

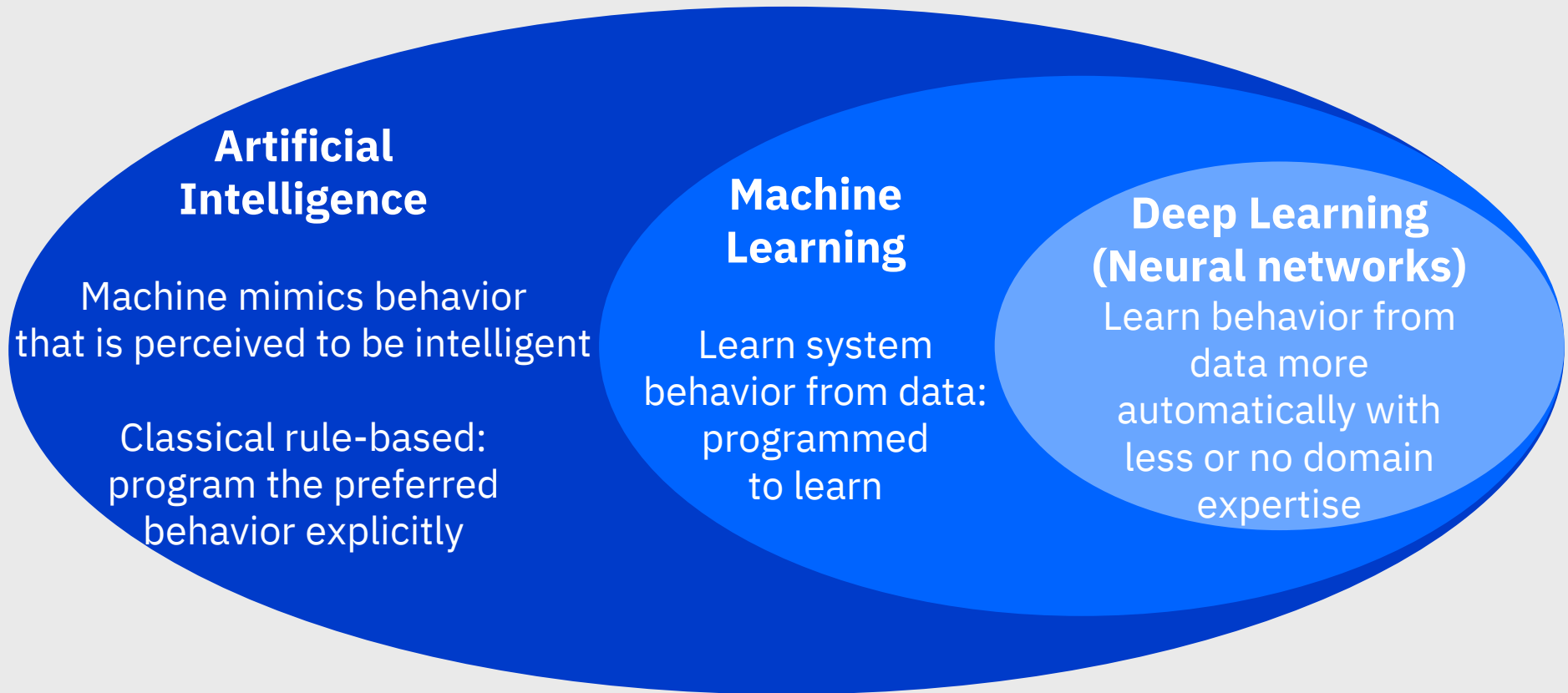
Multiplied data science with IBM deep learning platforms - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

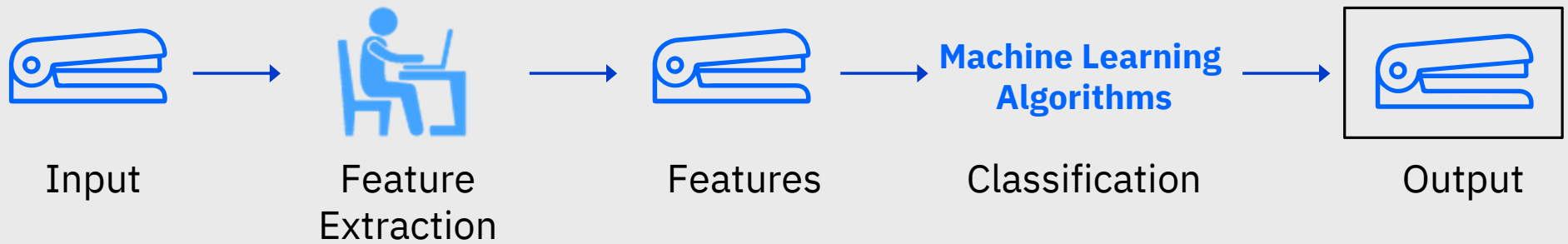
Deep learning development
and computation platforms:

From limited deep learning HW and SW setups
to productized versatile solutions and scalability

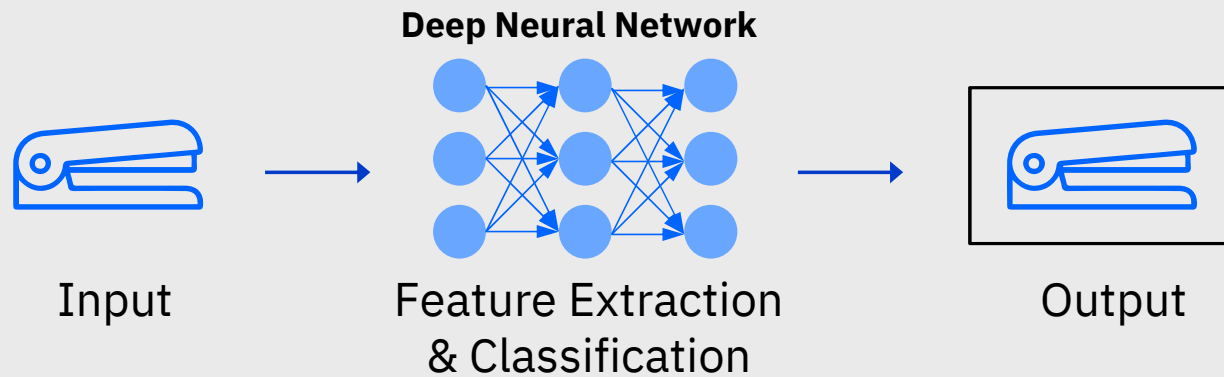
AI Landscape – tools for digital automation



