

Deep Learning e Visão Computacional com CNTK

VITOR MERIAT
@VITORMERIAT



Visual Studio Summit

#VSSUMMIT

About me

Vitor is a computer scientist who is passionate about creating software that will positively change the world we live in.

Currently, he works as **Data Scientist and Machine Learning Engineer** at **ESX**, where he is helping to shape new disruptive services based in Cloud Computing. Data science enthusiast, he works with Big Data projects, data analytics and **Microsoft MVP Azure**.



vitormeriat.com.br



linkedin.com/in/vitormeriat



twitter.com/vitormeriat



github.com/vitormeriat



youtube.com/vitormeriat



Visual Studio Summit



Microsoft



Agenda

O que é Visão Computacional e por que é uma assunto tão difícil;

Como Deep Learning e a Cloud Computing me ajudam com isso?





The researchers used 14 hours of footage of Barack Obama to produce their model.

SUBSCRIBE



0:34 / 1:25

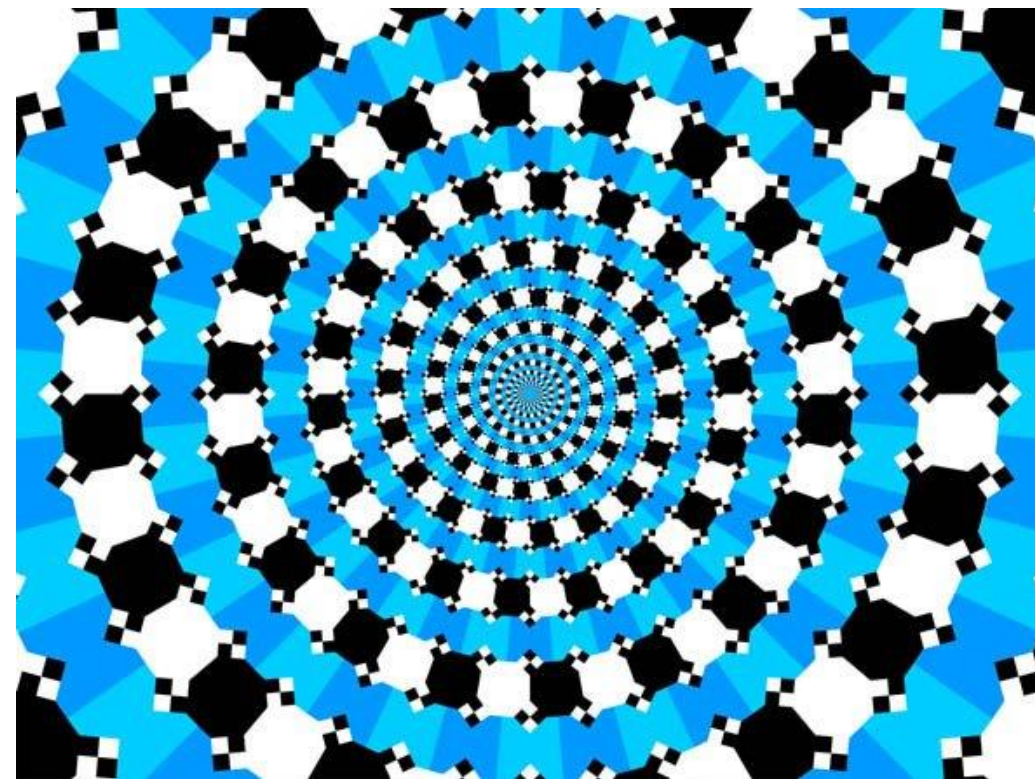


Visão é um problema inverso e mal-posto

- Pequenas variações na imagem levam a interpretações distintas;
- Visão é ambígua
- Ruído piora o problema

Visão Computacional é difícil

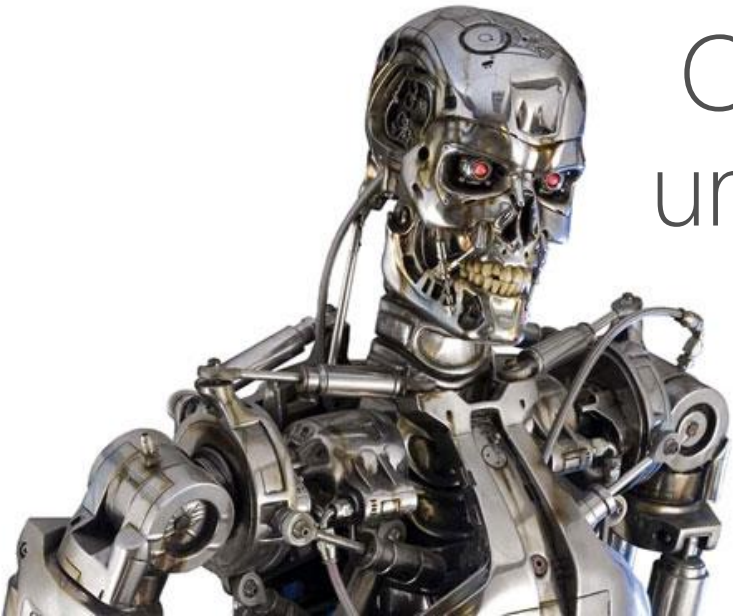
- Ambiguidades
- Distorções
- Paradoxos
- Ficções



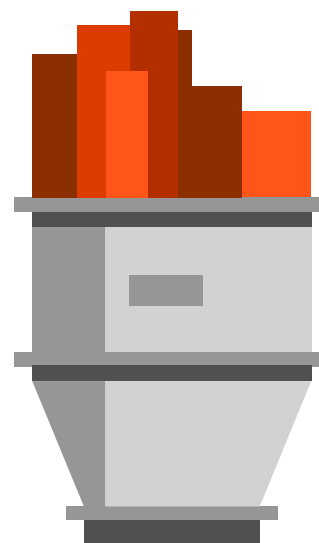
Eu quero ensinar para o meu computador
que isso é um gato

Eu sou
um gato

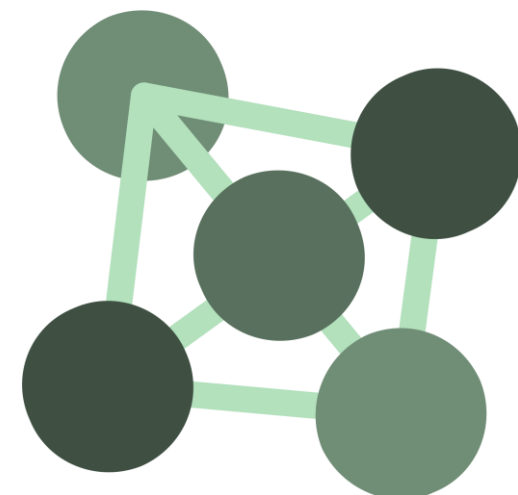
O que é
um gato?



Vamos ensinar esse computador

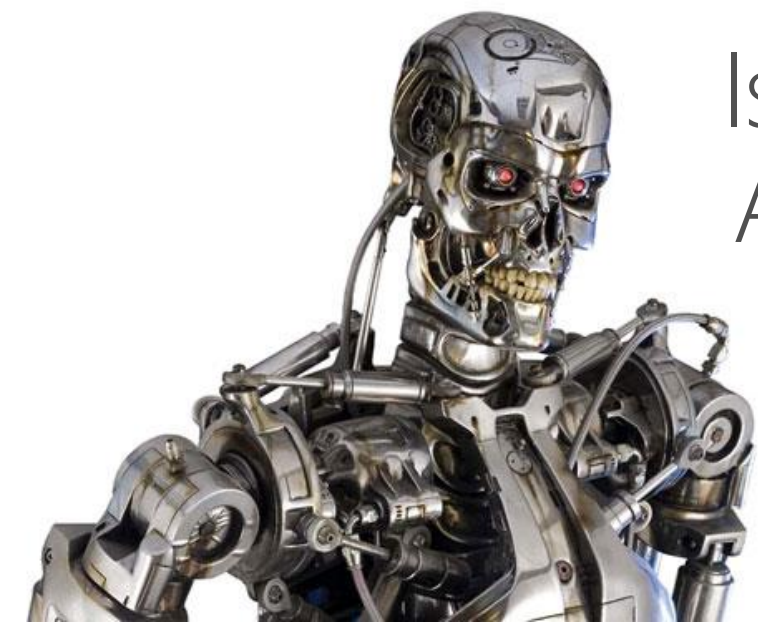


10101
01010
00100
11011
0101
010
00





Isso é um gato.
Agora eu já sei.



WTF?!?!



