# HPC, Big Data, AI: What are the new needs ? Are all infrastructure solutions equal ?



## Alain Cyr, PhD
*Infrastructure Solutions Architect*
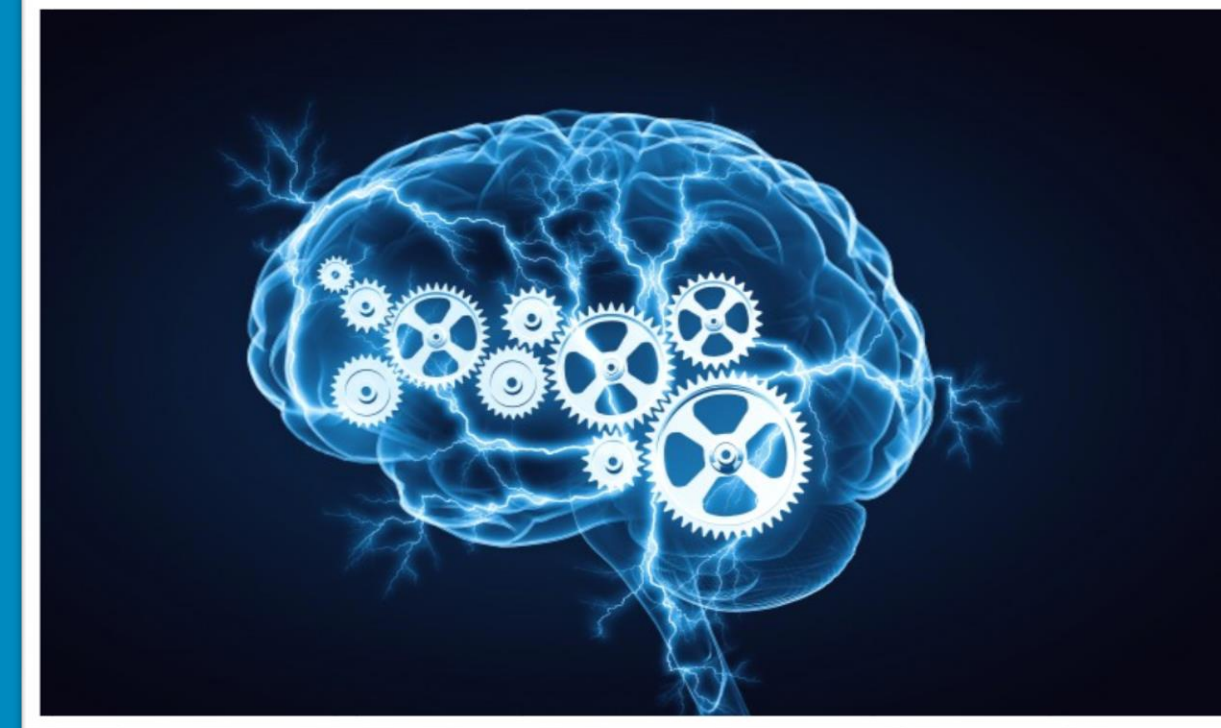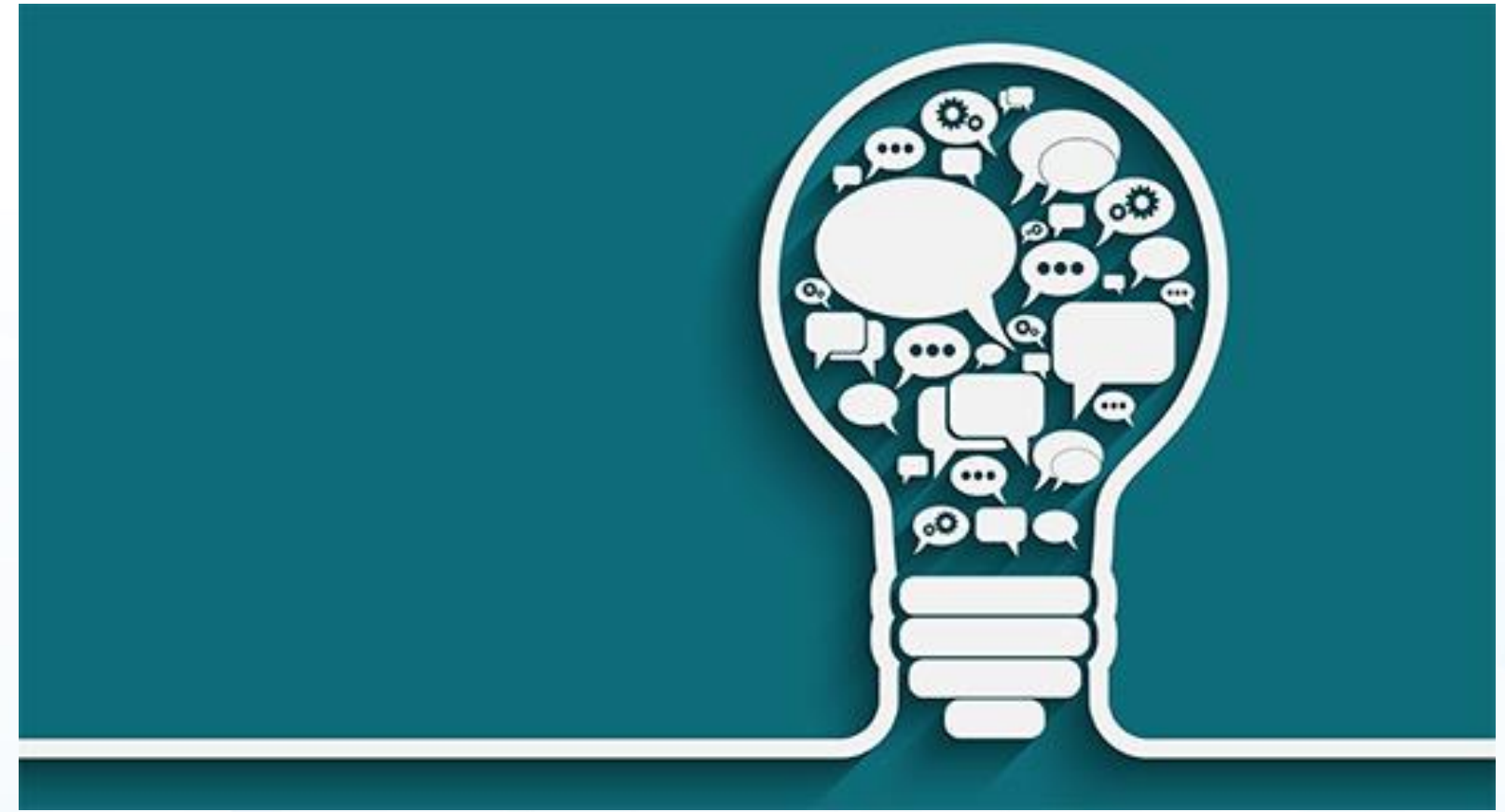*Montpellier IBM Client Center*

