# AI Edge
## How traffic patterns change in the age of AI

Johan Ervenius
Systems Engineer, Arista Networks

ARISTA

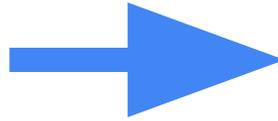# AI Acronyms

Training

Inference
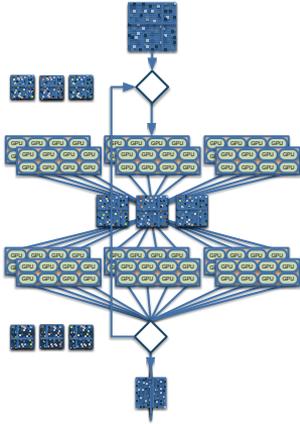
LLM

Generative AI

Agentic AI
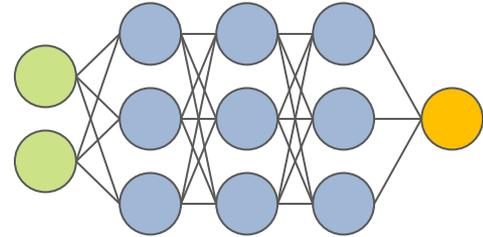
NLP

Neural Networks

Computer Vision

# AI Training vs Trained

Training

Trained
Model

Within the Data Center

ARISTA

# Use of the Trained model



Inference

Unknown data → Trained Model → New output

# Online Inference
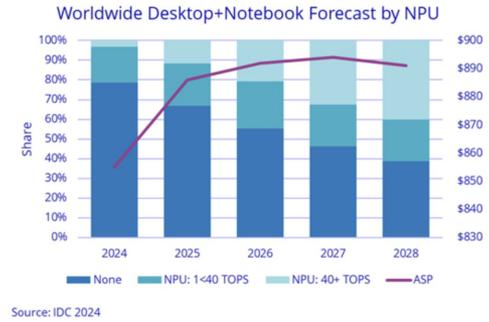


Trained Model

Data Center

ARISTA

# NPUs Everywhere Enabling Pervasive AI At The Edge

## Laptop/PC Edge



IDC: "…we see AI PCs (with NPUs) ramping up to **nearly two out of every three shipped in 2028**."
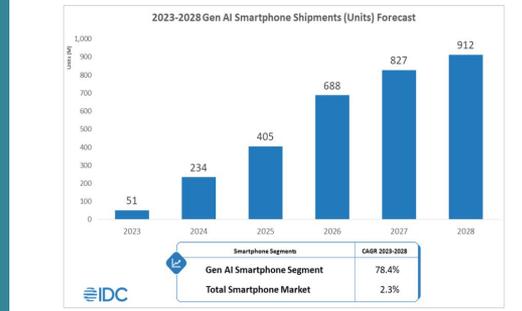
Worldwide Desktop+Notebook Forecast by NPU

**Increased Bandwidth & Power**

## Mobile Edge



By 2028, worldwide Gen AI smartphone shipments will reach 912 million units, resulting in a compound annual growth rate (CAGR) of 78.4% for 2023-2028.

2023-2028 Gen AI Smartphone Shipments (Units) Forecast

| Smartphone Segments | CAGR 2023-2028 |
|---|---|
| Gen AI Smartphone Segment | 78.4% |
| Total Smartphone Market | 2.3% |

**AI Application Performance at Edge**

## IoT Edge



By 2028, worldwide Gen AI smartphone shipments will reach 912 million units, resulting in a compound annual growth rate (CAGR) of 78.4% for 2023-2028.
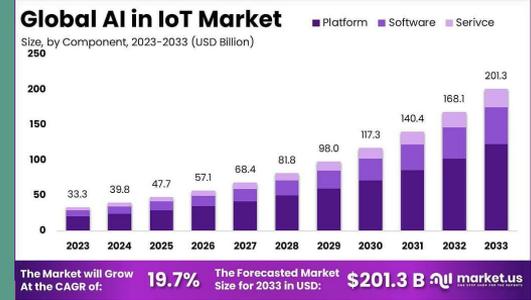
**Global AI in IoT Market**
Size, by Component, 2023–2033 (USD Billion)

The Market will Grow At the CAGR of: **19.7%** The Forecasted Market Size for 2033 in USD: **$201.3 B**

