

AI acceleration at the edge

Jonathan Lee

Chief Product Officer
ai.io

Alex White

Sales Applications Engineer
Intel

Michael Kleiner

VP Edge AI Solutions
OnLogic

Mohan Potheri

Cloud Solutions Architect
Intel

Allan Gagnon

Senior Solutions Architect
Arduino

AWS re:Invent 2023 - AI acceleration at the edge (AIM311)

 AWS Events
100K subscribers

[Subscribe](#)

 0 |  |  Share |  Clip |  Save | ...

53 views Dec 3, 2023 #AWSreInvent #AWSreInvent2023

AI has the power to transform and create innovative business models that generate real value. Companies can build models with minimal data and, with Intel domain-specific toolkits, facilitate deploying solutions at scale. Listen as industry experts and customers OnLogic, Arduino, and ai.io describe how they were able to use AI technology at the edge to achieve strong business outcomes in agriculture, manufacturing, and automotive by optimizing customer experiences, forecasting, and managing inventory. Come away from this session with tangible examples of Intel software and hardware optimizations that led to creative, successful AI solutions at the edge. This presentation is brought to you by Intel, an AWS Partner.

Learn more about AWS re:Invent at <https://go.aws/46iuzGv>.

Subscribe:

More AWS videos: <http://bit.ly/203zS75>

More AWS events videos: <http://bit.ly/316g9t4>

ABOUT AWS

Amazon Web Services (AWS) hosts events, both online and in-person, bringing the cloud computing community together to connect, collaborate, and learn from AWS experts.

AWS is the world's most comprehensive and broadly adopted cloud platform, offering over 200 fully featured services from data centers globally. Millions of customers—including the fastest-growing startups, largest enterprises, and leading government agencies—are using AWS to lower costs, become more agile, and innovate faster.

Agenda



1. Top Business Outcome for Edge AI
2. Bringing AI Everywhere
3. Compelling Edge AI Use cases
Edge Computer vision booth demo
4. Real World Successes
Featuring OnLogic, Arduino & ai.io
5. 4th Generation Intel Xeon for AI
6. The AI Pipeline of processing
7. Start building your AI solutions today

Edge Compute and AI: **Top Trends We're Seeing**

- By year-end 2026, 70% of large enterprises will have a documented strategy for edge computing, compared to fewer than 10% in 2023.
- CEO use of the term “digital” in their top business priorities has roughly doubled from 2018 to 2022.
- Digital data production at the edge is growing exponentially, creating the opportunity for deeper analysis, automation, AI /ML.*

Edge Compute and AI: **Top Trends We're Seeing**



“Enterprises need a strategy so they can overcome and learn from challenges, choose technologies, manage costs and enable successful digital transformation.”

Gartner

Building an Edge Computing Strategy,
Thomas Bittman, April 12, 2023

Bringing AI to the Edge

it
starts
with

intel



AI Continuum

Bringing AI Everywhere



Note: Intel Core Ultra integrates NPU low power inference engine from Meteor Lake onwards.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

Bringing AI Everywhere at the Edge

Move from development to deployment **fast**



Scalable hybrid AI approach

Leverage automatic processing of an AI workload using available/targeted system resources and accelerators on the edge, or in the cloud. OpenVINO and oneAPI allow developers to seamlessly transition between edge-to-cloud with a single code base.

Open tools to speed deployment

The OpenVINO toolkit optimizes deep learning inference deployment for hundreds of pretrained models across multiple Intel platforms. The Intel Geti Platform allows domain experts and data scientists to quickly build and train AI models.

Hundreds of deployment-ready solutions

Intel solutions include defect detection, worker and public safety, robotics, supply chain management, imaging diagnostics, and enhanced service delivery.



Bringing AI Everywhere to Unlock New Levels of Innovation

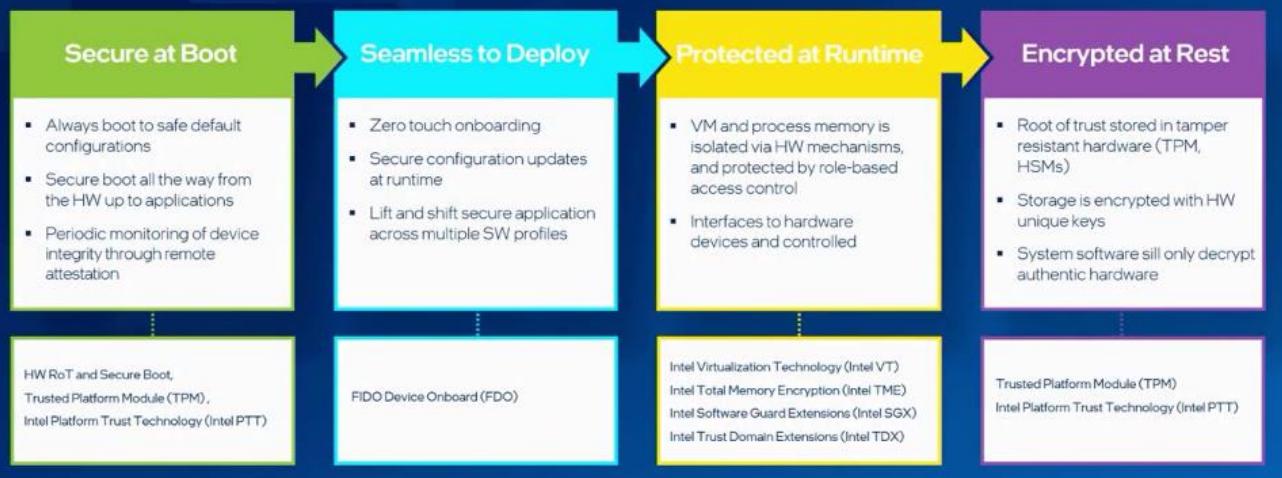
Education	Health	Finance	Retail	Government	Energy	Automotive	Manufacturing	Telco
Teacher Assistant	Drug Discovery	Algorithmic Trading	Product Promotion	Gov Services Chatbot	Energy Consumption Forecasting	Autonomous Car Development	Factory Automation	Personalized Customer Services
Student Study Buddy	Doctor Co-pilot	Customer Portfolio Assistant	Customer Interface and Sentiment Tool	Document Search Summarization	Operational Performance	Multi-language in car aid	Predictive Maintenance	Network Automation
Parent Chat Portal	Patient Family Chatbot	Risk / Credit Assessment	Image Shopping Aid	Live Language Translation	Energy Trading Assistant	Supply Chain Optimization	Precision Agriculture	Operational Performance

Secure the Edge

it
starts
with



Layering Security Foundation for Solution Developers



Edge Computer Vision

it starts with
intel.

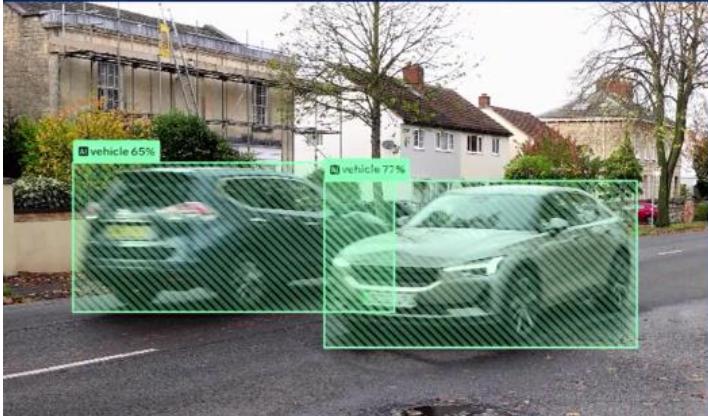


Where Will Computer Vision Applications Disrupt and Innovate?



Computer Vision scales into all sectors adding value to businesses and consumers alike

Computer vision at the edge



What is it?

Why does it matter

Where is it used?

ai.io

Data Solutions

aiScout

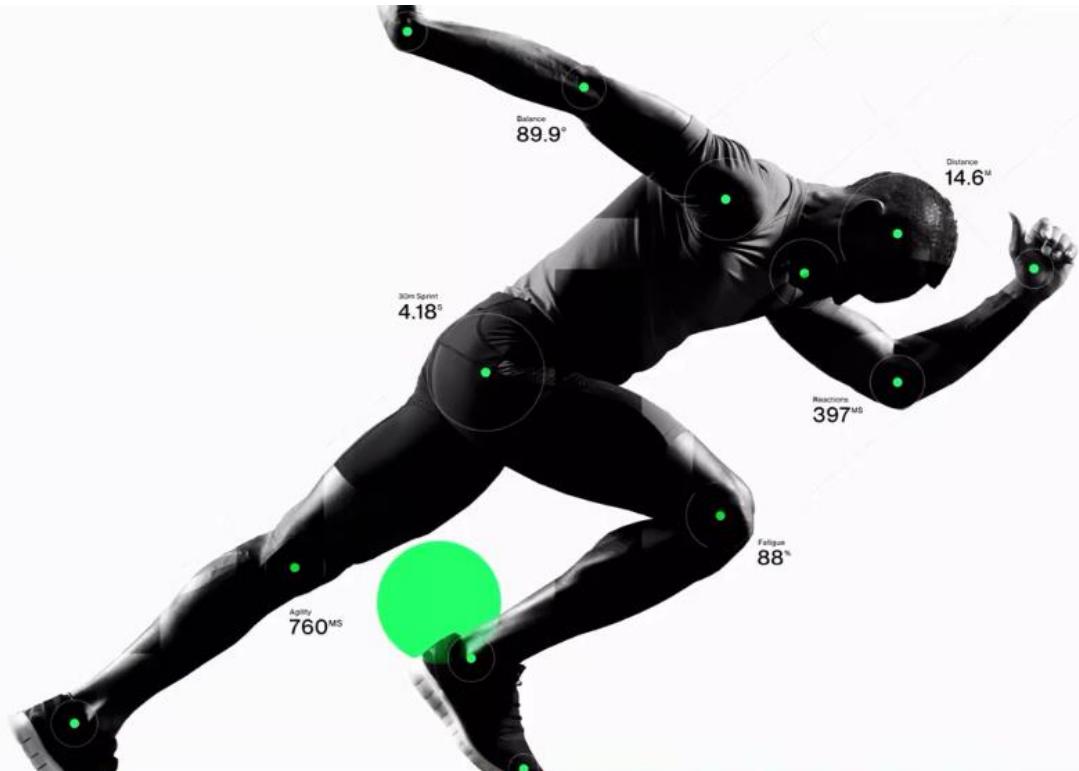
Talent Analysis & Development
App Platform

aiLabs

Elite Performance Labs

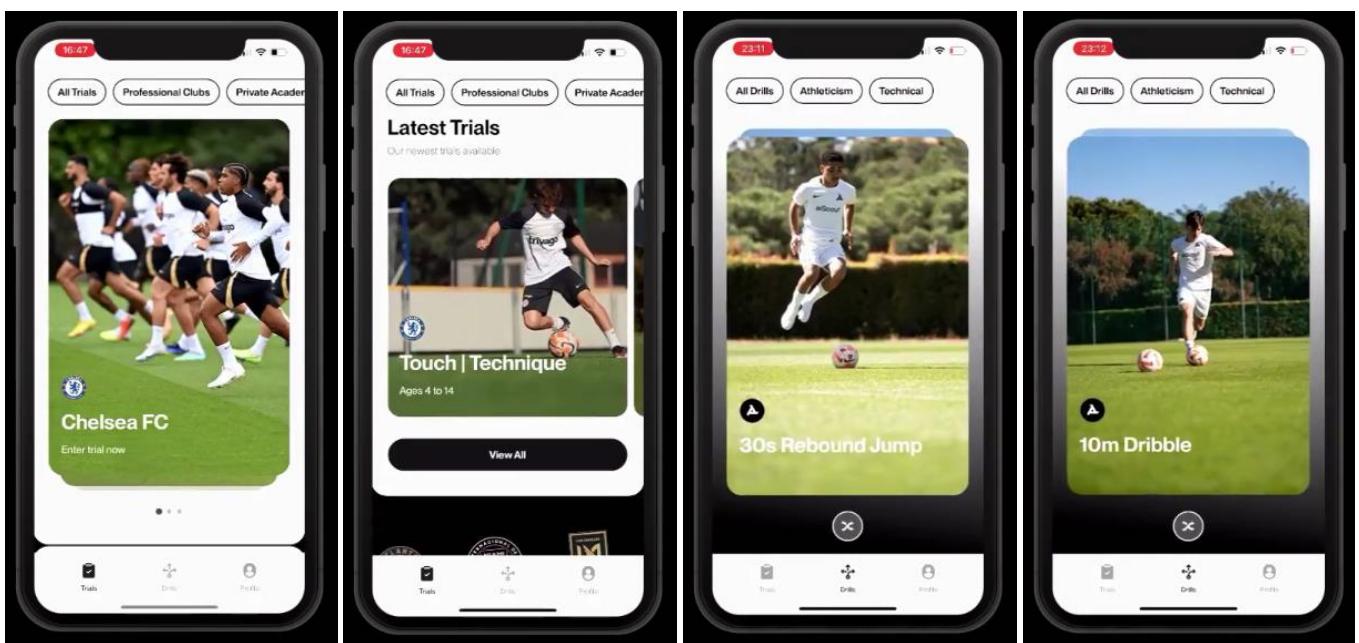
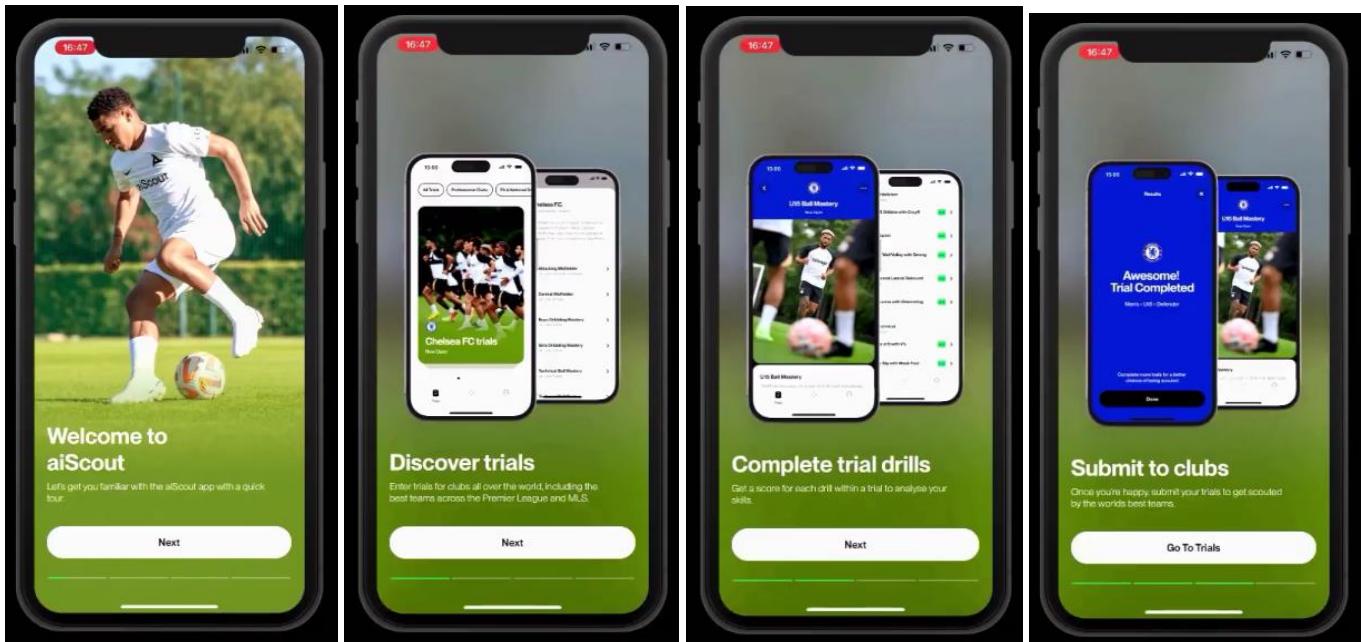
3DAT