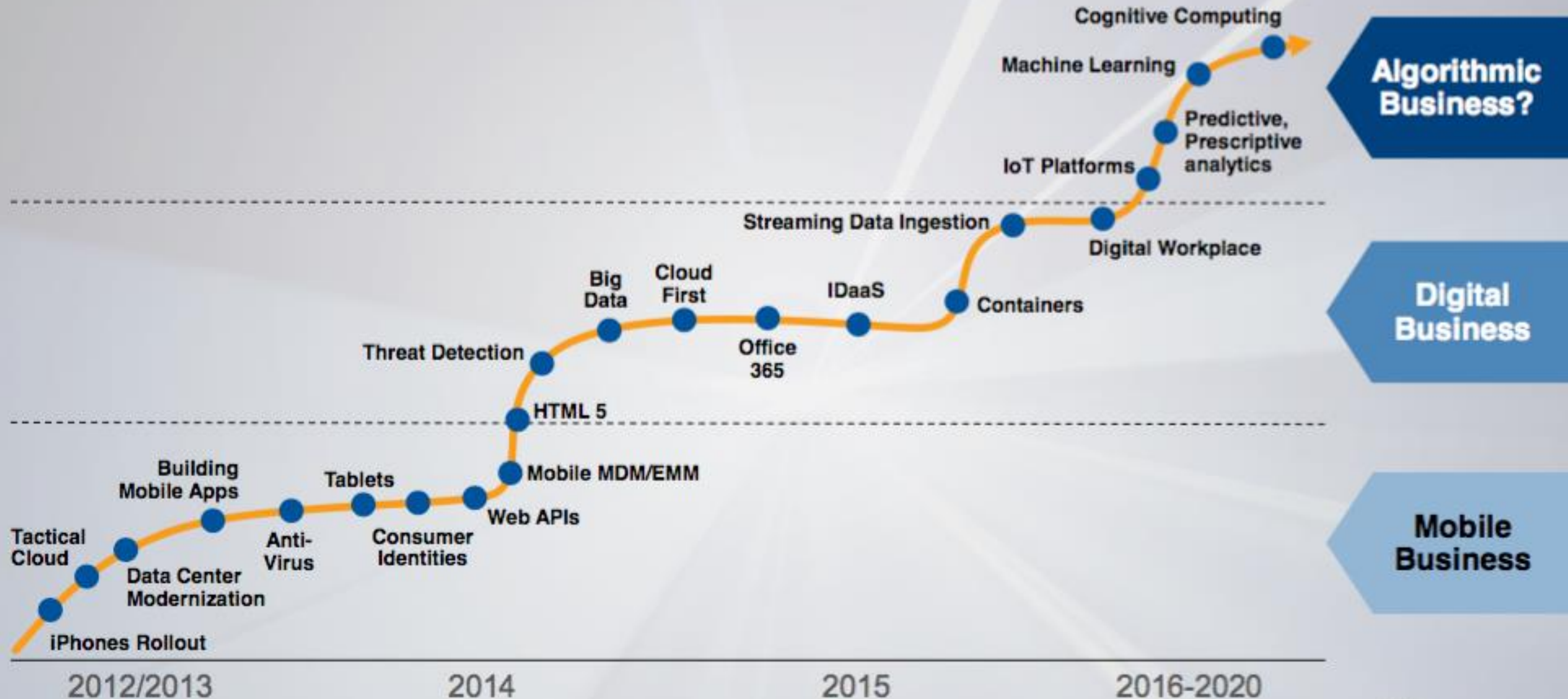cyralain@fr.ibm.com

@trollnyrd

POWER9

OpenPOWER™

# Agenda



- AI and Cognitive Use Cases

- Data driven needs: what are the right platforms ?