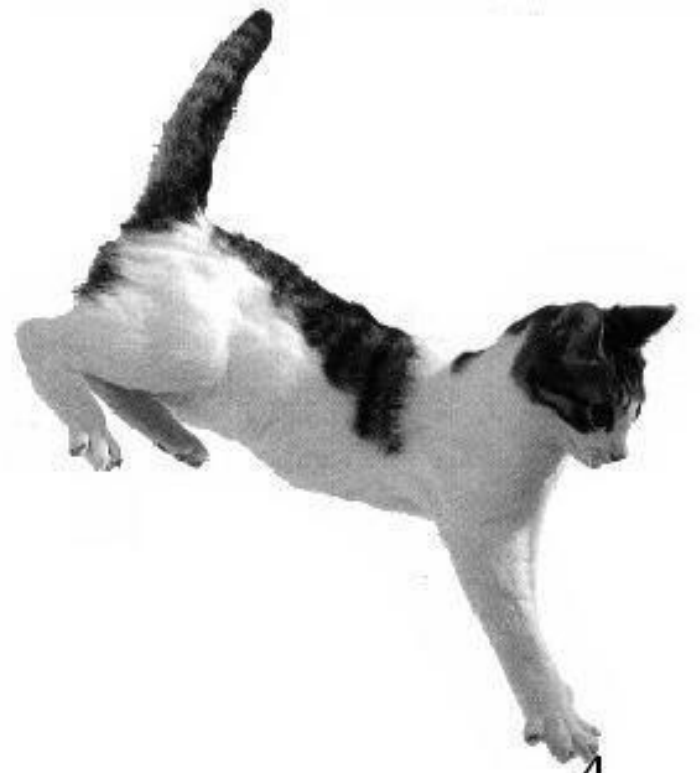
1



3

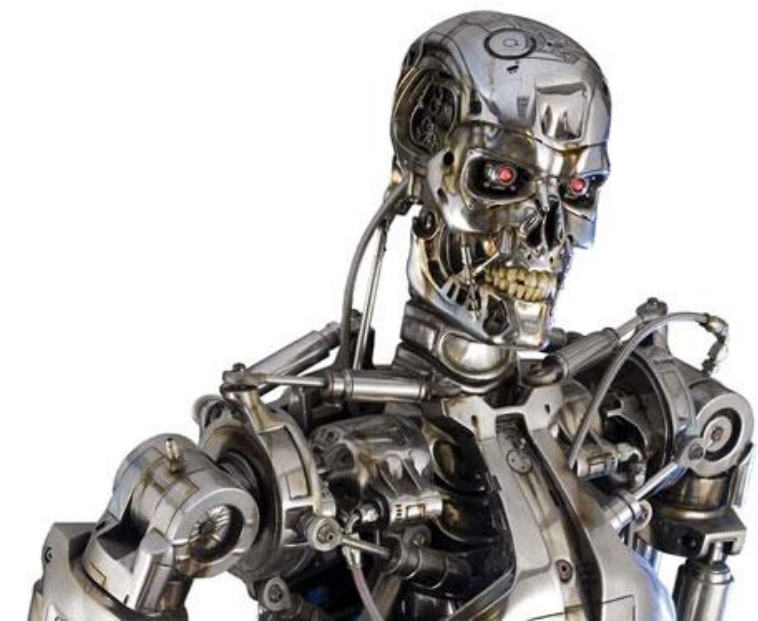


2



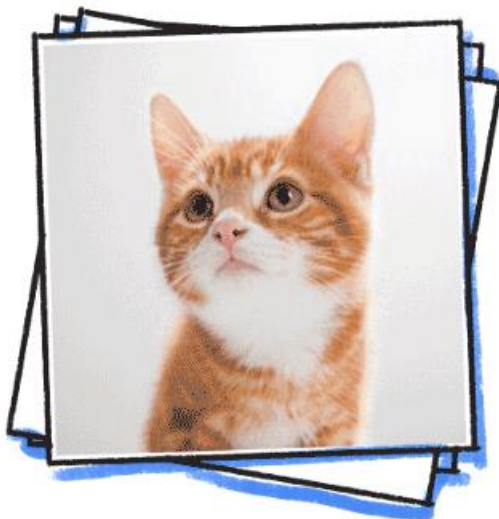
4

WTF?!?!



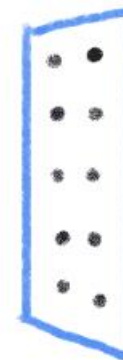
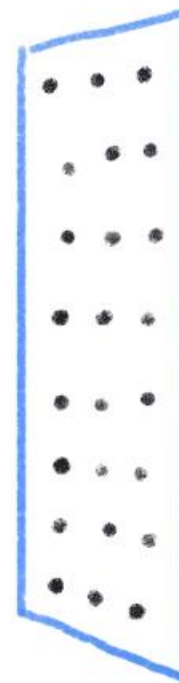
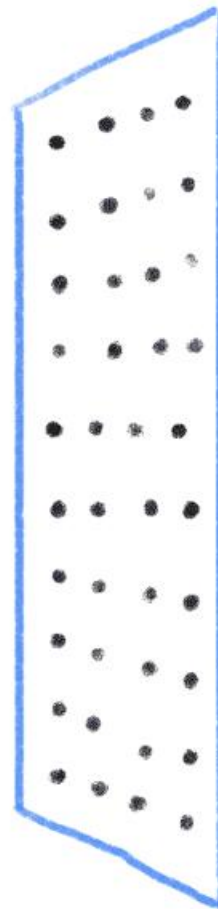
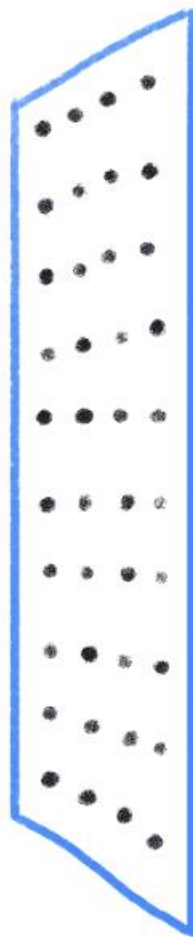
AngelBengal

