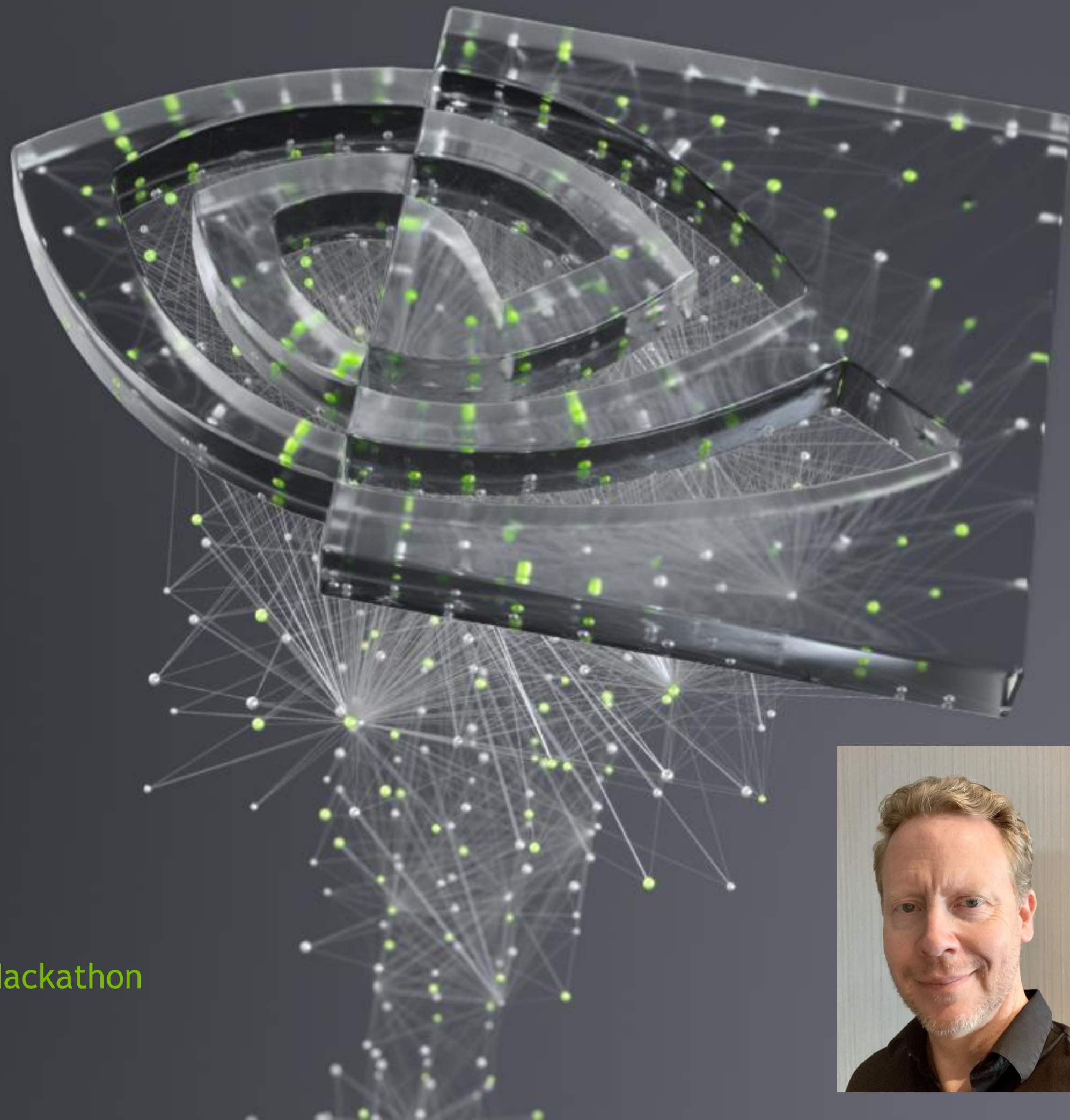




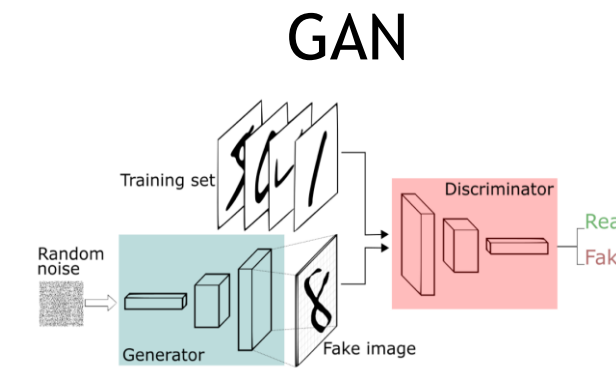
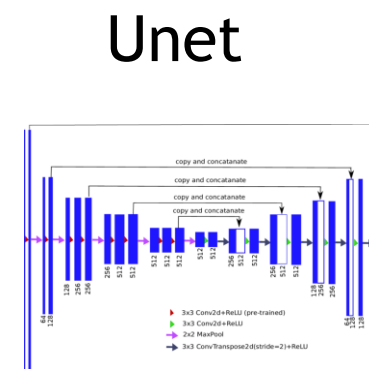
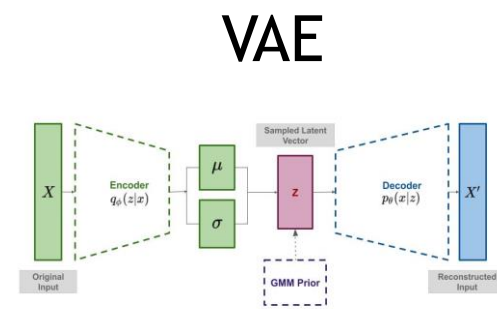
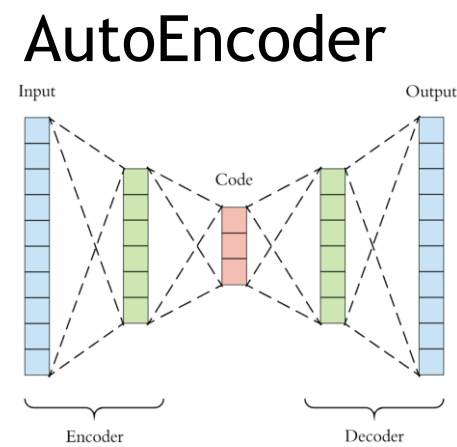
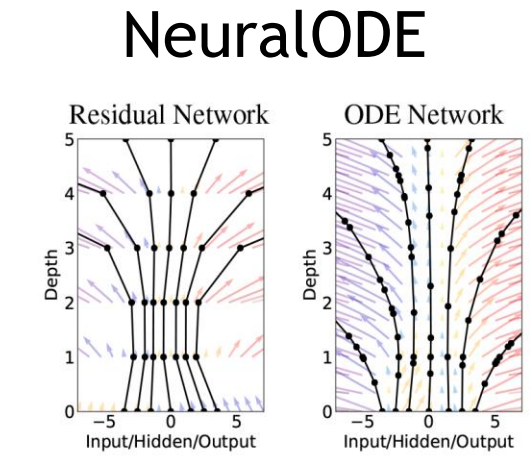
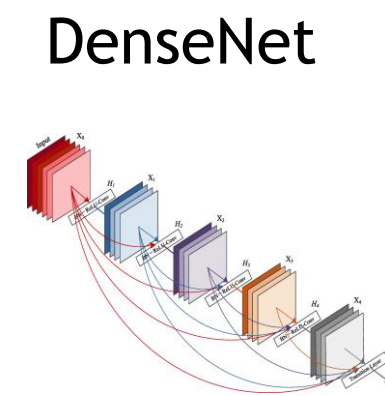
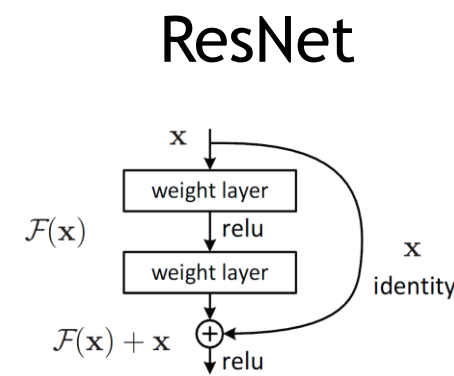
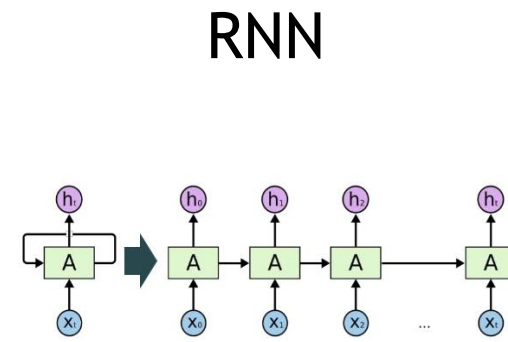
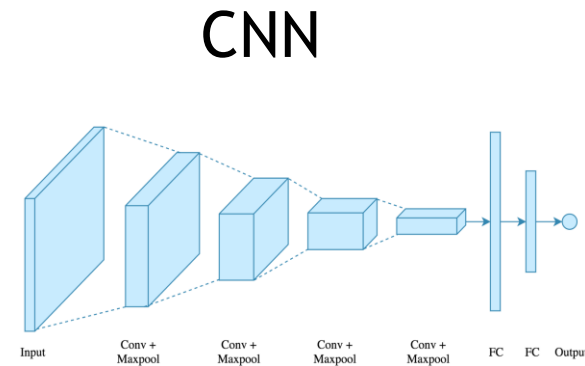
DEEP LEARNING ARCHITECTURES

David M. Hall, June 23, 2020

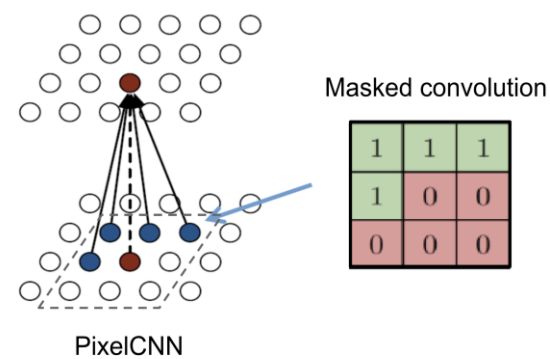
AI4ESS: AI for Earth System Science Workshop and Hackathon