Computer Vision &
Biomechanics Analysis

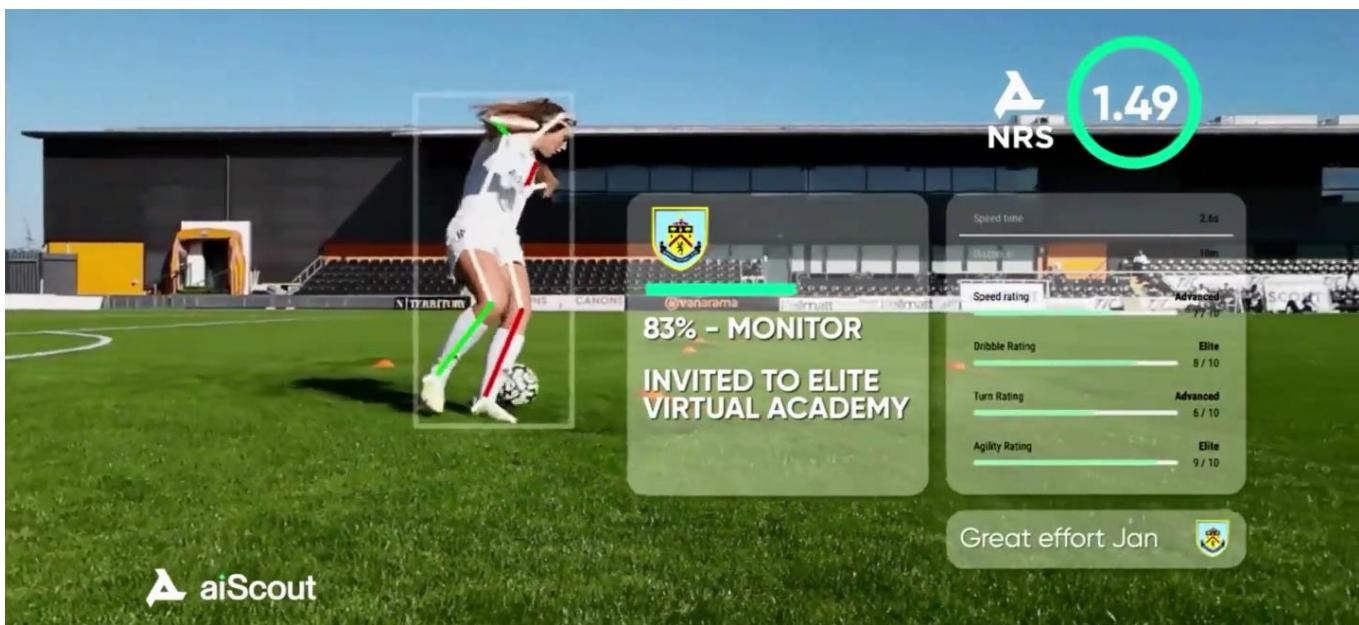


aiScout







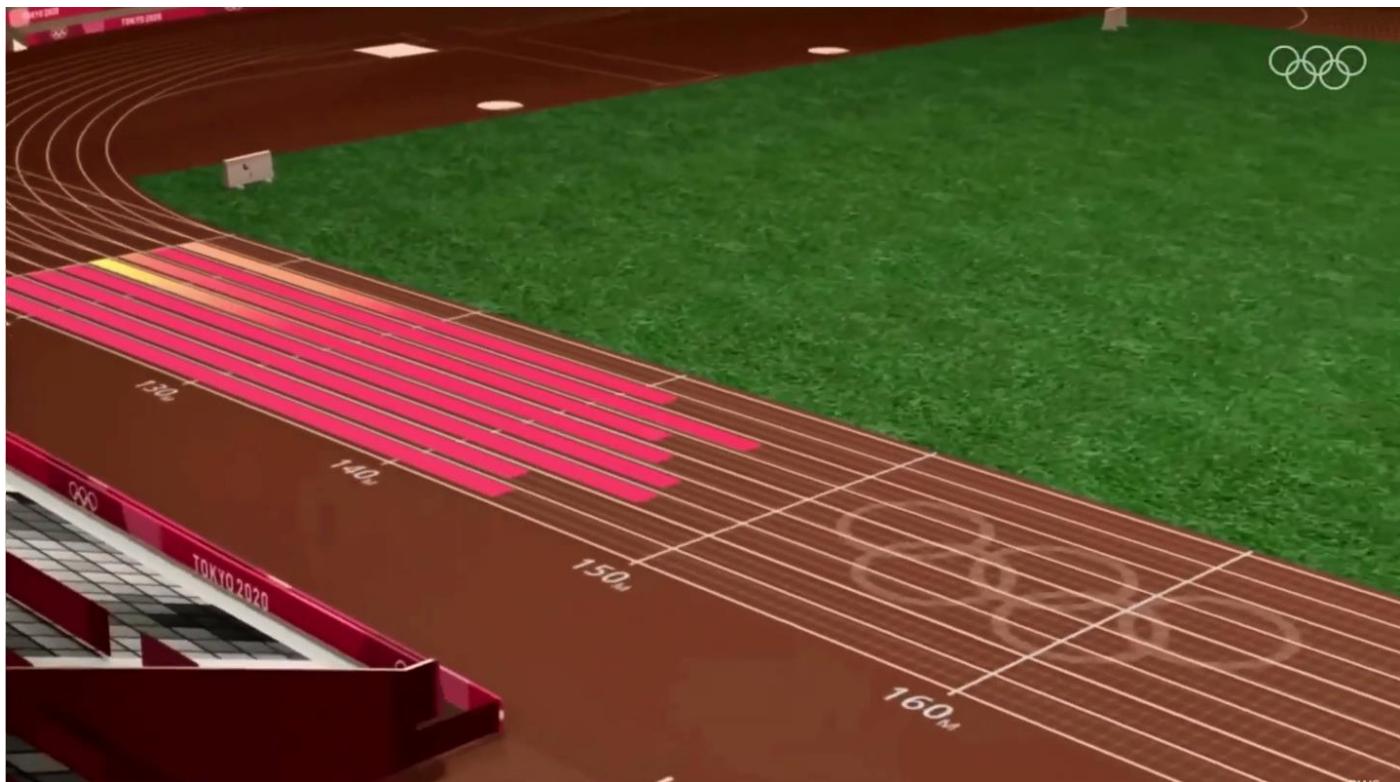




Δ

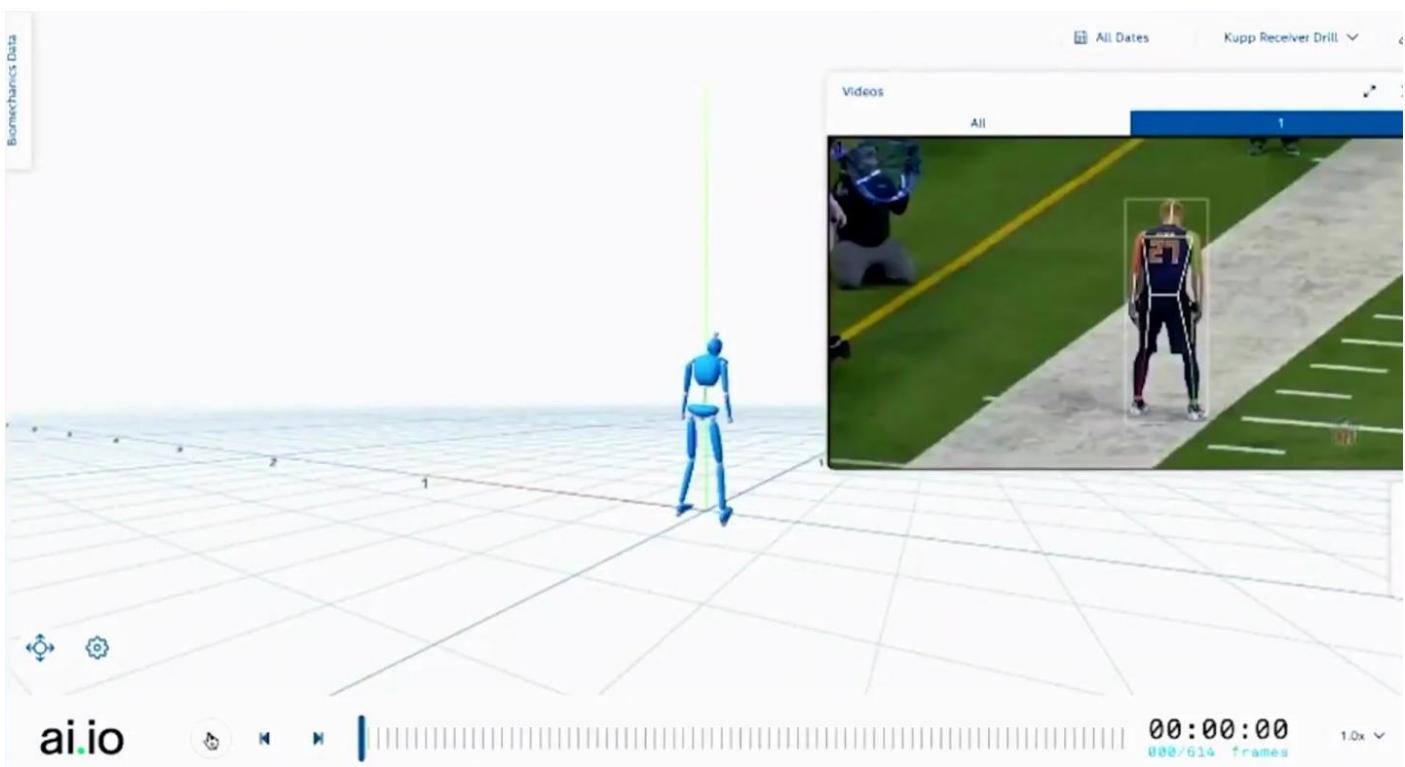
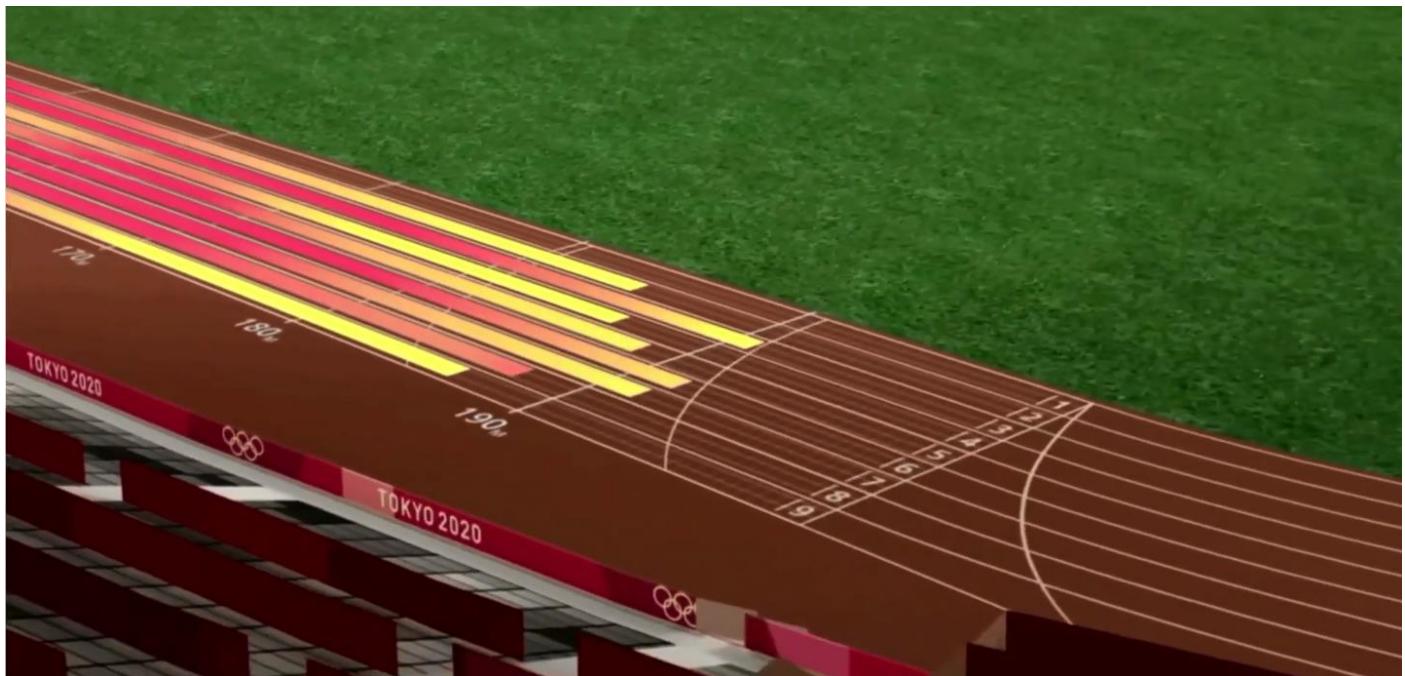
3DAT