CAT

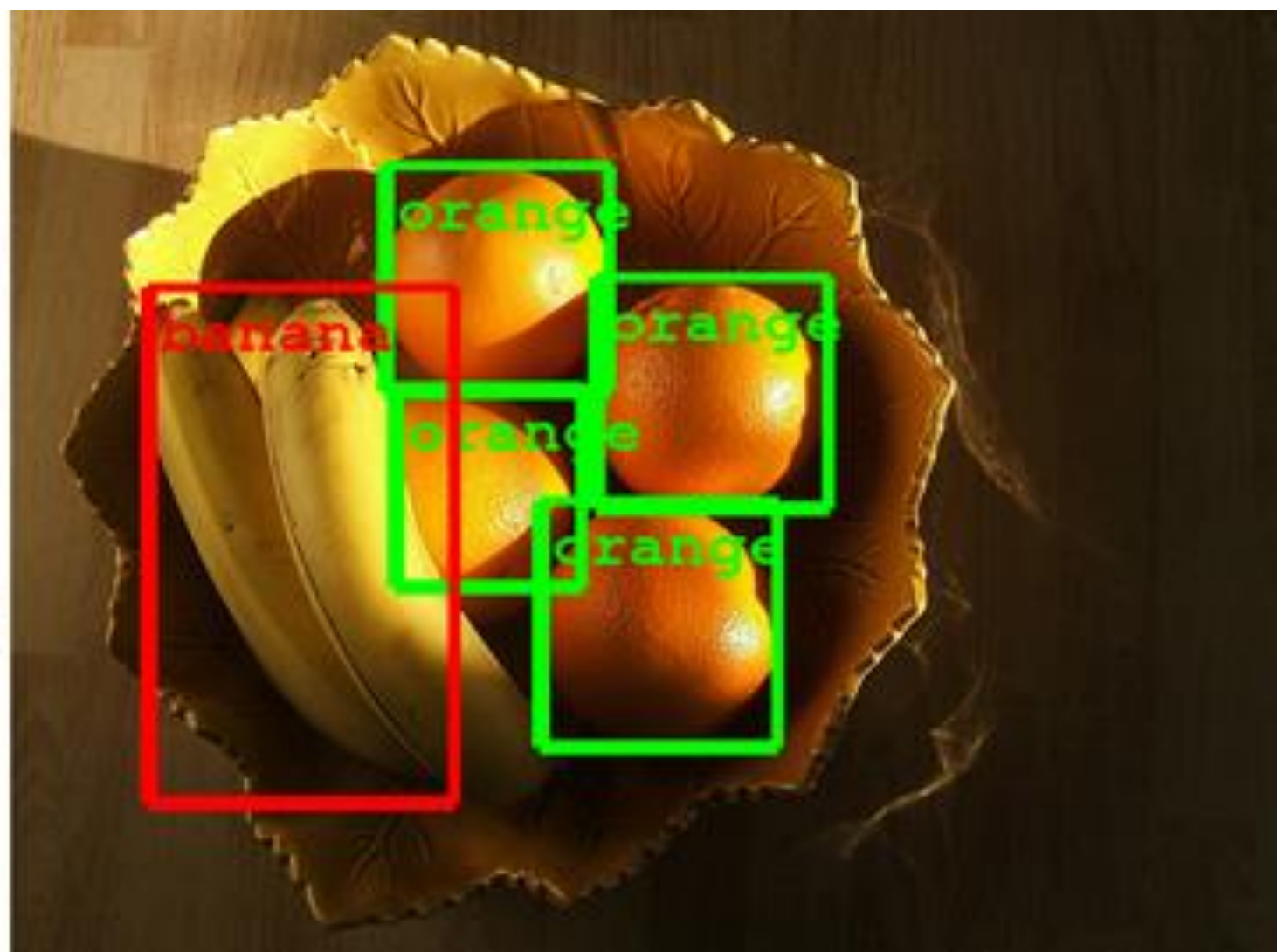
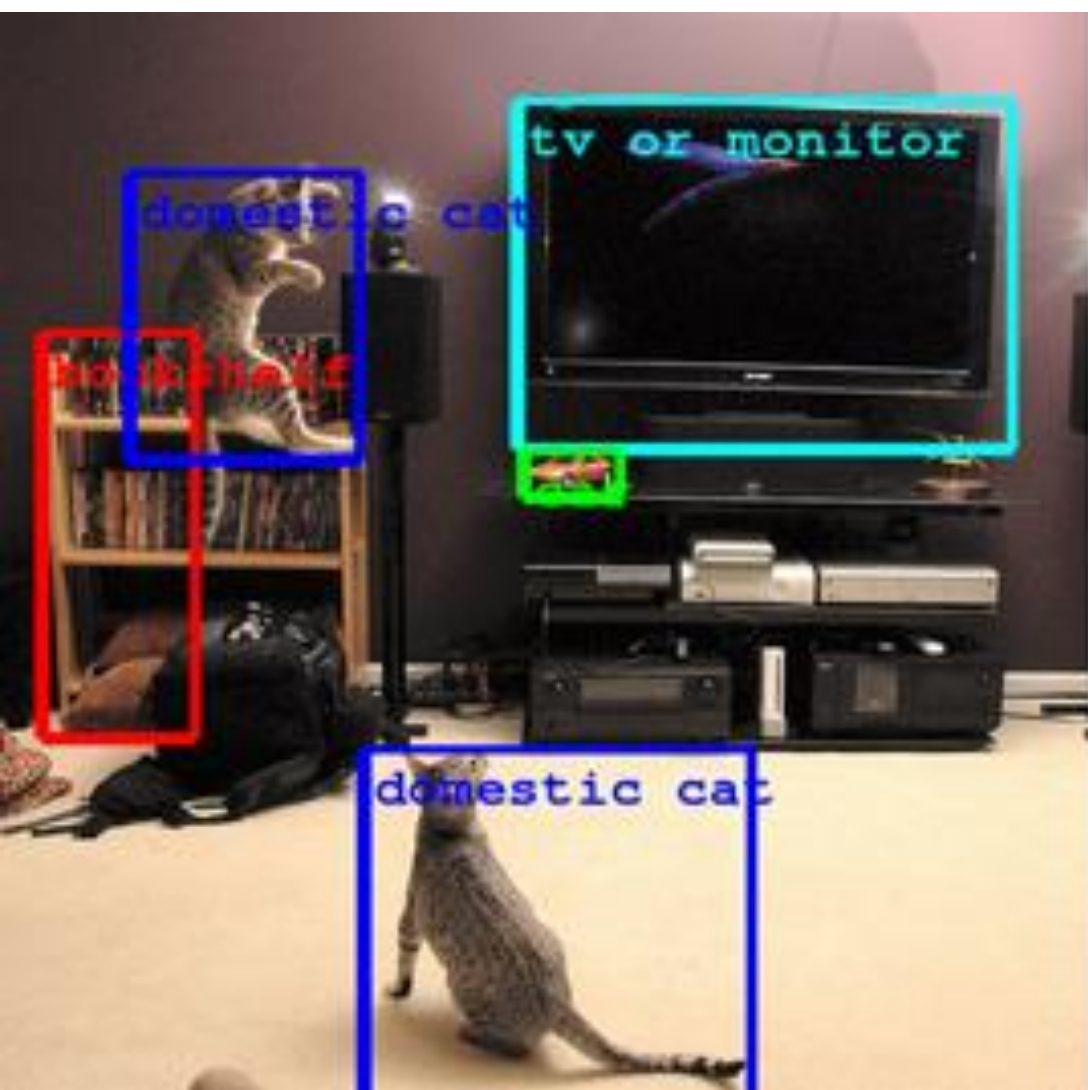


(LBELED
PHOTOS)

