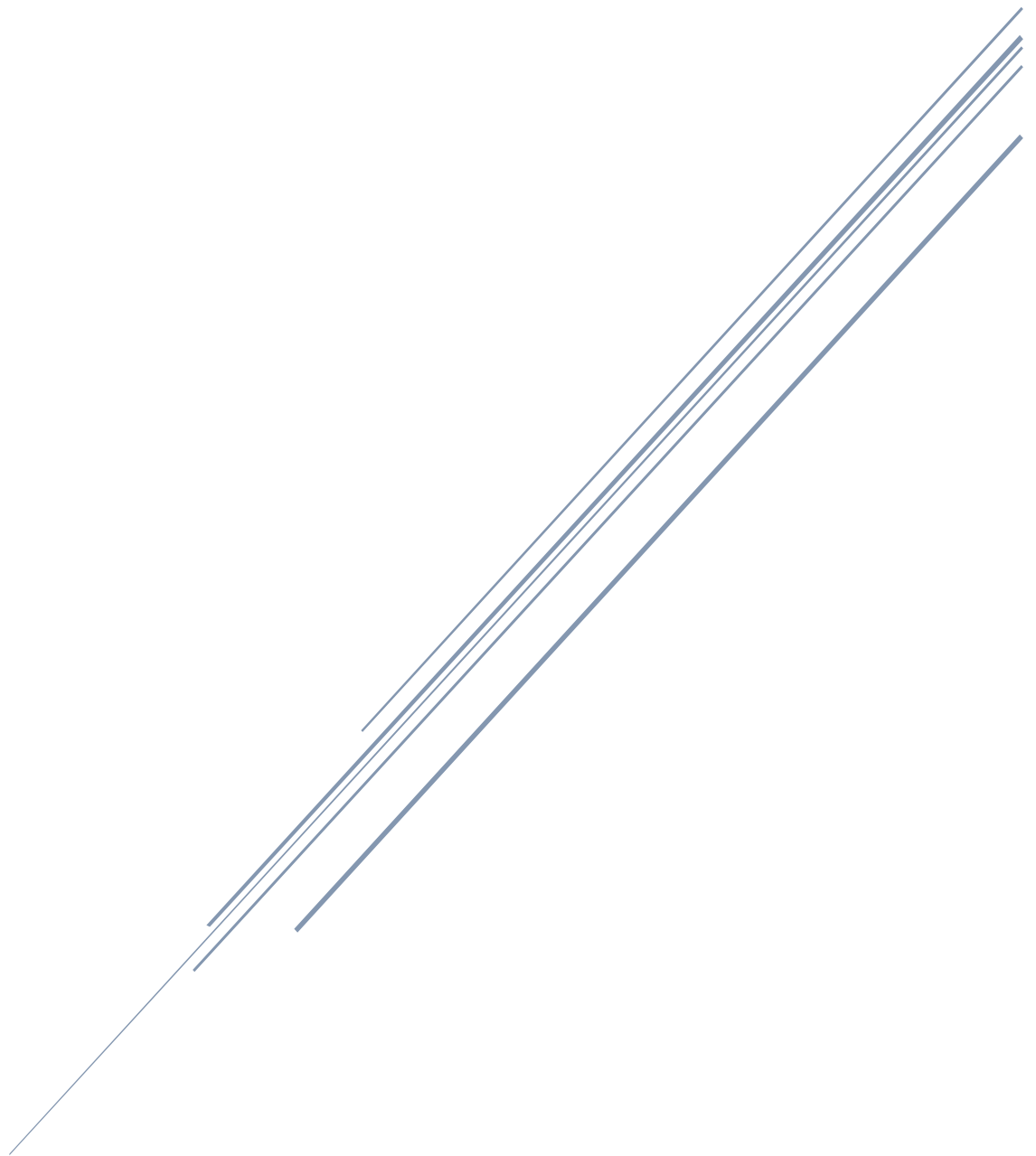


# DEEP LEARNING IN THE CLOUD

## Module 3



v-cardona

Using GPUs to Scale and Speed-up Deep Learning

## Contenido

Hardware accelerators.....	2
How does one use a GPU? .....	3
Deep learning in the cloud - PowerAI .....	3

# Hardware Accelerators

- **NVIDIA GPUs**
  - Software: CUDA
  - Card: GTX 1080, Tesla K80, Tesla P100, ....
  - GTX 1080(484 GB/s), K80(480 GB/s), P100(730 GB/s), ...
- **AMD GPUs**
- **TPUs (TensorFlow Processing Units)**
- **FPGAs**



You can use Google's TPUs, Nvidia GPU or even FPGA to accelerate your deep learning network computation time. These chips are particularly designed to support the training of neural networks, as well as the use of trained networks (that is, inference).

GPUs have two key limitations:

1. The first involves limited memory capacity. Yes, GPUs are very fast for data parallelism and, as such, we can take full advantage of their massive computing power. That said, we still need to store the data inside the GPU memory in order to access it and process it. You need a platform that can handle fast memory access in system memory, and also fast data exchange between GPUs.
2. The second limitation of GPUs is that you cannot easily buy these accelerators and embed them into your local machine. They're usually expensive and there are some dependencies and incompatibilities, which is the same as most hardware. Also, sometimes, you need a number of GPUs to handle your big datasets. So, these accelerators are not readily accessible, at least not for now.

# Hardware Accelerators

- **TPUs, GPUs, FPGAs**

- Training
- Inference



- **Limitations of GPUs:**

- **Limited Memory capacity**
  - Not practical for very large datasets
  - Alternative: reading data from system memory (overhead)
- **Accessibility**
  - Expensive, dependencies and incompatibilities

How does one use a GPU?

## Required Hardware

### 1. A laptop with an embedded GPU

- Not enough to solve real deep learning problems
- Needs to scale down the dataset

### 2. Using a GPU on a cloud service

- IBM cloud, AWS or Google cloud
- You can customize it
- Needs to move your data in the cloud
- Options for Single-GPU and Multi-GPU

### 3. Using a GPU cluster in the cloud

- IBM cloud

### 4. Using a GPU cluster on-premises

- Keep your data locally
- Cost- effective
- Perfect for sensitive data
- IBM PowerAI

## Deep learning in the cloud - PowerAI

There are various software frameworks for building and training a Deep Learning model, such as TensorFlow, Caffe, Torch, Chainer, and Theano. These frameworks can take advantage of graphical processing units (or GPUs) to accelerate the training or inference process, but need different types of extensions to work on GPUs, for example, CUDA Deep Neural Network, and nvCaffe libraries. IBM PowerAI is a package of software distributions for these types of software.

## What is powerAI?