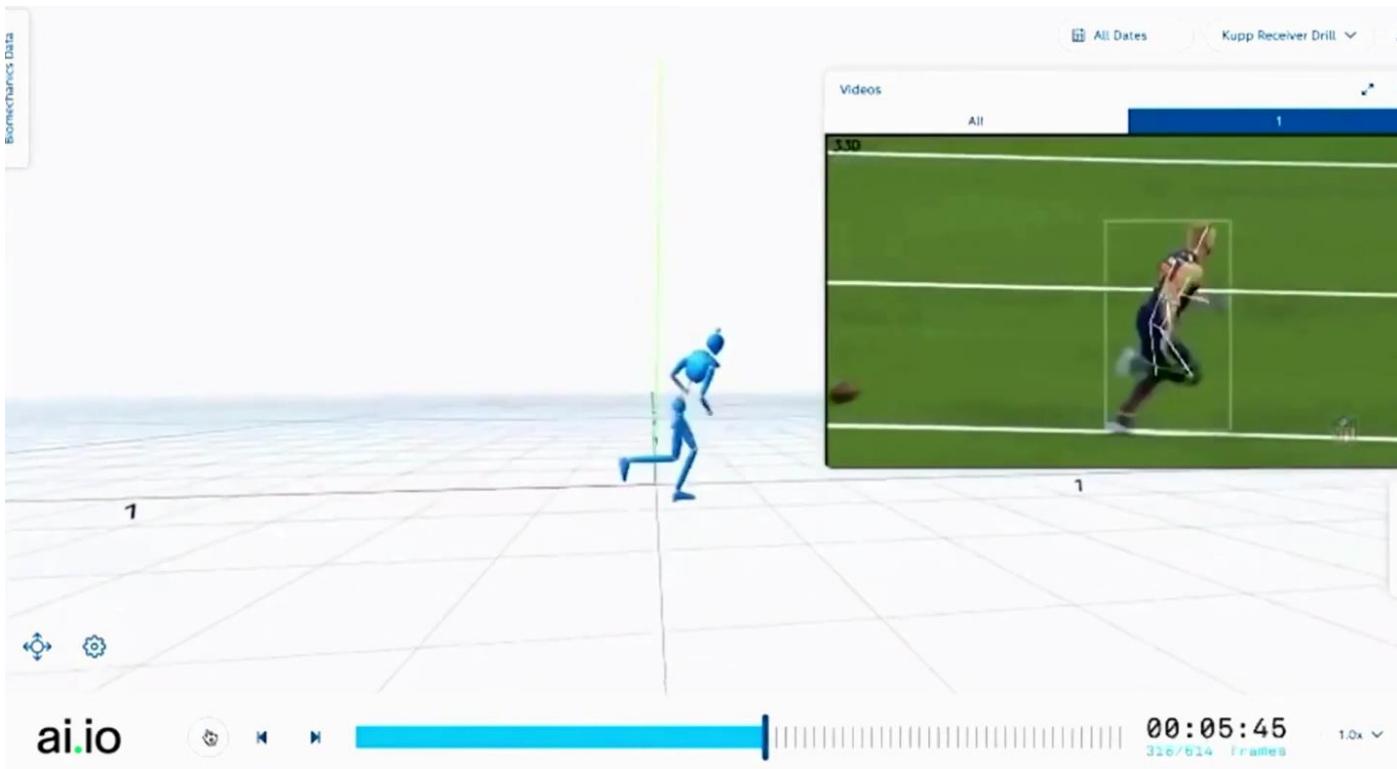
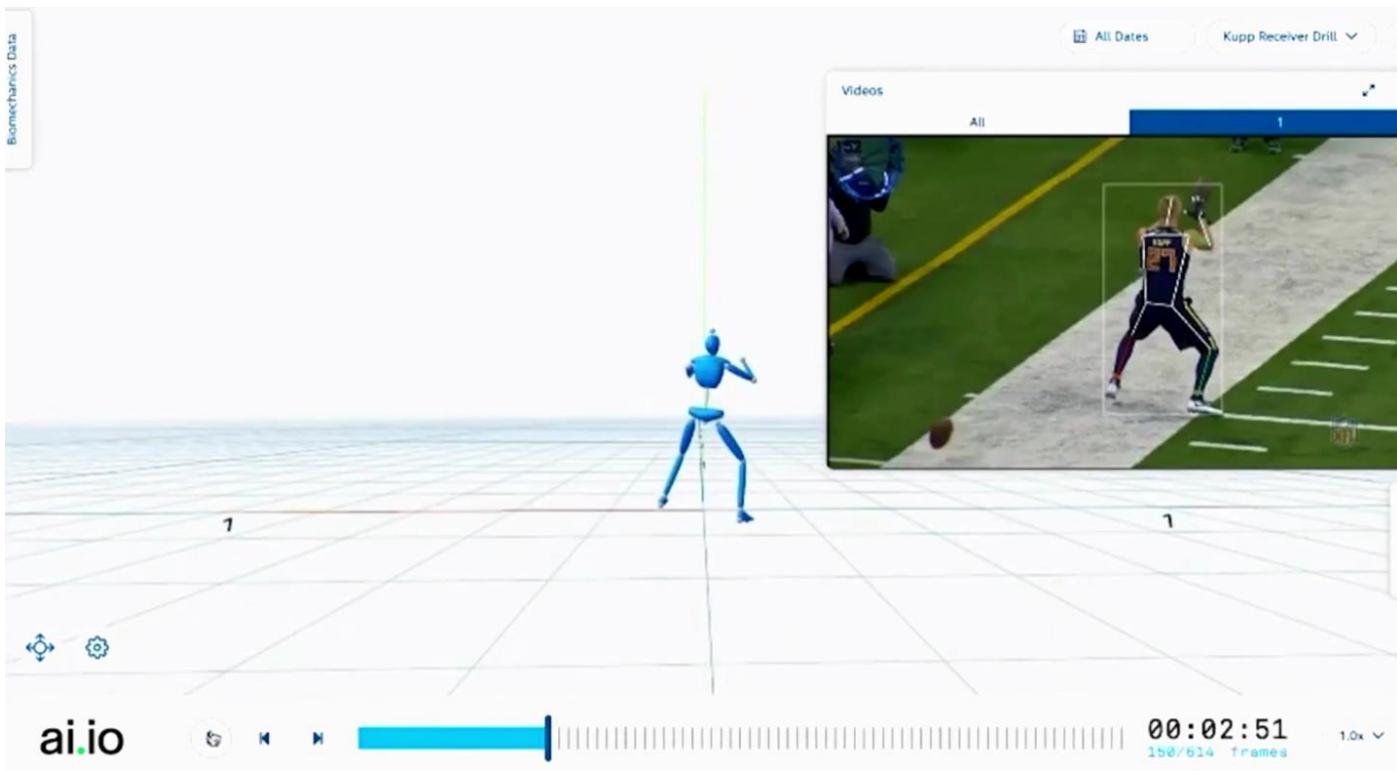
3D Athlete Tracking

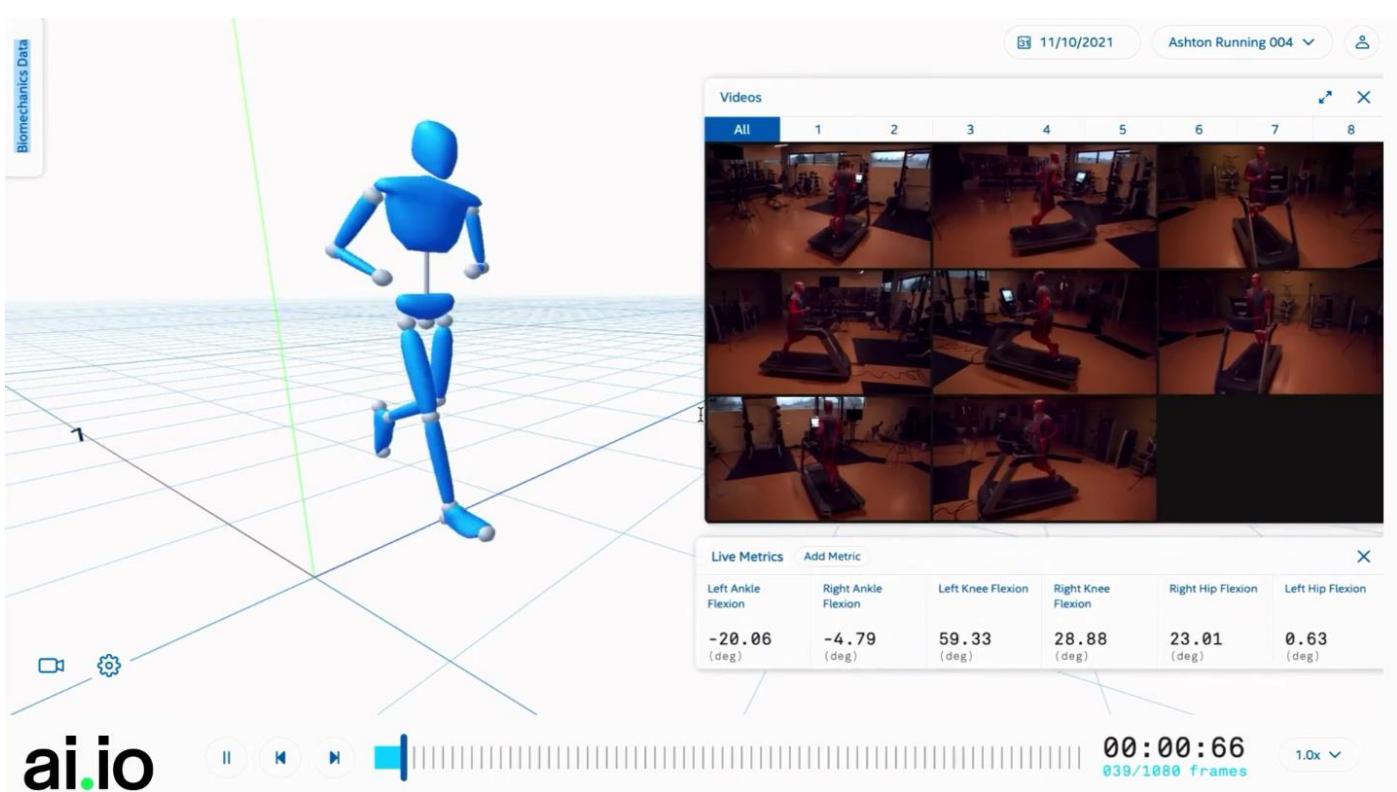
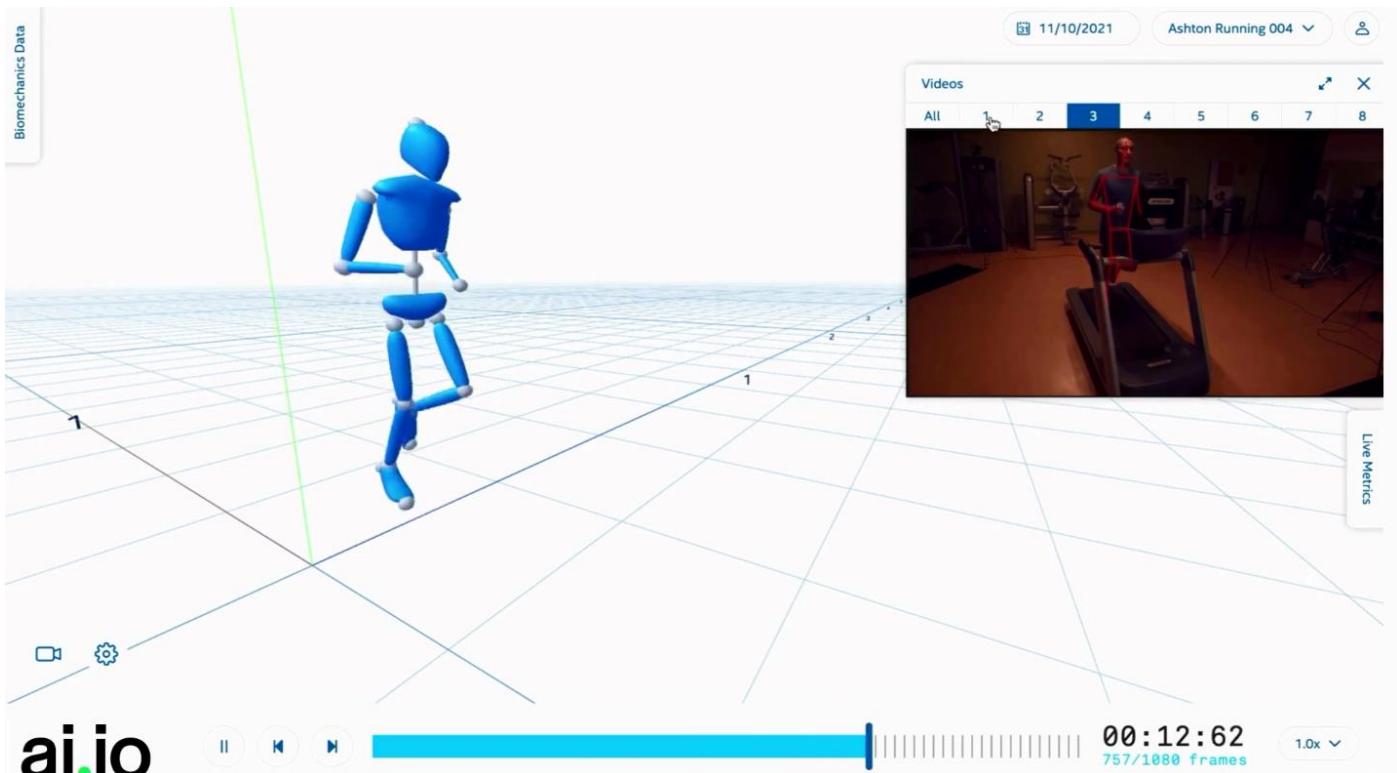