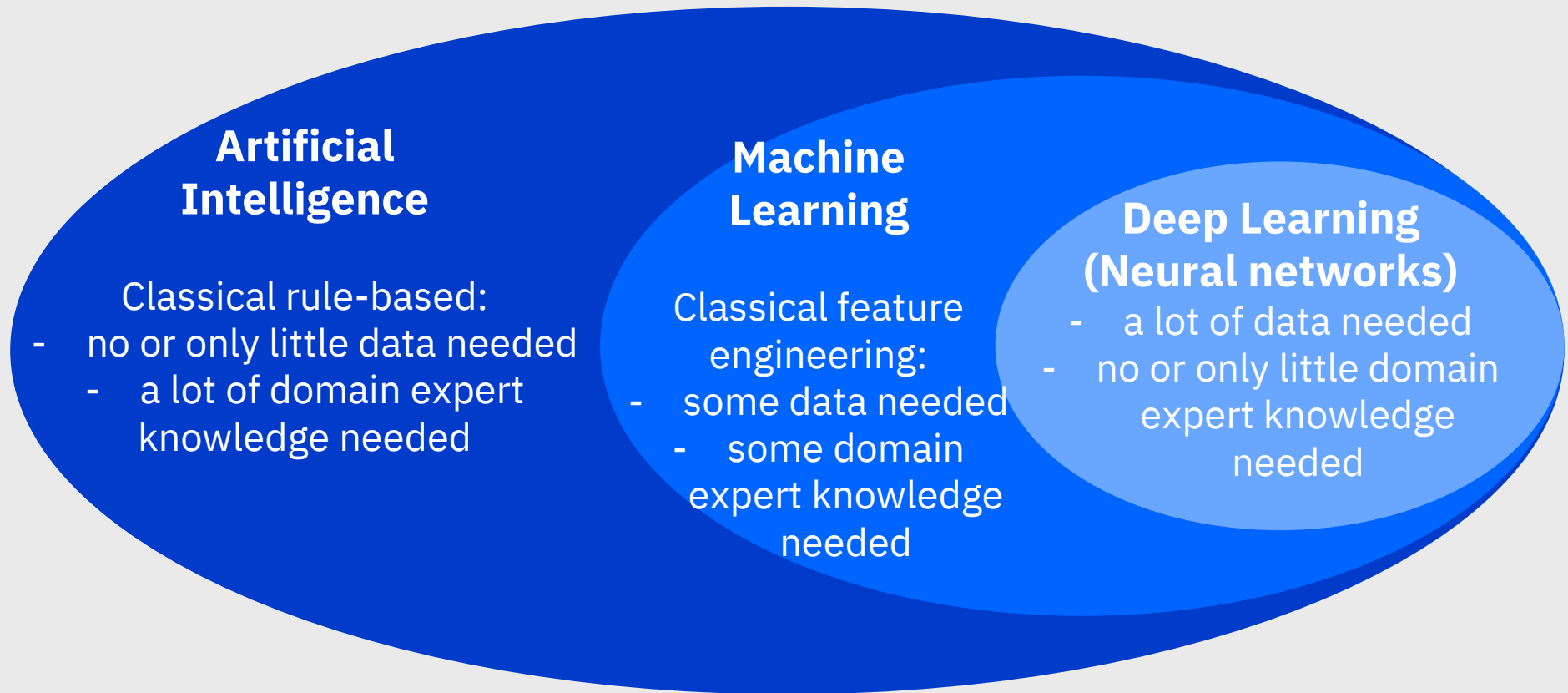
Machine Learning



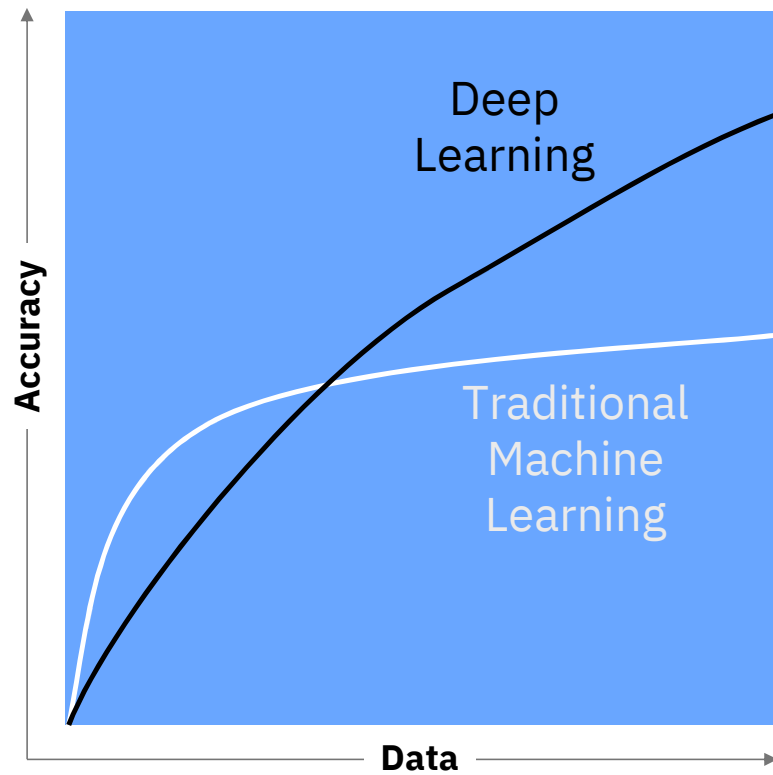
Deep Learning



AI pre-requisites: *application domain expert* knowledge and data



Deep Learning Has Revolutionized Machine Learning



Deep Learning Popularity Growing Exponentially



Source: Google Trends. Search term "Deep Learning"

Deep learning starting point..

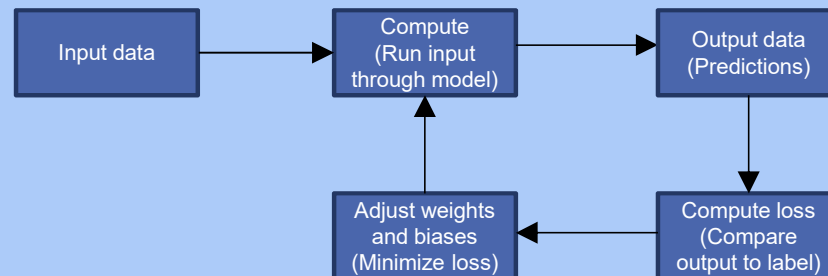
Typical target:

Task automation development through supervised learning by input-output mapping approximation available using neural nets

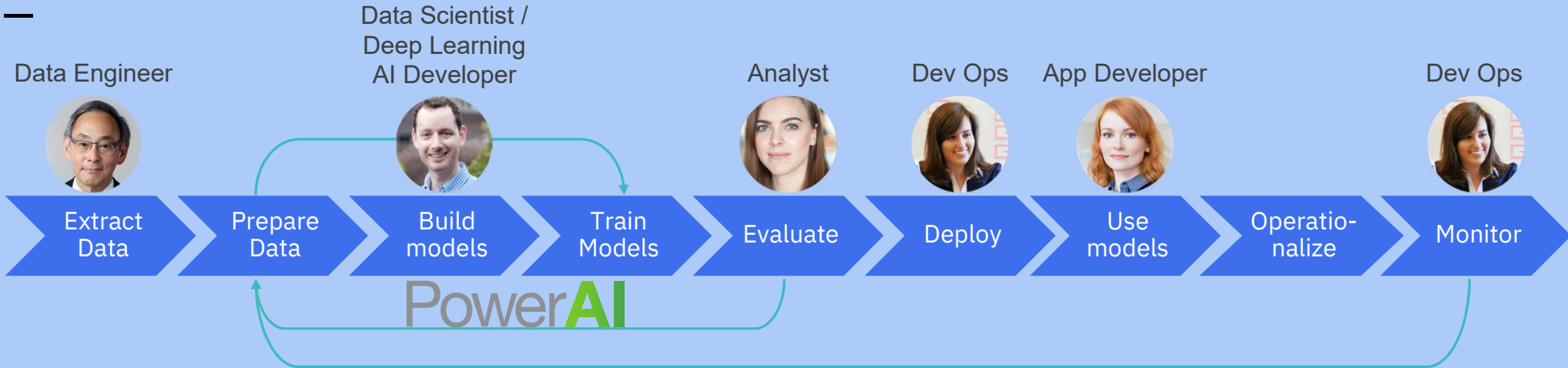
- X=medical image/signal/variables, Y=disease class/severity/prognosis
- X=customer feedback/preferences/.., Y=service action/campaign.., etc.



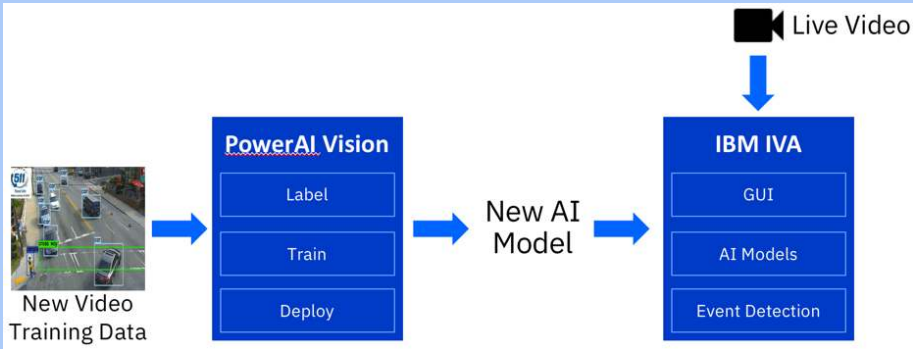
Calculation and data(/IO) intensive computing => GPU-training (thousands of cores) on high bandwidth systems



Systematic data science: deployment, integration and continuous AI development/model improvement



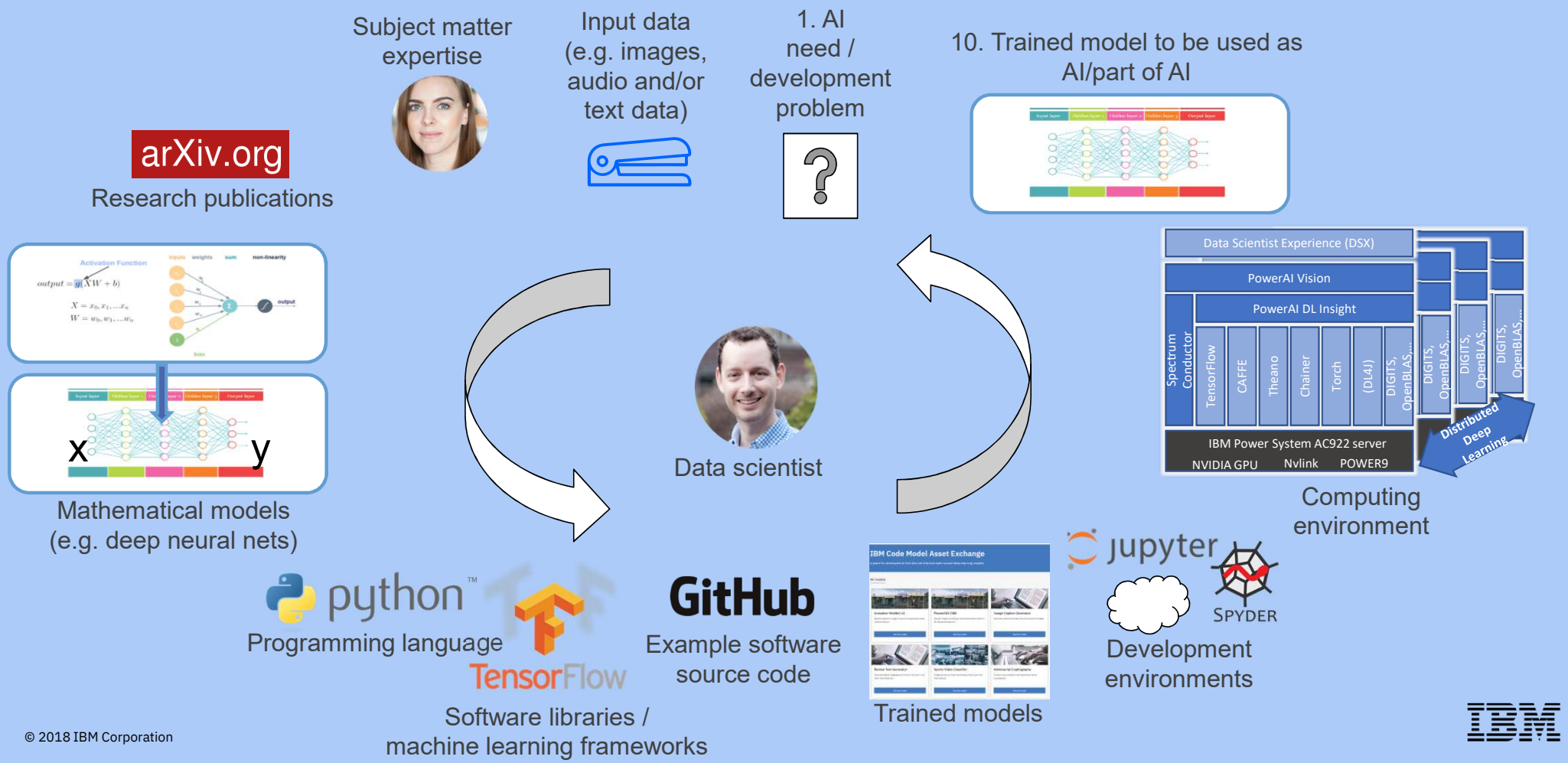
Interested parties need to be able to consume the trained model and the model needs to be supported