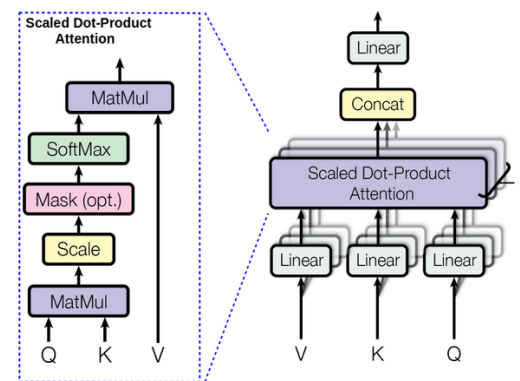
DEEP LEARNING MODEL ZOO



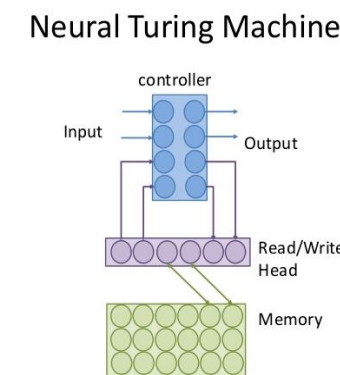
Partial Convolutions



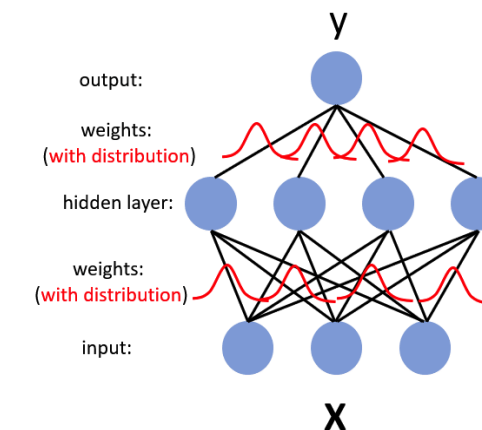
Transformer



Neural Turing Machine

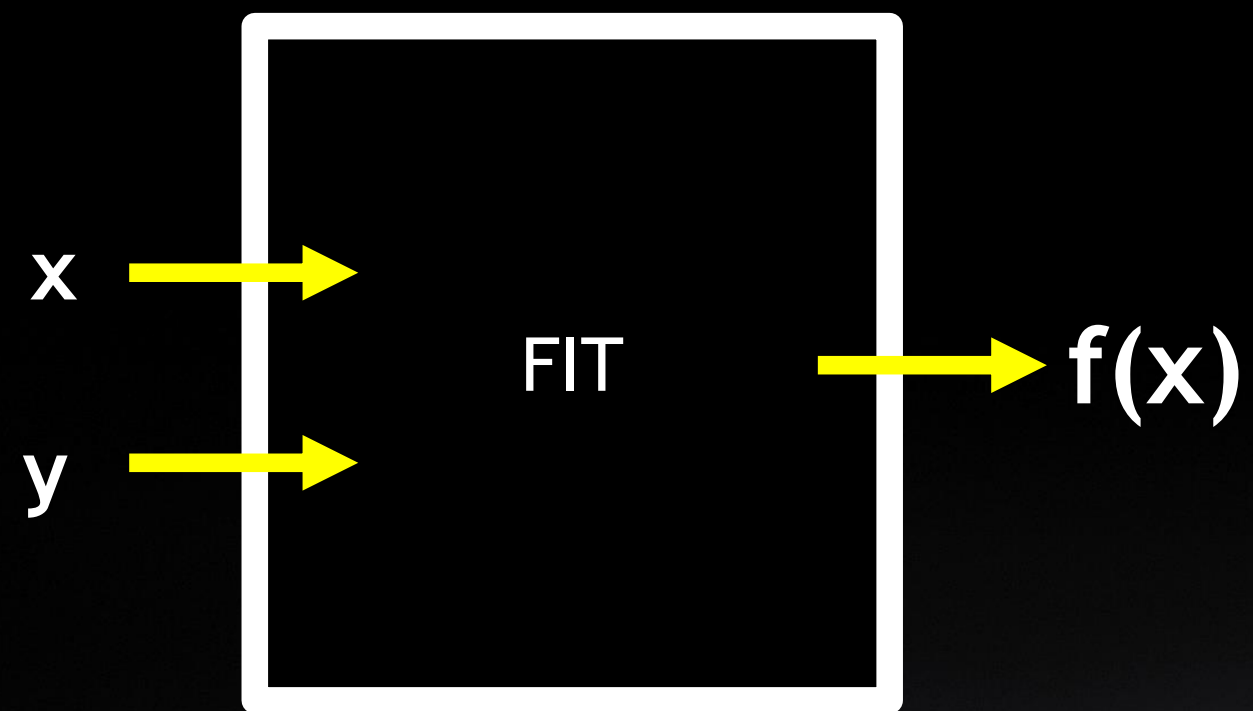


Bayesian Neural Net

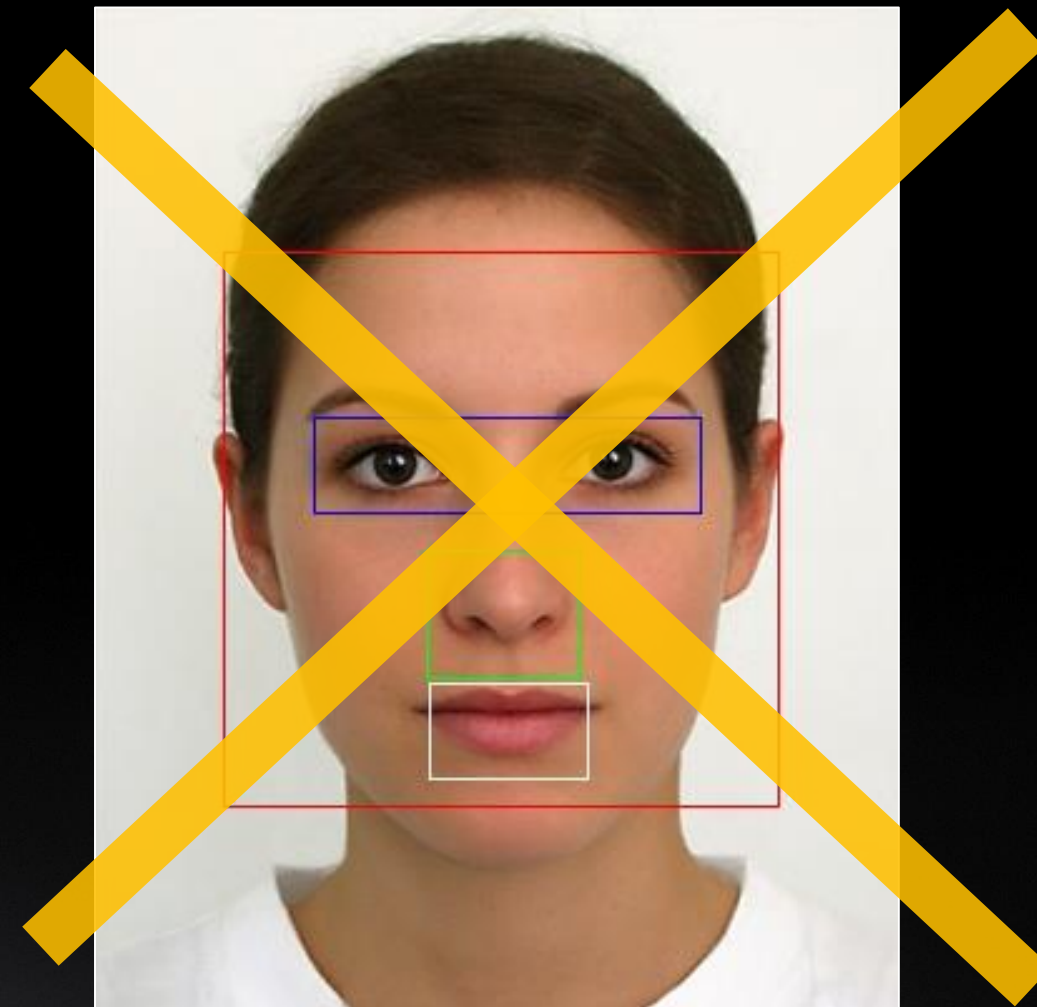


DEEP LEARNING'S CENTRAL PREMISE

LEARN FUNCTIONS
FROM DATA



FEATURE ENGINEERING
NOT REQUIRED



THE MASTER ALGORITHM

Universal Machine Learning is an Ideal, Not Yet a Reality

TOPIC

*"How to build
AGI and to use
it safely"*

*"How to Solve
Global Climate
Change"*



DISSERTATION



TRACKING THE STATE OF THE ART

Arxiv Sanity Preserver

Arxiv Sanity Preserver User: Pass: [Login or Create](#) [Fork me on GitHub](#)
Built in spare time by @karpathy to accelerate research.
Serving last 110403 papers from cs.[CV|CL|LG|AI|NE]/stat.ML

Search:

most recent | top recent | top hype | friends | discussions | recommended | library

Only show v1 | Last day | Last 3 days | Last week | Last month | Last year | All time

Top papers based on people's libraries:

Attention Is All You Need 1706.03762v5 pdf
Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin
12/6/2017 (v1: 6/12/2017) cs.CL | cs.LG
15 pages, 5 figures

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

NIPS 2016 Tutorial: Generative Adversarial Networks 1701.00160v4 pdf
Ian Goodfellow
4/3/2017 (v1: 12/31/2016) cs.LG
v2-v4 are all typo fixes. No substantive changes relative to v1

This report summarizes the tutorial presented by the author at NIPS 2016 on generative adversarial networks (GANs). The tutorial describes: (1) Why generative modeling is a topic worth studying, (2) how generative models work, and how GANs compare to other generative models, (3) the details of how GANs work, (4) research frontiers in GANs, and (5) state-of-the-art image models that combine GANs with other methods. Finally, the tutorial contains three exercises for readers to complete, and the solutions to these exercises.

Mask R-CNN 1703.06870v3 pdf
Kaiming He, Georgia Gkioxari, Piotr Dollár, Ross Girshick
1/24/2018 (v1: 3/20/2017) cs.CV
open source; appendix on more results

We present a conceptually simple, flexible, and general framework for object instance segmentation. Our approach

Papers With Code

Browse State-of-the-Art
2718 leaderboards • 1579 tasks • 2440 datasets • 24380 papers with code
[Follow on Twitter for updates](#)

Computer Vision

- Semantic Segmentation**: 59 leaderboards, 963 papers with code
- Image Classification**: 136 leaderboards, 802 papers with code
- Object Detection**: 79 leaderboards, 684 papers with code
- Image Generation**: 100 leaderboards, 330 papers with code
- Pose Estimation**: 76 leaderboards, 318 papers with code

[See all 789 tasks](#)

Natural Language Processing

- Machine Translation**: 45 leaderboards, 626 papers with code
- Language Modelling**: 14 leaderboards, 603 papers with code
- Question Answering**: 53 leaderboards, 562 papers with code
- Sentiment Analysis**: 37 leaderboards, 395 papers with code
- Text Classification**: 50 leaderboards, 237 papers with code

[See all 297 tasks](#)

Medical

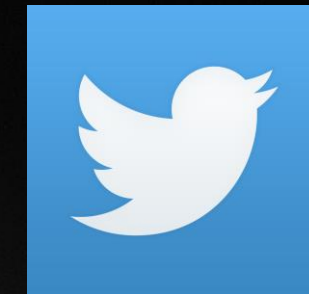
- Medical Image Segmentation**: 52 leaderboards, 84 papers with code
- Drug Discovery**: 13 leaderboards, 53 papers with code
- Lesion Segmentation**: 5 leaderboards, 49 papers with code
- Brain Tumor Segmentation**: 7 leaderboards, 27 papers with code
- Brain Segmentation**: 1 leaderboard, 22 papers with code

[See all 186 tasks](#)

Methodology

- Representation Learning**: 7 leaderboards, 557 papers with code
- Transfer Learning**: 8 leaderboards, 505 papers with code
- Word Embeddings**: 459 papers with code
- Domain Adaptation**: 30 leaderboards, 366 papers with code
- Data Augmentation**: 334 papers with code

[See all 124 tasks](#)





AGENDA

Deep Learning Basics

Fully Connected Networks

CNNs

ResNets

Encoder-Decoders

Masked Convolutions

Generative Models

Transformers

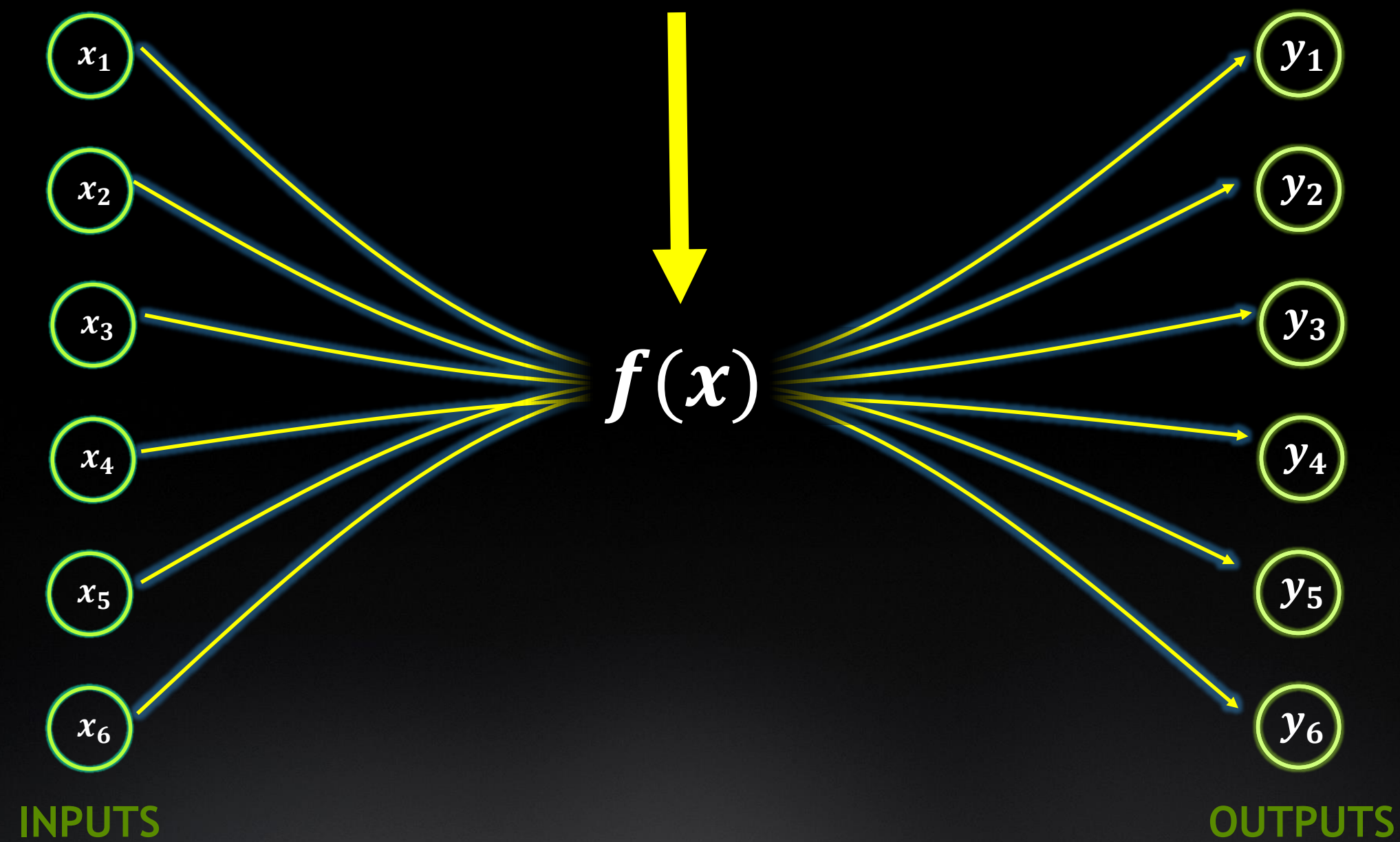
AutoML



DEEP LEARNING BASICS

REVERSE-ENGINEER FUNCTIONS FROM EXAMPLES

Find this, automatically



IT'S A NEW WAY TO BUILD SOFTWARE

TEMP, PRESSURE, MOISTURE



PROBABILITY OF RAIN

HAND-WRITTEN FUNCTION

```
Function1(T,P,Q)
update_mass()
update_momentum()
update_energy()
do_macrophysics()
do_microphysics()
y = get_precipitation()
return y
```

