# ARTIFICIAL INTELLIGENCE
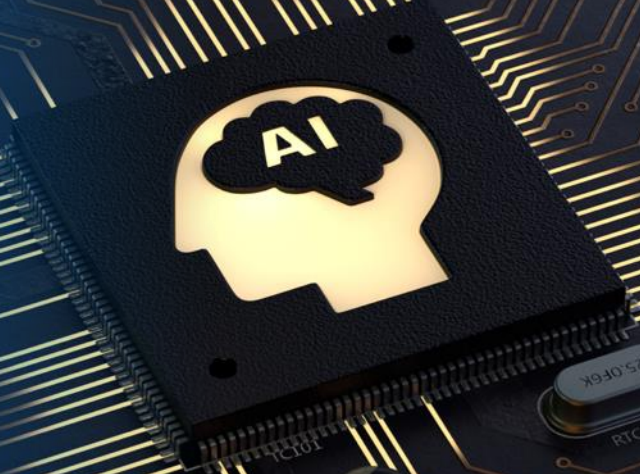
September 2019 - ITU

Paul Haines, Thailand Account Executive
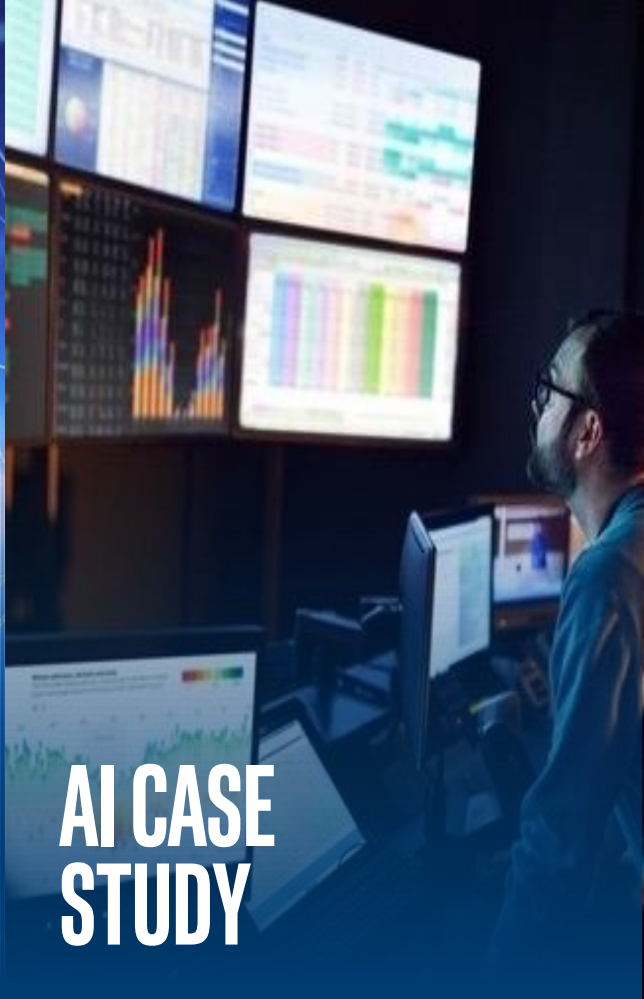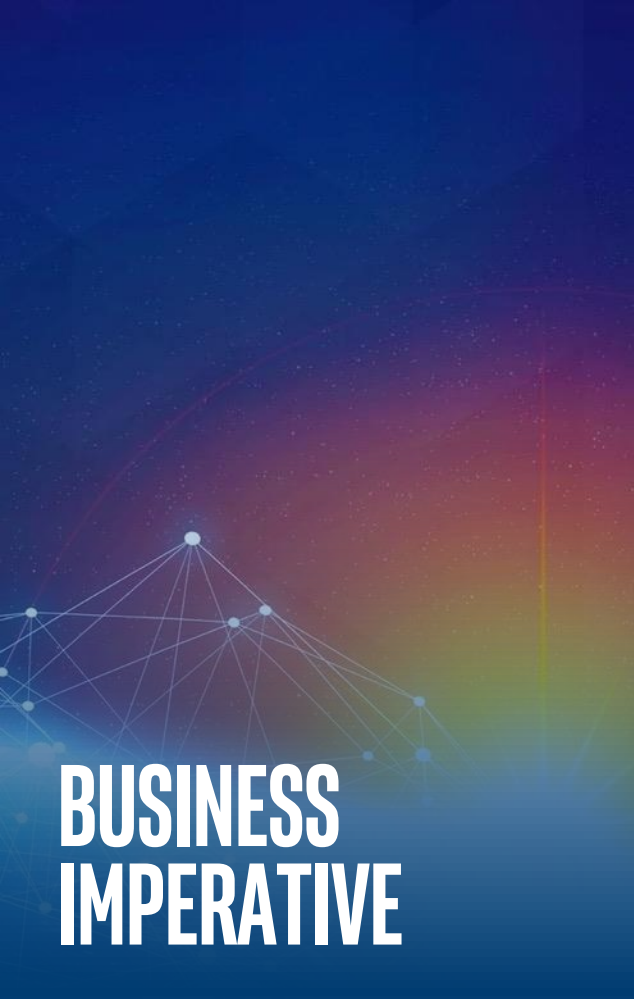
BUSINESS IMPERATIVE