- Entreprise class offering for AI: PowerAI

# The Path to the Digital Business and Beyond

Cognitive Computing

Machine Learning

**Algorithmic Business?**

Predictive, Prescriptive analytics

IoT Platforms

Streaming Data Ingestion

Digital Workplace

Big Data

Cloud First

IDaaS

**Digital Business**

Threat Detection

Office 365

Containers

HTML 5

Building Mobile Apps

Tablets

Mobile MDM/EMM

Web APIs

**Mobile Business**

Tactical Cloud

Anti-Virus

Consumer Identities

Data Center Modernization

iPhones Rollout

2012/2013      2014      2015      2016-2020

Gartner

# Radiologists

**Overloaded** with medical imaging data.

Eye Fatigue.
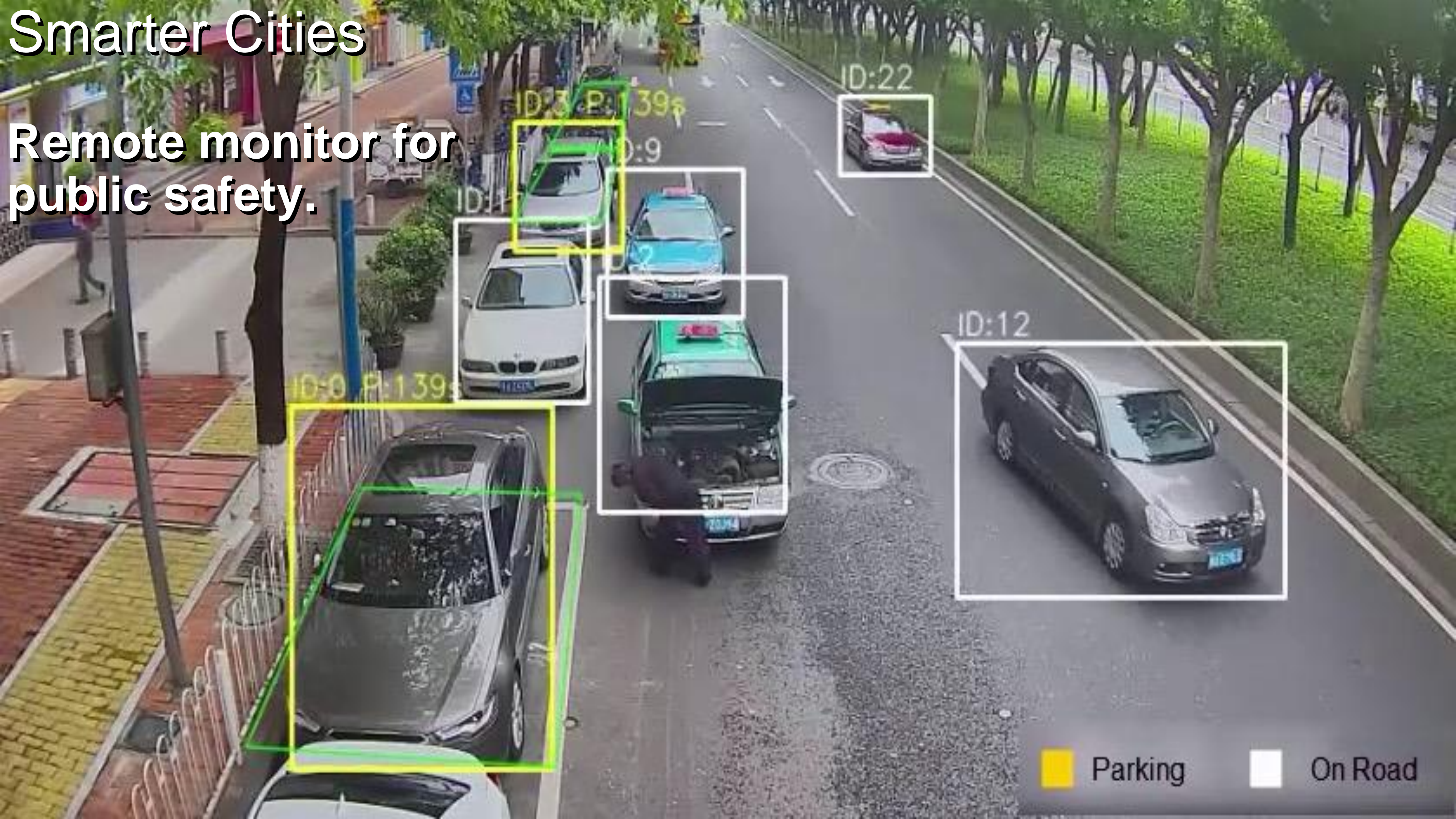Missed Diagnoses.
Radiologists are **scarce**.

Shape

Boundary

Attenuation

Technology understands morphology.
91% accuracy cancerous determination.

Holy grail? Premalignant lesions

save
time
money
lives

# 25 gigabytes
# of data per hour

is generated by a
connected car.

# 90% of cars will
# be connected by 2020.

**80 million**
wearable health
devices will
be available by
2017.

# 153 exabytes
# of healthcare
# data generated by

devices in 2013.

Increasing to **2,314**
**exabytes** in 2020.

