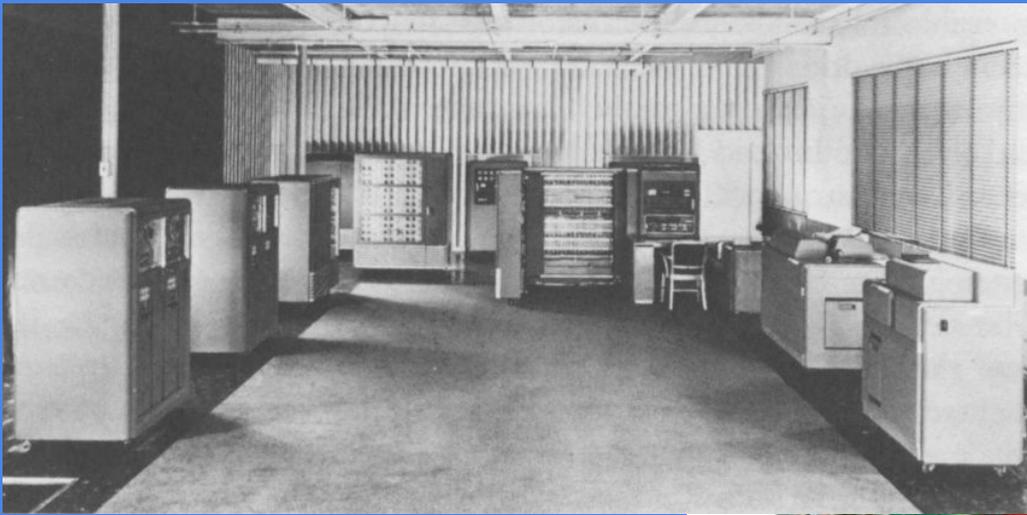




The development of AI thanks to advancements in hardware

Iker García Ferrero
Eritz Yerga Gutierrez

Motivation



IBM 701



IBM Minsky

Initial optimism

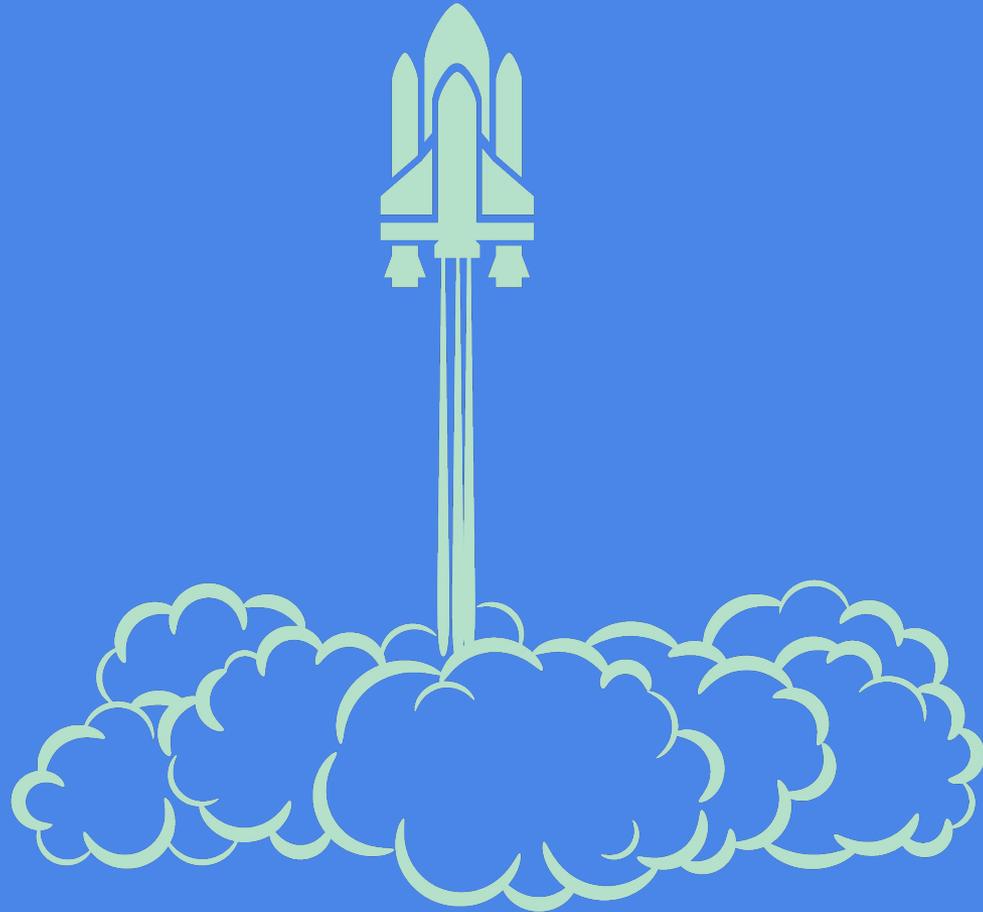
“ I’m not trying to impress you or leave you amazed, but the easiest way I have to sum it up is to say that at the moment machines capable of thinking, learning and creating exist in the world. Furthermore, their ability to do that will quickly increase until (in a near future) the magnitude of problems they will be able to solve will be on par with human mind itself on the same tasks. ”

- Herbert Simon. 1957.



February 24, 1956. Arthur Samuel demonstration of his program in an IBM 701 computer

AI is skyrocketing. Why now?



Advancements in CPUs

169,200 TIMES FASTER



APPLE II
1977

500,000 FLOPS
US\$1298



Ryzen 7 1800X
2017

84,600,000,000 flops
Full System: US\$800

1982



APPLE II plus

2017



APPLE iMAC 21.5"

WHAT'S DIFFERENT?