INTEL® AI

AI CASE STUDY

BUSINESS
IMPERATIVE

INTEL® AI

AI CASE
STUDY

# The AI Mandate

"
AI technologies are evolving fast and growing increasingly **critical** to firms' ability to win, serve, and retain customers.
"

**FORRESTER**

"
...strategic technologies for 2019 with the potential to drive significant **disruption** and deploy **opportunity** over the next five years
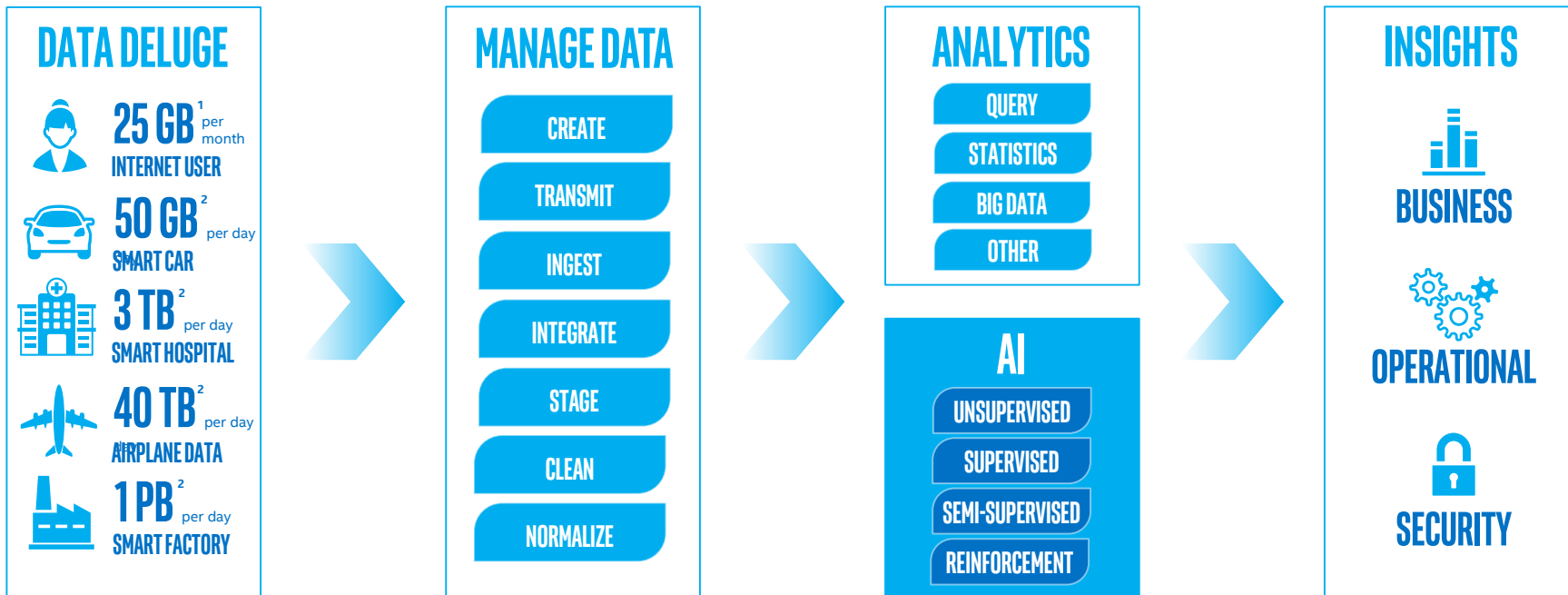"

**GARTNER**

"
...**70%** of CIOs will aggressively apply data and AI to IT operations, tools, and processes by 2021.
"

**IDC**

## THE TIME TO BEGIN AI ADOPTION IS NOW

# Why AI?

## DATA DELUGE

- **25 GB** [1] per month — INTERNET USER
- **50 GB** [2] per day — SMART CAR
- **3 TB** [2] per day — SMART HOSPITAL
- **40 TB** [2] per day — AIRPLANE DATA
- **1 PB** [2] per day — SMART FACTORY