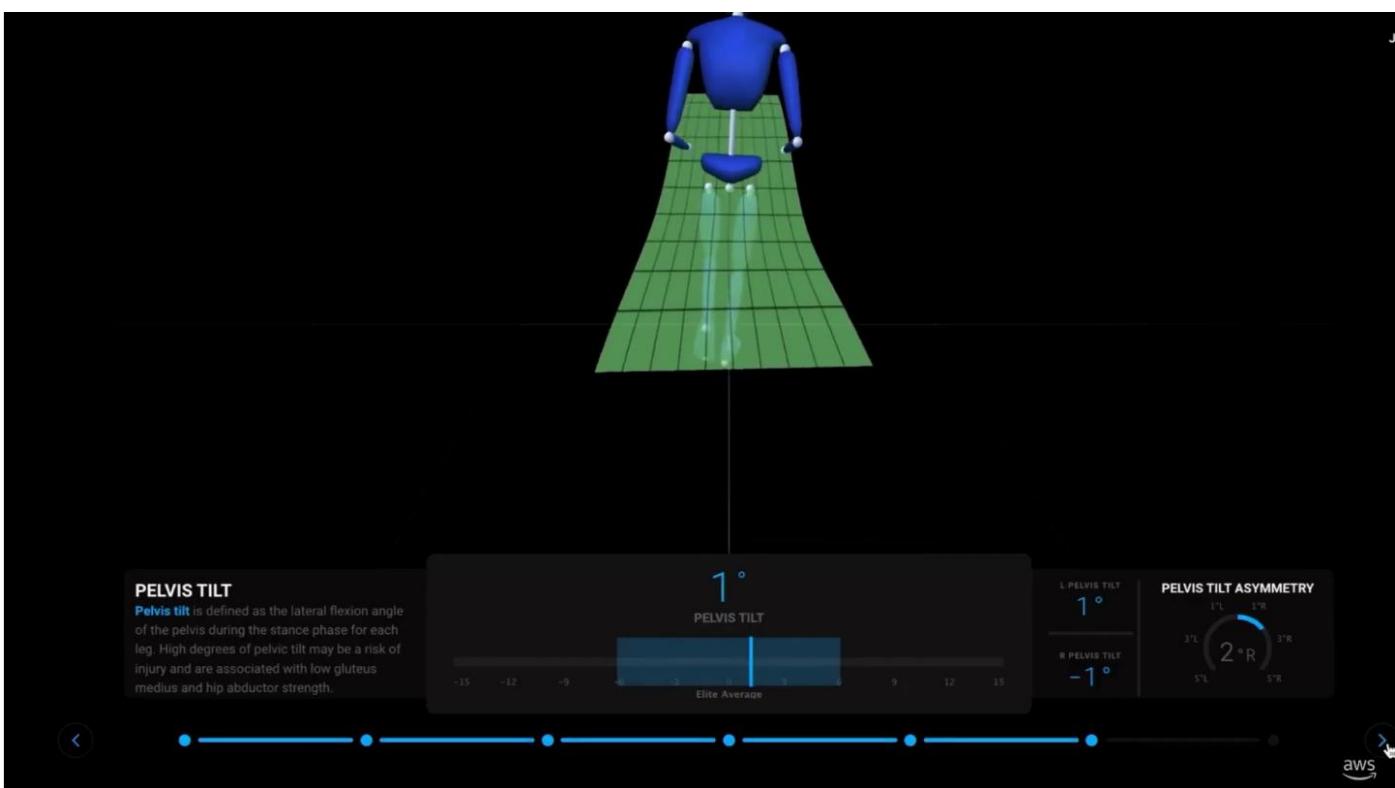
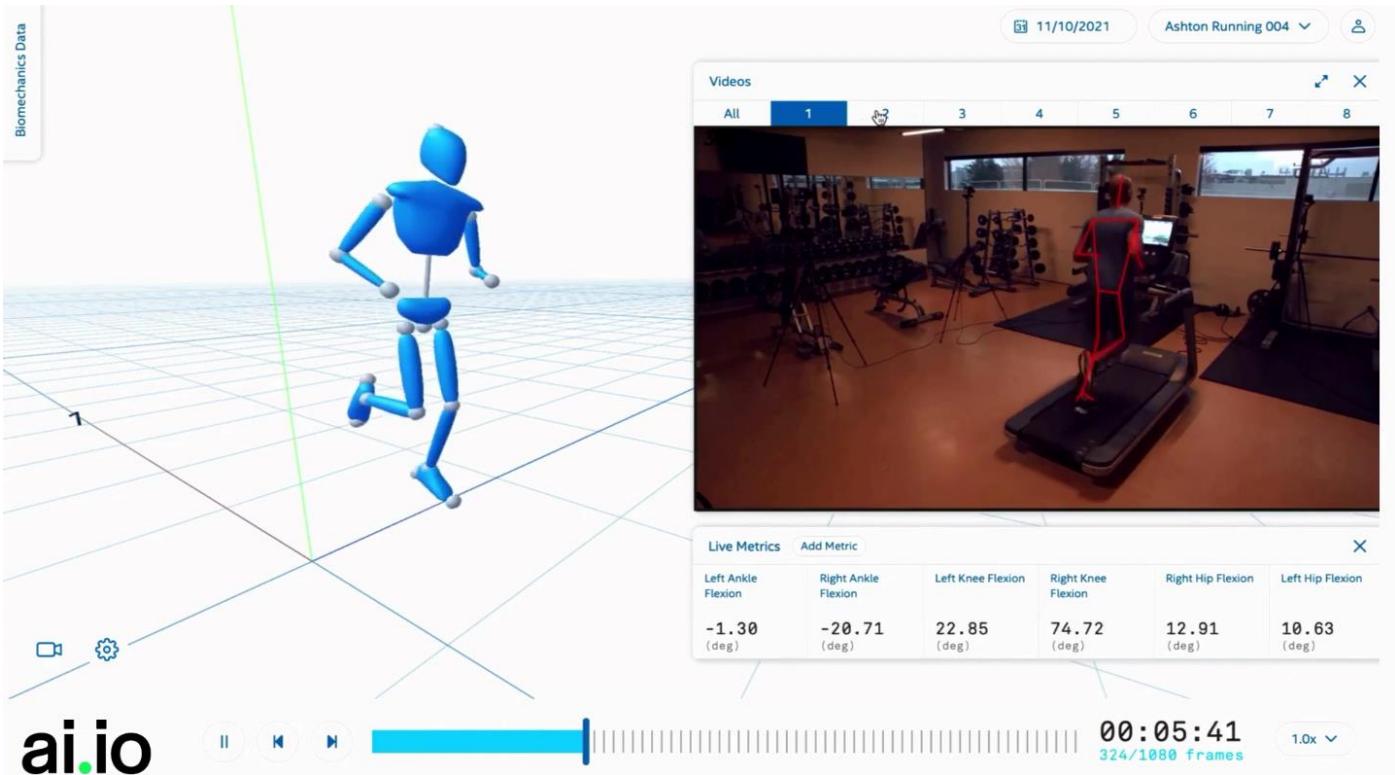
SPEED	
RICHARDS	37.3 KM/H
LYLES	38.3 KM/H
BROWN	36.6 KM/H
KNIGHTON	37.1 KM/H
DE GRASSE	38.2 KM/H
BEDNAREK	37.7 KM/H
FAHNBULEH	38.5 KM/H
DWYER	37.2 KM/H

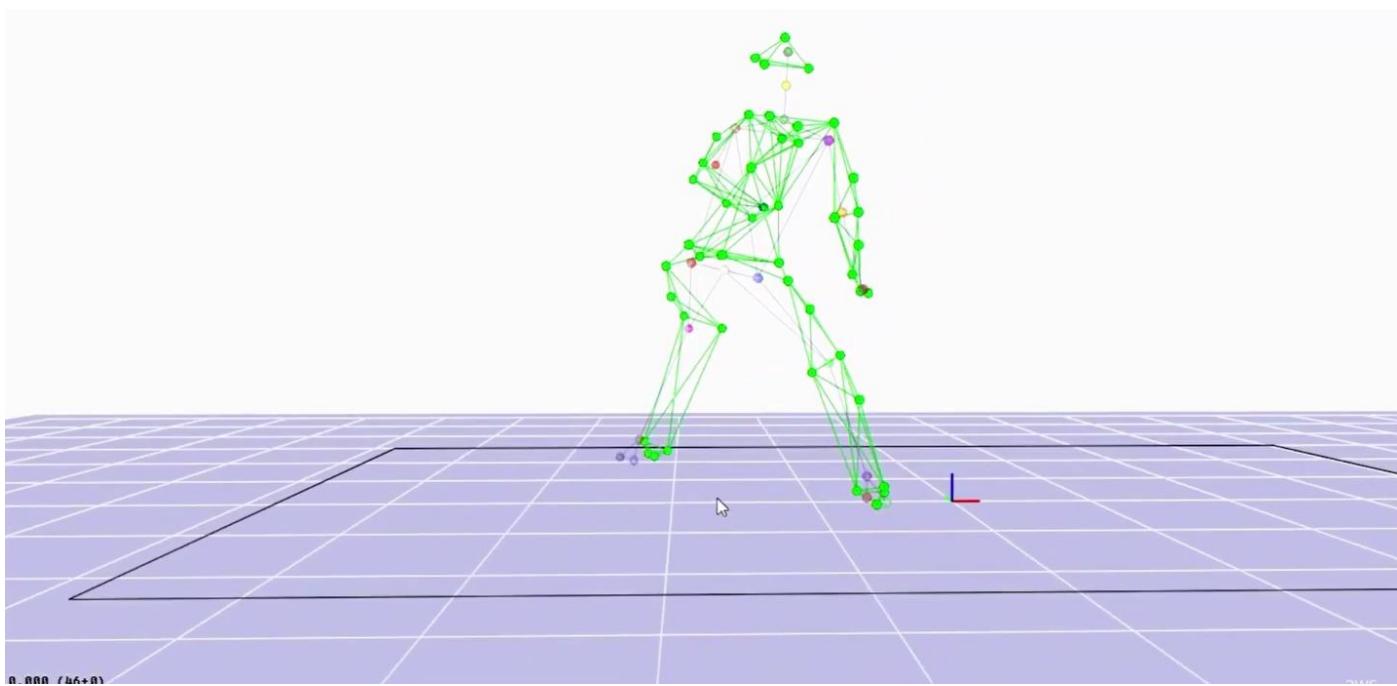
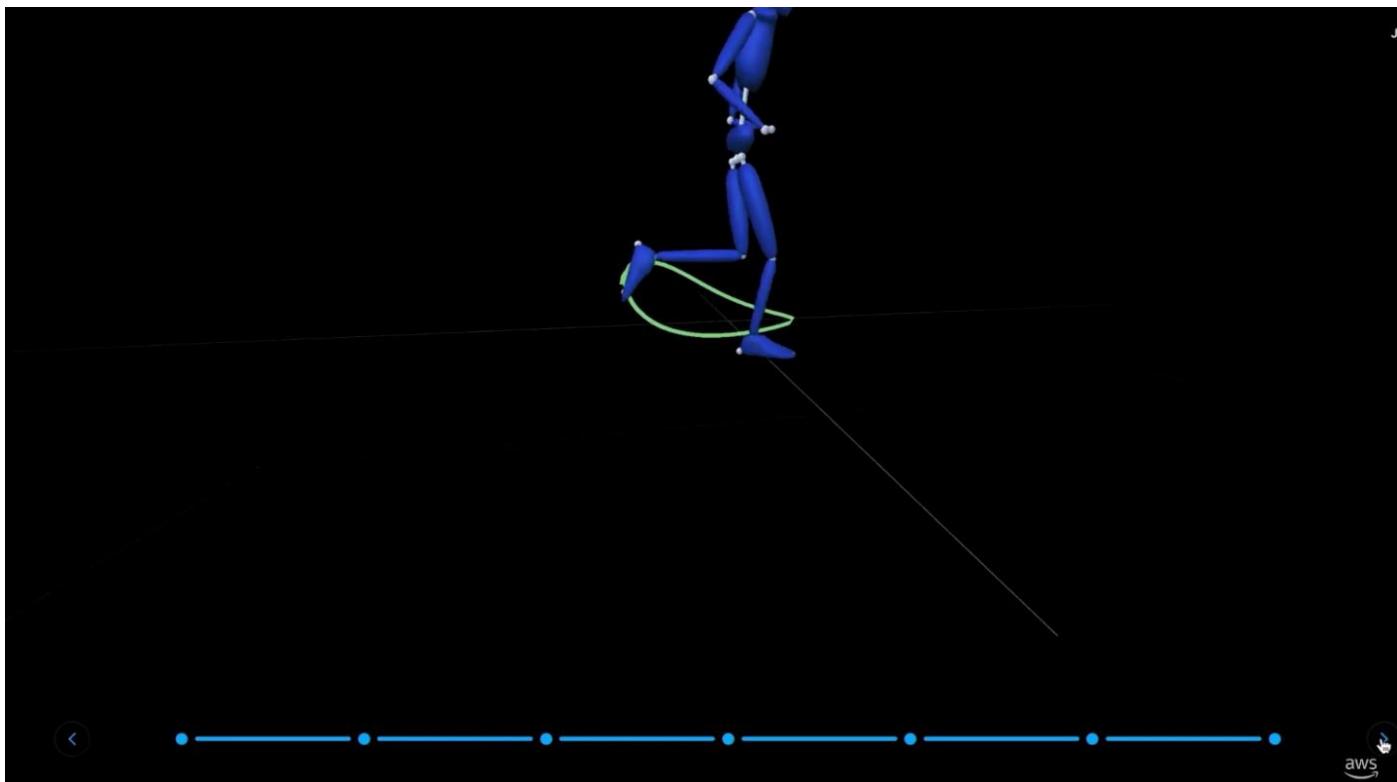
140_M 150_M 160_M