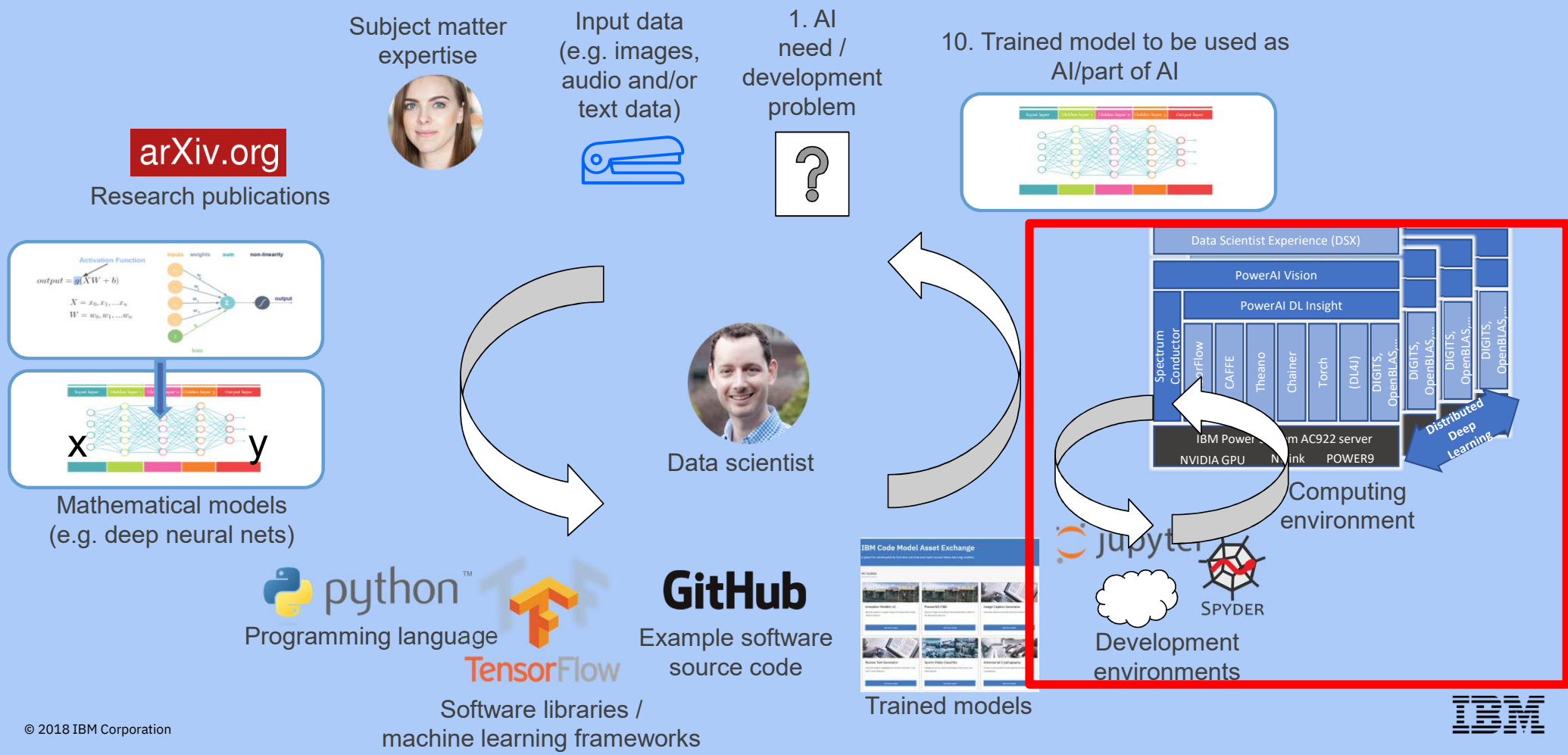
https://medium.com/@sumitg_16893/deep-insights-with-ai-for-video-analytics-5464fd30ebe1



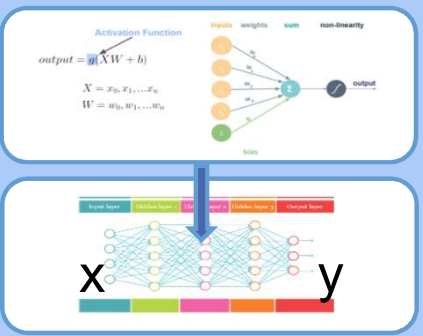
Machine learning and deep learning in AI development - dependencies



Machine learning and deep learning in AI development - dependencies



arXiv.org
Research publications



Mathematical models (e.g. deep neural nets)

python
Programming language

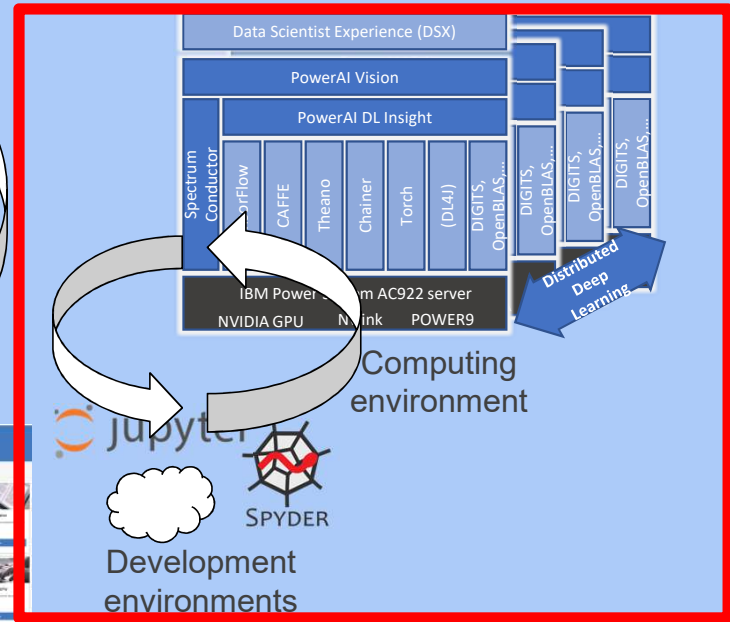
TensorFlow

Software libraries / machine learning frameworks

GitHub
Example software source code

IBM Code Model Asset Exchange

Trained models



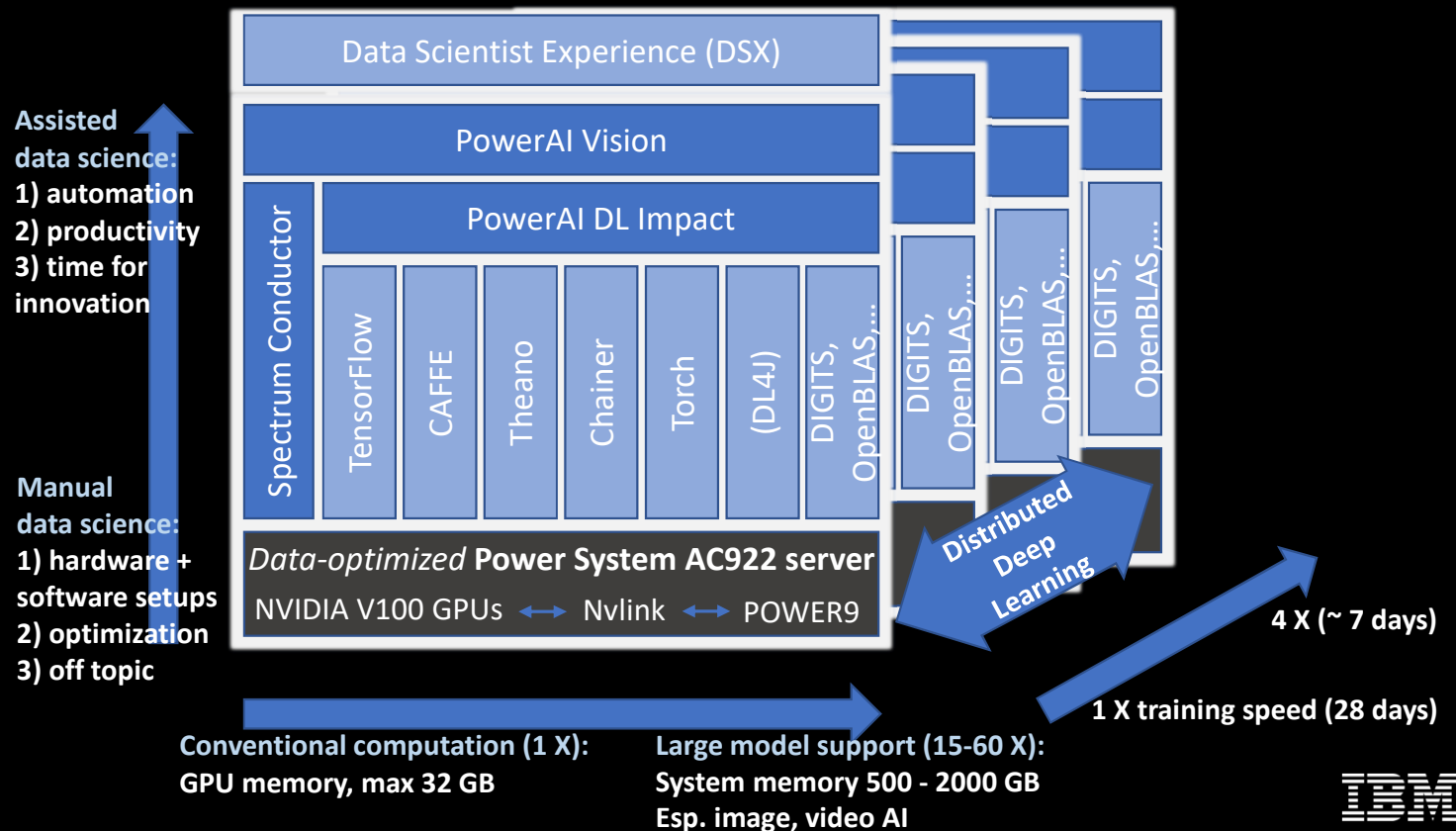
Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

Options for deployment

- a) Self-operated system or cluster suitable for IT + user departments
- b) Services from IBM business partner ecosystem and networks

PowerAI package & tools – 3 axes of differentiation and productivity



Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect



Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

SW stack - use the level you need:

- 1) Full data science environment with GUI (DSX) (and/or cluster management IBM Cloud Private local installation)
- 2) Dedicated computer vision solution for low threshold DL adoption (PowerAI Vision)
- 3) DL experiment automation, e.g. hyperparameter optimization (PowerAI DL Impact)