## MANAGE DATA

- CREATE
- TRANSMIT
- INGEST
- INTEGRATE
- STAGE
- CLEAN
- NORMALIZE

## ANALYTICS

- QUERY
- STATISTICS
- BIG DATA
- OTHER

## AI

- UNSUPERVISED
- SUPERVISED
- SEMI-SUPERVISED
- REINFORCEMENT
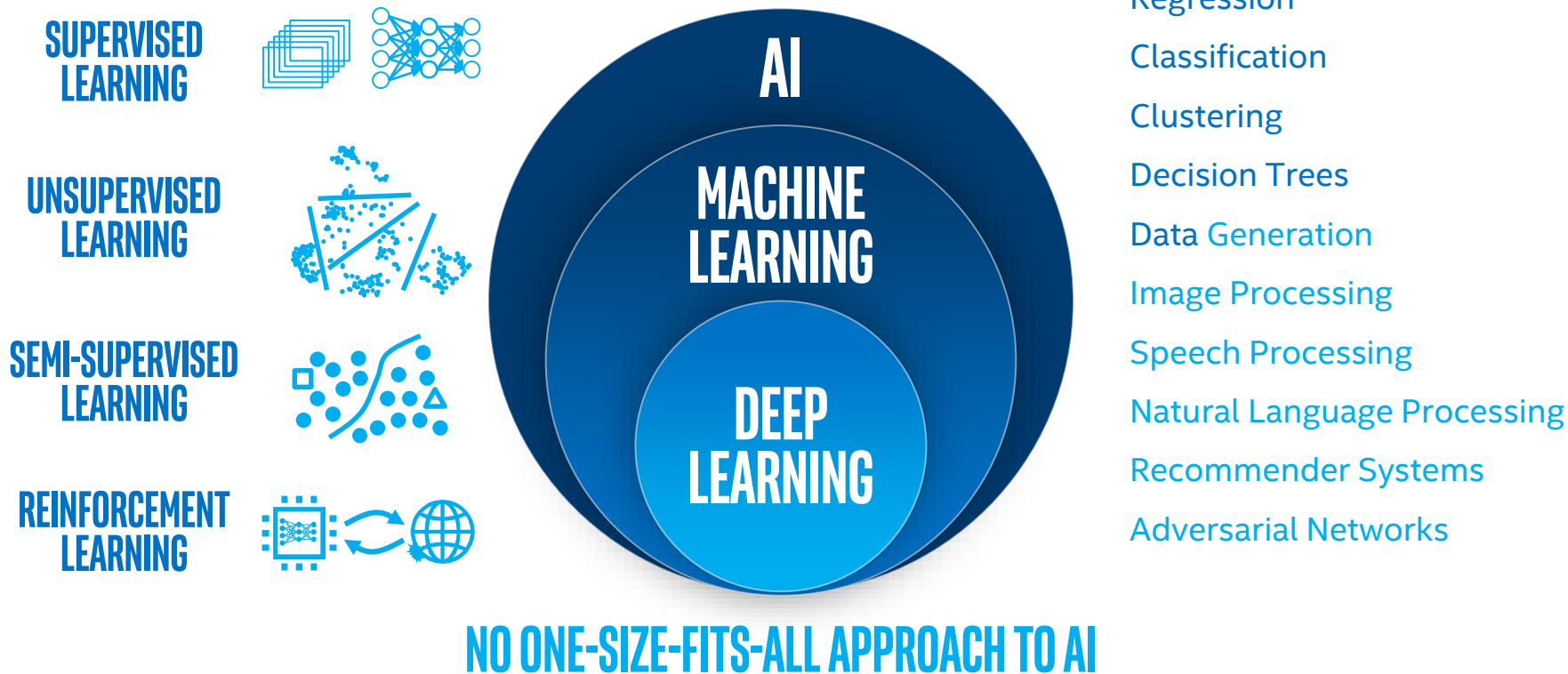
## INSIGHTS

- BUSINESS
- OPERATIONAL
- SECURITY

## EXTRACT VALUABLE INSIGHTS FROM DATA

1. Source: http://www.cisco.com/c/en/us/solutions/service-provider/vni-network-traffic-forecast/infographic.html
2. Source: https://www.cisco.com/c/dam/m/en_us/service-provider/ciscoknowledgenetwork/files/547_11_10-15-DocumentsCisco_GCI_Deck_2014-2019_for_CKN__10NOV2015_.pdf

# What is AI?

**SUPERVISED LEARNING**

**UNSUPERVISED LEARNING**

**SEMI-SUPERVISED LEARNING**

**REINFORCEMENT LEARNING**

AI

MACHINE LEARNING

DEEP LEARNING

Regression

Classification

Clustering

Decision Trees

Data Generation

Image Processing

Speech Processing

Natural Language Processing

Recommender Systems

Adversarial Networks

**NO ONE-SIZE-FITS-ALL APPROACH TO AI**

# AI Closer Look

## MACHINE LEARNING

*Algorithms designed to deploy better insight with more data*

**Regression** (Linear/Logistic)

**Classification** (Support Vector Machines/SVM, Naïve Bayes)

**Clustering** (Hierarchical, Bayesian, K-Means, DBSCAN)

**Decision Trees** (RandomForest)

**Extrapolation** (Hidden Markov Models/HMM)

**More…**

## DEEP LEARNING

*Neural networks used to infer meaning from large dense datasets*

**Image Recognition** (Convolutional Neural Networks/CNN, Single-Shot Detector/SSD)

**Speech Recognition** (Recurrent Neural Network/RNN)

**Natural Language Processing** (Long-Short Term Memory/LSTM)

**Data Generation** (Generative Adversarial Networks/GAN)

**Recommender System** (Multi-Layer Perceptron/MLP)

**Time-Series Analysis** (LSTM, RNN)

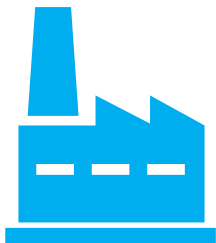**Reinforcement Learning** (CNN, RNN)

