



Can HPC Technologies Lead the way for AI?

Adnan Khaleel

Global Sales Strategy for HPC & AI

Adnan.Khaleel@dell.com

Sept 2018

AI customers demonstrate the value



MasterCard

Turn 160M transactions/
hour with 1.9M rules to
protect against fraud



Caterpillar

Autonomous
mining for safety



Zenuity

4.4Pb of data/mo
~50 simulations/ hr
safe autonomous driving



MIT Lincoln Labs

1+ petaFLOP to enable
robotic vehicles, cyber
security, bioinformatics

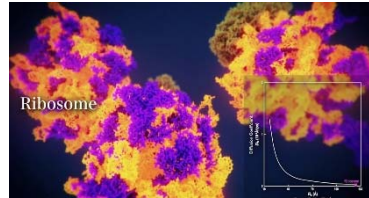


Develop, test and run
innovative algorithms



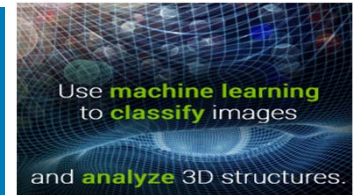
ZIFF.ai

AI startup accelerates
training models with 10s
of millions of images



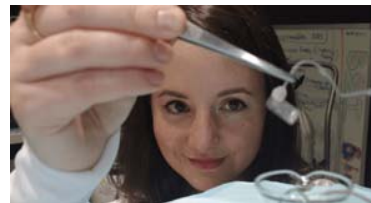
Simon Fraser University

Analyze DNA of microbes
and stop outbreaks



CSIRO

More than 1,800
patents in science
and technology



AeroFarms

390x more productive
than traditional field
farming



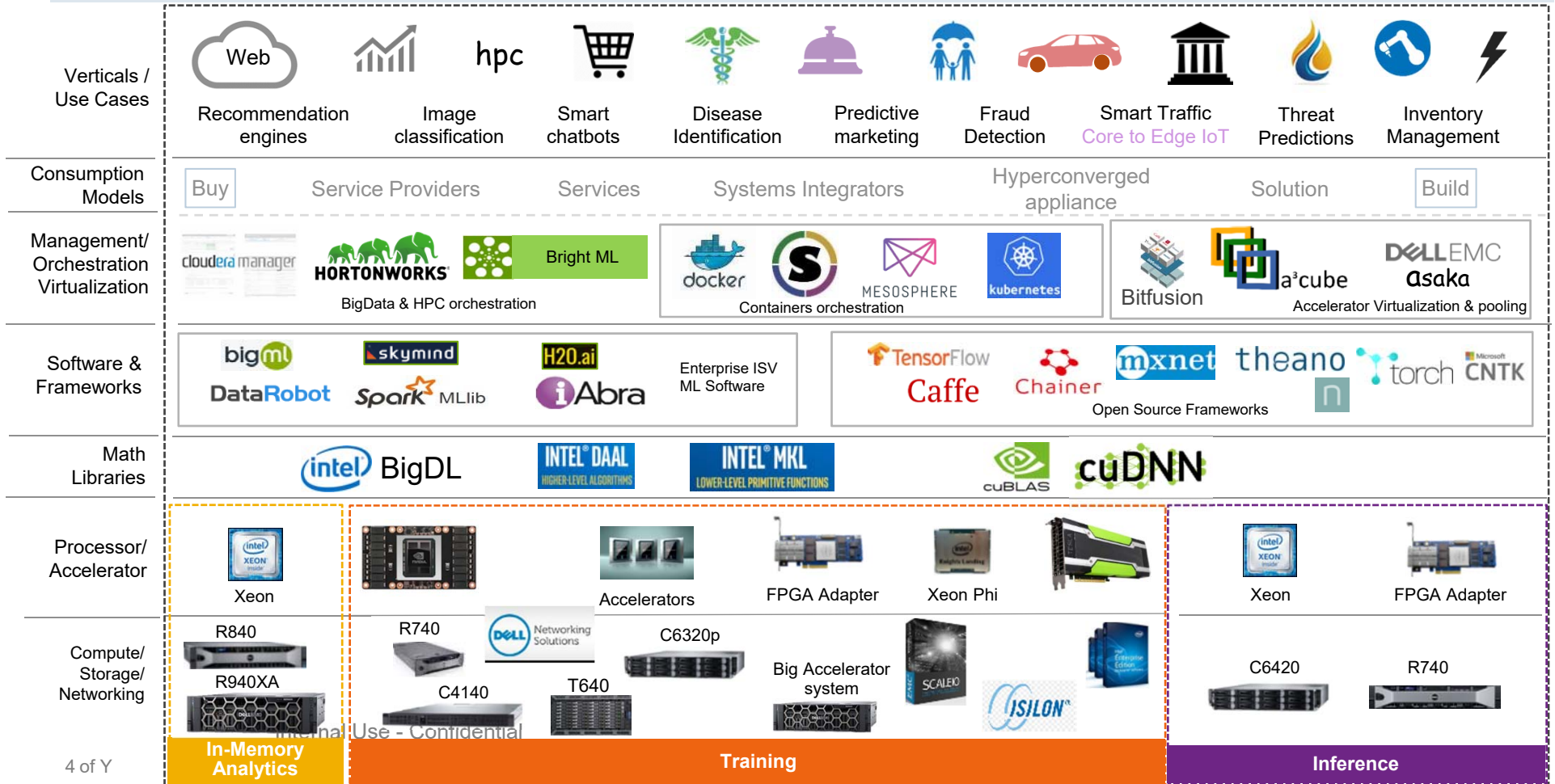
CU

300% faster
performance

Motivations for this discussion

- AI is still in its early days as computational technology; lots of dynamism
- Advise CIOs/CTOs and our PoV of how we see the technology evolving
 - Also from a future product planning perspective
 - What can we make an educated guess
- Maximize ROI on current investments
- What will be the right way to “scale” this problem size and or speed?
- It’s easier to come up with an innovative HW architecture than to “actually” program it
- What can we apply from HPC lessons learnt to this domain?
 - Interconnects, memory, scale-out within nodes vs across nodes
 - Are nodes with 8, 16, 32 GPUs useful and how far do we go?

