



EMERGE 2019

ART OF THE POSSIBLE

April 23, 2019

Transforming Through AI

Anthony Dina

Director of Data Analytics / AI, Dell Technologies



Agenda

1 The compelling need for AI

2 The infrastructure view of AI

3 AI at the Edge

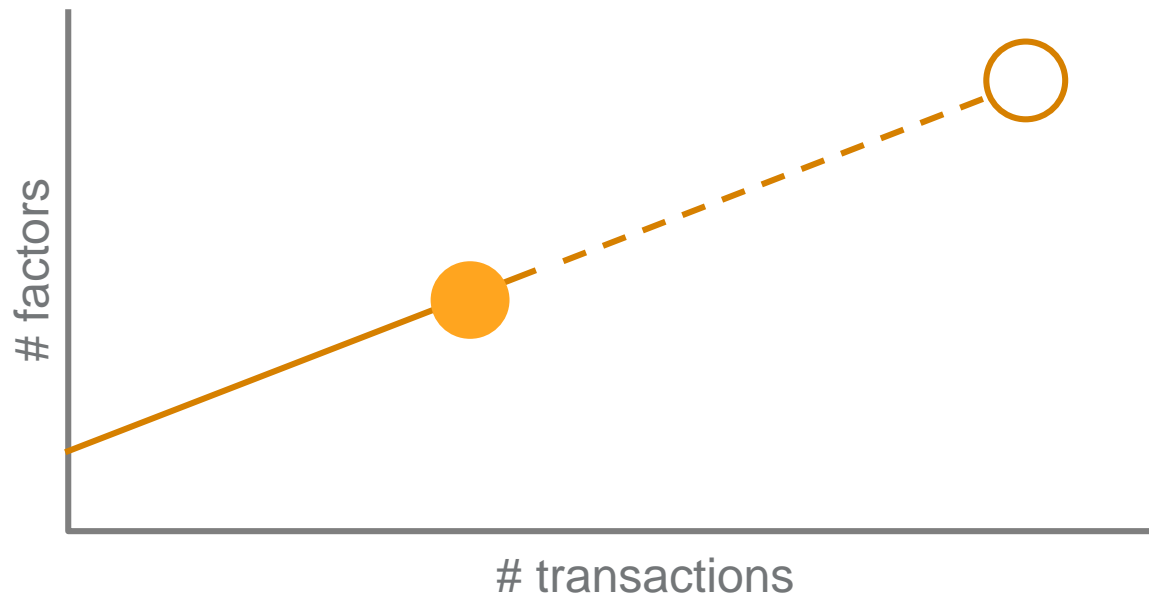
4 What's next



Nick Curcuru
Vice President of Data Analytics and
Cyber Security, Master Card

BUSINESS PROBLEMS PUSH US INTO AI





Fraud is estimated **\$24.26 billion in losses** in 2019 - is projected to rise to **\$34.66 billion in 2022** ¹
Card-not-present fraud is now **81%** more likely than card-present ²