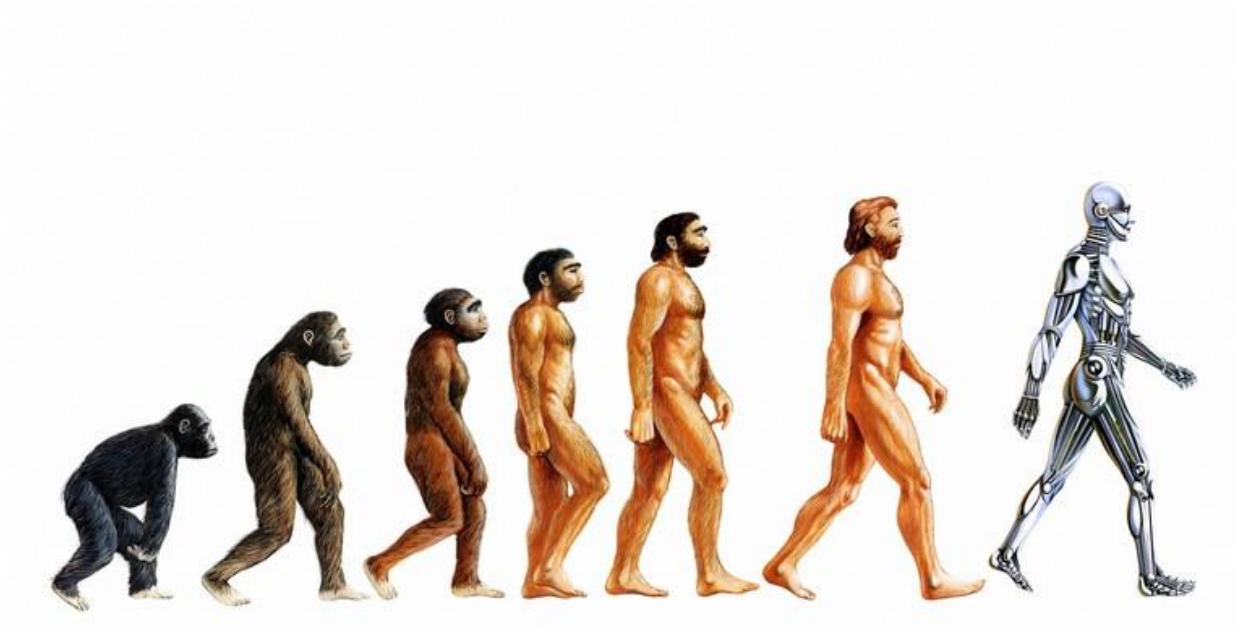
DOG



OUTPUT

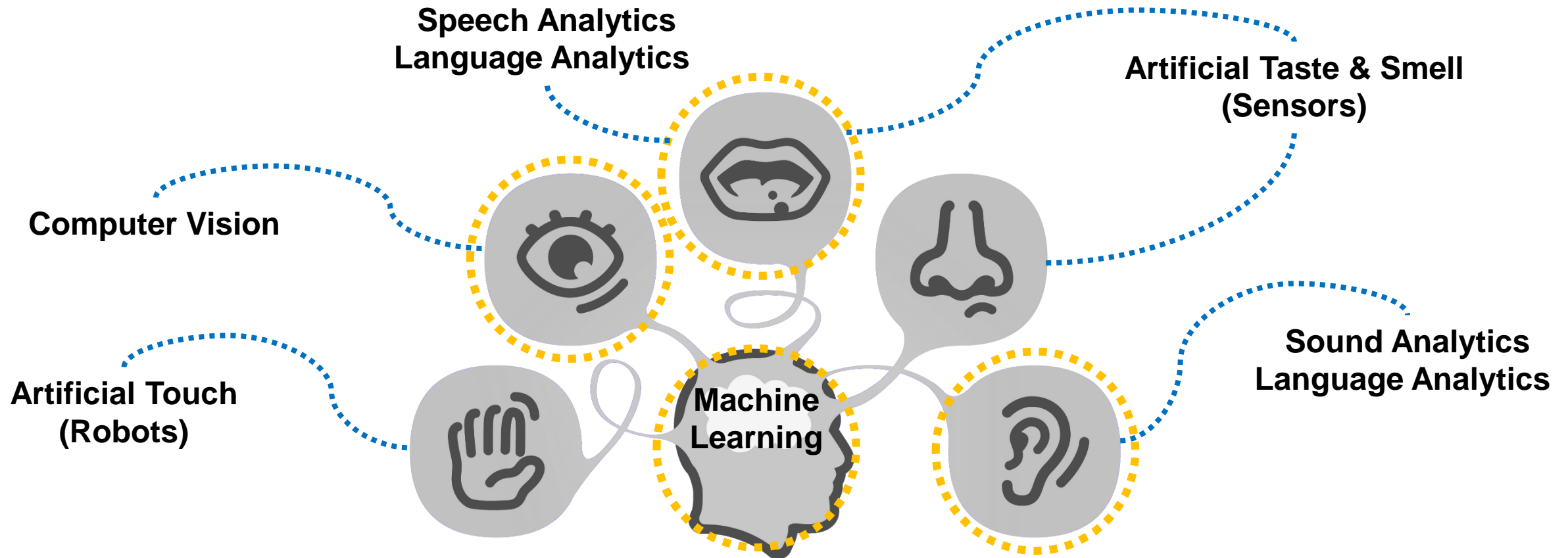


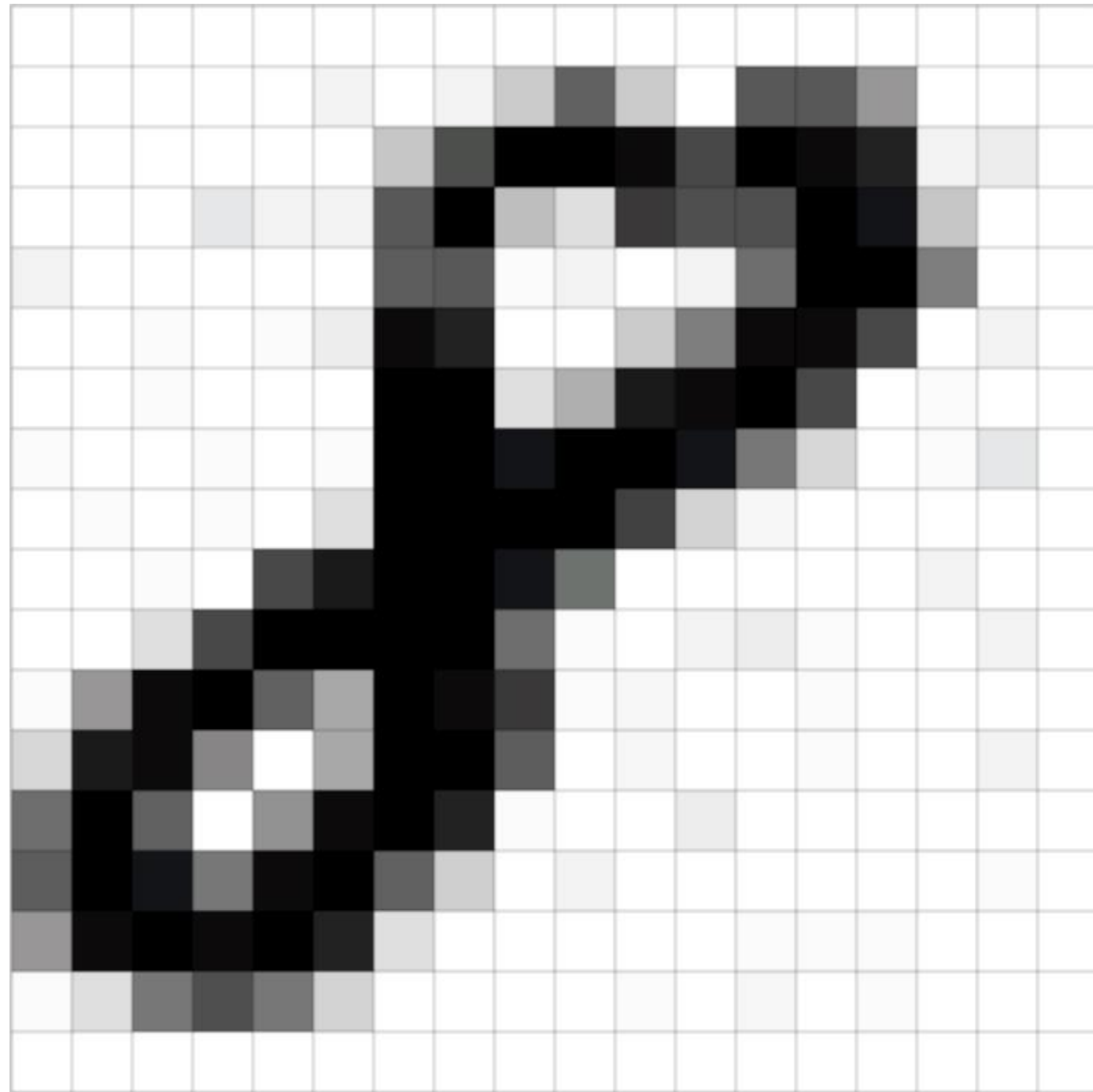
Artificial Intelligence Machine Learning Deep Learning Cognitive Computing



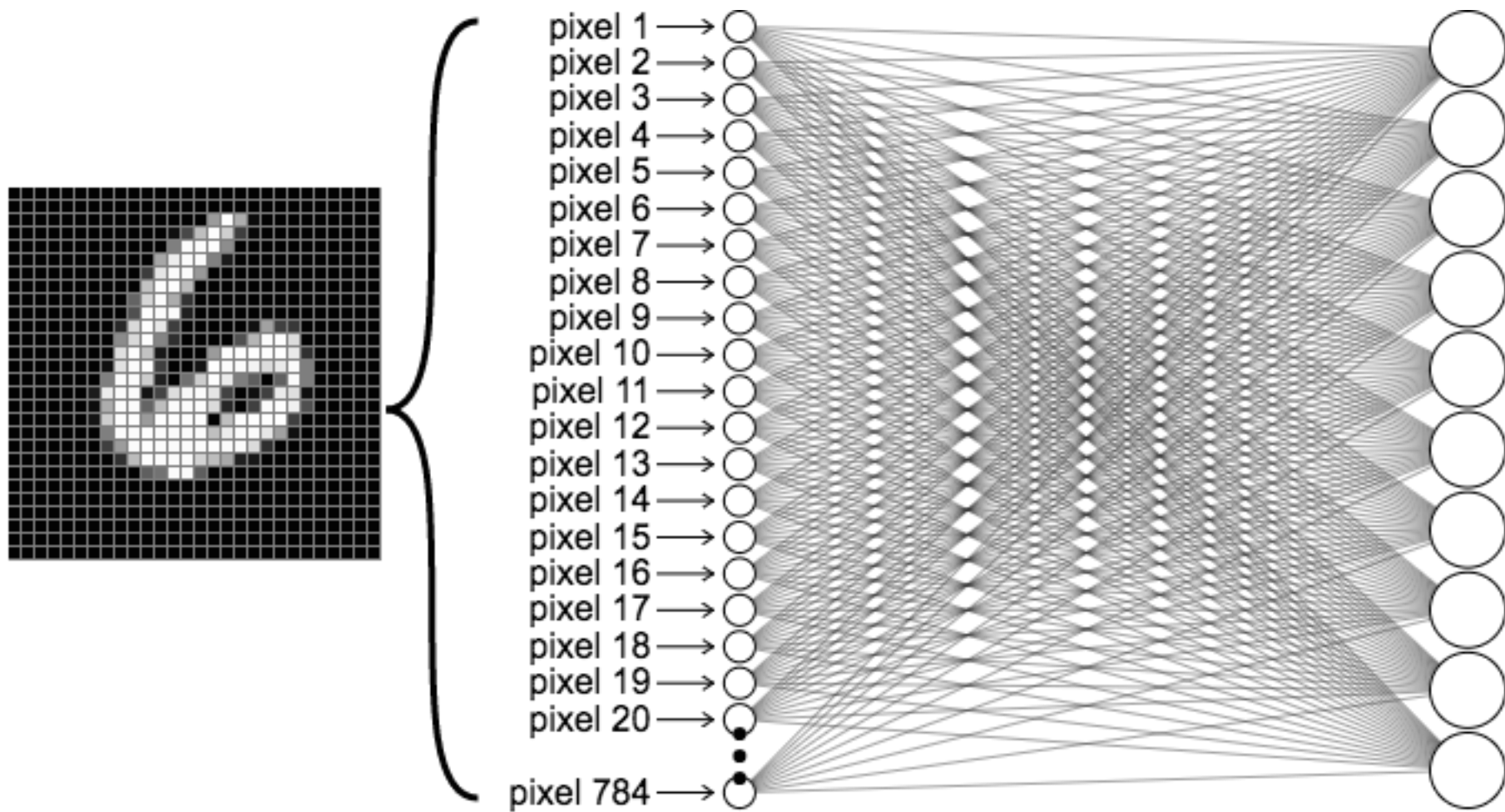
O que é computação cognitiva?