## 2.5
## quintillion
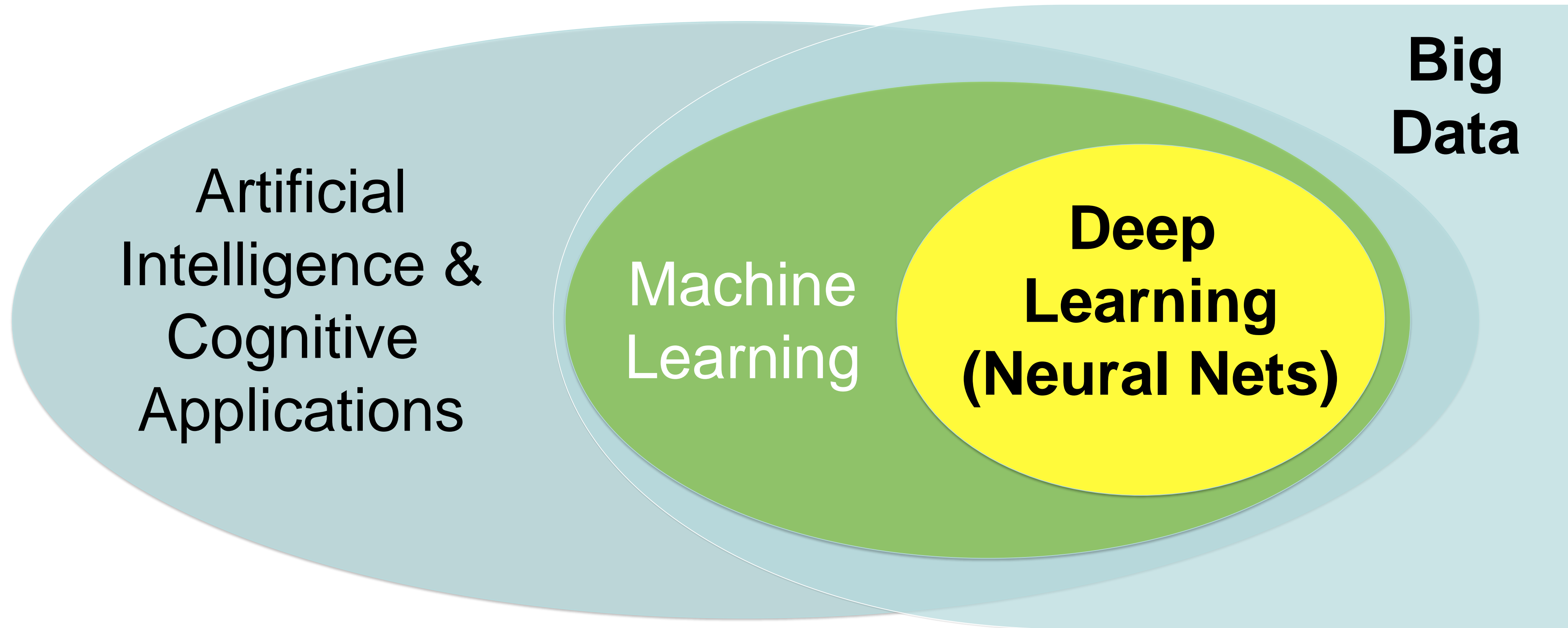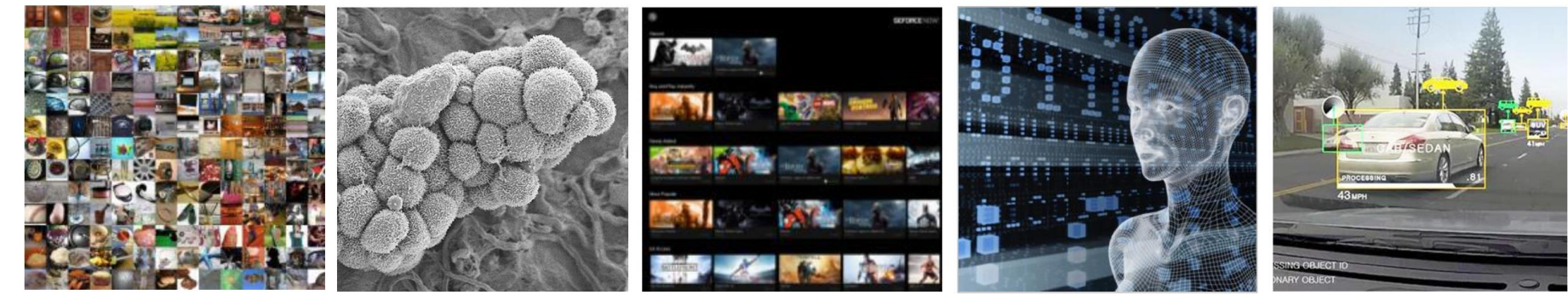## bytes of data

generated daily
by connected
machines.

There
will be
**28 times**
**more**
**sensor-**
**enabled**
**devices**
**than**
**people**
by the
year 2020.

## 1.7 megabytes
## of data per
## second

generated by
every human
being on the
planet by 2020.

# The Landscape is Evolving



**Big Data**

Artificial Intelligence & Cognitive Applications

Machine Learning

**Deep Learning (Neural Nets)**

# Gartner Views…

## Why DL/ML/AI Now?



- Three factors driving new impetus in deep learning:
  - Big data — large-scale training data
  - Algorithmic innovations
  - Highly parallel compute infrastructure
- Core computation in training DNNs: Dense matrix × vectors
  - A highly parallel workload
  - Experiment with different training models for their dataset
  - Minimize training times
- Different parts of the workflow have varying requirements

**Training in deep learning is compute and communication intensive.**

**Gartner.**

# Cognitive/AI applies across all industries

**AUTOMOTIVE**
Self-driving cars, Driver safety, Insurance

**COMMUNICATIONS**
Location-based advertising, Speech processing

**CONSUMER PACKAGED GOODS**
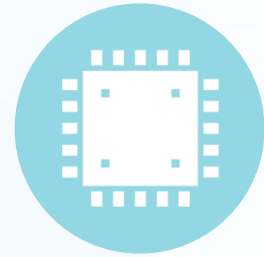Sentiment analysis of what's hot, product positioning