48 KILOBYTES	MEMORY	8 GIGABYTES	166,666x MORE
1 MHZ	CPU SPEED	1.6 GHZ	1600x FASTER
140 KILOBYTES	STORAGE	1 TERABYTE	7,508,684x MORE
\$95 (WITH DRIVE) \$4353 (2017 DOLLARS)	PRICE	\$1099	3.9x LESS



FORD MUSTANG 4CYL



FORD MUSTANG ECOBOOST 4CYL

22 MPG	MPG	20 MPG	0.9x LESS
16.2 SEC	0-60 TIME	5.5 SEC	2.9x FASTER
88 HP	HORSEPOWER	310 HP	3.5x MORE
\$6789 \$17,437 (2017 DOLLARS)	PRICE	\$26,195	1.5x MORE



IF CARS DEVELOPED AT THE PACE OF COMPUTERS, A 2017 FORD MUSTANG ECOBOOST WOULD HAVE:

660,764,192 HP
GET FROM 0-60 IN 0.0034 SEC
GET ABOUT 3,666,652 MPG AND COST \$4,471

BUT,

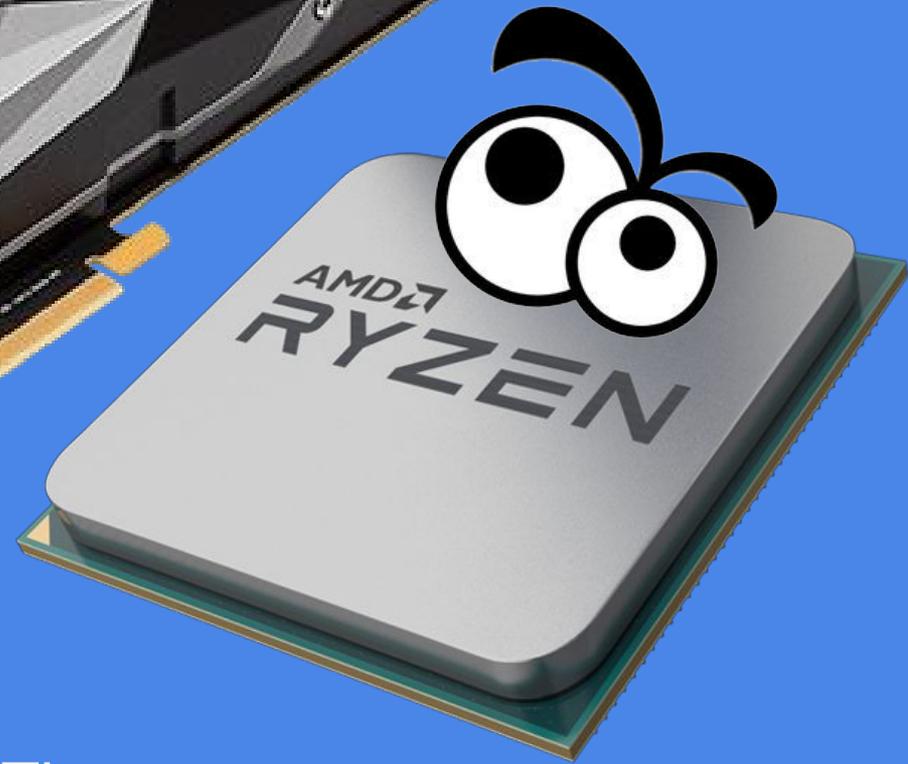
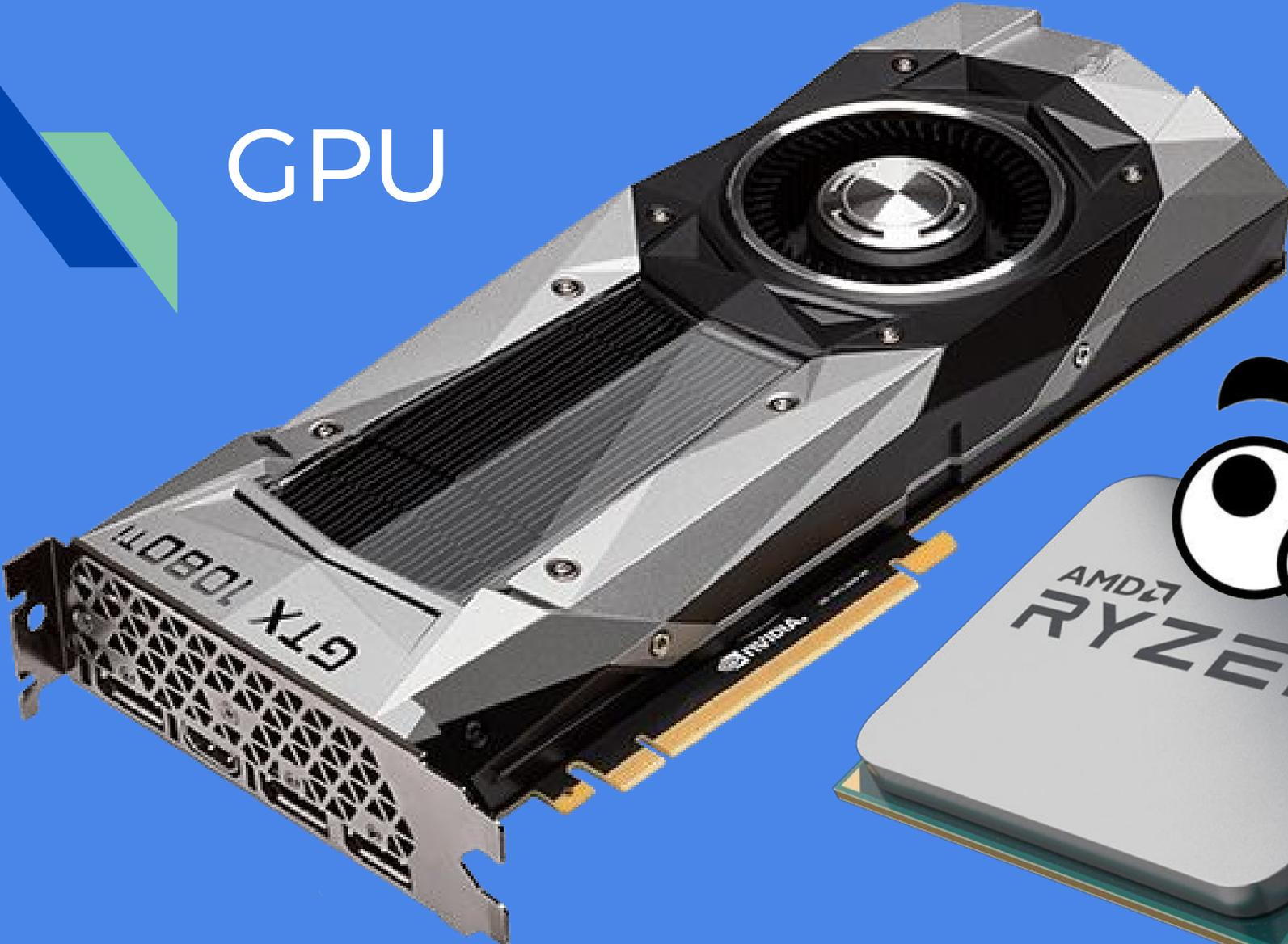
IF COMPUTERS DEVELOPED AT THE PACE OF CARS, A 2017 APPLE iMAC 21.5" WOULD HAVE:



43.2K OF RAM MEMORY,
WOULD RUN AT 2.9MHz
WOULD BE ABLE TO STORE 490K OF DATA
AND WOULD COST \$6,529

IT'S SORT OF A STUPID COMPARISON. BUT IT'S INTERESTING.

