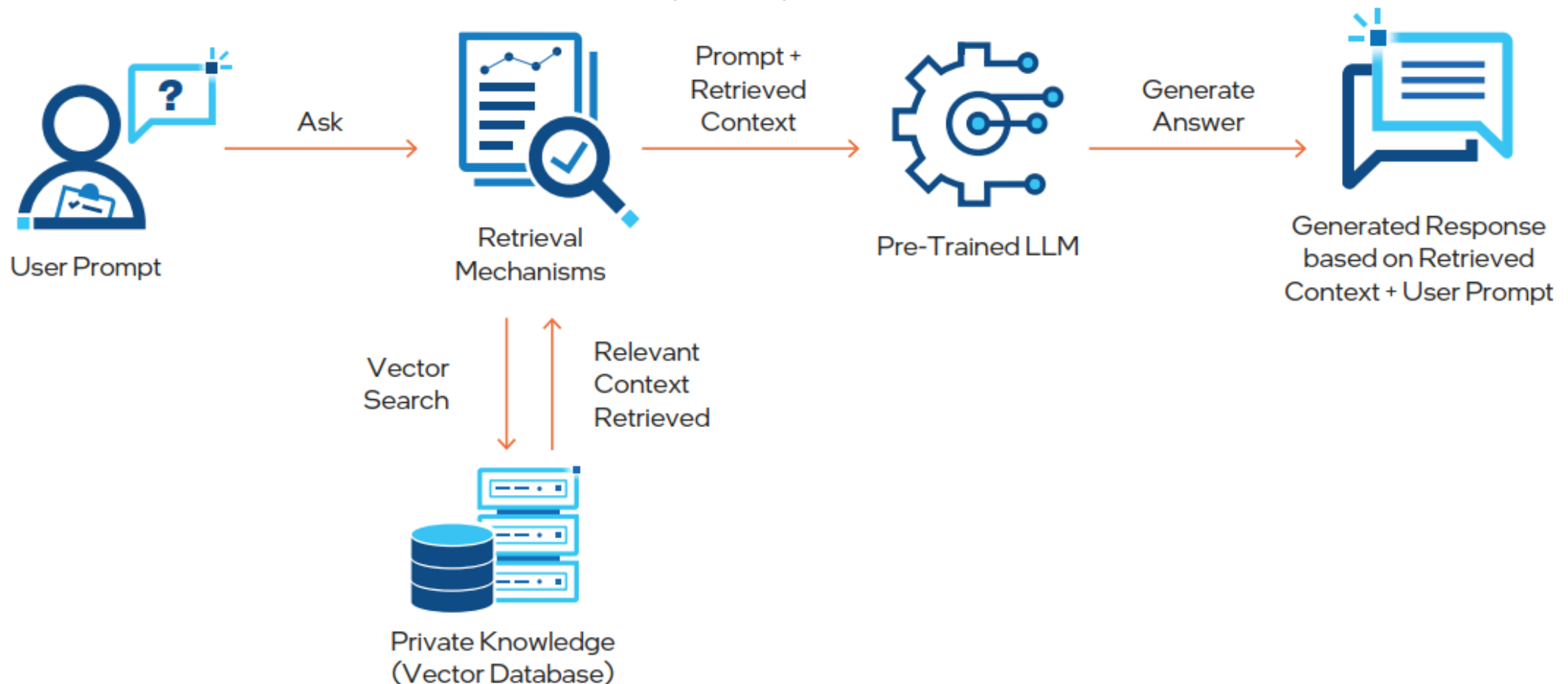
# Enterprise Intelligence with LLMs using RAG

Connect knowledge bases to LLMs with Retrieval Augmented Generation (RAG)

