



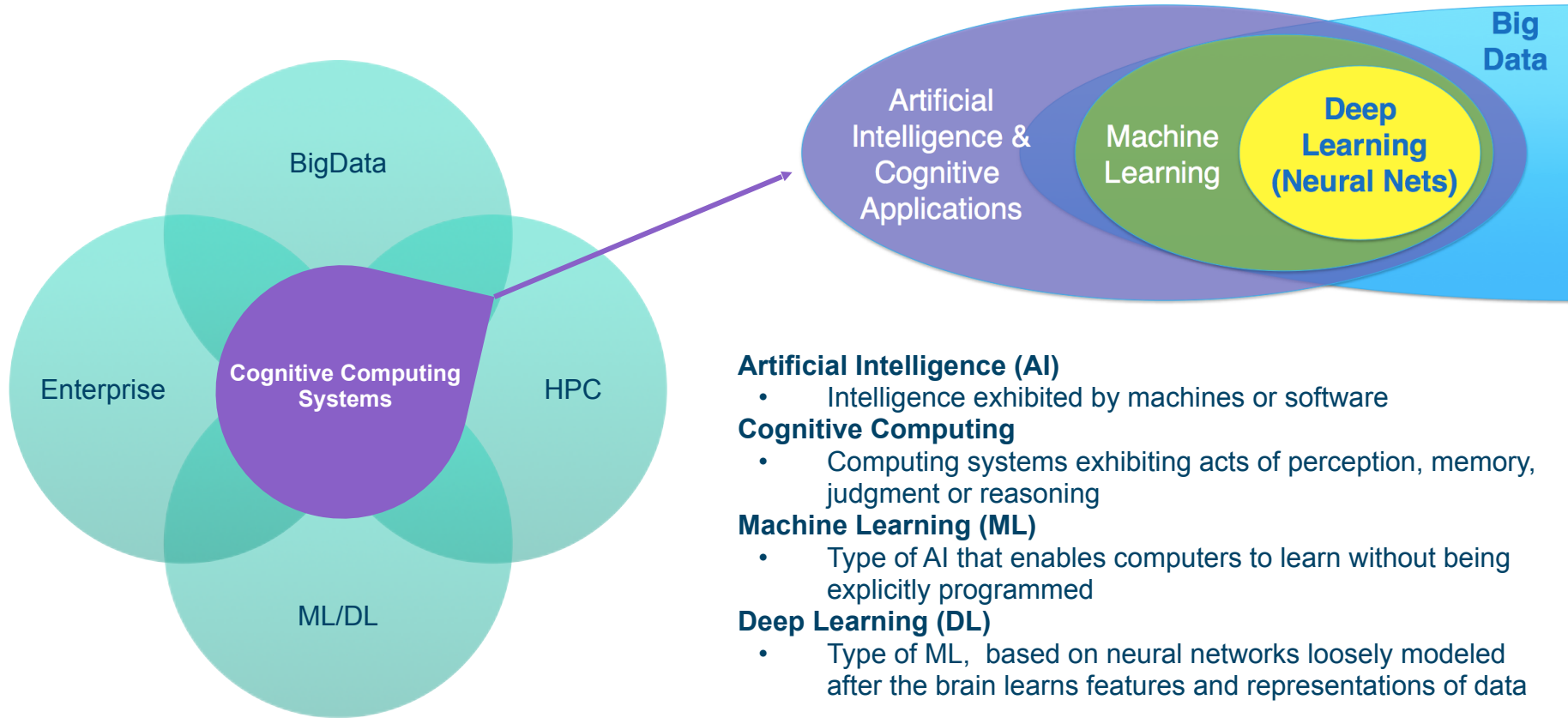
flattening the Deep Learning time to value curve

Franz Bourlet
POWER Systems Technical Sales
IBM Belgium & Luxembourg
Franz_Bourlet@be.ibm.com

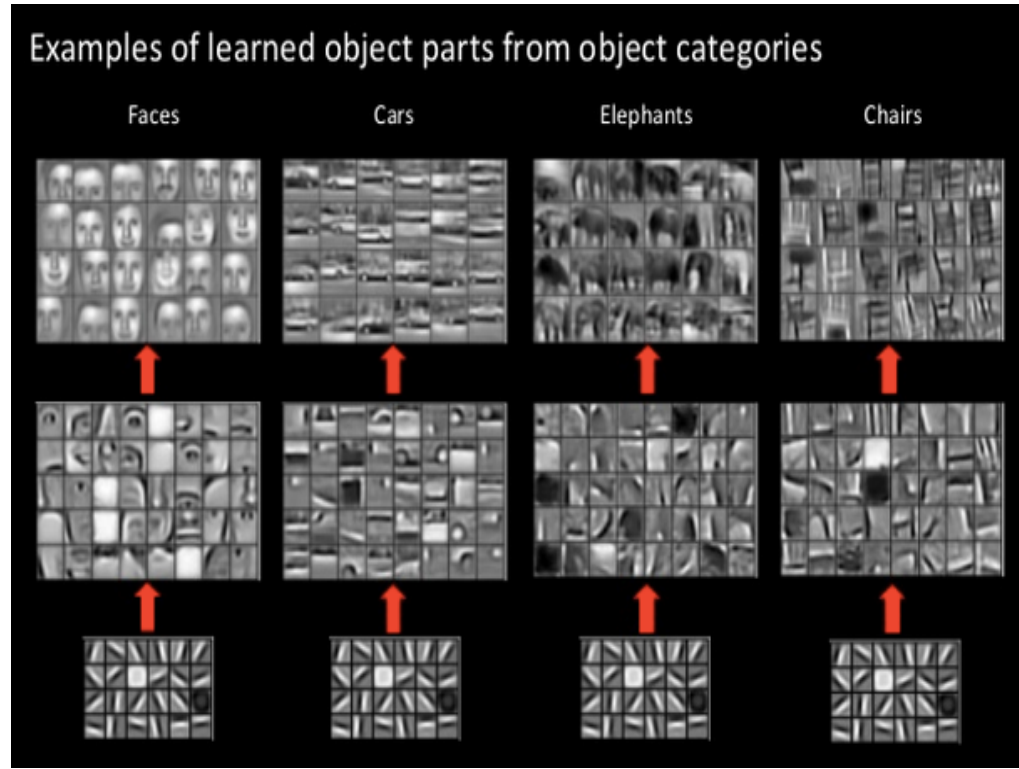
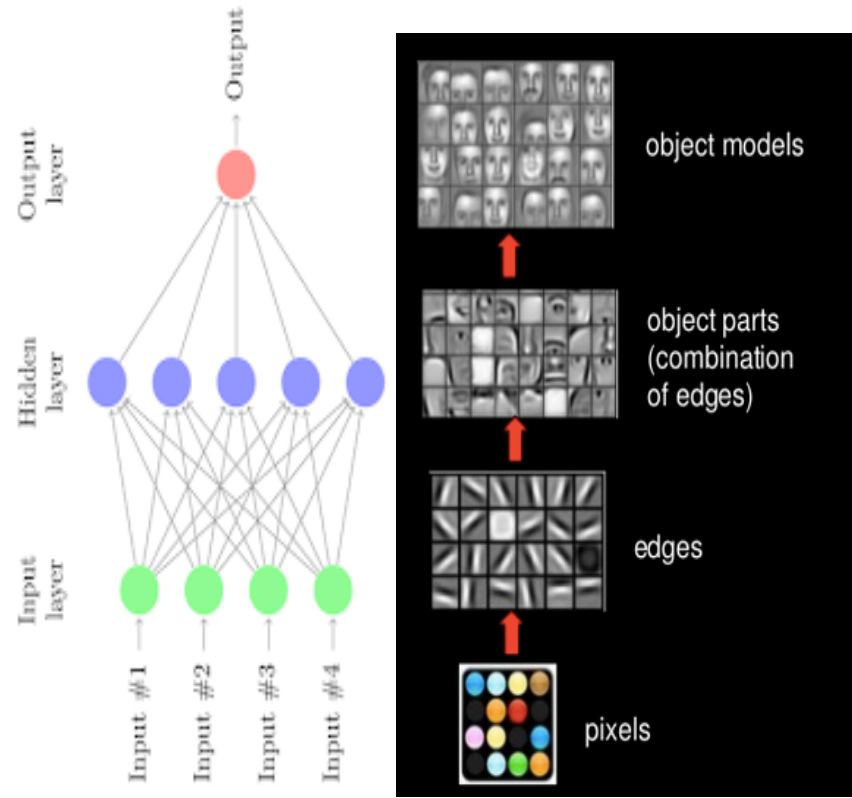


Agenda

- **Introduction to machine learning / deep learning**
- **Watson vs PowerAI**
- **POWER for cognitive solutions**
- **Cognitive solutions on POWER**
- **PowerAI demo**
- **Use cases / references**
- **AI in Belux**



Deep Learning....Under the Hood



What is Watson?

- **It is a brand**
- **Multiple products in Watson family**
 - A diverse set of cognitive technologies covering AI, machine learning, and deep learning
 - Evolve from the Jeopardy Game and cancer research applications
 - Organic and acquired technologies
- **Cloud only**

The Watson Data Platform Portfolio

Legend:
Available Now
Generally Available in December
Generally Available in 1H 2018

Data Science Experience

Collaborative data science capabilities

Data Refinery

Data preparation & integration

Lite Plan available **today**. Standard Plan GAs in 1Q18

Dashboards

Dashboarding & reporting

Analytics Engine
Hadoop & Spark

Watson Machine Learning
ML processing & management

Streaming Analytics
Analytics for data-in-motion

Message Hub
Kafka

Lift
Ground-to-cloud transfer

IBM Cloud Object Storage

Cloudbant

Compose
(Mongo, Graph, Redis, Postgres, Elastic Search, ScyllaDB, MySQL, RabbitMQ, etcd, RethinkDB)

DB2 Warehouse on Cloud

DB2 on Cloud

IBM Cloud SQL Query

Data Catalog**
Data & analytical asset organization with governance

** Data Catalog PNs in SQO **now**

Watson & Cloud Overall Architecture

The IBM Watson and Cloud Platform

Application

Watson
OncologyWatson
Cyber
Security

Weather

Watson
Explore +
DiscoverWatson
Virtual
AgentWatson
Compare
+ ComplyGBS +
GTS
Industry
Solutions

+ more...



Industry Solutions –
BP/ISV/Client play for
PowerAI

AI

Conversation

Discovery

Compare
+ ComplyKnowledge
Query

Tone Analysis

Personality
InsightsVisual
Recognition

Speech

Document
ConversionNat. Language
UnderstandingNat. Language
Classifier

+ more...



Cognitive Services
(comparable to what can
be built from open source
frameworks that are part
of Power AI) – API delivery
model

Data

Crawl



Cleanse



Enrich



Store



Analyze



The Watson Data Platform

Data Enrichment &
Storage services –
complimentary to PowerAI

Cloud

Dev Services

Containers

Messaging

Blockchain

Logging

+ more...

Infrastructure

Storage

Compute

Physical Network

Infrastructure
Mgmt

+ more...



Infrastructure – Only
cloud but 2 flavors –
Public - Multi-
tenant, share all
Dedicated – VPC-
driven isolation

Watson Data Platform : an integrated, unified self-service experience

- Shop for data
- Manage policies
- Shape data
- Build dashboards
- Auto model building
- Build ML flows
- Auto-optimize models
- Develop notebooks
- Streaming pipelines
- Build data apps

The screenshot displays the IBM Watson Data Platform interface. The top navigation bar includes 'Projects', 'Tools', 'Catalog', 'Data Services', and 'Community'. The 'Tools' and 'Catalog' tabs are highlighted with a green box. Below the navigation bar, the breadcrumb trail shows 'My Projects > Recent Sales Performance Dr...'. The main content area is divided into three sections: 'Notebooks', 'Streams flows', and 'Models'. Each section has a '+ New' button highlighted with a green box. A large blue semi-transparent box with white text is overlaid on the 'Notebooks' section, stating 'Analyze data and build data products from within a single environment'.

Notebooks

NAME	SHARED	SCHEDULED	STATUS	LANGUAGE	LAST EDITOR	LAST MODIFIED	ACTIONS
Select GPS Events				Python 2.7	paul taylor	2 Nov 2017	
Retail Sales Analysis				Python 2.7	Dirk deRoos	2 Nov 2017	
Machine Learning R				R 3.3.2	John Emmert	6 Nov 2017	

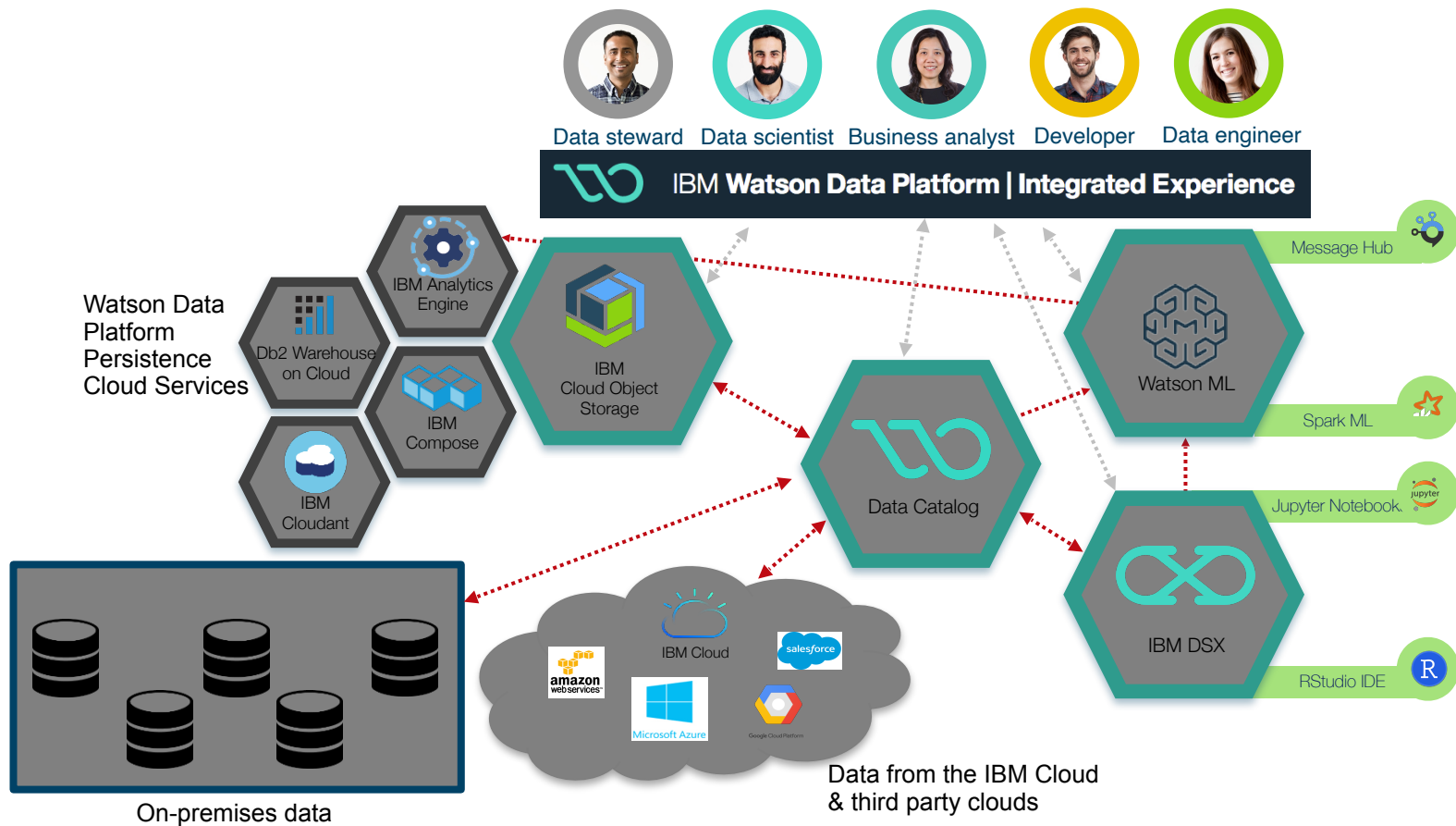
Streams flows

NAME	MODIFIED BY	LAST MODIFIED	ACTIONS
Streaming Tweets	jmemmert@us.ibm.com	8 Nov 2017, 11:48:06 am	