Machine Learning & Deep Learning eco-system – solving real world problems



Internal Use - Confidential

Dell EMC Ready Solutions for Machine & Deep Learning – Clearly Scale-out

AI, ML and DL Services



Deep Learning with NVIDIA

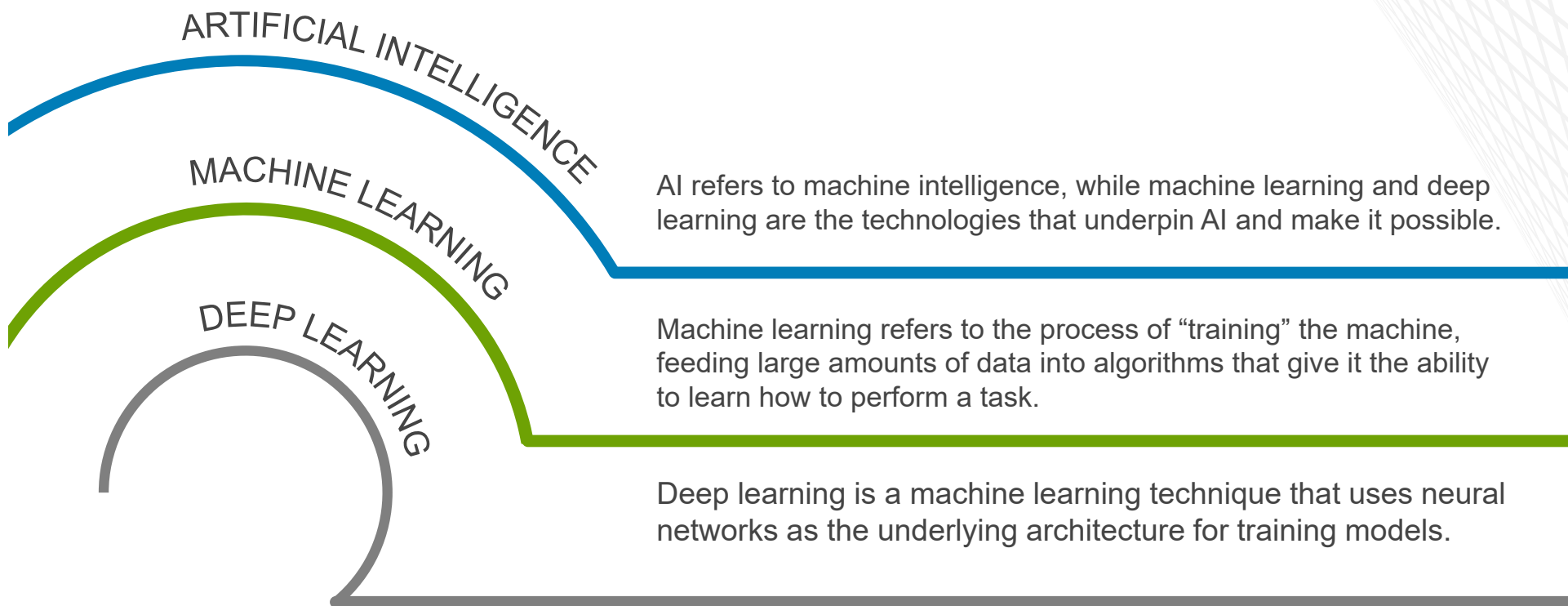