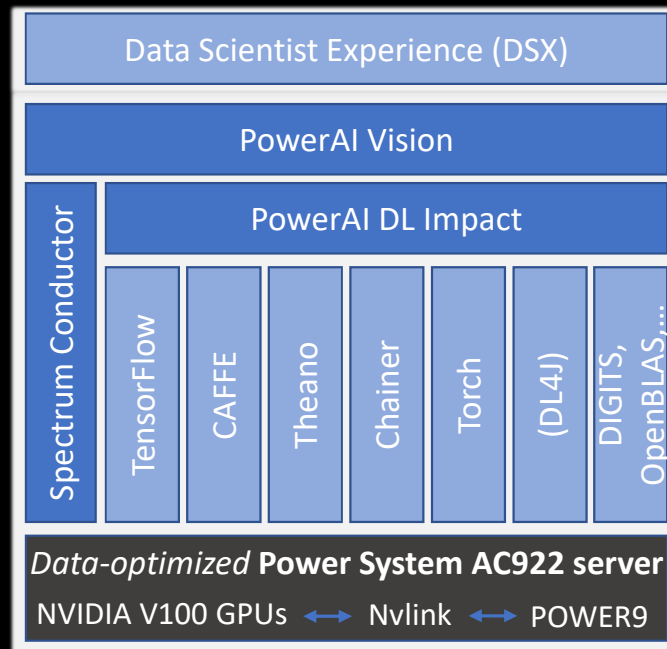
and/or

- 4) Conventional use of deep learning frameworks (TF, Caffe, Torch, Keras etc.) out-of-box (e.g. Linux command line, Python, tools of your choosing)

Assisted data science:
1) automation
2) productivity
3) time for innovation

Manual data science:
1) hardware + software setups
2) optimization
3) off topic

PowerAI package & tools – 1st axis: different options for usage



Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect



Deep Learning Impact (DLI), WML Accelerator (including model deployment)

New Dataset

Create a dataset from:

- LMDBs
- TensorFlow Records
- Images for Object Classification
- Images for Object Detection
- Images for Vector Output
- CSV
- Other

Trained Model Name

Trained Model Name	State	Dataset
Cifar10-20171029143150	Finished	Cifar10-datas
Cifar10-20171021133946	Finished	Cifar10-datas
Cifar10-20171021133732	Finished	Cifar10-datas
Cifar10-20171021133155	Finished	Cifar10-datas
Cifar10-20171021132644	Finished	Cifar10-datas
Cifar10-20171021120157	Failed	Cifar10-datas

Training Insights: Cifar10-20171029143150

Spark application used for training: app-20171030023251-0017-eeee0a87-422a-4440-bef2-...

Number of Workers: 1
CPUs per Worker: 1
Framework: CaffeOnSpark
Max Iteration: 3000

Time and Iteration

Scale: 5m30s
Number of Iterations: 2950

Loss

Loss: 2.4, 1.6, 0.8, 0.4

Accuracy

Accuracy: 0.9, 0.7, 0.5, 0.3, 0.1

Weight histogram

Weight: 1.5, 1.0, 0.5, 0.0, -0.5, -1.0

Tune Hyperparameters for model:Cifar10

A tuning task will launch multiple jobs to search hyperparameters. A new model will be created that contains the tuned hyperparameters.

* Name of tuning attempt: Cifar10- tuning-20171029185346
① The new model will be named as:Cifar10-tuning-20171029185346

* Hyperparameter search type: Random Search

* Framework: CaffeOnSpark

Tuning Parameter Settings

Input the parameters that will be tuned

* Optimizer (select at least 1):
 SGD
 AdaDelta
 AdaGrad
 Adam
 Nesterov
 RMSProp

* Learning rate range: []

* Weight decay range: []

* Momentum range: []

* Max batch size: []

Tuning Workload Setting

Input the parameters to control tuning jobs and process

* Number of workers: 1

GPUs per worker: 1

* Max iterations: []

* Total tuning jobs number: 10

* Max tuning jobs in parallel: 2

* Max running time(minutes): 60

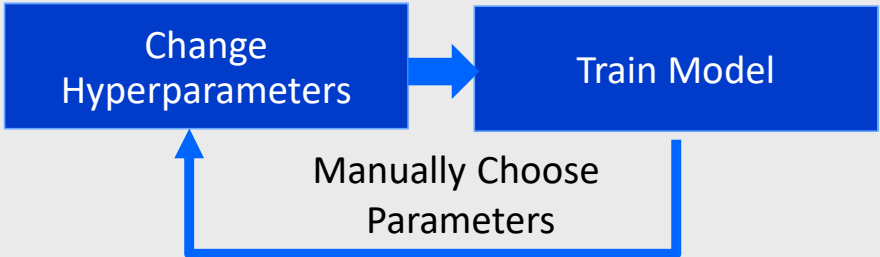
Start Tuning **Cancel**

<http://www.redbooks.ibm.com/abstracts/sg248409.html?Open>

Auto Hyper-Parameter Optimization (HPO) in WML Accelerator / PowerAI Enterprise

Run Model Training 100s of Times

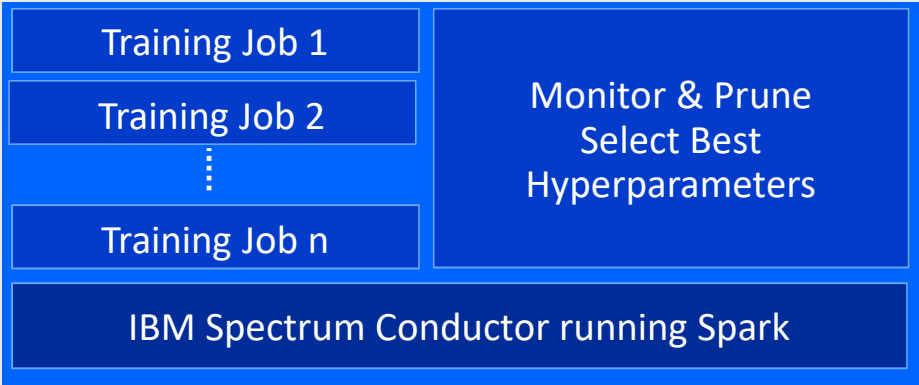
Manual Process
Can take weeks



Lots of Hyperparameters:

Learning rate, Decay rate, Batch size, Optimizers (Gradient Descent, Momentum, ..)

Auto-Hyperparameter Optimizer (Auto-HPO)
Done in Hours



Auto-HPO has 3 search approaches

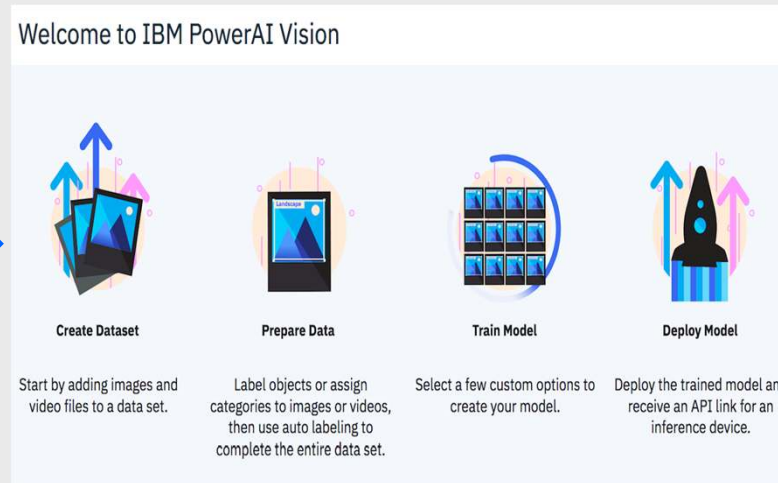
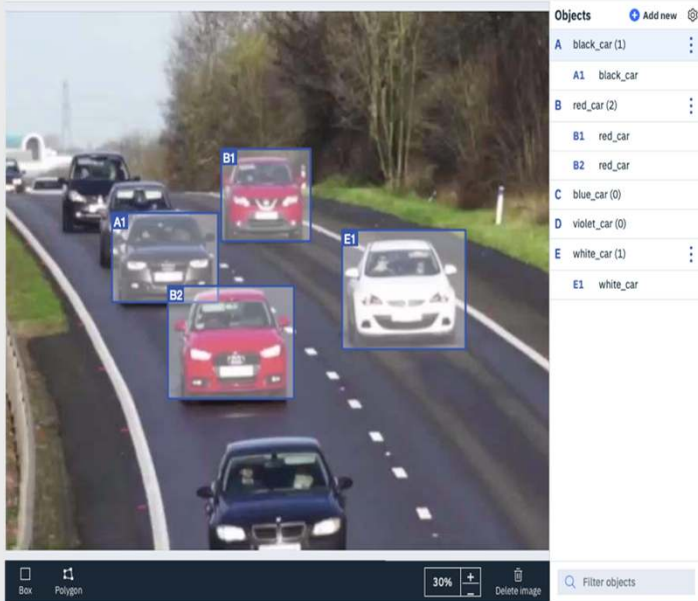
Random, Tree-based Parzen Estimator (TPE), Bayesian

PowerAI Vision: “Point-and-Click” AI for Images & Video

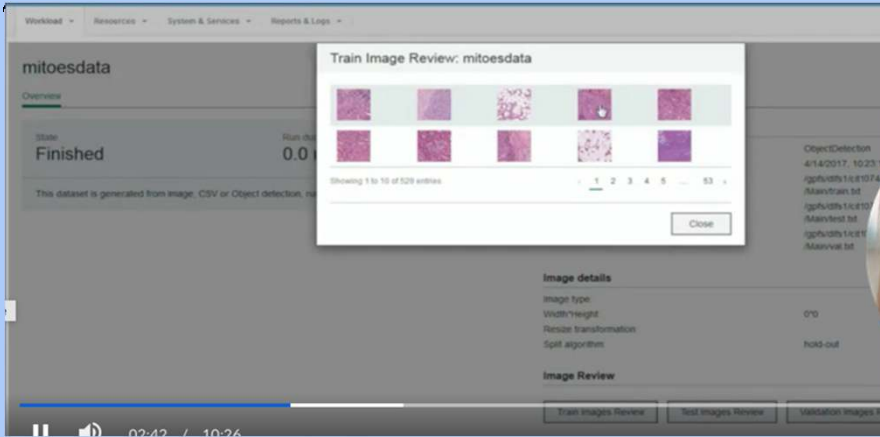
Label Image or Video Data

Auto-Train AI Model

Package & Deploy AI Model



... hiding all together the machine learning details
(getting greater group involved in AI development)