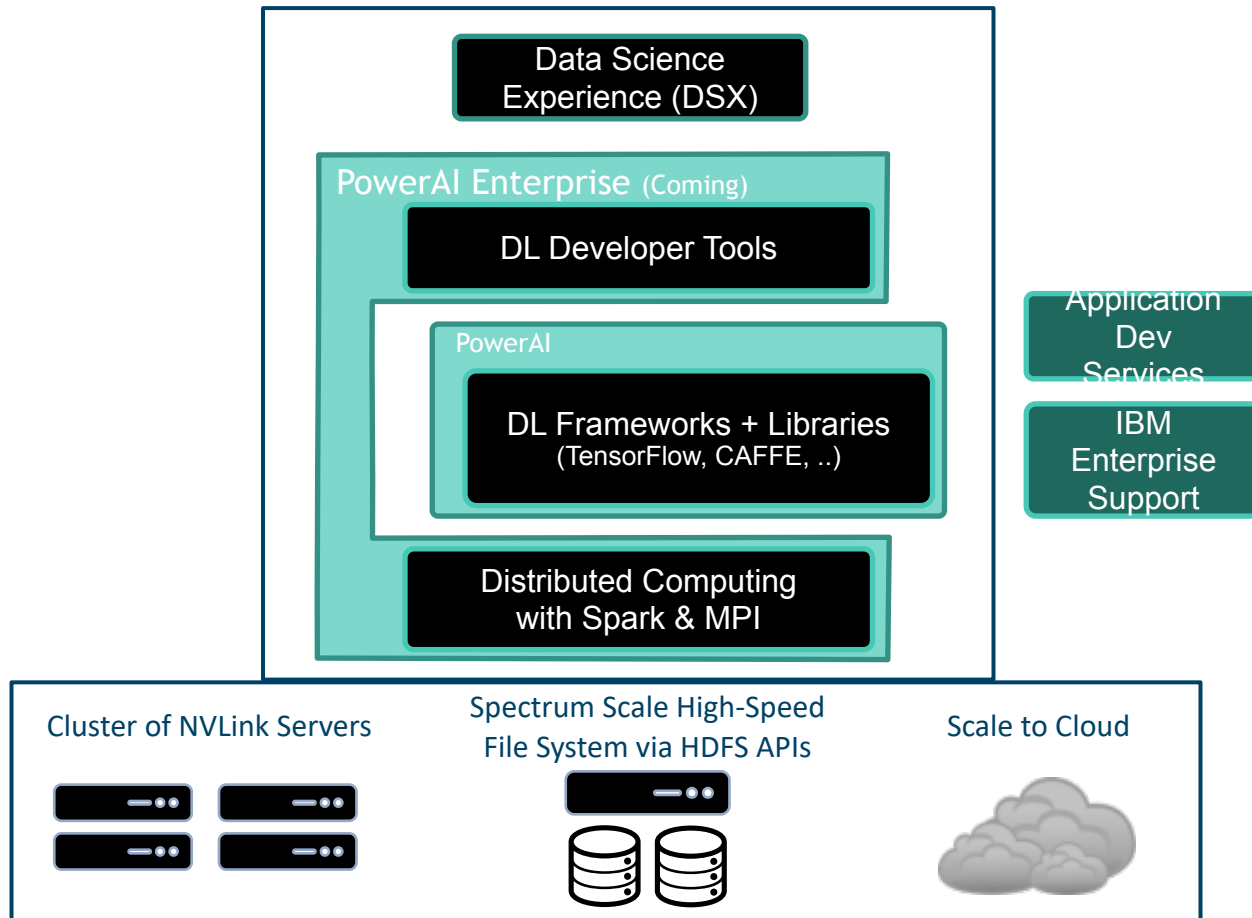
Models

NAME	STATUS	RUNTIME	LAST MODIFIED	ACTIONS
Repellent model	trained	spark-2.0	30 Oct 2017	

Example Deployment Architecture: Data Science



PowerAI Enterprise: Enhancing Developer Experience (& DSX integration)



Positioning --- Infrastructure Location

Watson AI

Cloud

Mainly looking at cloud deployment, minimal data privacy issues, i.e., data can be moved off-premise

Training, validation data can be hosted in a cloud environment

Does not envisage an issue with scaling on the cloud

Use case requires minimal data movement, i.e, for model re-training or the model used is well established

Use cases : chatbots, standard image classification, product recommendation, etc.

PowerAI

On-prem

Mainly looking at on-prem deployment, strong requirement that data does not move offsite

Training, validation data needs to be kept in-house due to regulatory, competitive or other reasons

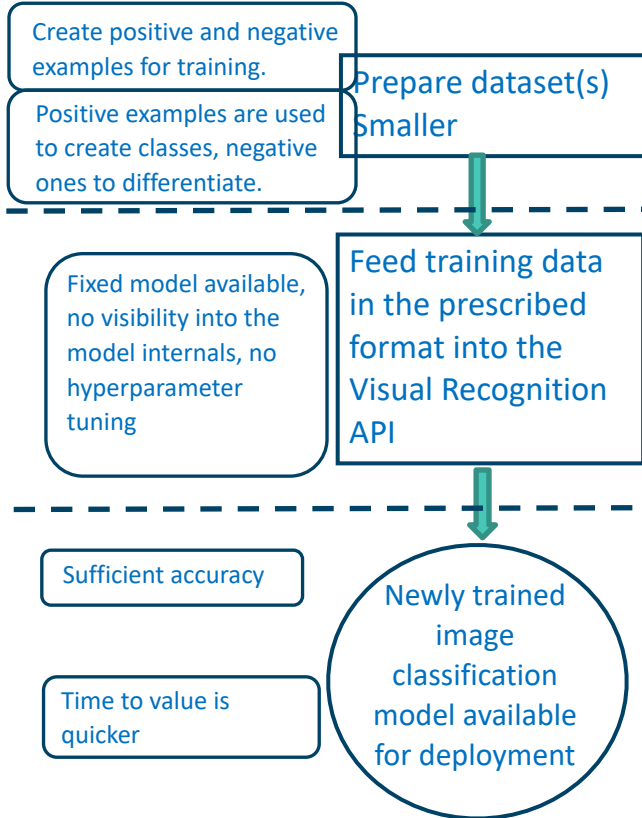
Multiple training runs need to be carried out on the model as focus on training is higher

Large amount of data for training and training data is updated frequently

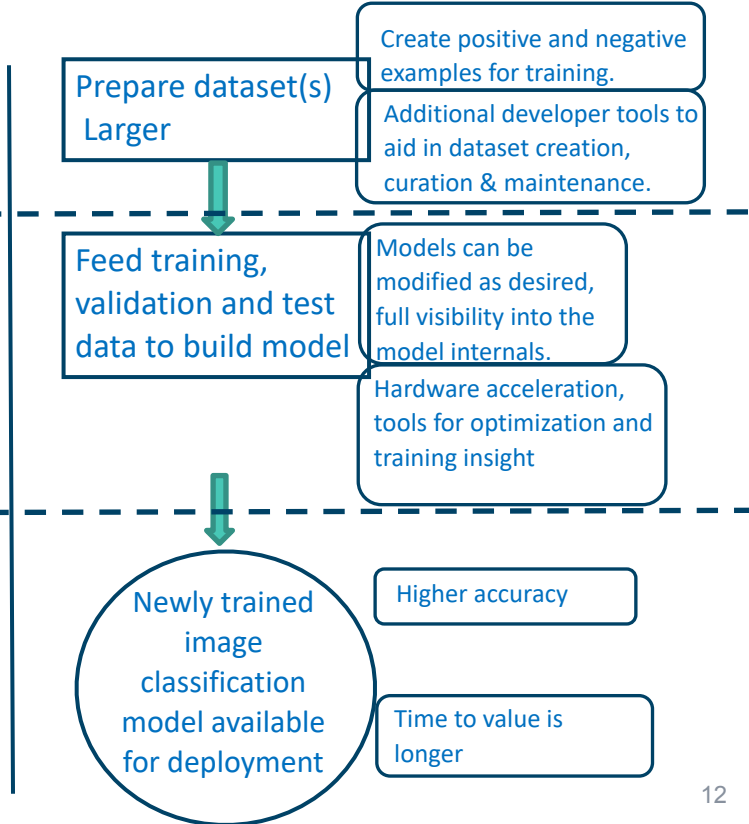
Use cases : fraud detection, credit analysis, image analysis for security, etc.

Positioning --- Model Training

Watson-based flow to create a custom image classification model – Focused on solution building



PowerAI-based flow to create a custom image classification model – build the classifier from ground up



Positioning --- Data

Watson AI

Unstructured Data

Unstructured data - systems of engagement

Text, audio, images, video

Use cases generally center around add on to business functions: chatbots, translation, visual recognition

Use data about what people say

PowerAI

Unstructured and Structured Data

Unstructured data - systems of engagement
Structured data - systems of records

Text, audio, images, video
Transactional data, warehouse data

Besides Watson AI use cases, it is also used to improve core business functions: churn reduction, fraud detection, product promotion

Use data about what people do

Positioning --- Targeted Users

Watson AI

Application Developers

Easier to use. Hide data science complexity.

Easier to retrain

Build applications quickly

PowerAI

Data Scientists

Full data science capabilities

More optimized results with retraining

Slice and dice data to get optimal results
Longer training effort

Positioning --- Infrastructure Management

Watson AI

PowerAI

Cloud

Is not keen on visibility into the infrastructure & its management (compute, storage, etc.)

Client is not invested in building a data center and does not plan to

Does not envisage an issue with scaling on the cloud

View the infrastructure management piece as a significant barrier to adoption

Is content with the performance provided by standard configurations (including accelerators) available in the cloud

No push from externalities like regulations, privacy or data as competitive advantage to make the move to an on-prem solution

On-prem

Is very keen to have complete visibility into infrastructure & its management

Already has a data center or is in the process of building one

Considering the impact of scaling on the cloud

Infrastructure management is not a significant barrier to adoption

Is looking to understand the performance benefit provided by non-standard configurations with on-prem infrastructure

Does view regulations, data privacy and/or data as competitive advantage as these issues drive a move to an on-prem solution

THE AI ERA IS HERE.



We continue to experience exponential growth of data and data sources.



Computing has moved beyond 'post CPU only' era, giving us vast computational power that was not accessible before.



CIOs are evolving from 'chief information officer' to 'chief intelligence officer' and the data science organization has continued to gain power and influence.

Right now, your infrastructure is putting up

ROADBLOCKS



Not equipped for
enterprise-level data
volumes



Blocks to
acceleration



Servers not specifically
designed for AI
workloads



Not able to
easily scale

AI DEMANDS A DIFFERENT TYPE OF SYSTEM

IBM Power Systems provides the cutting-edge advances in AI that data scientists demand, and the critical reliability that IT needs.



IBM Power Systems AC922



*Architecture designed
for the AI era with
advanced IO interfaces*

*AI at unrivaled scale with what
will likely become the worlds
most powerful supercomputer*

*Extends on heritage of
performance leadership
across AI, HPC and
accelerated DBs*

Best Server for Enterprise AI

AC922



An Acceleration Superhighway

Unleash accelerated computing potential in the post CPU-only era



Designed for the AI Era

Architected for the modern analytics and AI workloads that fuel insights



Delivering Enterprise-Class AI

Cutting-edge AI innovation data scientists desire, with dependability IT requires



An Acceleration Superhighway

Unleash accelerated computing potential in the post CPU-only era

Fastest Accelerator Performance

NVIDIA 2nd generation NVLink for Power9 and PCIe Gen4 deliver performance advantage vs. PCIe gen3

Simplest path to Acceleration

Coherence simplifies path to acceleration by abstracting data movement and locality

Most Efficient use of Accelerators

Share resources across CPUs and GPUs, while reducing bottlenecks to more fully-utilize acceleration investments

Developed via Collaborative Effort

Benefits of Industry collaboration are etched in silicon via NVLink and OpenCAPI

Designed for the AI Era

Exceptional data-intensive architecture for the types of AI and Analytics workloads that are fueling the AI Era

Exceptional Next-gen Design

Architecture ahead of x86 servers and better-equipped to handle today's data-intensive HPC and analytics workloads

More Data. Faster Data.

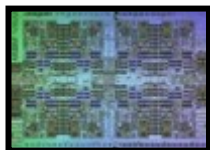
Up to 5.6x more I/O throughput and memory bandwidth, 4x threads*. 4th generation PCIe and 2nd generation NVLink, 2nd generation of CAPI

AI at Unrivaed Scale

AC922 is the backbone of CORAL Summit, meeting milestones to deliver 200+PetaFlops of HPC and 3 ExaFlops of AI as a service performance.

* 5.6x I/O bandwidth claim based on NVIDIA measurement test conducted on a Xeon E5-2640 V4 +P100 vs Power9 + V100 (12 GB/s vs 68 GB/s rated)

An Acceleration Superhighway: POWER 9 is IBM's Latest Processor

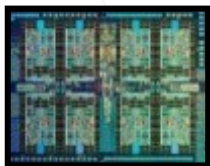


POWER7
45 nm

Enterprise

- 8 Cores
- SMT4
- eDRAM L3 Cache

1H10

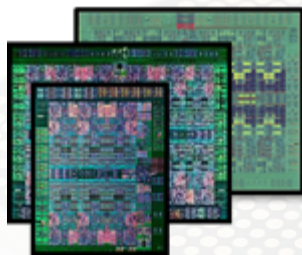


POWER7+
32 nm

Enterprise

- 2.5x Larger L3 cache
- On-die acceleration
- Zero-power core idle state

2H12



POWER8
Family 22nm

Enterprise & Big Data Optimized

- Up to 12 Cores
- SMT8
- CAPI Acceleration
- High Bandwidth GPU Attach

1H14 – 2H16



POWER9 Family
14nm

Built for the Cognitive Era

- Only processor with NVLink, PCIe Gen 4 advanced IO interfaces and coherence
- Premier Platform for Accelerated Computing
- Processor Family with Scale-Up and Scale-Out Optimized Silicon

2H17 – 2H18+

POWER9 Processor Family

Core Count / Size

SMP scalability / Memory subsystem

Scale-Out – 2 Socket Optimized

Robust 2 socket SMP system

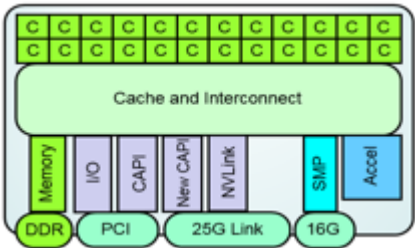
Direct Memory Attach

- Up to 8 DDR4 ports
- Up to 170 GB/s memory BW
- Commodity packaging form factor

SMT4 Core

24 SMT4 Cores / Chip

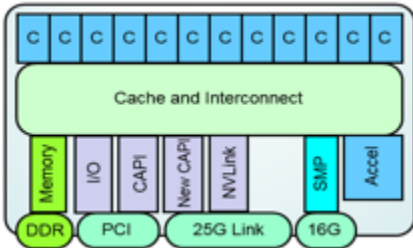
Linux Ecosystem Optimized



SMT8 Core

12 SMT8 Cores / Chip

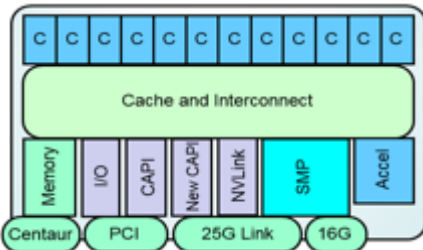
PowerVM Ecosystem Continuity



Scale-Up – 4+-Socket Optimized

Scalable System Topology / Capacity

- Large multi-socket
- Buffered Memory Attach
- 8 Buffered channels
- Up to 230 GB/s memory BW



POWER9

An acceleration superhighway.

The only processor specifically designed for the AI era.