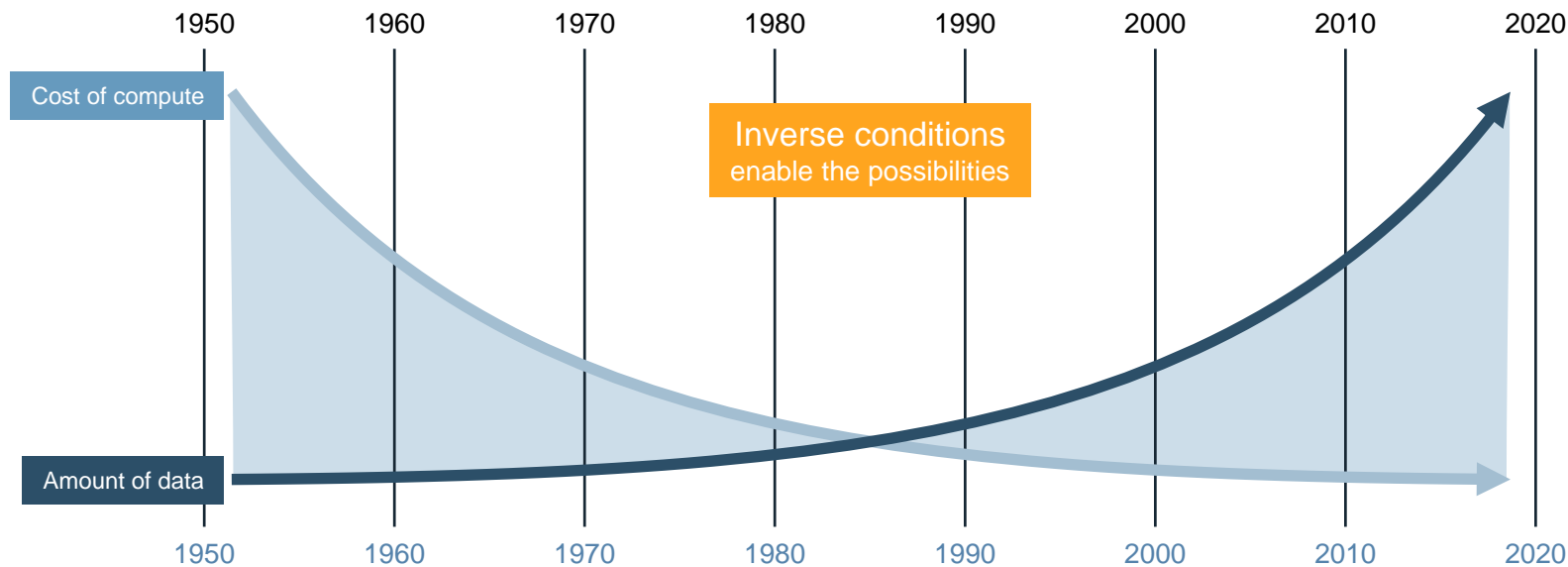
¹ Nilson Report, a newsletter

² Javelin Strategy & Research

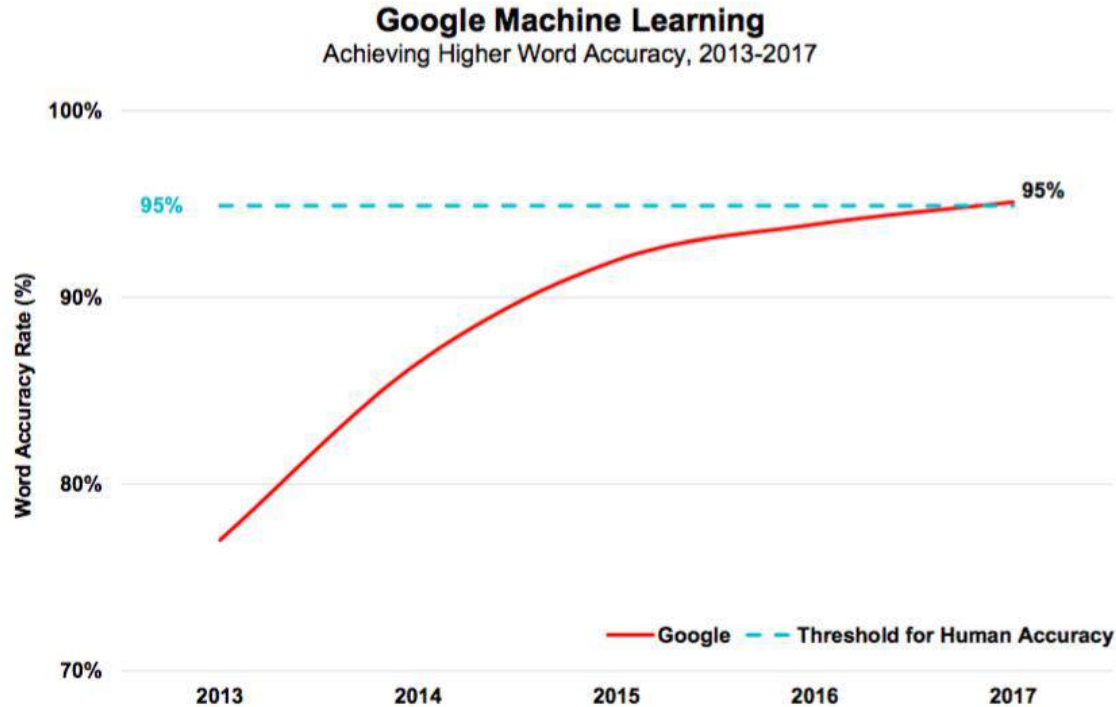
AI FITS WITH FEDERAL NEEDS AS WELL

 Healthcare	 Financial	 National Security	 Emergency Management
Case Automation Fraud Detection/Prevention	Tax Audit Fraud Detection/Prevention Investment predictions Customer service Network security	Facial recognition Video surveillance