**Increased exposure to perimeter breaches / data loss**

ARISTA

# AI in Retail

## AI Applications

Customer Experience

Video analytics

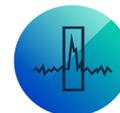Inventory Optimization

## Network Requirements

Traffic Asymmetry
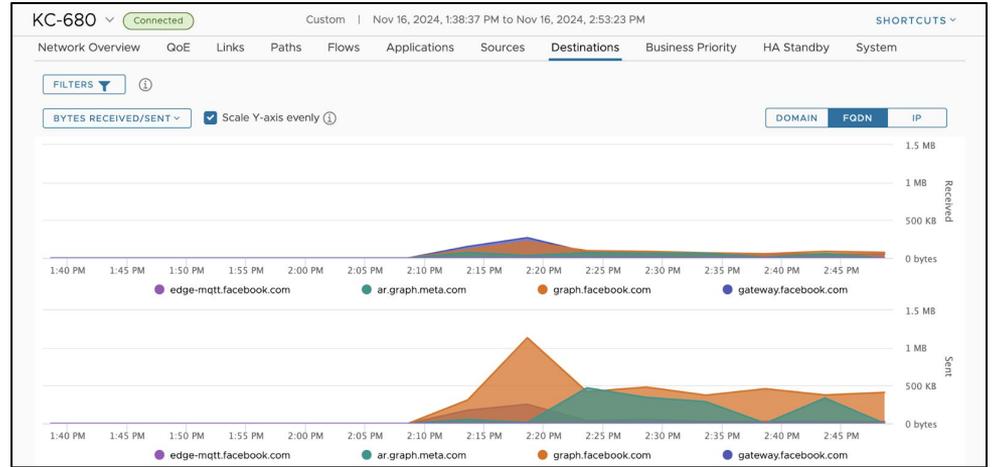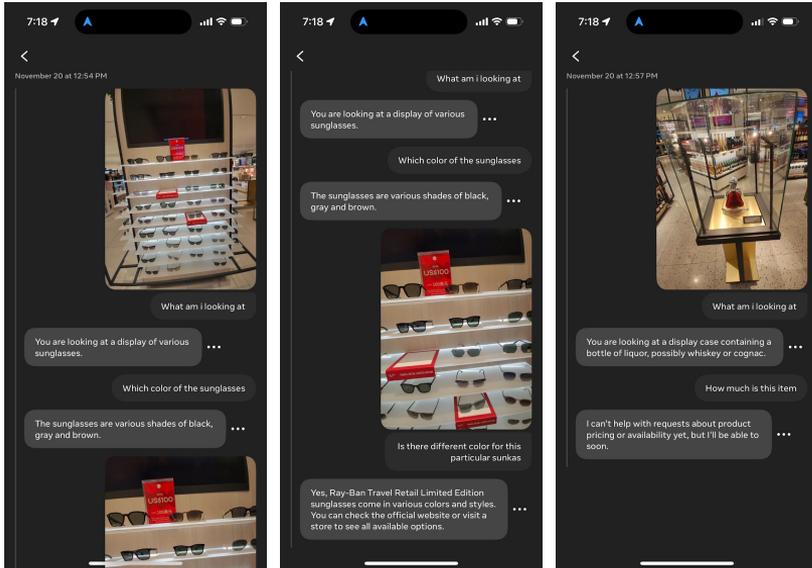10:1 upload to download

Increased Bandwidth Requirements

Latency Sensitive

Bursty Traffic

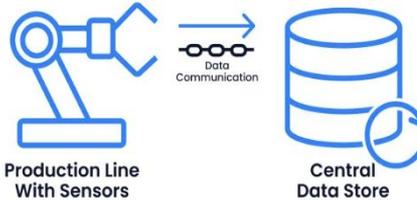# AI Traffic Is Now Driving New Requirements

# AI in Manufacturing

## AI Applications

AI analyzes sensor data from machinery to predict equipment failures

AI algorithms identify defects in real-time

AI predicts demand, optimizes inventory levels



Production Line With Sensors

Data Communication

Central Data Store

## Network Requirements

Traffic Asymmetry
Increased upload

Increased Bandwidth Requirements

Latency Sensitive

Bursty Traffic

ARISTA

# AI in Healtcare

## AI Applications

Robotic Surgery

Remote patient
monitoring

Supporting with
diagnostics

Clinical support



## Network Requirements

Traffic Asymmetry
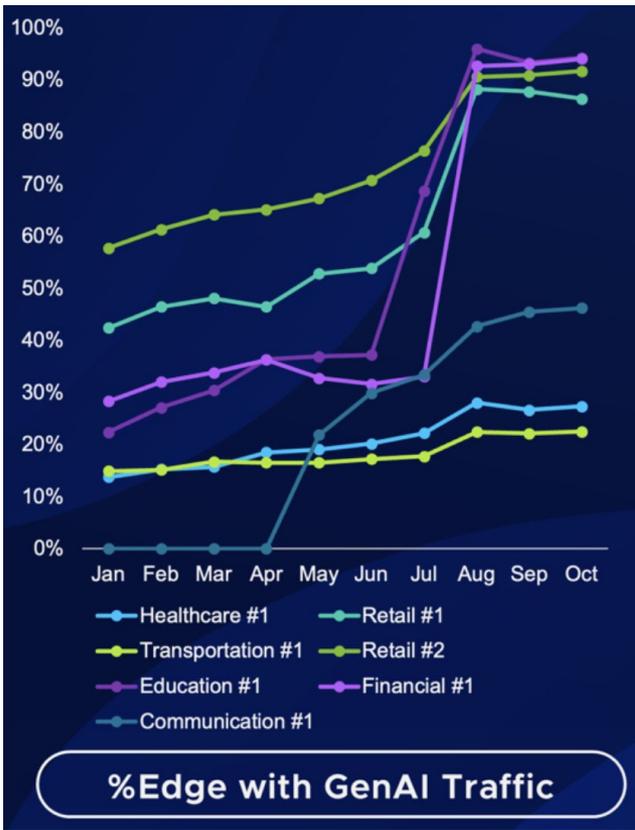Increased upload to
download