4x

Threads per
core vs x86

9.5x

Up to 9.5x more I/O
bandwidth than x86

2.6x

More RAM
possible vs. x86

1st

CPU to deliver
PCIe gen 4

An Acceleration Superhighway: POWER9 offers a variety of Acceleration Options

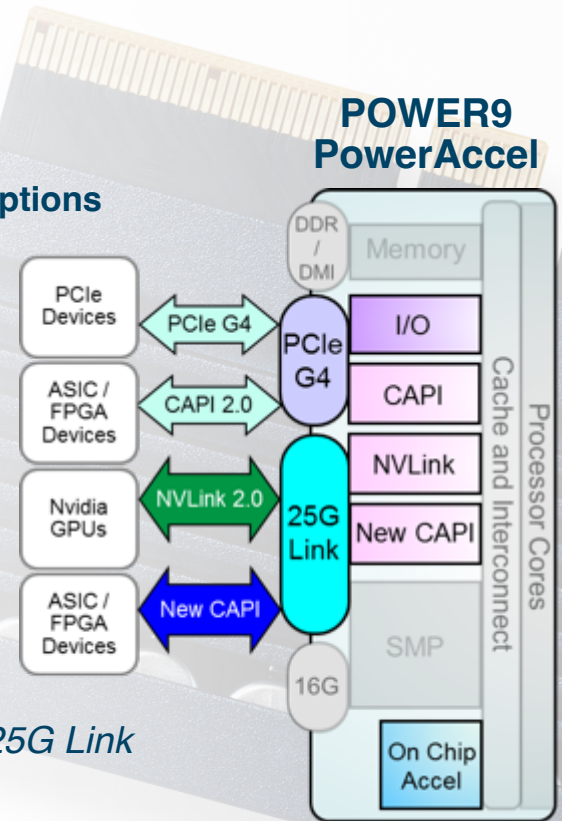
- Extreme Processor / Accelerator Bandwidth and Reduced Latency
- Coherent Memory and Virtual Addressing Capability for all Accelerators
- OpenPOWER Community Enablement – Robust Accelerated Compute Options

State of the Art I/O and Acceleration Attachment Signaling

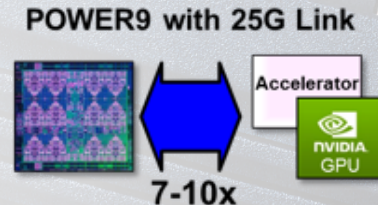
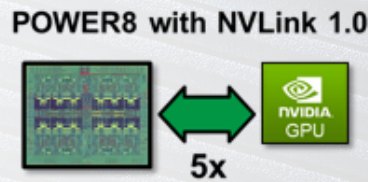
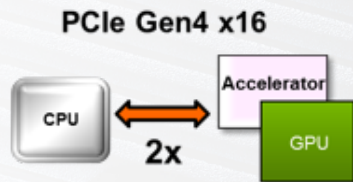
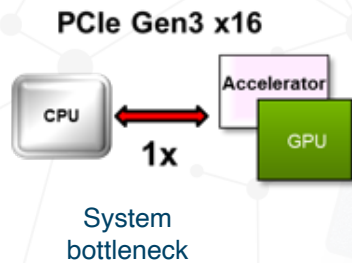
- **PCIe Gen 4** x 48 lanes – 192 GB/s duplex bandwidth
- **25G Link** x 48 lanes – 300 GB/s duplex bandwidth

Robust Accelerated Compute Options with OPEN standards

- **On-Chip Acceleration** – Gzip x1, 842 Compression x2, AES/SHA x2
- **CAPI 2.0** – 4x bandwidth of POWER8 using *PCIe Gen 4*
- **OpenCAPI 3.0** – High bandwidth, low latency and open interface using *25G Link*
- **NVLink 2.0** – Next generation of GPU/CPU bandwidth and integration



An Acceleration Superhighway: POWER9 Introduces Acceleration Innovations



Extreme CPU/Accelerator Bandwidth

Seamless CPU/Accelerator Interaction

- Coherent memory sharing
- Enhanced virtual address translation

Broader Use of Heterogeneous Compute

- Designed for efficient programming models
- Accelerate complex analytic / cognitive applications

Designed for Great Supercomputing & AI Leaders

Unprecedented performance and application gains with **advanced IO interfaces** integrated into **the NEW P9 processor** delivering **capabilities** not available on x86

Advanced IO interfaces:

- **2nd Generation CPU - GPU NVLink:** ~5.6X the CPU-GPU bandwidth compared to x86
- **PCIe Gen4/CAPI 2.0:** First to market with PCIe Gen 4 with 2x improvement over PCIe Gen 3 and next gen CAPI for coherent device attachment

Introducing Coherence: Treat system memory just like GPU memory enabling game changing simplification of programming and larger model sizes

Burst Buffer: Start/finish job data staging with high performance storage adapter resulting in significant improvement of computing efficiency

GPU Direct: Peer to peer communication within a cluster for remote data access and transfer (RDMA)



Configuration System Details for 4Q GA

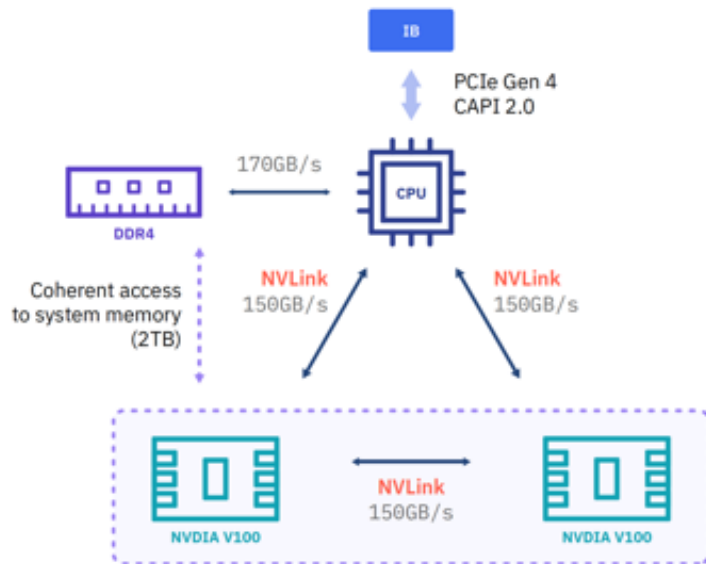
- MTM: 8335-GTG - Specific 4Q Feature Availability:
 - 1.6 TB NVMe high performance storage adapter
 - 2 SFF SATA: HDD (Max 4TB); SSD (Max 3.84TB)
 - IO: EDR InfiniBand, Quad ENET (2x1/2x10 GB), Quad ENET (4x1 GB), 100 GB ENET
 - RHEL 7.4 for P9
 - Air cooled only version available (max 4 GPU's)
 - OpenBMC

* Results are based on IBM Internal Measurements running the CUDA H2D Bandwidth Test

* Hardware: Power AC922; 32 cores (2 x 16c chips), POWER9 with NVLink 2.0; 2.25 GHz; 1024 GB memory, 4xTesla V100 GPU; Ubuntu 16.04. S822LC for HPC; 20 cores (2 x 10c chips), POWER8 with NVLink; 2.86 GHz; 512 GB memory, Tesla P100 GPU

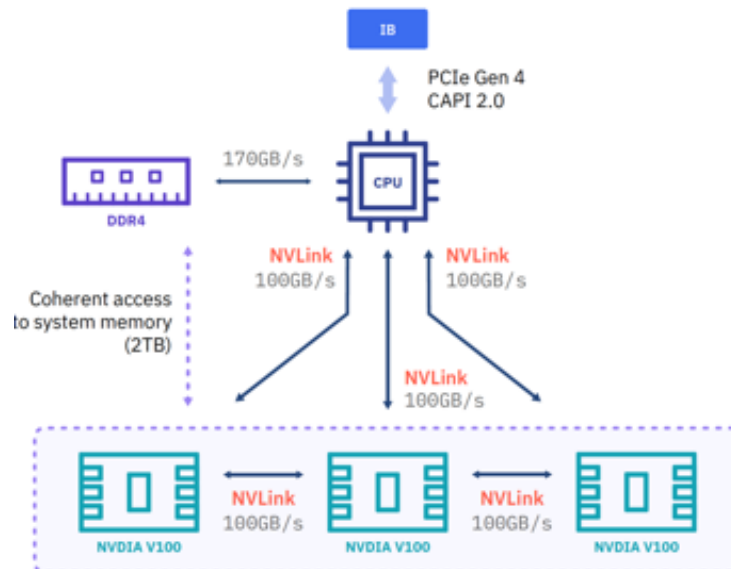
* Competitive HW: 2x Xeon E5-2640 v4; 20 cores (2 x 10c chips) / 40 threads; Intel Xeon E5-2640 v4; 2.4 GHz; 1024 GB memory, 4xTesla V100 GPU, Ubuntu 16.04

4 GPUs - Air (4Q'17)/Water Cooled (Coming)



- Up to 4 GPUs, air/water cooled options
- 150GB/s of bandwidth from CPU-GPU

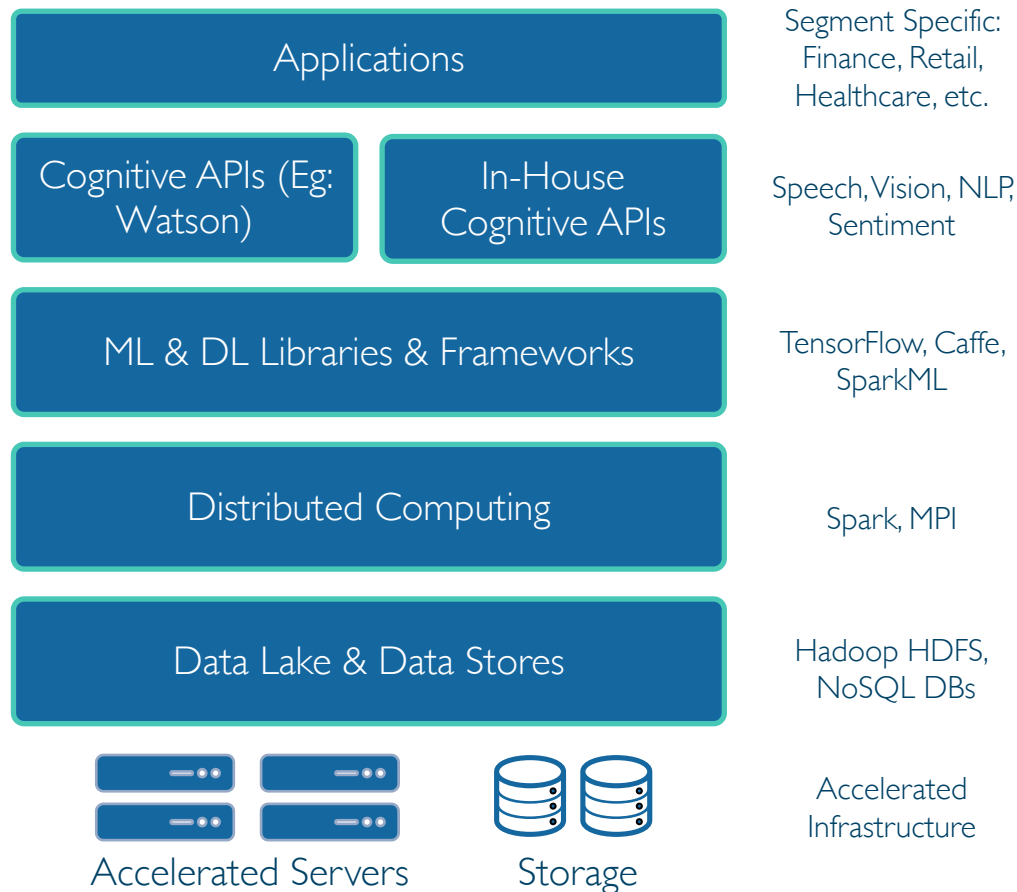
6 GPUs - Water Cooled (coming)



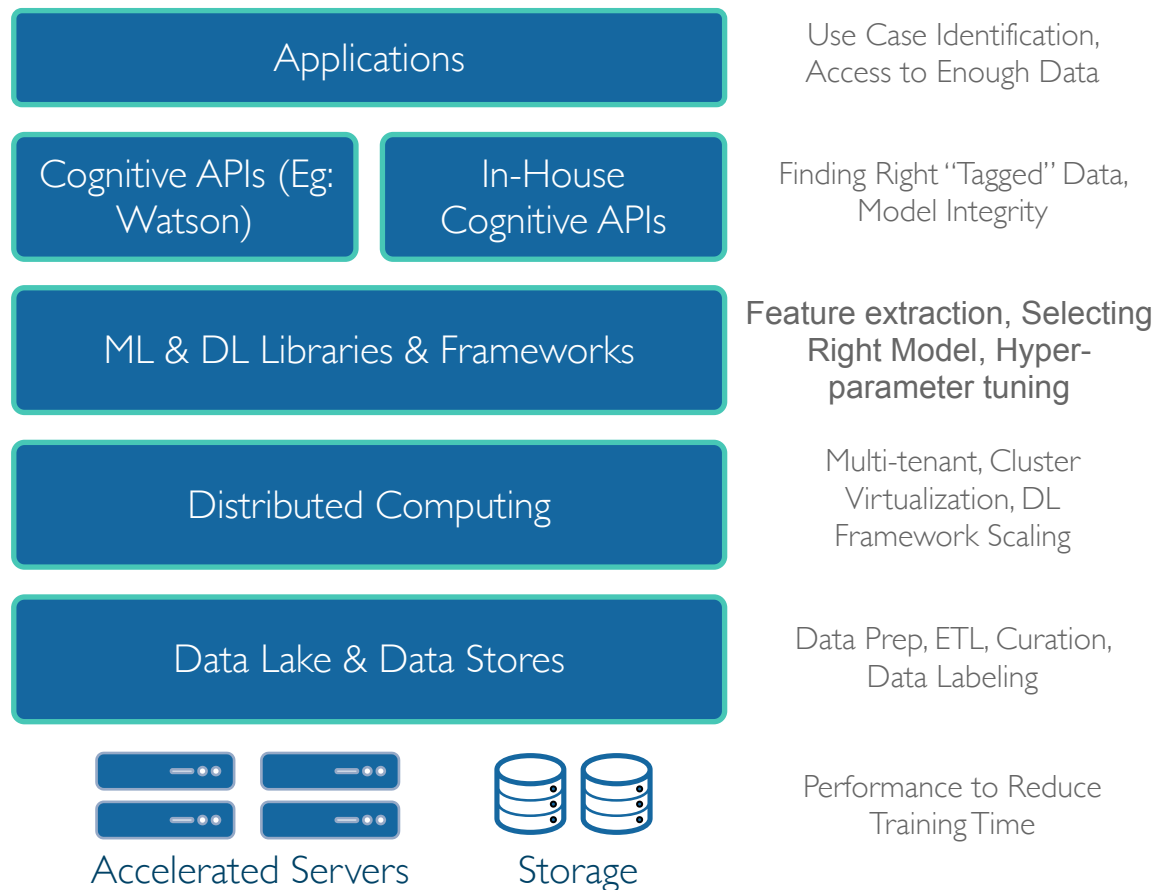
- Up to 6 GPUs, water cooled only
- 100 GB/s of bandwidth from CPU-GPU

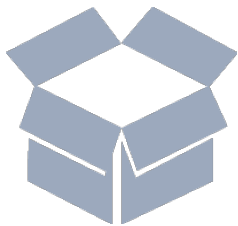
- Coherent access to system memory
- PCIe Gen 4 and CAPI 2.0 to InfiniBand
- Water cooled options available in 2Q'18