Cyber security Satellite imagery	Emergency Services Event Prediction

YES AI AGAIN, BUT THIS TIME...



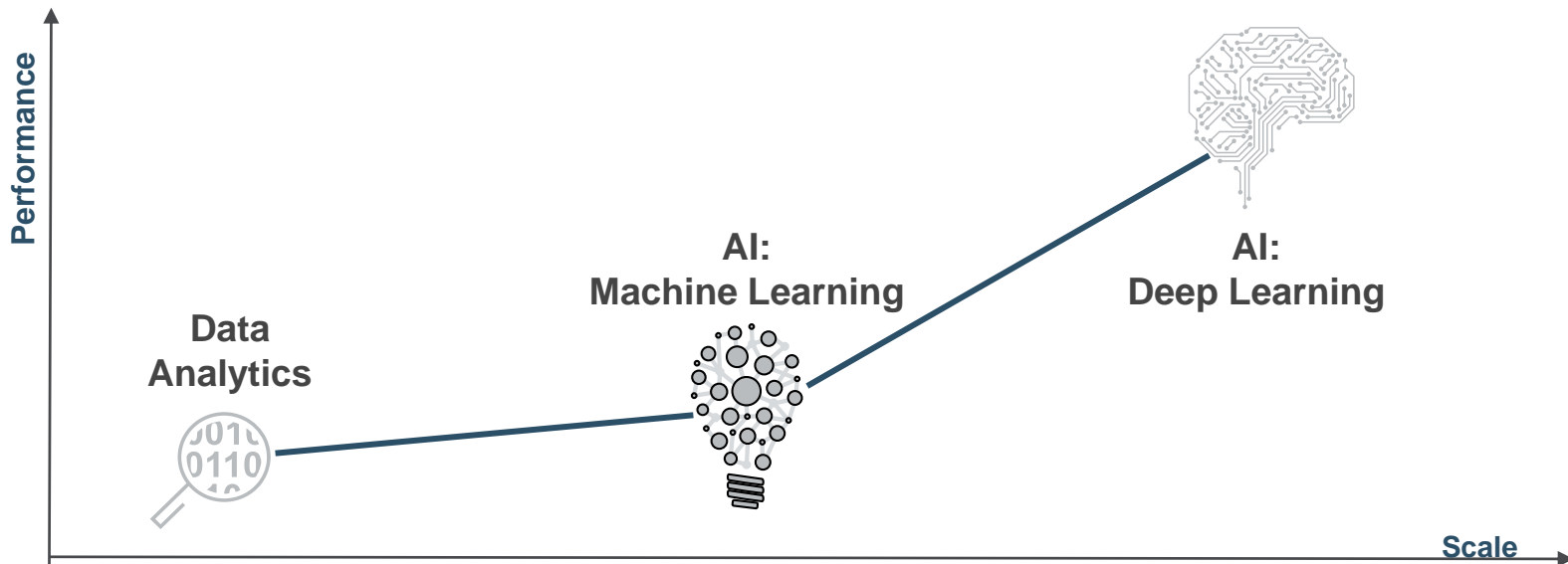
AI IS ABOUT ACCURACY AT SCALE



Source: Google (5/17)
Note: Data as of 5/17/17 and refers to recognition accuracy for English language. Word error rate is evaluated using real world search data which is extremely diverse and more error prone than typical human dialogue.