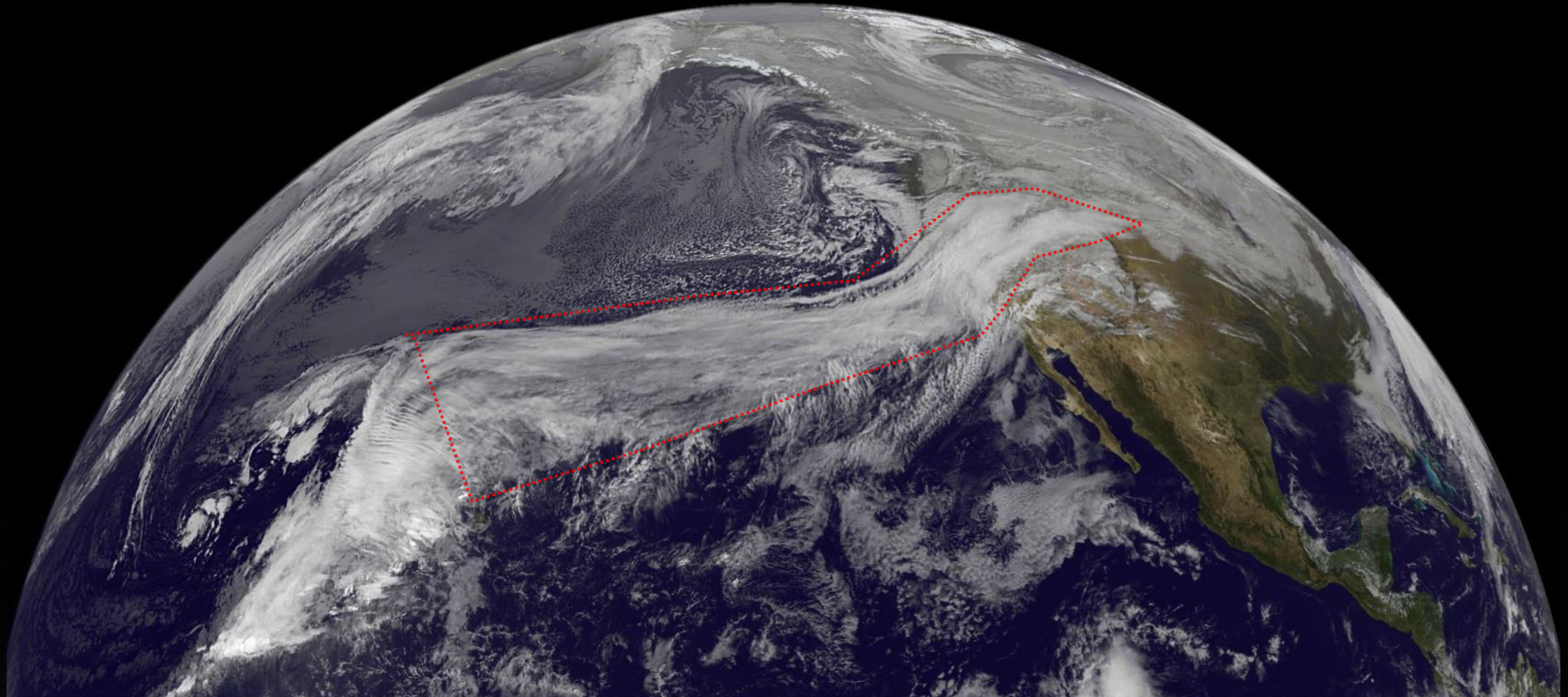
Convert expert knowledge into a function

LEARNED FUNCTION

```
Function1(T,P,Q)
A = relu( w1 * [T,P,Q] + b1)
B = relu( w2 * A + b2)
C = relu( w3 * B + b3)
D = relu( w4 * C + b4)
E = relu( w5 * D + b5)
y = sigmoid(w6 * E + b6)
return y
```