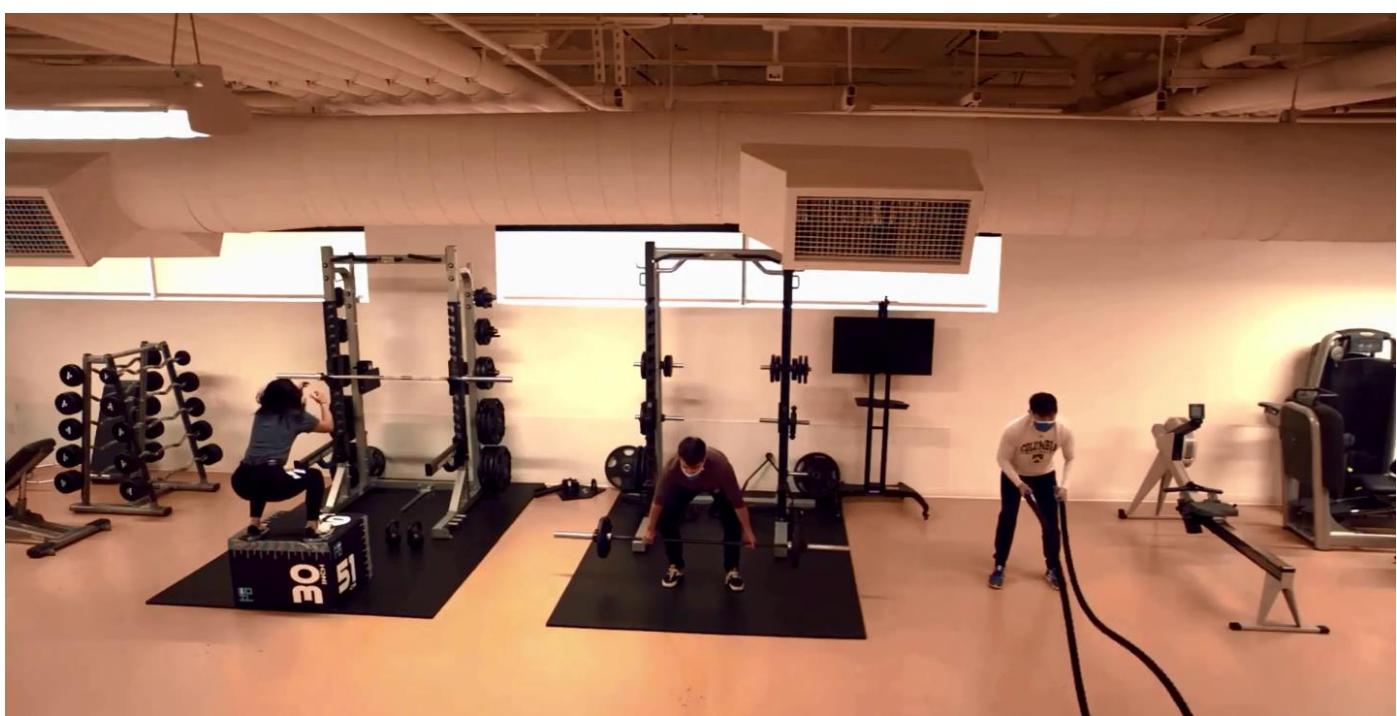
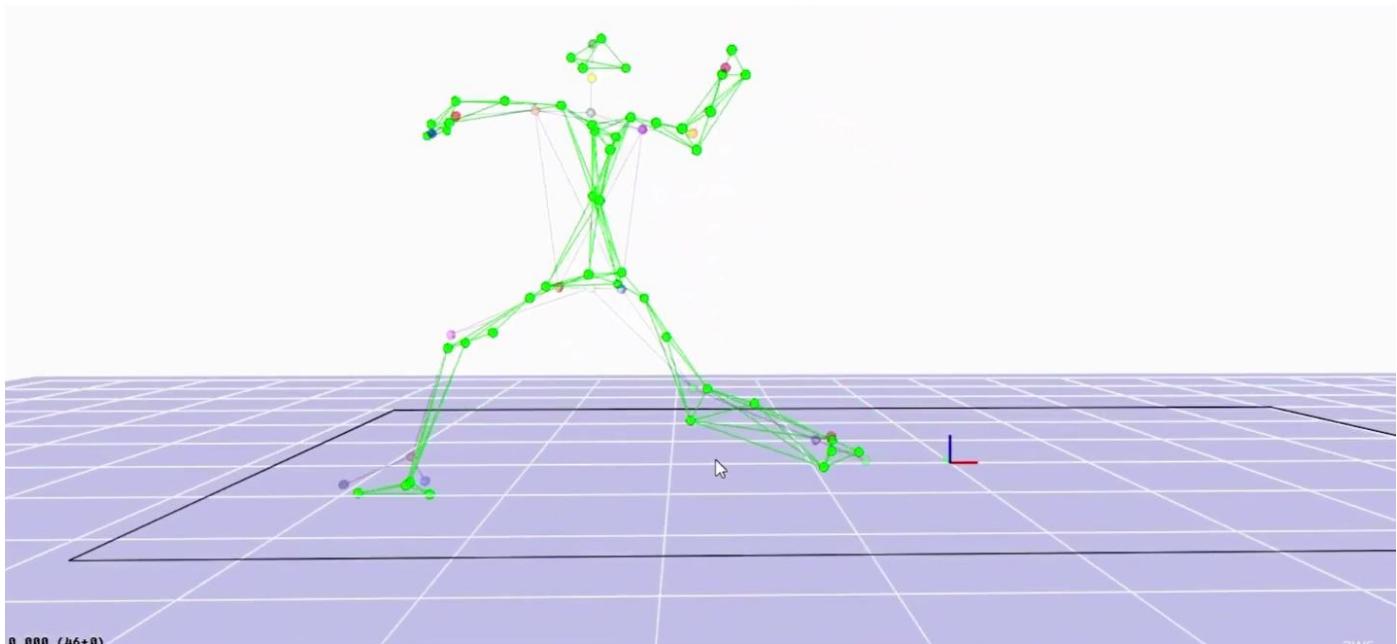