Layers in AI Infrastructure Stack

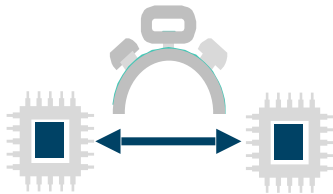


Challenges in Building Cognitive Solutions

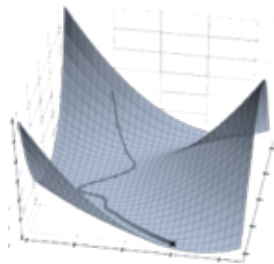




**enterprise-ready
software distribution
built on open source**

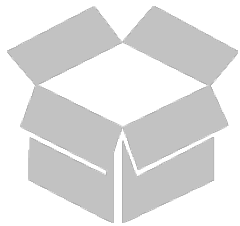


**performance
faster training times
for data scientists**



**tools for ease
of development**

IBM Power**AI**



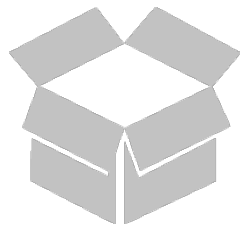
**Enterprise-Ready
Software Distribution
Built on Open Source**

fast, easy deployment



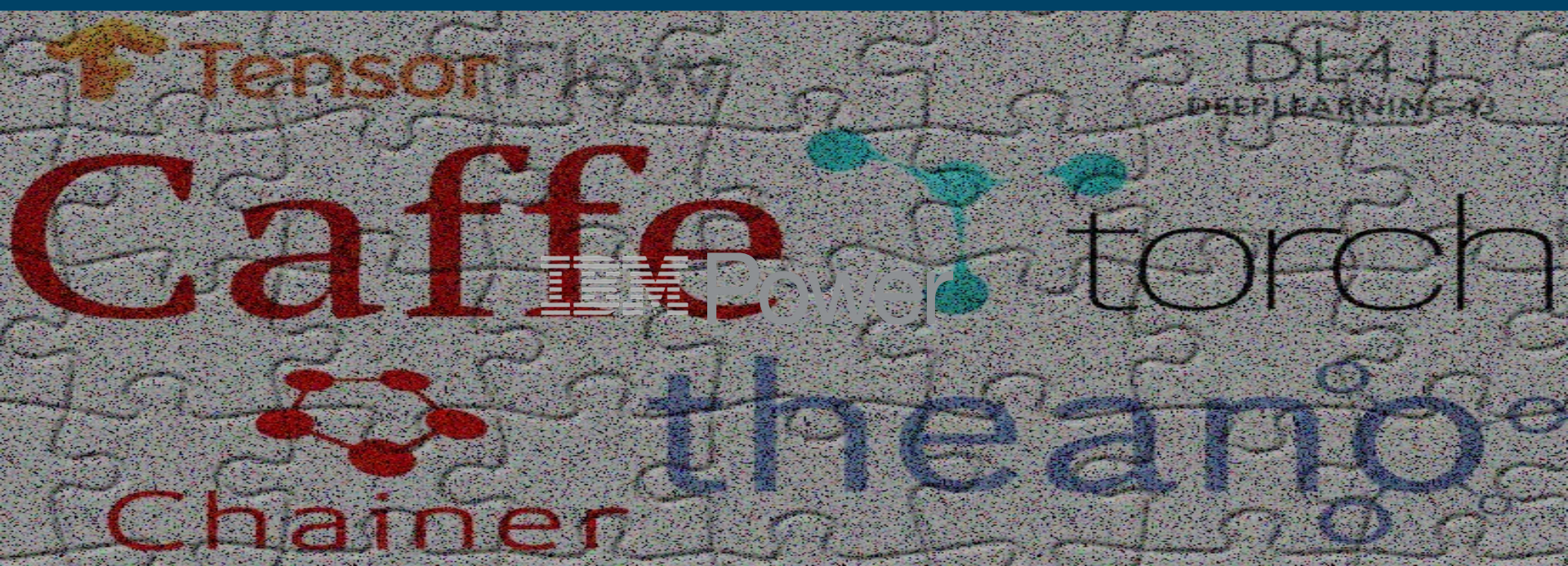
NOT PowerAI!

PowerAI!



**Enterprise-Ready
Software Distribution
Built on Open Source**

***precompiled and current
open source frameworks***



PowerAI v4.0 - Deep Learning Software Distribution

Deep Learning
Frameworks

Caffe

NVCaffe

IBMCaffe

Torch

TensorFlow

Distributed
TensorFlow

Theano

Chainer

Supporting
Libraries

OpenBLAS

Bazel

Distributed
Communications

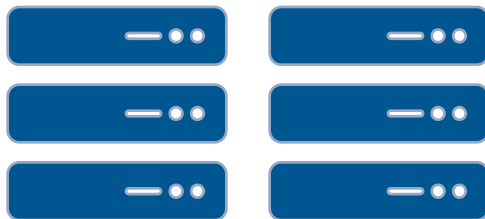
NCCL

DIGITS

Power

Accelerated Servers
and Infrastructure
for Scaling

IBM Cluster of NVLink Servers

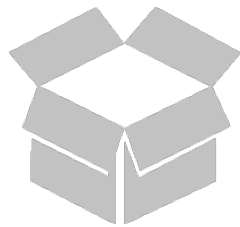


Spectrum Scale:
High-Speed Parallel File
System



Scale to
Cloud



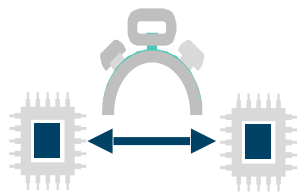


**Enterprise-Ready
Software Distribution
Built on Open Source**

***available enterprise
support for the entire stack***

IBM Power

SUPPORT

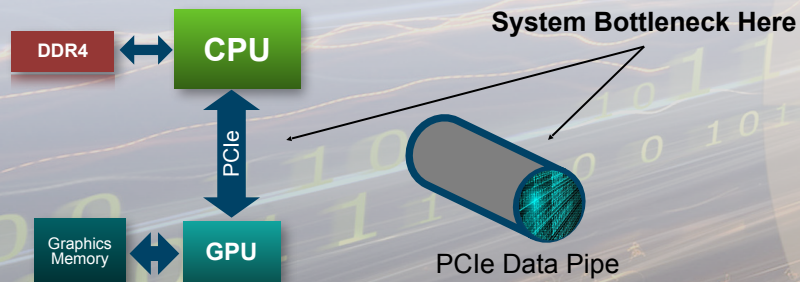


**Performance...
Faster Training
and Inferencing**

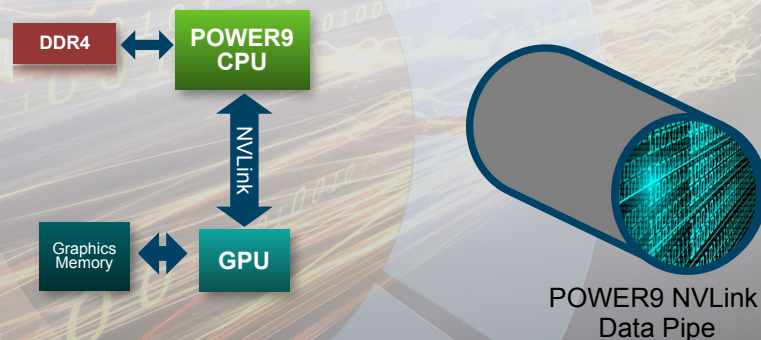
***unique innovation through
OpenPower collaboration***



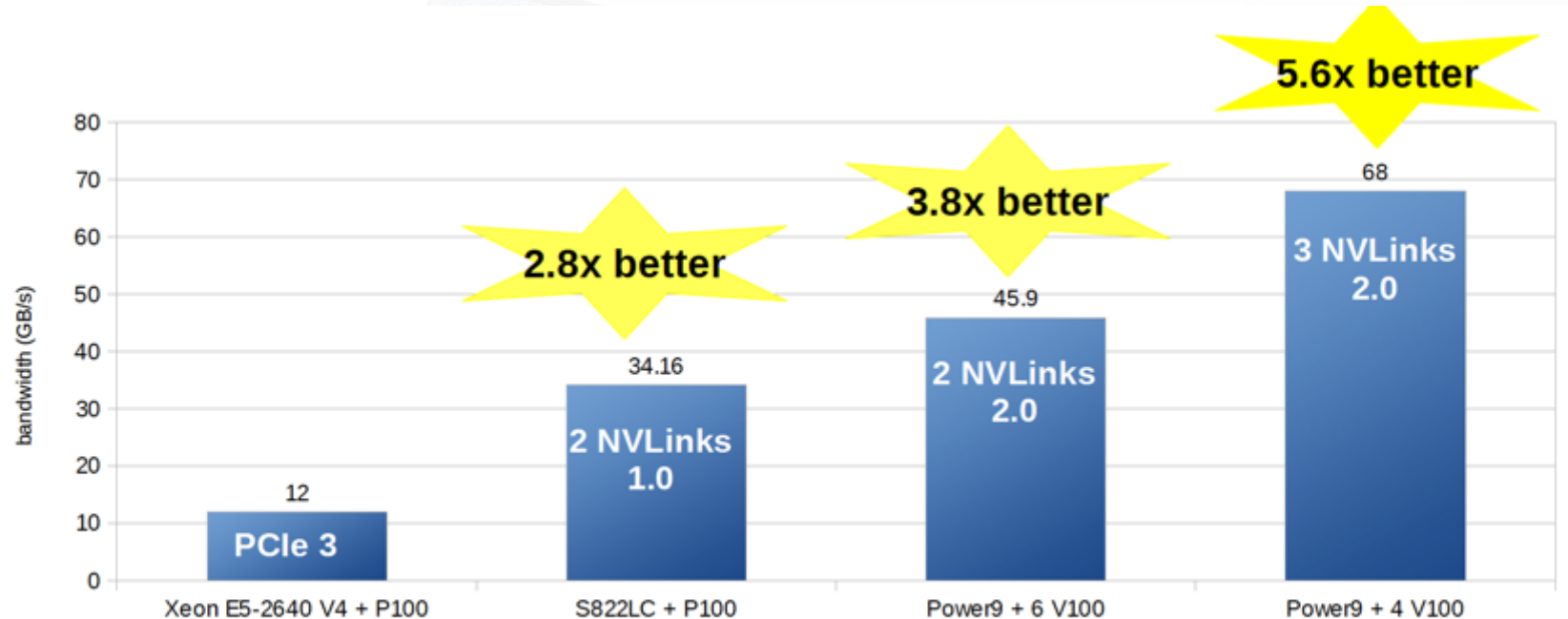
THE SYSTEM BOTTLENECK SHIFTS TO PCI-EXPRESS



***POWER9 with NVLink
delivers 5.6X the bandwidth***



The NVLink difference CPU-GPU

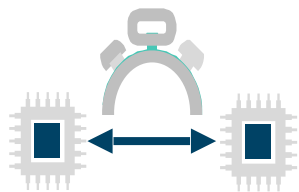


- P9 with 2nd Gen NVLink enables 5.6x faster data movement from CPU-GPU in 4 GPU system
- In 6 GPU system bandwidth is minimally reduced but balanced by higher compute capability

Results are based on IBM Internal Measurements running the CUDA H2D Bandwidth Test

Hardware: Power AC922; 32 cores (2 x 16c chips), POWER9 with NVLink 2.0; 2.25 GHz; 1024 GB memory; 4xTesla V100 GPU; Ubuntu 16.04. S822LC for HPC; 20 cores (2 x 10c chips), POWER8 with NVLink; 2.86 GHz; 512 GB memory; Tesla P100 GPU

Competitive HW: 2x Xeon E5-2640 v4; 20 cores (2 x 10c chips) / 40 threads; Intel Xeon E5-2640 v4; 2.4 GHz; 1024 GB memory; 4xTesla V100 GPU; Ubuntu 16.04



Performance... Faster Training and Inferencing

*faster training times
for data scientists*

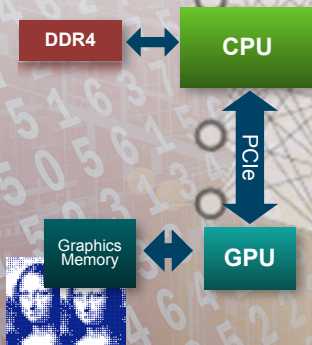
Distributed Deep Learning



Traditional Model Support

(Competitors)

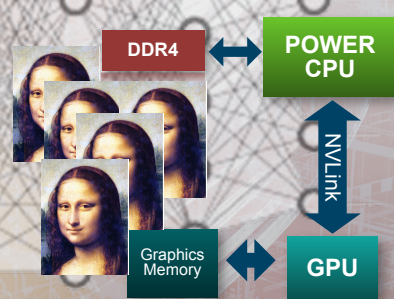
Limited memory on GPU forces trade-off in model size / data resolution



Large Model Support

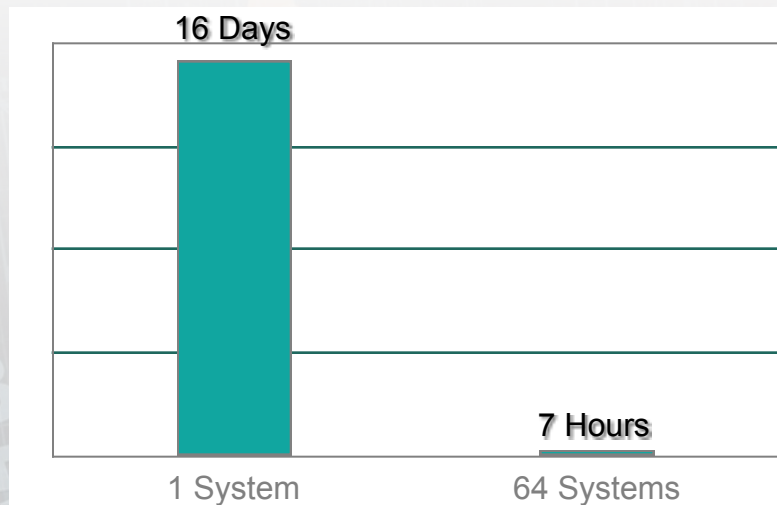
(PowerAI)