**FINANCIAL SERVICES**
Risk, fraud, surveillance, product opportunities

**EDUCATION & RESEARCH**
Interactive learning

**HIGH TECHNOLOGY / INDUSTRIAL MFG.**
Robotics, Mfg. quality, Warranty analysis

**LIFE SCIENCES**
Drug reactions, drug discovery

**MEDIA/ENTERTAINMENT**
Viewers / advertising effectiveness

**ON-LINE SERVICES / SOCIAL MEDIA**
Dialogue, image processing, sentiment

**HEALTH CARE**
Patient monitoring, diagnostics

**OIL & GAS**
Exploration, simulation efficiency

**RETAIL**
Consumer sentiment Demand forecasting

**TRAVEL & TRANSPORTATION**
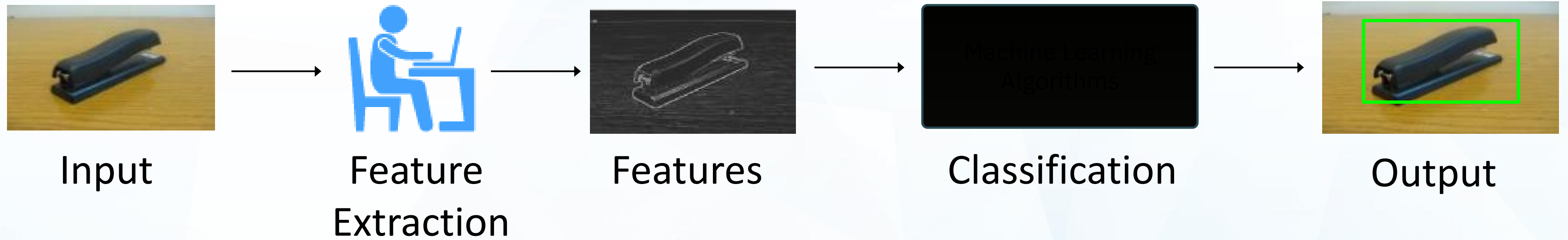Traffic and safety management

**UTILITIES**
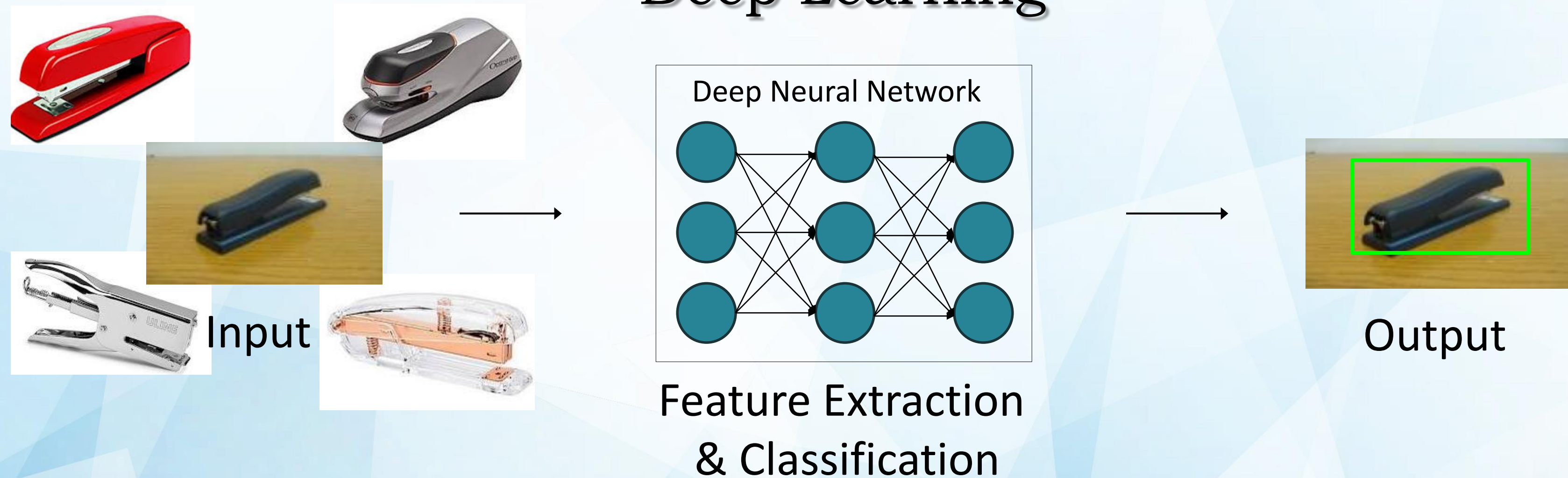Smart Meter analysis for network capacity,

**LAW ENFORCEMENT & DEFENSE**
Threat analysis - social media monitoring, photo analysis
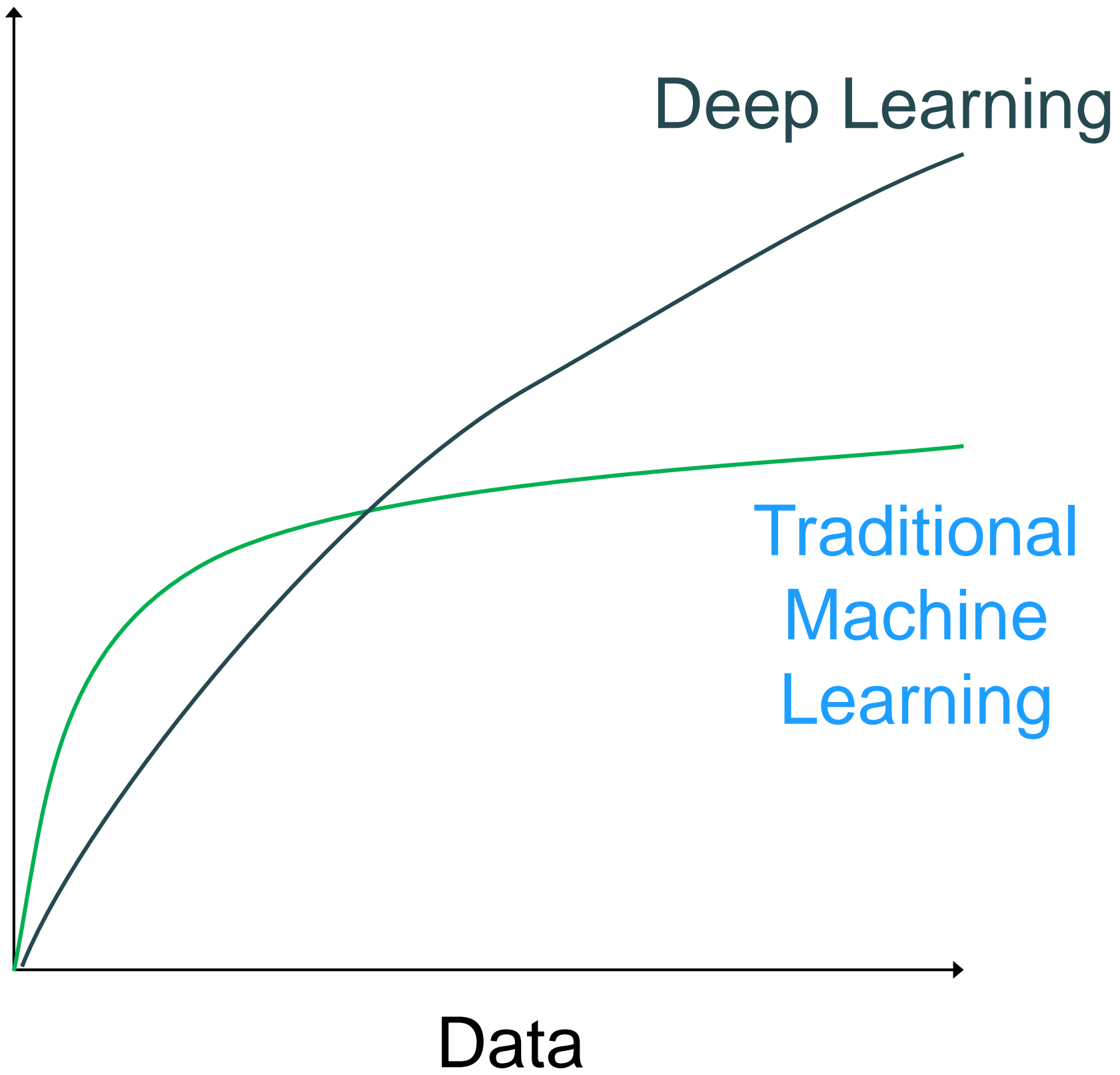
# Machine Learning