Deep Learning with Intel



Machine Learning with Hadoop

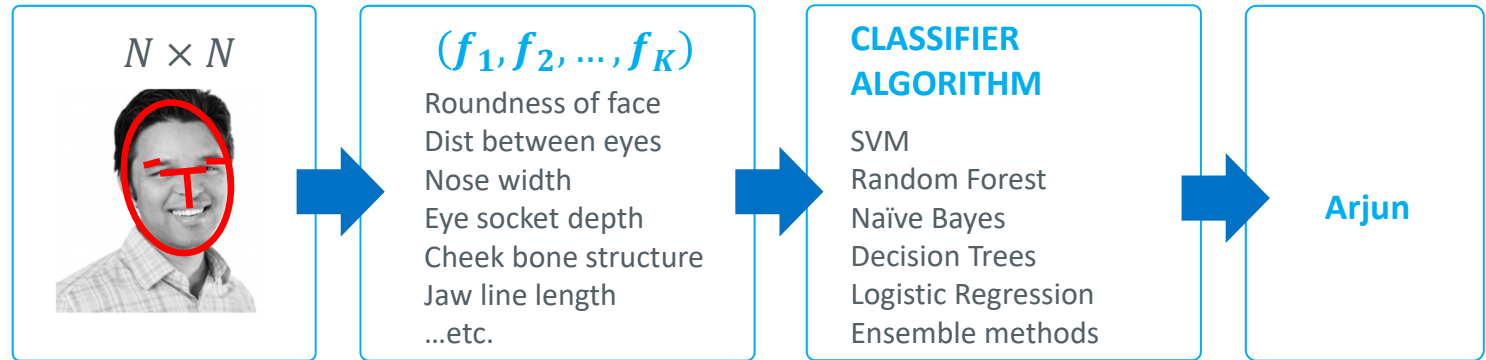
What is AI, machine learning and deep learning?



Quick Primer

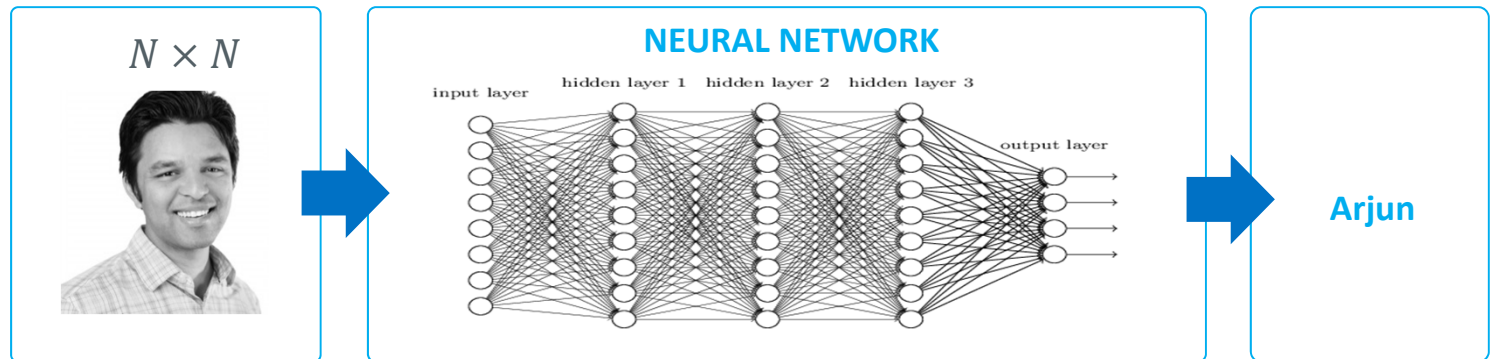
Machine Learning

How do you engineer the best features?