# Running LLM + RAG with OpenVINO™ and LangChain on iGPU for the edge

**Enable enterprise intelligence through knowledge-based search**
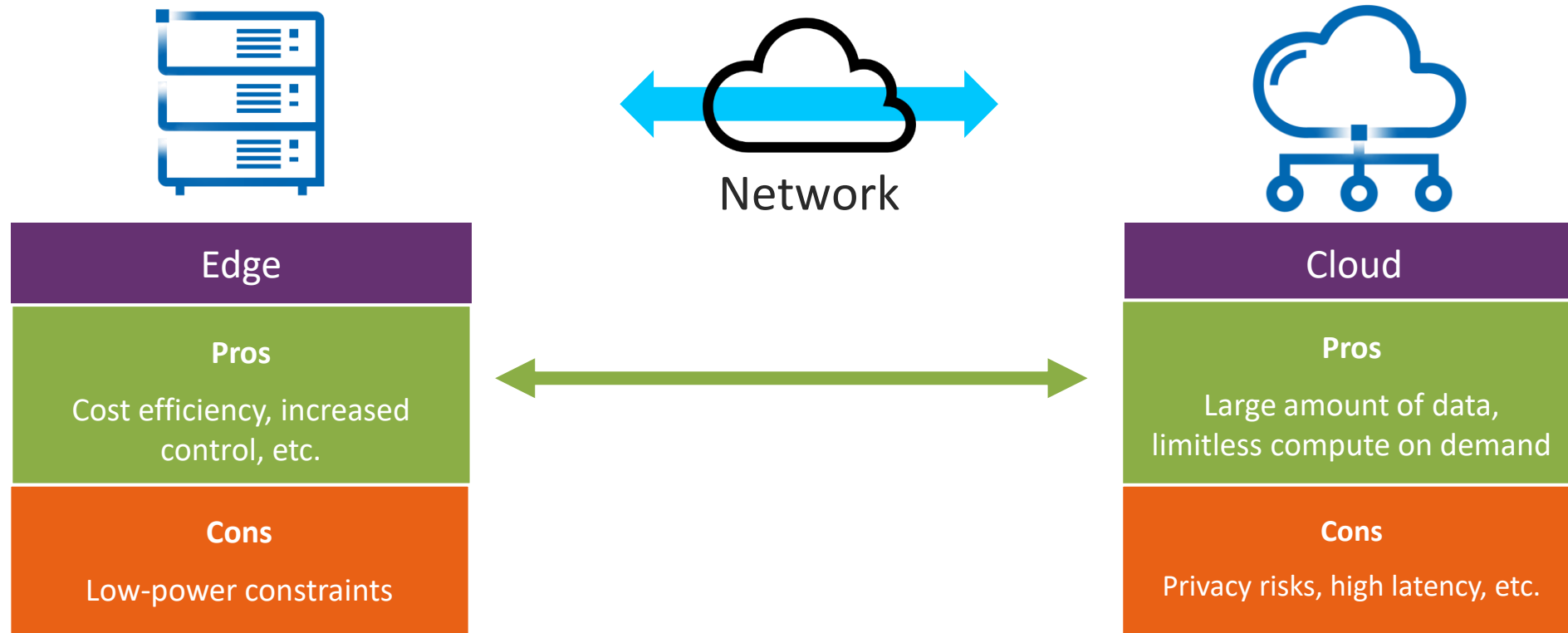
# Enterprise Data Protection at the Edge



Intel® Software Guard Extensions (Intel® SGX)

Secure Access Service Edge (SASE)