Input → Feature Extraction → Features → Classification → Output

# Deep Learning



Deep Neural Network

Input → Feature Extraction & Classification → Output

# Deep Learning Has Revolutionized Machine Learning

Accuracy



Deep Learning

Traditional Machine Learning

Data

# of Searches for Deep Learning from 2011 to 2017



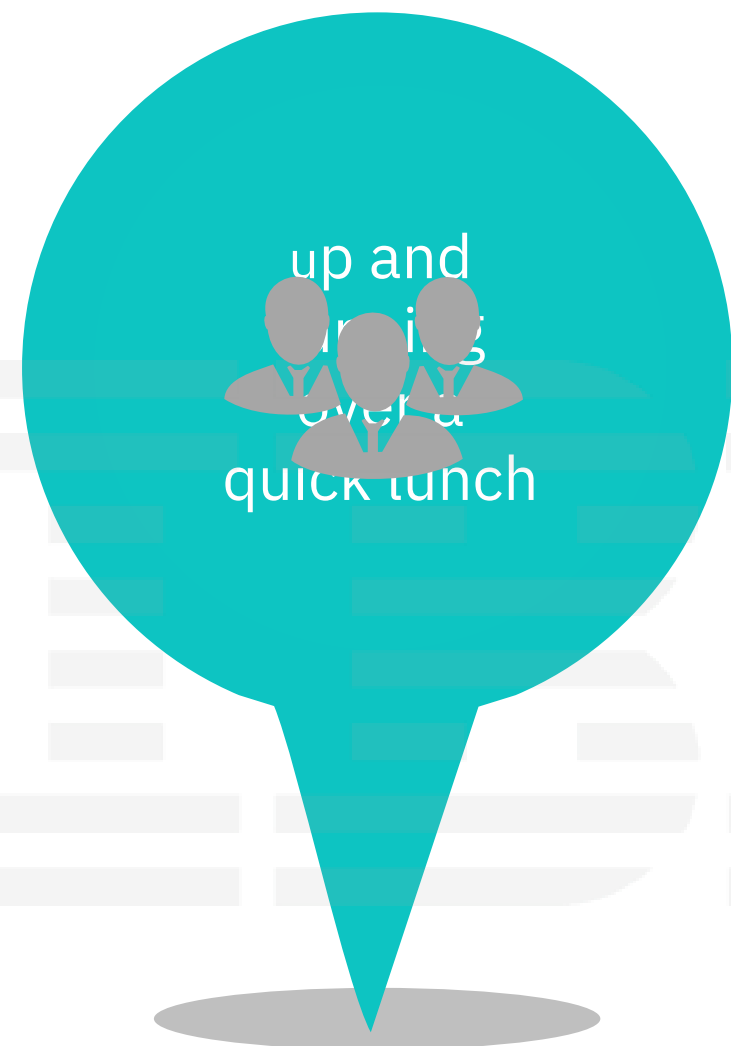Source: Google Trends.  Search term "Deep Learning"
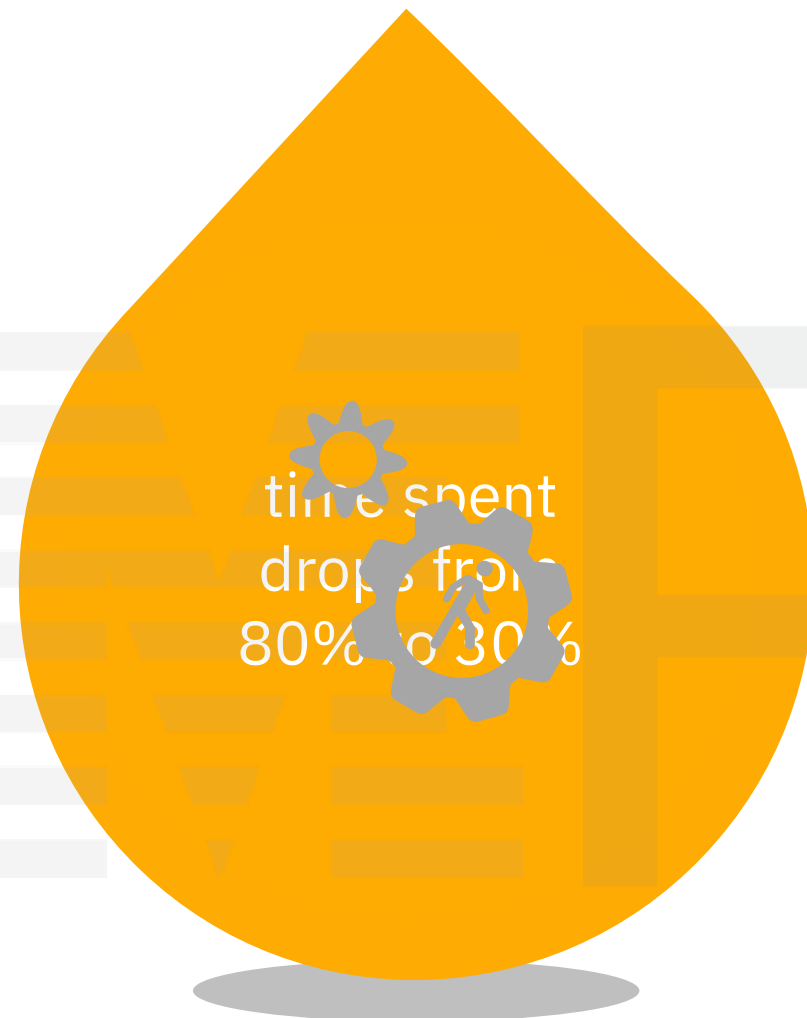
13