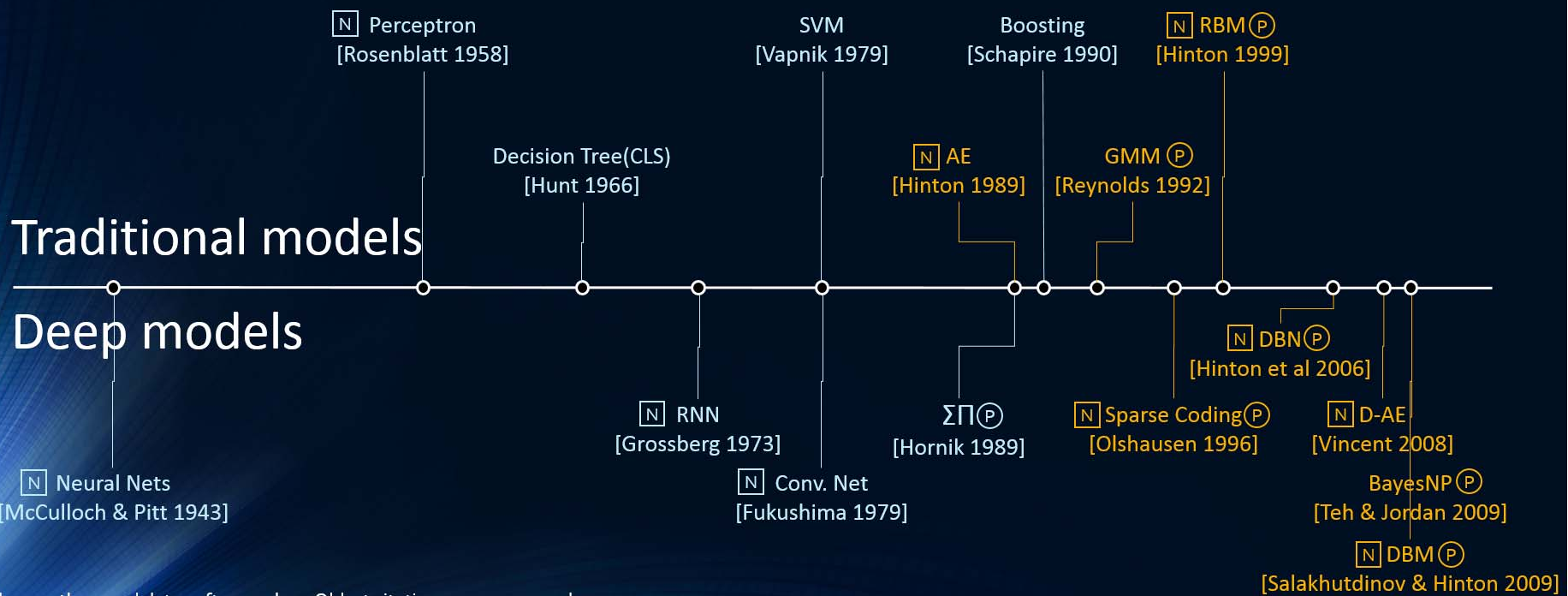
Deep Learning

How do you guide the model to find the best features?



Deep Learning evolution

- N Neural Network
- P Probabilistic Model
- Supervised learning
- Unsupervised learning



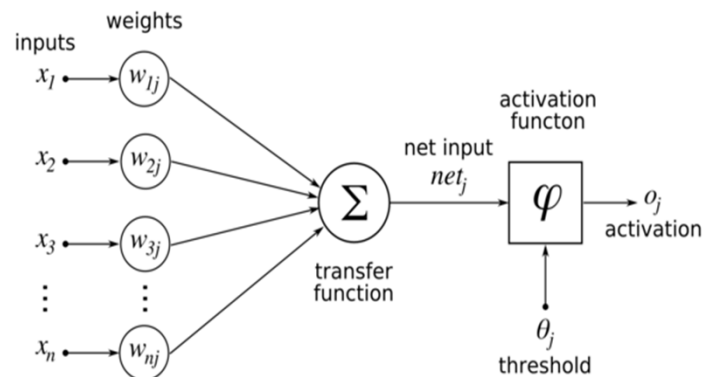
Algorithms authors and dates often unclear. Oldest citations were assumed
Classifications based on Yann LeCun's Deep Learning class at NYU – spring 2014.

Biologically Inspired Artificial Neurons



Neurons and biological brains

- Each Human has a 100 Billion neurons. Each Neuron connects to 10k other neurons ~ 1 Trillion Synapses
- Neurons send and receive signals. A neuron “fires” if the summation of its input is above a threshold



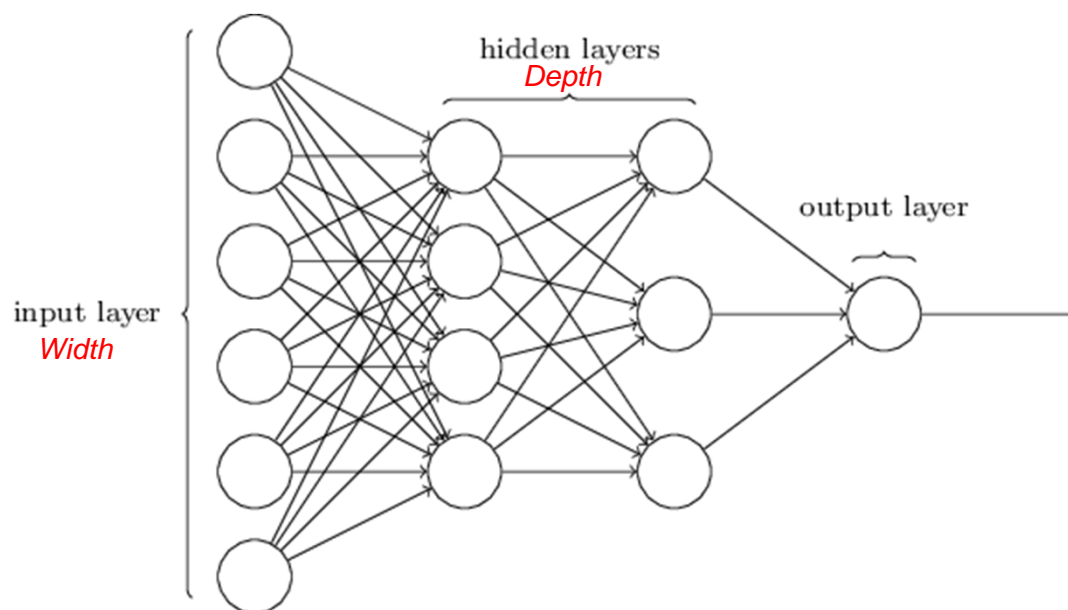
Simple Artificial Neuron or “Switch” or (Single Layer) Perceptron