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

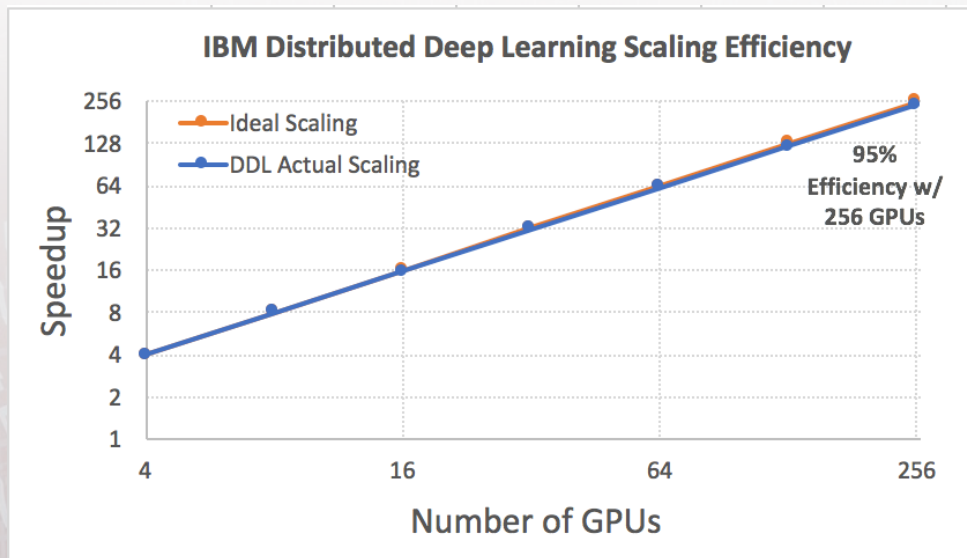


PowerAI Rel. 4 with Distributed Deep Learning

**16 Days Down to 7 Hours:
58x Faster**

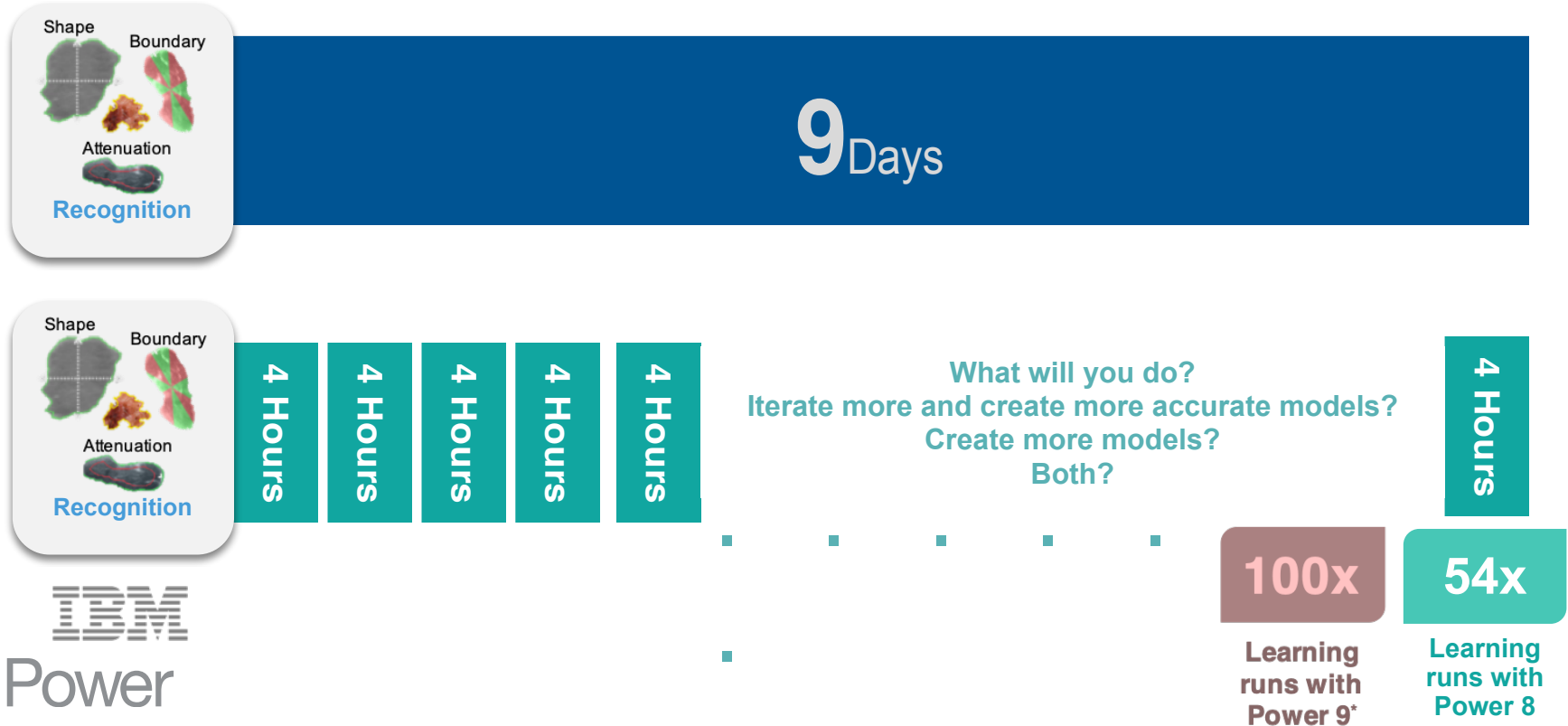


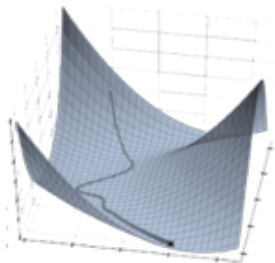
Near Ideal Scaling to 256 GPUs and Beyond



ResNet-101, ImageNet-22K, Caffe with PowerAI DDL, Running on Minsky (S822Lc) Power System

Acceleration training days become hours





**Tools for Ease
of Development**

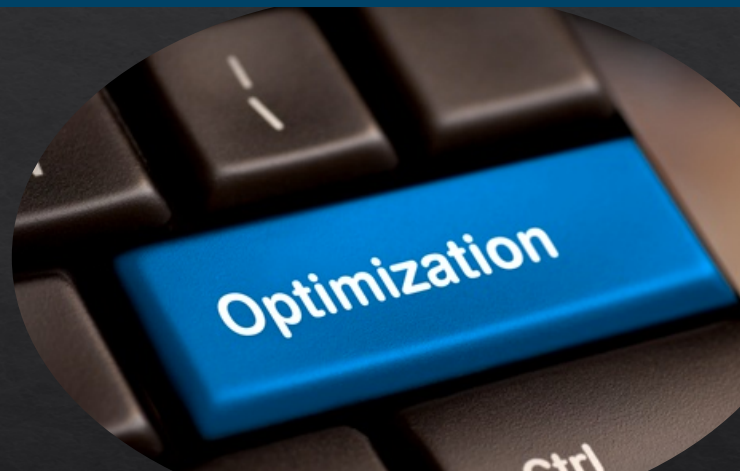
**rich advisory and building
toolsets to flatten
time to value**



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
+++

“IBM Power is a great cognitive platform if not the best out there. The IBM Power team identified the need for and implemented acceleration before anyone else in the industry and are already on their third generation with the highest speed accelerator interconnects and coherent architecture that can share main memory with accelerators.”

Forbes

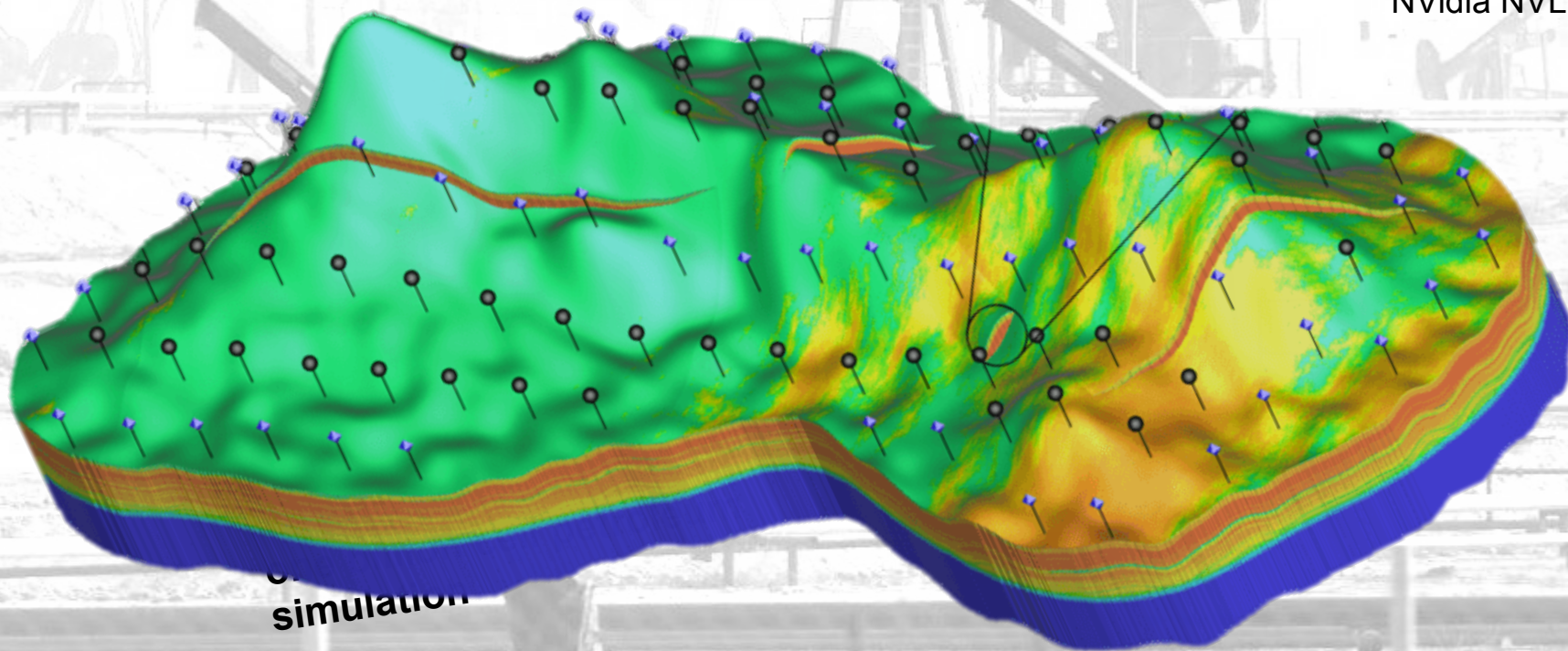
**IBM's New
PowerAI
Features
Demonstrate
Enterprise AI
Leadership
...Again**

Oil & Gas Billion Cell Reservoir calculation

IBM Power

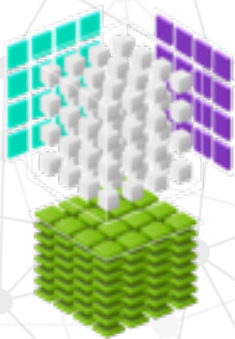
22,000 computer node cluster
716,000 Intel CPUs

30 computer node cluster
60 Power CPUs
120 NVIDIA Tesla P100 GPUs
NVidia NVLink



CORAL – Summit and Sierra on day 1

AI Exascale
Today



Order of Magnitude
Leap in
Computational Power

200 PF

20 PF

3+EFLOPS

Tensor Ops

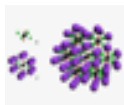
10X

Perf Over Titan

5-10X

Application Perf Over Titan

Real,
Accelerated
Science



DIRAC



HACC



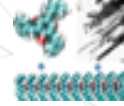
NUCCOR



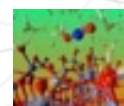
RAPTOR



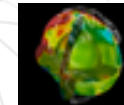
FLASH



LSDALTON



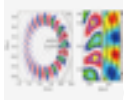
NWCHEM



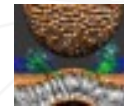
SPECFEM



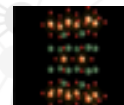
ACME



GTC



NAMD



QMCPACK



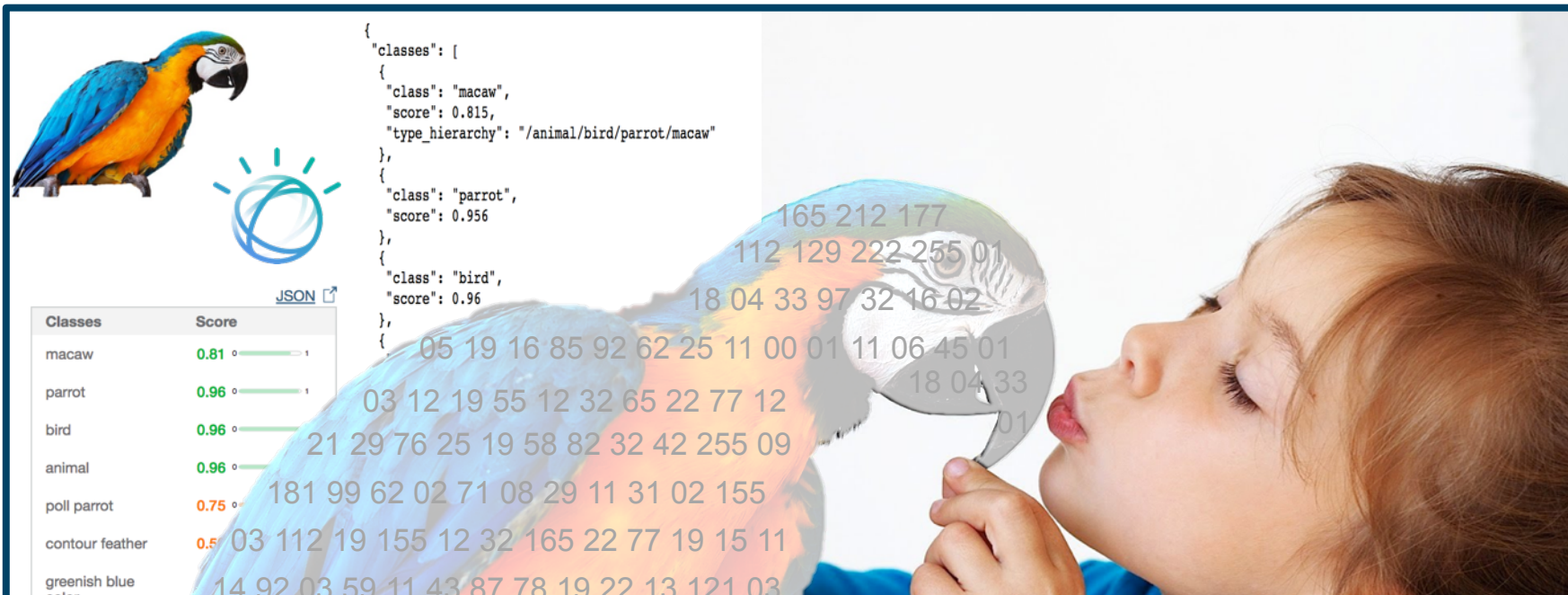
XGC



Deployment at Summit

- Significant application performance over Titan (AMD/NVIDIA) Achieved with $\frac{1}{4}$ the number of servers
- Similar wins at :





IBM's new PowerAI tools automate image recognition

New AI Vision software will make image recognition easier and faster for developers

By Agham Shah, U.S. Correspondent, IDG News Service

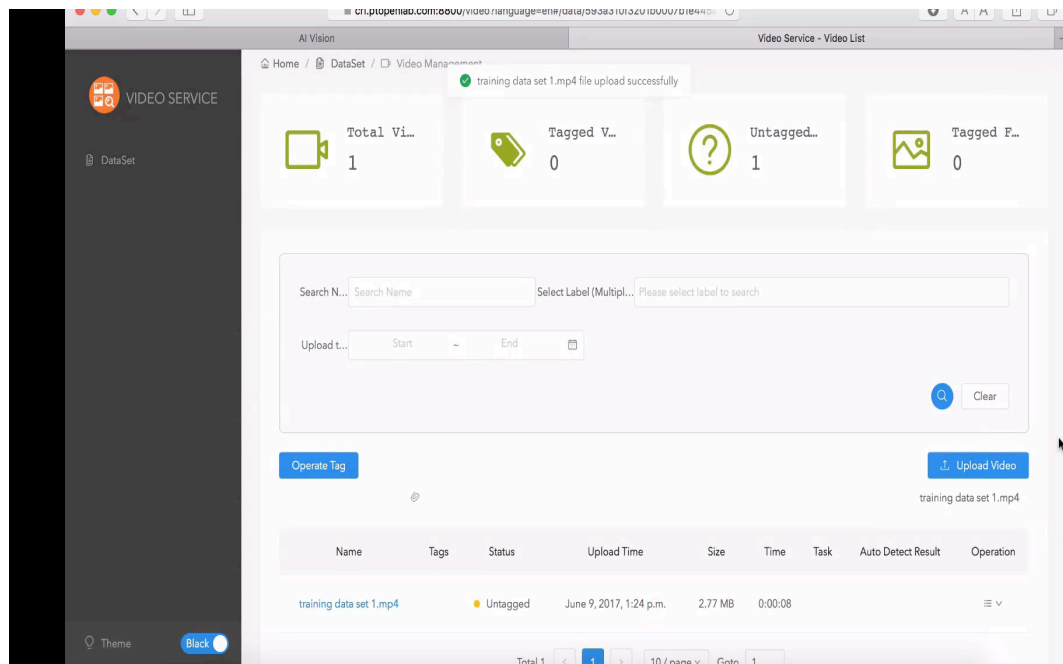
PCWorld

AI Vision toolset

Experts only becomes
beginner knowledge
requirement to build
image-based neural nets

Tooling lets non-techies
label data – brings expertise
to algorithm from LOB and
mitigates errors

Choose best model and
framework to apply based
on data set



**Easily upload multiple videos which will be used to create
labels for objects**

IBM Data Science Experience

~~In my dreams~~

I'm coding in
an open data science
framework,
running on Spark and
Power

...in minutes

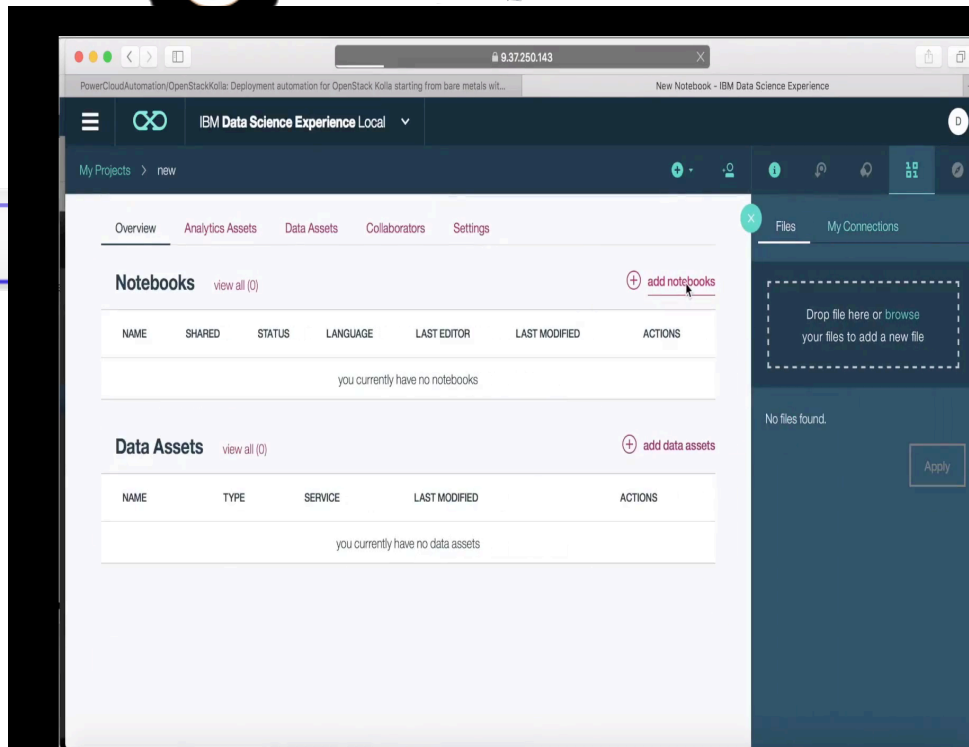
Spark

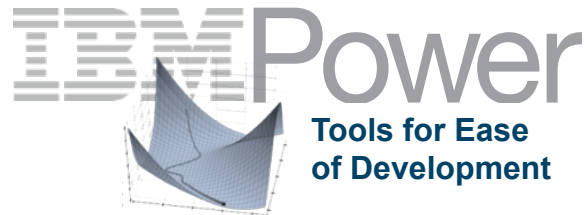
Learn

Create

Collaborate

POWER





Build and deploy with mouse clicks

Built training set with a mere 3 MB of video

Deployed as a REST-API with a mouse clicks

Data is inferenced with a single mouse click

The screenshot shows the IBM AI Vision console interface. A 'Please Confirm' dialog box is open, asking for confirmation to deploy a specific trained model. The background shows a table of trained models with columns for Id, Usage, Categories, Accuracy, Created At, and Operation. The 'Deploy API' button is highlighted in the dialog box.