**More…**

## REASONING

*Hybrid of analytics & AI techniques designed to find meaning in diverse datasets*

**Associative Memory** (Intel® Saffron AI memory base)

← **See also:** machine & deep learning techniques

**More…**

# Which Approach is Best?

**SMART FACTORY**

| QUESTION | METHOD | APPROACH |
|---|---|---|
| How many parts should we manufacture? | Historical supply and demand analysis | Statistical Analytics |
| What will our production yield be? | Algorithm learns which variables correlate to yield | Machine Learning *(Unsupervised)* |
| Which parts have visual defects? | Algorithm learns to identify defects in images | Deep Learning *(Supervised)* |
| Can my robotic arm learn to get better? | Algorithm that acts and adapts based on feedback | Deep Learning *(Reinforcement)* |

## CHOOSE THE RIGHT AI APPROACH FOR YOUR CHALLENGE

# AI Solutions in Every Market

**AGRICULTURE**
Achieve higher yields and increase efficiency

**ENERGY**
Maximize production and uptime

**EDUCATION**
Transform the learning experience

**GOVERNMENT**
Enhance safety, research, and more

**FINANCE**
Turn data into valuable intelligence

**HEALTH**
Revolutionize patient outcomes

**INDUSTRIAL**
Empower truly intelligent Industry 4.0

**MEDIA**
Create thrilling experiences

**RETAIL**
Transform stores and inventory

**SMART HOME**
Enable homes that see, hear, and respond

**TELECOM**
Drive network and operational efficiency

**TRANSPORT**
Automated driving

## OUR PARTNERS ARE DRIVING REAL-WORLD VALUE WITH INTEL® AI

BUSINESS
IMPERATIVE

INTEL® AI

AI CASE
STUDY

# AI Opportunities are Diverse



**DEVICE**

**INTELLIGENT EDGE**

**MULTI-CLOUD**

| THINGS | ON PREMISE | REGIONAL DATA CENTERS | PUBLIC/HYBRID/PRIVATE DC |

Personal, industrial & community devices

Data analysis & action & gateways

Local private, cloud & telco

Large scale centralized

# Intel® AI Strategy



## VIBRANT COMMUNITY

Drive innovative use cases

Pioneer leading-edge AI

Fuel the ecosystem

## INDUSTRY & OPEN SOFTWARE

Optimize customer software

Unify APIs across Intel

Empower developers

## PLATFORM WITH BEST HARDWARE

Extend the CPU

Lead in acceleration

Build a common platform

# Accelerate Your AI Journey with Intel

**DISCOVER**
Get started faster with community support

**DATA**
Tame the deluge with a modern data layer

**MODEL**
Speed up development with open AI software

**DEPLOY**
Deliver on the best AI hardware for your needs

## COMMUNITY

**CONSULT** — Intel® AI

**PARTNER** — AI Builders / AI In Production

**LEARN** — AI Developer Program

## SOFTWARE

**DATA MANAGEMENT** — Choice of 50+ Optimized Tools for Data Preparation

**MACHINE LEARNING** — Intel® DAAL / Intel® Distribution for Python* / ANALYTICS ZOO

**DEEP LEARNING** — TensorFlow / BigDL / ONNX / OpenVINO / mxnet / PyTorch / MODEL ZOO / & More

## HARDWARE

**MOVE** — intel Silicon Photonics / intel Omni-Path Fabric / intel Ethernet

**STORE** — intel OPTANE DC PERSISTENT MEMORY / intel OPTANE DC SOLID STATE DRIVE

**PROCESS** — intel XEON PLATINUM / FPGA / GPU / intel NERVANA inside / intel MOVIDIUS inside

# Get Started Faster
## with community support

## CONSULT

intel AI

Consult with your Intel and supplier representative(s) to learn more

Visit: plan.seek.intel.com/SMARTForm_ICS

## PARTNER

### AI BUILDERS
*(Cloud to Device)*

### AI IN PRODUCTION
*(IOT Edge/Device)*

Partner with an Intel® AI provider and/or access a catalog with >100 solutions

Visit: builders.intel.com/ai
software.intel.com/ai-in-production

## LEARN

### AI DEVELOPER PROGRAM

Learn AI skills with the FREE¥ Intel® AI Developer Program, including cloud access

Visit: software.intel.com/ai

# Tame the Deluge
## with a modern data layer

See also:
Analytics Gold Deck

DATA

**DATA MANAGEMENT**

| CREATE | TRANSMIT | INGEST | INTEGRATE | STAGE | CLEAN | NORMALIZE |

**ANALYTICS?**
Analytics Gold Deck

**END-TO-END:**

SAP*, Microsoft*, Oracle*, SAS*, Cloudera*, IBM*...