- Has a number of “weighted” inputs
- Considers sum of weighted inputs and “fires” an output based on it’s activation function

Multi Layer Perceptrons or Artificial Neural Networks (ANNs)

Group together multiple Artificial Neurons into a Network

- Single Layer Perceptrons can be used to represent linearly separable functions
- However many problems are not linearly separable, may be represented as 2D or 3D spaces (*Lippmann in the 1987 paper “An introduction to computing with neural nets”*)
- *Generally arrive at these parameters through experimentation*
- *Architecture is based on the layout of various layers E.g. AlexNet, ImageNet*

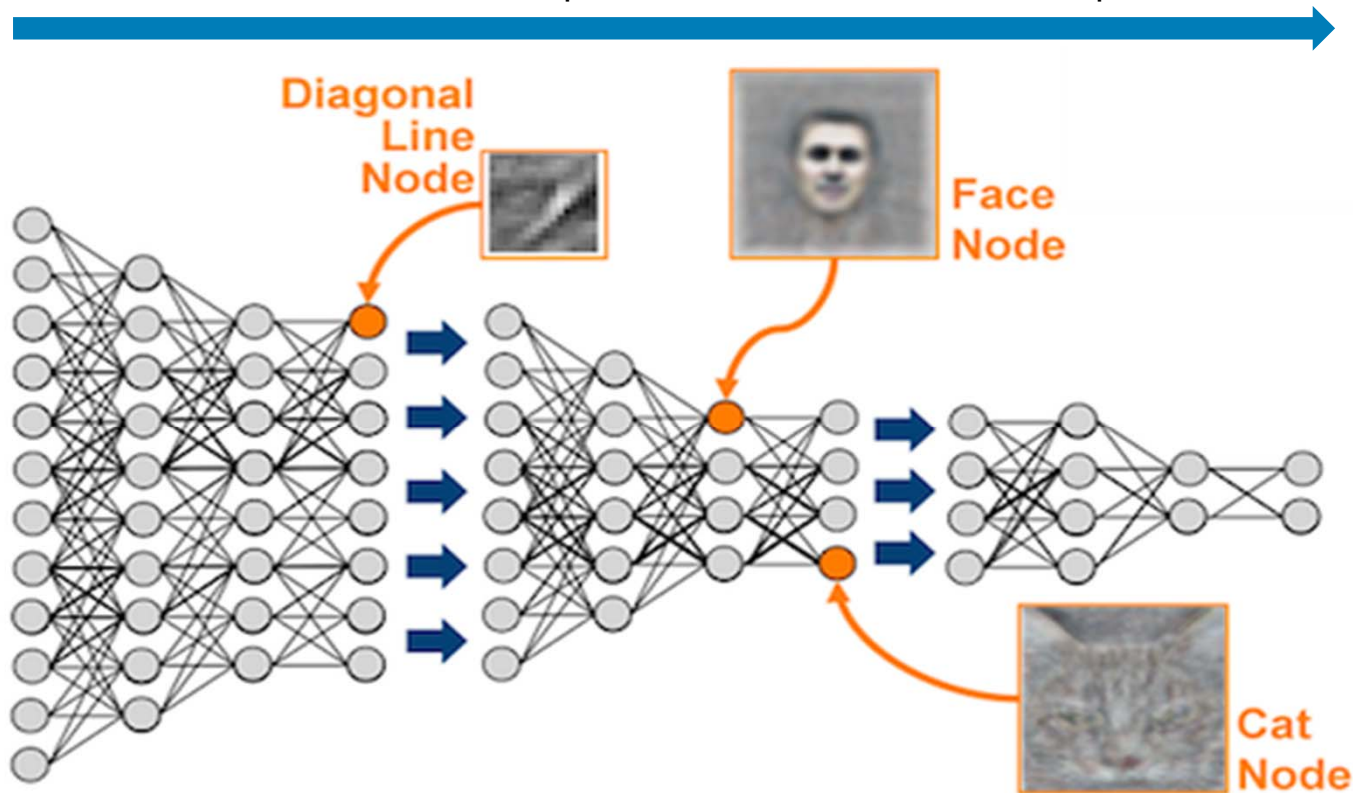


Deep Neural Network – How do they work?