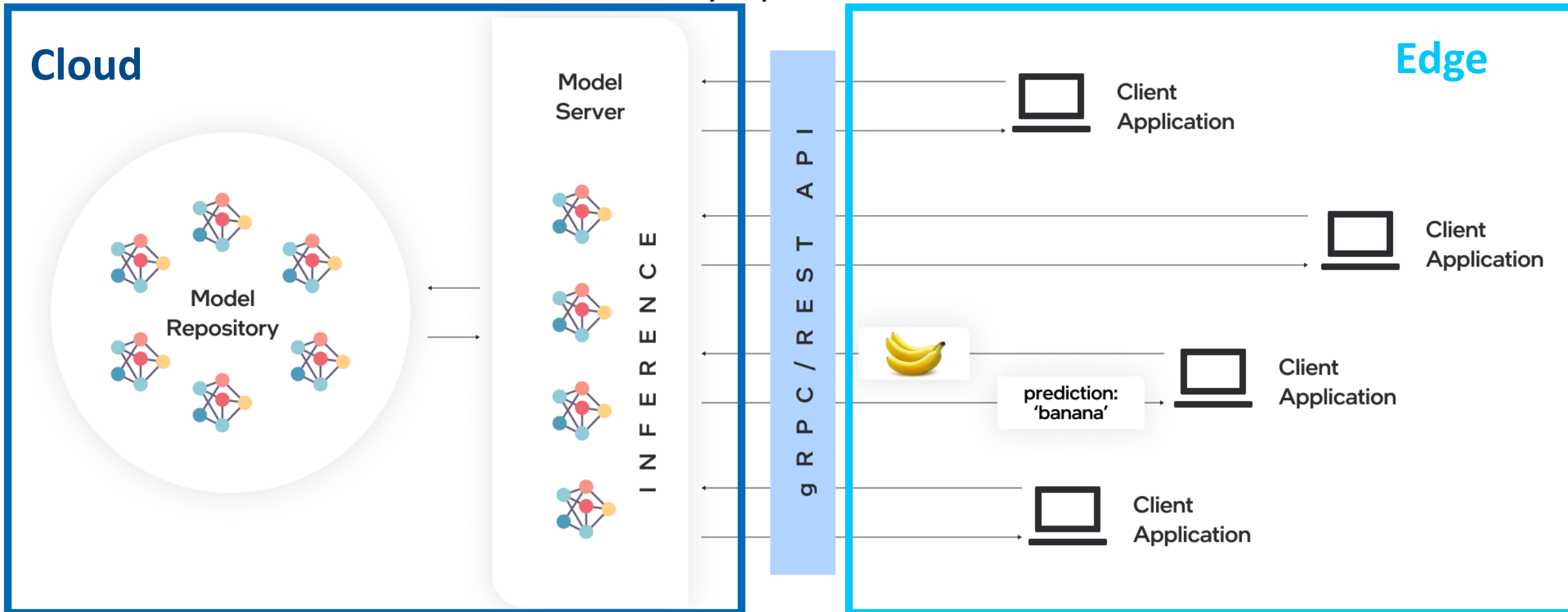Intel® QuickAssist Technology (Intel® QAT)

# Edge to Cloud Paradigms

# Edge to Cloud: Flexibly Using Compute



Network

| Edge |
|------|
| **Pros** |
| Cost efficiency, increased control, etc. |
| **Cons** |
| Low-power constraints |

| Cloud |
|-------|
| **Pros** |
| Large amount of data, limitless compute on demand |
| **Cons** |
| Privacy risks, high latency, etc. |

# Edge to Cloud with OpenVINO™ Model Server

Move workloads across the edge and cloud
Powered by OpenVINO™ Runtime

© 2024 Intel

# LLM Assistants: OpenVINO™ Model Server with INT8 Compression



**Deploying a Quantized Tiny-llama model across client and server**

# Intel®'s AI Hardware Portfolio



## Edge AI Platforms

Partner edge platforms using Intel® Arc™ GPU

## Cloud Platforms

## Client Platforms

Desktops

Laptops

# Conclusion

- AI at the edge is transforming enterprise intelligence

- But not without several challenges: scalability, setup, AI performance, etc.

- At Intel®, we see the full end-to-end stack as key for optimizing AI at the edge, and across the cloud to edge

- OpenVINO is an open standard, ready-to-use for building AI and Gen AI

- Try It Yourself: openvino.ai

# Resources

## Resources

- [openvino.ai](openvino.ai)

- [intel.com/edgeai](intel.com/edgeai)

- Demos: [intel.com/openvinonotebooks](intel.com/openvinonotebooks)

- [Enterprise Security Solutions at the Edge](#) with Intel

## 2024 Embedded Vision Summit

May 23[rd] (12:00 pm – 12:30 pm)

"Identifying and Mitigating Bias in AI"

May 23[rd] (1:30 pm – 2:00 pm)

"Intel's Approach to Operationalizing AI in the Manufacturing Sector"