KP INTERNET TRENDS 2017 | PAGE 48

ACCURACY AT SCALE HAS A COST



Analytics Requirements

Compute: **CPU**
Throughput: **MB/s**
Concurrency: **10s**
Scale: **GBs to TBs**
Data Type: **Structured**

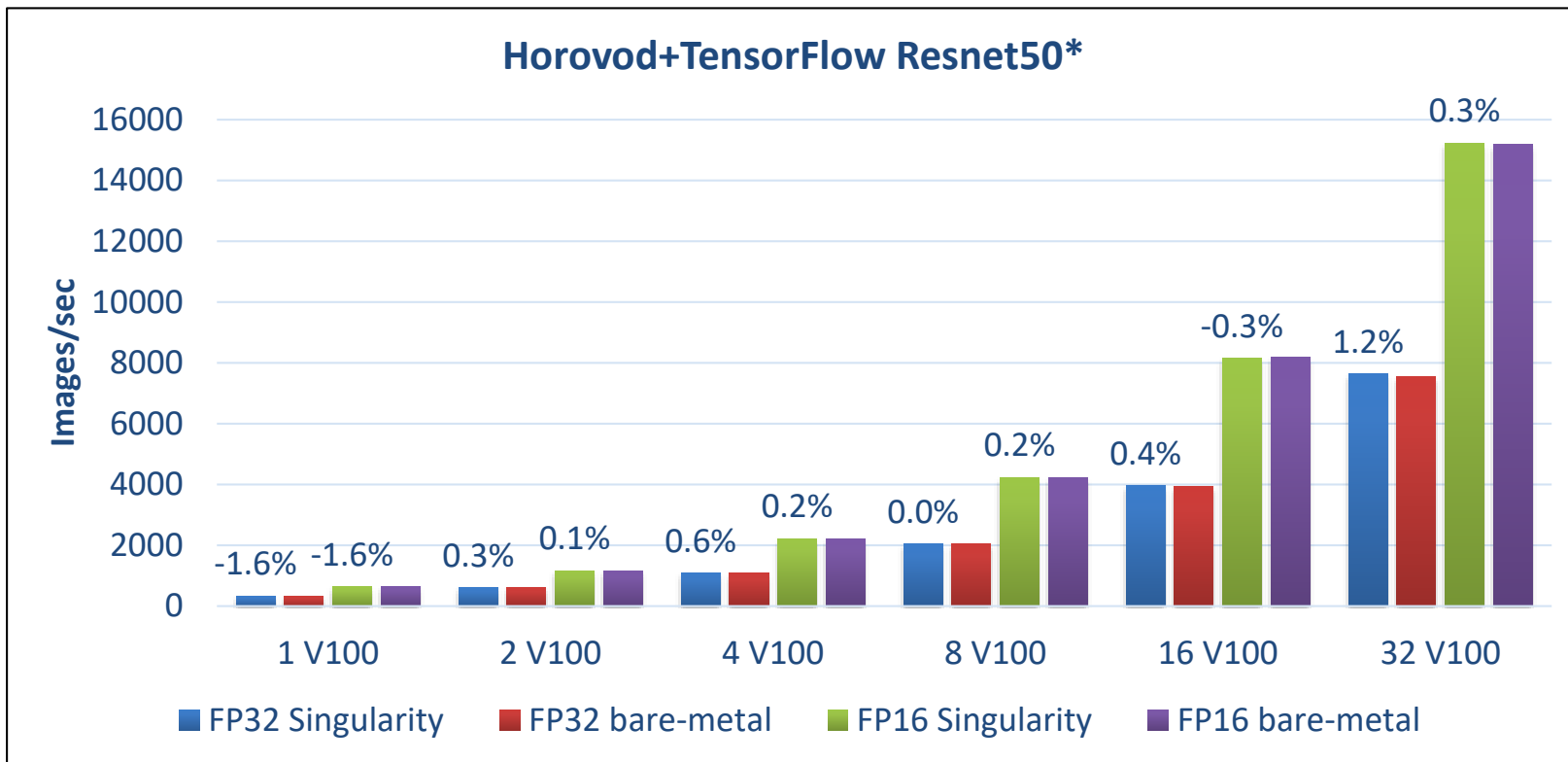
ML Requirements

Compute: **CPU & GPU**
Throughput: **GB/s**
Concurrency: **100s–K's**
Scale: **10s TB's–PB's**
Data Type: **Semi-structured, labeled**