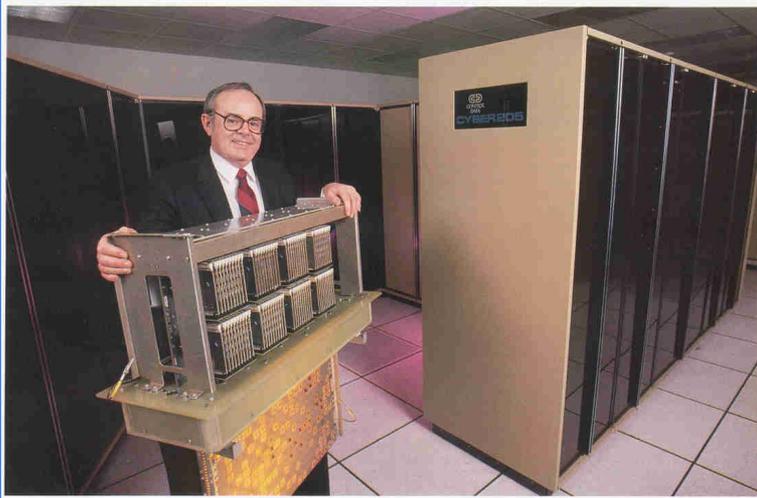
JALOPNIK

 GPU

GTX 1080 Ti
11000000000000 flops
130 times faster than the Ryzen 7 1800X

High computational power available to everyone thanks to GPUs

Millions of researchers, projects, applications... around the world



ETA10

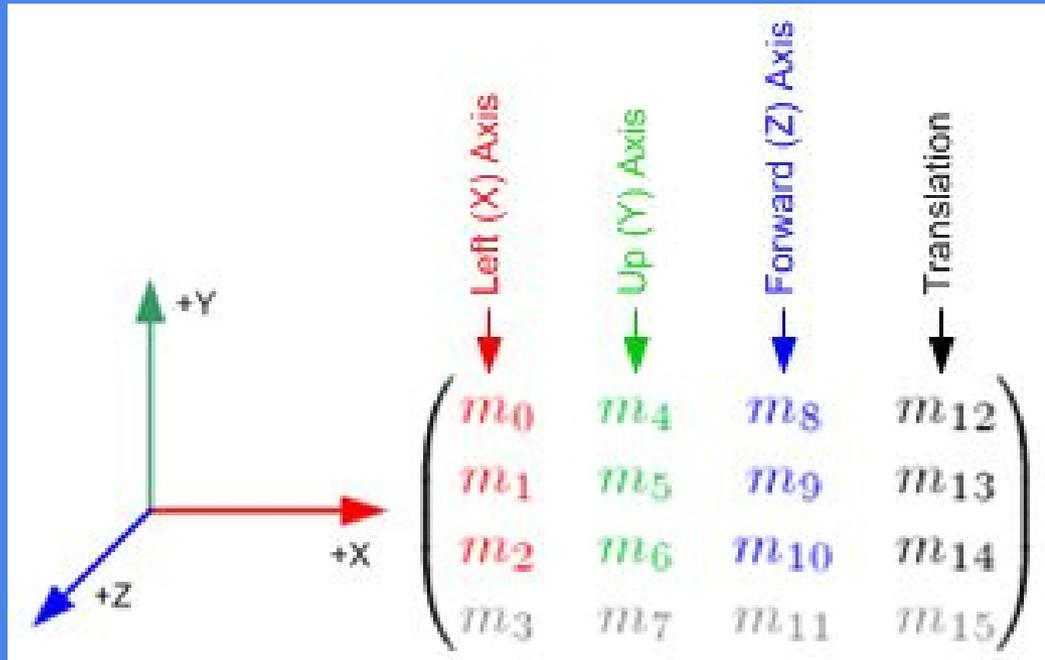
- 1987 Supercomputer
- 10 GLOPS
- High power consumption
- Liquid nitrogen cooling
- Very expensive
- Only available for few universities and research centers, such as Florida State University or Johnson Space Center



NVIDIA GTX 1050

- 2016 entry level gpu
- 1800 GLOPS
- Less than 75 watts
- Passive cooling
- US\$ 110
- Available to everyone. You can buy one right now and be using it tomorrow morning.

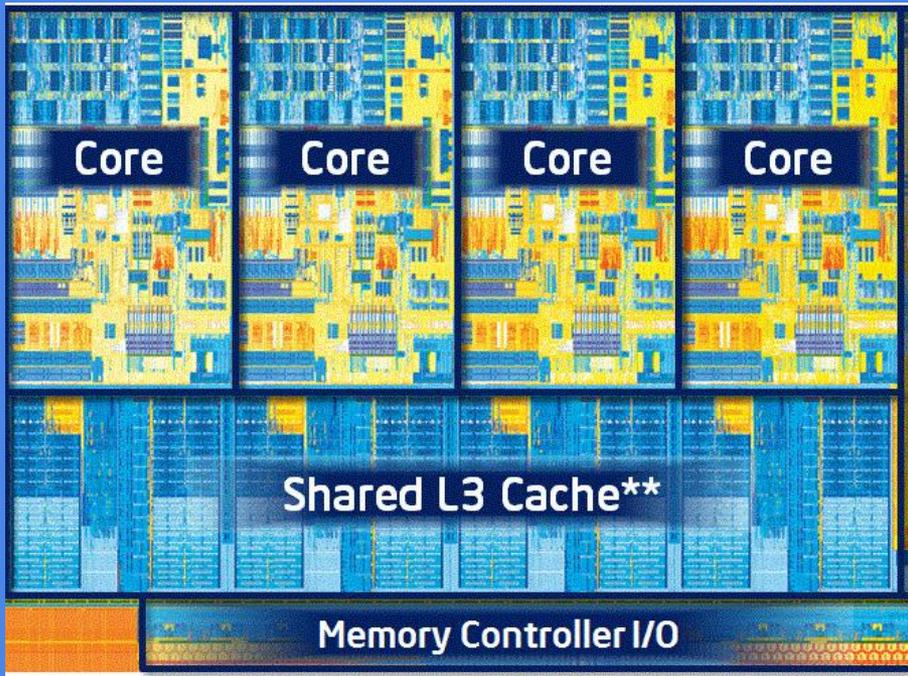
Why GPUs?