Área da Computação que estuda e desenvolve algoritmos que “tentam” imitar o processo da Inteligência Humana, de aprender, sentir, identificar e resolver problemas.





[illegible]



Escalar

```
s = np.array(8)  
s.shape = ()
```

[8]

Vetor

```
v = np.array([1,2,3])  
v.shape = (3,)
```

[1,2,3]

Matriz

```
m = np.array([[1,2,3], [4,5,6], [7,8,9]])
```

```
m.shape = (3, 3)
```

[1,2,3]

[4,5,6]

[7,8,9]

Tensor

```
t = np.array([[[[1],[2]],[[3],[4]],[[5],[6]]],  
              [[7],[8]],[[9],[10]],[[11],[12]]],  
              [[13],[14]],[[15],[16]],[[17],[17]]])
```

```
t.shape = (3, 3, 2, 1)
```

't'
'e'
'n'
's'
'o'
'r'

Vetor

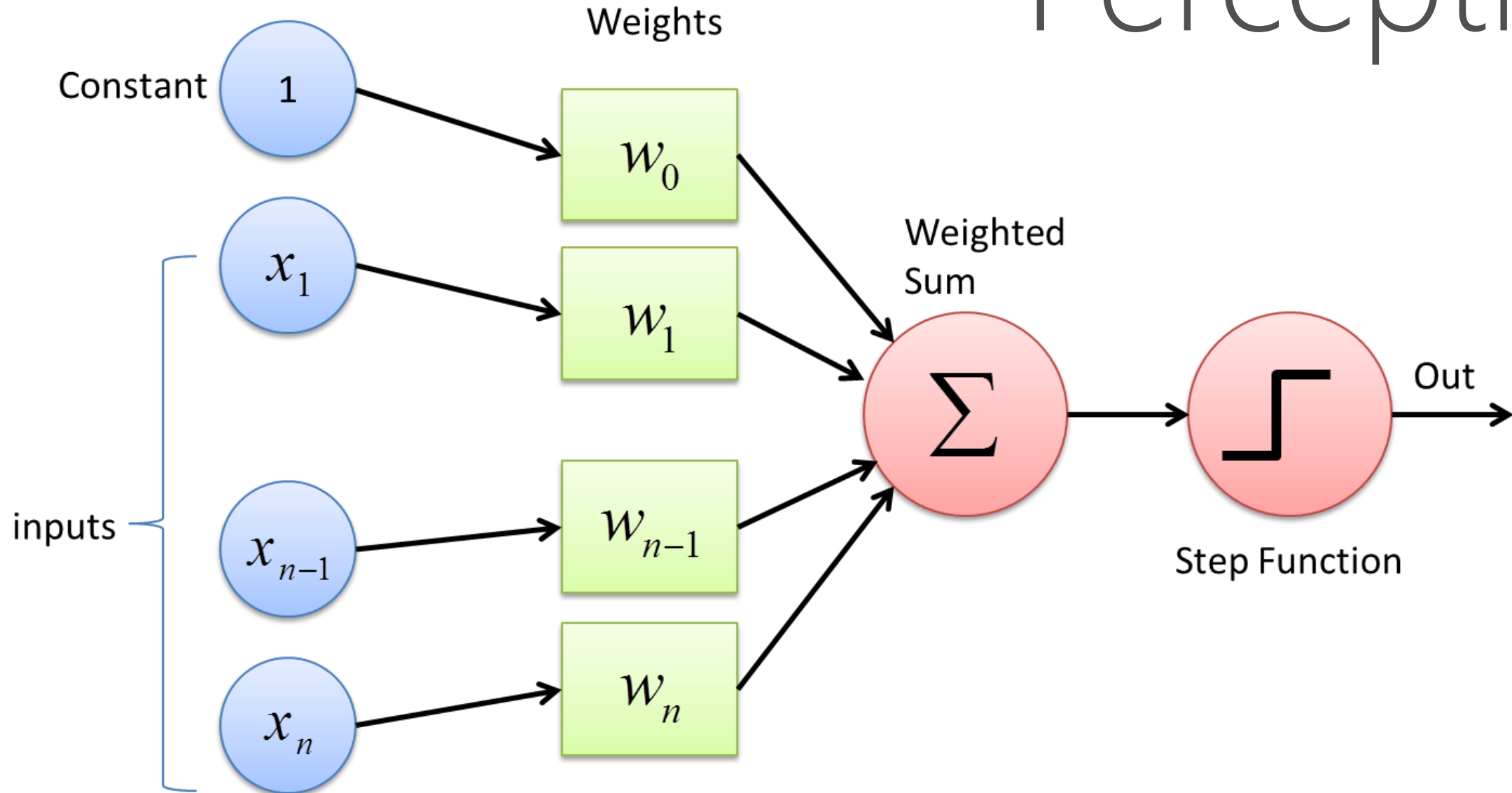
3	1	4	1
5	9	2	6
5	3	5	8
9	7	9	3
2	3	8	4
6	2	6	4