Reverse-engineer a function from inputs / outputs

COMPLEX PHENOMENA ARE BEST DESCRIBED IMPLICITLY

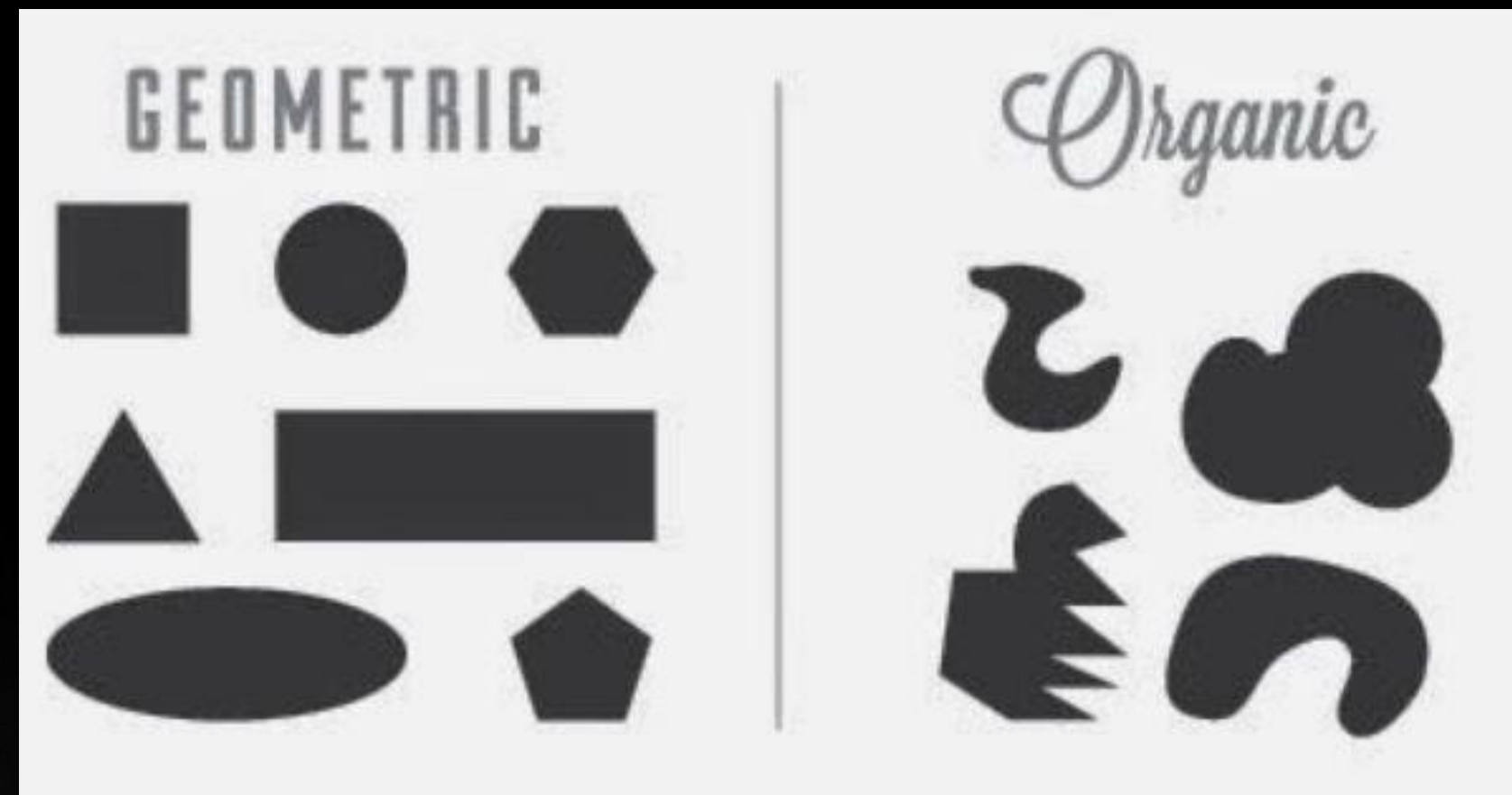


EXAMPLE: ATMOSPHERIC RIVER

FORWARD AND REVERSE ENGINEERING ARE COMPLIMENTARY

SOFTWARE DEVELOPMENT

ENGINEERED
PROGRAMMED
LABOR INTENSIVE
EXPLICIT
EXPLAINABLE
HEURISTIC
SIMPLE
FROM EXPERTISE

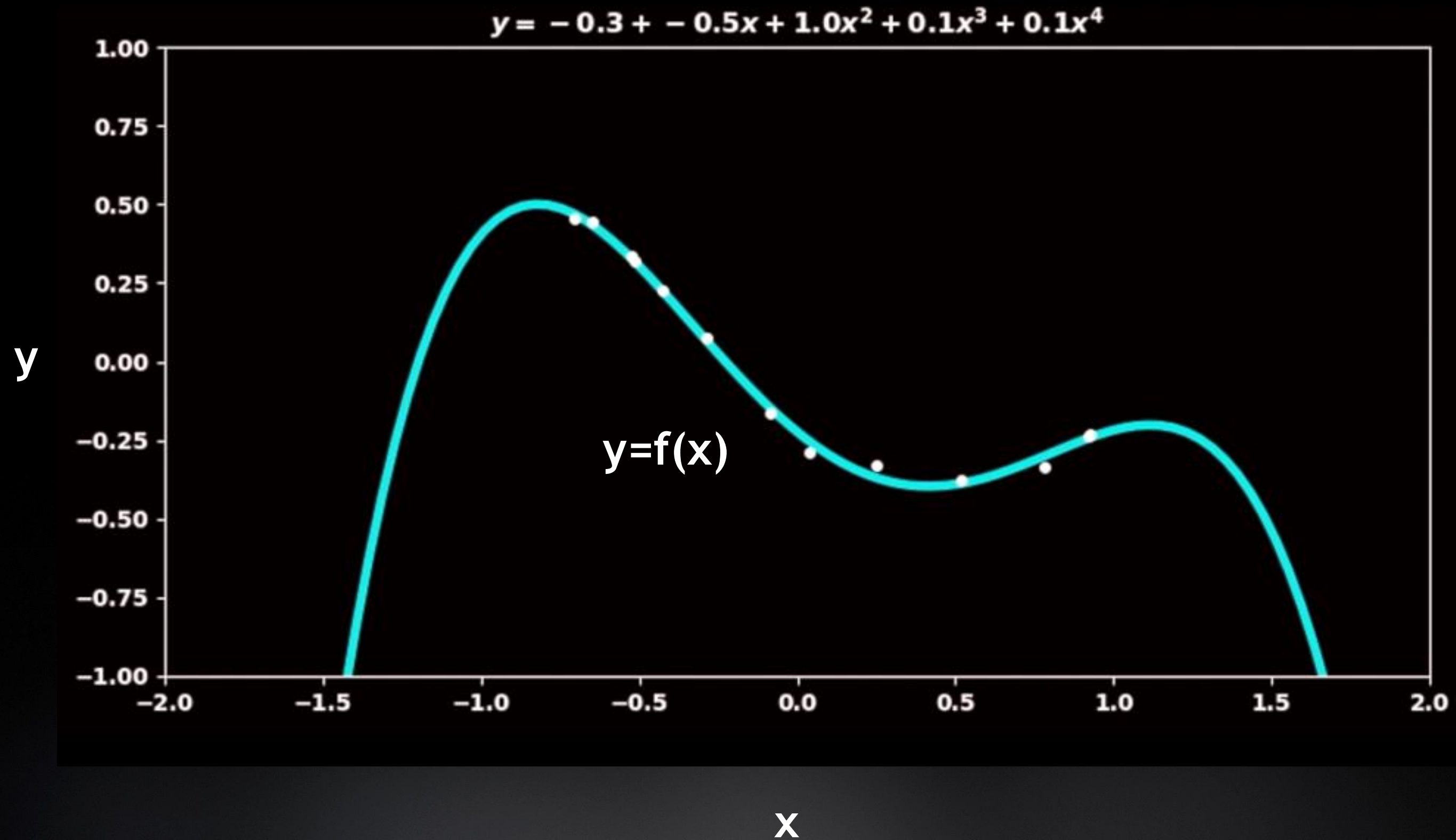


MACHINE LEARNING

REVERSE ENGINEERED
LEARNED
AUTOMATIC
IMPLICIT
SUBTLE
REALISTIC
COMPLEX
FROM EXAMPLES

For best results, combine them

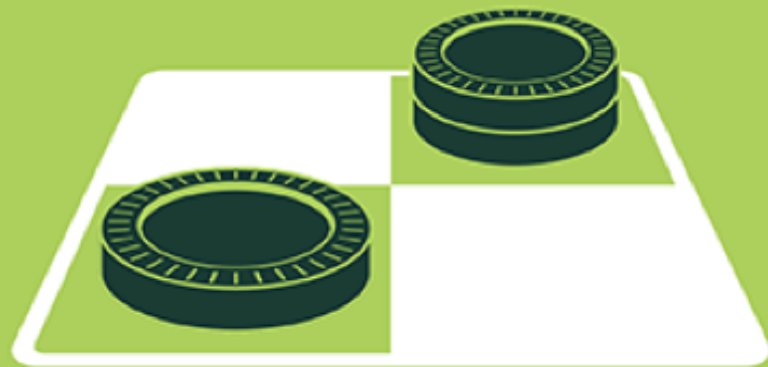
MACHINE LEARNING IS CURVE FITTING, GENERALIZED



AI, MACHINE LEARNING, DEEP LEARNING

ARTIFICIAL INTELLIGENCE

EXPERT SYSTEMS
EXECUTE HAND-WRITTEN ALGORITHMS AT HIGH SPEED



MACHINE LEARNING

TRADITIONAL ML
LEARN FROM EXAMPLES USING HAND-CRAFTED FEATURES



DEEP LEARNING

LEARNS BOTH OUTPUT AND FEATURES FROM DATA



1950's

1960's

1970's

1980's

1990's

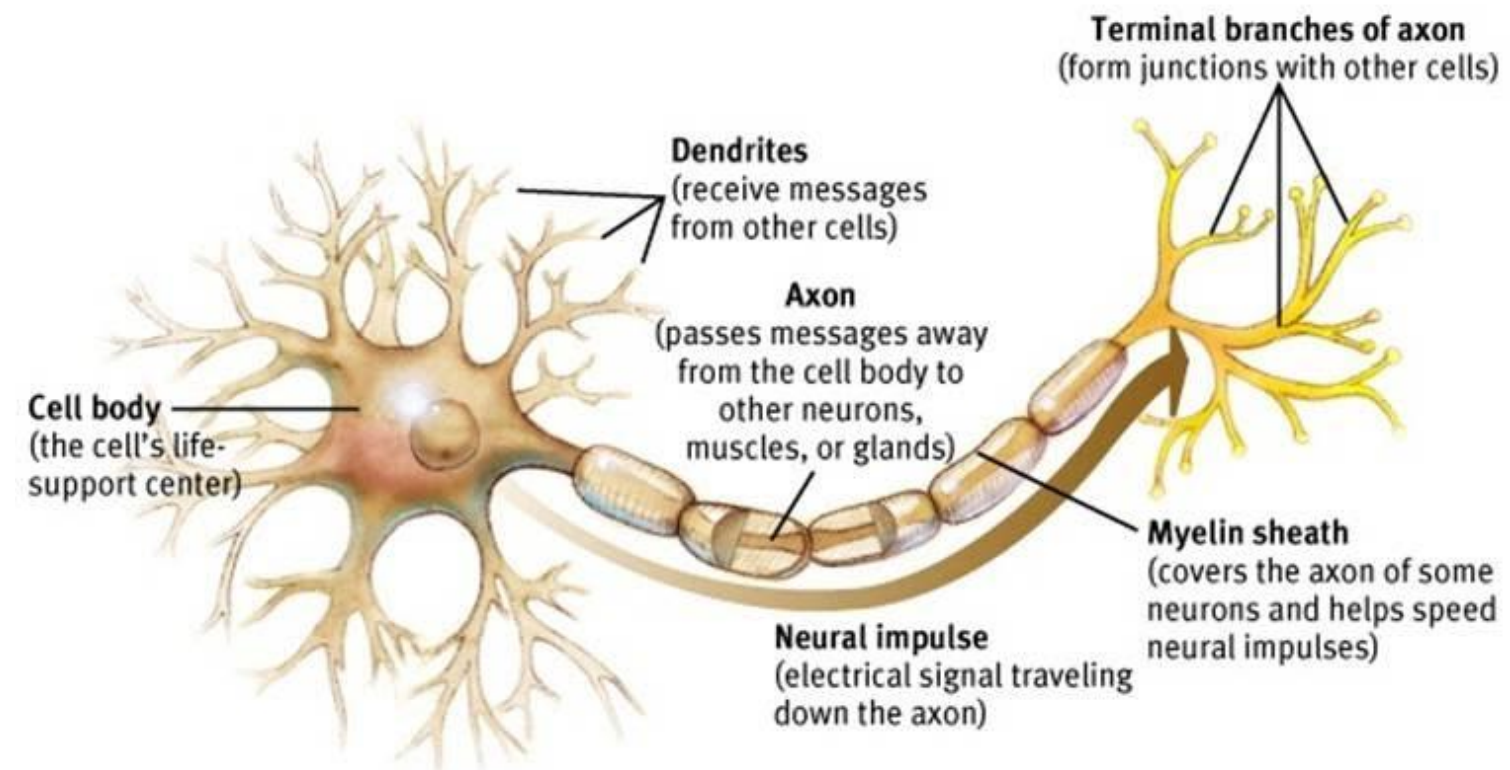
2000's

2010's

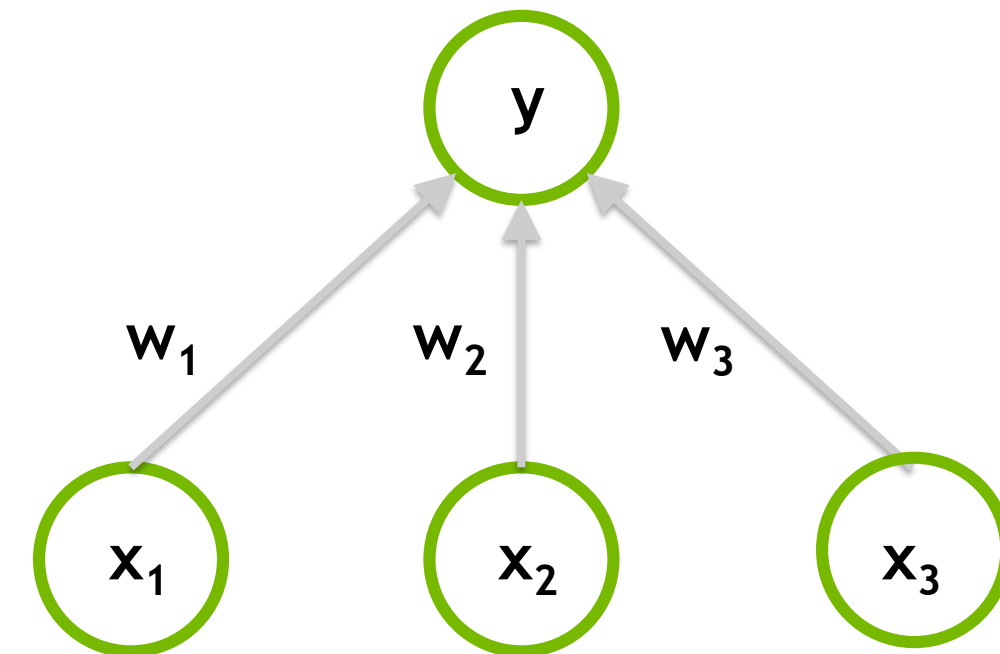
ARTIFICIAL NEURONS

Are simple equations with a set of adjustable parameters

Biological neuron



Artificial neuron

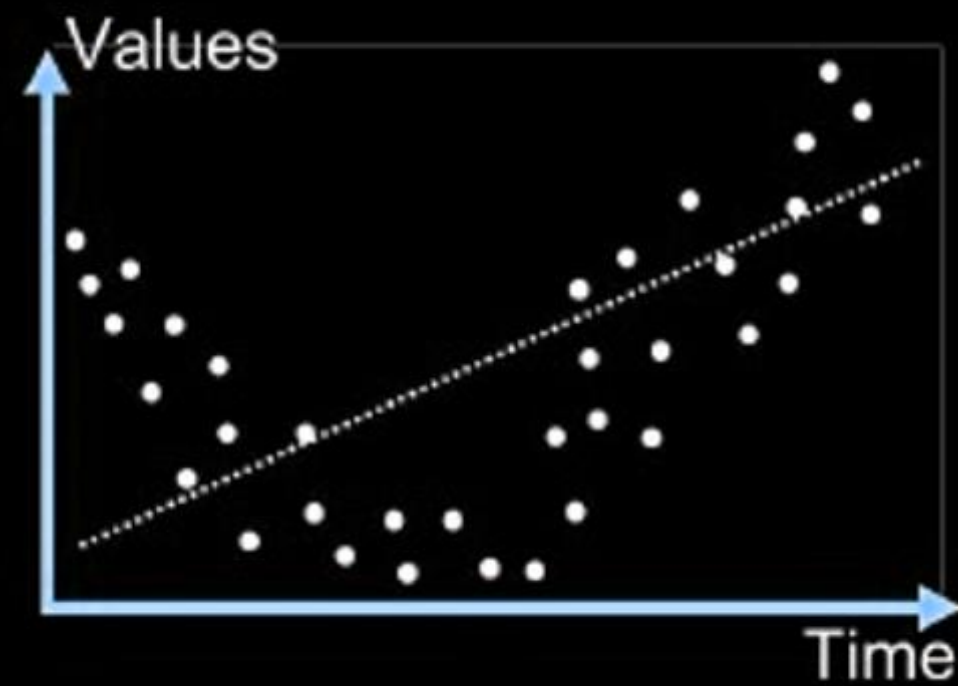


$$y = f(w_1x_1 + w_2x_2 + w_3x_3)$$

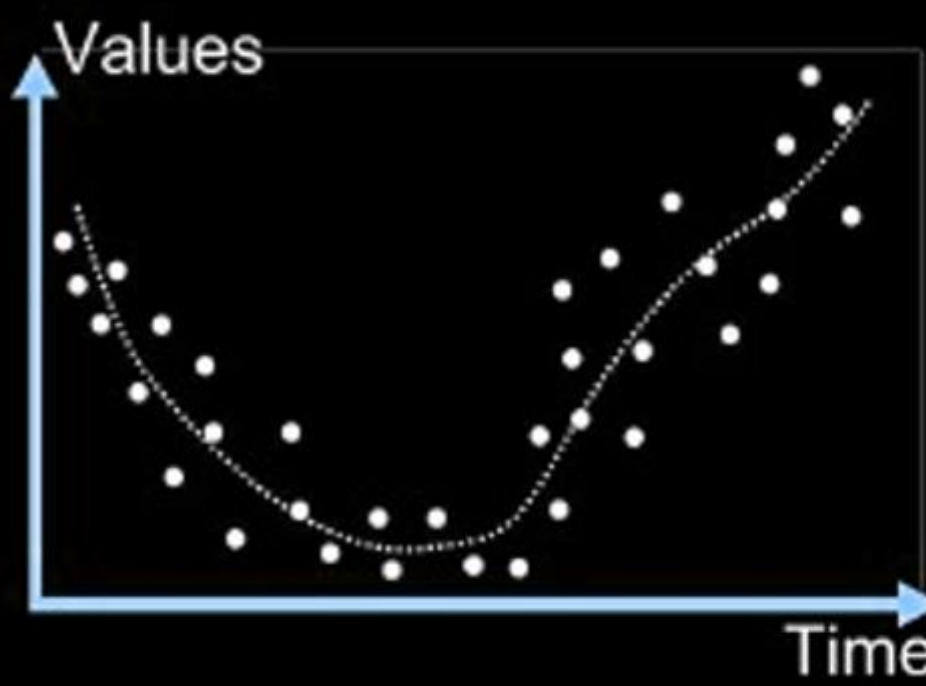
<https://towardsdatascience.com/the-differences-between-artificial-and-biological-neural-networks-a8b46db828b7>