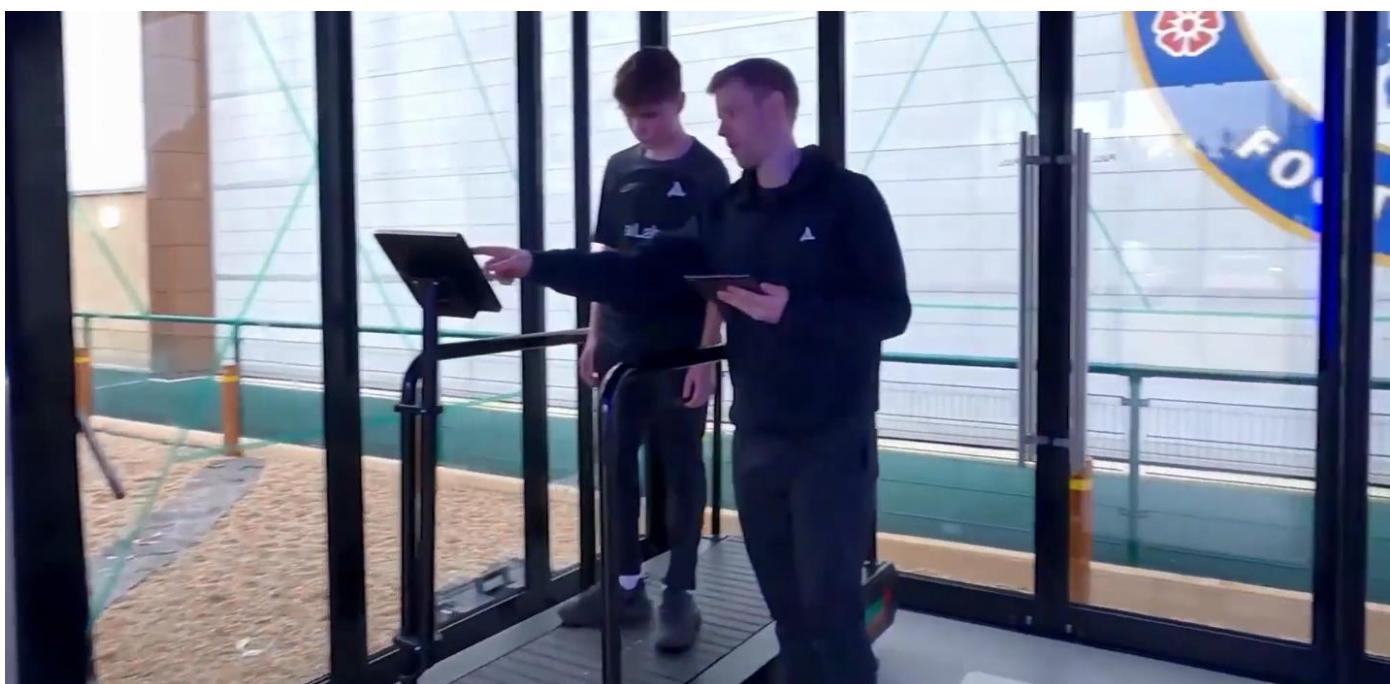


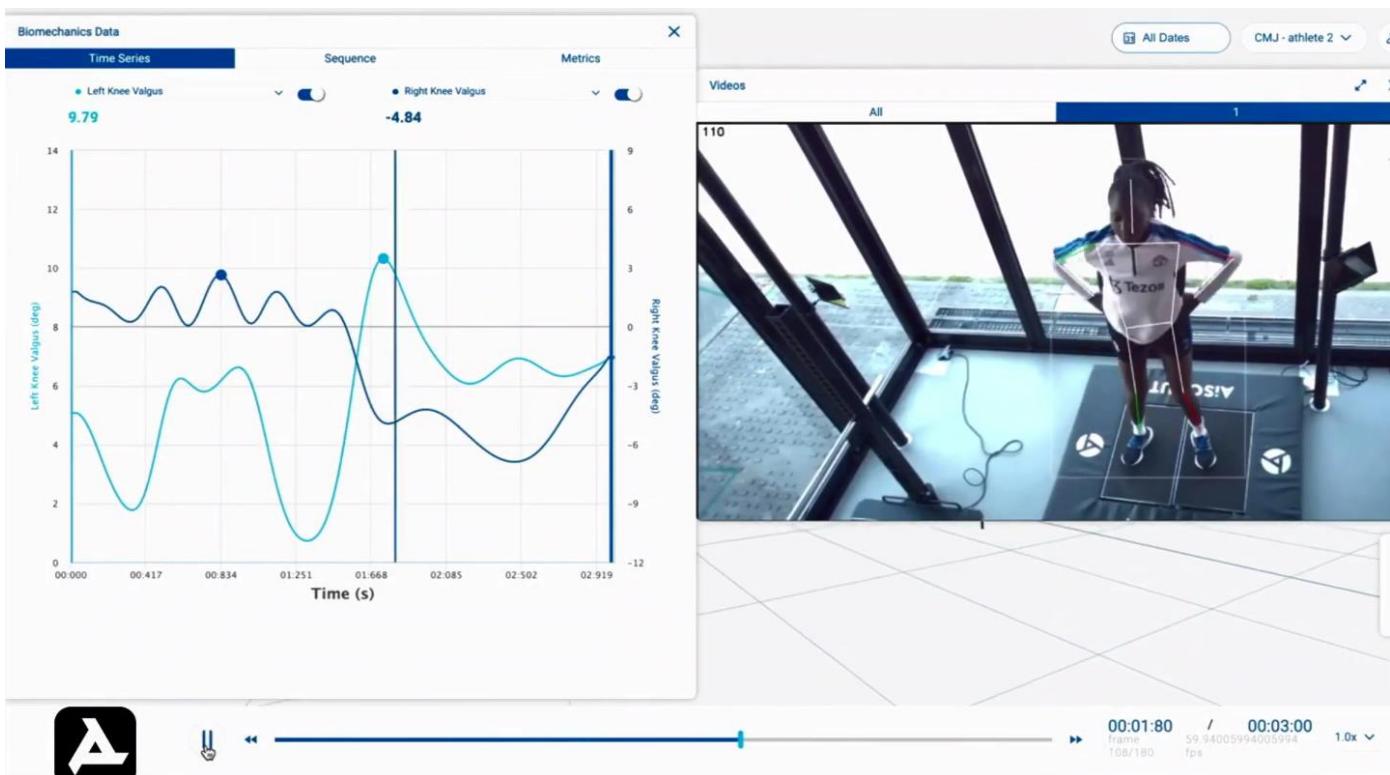
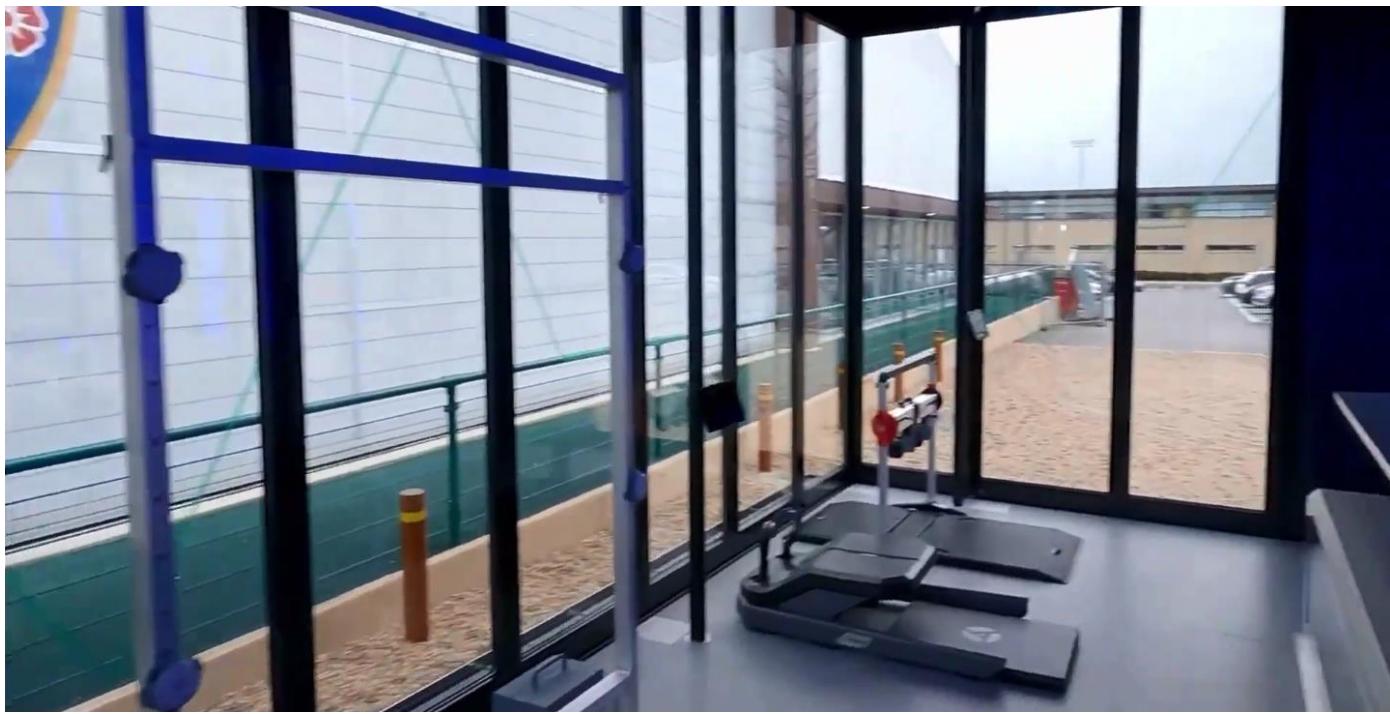


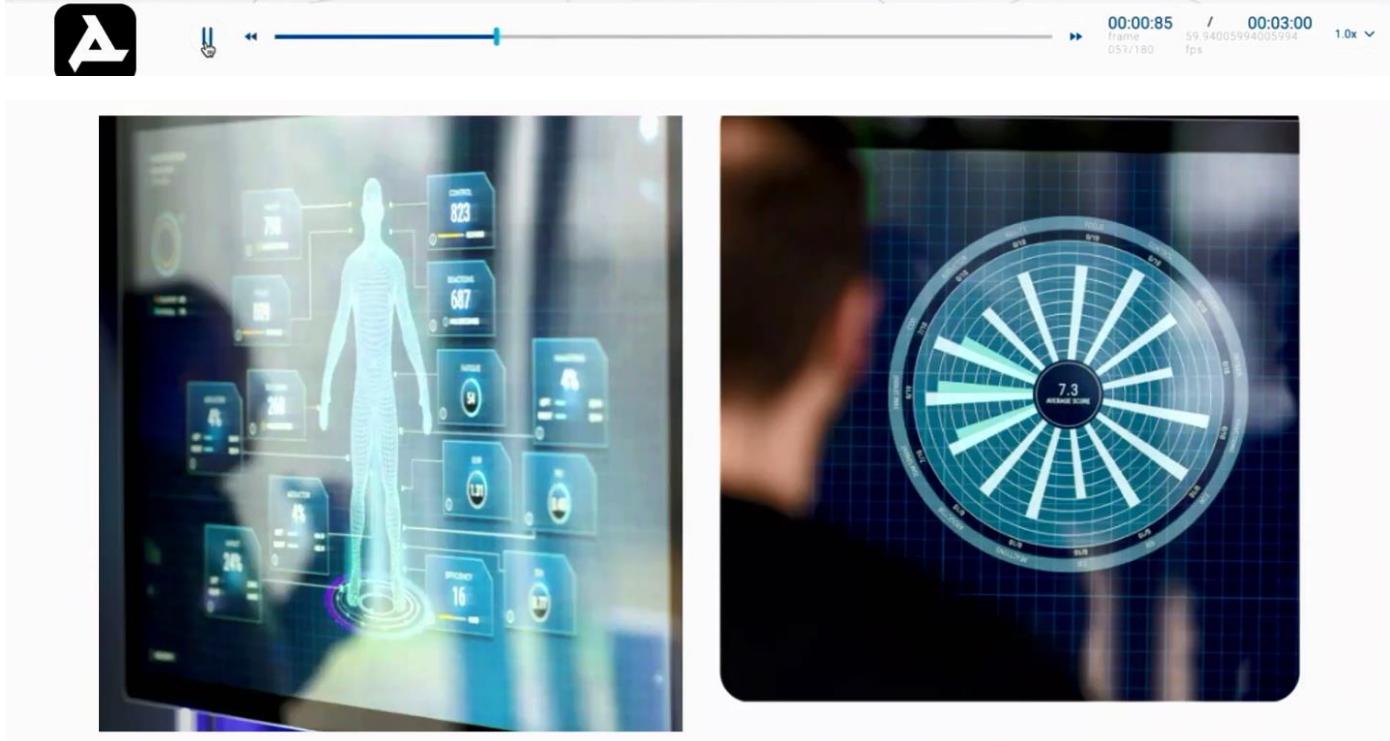
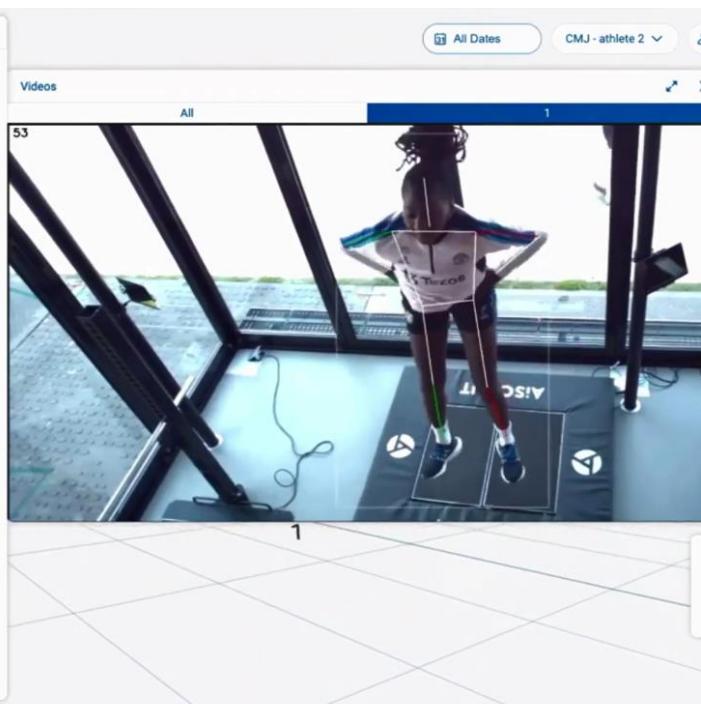
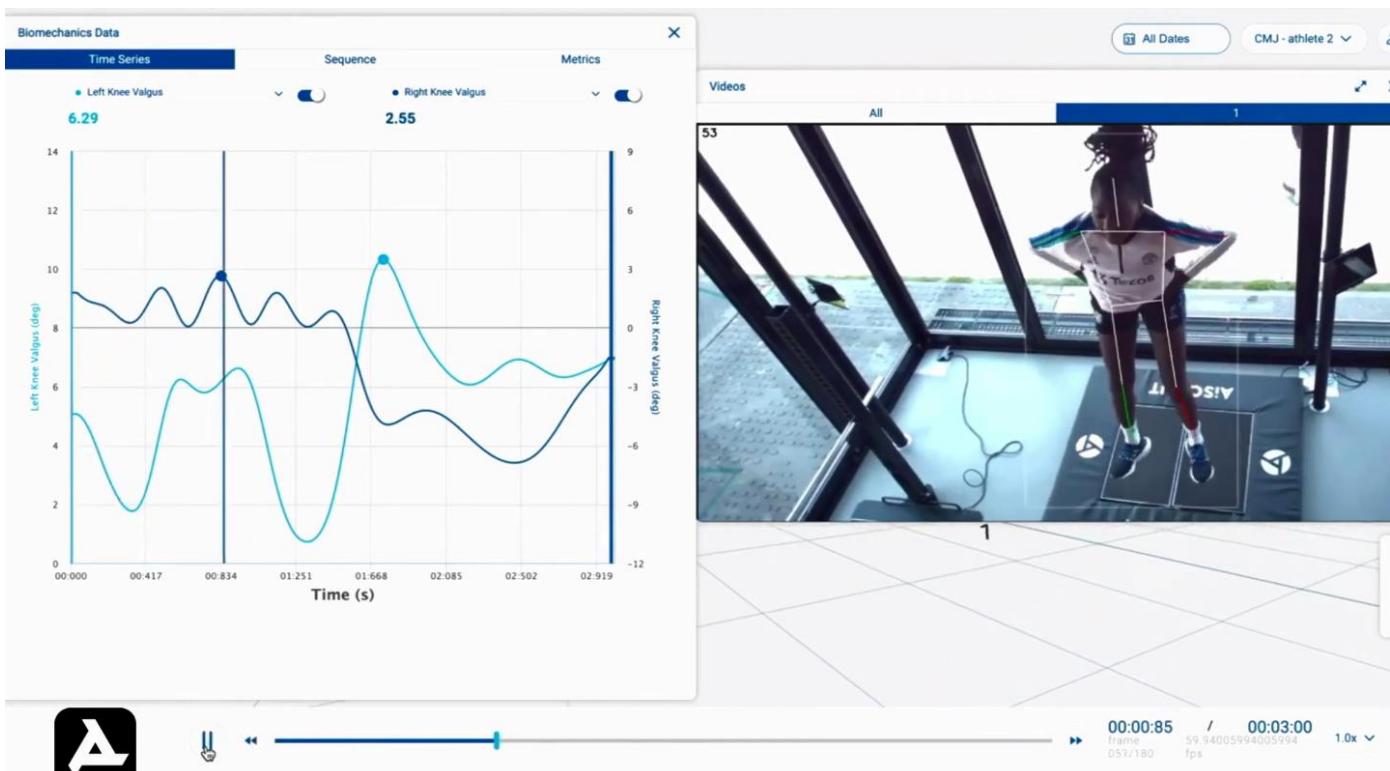