# More iPod than iPhone X

We are here!

deep learning timeline

# Realizing business value comes with its challenges

**DATA PREPARATION**
most time spent here

**DEPLOY & INFER**
requires different skills

up and running over a quick lunch

time spent drops from 80% to 30%

9 days work becomes 4 hours ... more models

Single click to deploy

Iterate faster and do it again

**UP & RUNNING**
weeks to months

**BUILD, TRAIN, OPTIMIZE**
very iterative

**MAINTAIN ACCURACY**
experience all that pain again

# Houston, we have a problem...

# Cognitive Systems are built with optimized hardware and software



**Dev Ecosystem**

| Industry Solutions |
| Partner Software |
| Open Source Software |
| Optimized Libraries |

P8 〉 P9 〉 P10

| Open Accelerator Interfaces |
| Accelerator Roadmaps |

## Not Just About Hardware Design

It's about co-optimized

hardware + software

which *just work* for Machine Learning, Deep Learning, and AI

IBM

# IBM Power Systems: open to the core

## OpenPOWER

**>340 members**

Mellanox
TECHNOLOGIES

XILINX

NVIDIA.

## Open Source Workloads

mongoDB.

Spark
APACHE

hadoop

EDB
POSTGRES

cassandra

## OpenCAPI

AMD

Hewlett Packard
Enterprise

Google

DELL

## Open Frameworks

torch

theano

Chainer

TensorFlow ™

Caffe

Making machine learning and AI more affordable

# *This* is what a r-evolution looks like



OpenPOWER™

## Implementation / HPC / Research

ASU Arizona State University · ASTRI · Agency for Science, Technology and Research Singapore · CAMBRIDGE CARES · Carnegie Mellon University · cfms · CINECA · GSIC Global Scientific Information and Computing Center · Hartree Centre · ICCS-NTUA · iit-b · icm

JÜLICH FORSCHUNGSZENTRUM · Lawrence Livermore National Laboratory · LSU Louisiana State University · UNIVERSITY OF MICHIGAN · NUS National University of Singapore · NANYANG TECHNOLOGICAL UNIVERSITY · OAK RIDGE National Laboratory · OSU Oregon State University · OVH.COM · rackspace the #1 managed cloud company · MAISON DE LA SIMULATION de Champagne-Ardenne

RICE Unconventional Wisdom · rzg IPP RECHEN-ZENTRUM GARCHING · SASTRA UNIVERSITY · Sandia National Laboratories · SDSC San Diego Supercomputer Center · Symbiosis Institute of Computer Studies & Research · TEES Texas A&M Engineering Experiment Station · UF UNIVERSITY of FLORIDA · TACC

USC UNIVERSITY OF SOUTHERN CALIFORNIA · 清华大学 Tsinghua University · UNIVERSITY OF ARKANSAS · UNICAMP · UNIVERSITY OF OREGON · ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΑΤΡΩΝ · WASEDA UNIVERSITY

## Software

American Megatrends · BYOSOFT 百敖软件 · FIXSTARS · SuSE · redhat · freeBSD · Google · SANMARCO INFORMATICA · gpudb · 红旗 Linux · RedHadoop · redislabs · synerscope · Groupe T2i innovation informatique · ubuntu Supported by Canonical

## System / Integration

ARROW ELECTRONICS, INC. · ASETEK · AVNET · BULL an atos company · CIARA Speed | Performance | Passion for Innovation · Cirrascale · CSPi · DRC A Security First Corp Company · E4 COMPUTER ENGINEERING · IBM · INSPUR 浪潮 · Mark III Systems

Microway Technology you can count on since 1982. · NEC · Neu Cloud Oriental 新云东方 · OCF · One Stop Systems · PENGUIN COMPUTING · RAPTOR ENGINEERING · rikor · RTDS Technologies · STACK VELOCITY · ТЕХНОПРОМ · UNISOURCE · YADRO · 中太数据 ZOOM NETWORKS

## I/O / Storage / Acceleration

ALGO-LOGIC · ALPHA DATA · AVAGO TECHNOLOGIES · BittWare FPGA COMPUTING SYSTEMS · blazegraph · BLUEBEE High Performance Genomics · BROCADE · Chelsio Communications Accelerate · CONVEY COMPUTER · DataDirect NETWORKS · edico genome · EVERSPIN TECHNOLOGIES The MRAM Company™ · FUSION-io

HGST a Western Digital company · HITACHI Inspire the Next · Inphi Think fast. · Interface Masters TECHNOLOGIES Innovative Network Solutions · MAXELER Technologies · Mellanox TECHNOLOGIES · MEMBLAZE · Micron · Microsemi · MYRICOM Network Products · Nallatech · NVIDIA

PMC · QLOGIC · SAMSUNG · SanDisk · Semptian 恒扬科技 · SK hynix · SOLARFLARE · XILINX

## Boards / Systems

Celestica · 创和通讯 · IBM · Inventec · msi · PET 柏飞电子 · TYAN · wistron

## Chip / SOC

IBM · IDT · Infineon · POWERCORE · SYNAPSE design · VeriSilicon

# Cross community collaboration is essential

| OPENSTACK | LINUX | OPEN COMPUTE | OPENPOWER |
|-----------|-------|--------------|-----------|

**OPENSTACK**
- Orchestration
- Identity
- Compute/Network/Storage Fabric

**LINUX**
- Applications
- Kernel
- Drivers

**OPEN COMPUTE**
- Mechanical
- Electrical
- Firmware
- Protocols
- Interconnects
- Memory/CPUs

# FROM HERE TO AI

Our POWER9 servers and solutions are built to crush today's most advanced data applications – from the mission critical applications you run today to the next generation of AI workloads.