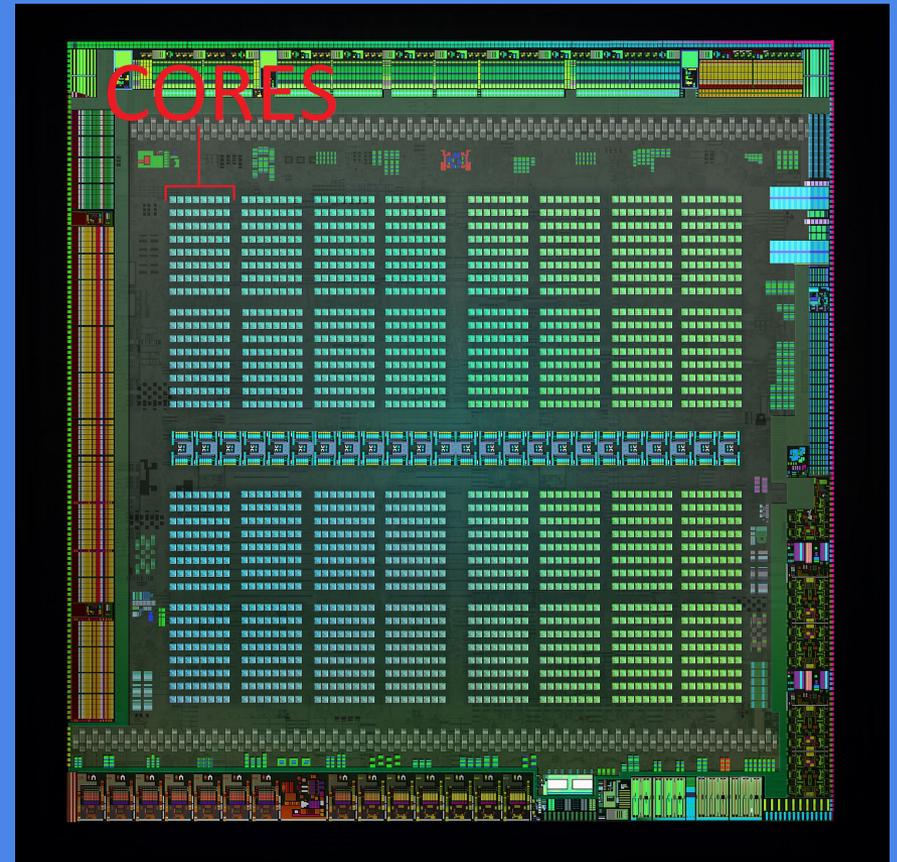
World's First GPU



GeForce 256: “A single-chip processor with integrated transform, lighting, triangle setup/clipping, and rendering engines that is capable of processing a minimum of 10 million polygons per second”
- NVIDIA 1999



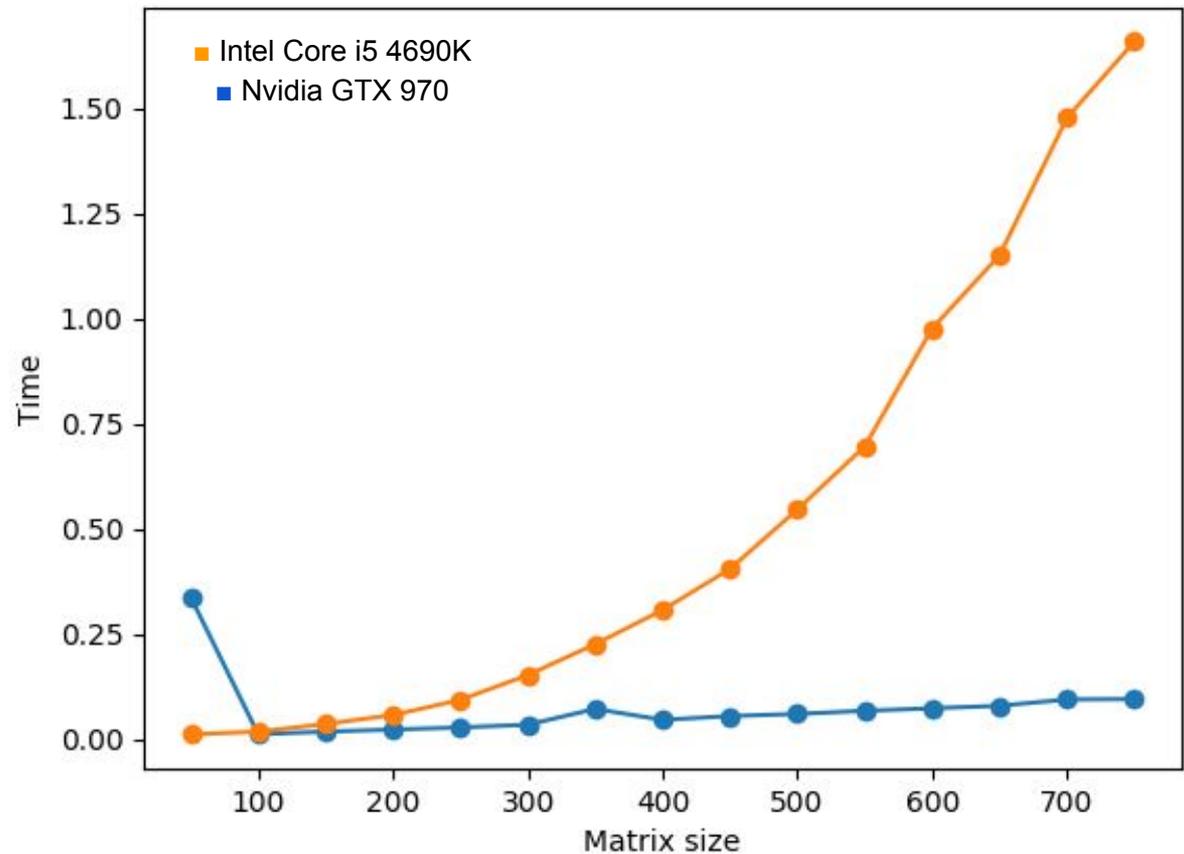
Intel i7. (4 Cores)