**TRENDING:**

C3IoT*
Thoughtspot*
Streamsets*
Confluent* (Kafka*)...

BlueData*
MemSQL*
RedisLabs*
Cassandra*
Pandas*...

Aerospike*
MarkLogic*
Splunk*
Spark*
SKLearn*...

**AI?**
See Next Slide>>

**SOLUTIONS:**

30+ AAI & HPC Solutions (Genomics, ICPD, Splunk...)

**CPU**

**Visit:** www.intel.com/analytics

# Speed Up Development
## with open AI software



**MACHINE LEARNING** — **DEEP LEARNING**

| | MACHINE LEARNING | DEEP LEARNING |
|---|---|---|
| **TOOLKITS**<br>App Developers | ANALYTICS ZOO | MODEL ZOO · OpenVINO™ |

| | **Intel® Data Analytics Acceleration Library** (DAAL) | **Intel® Distribution for Python*** (Sklearn*, Pandas*) | **R** (Cart, Random Forest, e1071) | **Distributed** (MlLib on Spark, Mahout) | **Frameworks** | **Intel Tools** |
|---|---|---|---|---|---|---|
| **LIBRARIES**<br>Data Scientists | | | | | TensorFlow* · BigDL · Caffe · ONNX · mxnet · PyTorch<br>*More framework optimizations in progress…* | NAUTA<br>RL Coach<br>NLP Architect<br>NN Distiller |

| | KERNELS | | |
|---|---|---|---|
| **KERNELS**<br>Library Developers | Intel® Math Kernel Library (Intel® MKL) | Intel® Machine Learning Scaling Library (Intel® MLSL) · Intel® Math Kernel Library for Deep Neural Networks (Intel® MKL-DNN) | nGraph |

**CPU** — **CPU ▪ GPU ▪ FPGA ▪ ACCELERATOR**

**Visit:** [www.intel.ai/technology](www.intel.ai/technology)

# Deploy with Unprecedented
## AI hardware choice

## MOVE FASTER

INTEL® SILICON PHOTONICS

INTEL® ETHERNET

INTEL® OMNI-PATH FABRIC

## STORE MORE

INTEL® OPTANE™ SSD

INTEL® OPTANE™ DC
PERSISTENT MEMORY

## PROCESS EVERYTHING

CPU

GPU, FPGA

ACCELERATORS

**Visit:** www.intel.ai/technology

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.
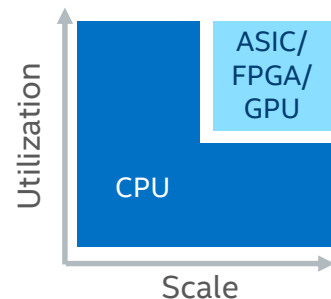*Other names and brands may be claimed as the property of others

# Intel® AI Hardware



DEVICE    INTELLIGENT EDGE    MULTI-CLOUD

## OPTIMIZED FRAMEWORKS & SOFTWARE

| CPU | GPU | FPGA | ASIC | | |
|-----|-----|------|------|------|------|
| intel CORE 10TH GE / XEON PLATINUM inside | intel GPU FUTURE | intel FPGA | intel NERVANA inside NNP-I COMING SOON | intel NERVANA inside NNP-T COMING SOON | intel MOVIDIUS inside |

**← WORKLOAD BREADTH        AI SPECIALIZATION →**

| Multi-Purpose Foundation for AI | Data-Parallel Media, Graphics, HPC & AI | Multi-Function & Real-time Deep Learning Inference | Deep Learning Inference | Deep Learning Training | Media & Vision DL Inference at the Edge |

*Visit:*

# Intel® AI Use Cases

## CPU
### Intel® Xeon® Scalable Processors

**MULTI-CLOUD**

| | | |
|---|---|---|
| **JD.com*** | **HYBRID ANALYTICS + AI** Fast time-to-solution on Spark* with MlLib & BigDL | |
| **CERN*** | **HPC AND AI** Fast time-to-solution for deep learning in classic workflows | |
| **Novartis*** | **LARGE DL TRAINING** Fast DL training for large image recognition in drug discovery | |
| **Taboola*** | **DEEP LEARNING INFERENCE** High throughput real-time recommendation (billion items) | |
| **Ziva*** | **MACHINE LEARNING** Animating movie creatures using machine learning techniques | |

## CPU
### Intel® Xeon® Scalable Processors

**INTELLIGENT EDGE**

| | |
|---|---|
| **GE Health*** | **DEEP LEARNING INFERENCE** Low TCO for image recognition in CT scanner for radiology |

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.
*Other names and brands may be claimed as the property of others.