Id	Usage	Categories	Accuracy	Created At	Operation
e4891e67-a061-4778-94cb-e447d137b431	Object Detection	pedestrian, cyclist, car	0.95000	2017-06-08 22:13:11	Deploy Actions
4cb1ca8b-aedf-46ab-b484-9d2d93c3db97	Image Classification	Acridotheres, Acrocephalus...	0.95000	2017-06-06 11:01:22	Deploy Actions
16e4d003-8f1d-4aa7-b0c7-8ff53573dd48	Image Classification	1, 2, 3	0.96875	2017-06-06 11:00:01	Deploy Actions
9d875b90-d8ae-44b6-ad40-8662b2bb1e7a	Image Classification	fire, no_fire	1.00000	2017-06-05 16:53:18	Deploy Actions
d292875f-a0be-49b2-8f6d-e373470812a2	Object Detection	fire, smoke	0.14617	2017-05-31 11:09:58	Deploy Actions
c83fb544-9907-4fc1-81ff-2ddea10401d6	Object Detection	car, motor	0.57332	2017-05-26 12:30:34	Deploy Actions
cca5c9a0-38c5-4c23-b82d-89ebe3d96718	Object Detection	car	0.96537	2017-05-24 21:01:08	Deploy Actions
37f11db0-ee5e-486c-8622-989a79dcbcb24	Object Detection	car, moto	0.87596	2017-05-16 16:56:11	Deploy Actions
1e96cb5e-8883-4d6c-9516-b896d00ae47b	Object Detection	整体, 头部	0.87596	2017-05-16 12:40:13	Deploy Actions
f05a5016-a77f-4057-a63f-396e7ee35a91	Image Classification	Acridotheres, Acrocephalus...	0.88542	2017-05-07 17:43:51	Deploy Actions

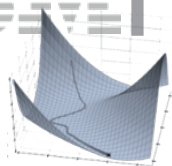
Total: 20, Page Count: 2

Deep Learning

What you and I (our brains) do without even thinking about it.....we recognize a bicycle



IBM Power

Tools for Ease
of Development

Point and click label tuning

10x faster to
create a large
labeled data
set compared
to traditional
methods

AI Vision

Video Service - Video List

1943895@qq.com

My Trained Models

Id	Usage	Categories	Accuracy	Created At	Operation
e4891e67-a061-4778-94cb-e447d137b431	Object Detection	pedestrian, cyclist, car		2017-06-08 22:13:11	Actions
4cb1ca8b-aedf-46ab-b484-5d2d93c3dbe7	Image Classification	Acridotheres, Acrocephal...	0.95000	2017-06-06 11:01:22	Deploy Actions
18e4d003-811d-4aa7-b0c7-8ff53573dd48	Image Classification	1, 2, 3	0.96875	2017-06-06 11:00:01	Deploy Actions
9d875b90-d6ae-44b6-ad40-8662b2bb1e7a	Image Classification	fire, no_fire	1.00000	2017-06-05 16:53:18	Deploy Actions
d292875f-a0be-49b2-8f6d-e373470812a2	Object Detection	fire, smoke	0.14617	2017-05-31 11:09:58	Deploy Actions
c83fb544-9907-4fc1-81ff-2dd6a10401d6	Object Detection	car, motor	0.57332	2017-05-26 12:30:34	Deploy Actions
cca5c9a0-38c5-4c23-b82d-89ebc3d96718	Object Detection	car	0.96537	2017-05-24 21:01:08	Deploy Actions
37f11db0-ee5e-486c-8622-989a79dbcb24	Object Detection	car, moto	0.87596	2017-05-16 16:56:11	Deploy Actions
1e96cb5e-8883-4d5c-9516-b896d00ae47b	Object Detection	整体, 头部		2017-05-16 12:40:13	Deploy Actions
f05a5016-a77f-4057-a63f-396e7ee35a91	Image Classification	Acridotheres, Acrocephal...	0.88542	2017-05-07 17:43:51	Deploy Actions

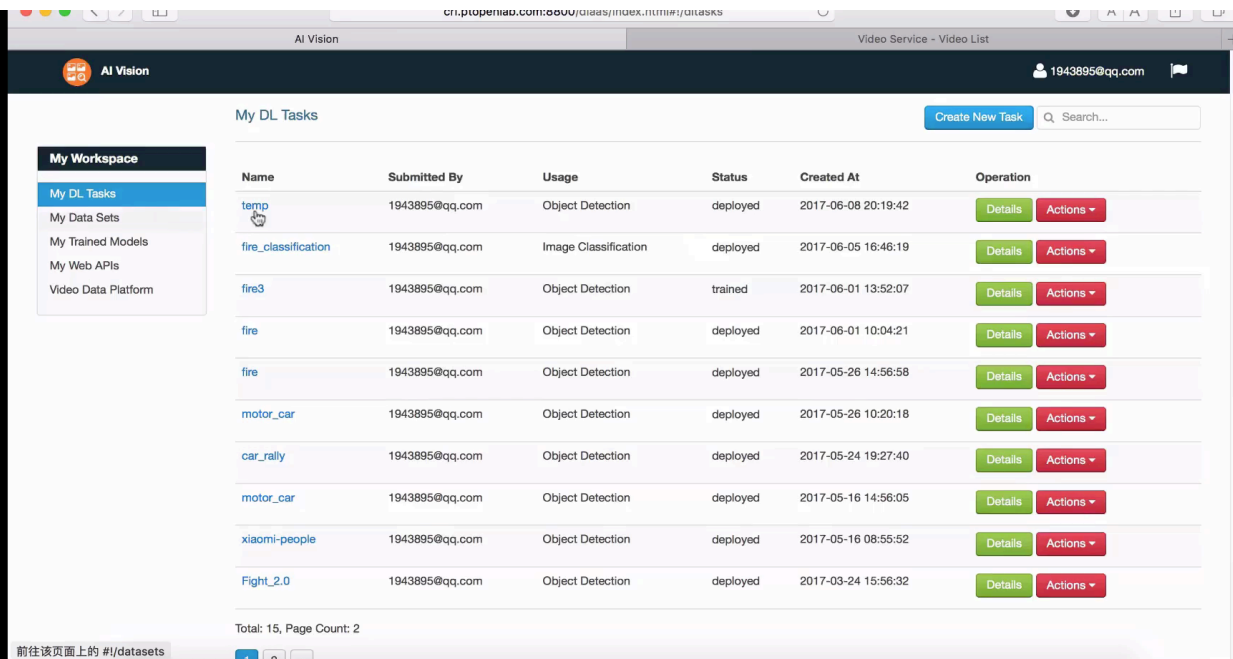
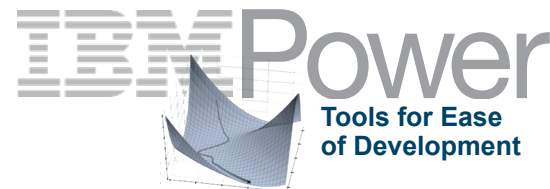
Total: 20, Page Count: 2

1 2

前往该页面上的 #/datasets

AI Vision toolset

AI Vision assistants help to rapidly iterate over models and makes suggestions along the way to improve a model's accuracy



The screenshot shows the 'My DL Tasks' page in the IBM AI Vision toolset. The page has a dark header with the 'AI Vision' logo and a user profile. A sidebar on the left lists 'My Workspace' items: 'My DL Tasks' (selected), 'My Data Sets', 'My Trained Models', 'My Web APIs', and 'Video Data Platform'. The main content area displays a table of tasks with columns: Name, Submitted By, Usage, Status, Created At, and Operation. The table lists 15 tasks, with the first 10 visible. The 'temp' task is highlighted. At the bottom, there is a pagination bar showing 'Total: 15, Page Count: 2' and a link to '前往该页面上的 #1/datasets'.

Name	Submitted By	Usage	Status	Created At	Operation
temp	1943895@qq.com	Object Detection	deployed	2017-06-08 20:19:42	Details Actions
fire_classification	1943895@qq.com	Image Classification	deployed	2017-06-05 16:46:19	Details Actions
fire3	1943895@qq.com	Object Detection	trained	2017-06-01 13:52:07	Details Actions
fire	1943895@qq.com	Object Detection	deployed	2017-06-01 10:04:21	Details Actions
fire	1943895@qq.com	Object Detection	deployed	2017-05-26 14:56:58	Details Actions
motor_car	1943895@qq.com	Object Detection	deployed	2017-05-26 10:20:18	Details Actions
car_rally	1943895@qq.com	Object Detection	deployed	2017-05-24 19:27:40	Details Actions
motor_car	1943895@qq.com	Object Detection	deployed	2017-05-16 14:56:05	Details Actions
xiaomi-people	1943895@qq.com	Object Detection	deployed	2017-05-16 08:55:52	Details Actions
Fight_2.0	1943895@qq.com	Object Detection	deployed	2017-03-24 15:56:32	Details Actions

Total: 15, Page Count: 2

[前往该页面上的 #1/datasets](#)

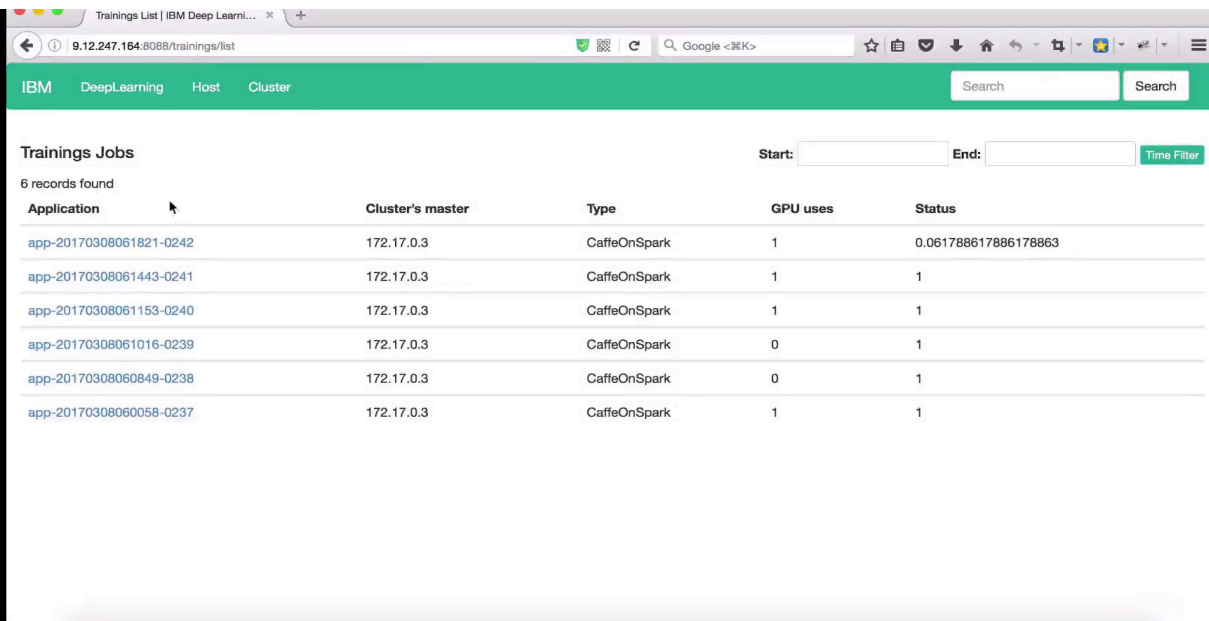
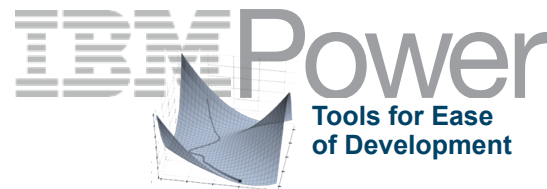
Create a new task to train the model

Real time monitoring of hyper parameters

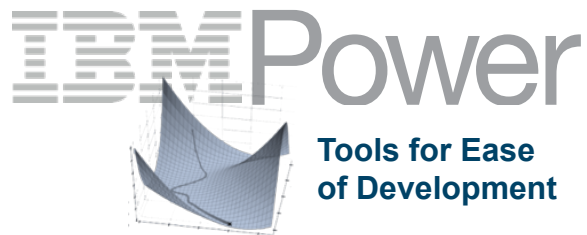
Expert optimization advice for hyper parameter selection and tuning

Traffic light alerting for required parameter optimization with early stop advice and more

CPU, GPU, memory utilization info, comms overhead, +++

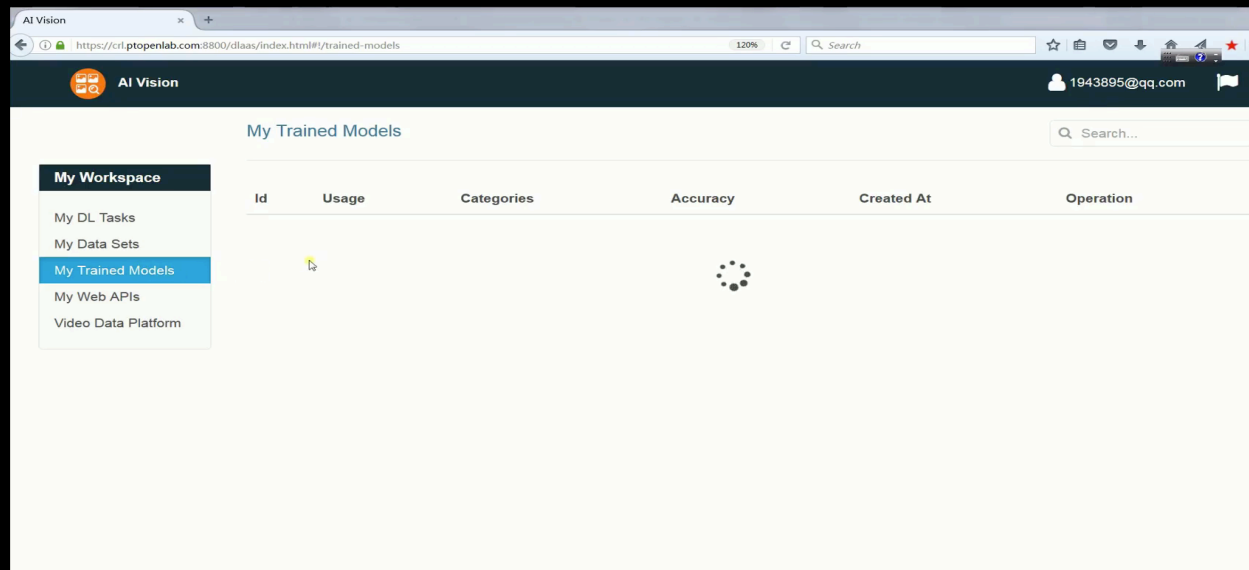
A screenshot of a web browser displaying the "Trainings List" page of the IBM Deep Learning Insight dashboard. The browser's address bar shows the URL "9.12.247.164:8088/trainings/list". The dashboard has a green header with navigation links for "IBM", "DeepLearning", "Host", and "Cluster". Below the header, there's a section titled "Trainings Jobs" with a "6 records found" status. To the right of this section are input fields for "Start:" and "End:" and a "Time Filter" button. The main content is a table with the following columns: "Application", "Cluster's master", "Type", "GPU uses", and "Status". The table lists six training jobs, all of which are "CaffeOnSpark" type and running on the "172.17.0.3" master. The "GPU uses" column shows values of 1, 1, 1, 0, 0, and 1 for the respective jobs. The "Status" column shows values of "0.061788617886178863", "1", "1", "1", "1", and "1".