ADJUST MODEL CAPACITY TO FIT YOUR DATA

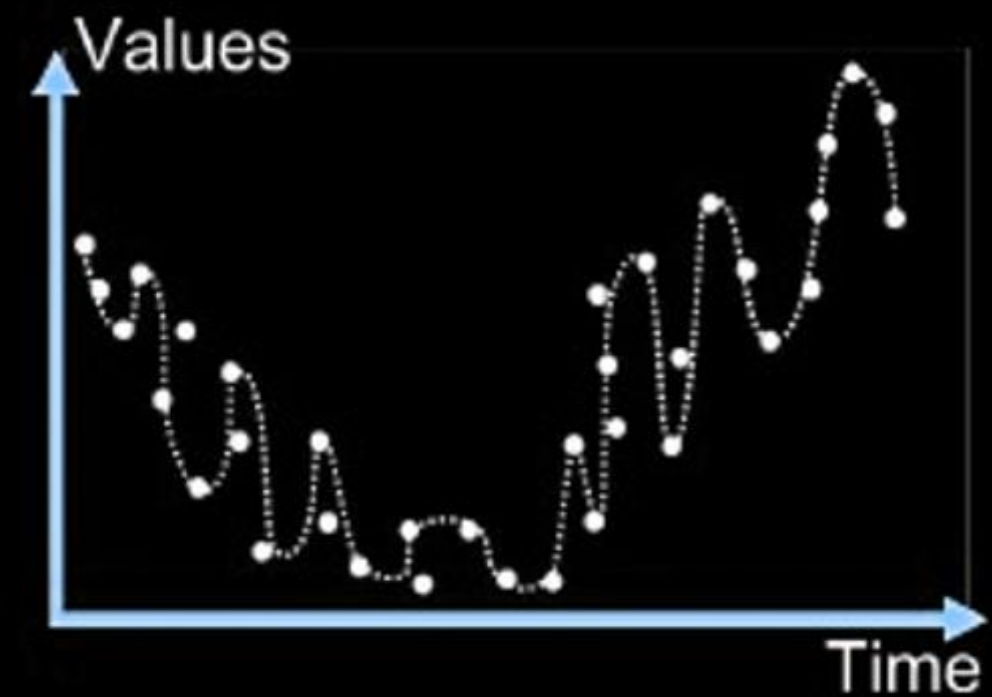
A good model is one that generalizes to new data



UNDERFIT



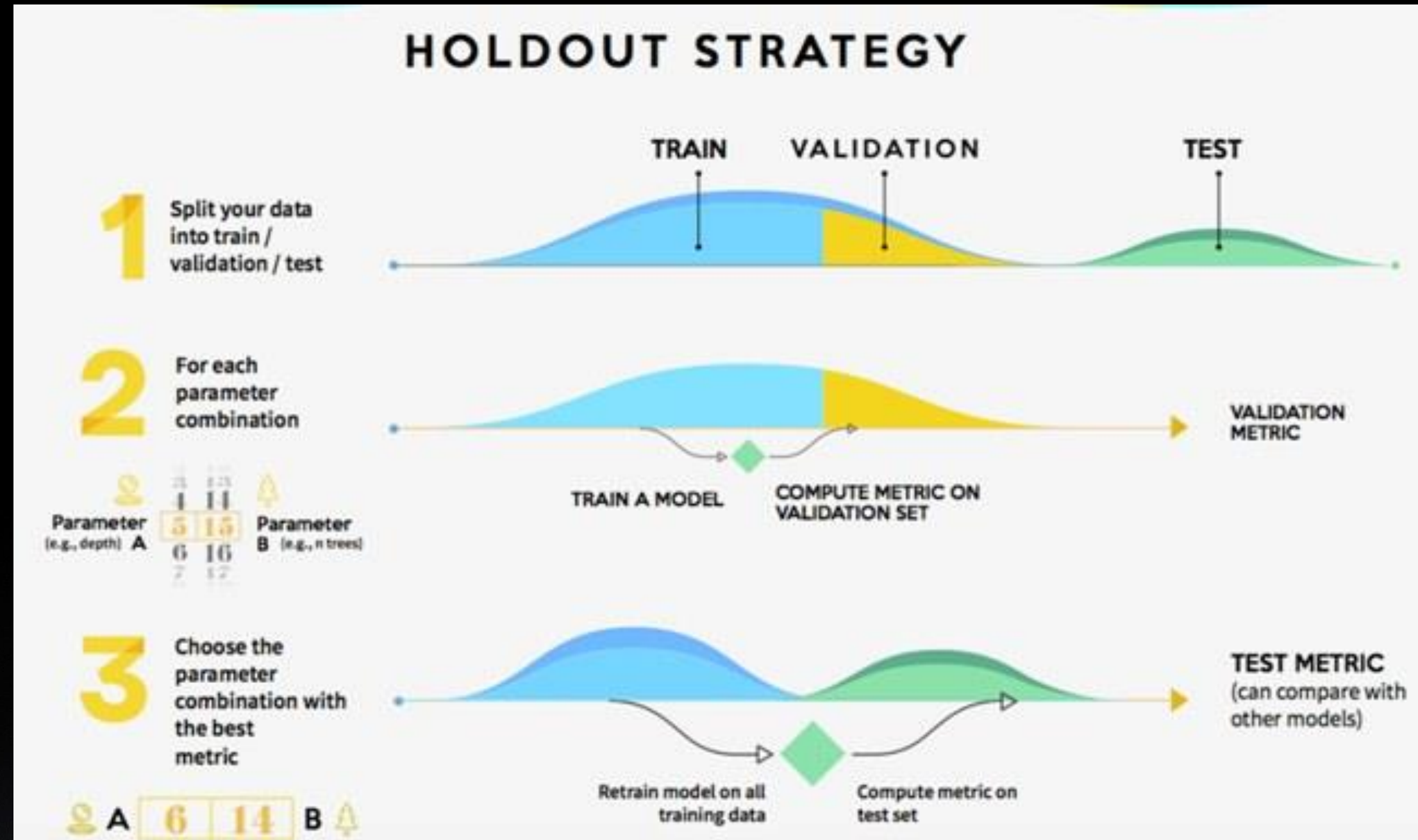
GOOD FIT



OVER FIT

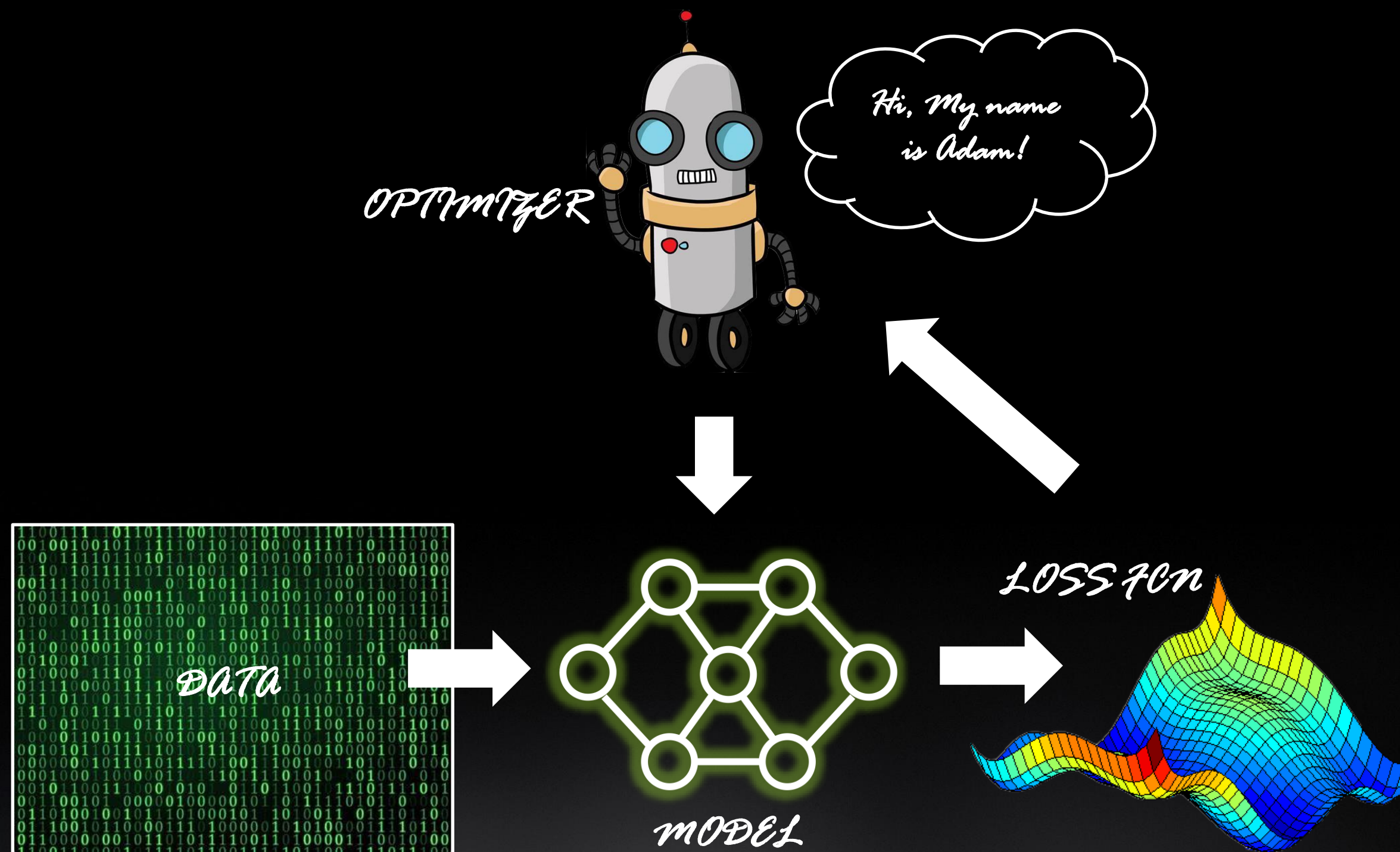
KEEP TRAINING AND TEST DATA SEPERATE

KEEP TEST, TRAINING, AND VALIDATION DATA SEPERATE



TRAINING

DATA, MODEL, LOSS, AND OPTIMIZER



TRAINING: SEARCHING FOR A GOOD SOLUTION

Model training is a form of search, performed by the optimizer.

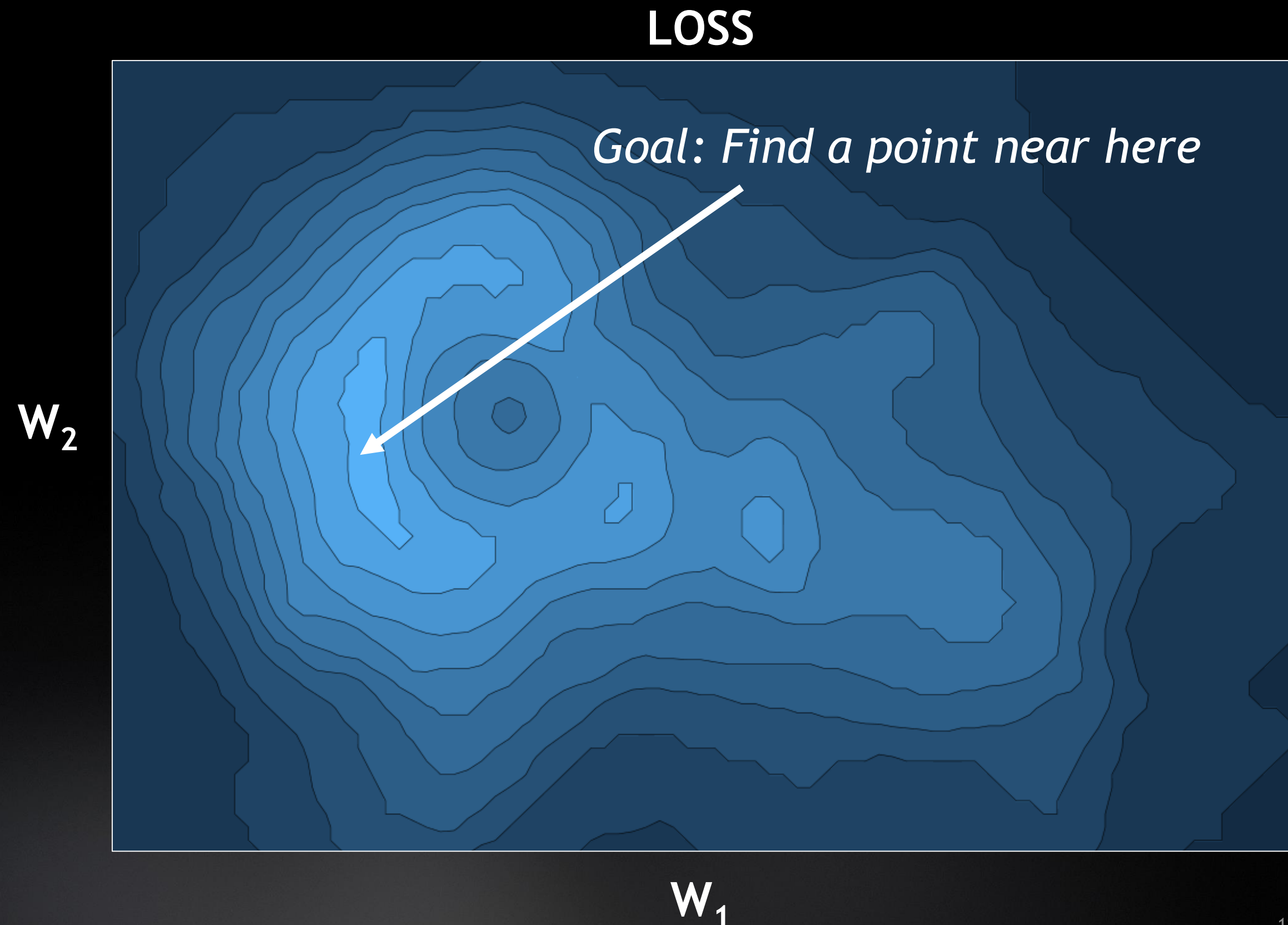
Adjust W_1, W_2 to minimize loss

COULD USE:

- Grid Search
- Evolutionary Algorithms
- Conjugate Gradient
- Newton's method
- Other 2nd order methods

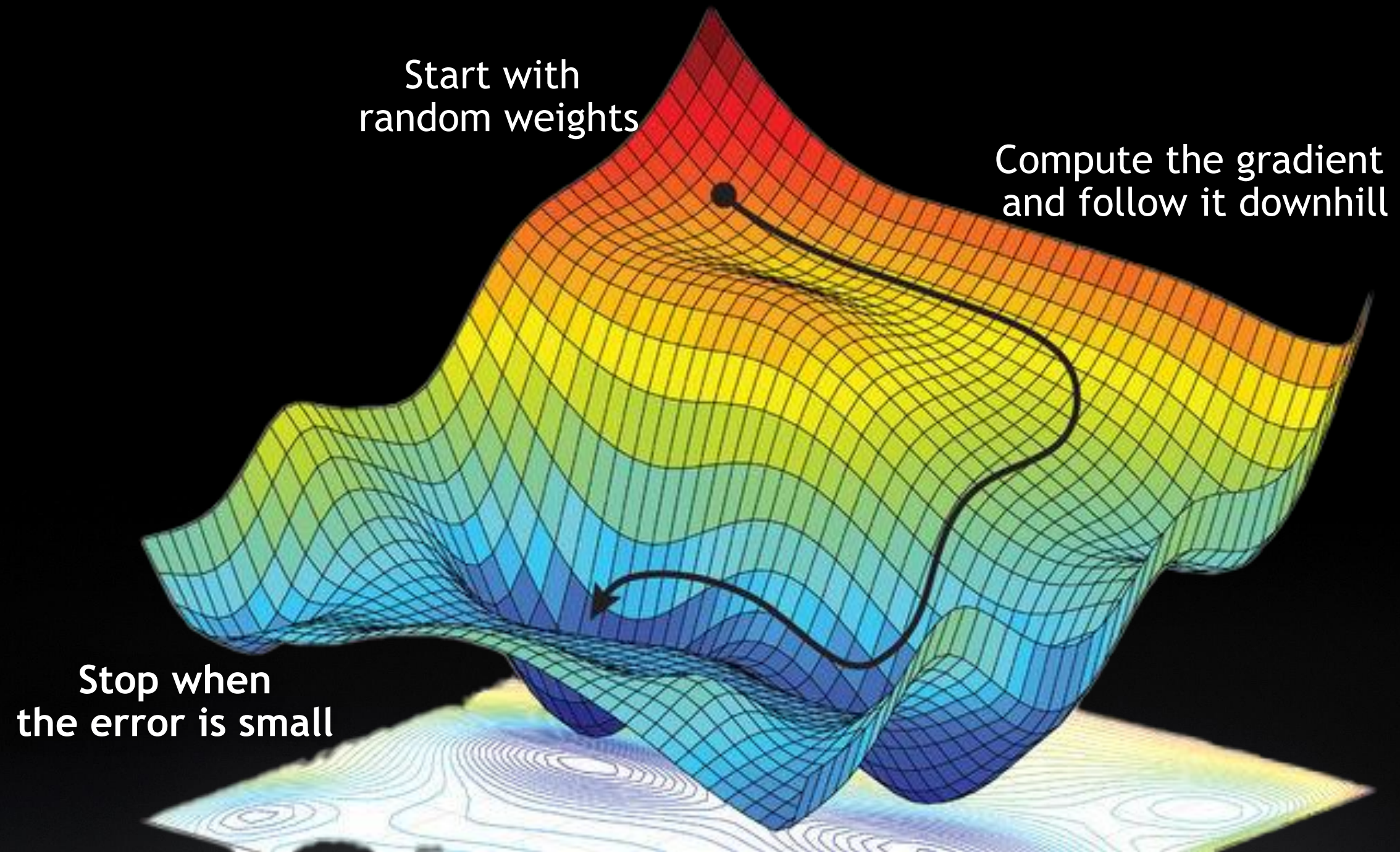
ACTUALLY USE:

- Gradient Descent



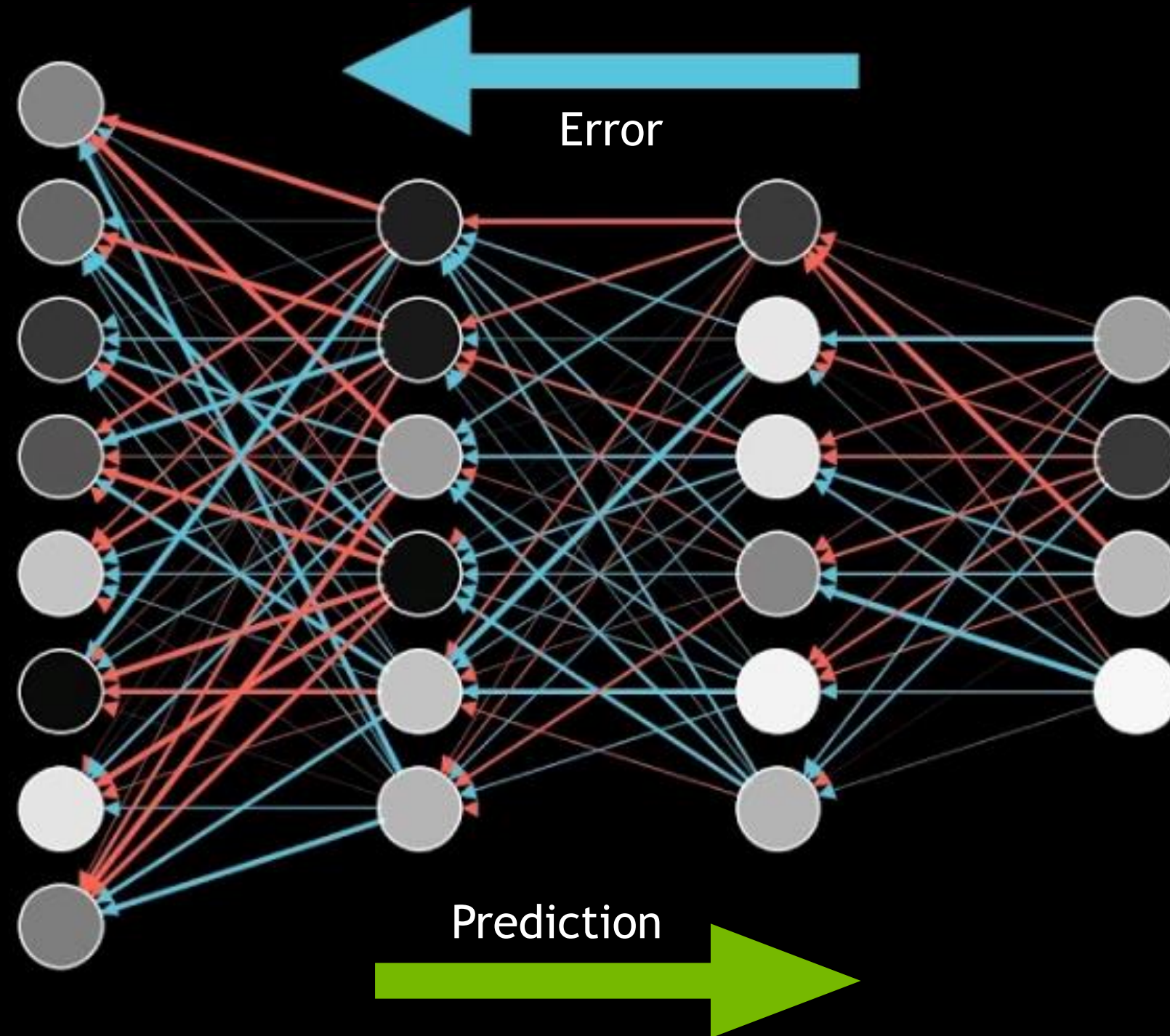
GRADIENT DESCENT

Finding as solution is as easy as falling down a hill



BACKPROPAGATION

Compute the gradient, by efficiently assigning blame



AUTOGRAD

Let a framework keep track of your gradient, so you don't have to

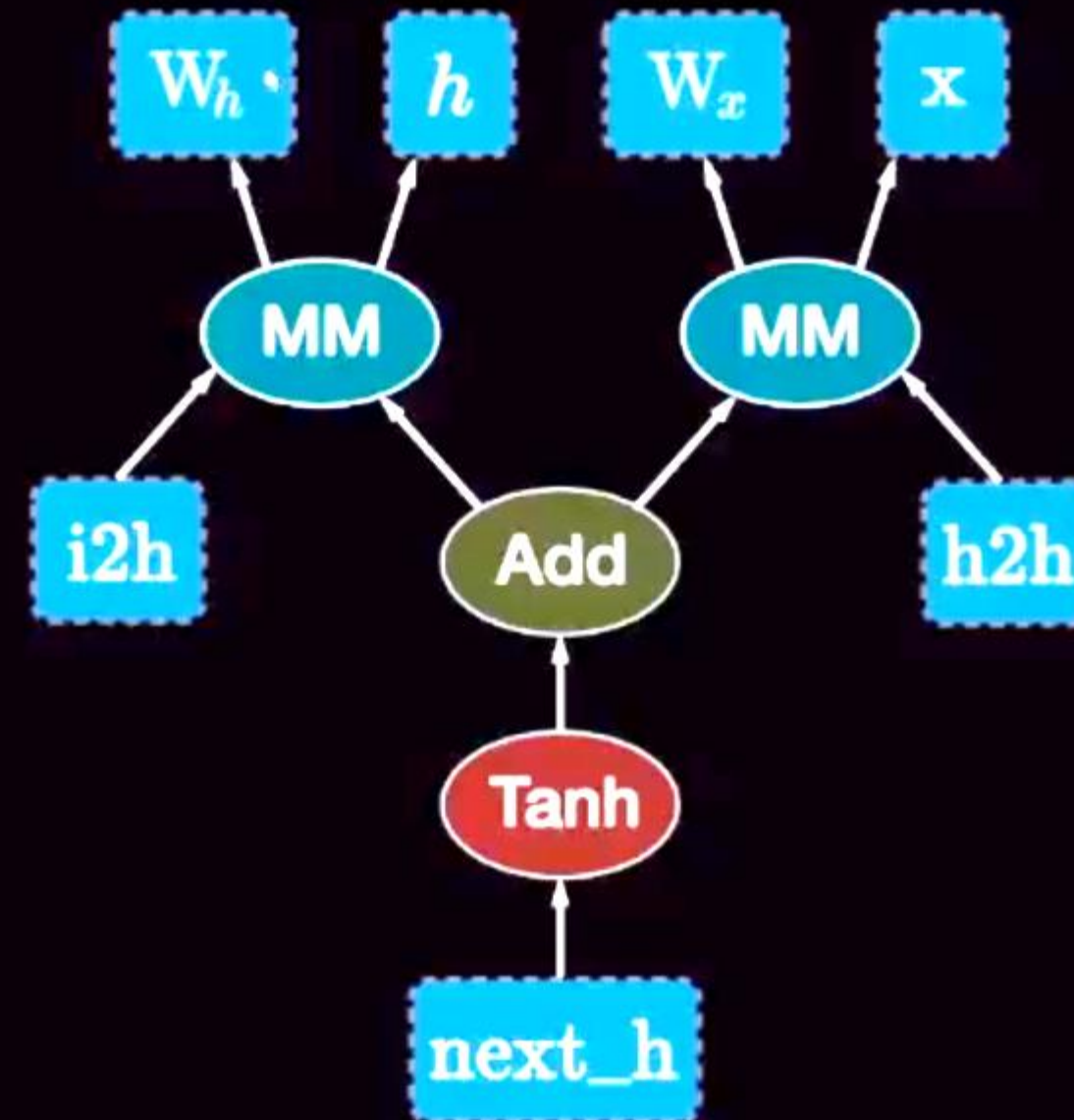
PyTorch Autograd

```
from torch.autograd import Variable

x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))

i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()

next_h.backward(torch.ones(1, 20))
```