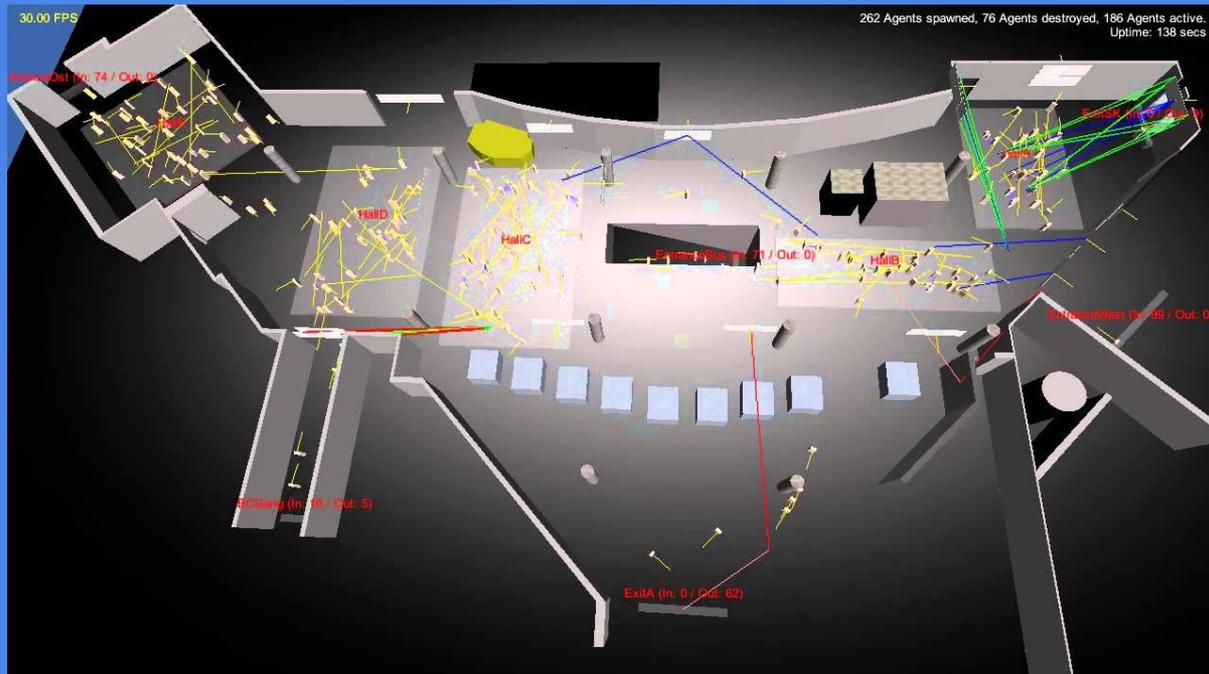
Nvidia GTX 1080 (2560 Cores)

Intel Core i5 4690K a 4 core CPU at 4GHz VS NVIDIA GTX 970 a 1664 core GPU at 1,350GHz

```
1 from __future__ import print_function
2 import matplotlib
3 import matplotlib.pyplot as plt
4 import tensorflow as tf
5 import time
6
7 def get_times(maximum_time):
8
9     device_times = {
10         "/gpu:0": [],
11         "/cpu:0": []
12     }
13     matrix_sizes = range(50, 50000, 50)
14
15     for size in matrix_sizes:
16         for device_name in device_times.keys():
17
18             print("##### Calculating on the " + device_name + " #####")
19
20             shape = (size, size)
21             data_type = tf.float16
22             with tf.device(device_name):
23                 r1 = tf.random_uniform(shape=shape, minval=0, maxval=1, dtype=data_type)
24                 r2 = tf.random_uniform(shape=shape, minval=0, maxval=1, dtype=data_type)
25                 dot_operation = tf.matmul(r2, r1)
26
27
28             with tf.Session(config=tf.ConfigProto(log_device_placement=True)) as session:
29                 start_time = time.time()
30                 result = session.run(dot_operation)
31                 time_taken = time.time() - start_time
32                 print(result)
33                 device_times[device_name].append(time_taken)
34
35             print(device_times)
36
37             if time_taken > maximum_time:
38                 return device_times, matrix_sizes
39
40
41 device_times, matrix_sizes = get_times(1.5)
42 gpu_times = device_times["/gpu:0"]
43 cpu_times = device_times["/cpu:0"]
44
45 plt.plot(matrix_sizes[:len(gpu_times)], gpu_times, 'o-')
46 plt.plot(matrix_sizes[:len(cpu_times)], cpu_times, 'o-')
47 plt.ylabel('Time')
48 plt.xlabel('Matrix size')
49 plt.show()
```