Application	Cluster's master	Type	GPU uses	Status
app-20170308061821-0242	172.17.0.3	CaffeOnSpark	1	0.061788617886178863
app-20170308061443-0241	172.17.0.3	CaffeOnSpark	1	1
app-20170308061153-0240	172.17.0.3	CaffeOnSpark	1	1
app-20170308061016-0239	172.17.0.3	CaffeOnSpark	0	1
app-20170308060849-0238	172.17.0.3	CaffeOnSpark	0	1
app-20170308060058-0237	172.17.0.3	CaffeOnSpark	1	1

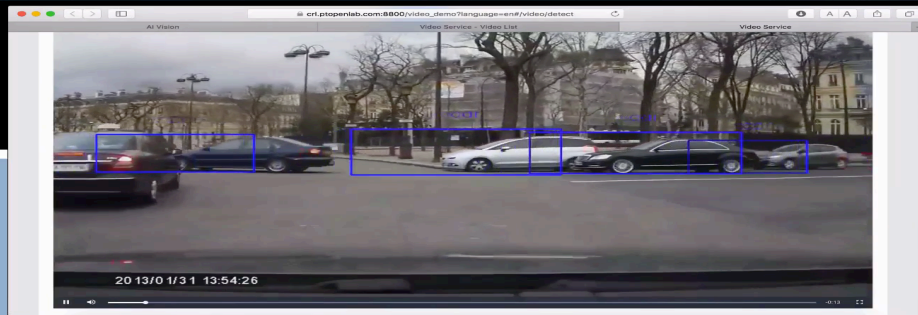


Point and click model deployment

Hardcore coding
days become
clicks to
expose the
model as a
REST-API
endpoint
and scored
from anywhere



IBM Power



Real-time detection of objects for which the Deep Learning model was trained (Car, Motor & Pedestrian)

AI Vision in Data Center

DATA PREPARATION

most time
spent here

up and
running
over a
quick
lunch

time spent
drops from
80% to
30%

9 days work
becomes
hours
more
models

DEPLOY & INFER

requires different
skills

Assign
pairs for
selection and
turning

Iterate faster
and more
again

UP & RUNNING

weeks to months

BUILD, TRAIN, OPTIMIZE

very iterative

MAINTAIN ACCURACY

experience all that
pain again

IBM POWER AC922 pricing

65K street price, including :

- 2 X POWER9 (40 cores)
- 4 X Tesla V100
- 1TB RAM
- NVLink everywhere

**Comparable to DGX Station
Superior performance**

Superb value for the money!



Deep Learning / AI Enterprise Use Cases



AUTOMOTIVE

Auto sensors
reporting location,
problems



COMMUNICATIONS

Location-based
advertising



CONSUMER PACKAGED GOODS

Sentiment analysis of
what's hot, problems



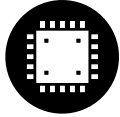
FINANCIAL SERVICES

Risk & portfolio analysis
New products



EDUCATION & RESEARCH

Experiment sensor analysis



HIGH TECHNOLOGY / INDUSTRIAL MFG.

Mfg. quality
Warranty analysis



LIFE SCIENCES

Clinical trials



MEDIA/ENTERTAINMENT

Viewers / advertising
effectiveness



ON-LINE SERVICES / SOCIAL MEDIA

People & career matching



HEALTH CARE

Patient sensors,
monitoring, EHRs



OIL & GAS

Drilling exploration
sensor analysis



RETAIL

Consumer sentiment



TRAVEL & TRANSPORTATION

Sensor analysis for
optimal traffic flows



UTILITIES

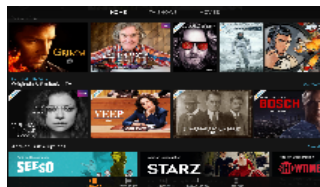
Smart Meter analysis
for network capacity,



LAW ENFORCEMENT & DEFENSE

Threat analysis - social
media monitoring, photo
analysis

Deep Learning in Industries



Automotive and Transportation

- Autonomous driving:
- Pedestrian detection
- Accident avoidance

Auto, trucking, heavy equipment, Tier 1 suppliers (Hyundai, Toyota, Komatsu, General Motors, Volvo)

Broadcast, Media and Entertainment

- Captioning
- Search
- Recommendations
- Real time translation
-

Consumer facing companies with large streaming of existing media, or real time content

Consumer Web, Mobile, Retail

- Image tagging
- Speech recognition
- Natural language
- Sentiment analysis

Hyperscale web companies, large retail (Google photos, Twitter, Woolworths, Aeon)

Security and Public Safety

- Video Surveillance
- Image analysis
- Facial recognition and detection

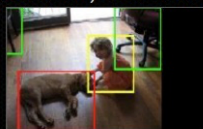
Local and national police, public and private safety/security (ADT, IViz, Pinkerton, Sentry)

Medicine and Biology

- Drug discovery
- Diagnostic assistance
- Cancer cell detection

Pharmaceutical, Medical equipment, Diagnostic labs (Takeda, Asian Pharma, Pfizer)

Image Classification, Object Detection, Localization



Speech & Natural Language Processing



Face Recognition



Medical Imaging & Interpretation



Deep Learning in Banking Industry

1. Predictive Chat Boots for Customer Support
2. Customer recommendations
3. Fraud Detection
4. Algorithmic trading
5. Credit Risk



<http://ieeexplore.ieee.org/document/7359417/>



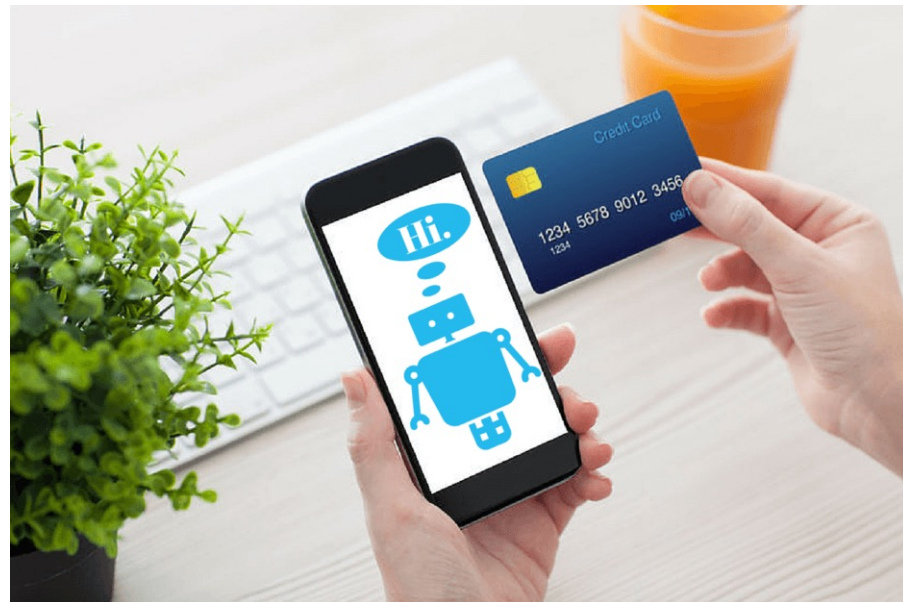
Contextual Chat Boot

What is provided:

- Ability to search for customer financial data
- Provide excellent answers to the top 5 customer support questions
- Handle other questions reasonably well
- Small talk on a basic level

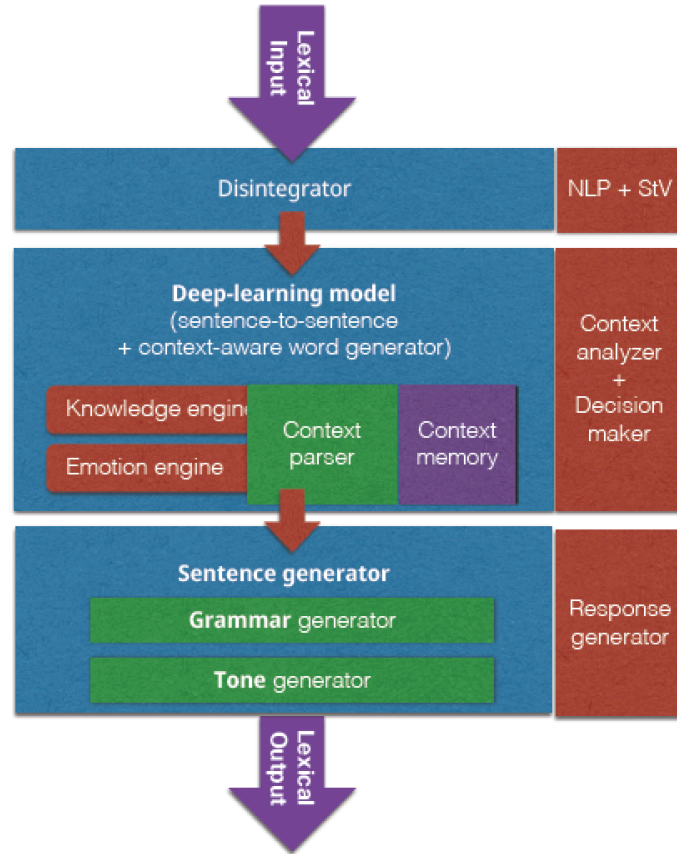
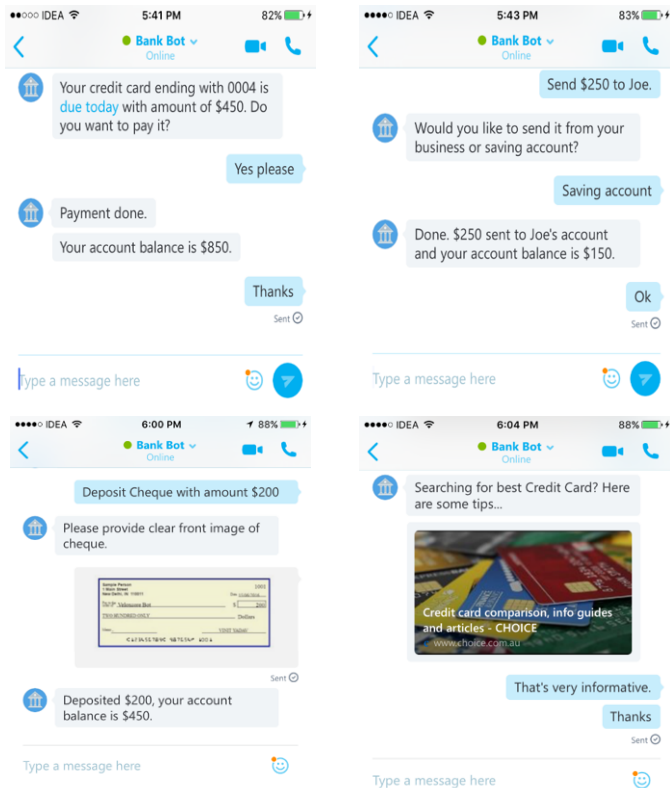
Attributes:

- Enhance
- Understand
- Knowledge
- Decide





Contextual Chat Boot

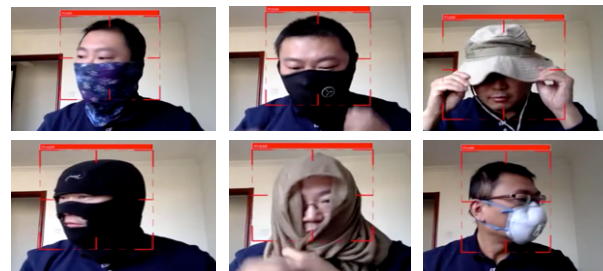




Fraud Detection – Masked Face Detection at ATM's

Challenge

- Unattended ATMs become target of crimes
- Masked face detection could help recognize potential criminal actions and then trigger alarm or limit function
- Traditional pattern identification algorithms are not so effective to resolve such many diversified and changing possibilities



Benefits

- Masked Face detection for ATM proved production ready and significantly improve the security of ATMs and banks

Solution

- A cluster of IBM Minsky servers (S822LC) with Nvidia P100 GPU and NVLink
- PowerAI stack for the deep learning framework, running over Spectrum Conductor with Spark and Spectrum Scale
- Modeling services of face occlusion detection thru Caffe
- Sampled 150+ face occlusion videos to generate 1500+ images as the training and testing dataset
- Real time video recording and auto-detect face occlusion, support multiple ways of recognition simultaneously



EU GDPR discovery using Deep Learning

- **Use Case:** *How to find and identify GDPR data of single individuals from multiple data sources?*
- **Addressable Market:** In general (a) public authorities, (b) organizations that engage in large scale systematic monitoring, or (c) organizations that engage in large scale processing of sensitive personal data. The GDPR not only applies to organisations located within the EU but it will also apply to organisations located outside of the EU if they offer goods or services to, or monitor the behaviour of, EU data subjects. It applies to all companies processing and holding the personal data of data subjects residing in the European Union, regardless of the company's location.
- **Offerings:** ELINAR / PowerAI based solution for GDPR discovery using text mining and deep Learning



JANI WAHLMAN

INFO

TEST AI

TRAIN AI

V1.2-54-G3591352



Upload a file

RECEIPT _ _ X

2013/2017 Receipt LinkedIn

https://www.linkedin.com/company/receipt/3375260618/?printReceipt=true 1/1

First name

LinkedIn ~~releases~~ company ~~Cardner~~ House ~~Wilcox~~ ~~3375260618~~ ~~Ireland~~ VAT : IE9740425P Billed to: ~~MAKAS~~ ~~100000~~ ~~solinnmarkat~~ 28 28 ~~Finland~~ Date ~~3/27/2017~~ Method : MasterCard ~~releases~~ Receipt # ~~357830778~~ Invoice # ~~3375260618~~ VAT # FI10267147 Item Description Rate Quantity ~~3375260618~~ 1 Sponsored Updates - Predictive ~~3375260618~~ - €45.33 Subtotal: €45.33 VAT : 0.00 % €0.00 Invoice : €45.33 ~~3375260618~~ €45.33 Balance : €0.001. You may be required to account for VAT under the reversecharge procedure according to the local VAT rules in your country .