WHAT YOU NEED TO MAKE IT WORK

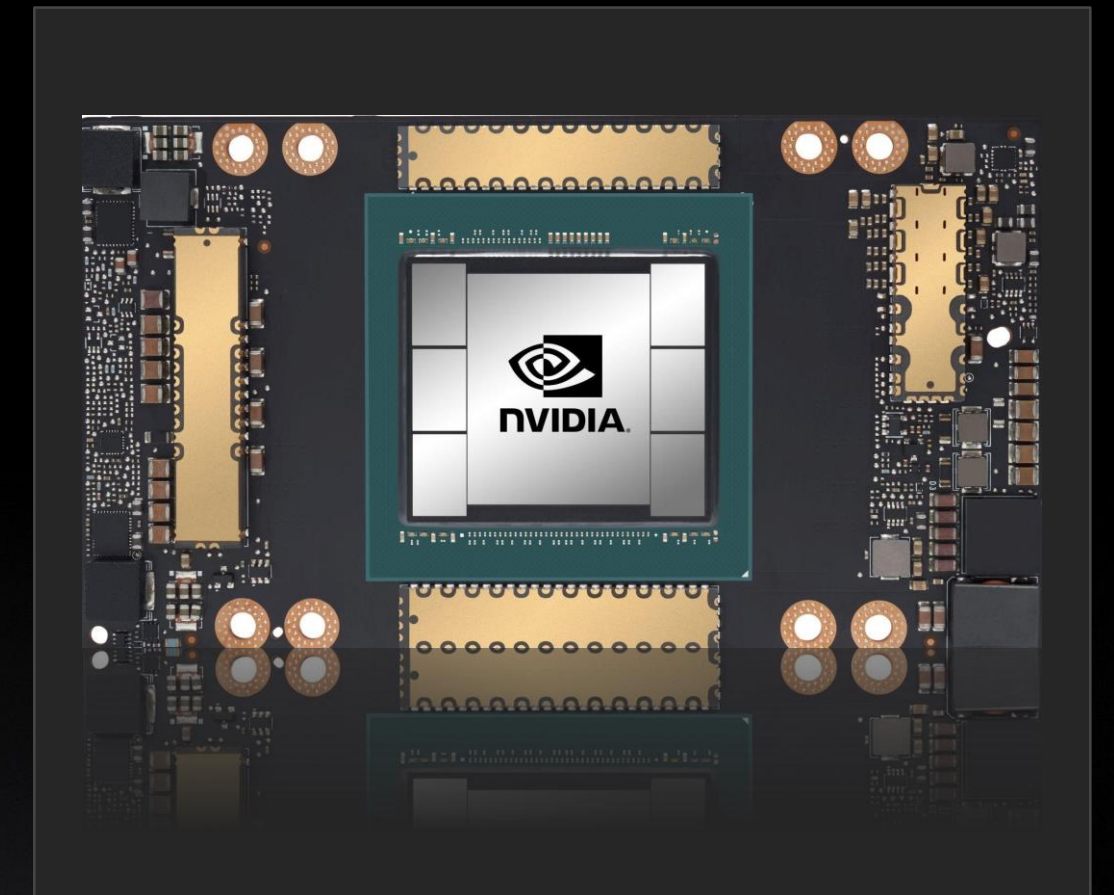
You need three main ingredients (and some skill)



LARGE QUANTITIES OF DATA



ML FRAMEWORK



GPU ACCELERATOR

DEEP LEARNING FRAMEWORK

Many frameworks to choose from (but not Fortran)

 PyTorch



Python

 mxnet



C++

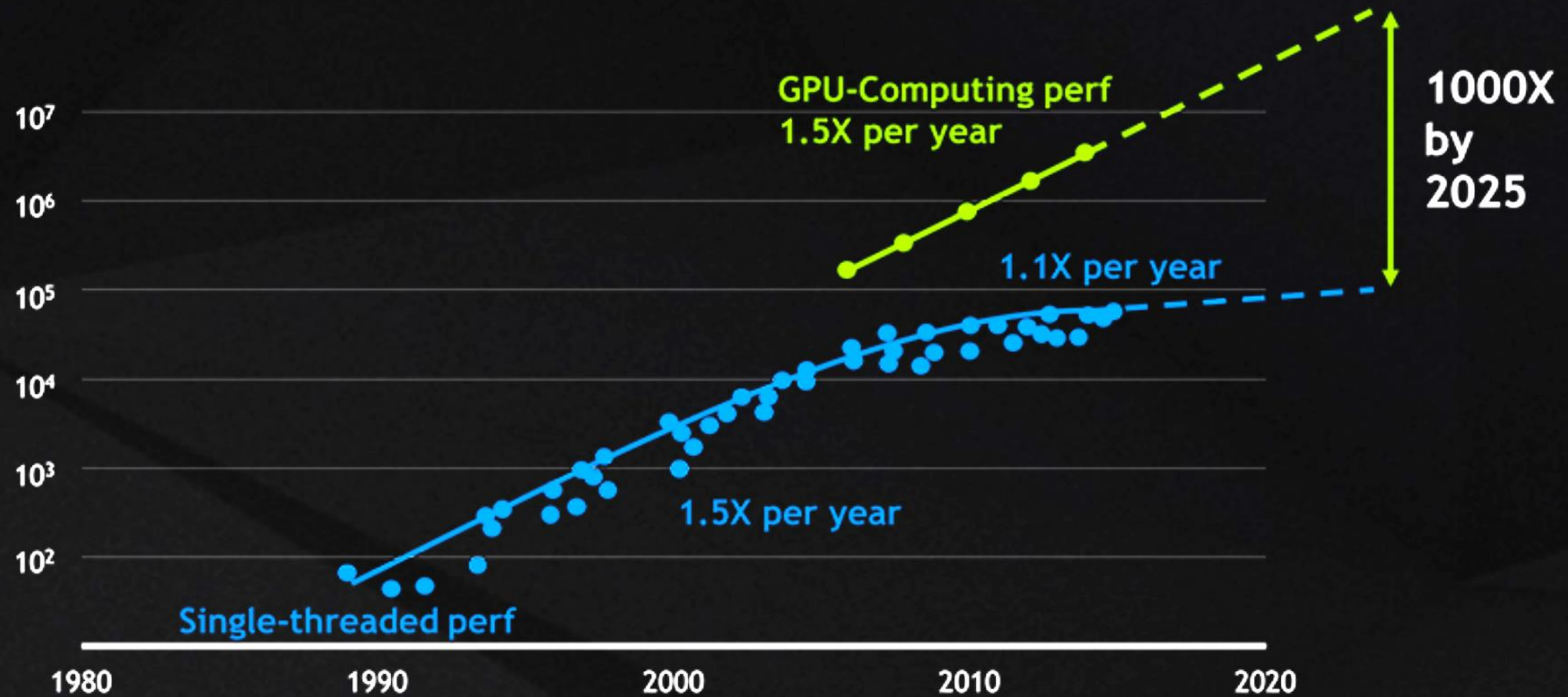


Julia



GPUS MAKE MACHINE LEARNING PRACTICAL

Train in a day? Or a month?



40 Years of Microprocessor Trend Data

LEARNED FUNCTIONS ARE GPU ACCELERATED

Next level software. No porting required.



HOW CAN I GET ACCESS TO A POWERFUL GPU?

Many way to take advantage of NVIDIA GPUs for Deep Learning



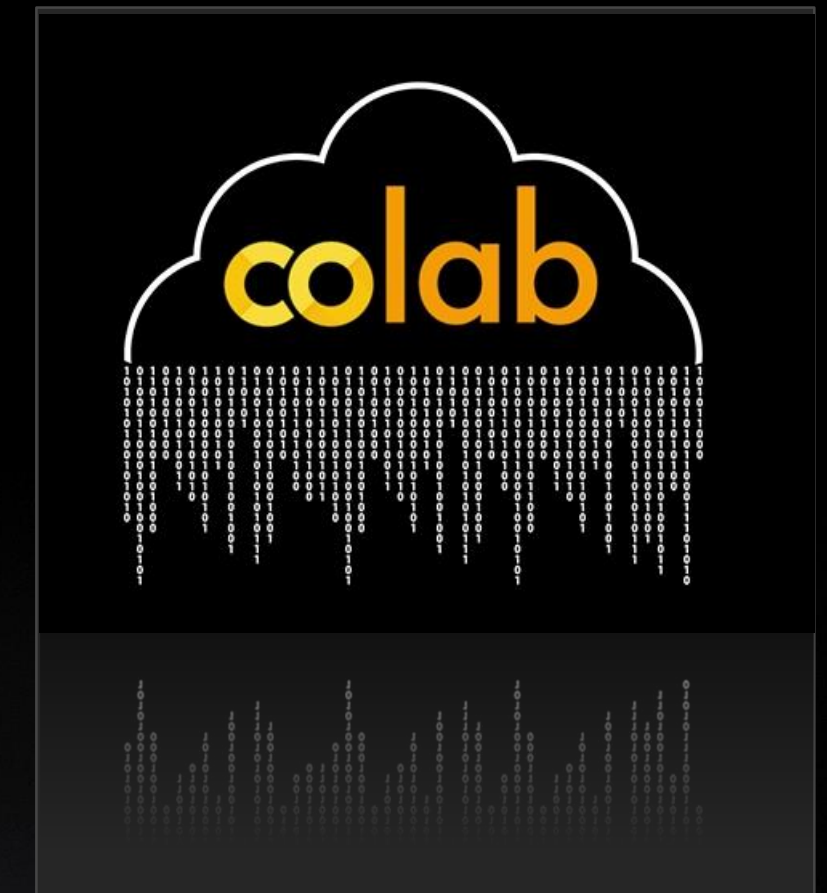
NVIDIA Quadro
Laptop or Workstation



Cloud Computing Services
(Free hours to start)



National Supercomputers
(Apply for compute)



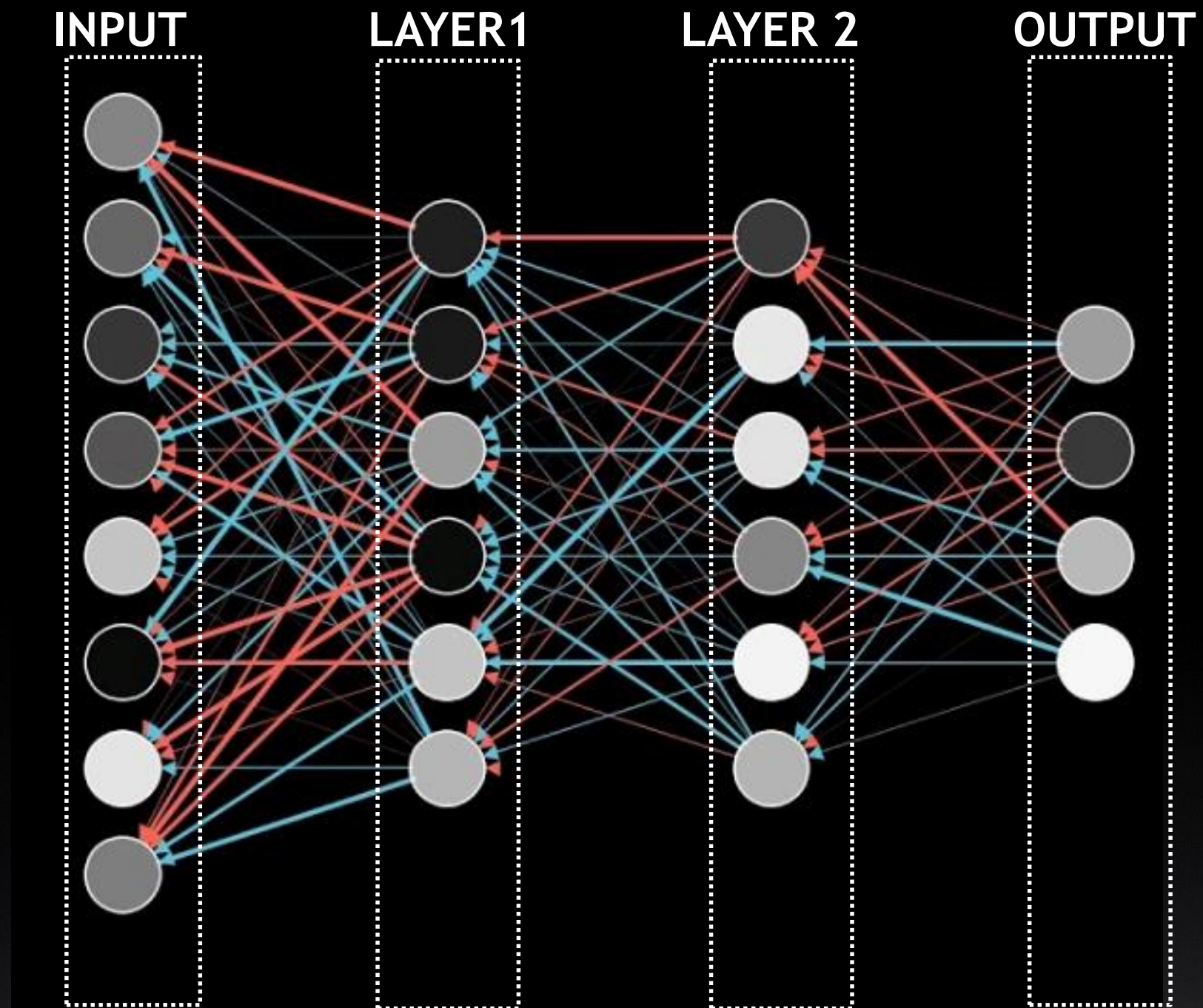
Google Colab
(1 Free NVIDIA GPU)