| Mission Critical Workloads | Big Data Workloads | Enterprise AI Workloads |
|---|---|---|

S922/S914/S924,
H922/H924,
L922

E950/H950

E980/H980

LC922/LC921

AC922

Core Infrastructure

Next Gen AI
Workloads

enterprise-ready
software distribution
built on open source

performance
faster training times
for data scientists

tools for ease
of development

IBM PowerAI

# Enterprise Ready Build on Open Source

## IBM PowerAI Platform

### PowerAI Software Distribution

**Deep Learning Frameworks**

| Caffe | NVIDIA Caffe | IBM Caffe | torch |
|---|---|---|---|
| TensorFlow | theano | Chainer | |

**Supporting Libraries**

| DIGITS | OpenBLAS | Distributed Frameworks | Bazel | NCCL |
|---|---|---|---|---|

### IBM Power System for HPC, with NVLink

Breakthrough performance for GPU accelerated applications, including Deep Learning and Machine Learning.

**Tools for Ease of Development**

**rich advisory and building toolsets to flatten time to value**

AI Vision
rich toolset image
recognition neural
networks

automated deep learning
toolkit data preparation

DL Insight toolkit supports
auto-training runs for
hyper parameter tuning
+++

# Performance… Faster Training and Inferencing

# faster training times for data scientists

## Distributed Deep Learning



## Traditional Model Support → Large Model Support (LMS)

(Competitors)

Limited memory on GPU forces trade-off in model size / data resolution

(PowerAI)

Use system memory and GPU to support more complex models and higher resolution data

# COMMUNICATION PATHS



**Power AI DDL:** Fully utilize bandwidth for links within each node and across all nodes
→ Learners communicate as efficiently as possible

# Deep Learning Impact: Monitor, Adviser and Optimizer

**Deep Learning Applications**

- Image recognition
- Object detection
- Translation
- Others

**Optimize** →

**Hyper-parameters Optimizer**

- optimizing procedure parameters
- criterion and model parameters

app-20160824141112-1310

Cluster Master          172.17.0.11
Start Time              2016-08-24T14:11:15.175Z

Static parameters ▾

| Hyper Parameters | | Layer information | |
|---|---|---|---|
| weight_bin_range | 7 | L1 | data |
| debug_info | true | L2 | data |
| weight_scale | 20 | L3 | conv1 |
| snapshot_format | HDF5 | L4 | pool1 |
| snapshot_prefix | "cifar10_quick" | L5 | ip1 |
| max_iter | 2000 | L6 | ip2 |
| gradient_scale | 50 | L7 | accuracy |
| display | 10 | L8 | loss |
| base_lr | 0.0001 | | |
| snapshot | 4000 | | |
| solver_mode | GPU | | |
| test_iter | 100 | | |
| test_interval | 50 | | |
| net | "cifar10_quick_train_test.prototxt" | | |
| weight_decay | 0.004 | | |
| lr_policy | "fixed" | | |
| momentum | 0.9 | | |

**Monitor** →

**Real-time monitor for running application**

- learning curve
- weight/gradient/activation histogram and norm
- worst cases of training samples

**Advice** →

**Real-time adviser for running application**

check the training process

- overflow
- underfitting
- overfitting
- divergence
- convergence

Problems Analysis and Suggestion

early stop: exception found
underfitting detected
add more hidden layers
sync_interval=20

Learning Curves

30

# IBM PowerAI Platform

## PowerAI Software Distribution: Optimized for Power

| Deep Learning Frameworks & Enhancements | TensorFlow | Caffe | IBM Caffe | Watson APIs |
| --- | --- | --- | --- | --- |
| | IBM Research Distributed Deep Learning | Power Systems Large Model Support | AI Vision Tools | |
| Supporting Capabilities And Libraries | Distributed Frameworks | IBM Research AI Vision Runtime | IBM Spectrum Conductor | IBM Data Science Experience |
| | NVIDIA DIGITS | OpenBLAS | Bazel | NVIDIA NCCL |
| IBM Services And Support | IBM Entire Stack Support | IBM Research Pioneering AI Research | COGNITIVE CLASS Education & Certification | Power Systems Optimization and testing |

## IBM Power Accelerated Servers: Ideal for PowerAI

| IBM Services And Support | Acceleration Superhighway | Designed for The AI era | Enterprise Grade | POWER9 Performance |
| --- | --- | --- | --- | --- |

# Smarter & Safer Cities

Near miss at intersections

Monitor and Impose regulations

# More Use Cases

# Programmatic Approach: Using Jupyter Notebooks

# IBM AI & Data Science Workbench

Democratizing ML/DL

Operationalizing ML/DL

Fit-for-purpose

Integration with Watson functions



**IBM Data Science** Experience

Model Builder    SPSS Flows    jupyter    R Studio

Machine Learning Runtimes    Deep Learning Runtimes

APACHE Spark    scikit learn    Spark MLlib    dmlc XGBoost    IBMPowerAI    ANACONDA    H2O.ai

Cloud Infrastructure as a Service

docker    kubernetes    nVIDIA

Model Lifecycle Management

Cognitive Systems    Data Stores

No idea of required load, or wildly elastic?
PUBLIC CLOUD

Staying in NYC for a weekend?
PUBLIC CLOUD

Staying in NYC 6 months?
PRIVATE or PUBLIC?

Staying in NYC 3 years?
PRIVATE CLOUD

**Elastic Workloads**
Spin up and down resources on the public cloud

25%

75%

**Predictable Workloads**
Lower costs with private cloud infrastructure

balance **owning** and **renting** for today's enterprise workloads

# What's next ?

- **Getting PowerAI ?**   **http://ibm.biz/powerai**


- **Testing PowerAI in the Cloud ?**
  **https://power.jarvice.com/**


- **Need for any further reference or contact ?**
  **https://www-03.ibm.com/systems/uk/power/hardware/hpc/outthink.html**

# So, What's really next ?

# Questions ?

**Alain Cyr**

*cyralain@fr.ibm.com*

IBM

*Power & Storage Business Development*
*MOP IBM Client Center*