DL Requirements

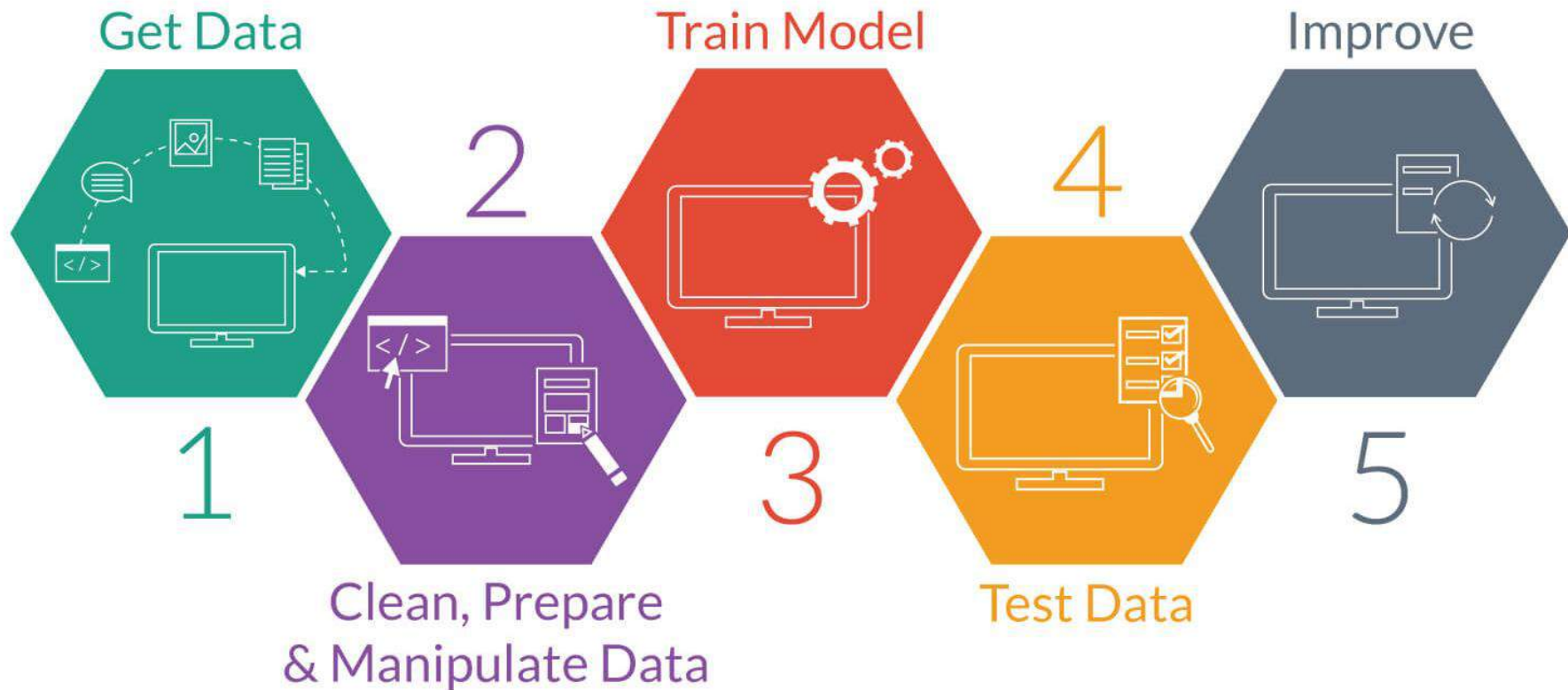
Compute: **10's-1000's GPU, FPGA, etc.**
Throughput: **10's-100's GB/s**
Concurrency: **1000s – M's**
Scale: **100s TB's to 10s PB's**
Data Type: **Unstructured, unlabeled**