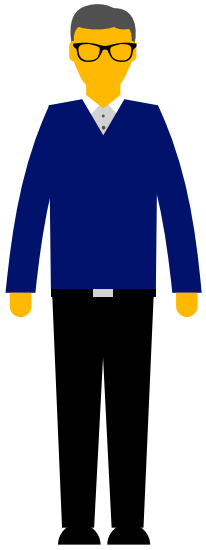
Matriz

2	1	8	8	1	8
2	8	4	5	9	0
2	3	5	3	6	0
7	4	7	1	3	5

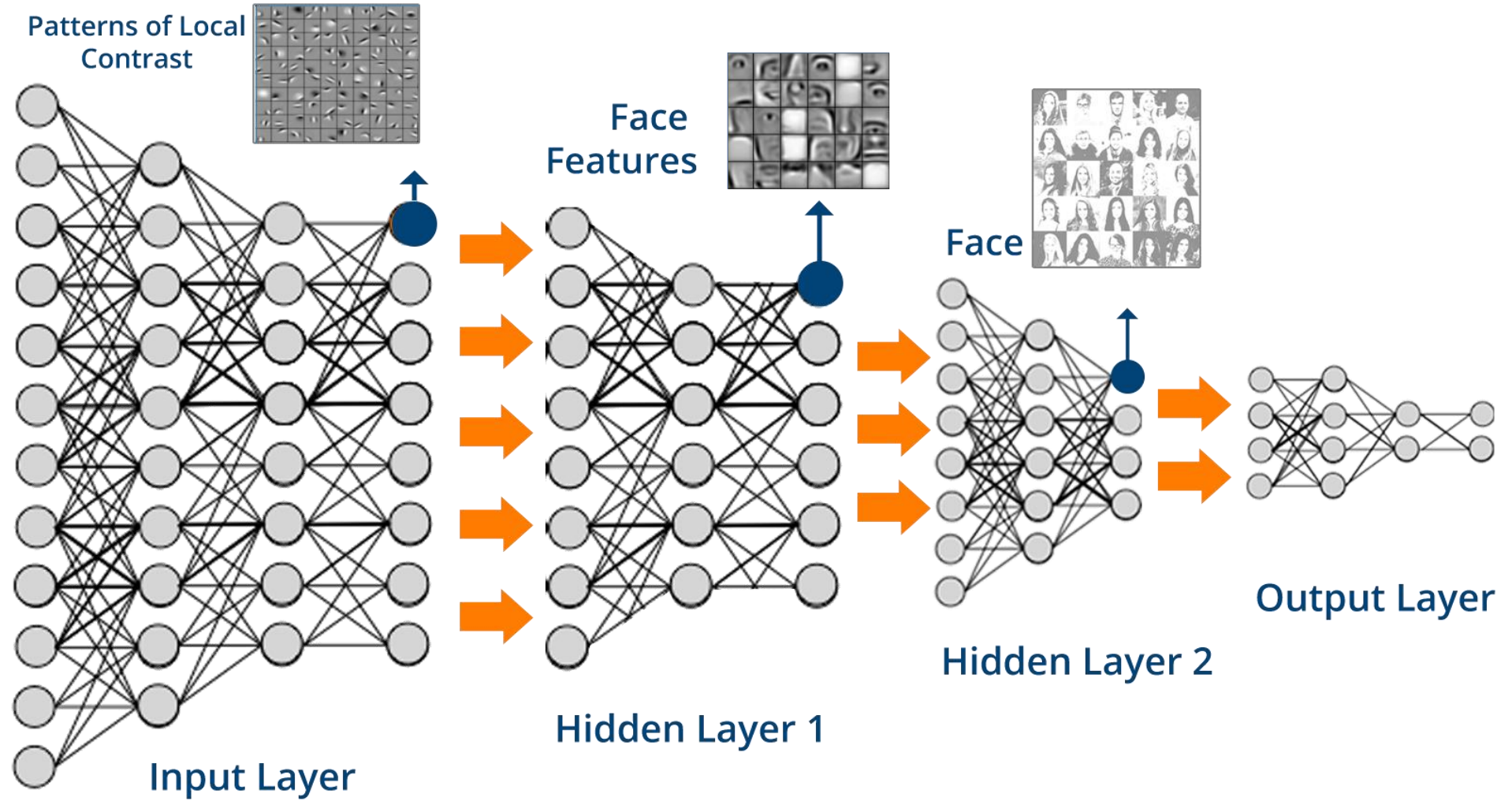
Tensor

Perceptron

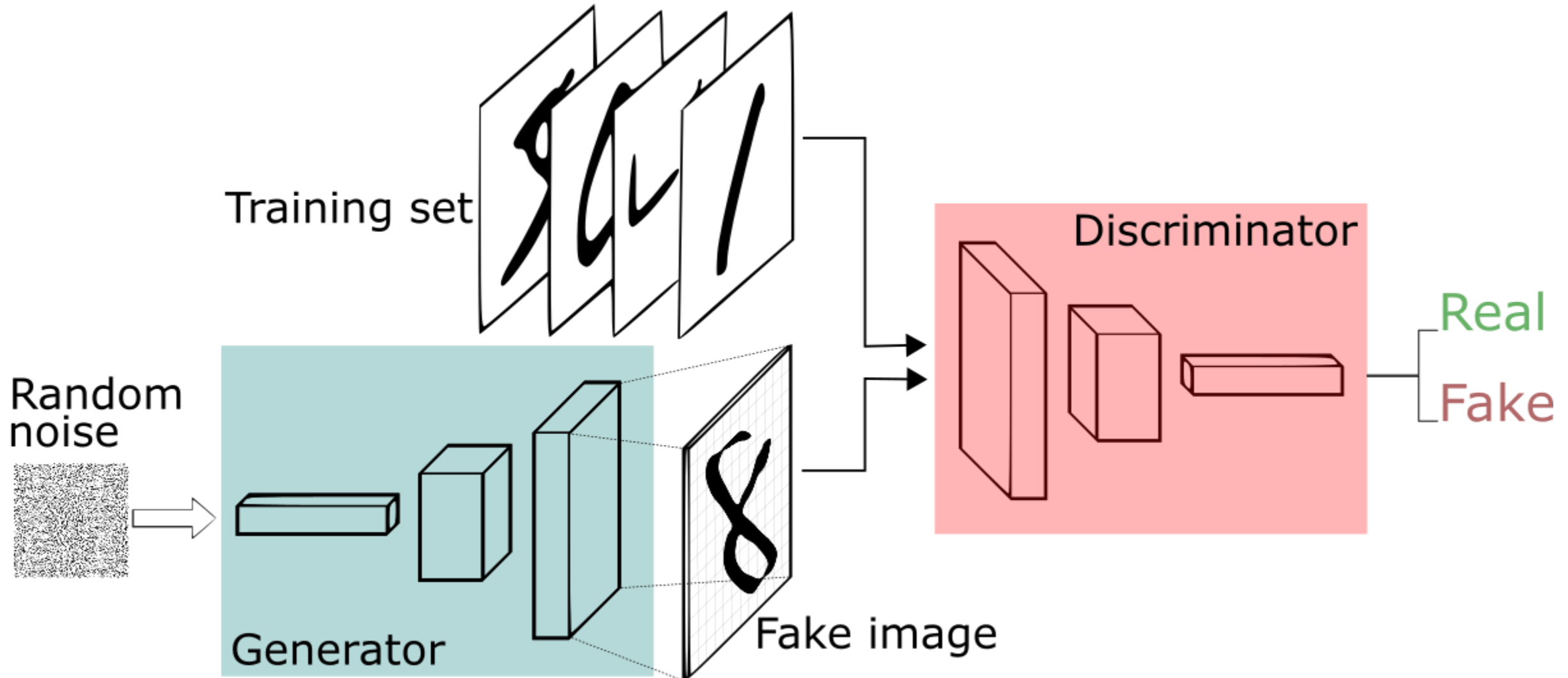




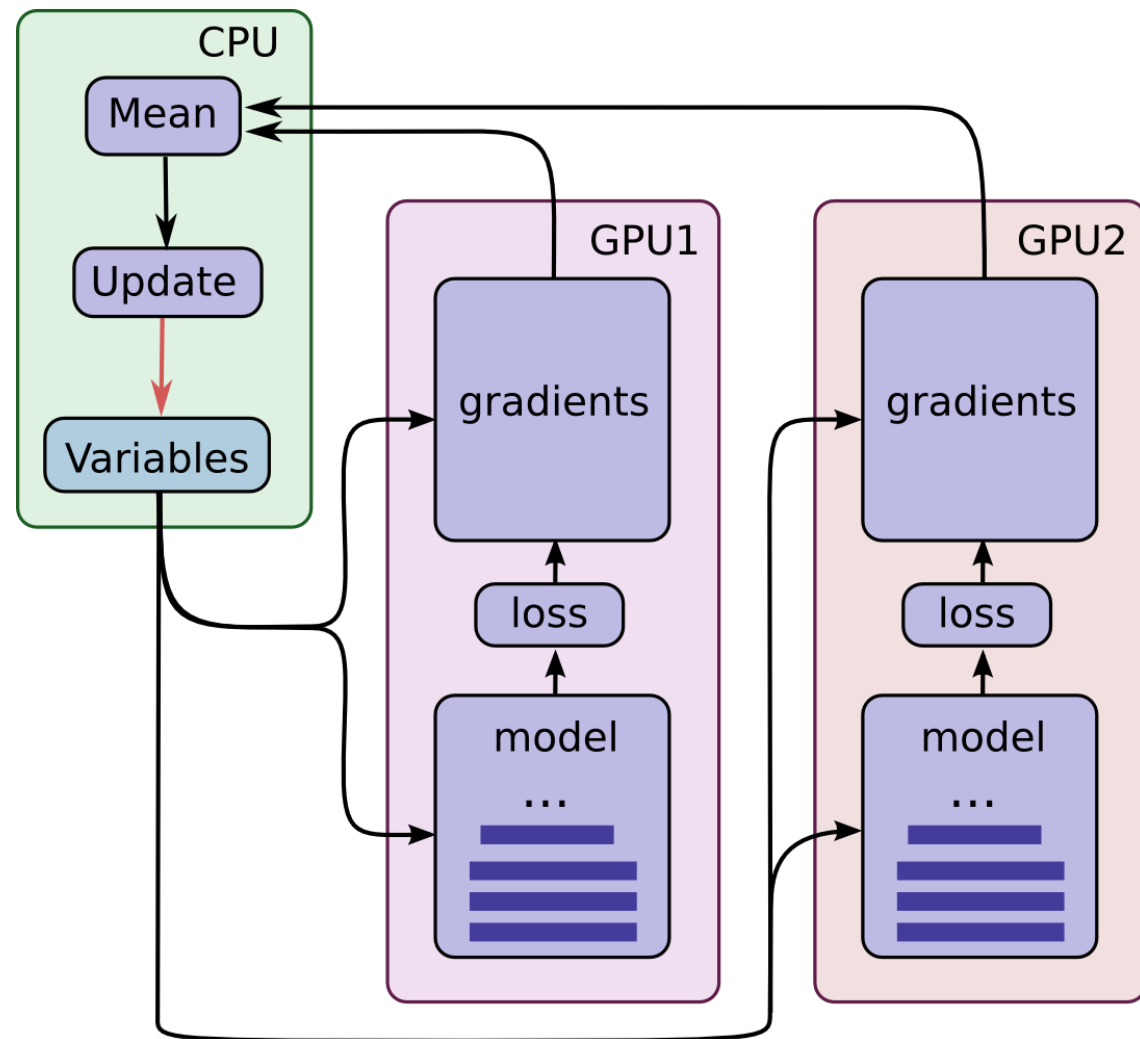
What?