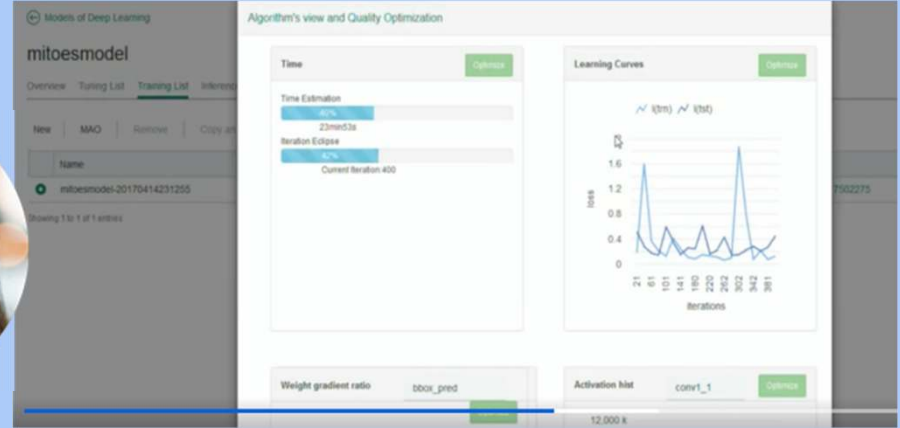
Subject matter experts



Application developers

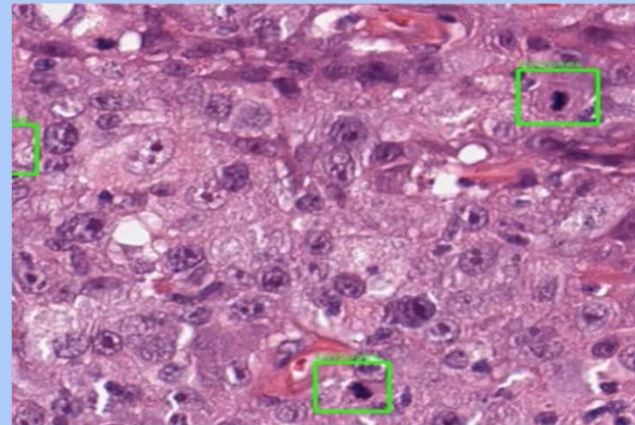


Data scientists



View Inference:

01_11.jpg	mitoses	99.5%	(1339,1061)(1427,1147)	View Result Image
	mitoses	99.3%	(1135,873)(1219,957)	
	mitoses	99.0%	(1574,861)(1659,945)	
01_21.jpg	mitoses	92.2%	(1737,270)(1830,359)	View Result Image
04_28.jpg	mitoses	99.0%	(369,533)(454,617)	View Result Image
	mitoses	98.8%	(1445,1355)(1530,1438)	
	mitoses	98.5%	(908,563)(992,646)	
	mitoses	96.5%	(419,197)(504,279)	
	mitoses	96.4%	(981,842)(1074,932)	
	mitoses	84.1%	(710,712)(800,800)	
04_18.jpg	mitoses	99.6%	(339,1707)(423,1791)	View Result Image
	mitoses	97.2%	(83,1597)(170,1683)	
	mitoses	87.1%	(936,1599)(1010,1973)	

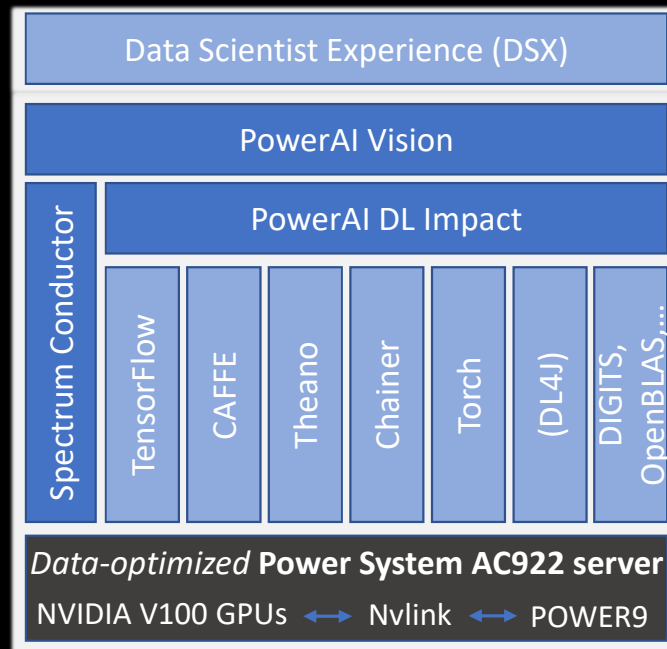


Tumor Proliferation Assessment – mitosis detection
Images from electron-microscope
Size of image - 70K * 60K

Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

PowerAI package & tools – 2nd axis: go beyond GPU memory limits



Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect

Conventional computation (1 X):
GPU memory, max 32 GB

Large model support (15-60 X):
System memory 500 - 2000 GB
Esp. image, video AI

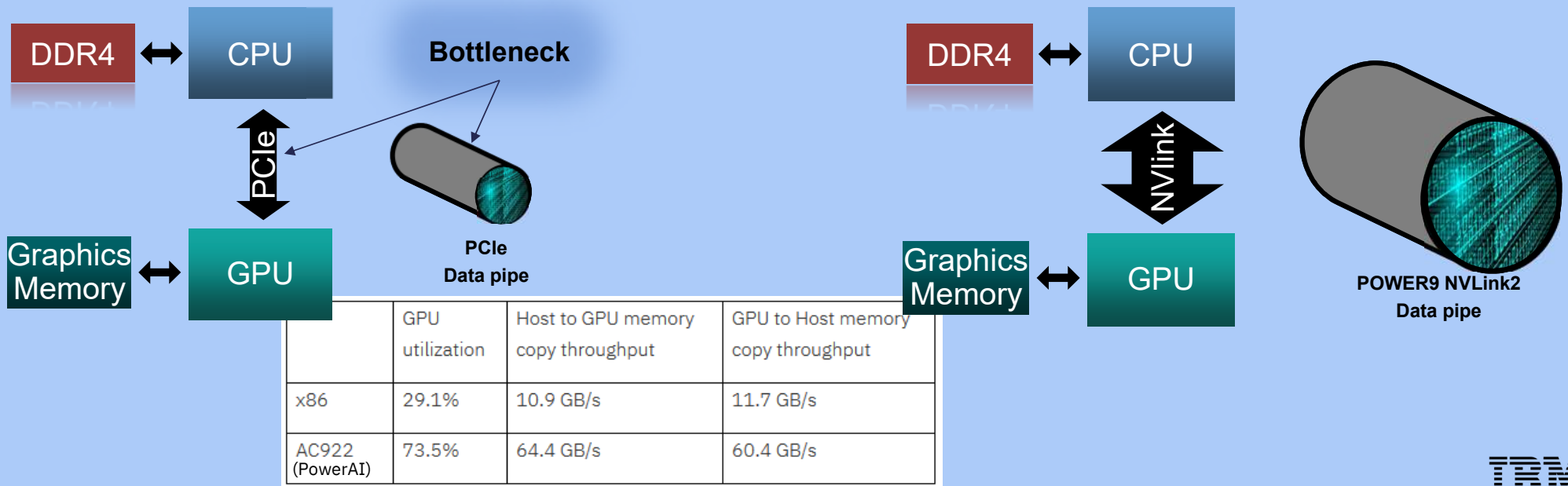


Large Model Support (LMS) on PowerAI (only) 1/2

GPU memory (max 8-16/32 GB) has to host active experiment (including neural net parameters) as well training data

System main memory much larger in server systems (128 GB – 2000 GB)

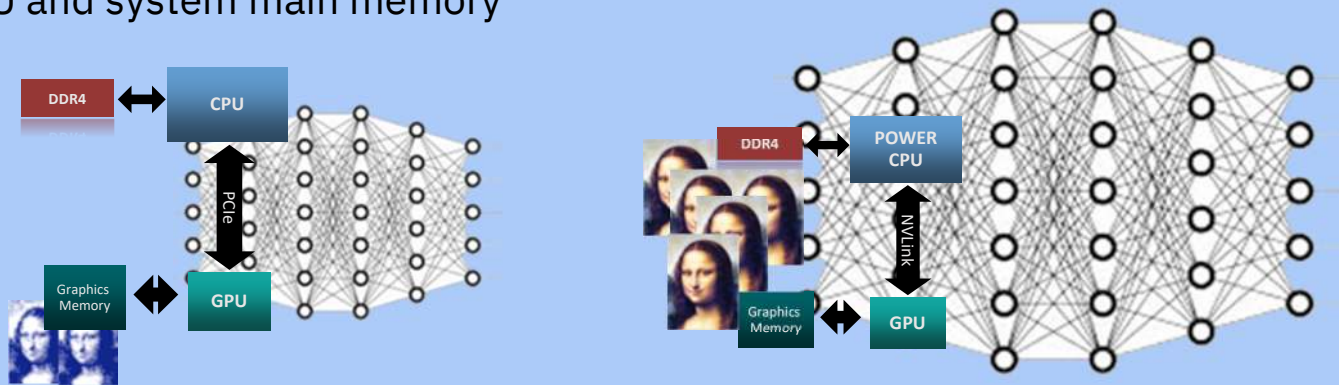
NVLINK provides both fast GPU interconnection and in Power servers also fast connection between GPUs and system main memory - IO-advantage in itself for feeding training data to GPUs



Large Model Support (LMS) on PowerAI (only) 2/2

“Large Model Support” (LMS) – larger experiments than GPU memory:

- Caffe: swap chunks of data between GPU and system main memory
- Tensorflow: computation graph rewriting in order to swap results between GPU and system main memory



<https://github.com/ibmsoe/caffe/tree/master-lms>

<https://developer.ibm.com/linuxonpower/2018/07/27/tensorflow-large-model-support-case-study-3d-image-segmentation/>