ACCURACY AT SCALE IS POSSIBLE: 32 V100 = 22.4X

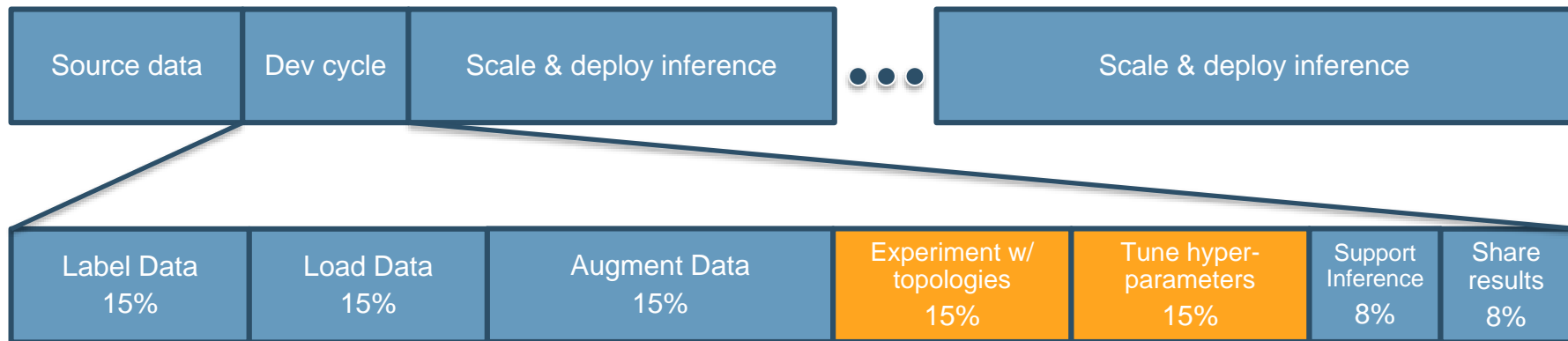


*Note: In FP32 mode, batch size: 128 per GPU; In FP16 mode, batch size: 256 per GPU
 Note: MPI Used for multi-node communication

AI PROCESS



AI LIFECYCLE



Data Collection

Managing the Stack

Data Cleansing

Training Models

What about AI at the Edge?

IS THIS THE EDGE?



IS THIS THE EDGE?



<https://latinousa.org/2017/02/01/us-customs-border-protection-publishes-specific-qa-trump-immigration-executive-order/>