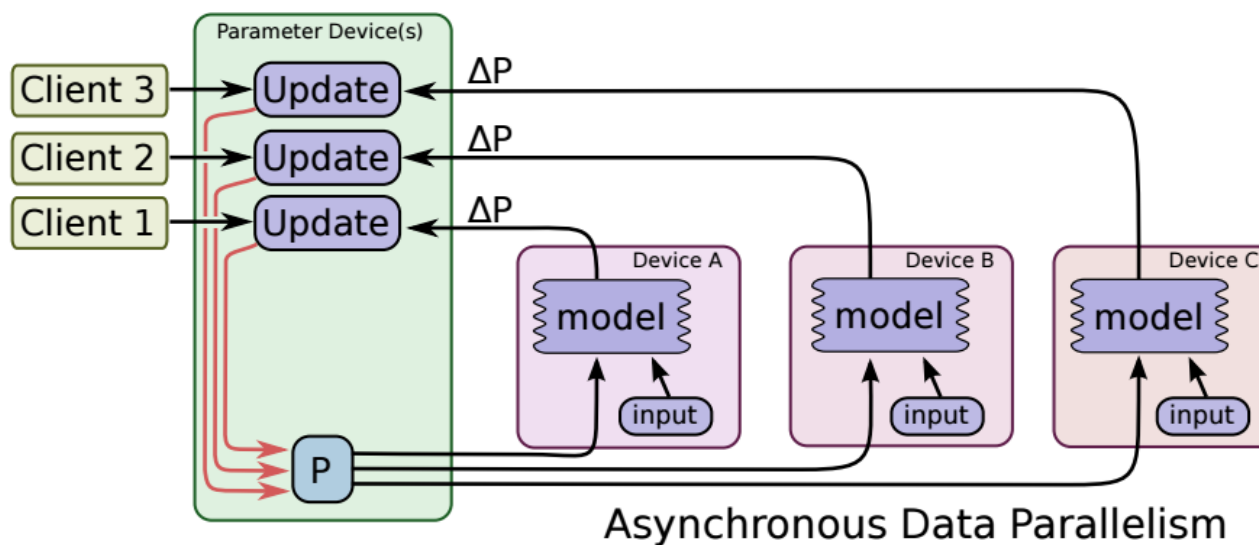
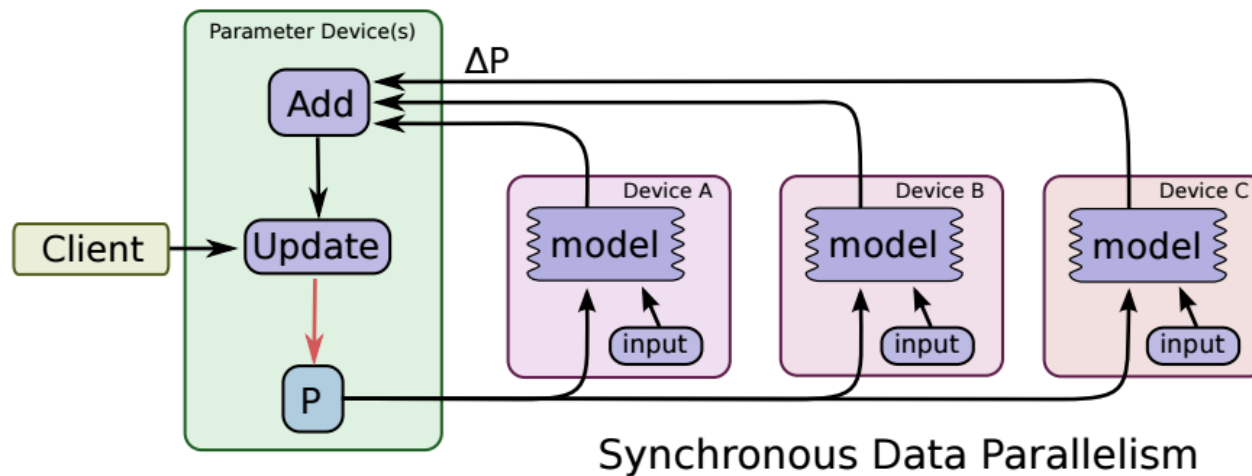
Generative Adversarial Networks



Multiple GPU Cards



Parallel Training



A faster, more efficient, more intelligent cloud

➔ The need for **SCALE**

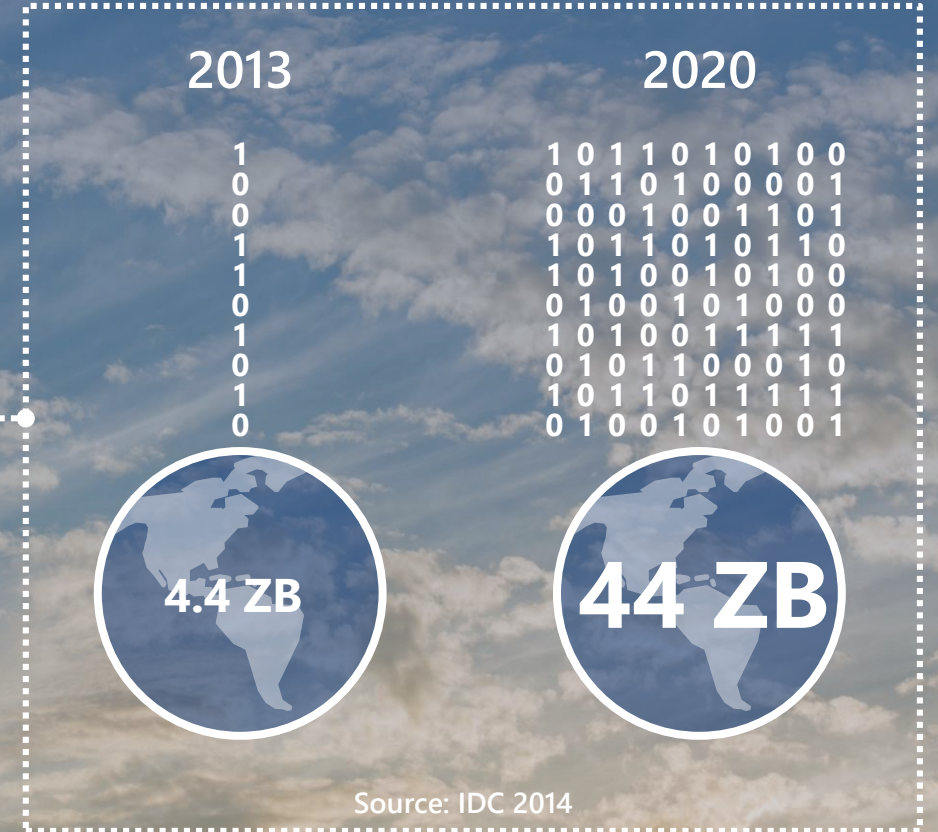
Data explosion: 2013 4.4 ZB - 2020 44 ZB
ML, DNN, AI are driving requirements up faster

➔ The need for **LOW-LATENCY**

Autonomous decision making
Real-time insights into connected devices
Interactive user experiences

➔ The need for **THROUGHPUT**

Cloud-scale services
Searches and recommendations (Indexing the Internet!)