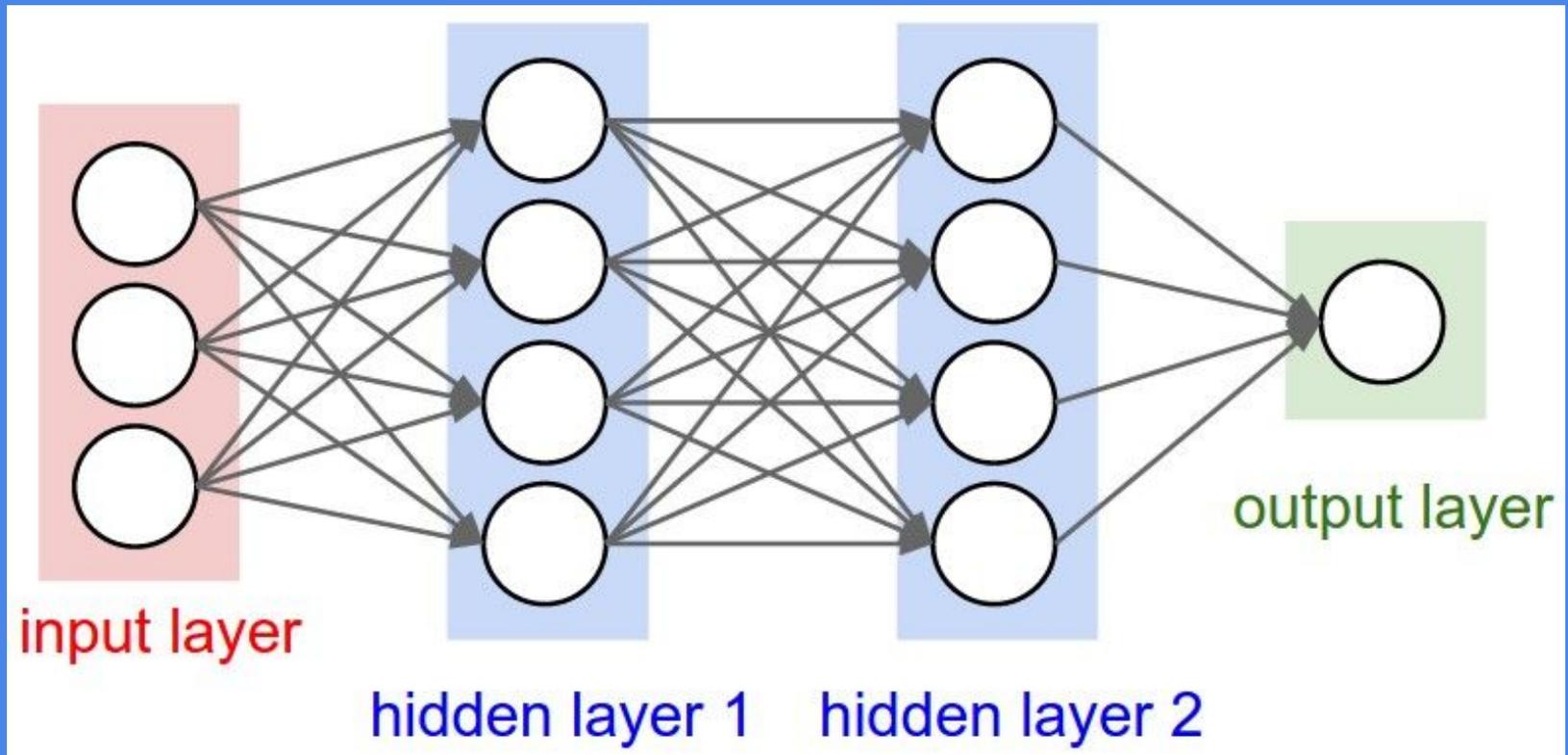


Ok, GPUs are faster for videogames, but, why are so good for AI?



Airport simulation using multi agent systems

Neural networks, deep learning...



High bandwidth, computer power isn't everything



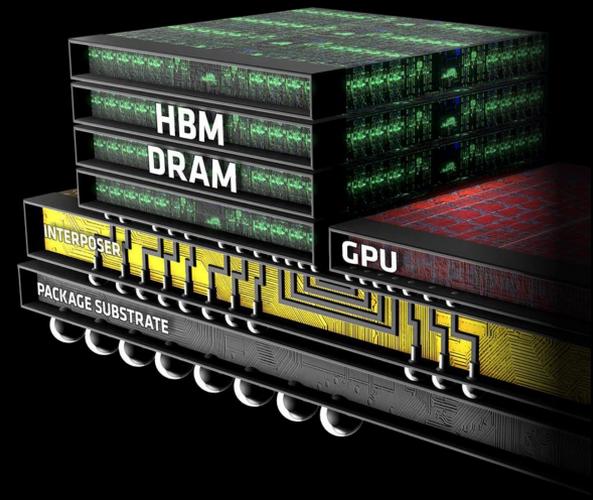
DDR4
Low latency
50GB/S

GRAPHICS TECHNOLOGY LEADERSHIP



HIGH BANDWIDTH MEMORY

- ▲ First in the Industry with High Bandwidth Memory (HBM) Technology
- ▲ 3D HBM DRAM Die Stack on Silicon Interposer
- ▲ >3X Performance/Watt Compared to GDDR5³
- ▲ >50% Power Savings Versus GDDR5⁴



10 | 2015 FINANCIAL ANALYST DAY | MAY 6, 2015

GDDR5, GGDR5X, HBM
High bandwidth
900GB/s

NEURAL NETWORK CPU vs GPU

TensorFlow. MNIST convolutional mode.

AMD Ryzen 7 1700 a 8 cores / 16 threads CPU at 3,2Ghz.