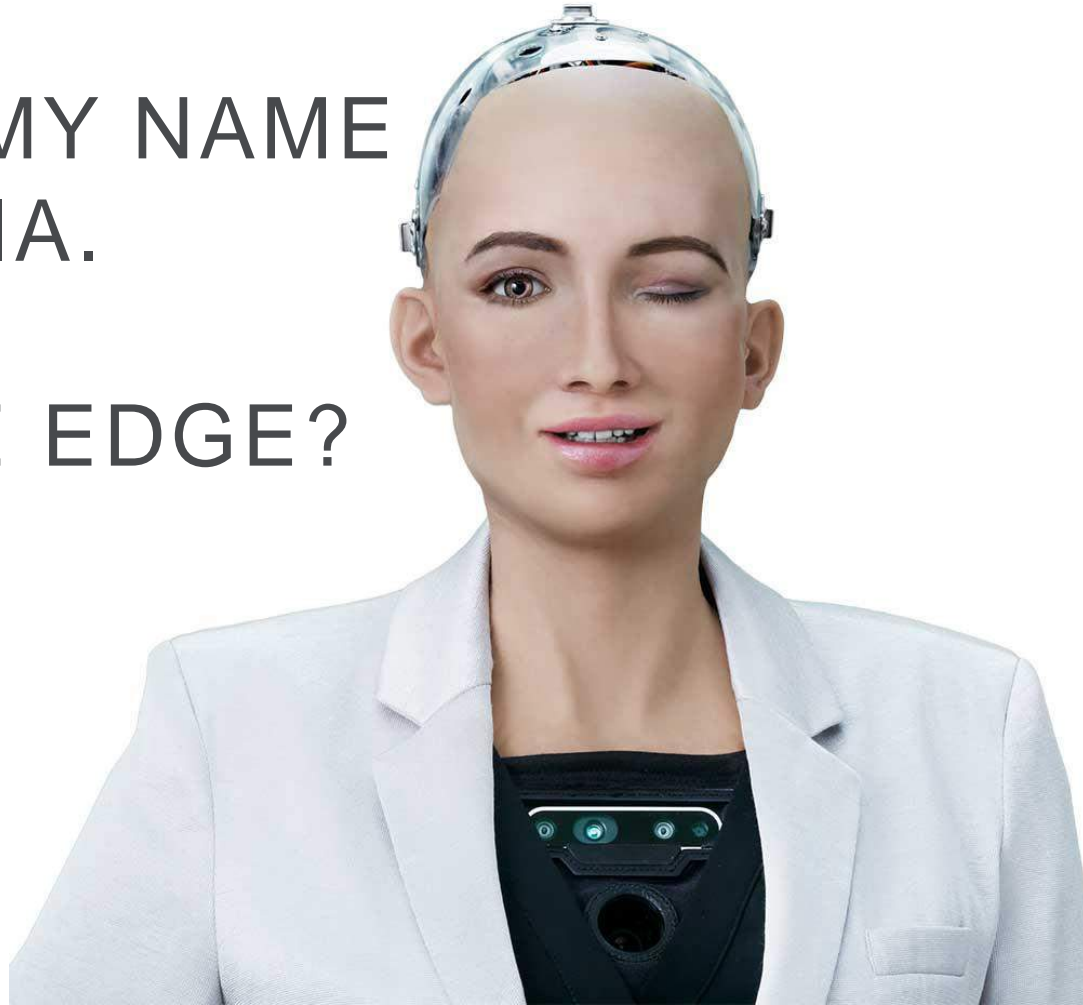
IS THIS THE EDGE?



<https://www.rferl.org/a/u-s-warship-donald-cook-heads-to-black-sea-for-security-operations-/29720003.html>

HELLO. MY NAME
IS SOPHIA.

AM I THE EDGE?



EDGE

CORE

CLOUD

Varies to <1ms

<5ms

<10-<40ms

~100ms

Factors at the edge:

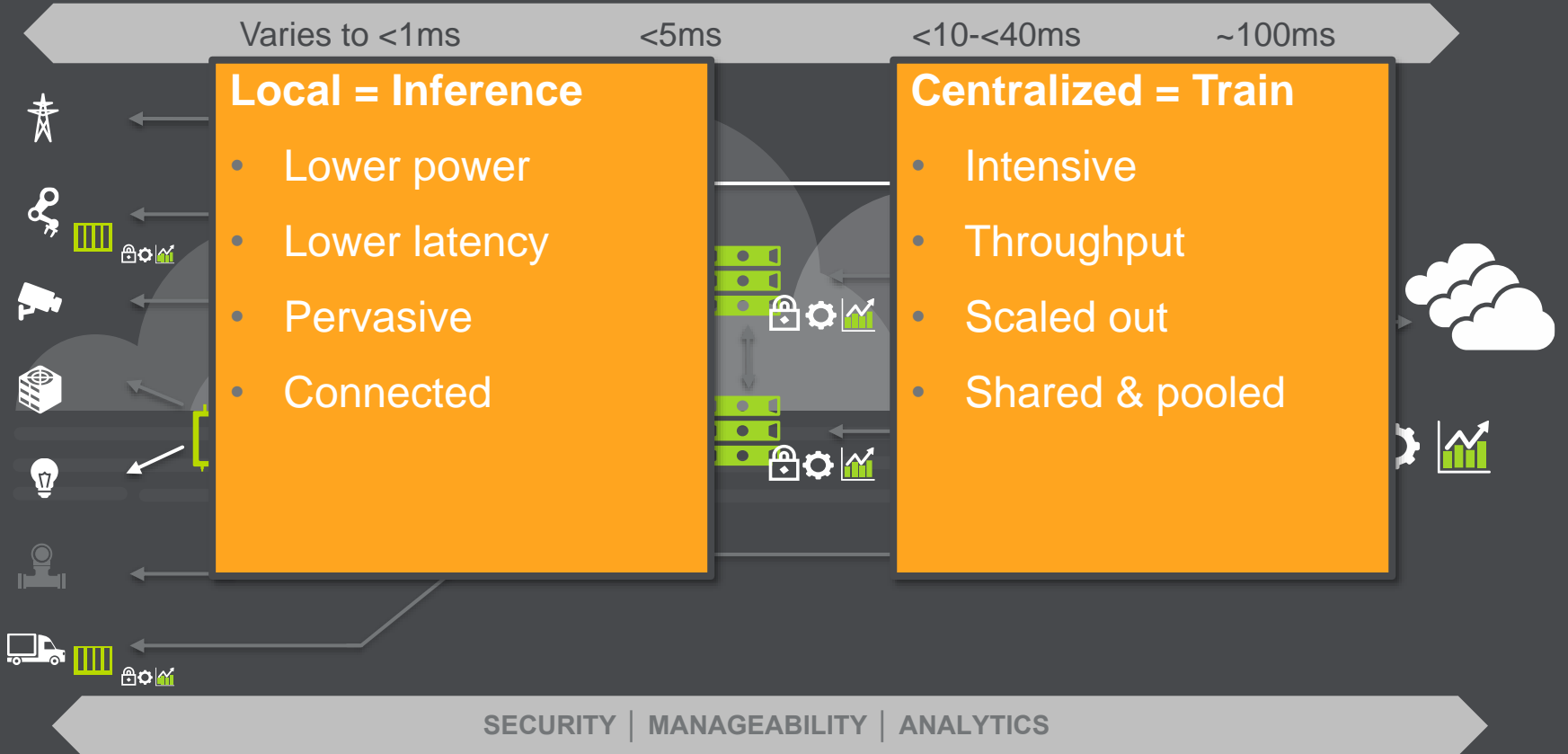
Power, Data; Models / logic; Comms; Security; Management