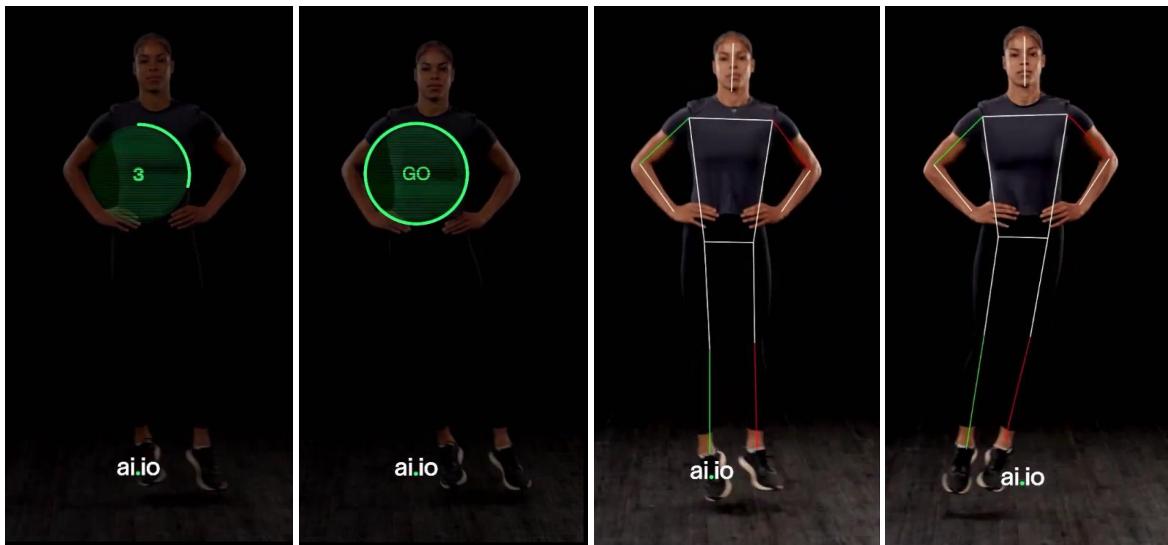












GPU Accelerated Computer Vision

Learn more on the Intel Booth

Visit us at our booth: Intel Booth #750



Intel Data Center GPU

FLEX SERIES

140

4

Media
Engines

75W

Power
Envelope

16

Ray Tracing
Units

16

Xe
Architecture

Half
Height

PCIe

Optimized
for lower TCO

5X

Media transcode throughput at half the power!
Intel Flex 140 GPU compared to NVIDIA A10

Media Delivery

30

1080p 60fps Stream

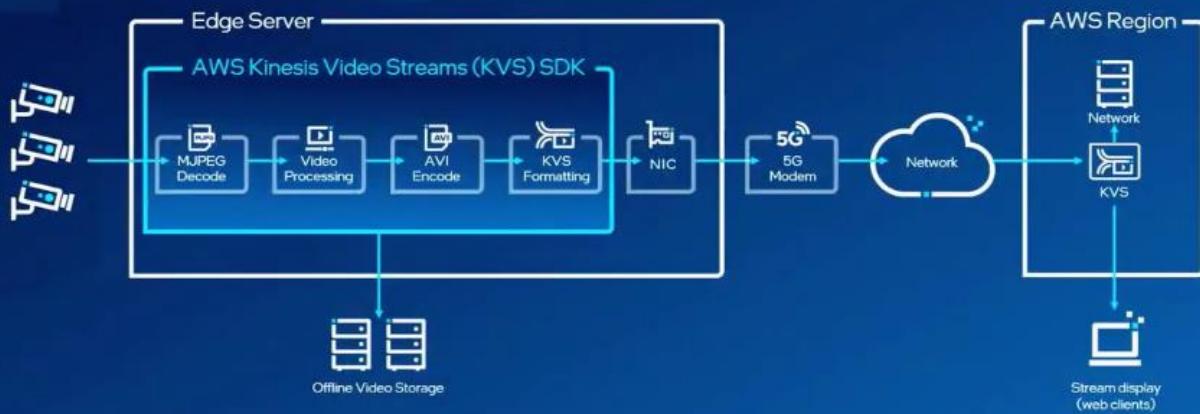
Cloud Gaming

40

720p 30fps Game Stream

See performance claims <https://edc.intel.com/content/www/us/en/products/performance/benchmarks/intel-data-center-gpu-flex-series/>

Solution Architecture





ONLOGIC

Let's Make It Possible



Industrial Computers

- Small Form Factor
- Fanless or Active Cooling
- 0 to 50°C Operating Temp.

Rugged Computers

- Resistant to Shock & Vibration
- Wide Power Input Range
- -40 to 70°C Operating Temp.

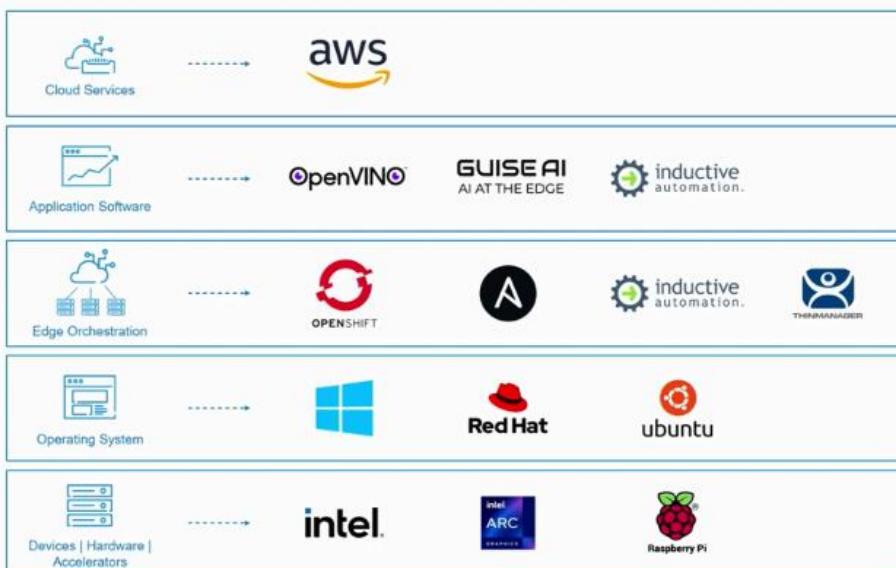
Panel PCs

- 8.4" to 24" Screen Sizes
- Resistive or Capacitive Screens
- Available with IP65 Front Bezels

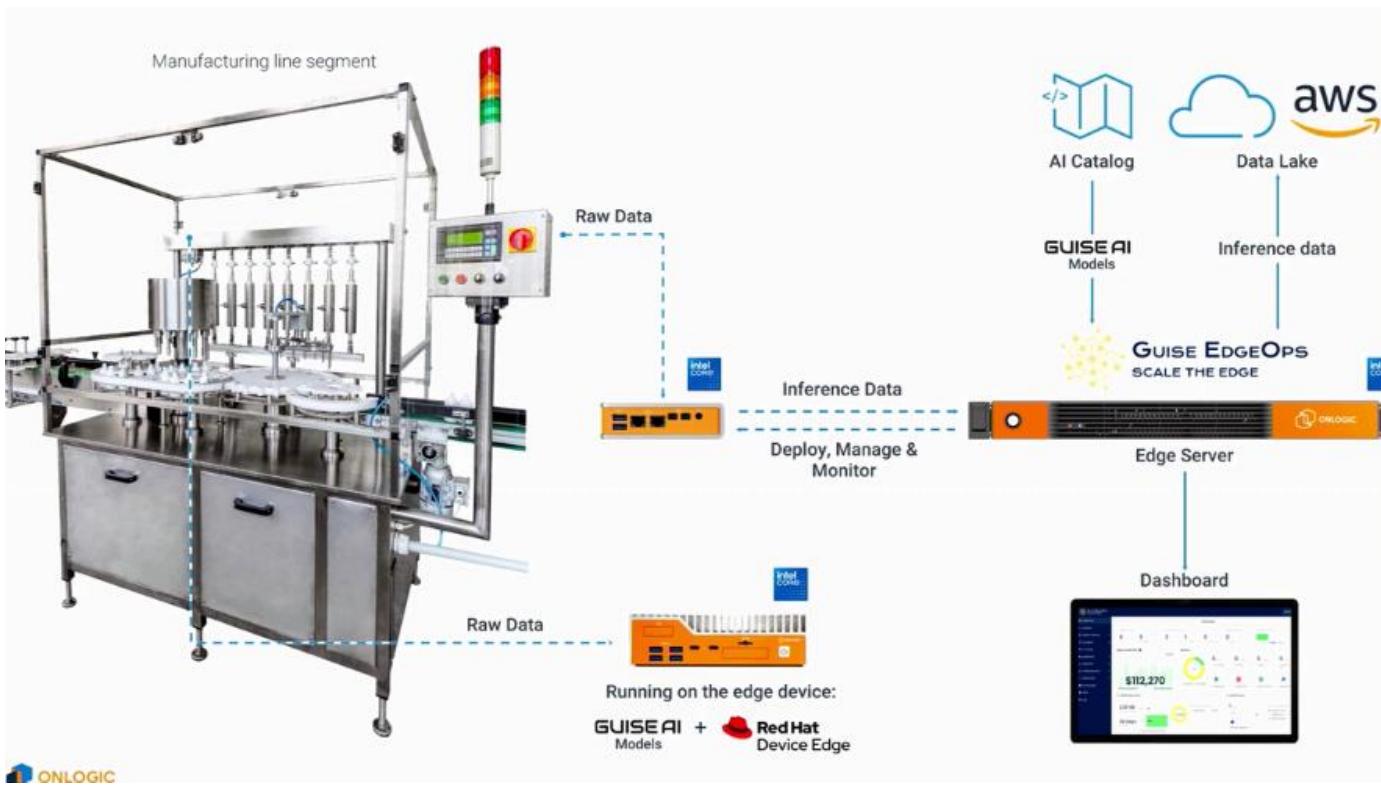
Edge Servers

- 1U to 4U Sizes
- Advanced Intel Processing
- Highly Configurable

Collaborations that make AI at the Edge easy



*Only a subset of our partners listed



OnLogic Demos at re:Invent

Intel Booth - KVS Demo:

Expo Hall - Venetian Conference Center, Booth 750

- Running Amazon Kinesis to process HD video streams and send data to the cloud.

Builders' Fair – Find My Ship:

Expo Hall - Venetian Conference Center

- Controlling AWS IoT Shadow updates along with view of Ultra-Wide Band Data visualization in 3D.

AWS IoT Kiosk – AWS IoT DeepRacer Vx Demo:

Expo Hall - Venetian Conference Center

- Showcasing Greengrass Runtime Software at the edge to get car data (via CAN Bus) and visualize it.

AWS Disaster Response DDIL activation:

Venetian Hotel Lobby - Disaster Response Jeep area

- Aggregating the GPS locations and metadata from 60 LoRaWAN GPS trackers deployed to re:Invent staff.

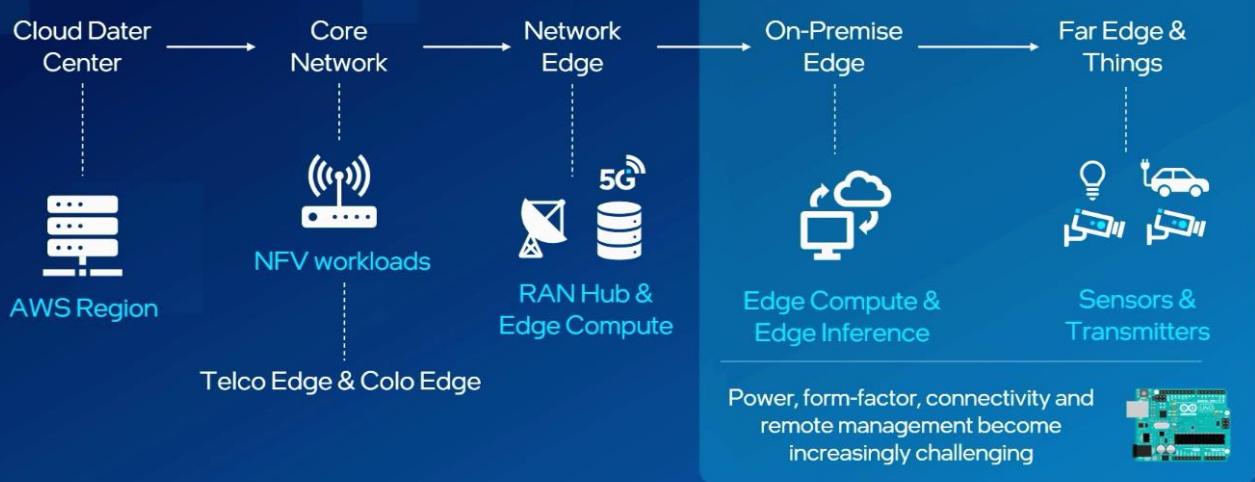
AWS GenAI Chess activation:

Rec Center – Mandalay Bay Convention Center

- Used for AWS IoT Shadow and Jobs update for each Chess move of the robots and to perform a Trust and Verify check between the physical chess board and the output of the GenAI model.