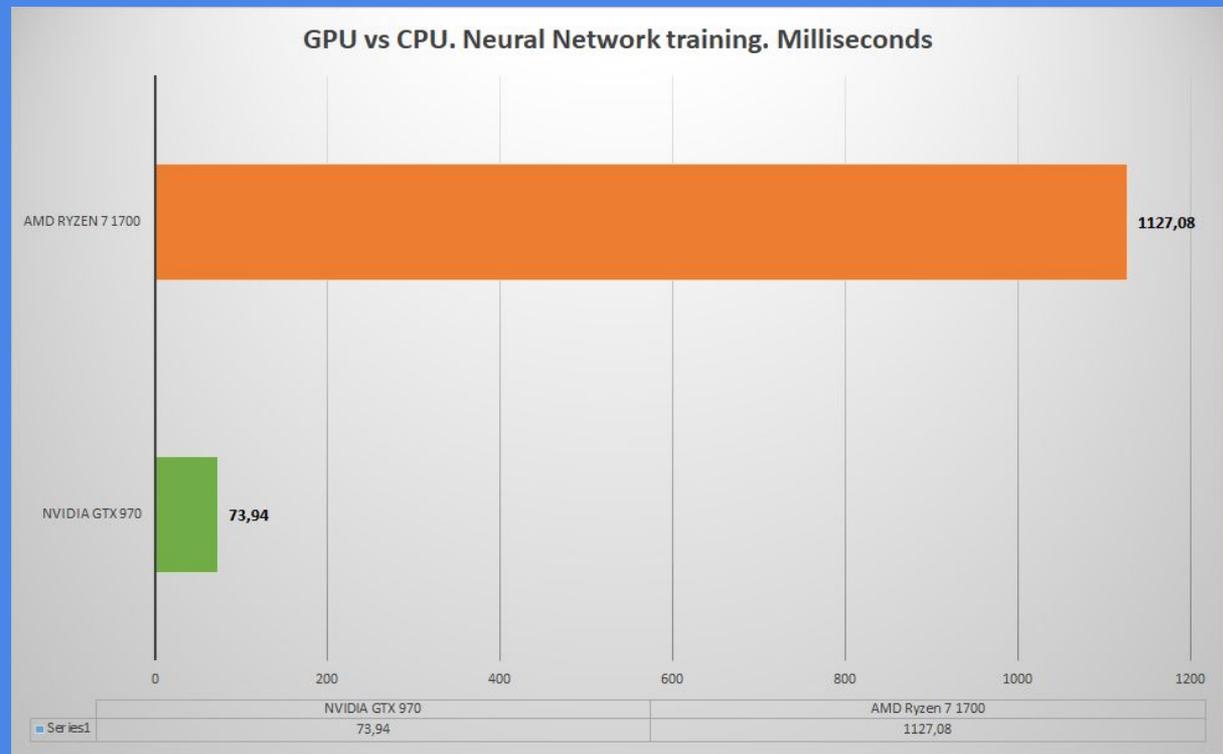
Released: 2017. Price: US\$ 320

NVIDIA GTX 970 a 1664 core GPU at 1,350GHz

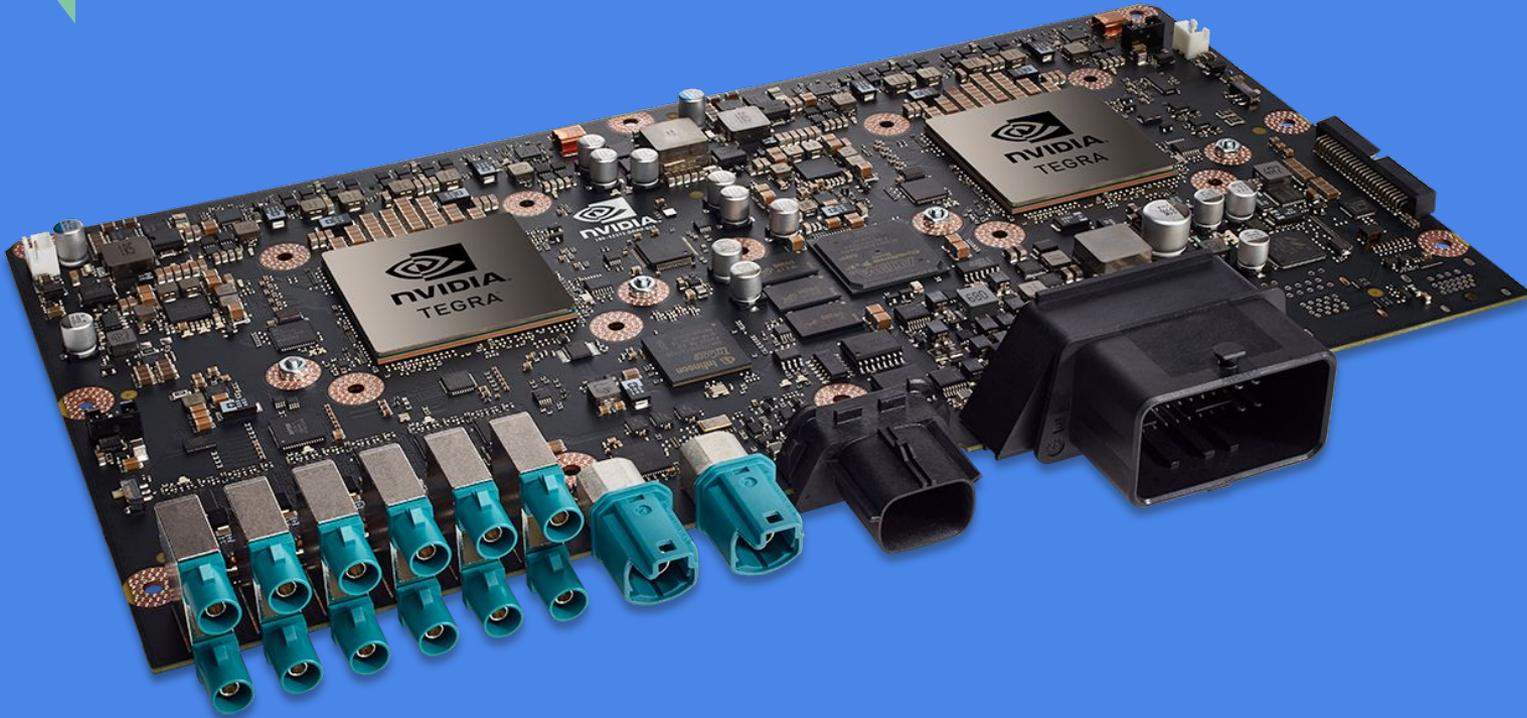
Released: 2014. Price; US\$ 329

15.24 times
faster

```
10 import numpy
11 from six.moves import urllib
12 from six.moves import xrange # pylint: disable=redefined-builtin
13 import tensorflow as tf
14
15 SOURCE_URL = "http://yann.lecun.com/exdb/mnist/"
16
17 # MNIST constants
18 IMAGE_SIZE = 28
19 NUM_CHANNELS = 1
20 PIXEL_DEPTH = 255
21 NUM_LABELS = 10
22 VALIDATION_SIZE = 1000 # Size of the validation set.
23 SEED = 66669 # Set to None for random seed.
24 BATCH_SIZE = 54
25 NUM_EPOCHS = 10
26 EVAL_BATCH_SIZE = 60
27 EVAL_FREQUENCY = 100 # Number of steps between evaluations.
28
29
30 FLAGS = None
31
32
33 def data_type():
34     """Return the type of the activations, weights, and placeholder variables."""
35     if FLAGS.use_fp16:
36         return tf.float16
37     else:
38         return tf.float32
39
40
41 def maybe_download(filename):
42     """Download the data from 'http://www.tensorflow.org/mnist/'
43     if not in 'mnist_data_dir' directory.
44     If file already exists, do nothing.
45     If file not found, download it from the source URL.
46     If file not found, print an error message.
47     If file not found, return None.
48     """
49     if not os.path.exists(filename):
50         url = SOURCE_URL + filename
51         with urllib.urlopen(url) as f:
52             data = f.read()
53             with open(filename, 'wb') as f:
54                 f.write(data)
55             print('Successfully downloaded', filename, size, 'bytes.')
56     return filename
57
58
59 def extract_data(filename, num_images):
60     """Extract the images into a 4D tensor [image index, y, x, channels].
61     Values are rescaled from [0, 255] down to [-0.5, 0.5].
62     """
63     print('Extracting', filename)
64     with gzip.open(filename) as bytestream:
65         bytestream.read(1)
66         buf = bytestream.read(DIMAGE_SIZE * DIMAGE_SIZE * num_images * NUM_CHANNELS)
67         data = numpy.frombuffer(buf, dtype=numpy.uint8).astype(numpy.float32)
68         data = (data - 127.5) / 127.5 # Pixel depth
69         data = data.reshape(num_images, DIMAGE_SIZE, DIMAGE_SIZE, NUM_CHANNELS)
70         return data
71
72
73 def extract_labels(filename, num_images):
74     """Extract the labels into a vector of 1000 labels IDs.
75     """
76     print('Extracting', filename)
77     with gzip.open(filename) as bytestream:
78         bytestream.read(1)
79         buf = bytestream.read(1 * num_images)
80         labels = numpy.reshape(buf, (num_images,)).astype(numpy.int64)
81     return labels
82
83
84 def main(_):
85     """Generate a data dataset that matches the dimensions of MNIST."""
86     data = numpy.zeros(
87         shape=(num_images, DIMAGE_SIZE, DIMAGE_SIZE, NUM_CHANNELS),
88         dtype=numpy.float32)
89     labels = numpy.zeros(shape=(num_images,), dtype=numpy.int64)
90     for image in xrange(num_images):
91         label = image % 10
```