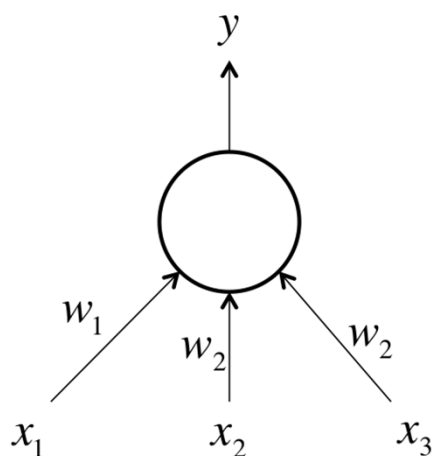
Basic Attributes

Compound Attributes

Complex Attributes



A Simple Example: Supervised Learning



$$y^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + w_3 x_3^{(i)}$$

$$z = \sum_k w_k x_k$$

$$y = \frac{1}{1 + e^{-z}}$$

$$z^{(1)} = w_1 x_1^{(1)} + w_2 x_2^{(1)} + w_3 x_3^{(1)}$$

$$z^{(2)} = w_1 x_1^{(2)} + w_2 x_2^{(2)} + w_3 x_3^{(2)}$$

$$z^{(3)} = w_1 x_1^{(3)} + w_2 x_2^{(3)} + w_3 x_3^{(3)}$$

Error criteria: Minimize the square error over all of the training examples

$$\text{Minimise } E = \frac{1}{2} \sum_i (t^{(i)} - y^{(i)})^2$$

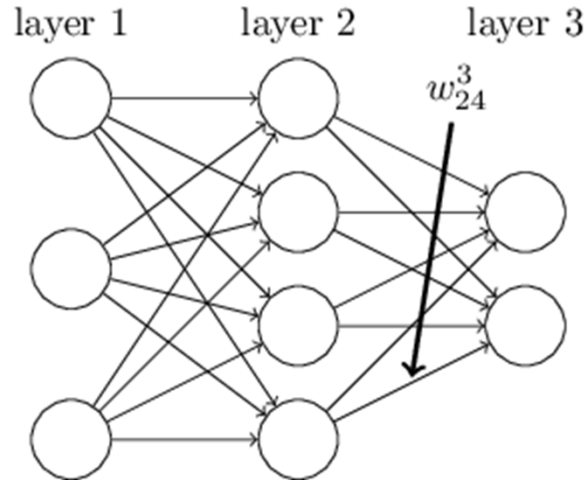
Where $t^{(i)}$ is the true answer for the i^{th} training example

To modify the weights:

$$\Delta w_k = \sum_i \epsilon x_k^{(i)} y^{(i)} (1 - y^{(i)}) (t^{(i)} - y^{(i)})$$

Training and Backpropagation (Supervised learning)

The backpropagation algorithm was originally introduced in the 1970s, but its importance wasn't fully appreciated until a famous 1986 paper by David Rumelhart, Geoffrey Hinton, and Ronald Williams