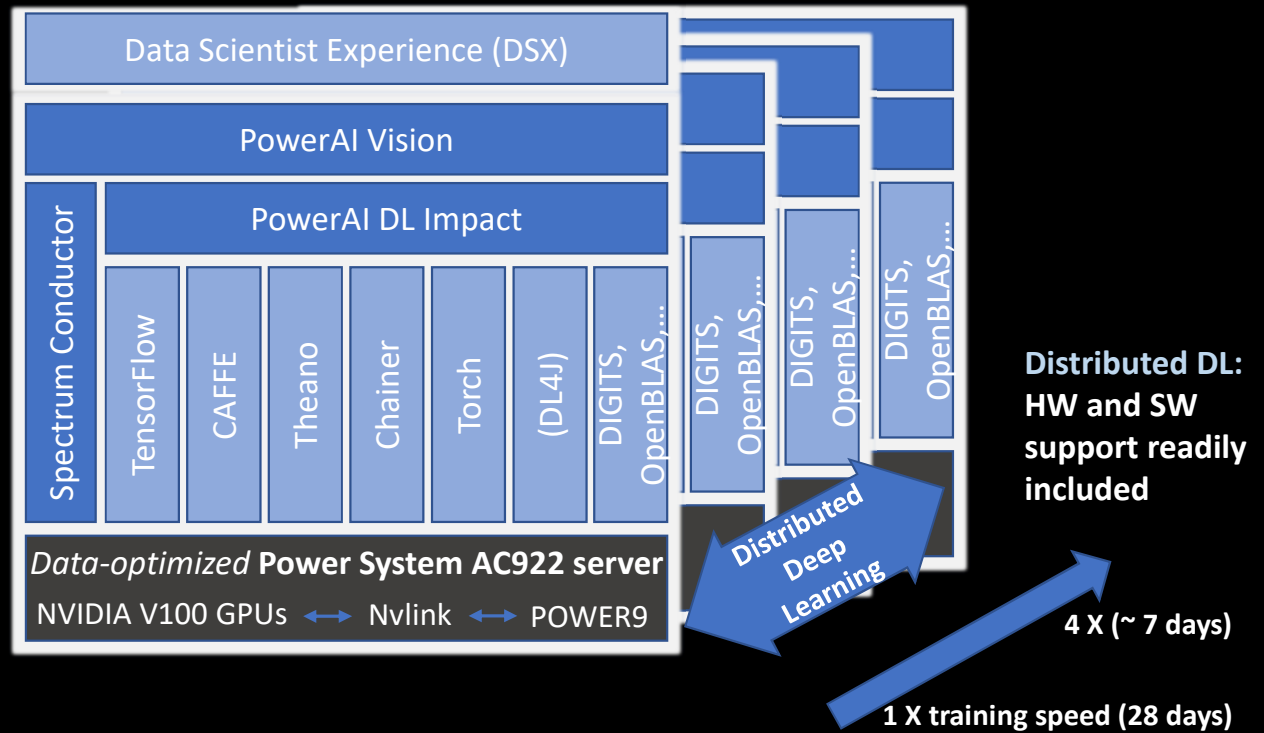
<https://arxiv.org/abs/1807.02037>

Technology preview of LMS Tensorflow support on PowerAI: /opt/DL/tensorflow/doc/README-LMS.md

Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

PowerAI package & tools – 3rd axis: easily divide workloads across cluster



Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect



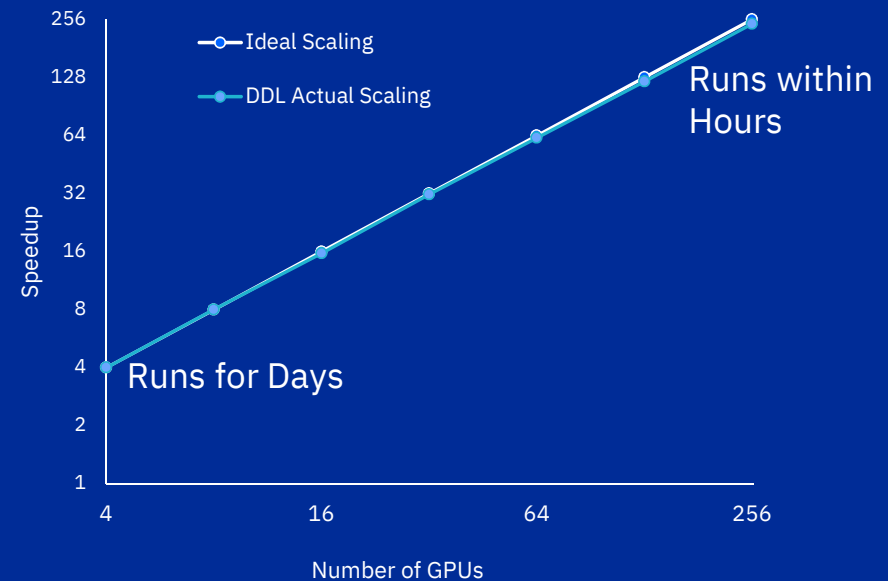
Distributed Deep Learning (DDL)

Deep learning training takes days to weeks

DDL in WML CE extends TensorFlow & enables scaling to 100s of servers

Automatically distribute and train on large datasets to 100s of GPUs

Near Ideal (95%) Scaling to 256 GPUs



ResNet-50, ImageNet-1K
Caffe with PowerAI DDL,
Running on S822LC Power System

Elastic Distributed Training (EDT) in WML Accelerator / PowerAI Enterprise

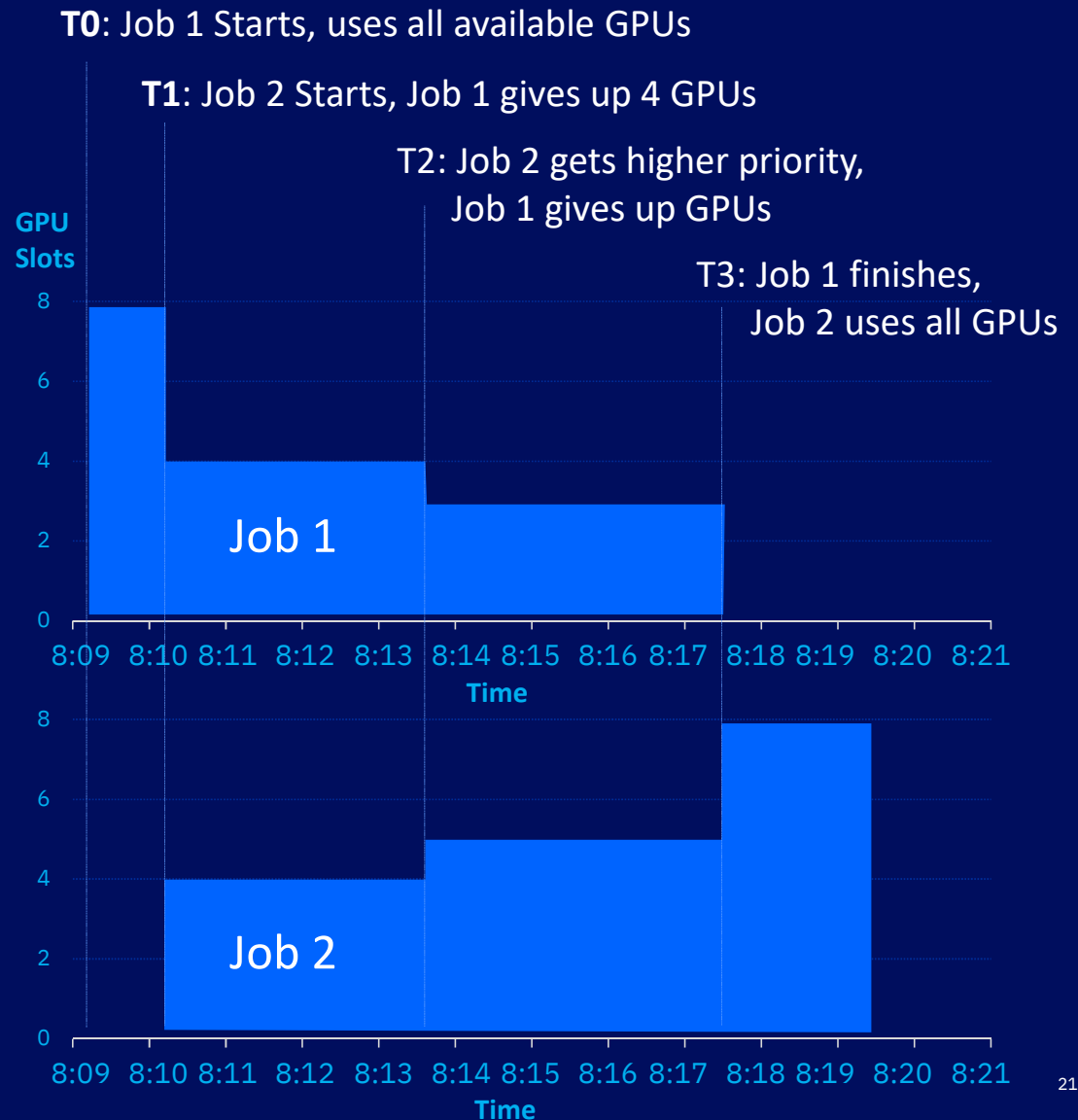
Dynamically Reallocates GPUs within
milliseconds

Increases Job Throughput and Server / GPU
Utilization

Works with Spark & AI Jobs

Works with Hybrid x86 & Power Cluster

2 Servers with 4 GPUs each: total 8 GPUs
Available Policies: Fair share, Preemption, Priority



Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

Options for deployment

- a) Self-operated system or cluster suitable for IT + user departments
- b) Services from IBM business partner ecosystem and networks

SW stack - use the level you need:

- 1) Full data science environment with GUI (DSX) (and/or cluster management IBM Cloud Private local installation)
- 2) Dedicated computer vision solution for low threshold DL adoption (PowerAI Vision)
- 3) DL experiment automation, e.g. hyperparameter optimization (PowerAI DL Impact)

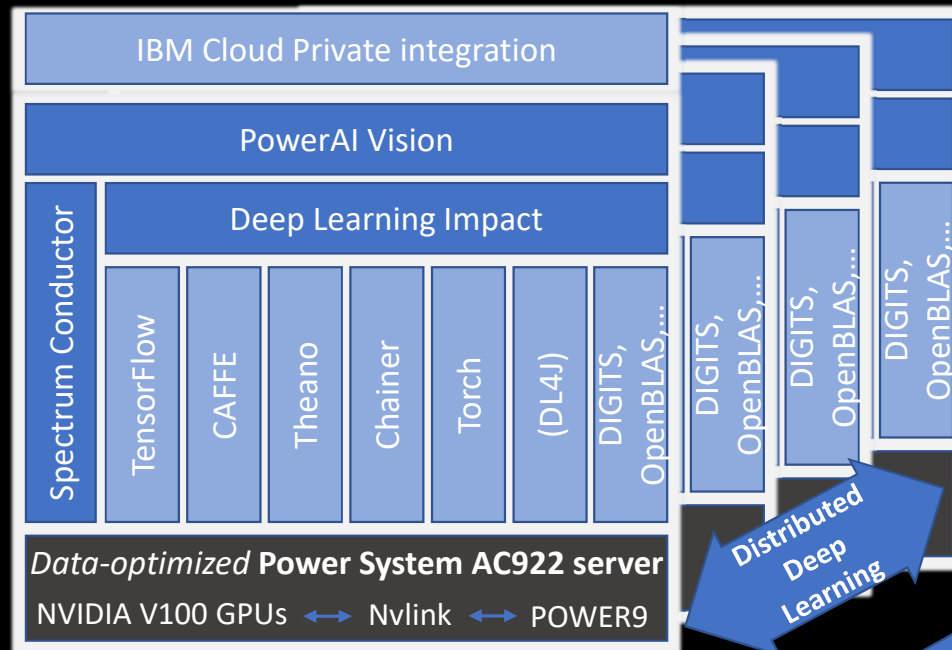
and/or

- 4) Conventional use of deep learning frameworks (TF, Caffe, Torch, Keras etc.) out-of-box (e.g. Linux command line, Python, tools of your choosing)