SECURITY | MANAGEABILITY | ANALYTICS

EDGE

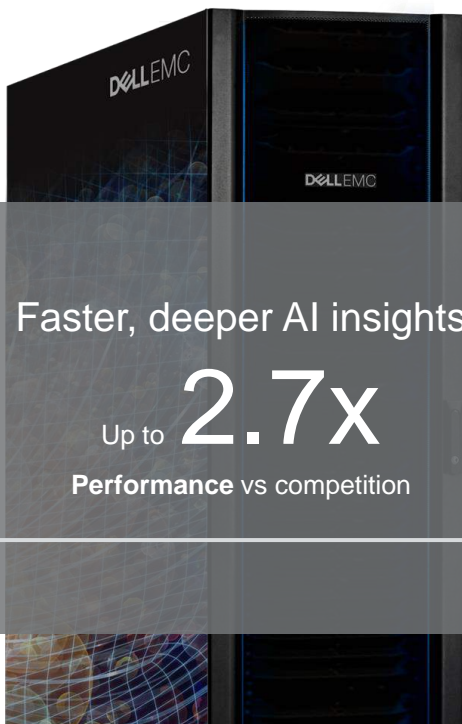
CORE

CLOUD



Ready Solutions for AI

Validated stack built to handle most demanding AI workloads



Simpler AI experience

30%

Improved data scientist productivity

Faster, deeper AI insights

Up to 2.7x

Performance vs competition

Proven AI expertise

98%

Lower training time

Where do we go from here?

WHERE DO WE GO FROM HERE?

Today

- World of clickstream
- E.g. chatbots; recommendation engines; etc.
- Lambda architectures: speed and batch
- Gateways / regional DC
- Train centrally; infer locally (Models pushed down)

Emerging

- World of smart IoT
- E.g. autonomous vehicle
- Speed layer is the only layer (Kappa)
- Decisions at ingest
- Mobile / Micro DC
- Train locally (Models built at edge then moved centrally)

The Possible

- World of meshed intelligence
- E.g. Swarm drones
- High speed P2P networks
- Continuously mobile DC
- Clusters built at edge
- Autonomous learning (models moved laterally)

D~~E~~LL EMC