# Intel® AI Use Cases

## FPGA
Intel® FPGA

| MULTI-CLOUD | **Microsoft*** | **REAL-TIME REC. ENGINE**<br>Real-time recommendations and more workload acceleration |
|---|---|---|
| | **Manjeera*** | **REAL-TIME TRANSCRIPTS**<br>Real-time transcription acceleration |
| | **JD.com*** | **TEXT RECOGNITION**<br>Faster time-to-market for custom CNN & LSTM for end-to-end text recognition |

| INTELLIGENT EDGE | **QNAP*** | **VISION INFERENCE**<br>Faster time-to-market for custom CNN workload with OpenVINO™ toolkit |
|---|---|---|
| | **NEC*** | **FACE RECOGNITION**<br>Faster time-to-market for custom CNN workload for surveillance and retail |
| | **Alibaba*** | **REAL-TIME VISION**<br>Real-time video encoding and decoding for smart city project |

## ASIC
Intel® Movidius™ Myriad™ X VPU

| INTELLIGENT EDGE | **HPE*** | **VISION AT THE EDGE**<br>Video analytics and DL inference in an edge server blade |
|---|---|---|

| DEVICE | **Hikvision*** | **VISION IN THE DEVICE**<br>Deep learning-based computer vision at low power |
|---|---|---|

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.
*Other names and brands may be claimed as the property of others.

# AI Compute Considerations

## WORKLOADS



What is my
workload profile?

## REQUIREMENTS



RESPONSE TIME
VERSATILITY **SECURITY** DATA
**SCALABILITY** COST
SIZE POWER
RELIABILITY **THROUGHPUT** EFFICIENCY
MANAGEABILITY
FLEXIBILITY **ACCURACY**

What are my use case
requirements?

## DEMAND



How prevalent is AI
in my environment?

Note: word cloud source is www.wordart.com
¥Free = available to download/access at no cost to qualified developers who are enrolled in the program
*Other names and brands may be claimed as the property of others.

# Bust the Deep Learning Myth



*"A GPU is required for deep learning..."* **FALSE**

- Most enterprises (---) use <u>CPU</u> for machine and deep learning needs

- Some early adopters (---) may reach a deep learning tipping point when acceleration is needed[1]

1"Most" of enterprise customers based on survey of Intel direct engagements and internal market segment analysis

# Deep Learning Use Case

**Source Paper:**

research.fb.com/
wpcontent/uploads/2017/12/hpca-
2018-facebook.pdf

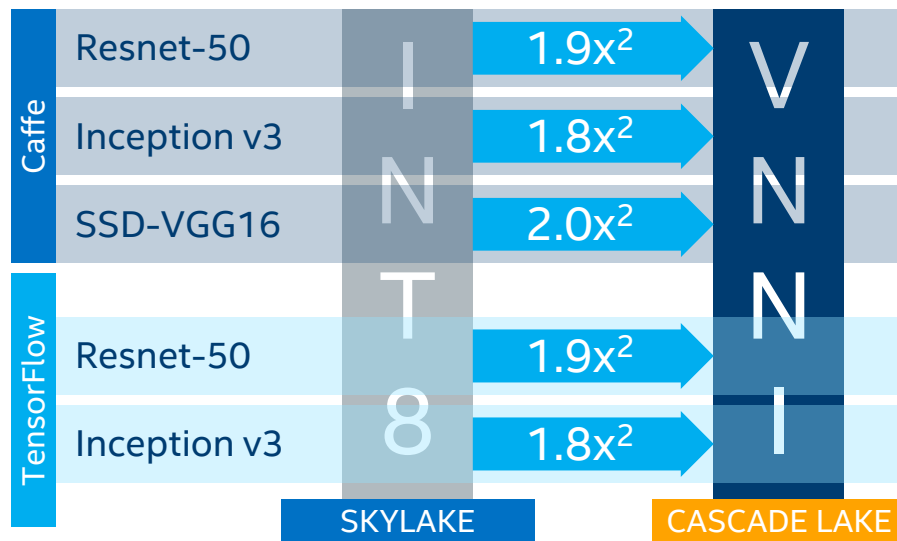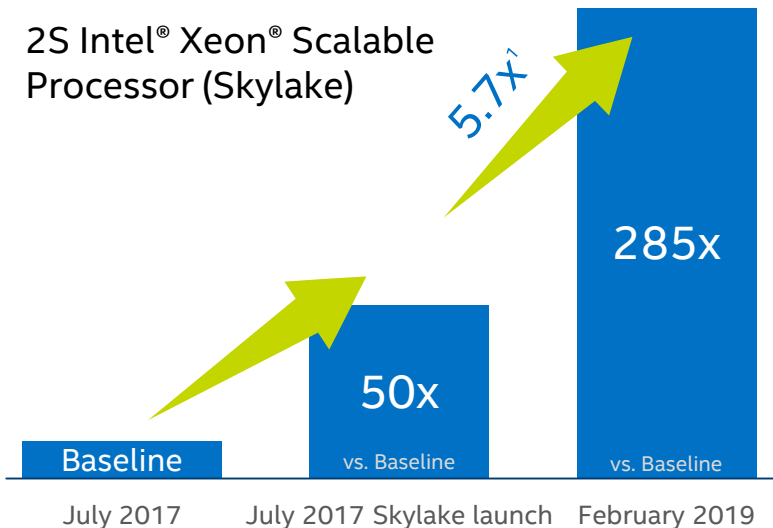| Services | Ranking Algorithm | Photo Tagging | Photo Text Generation | Search | Language Translation | Spam Flagging | Speech |
|---|---|---|---|---|---|---|---|
| **Model(s)** | MLP | SVM,CNN | CNN | MLP | RNN | GBDT | RNN |
| **Inference Resource** | CPU | CPU | CPU | CPU | CPU | CPU | CPU |
| **Training Resource** | CPU | GPU & CPU | GPU | Depends | GPU | CPU | GPU |
| **Training Frequency** | Daily | Every N photos | Multi-Monthly | Hourly | Weekly | Sub-Daily | Weekly |
| **Training Duration** | Many Hours | Few Seconds | Many Hours | Few Hours | Days | Few Hours | Many Hours |