Assisted data science:
1) automation
2) productivity
3) time for innovation

Manual data science:
1) hardware + software setups
2) optimization
3) off topic

PowerAI package & tools – 3 axes of differentiation and productivity



Distributed DL:
HW and SW support readily included

Distributed Deep Learning

4 X (~ 7 days)

1 X training speed (28 days)

Conventional computation (1 X):
GPU memory, max 32 GB

Large model support (15-60 X):
System memory 500 - 2000 GB
Esp. image, video AI

Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect



IBM Open Source Based AI Stack

Auto-AI software: PowerAI Vision, IBM Auto-AI

Watson Studio

WML CE

Data Preparation
Model Development
Environment



Watson Machine Learning

Watson ML Accelerator

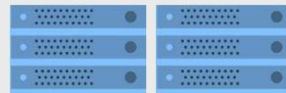
Watson ML CE

Runtime Environment
Train, Deploy, Manage Models



Watson OpenScale

Model Metrics, Bias,
and Fairness
Monitoring



Accelerated AC922
Power9 Servers



Storage
(Spectrum Scale ESS)

Previous Names:

WML Accelerator = PowerAI Enterprise
WML Community Ed. = PowerAI-base

Runs on x86 & other storage too

Available on Private Cloud or Public Cloud

Our Focus: Ease of Use & Faster Model Training Times

Watson ML Accelerator

Distributed Deep Learning (DDL)

Auto Hyper-Parameter Optimization (HPO)

Elastic Distributed Training (EDT) & Elastic Distributed Inference (EDI)

IBM Spectrum Conductor

Apache Spark, Cluster Virtualization, Job Orchestration

Watson ML

Model Management & Execution

Model Life Cycle Management

Watson ML Community Edition
WML CE

WML CE: Open Source ML Frameworks



PYTORCH



Chainer

Snap ML

Large Model Support (LMS)

DDL-16

Infrastructure Designed for AI



Power9 or x86 Servers with GPU Accelerators



Storage (ESS)

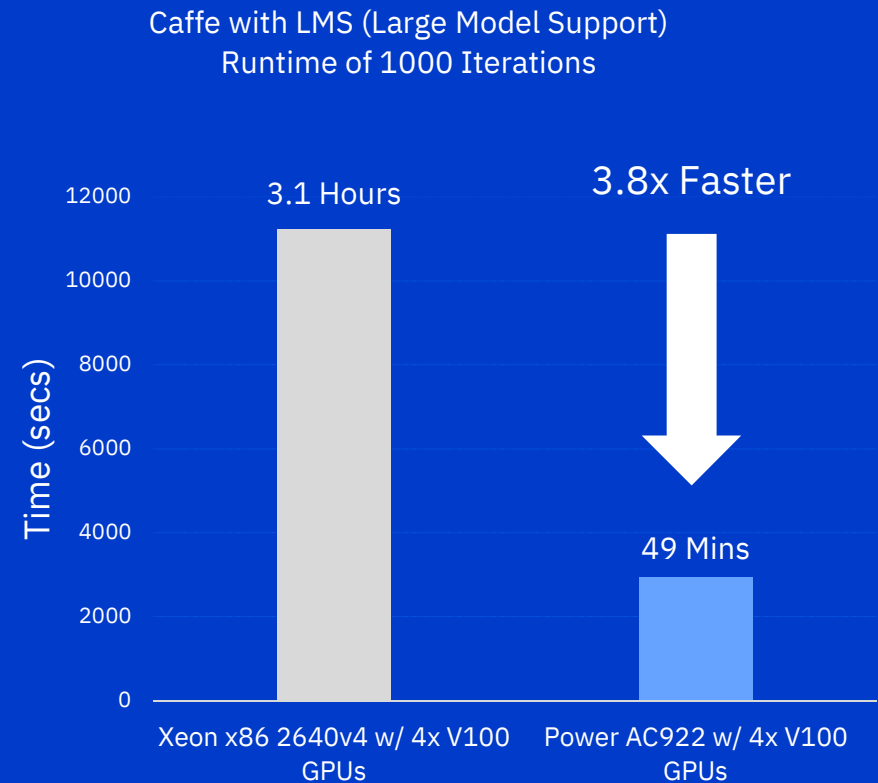
Large AI Models Train ~4 Times Faster

-

Time-to-Results Drops from Months to Weeks

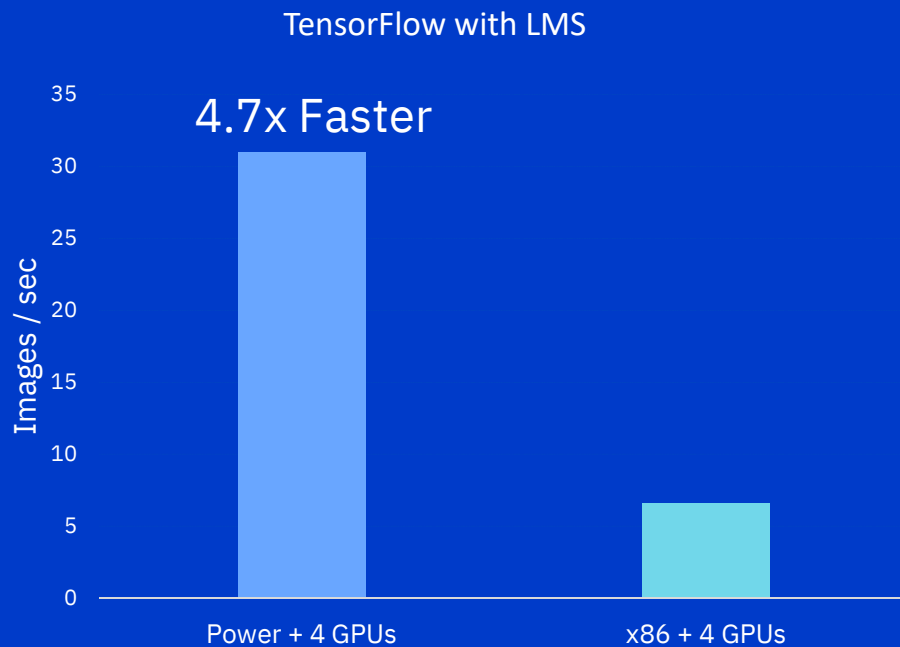
IBM POWER9 Servers with NVLink to GPUs
VS
x86 Servers with PCIe to GPUs

Detailed Benchmark Information is available



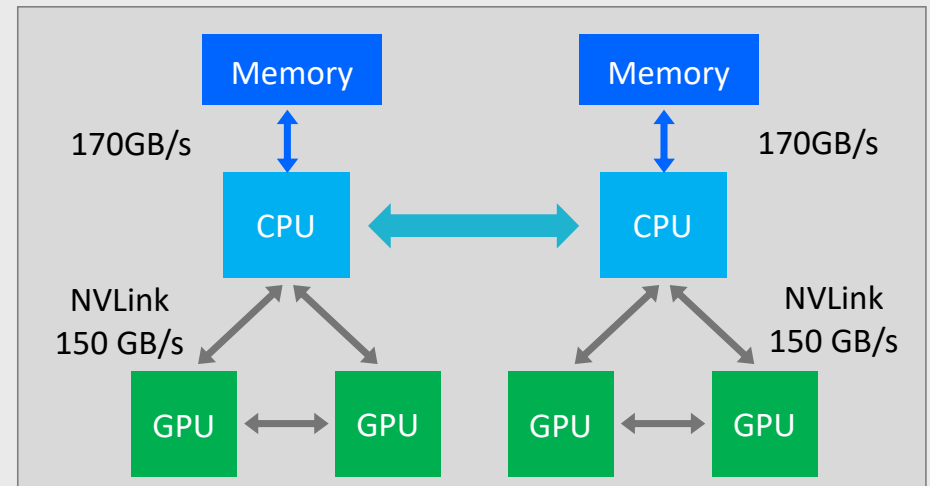
GoogleNet model on Enlarged
ImageNet Dataset (2240x2240)

Large Model Support (LMS) Enables Higher Accuracy via Larger Models



500 Iterations of Enlarged GoogleNet model on Enlarged ImageNet Dataset (2240x2240), mini-batch size = 15
Both servers with 4 NVIDIA V100 GPUs

Store Large Models & Dataset in System Memory Transfer One Layer at a Time to GPU



IBM AC922 Power9 Server
CPU-GPU NVLink 5x Faster than Intel x86 PCI-Gen3

THE US AGAIN HAS THE WORLD'S MOST POWERFUL SUPERCOMPUTER



The IBM-built Summit supercomputer is the world's smartest and most powerful AI machine. Its racks are connected by over 185 miles of fiber-optic cables.