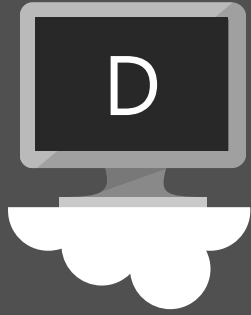


Source: IDC 2014

New Azure VM Sizes



Lowest Price



SSD Storage
Fast CPUs



New generation
of D family VMs



High memory and
Large SSDs



New A-Series



Compute Intensive



NVIDIA GPUs
K80 Compute



NVIDIA GPUs
M60 Visualization



Fastest CPU
IB Connectivity



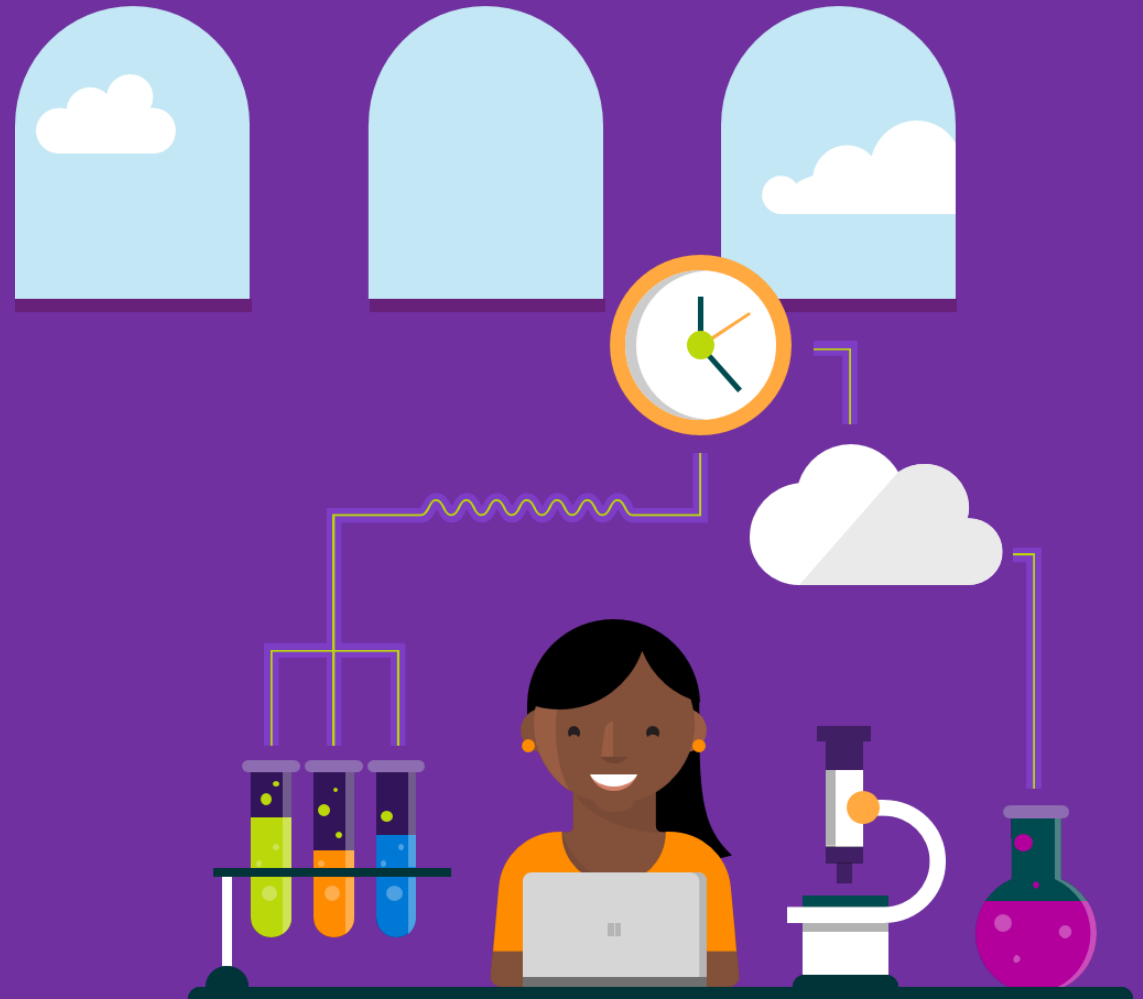
Large SSDs



SAP Large Instances

Azure Batch AI Training

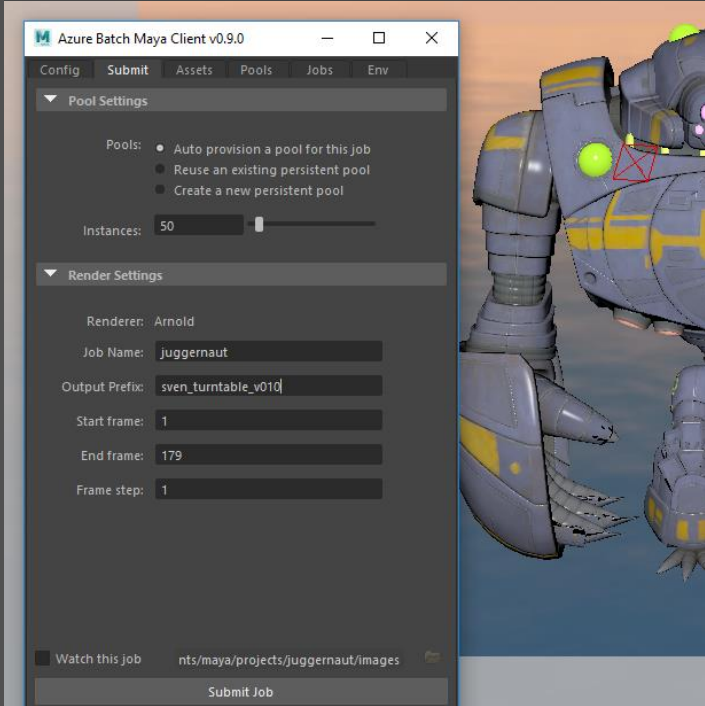
AI Training At
Scale with Azure



Announcing Azure Batch Rendering Service

Autodesk 3ds Max / Maya

Integrated Client Plugin



 Azure Batch

 AUTODESK MAYA

SOLIDANGLE

 3DS MAX

arnold



Monitoring
Reporting
Single bill



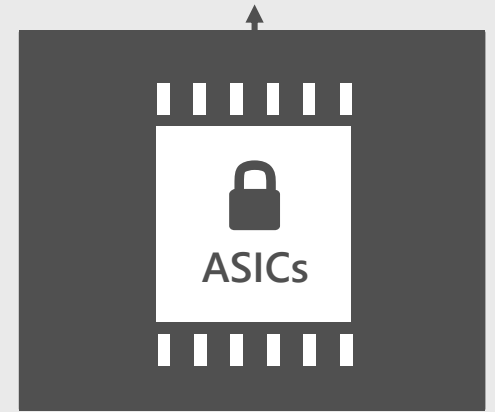
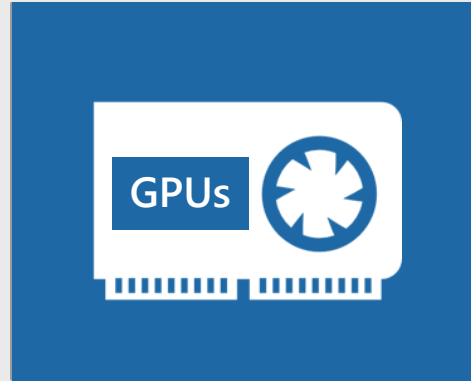
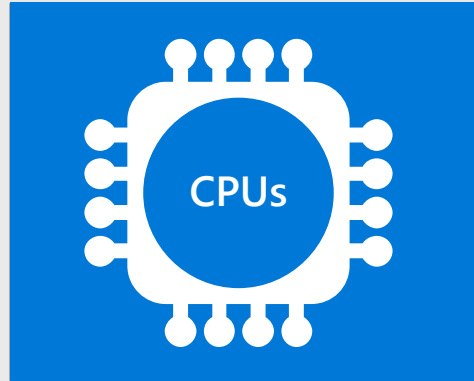
Silicon alternatives

TRAINING

CPUs and GPUs, limited FPGAs,
ASICs under investigation

EVALUATION

CPUs and FPGAs,
ASICs under investigation



The power of deep learning on FPGA

Performance

Tens to hundreds of TOPS of effective inference throughput at low batch sizes
Ultra-low latency serving on modern DNNs
>10X better than CPUs and GPUs
Scale to many FPGAs in single DNN service

Flexibility

FPGAs ideal for adapting to rapidly evolving ML
CNNs, LSTMs, MLPs, reinforcement learning, feature extraction, decision trees, etc.
Inference-optimized numerical precision
Custom binarized, ternarized, tiny precision nets
Sparsity, deep compression for larger, faster models

Scale

Microsoft has the world's largest cloud investment in FPGAs
Multiple Exa-Ops of aggregate AI capacity
We have built powerful DNN serving platform on our FPGA fabric

MUITO OBRIGADO!

VITOR MERIAT
@VITORMERIAT



Visual Studio Summit

#VSSUMMIT