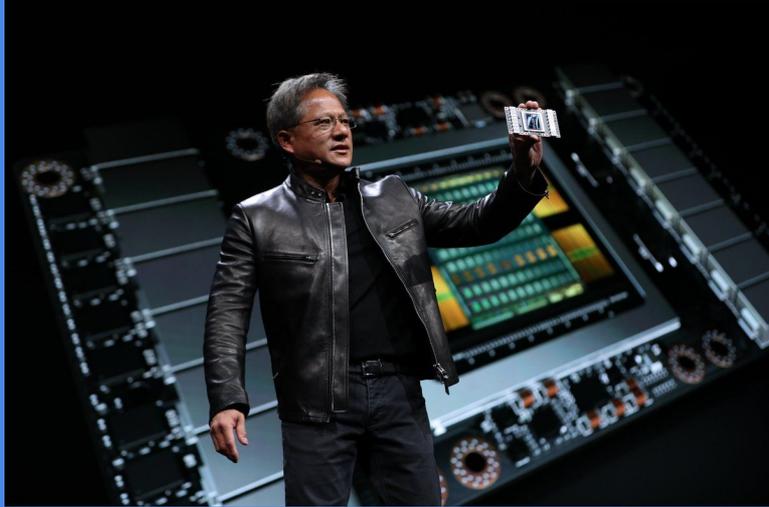


A world of possibilities \$\$ Real applications for AI \$\$



NVIDIA Parker, a computer for autonomous cars with a 256 core GPU

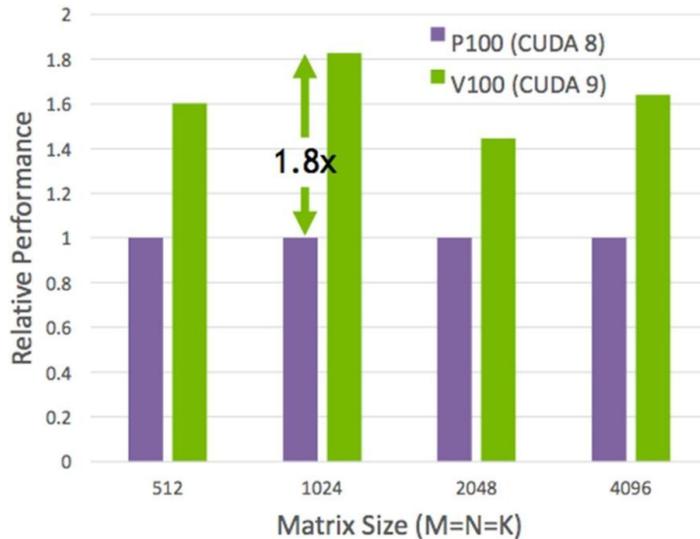
A vision of the future



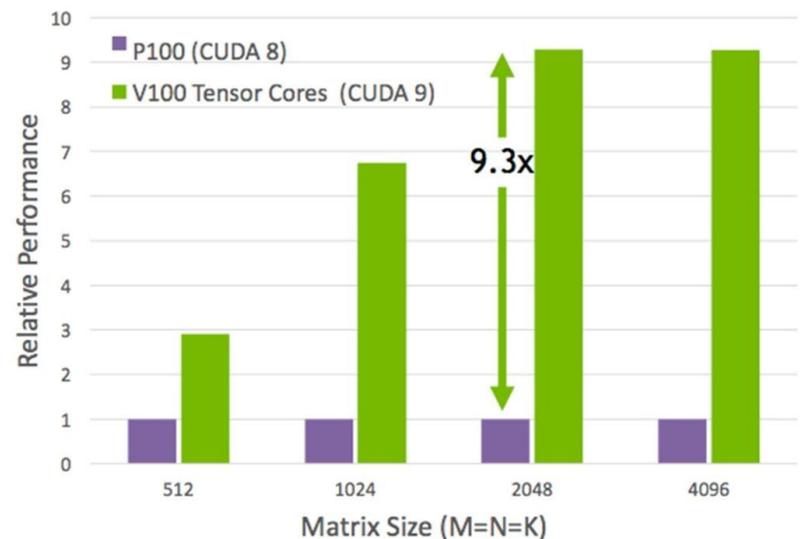
GPU PERFORMANCE COMPARISON

	P100	V100	Ratio
DL Training	10 TFLOPS	120 TFLOPS	12x
DL Inferencing	21 TFLOPS	120 TFLOPS	6x
FP64/FP32	5/10 TFLOPS	7.5/15 TFLOPS	1.5x
HBM2 Bandwidth	720 GB/s	900 GB/s	1.2x
STREAM Triad Perf	557 GB/s	855 GB/s	1.5x
NVLink Bandwidth	160 GB/s	300 GB/s	1.9x
L2 Cache	4 MB	6 MB	1.5x
L1 Caches	1.3 MB	10 MB	7.7x

cuBLAS Single Precision (FP32)



cuBLAS Mixed Precision (FP16 Input, FP32 compute)



New generation of supercomputers for AI



Nvidia DGX-1



IBM Minsky

Questions