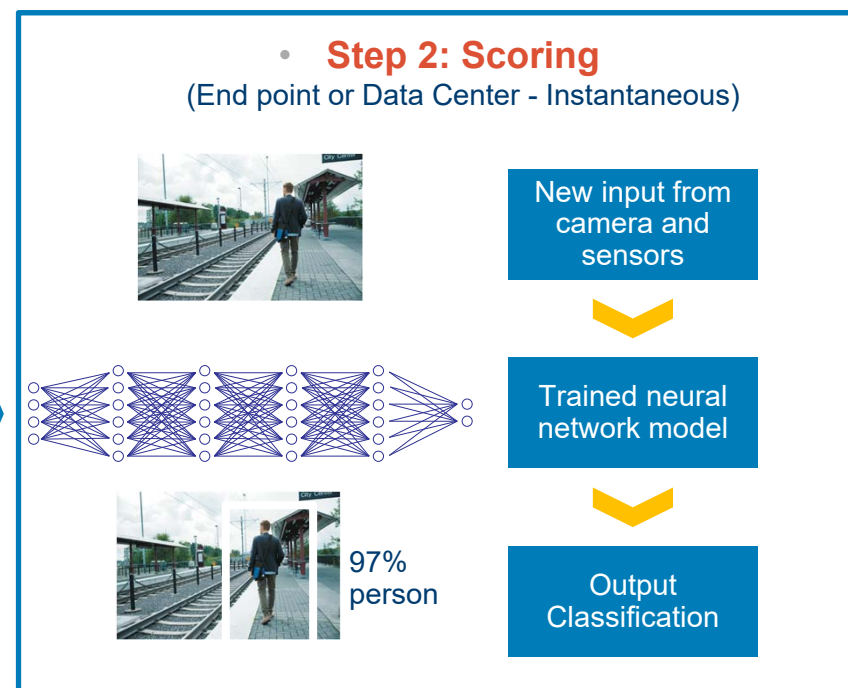
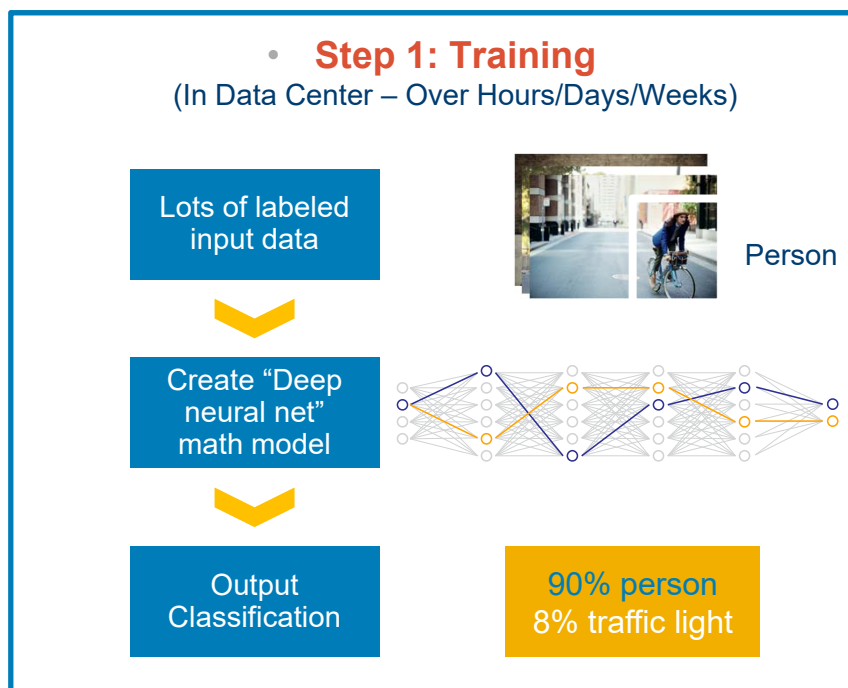


Fast Matrix Maths to the rescue

$$(W^{l+1})^T \delta^{l+1} = \begin{bmatrix} w_{11}^{l+1} & \dots & w_{m1}^{l+1} & \dots & w_{M1}^{l+1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{1k}^{l+1} & \dots & w_{mk}^{l+1} & \dots & w_{Mk}^{l+1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{1K}^{l+1} & \dots & w_{mK}^{l+1} & \dots & w_{MK}^{l+1} \end{bmatrix} \begin{bmatrix} \delta_1^{l+1} \\ \vdots \\ \delta_m^{l+1} \\ \vdots \\ \delta_M^{l+1} \end{bmatrix} = \begin{bmatrix} \sum_{m=1}^M \delta_m^{l+1} w_{m1}^{l+1} \\ \vdots \\ \sum_{m=1}^M \delta_m^{l+1} w_{mk}^{l+1} \\ \vdots \\ \sum_{m=1}^M \delta_m^{l+1} w_{mK}^{l+1} \end{bmatrix} \sim \delta^l$$

1. Do forward propagation (matrix operation)
2. Calculate Error of the final layer
3. Do backward propagation (matrix operation)
4. Update the weights

Deep Learning in practice



Internal Use - Confidential

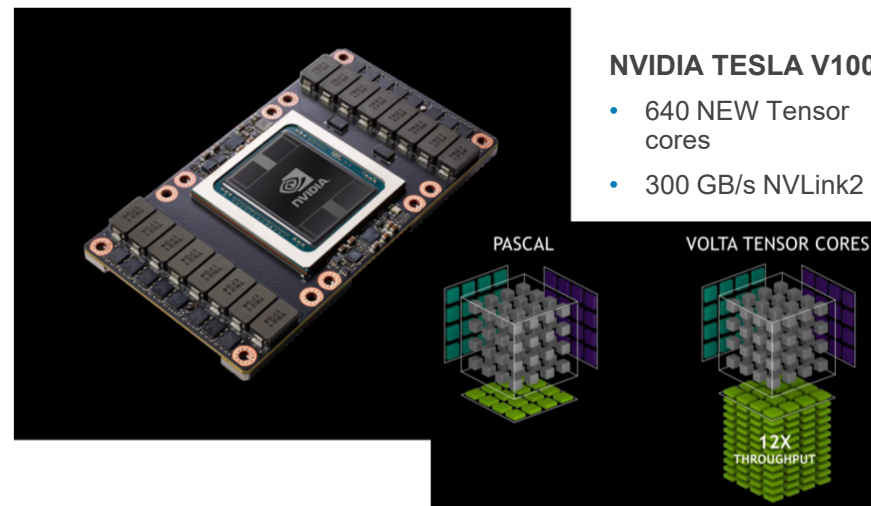
Scaling DL: Compute

Compute Efficiency

- Matrix operations are important (General Matrix to Matrix)
 - GEMM Processor friendly but need a lot of memory
- Floating point precision not critical – Int8, Int16 are sufficient
- Data transport:
 - High Bandwidth DRAM and memory hierarchy
- Dataflow architectures, memristors, Op AMPS