## LARGE CLOUD USERS EMPLOY CPU EXTENSIVELY FOR DEEP LEARNING

# Deep Learning Performance on CPU
## Hardware + software improvements for Intel® Xeon® processors



2S Intel® Xeon® Scalable Processor (Skylake)

5.7x[1]

285x

50x

Baseline

vs. Baseline

vs. Baseline

July 2017

July 2017 Skylake launch

February 2019

| Caffe | Resnet-50 | | $1.9x^2$ | |
| | Inception v3 | | $1.8x^2$ | |
| | SSD-VGG16 | | $2.0x^2$ | |
| TensorFlow | Resnet-50 | INT8 | $1.9x^2$ | VNNI |
| | Inception v3 | | $1.8x^2$ | |

SKYLAKE          CASCADE LAKE

BUSINESS
IMPERATIVE

INTEL® AI

AI CASE
STUDY

# Intel® AI Case Study



ACCELERATE YOUR AI JOURNEY

1 DISCOVER  2 DATA  3 MODEL  4 DEPLOY

# Intel® AI Case Study

**DISCOVER**

**IDENTIFY**

*Identify* prospects internally and using the 70+ AI solutions in Intel's portfolio; then assess business value of each one

**PRIORITIZE**

*Prioritize* projects based on business value and cost to solve with Intel guidance; choose industrial defect detection via DL[1]

**CONSIDER**

*Consider* ethical, social, legal, security and other risks and mitigation plans with Intel advisors prior to kickoff

**ORGANIZE**

*Organize* internally to get buy-in, support new development philosophy and grow developer talent via Intel® AI

**Value (H)**

**Cost (L)**

Corrosion

L — H

**AI DEVELOPER PROGRAM**

# Intel® AI Case Study



**DATA**

**INGEST**

*Ingest* streaming data from drones using a popular software tool among the many that run on the CPU

**STORE**

*Store* data in block storage (for high-performance) in a data lake with guidance from an Intel storage partner

**PREPARE**

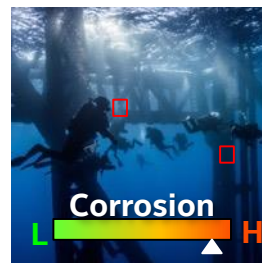*Prepare* data by performing cleanup and integration using popular software tools that run on the CPU

**ACT**

*Act* on the data using one of the many popular CPU tools for data analytics and visualization

011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100

**12 weeks**

| Source Data | Transmit Data | Ingest Data | Clean up Data | Integrate Data | Stage Data |
|---|---|---|---|---|---|
| 10% | 5% | 5% | 60% | 10% | 10% |

# Intel® AI Case Study

**MODEL** → **SETUP**

*Set up* compute environment; DL training (~7% of journey) acceleration <u>NOT</u> worthwhile due to high setup time & cost

**ACCELERATION ZONE**

Breakeven Threshold

**CPU ZONE**

intel XEON PLATINUM inside

DL Demand

Time

**YOU ARE HERE**

**MODEL**

*Model* development through training a deep neural network using an Intel-optimized DL framework

**TEST**

*Test* the deep learning model using a control data set to determine if accuracy meets requirements

**DOCUMENT**

*Document* the code, process, and key learnings for future reference

**12 weeks**

| **Train** (Topology Experiments) | **Train** (Tune Hyper-parameters) | **Test Inference** | **Document Results** |
|---|---|---|---|
| *30%* | *30%* | *20%* | *10%* |