High Bandwidth
Requirements

Very Latency
Sensitive

# Generative AI Consumption has Increased Significantly



%Edge with GenAI Traffic



GenAI volume/site trending up



Increase in upload traffic %

ARISTA

# AI Network Bandwidth Impact

**Before**
Upload/Download 1:10

Upload

Download

**2026**
Upload/Download 10:1

Upload

Download

# Enterprise WAN is getting increasingly complex

# Agentic AI interactions are on the rise

Amazon CodeWhisperer

Microsoft Copilot

Chat GPT

Google Gemini

GEN AI

GCP

AWS

IAAS / PAAS

LLM

DATA CENTER

SLM

BRANCH

LEASED LINES

SLM

BRANCH

LLM

DATA CENTER

SLM

BRANCH

SLM

BRANCH

LLM

DATA CENTER

SLM

BRANCH

LEASED LINES

LLM

DATA CENTER

SLM

BRANCH

LLM

DATA CENTER

SLM

BRANCH

BRANCH

TLS

TLS

TLS

TLS

TLS

TLS

TLS

LLM

LLM

ARISTA

# AI Edge Application Optimization



**What is on the network?**

Deep Application Recognition

**What paths are available?**

Secure Overlay

**How are the paths performing?**

Link Qualification

**What is the best path for the application?**
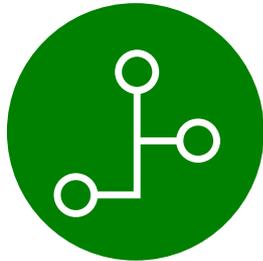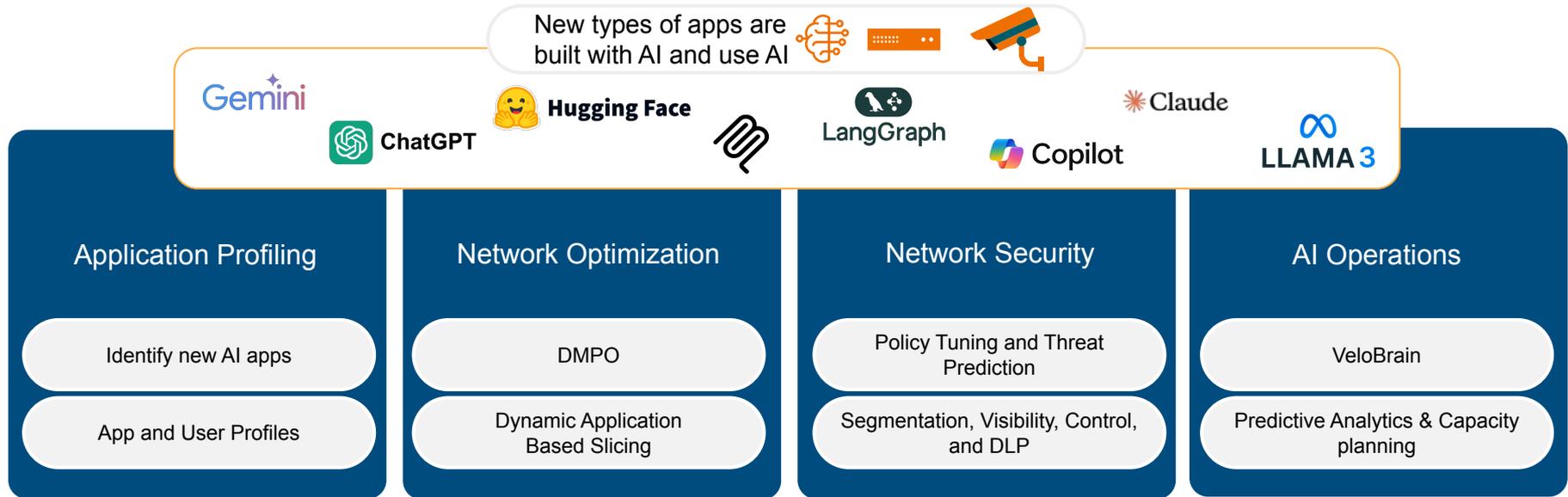
Application Steering

**Reacting to changed link conditions**

On-Demand Remediation & Aggregation

ARISTA

# Deliver an AI-Ready SD-WAN Architecture

You need a networking architecture designed to optimize performance, security, and scalability for distributed AI workloads at the edge - VeloRAIN from Arista



New types of apps are built with AI and use AI

Gemini
ChatGPT
Hugging Face
LangGraph
Claude
Copilot
LLAMA 3

| Application Profiling | Network Optimization | Network Security | AI Operations |
|---|---|---|---|
| Identify new AI apps | DMPO | Policy Tuning and Threat Prediction | VeloBrain |
| App and User Profiles | Dynamic Application Based Slicing | Segmentation, Visibility, Control, and DLP | Predictive Analytics & Capacity planning |

ARISTA

# Thank you!

ARISTA