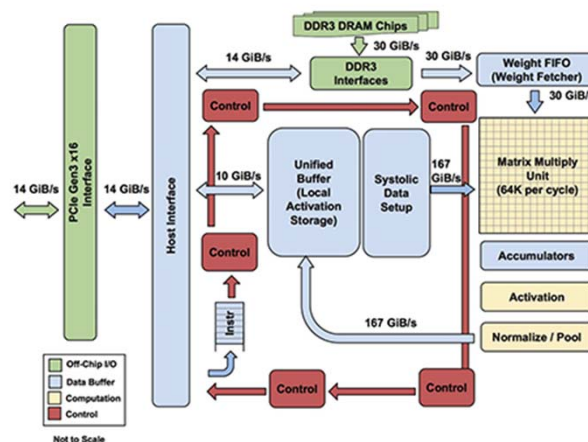
Scale Up vs Scale Out

- On-node vs off-node
- Faster off-chip interconnects
- Parallel Programming frameworks e.g. MPI
- BitFusion.io for GPU/FPGA Virtualization
- Distributed Deep Learning by IBM – GPUs over IB



NVIDIA TESLA V100

- 640 NEW Tensor cores
- 300 GB/s NVLink2



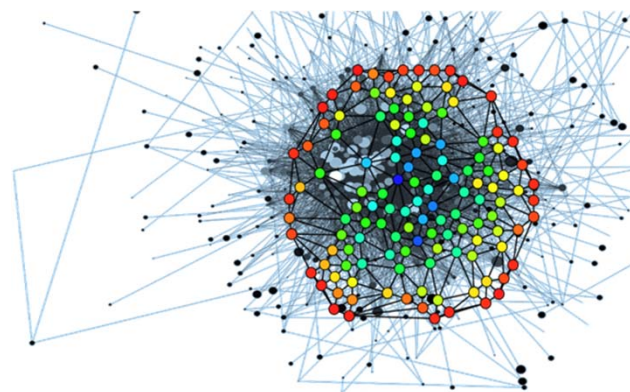
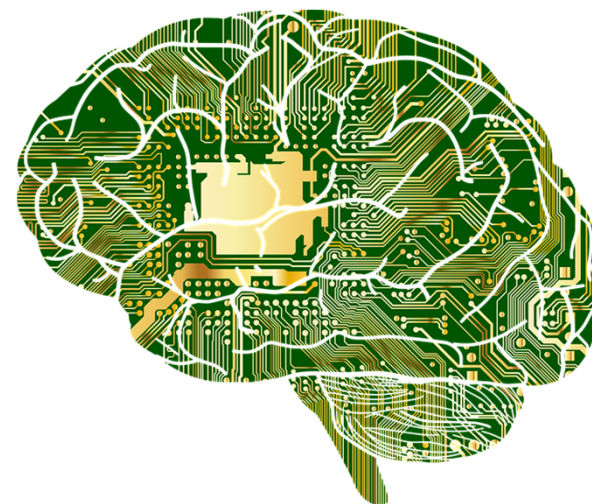
Google TPU

Matrix Multiplier Unit (MXU): 65,536 8-bit multiply-and-add units for matrix operations

Accelerating DL: Productivity

Software infrastructure:

- Partitioning problem space efficiently across nodes
- More efficient algorithms for doing the forward and back propagation (optimized algorithms and matrix computations)
- Distributed compute infrastructure that allows you to efficiently scale-out
- A fast (distributed) filesystem for managing data
- Resilient programming model, easy to deploy and maintain



Internal Use - Confidential

Summary

- Many of the techniques developed for HPC are already making their way into the latest DL processors
- Newer techniques unique to AI are also being developed, like optimized matrix calculators using memristors
- Scale-out and modularity will dominate provided we create the proper SW frameworks – higher speed interconnects will remain crucial
- Practicality of the technology will ultimately determine which technologies become mainstream

References

<https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>

M. Abadi et al. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. (2015). <http://www.tensorflow.org>

A. F. Aji and K. Heafield. 2017: Sparse Communication for Distributed Gradient Descent. (2017). arXiv:1704.05021

F. Akopyan et al. 2015. TrueNorth: Design and Tool Flow of a 65 mW 1 Million Neuron Programmable Neurosynaptic Chip. IEEE Transactions on ComputerAided Design of Integrated Circuits and Systems 34, 10 (2015), 1537–1557.

TAL BEN-NUN and TORSTEN HOEFLER, ETH Zurich, 2018: Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis (2018) arXiv:1802.09941v2



Questions or Feedback: Adnan.Khaleel@Dell.com