# Intel® AI Case Study

**DEPLOY** · **ARCHITECT** · **IMPLEMENT** · **SCALE** · **ITERATE**

*Architect* AI deployment with Intel® AI Builders

*Implement* AI in production environment

*Scale* to more sites and users as demand grows

*Iterate* on the models with new data over time

| Data Ingest | Drones |
| Media Store | Prepare Data |
| Model Store | Training |
| Inference | Label Store |
| Media Server | Service Layer |

## REMOTE DEVICES

Drone
Drone
Drone

**10 Drones**

*Real-time object detection and data collection*

Drone
Drone
Drone

**Per Drone**

| 1x | Intel® Core™ processor |
| 1x | Intel® Movidius™ VPU |

## MEDIA SERVER

Media Store
Media Store
Media Store

**110 Nodes**

8 TB/day per camera

10 cameras

3x replication

1-year retention

4 mgmt nodes

Media Store
Media Store
Media Store

**Per Node**

| 1x | 2S 61xx | 20x | 4TB SSD |

## MULTI-USE CLUSTER

Data Ingestion
Data Ingestion
Data Ingestion
Data Ingestion
Inference
Inference
Inference
Inference
Prepare Data
Prepare Data
Service Layer
Service Layer
Service Layer
Media Server
Media Server
Media Server

**4 Nodes**
*One ingestion per day, one-day retention*

**4 Nodes**
*20M frames per day*

**2 Nodes**
*Infrequent op*

**3 Nodes**
*Simultaneous users*

**3 Nodes**
*10k clips stored*

## DATA STORE

Model Store
Model Store
Model Store
Model Store
Label Store
Label Store
Label Store
Label Store

**4 Nodes**
*1-year of history*

**4 Nodes**
*Labels for 20M frames /day*

**Per Node**

| 1x | 2S 81xx |
| 5x | 4TB SSD |

## ADV. ANALYTICS

Training

**16 Nodes**
*Intermittent use 1 training/month for <10 hours*

Training

**Per Node**

| 1x | 2S 81xx |
| 1x | 4TB SSD |

## SOFTWARE

➤ OpenVino™ Toolkit
➤ Intel® MKL-DNN

➤ TensorFlow*
➤ Intel® Movidius™ Software Development Toolkit

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.
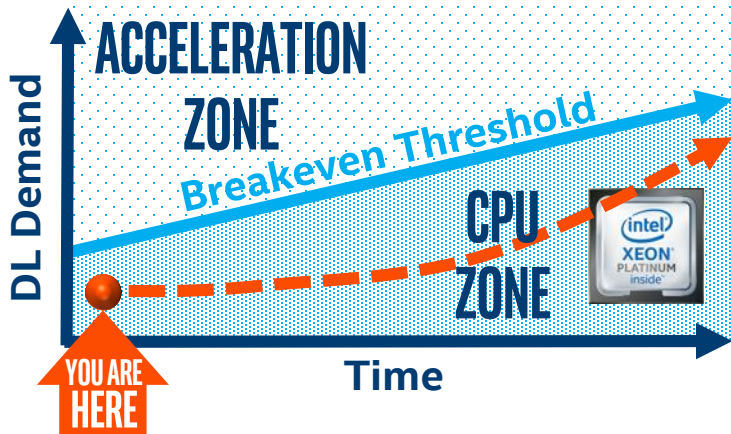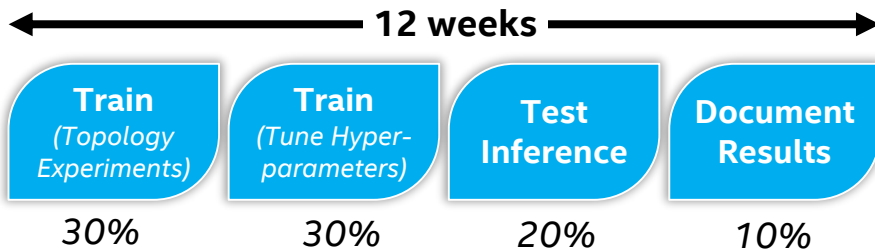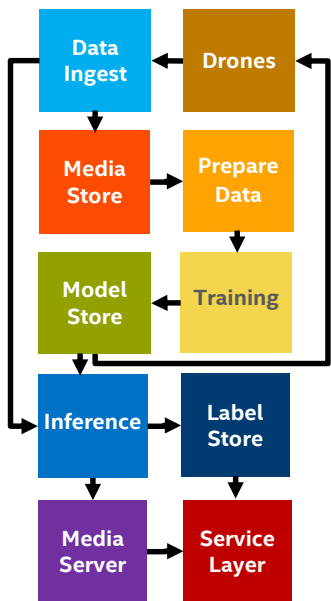*Other names and brands may be claimed as the property of others    Optimization Notice

THANK YOU