Expanding computer vision to the far edge



Arduino Cloud for Business Machine Learning Tools

Machine Learning Tools: Endless possibilities



Acquire valuable data securely from your devices and rapidly build custom dataset

Develop algorithms with ready-to-use digital signal processors and machine learning blocks

Validate and train Machine Learning models with real-time data

Build optimized embedded inference

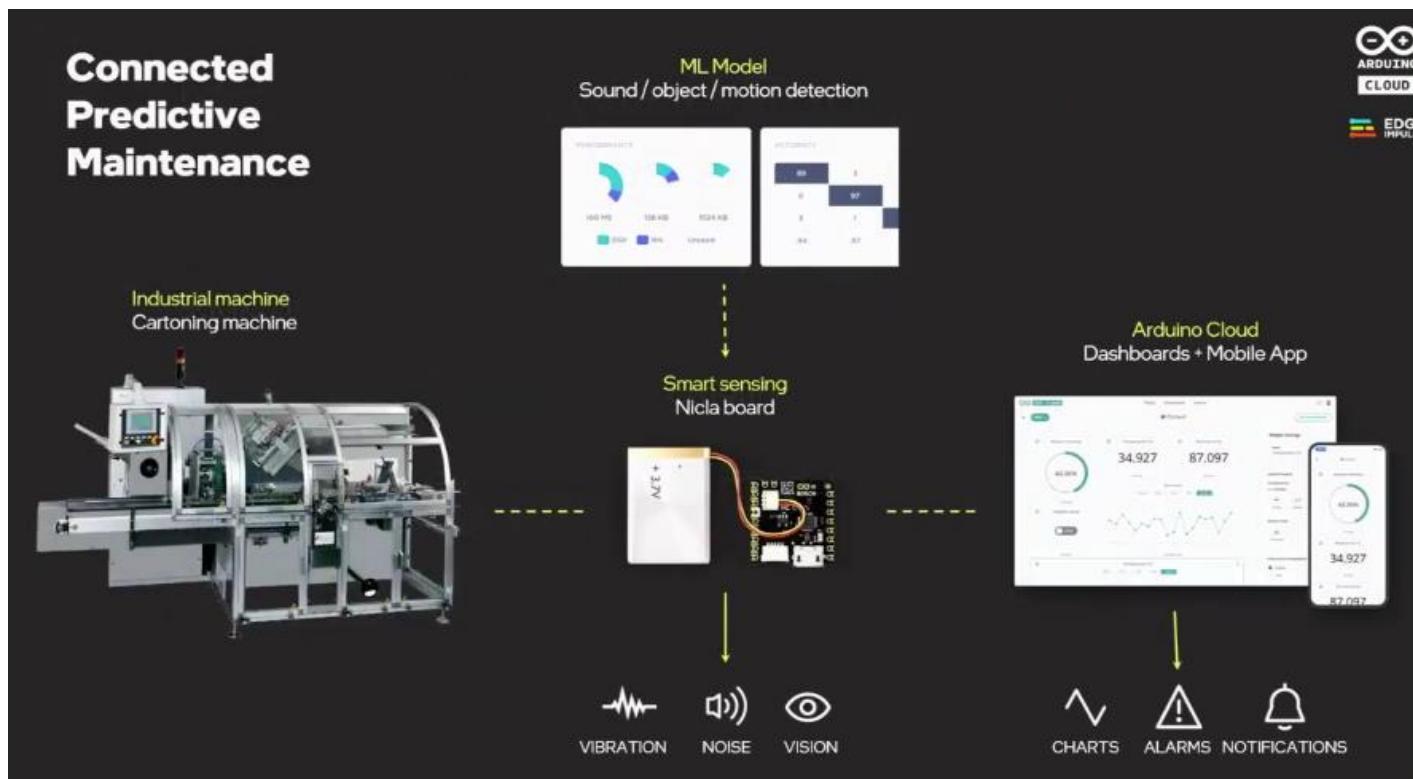
Machine Learning Tools at a Glance

Data collection - Build custom datasets by collecting sensor, audio or camera data straight from devices, files or cloud integrations.

ML model design - Start creating your model.

ML model test - Validate your model and picture how it will perform with real-world data.

ML model deploy - Deploy your model to any device.



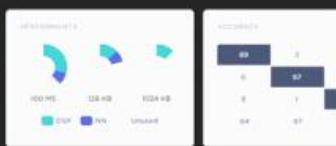
Component defect detection



Production line
Juice production



ML Model
Object detection



Defect detection
Nimble Vision



VISION

Arduino Cloud
Dashboards + Mobile App



CHARTS ALARMS NOTIFICATIONS

Component defect detection



Logistics
Automated inventory management



ML Model
Object detection



Amount / defect detection
Portenta H7 + Vision Shield



WiFi

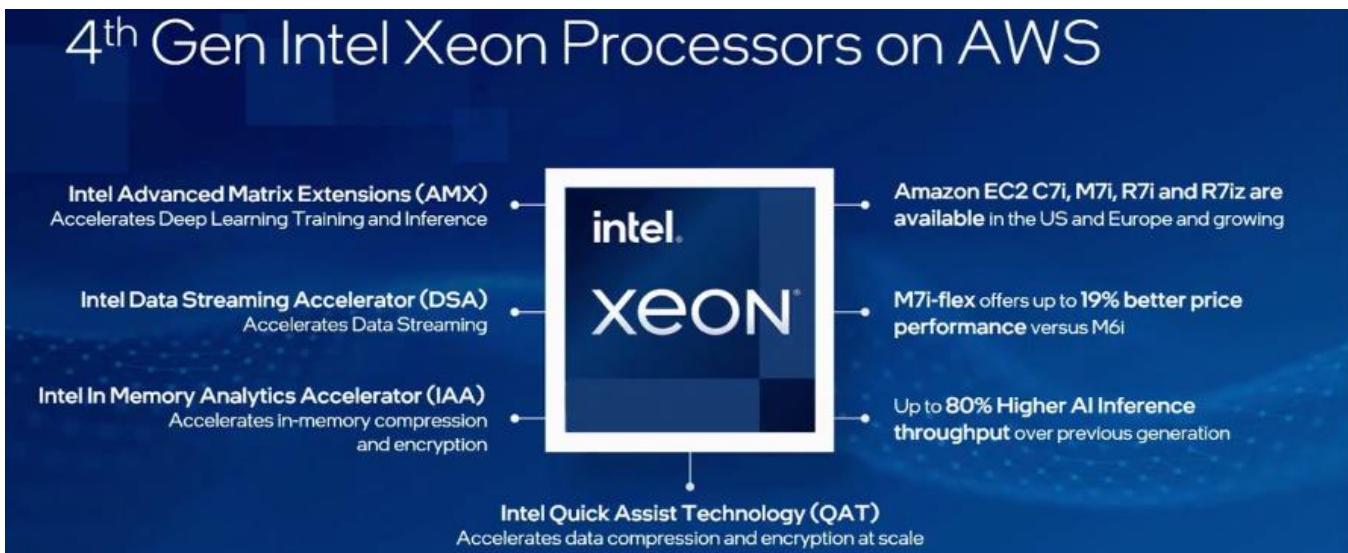
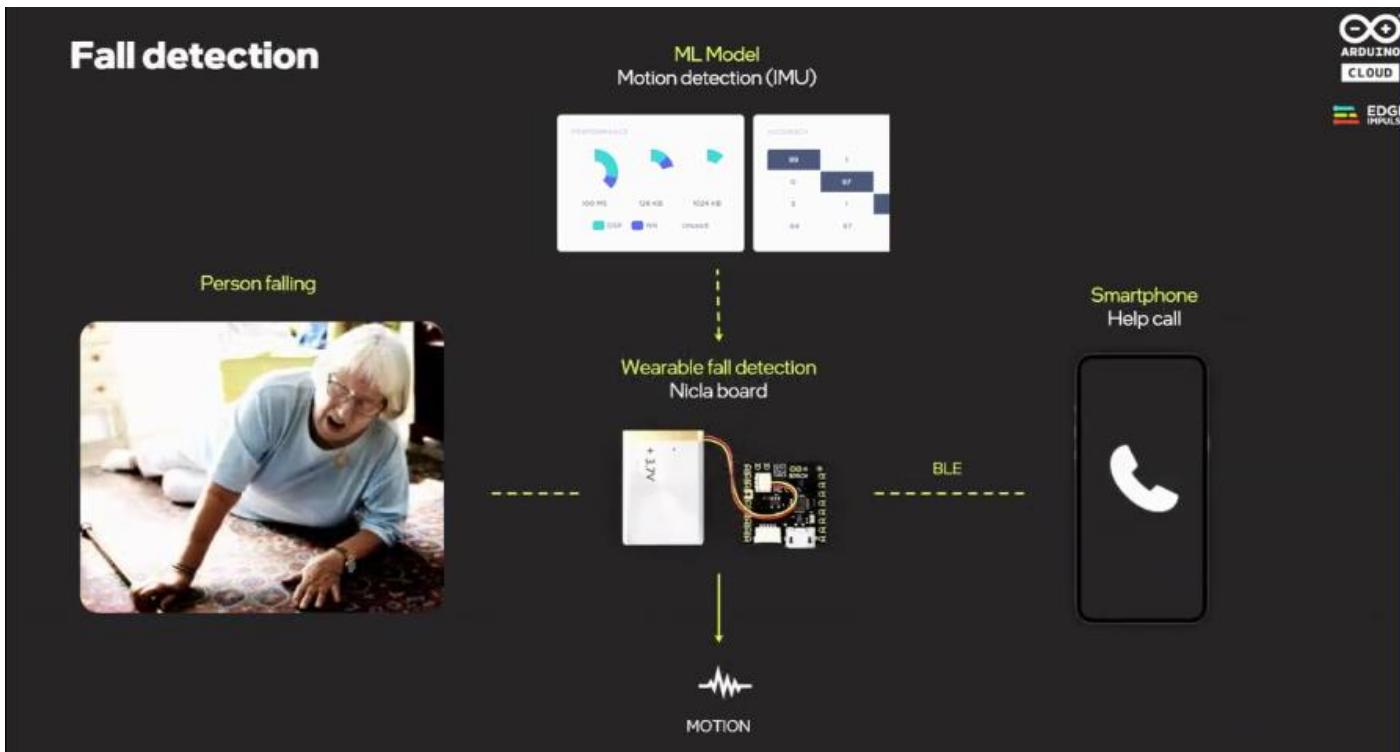


VISION

Arduino Cloud
Dashboards + Mobile App



CHARTS ALARMS NOTIFICATIONS



One-line Code Change = Huge Performance Gains

Engineer Data
~90x



```
import modin.pandas as pd
```

Create Machine Learning & Deep Learning Models

~38x



```
from sklearnex import patch_sklearn  
patch_sklearn()
```

Deploy

~3x



```
TF_ENABLE_ONEDNN_OPTS=1
```

Xeon AI and “Large” Language Models (LLMs)

Expert LM fine-tuned on a single task can outperform a multi-task model trained on 300 tasks



Race for one model to rule them all

“...we’re at the end of the era where it’s going to be these, like, giant, giant models,” - Sam Altman, OpenAI CEO

“More companies would be better served focusing on smaller, specific models that are cheaper to train and run.” - Clement Delangue, HuggingFace CEO

“small optimized models in healthcare are as accurate as large ones while being much more efficient.” - David Talby, Spark-NLP CEO



Models to fit every business need

Intel Xeon AI for LLMs

- Enable the most cost effective and ubiquitous approach
- Fine-tune and optimize inference models on Intel Xeon processors
- Leverage hundreds of Intel and 3rd party pretrained models

Summary



- Use Intel solutions to run edge AI efficiently and at scale
- Leverage Intel 4th Generation Xeon based instances with AMX capabilities for AI
- Enable Intel SW optimizations for data processing, training and inference
- Visit us at our booth: **Intel Booth #750**

Connect With Our Partners

Visit our partners' booths to see how they use Intel technology.

CDW	— Booth #305
HPE	— Booth #630
IBM	— Booth #930
SingleStore	— Booth #1586
World Wide Technology	— Booth #131
Nutanix	— Booth #132
SUSE	— Booth #501

Notices & Disclaimers

Intel is committed to the continued development of more sustainable products, processes, and supply chain as we strive to prioritize greenhouse gas reduction and improve our global environmental impact. Where applicable, environmental attributes of a product family or specific SKU will be stated with specificity. Refer to the 2022 Corporate Responsibility Report (p. 64) for further information.

Performance varies by use, configuration and other factors. Learn more at www.intel.com/PerformanceIndex.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Code names are used by Intel to identify products, technologies, or services that are in development and not publicly available. These are not "commercial" names and not intended to function as trademarks.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

Useful Links and Resources

AMX Enabled Amazon EC2 m7i Instances: <https://aws.amazon.com/about-aws/whats-new/2023/08/amazon-ec2-m7i-flex-m7i-instances/>

Public Proof Points for AI with 4th Gen Xeon: <https://www.intel.com/content/www/us/en/newsroom/news/4th-gen-intel-xeon-outperforms-competition-real-world-workloads.html>

OpenVino open-source Toolkit for Edge AI inference: <https://www.intel.com/content/www/us/en/newsroom/news/4th-gen-intel-xeon-outperforms-competition-real-world-workloads.html>

Leverage AMX to accelerate TensorFlow and PyTorch workloads: <https://www.intel.com/content/www/us/en/developer/articles/technical/accelerate-tensorflow-ml-performance-amx.html#gs.00me4r;> <https://www.intel.com/content/www/us/en/developer/articles/technical/accelerate-pytorch-training-inference-on-amx.html>

Optimizing Stable diffusion with Intel and Hugging Face: <https://huggingface.co/blog/train-optimize-sd-intel>

Performance Data for Intel AI Data Center Products: <https://www.intel.com/content/www/us/en/developer/topic-technology/artificial-intelligence/performance-4th-gen-56c-xeon-max.html>

Sustainability Technical Paper (Intel-Forrester): <https://www.intel.com/content/www/us/en/newsroom/news/4th-gen-intel-xeon-outperforms-competition-real-world-workloads.html#gs.00mfhg>