GENEVEVE MARTIN/DAK RIDGE NATIONAL LABORATORY

<https://www.wired.com/story/the-us-again-has-worlds-most-powerful-supercomputer/>



The Future of Accelerated Computing

Power System



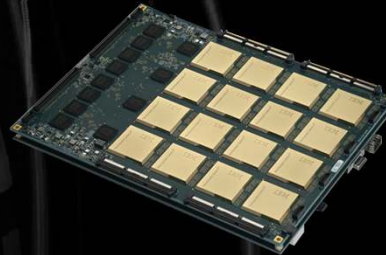
IBM PowerAI



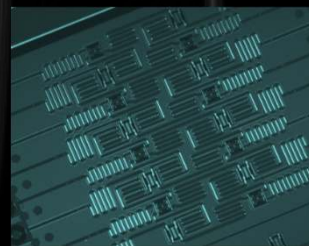
Today:
Graphic Processors
NVIDIA Volta



New:
FPGA
Xilinx Alveo



Tomorrow:
Neuromorphic Processors
IBM TrueNorth

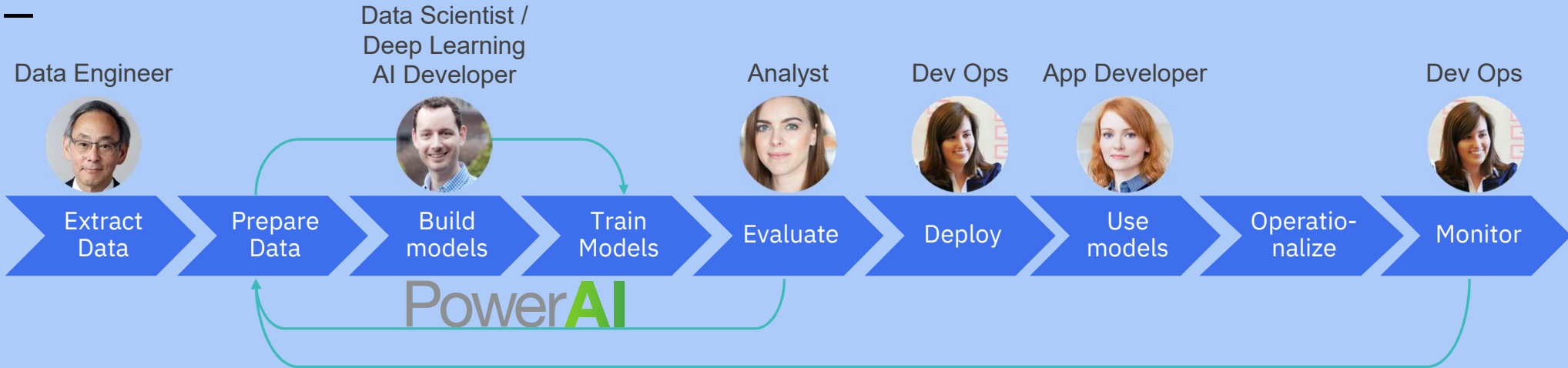


Future:
Quantum Computer
IBM Q

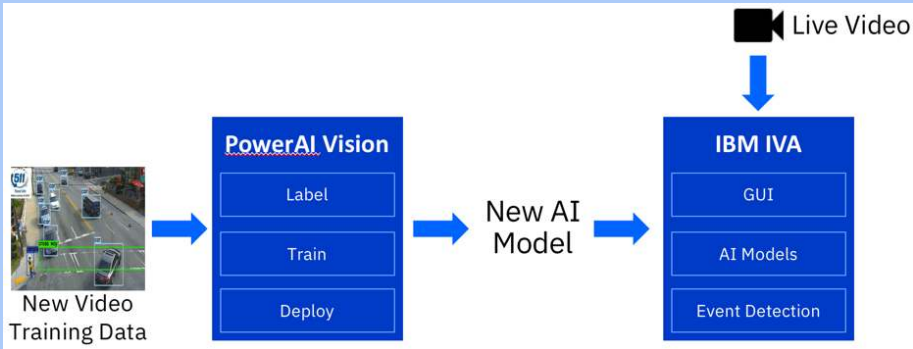
OpenPOWER™



Systematic data science: deployment, integration and continuous AI development/model improvement



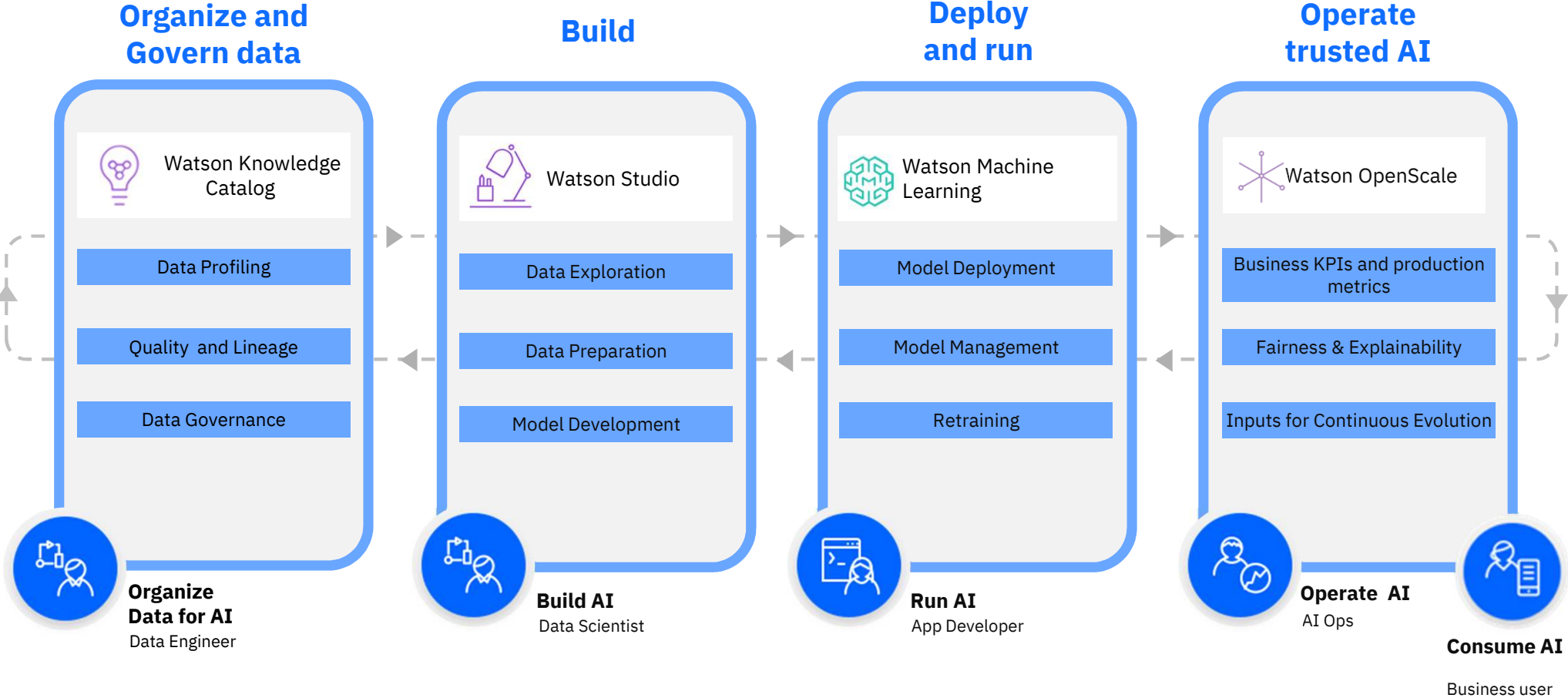
Interested parties need to be able to consume the trained model and the model needs to be supported



https://medium.com/@sumitg_16893/deep-insights-with-ai-for-video-analytics-5464fd30ebe1



Watson tool portfolio addresses all the stages of AI development and utilization



Thank you!

Questions?

Multiplying data science with IBM PowerAI - faster value creation from deep learning AI

"Data scientists don't train a model - they train thousands!"

Options for deployment

- a) Self-operated system or cluster suitable for IT + user departments
- b) Services from IBM business partner ecosystem and networks

SW stack - use the level you need:

- 1) Full data science environment with GUI (DSX) (and/or cluster management IBM Cloud Private local installation)
- 2) Dedicated computer vision solution for low threshold DL adoption (PowerAI Vision)
- 3) DL experiment automation, e.g. hyperparameter optimization (PowerAI DL Impact)

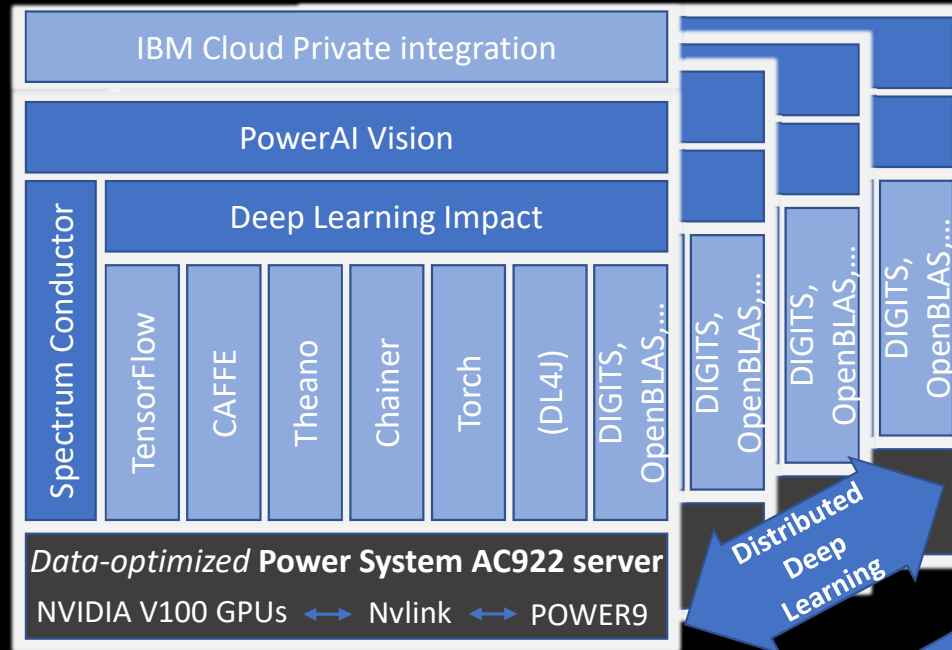
and/or

- 4) Conventional use of deep learning frameworks (TF, Caffe, Torch, Keras etc.) out-of-box (e.g. Linux command line, Python, tools of your choosing)

Assisted data science:
1) automation
2) productivity
3) time for innovation

Manual data science:
1) hardware + software setups
2) optimization
3) off topic

PowerAI package & tools – 3 axes of differentiation and productivity



Distributed DL:
HW and SW support readily included

4 X (~ 7 days)

1 X training speed (28 days)

Conventional computation (1 X):
GPU memory, max 32 GB

Large model support (15-60 X):
System memory 500 - 2000 GB
Esp. image, video AI

Jukka Remes, BDE, Dr (Tech.)
Teppo Seesto, Solution Architect