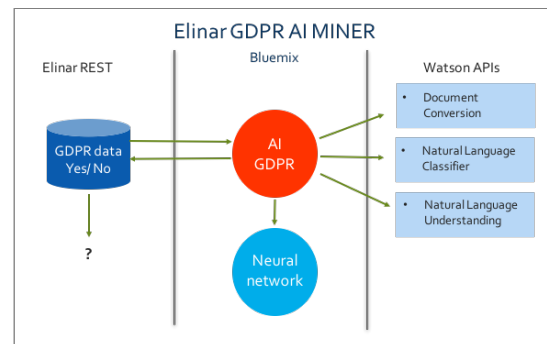
File: Receipt _ LinkedIn 45.pdf

AI: GDPR data found

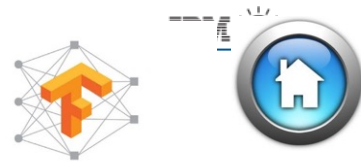
Language: English

Records

Record: 1



Insurance Price Optimization

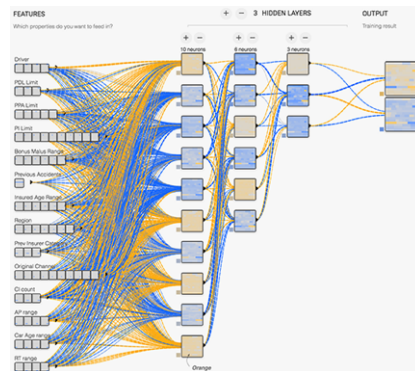


Approximately 7-10% of AXA's customers cause a car accident every year. Most of them are small accidents involving insurance payments in the hundreds or thousands of dollars, but about 1% are so-called large-loss cases that require payouts over \$10,000. As you might expect, it's important for AXA adjusters to understand which clients are at higher risk for such cases in order to optimize the pricing of its policies.

At the right, you can see there are about 70 values as input features including the following.

- Age range of the driver
- Region of the driver's address
- Annual insurance premium range
- Age range of the car

AXA entered these features into a single vector with 70 dimensions and put it into a deep learning model in the middle. The model is designed as a fully connected neural network with three hidden layers, with a [ReLU](#) as the activation function. This use case has end with a accuracy of 78%



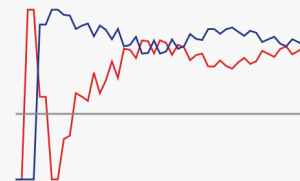
TEST PERFORMANCE

Large Loss Accuracy 0.783

Non-Large Loss Accuracy 0.785

Previous Best by Random Forest:

Large Loss Accuracy 0.386





90%
inspection
times



10X
number of
inspections



accident
risk
rate



IBM Power



KEPCO
KEPRI

AI as a Service – Fraud Surveillance

Financial services organisation evaluates high-performance PowerAI solution to combat fraud.

Winning Solution

- IBM Power System S822LC
- IBM Power AI Development Platform
- RHEL

Total Contract Value: \$92k

Competition: Nvidia DGX-1

Solution Benefits

- The shared, standardised & agile on prem Cloud AI PaaS Service offers CS LoBs cost benefits & time to market competitive advantage for new innovative AI services
- Significantly reduce the cost of Fraud & Fraud Surveillance & Detection
- Time to value with one click installation
- Level 3 support based on Enterprise grade AI development platform
- Significantly improved model accuracy enabled by large ML/DL model support leveraging IBM unique hardware acceleration architecture
- Significantly reduced model training times leveraging unique Distributed Deep Learning capability & graphical monitoring of training accuracy enabled by Deep Learning Impact
- Most cost effective Distributed Computing Framework – IBM Spectrum Computing

About the customer

Credit Suisse is one of the world's leading financial services providers and offers clients its combined expertise in the areas of private banking, investment banking and asset management. Credit Suisse provides advisory services, comprehensive solutions and innovative products to companies, institutional clients and high-net-worth private clients globally, as well as to retail clients in Switzerland. Credit Suisse is headquartered in Zurich and operates in about 50 countries worldwide.

Business challenge and solution

Credit Suisse has been providing low latency HPC Monte Carlo Simulation PaaS services based on IBM's market leading workload scheduling & shared resource management middleware platform (Symphony); & Data Analytics PaaS services based on Hadoop/ MapReduce to the Capital Markets Trading LoBs for many years; to support Risk Analytics, Pricing & FRTB driven XVA calcs . These mission critical services are deployed on a global 60,000 + core commodity computing grid cluster, with a few nodes using GPU hardware acceleration for specific Trading Risk analytics use cases. Within the last 6 months, the CS SAMI organisation has extended this PaaS capability to include a ML/DL AI PaaS service for Fraud Surveillance & Detection of all Trading email communication, using NLP algorithms based on the Theano framework. The challenges CS are facing is to scale the ML/DL AI PaaS service efficiently, to meet the performance demands of an increasing number of ML/DL uses across the firm, from Trading Risk Analytics & Market Predictions, extending the Fraud use case to include voice to text; to many operational efficiency use cases. The solution undergoing evaluation is IBM's enterprise grade AI application development framework, PowerAI, Vision AI, Distributed Deep Learning & Deep Learning Impact, leveraging IBM's unique hardware acceleration platform Power8 Minsky, based on Nvidia GPUs & NVLink technology.

Why did the Client choose IBM Systems

CS have selected IBM's AI platform for evaluation as the AI platform of choice, based on Vision & Strategy alignment with IBM's highly differentiated SDI Vision & Strategy. This leverages a common middleware workload scheduling & shared resource management platform, IBM's Spectrum Computing Framework. This uniquely manages not only traditional compute intensive low latency HPC & data intensive analytics workloads, but also large scale ML/DL AI model workloads, leveraging IBM's unique Power8 Minsky hardware acceleration capability, allowing CS to meet the wide range of AI use case performance demands cost effectively.

Business Partner: Recarta

**Winning Solution**

IBM Power AI AC922

Total Contract Value:

50k\$

Competition:

x86

Solution Benefits:

- 2x performance over x86
- NVLink technology with communication between CPUs and GPUs

Business Partner:

UMB

Use Case Contact: Rene Bersier

About the customer

LGT Group is the private banking and asset management group of the princely House of Liechtenstein. Originally known as The Liechtenstein Global Trust, LGT is the largest family-owned private wealth and asset manager in the world, wholly owned by the Prince of Liechtenstein Foundation. LGT is headquartered in Vaduz, Liechtenstein

Business challenge and solution

LGT Bank want to handle their data faster and want to get more out of their data to save time and use also new ways of AI. They have seen our new announcement from the PowerAI AC922 and liked the idea to start with deep learning to program neural networks and to analyse their data faster. Their idea is to analyse their data and also to learn more about deep learning. Specific use cases will be defined in March/April.

Why did the client choose IBM Systems

LGT Bank have bought one E880C and get the PowerAI included. With the new Announcement they understand how easy it is to set up Open Source Frameworks, which are pre-packaged and easy to implement. They like the fact that large Deep Learning jobs can be clustered across several servers with PowerAI. The key for them was the benefit in the NVLink technology, which creates a direct, ultra-fast connection between CPU and GPU, with up to 5.6 x faster communication. With that they get a faster output of their data.

OTP Bank Hungary

GPU accelerated server win

Financial services institution evaluates GPU accelerated computing in several use cases

Winning Solution

2 pcs of AC922 servers, each with:

- 32 core Power9 processor
- 1 TB memory,
- 4x NVIDIA Volta GPU,
- 100Gb EDR IB connection

Total Contract Value: \$ 110,000

Competition

- x86 vendors

Solution Benefits

- NVLink integrated GPU cards offer superior performance in the same price range as x86 competitors

Business Partner: InterComputer

Use-case contact: Jozsef Suranyi
(IBM HU)

About the customer

Owing to economic and legal considerations, OTP Group provides its universal financial services through several subsidiaries. In Hungary, traditional banking operations are performed by the Bank while specialized services, including car leasing and investment funds are developed and offered by the Bank's subsidiaries. OTP Bank has completed several successful acquisitions in the past years, becoming a key player in the region. OTP Group currently operates in Bulgaria (DSK Bank), Croatia (OTP banka Hrvatska), Romania (OTP Bank Romania), Serbia (OTP banka Srbija), Slovakia (OTP Banka Slovensko), Ukraine (CJSC OTP Bank), Russia (OAO OTP Bank, Donskoy Narodny Bank) and Montenegro (Crnogorska komercijalna banka AD) via its subsidiaries.

Business challenge and solution

OTP Bank has a small, independent special "Research and Innovation" IT team to test and evaluate leading IT technologies. They purchased a Power8 Minsky server in 2016 and tested under several workloads, ranging from in-house developed GPU assisted algorithms to GPU database for relational or non-relational data. Convinced by the superior performance of the server, they ordered the Power9 follow-on within two weeks of the announcement. They will further test MapD Core GPU db on Power and also porting Ethereum blockchain technology to Linux on Power.

Why did the Client choose IBM Systems

Designed to deliver extremely high performance for resource-intensive workloads, the IBM Power System AC922 is the only architecture with the latest NVIDIA NVLink technology, creating a direct, ultra-fast connection between CPU and GPU. AC922 is available on the same price range as brand x86 servers, but delivers superior performance and technical design.

Internal use only, OTP is not an approved reference.

AI in Belux



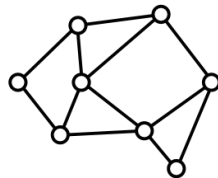
- AI Brussels meetup : 1175 members
- AI Mons : 179 members
- AI Ghent : 50 members
- Machine learning Luxembourg : 339 members
- Data science Luxembourg : 1605 members
- Applied deep learning (Antwerp) : 143 members
- TensorFlow Belgium : 892 members
- **AI 4 business conference** in Brussels on February 27, IBM presence

PowerAI in Europe

- Nordics : nearly 700 members in total (4 cities)
- London : 587 members
- Brussels : created 21/2

RoboVision

Located in Ghent



ROBOVISION
DEEP LEARNING APPLIED

12 employees, all data scientists / AI experts

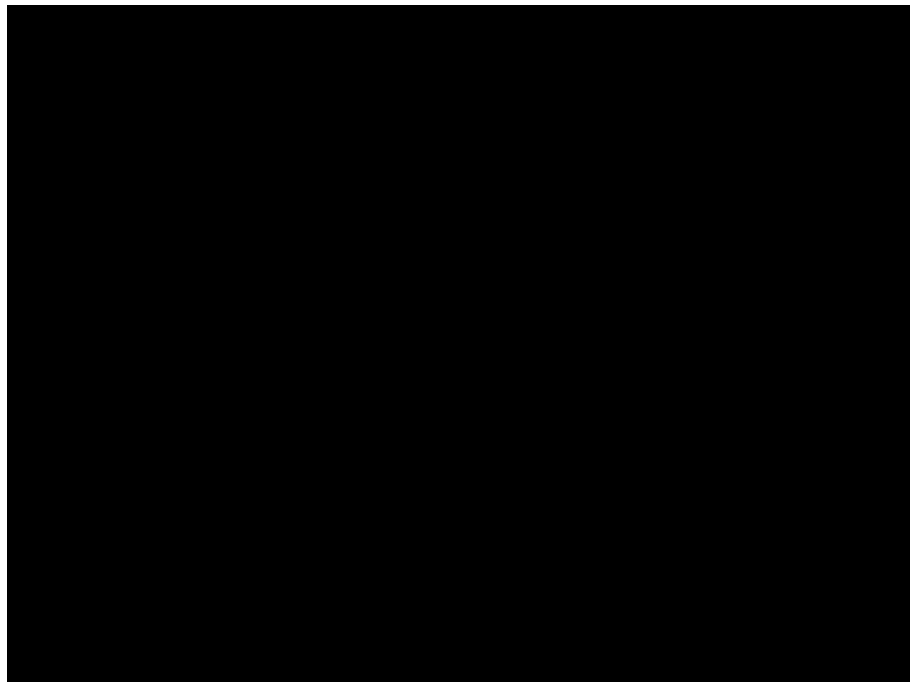
Use cases in manufacturing, finance,
agriculture, retail

References : Audi, BNP Paribas Fortis

Partnering with IBM since Q2 2017

POC in Montpellier/Poughkeepsie in Q3

On premise POC in Q4/Q1 2018



11th December 2017 | Written by Brytlyt

Brytlyt GPU database smashes benchmark record again, this time using IBM Minsky Hardware

In an independent benchmark by industry expert, Brytlyt's GPU Database outperformed all other vendors by a factor of four or more with its PostgreSQL fork tapping into the super computing power of IBM Minsky Hardware.

Based on PostGreSQL

Easy to integrate into existing landscape

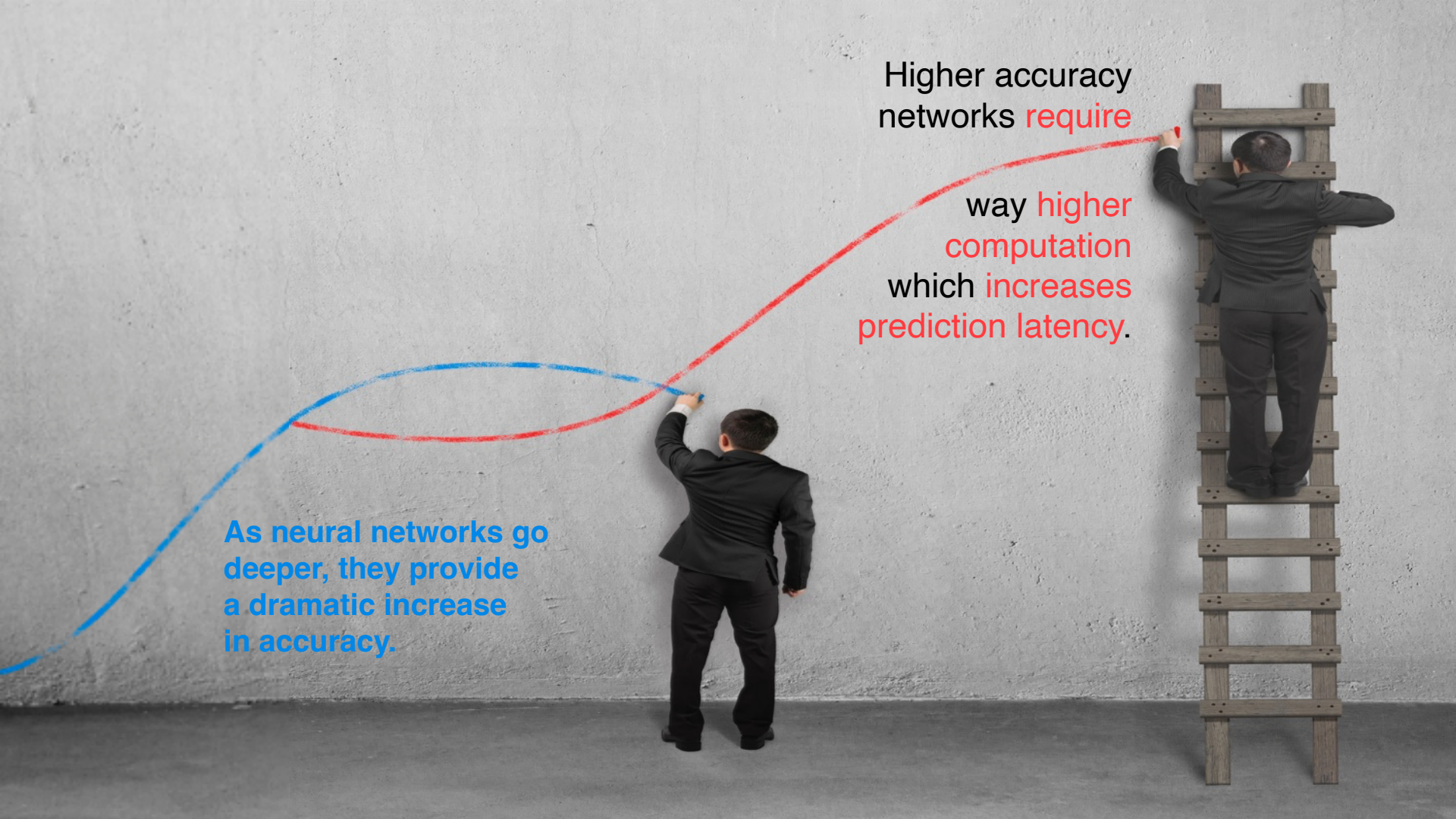
Rich functionality

Results based on (sub-optimal) POWER8 hardware

<https://www.brytlyt.com/blog/brytlyt-gpu-database-smashes-benchmark-record/>

[http://tech.marksblogg.com/billion-nyc-taxi-rides-brytlytdb-ibm-minsky.html?
utm_source=brytlyt%20website&utm_medium=MinskyBenchmark](http://tech.marksblogg.com/billion-nyc-taxi-rides-brytlytdb-ibm-minsky.html?utm_source=brytlyt%20website&utm_medium=MinskyBenchmark)





As neural networks go deeper, they provide a dramatic increase in accuracy.

Higher accuracy networks **require** way higher computation which **increases** prediction latency.