**FULLY CONNECTED
NETWORKS**
(MULTI-LAYER PERCPTRONS)

FULLY CONNECTED NETWORKS

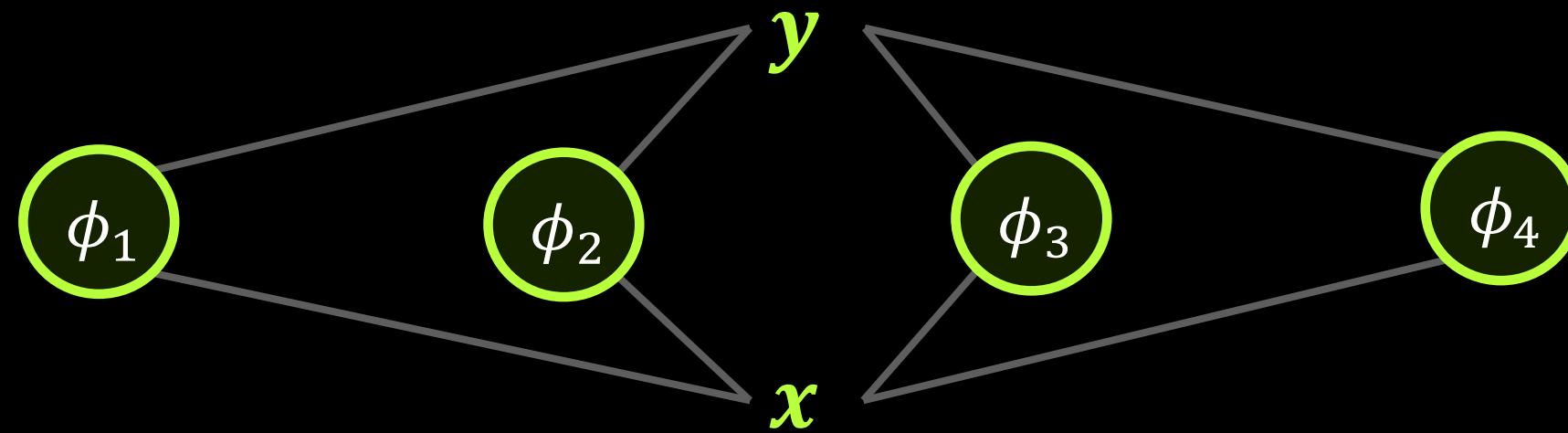
A given neuron is connected to every neuron in the previous layer



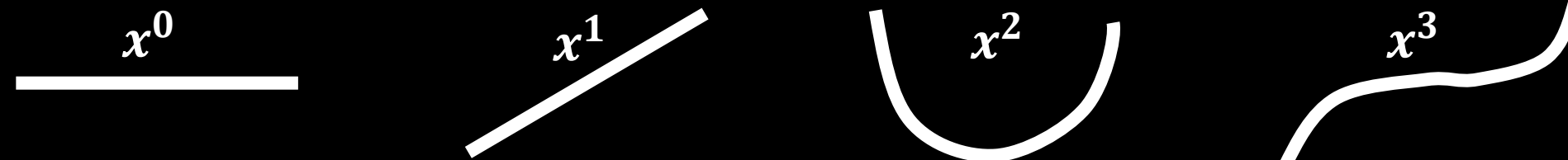
SINGLE LAYER NEURAL NETWORKS

A series expansion over basis functions ϕ .

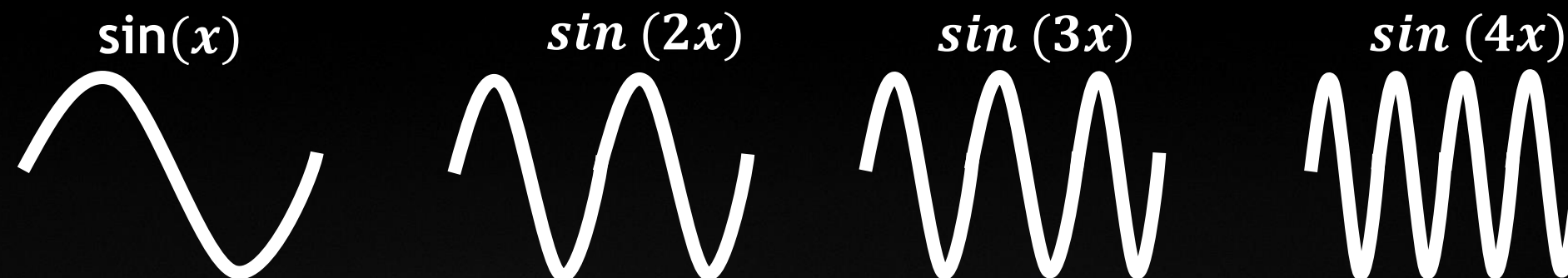
$$y = \sum_i w_i \phi_i(x + b_i)$$



TAYLOR SERIES



FOURIER SERIES

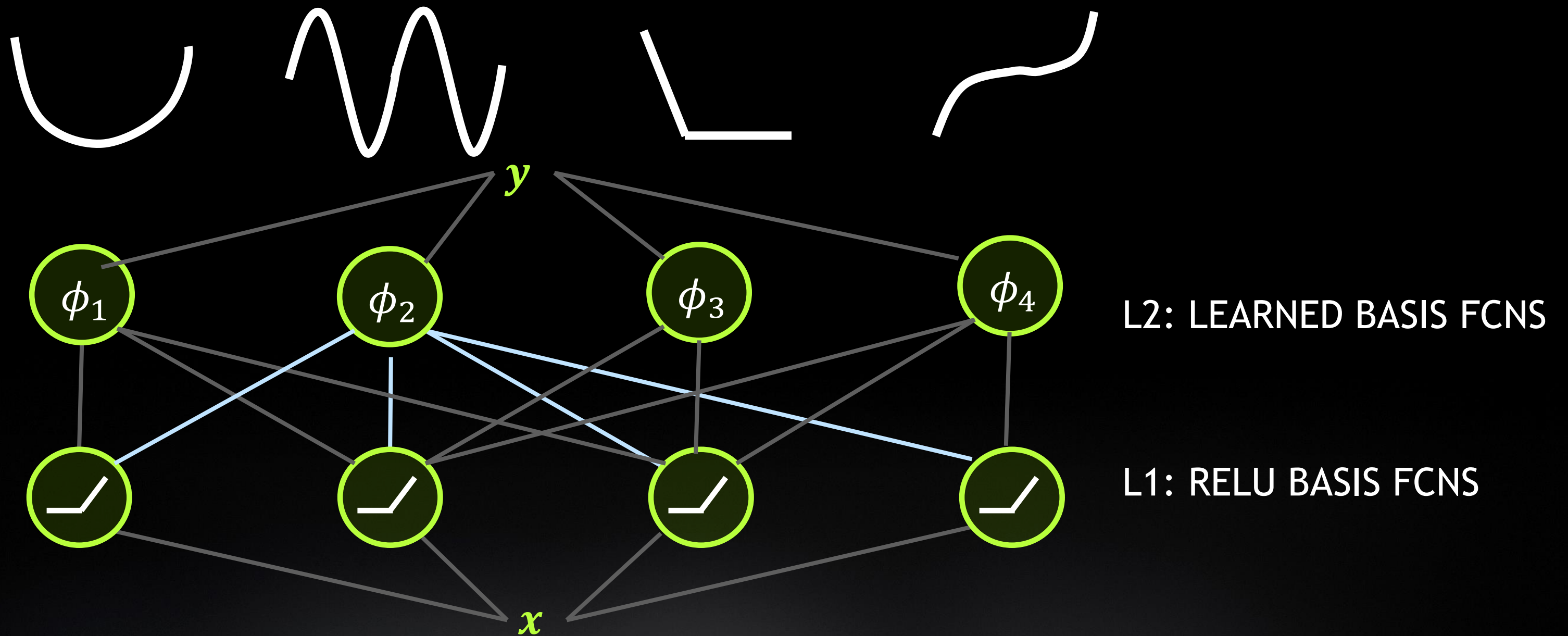


RELU



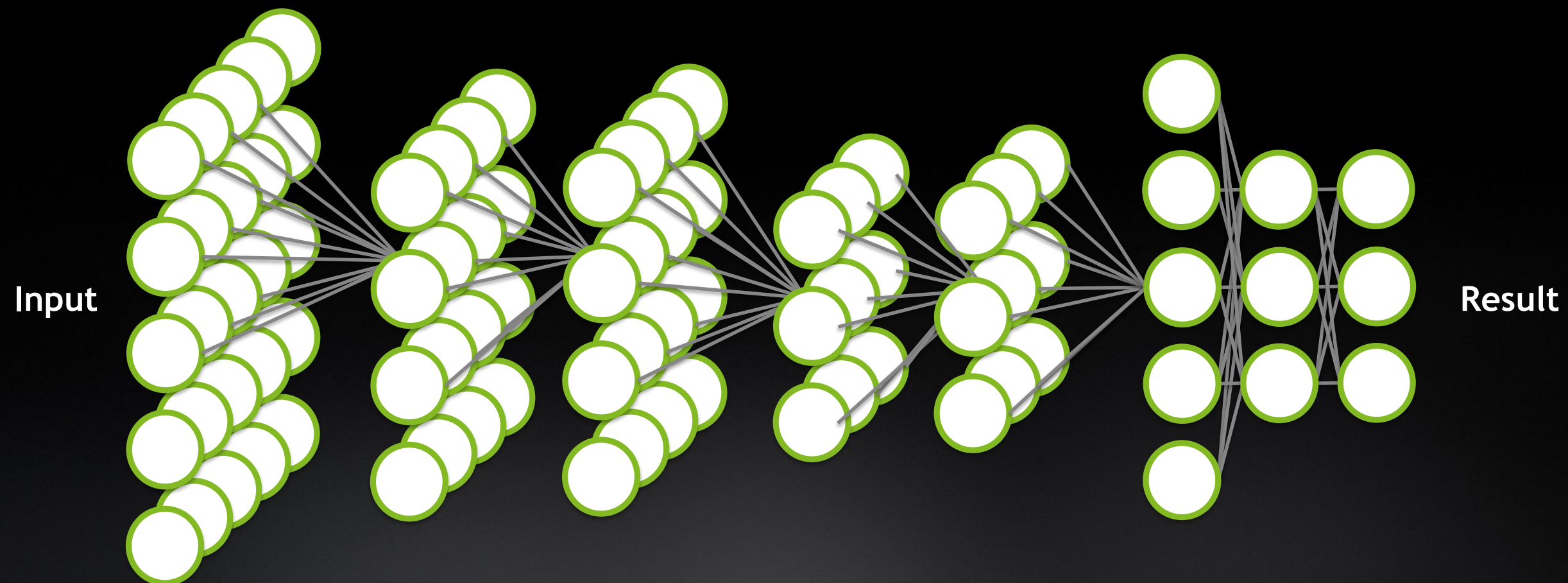
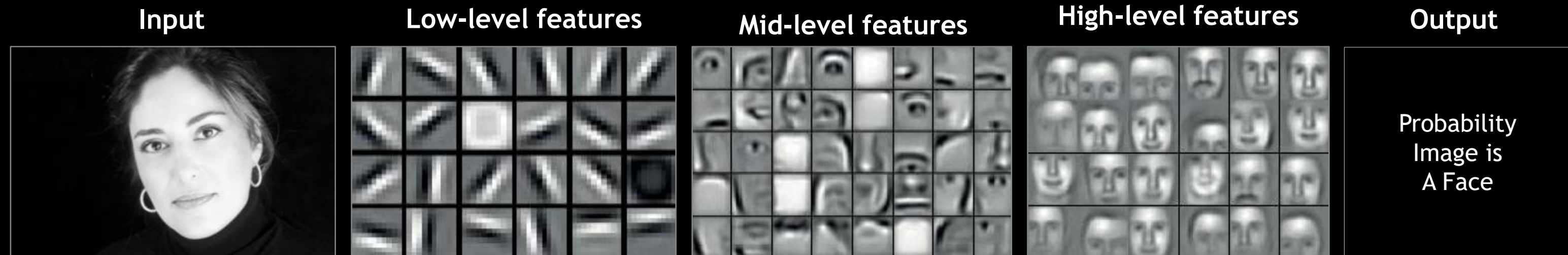
TWO LAYER NEURAL NETWORKS

Learn the function and the basis functions at the same time



DEEPER NEURAL NETWORKS

More layers allows for more levels of abstraction



Large Scale Visual Recognition Challenge 2012



The Imagenet competition: Automatically classify images from 1000 different categories



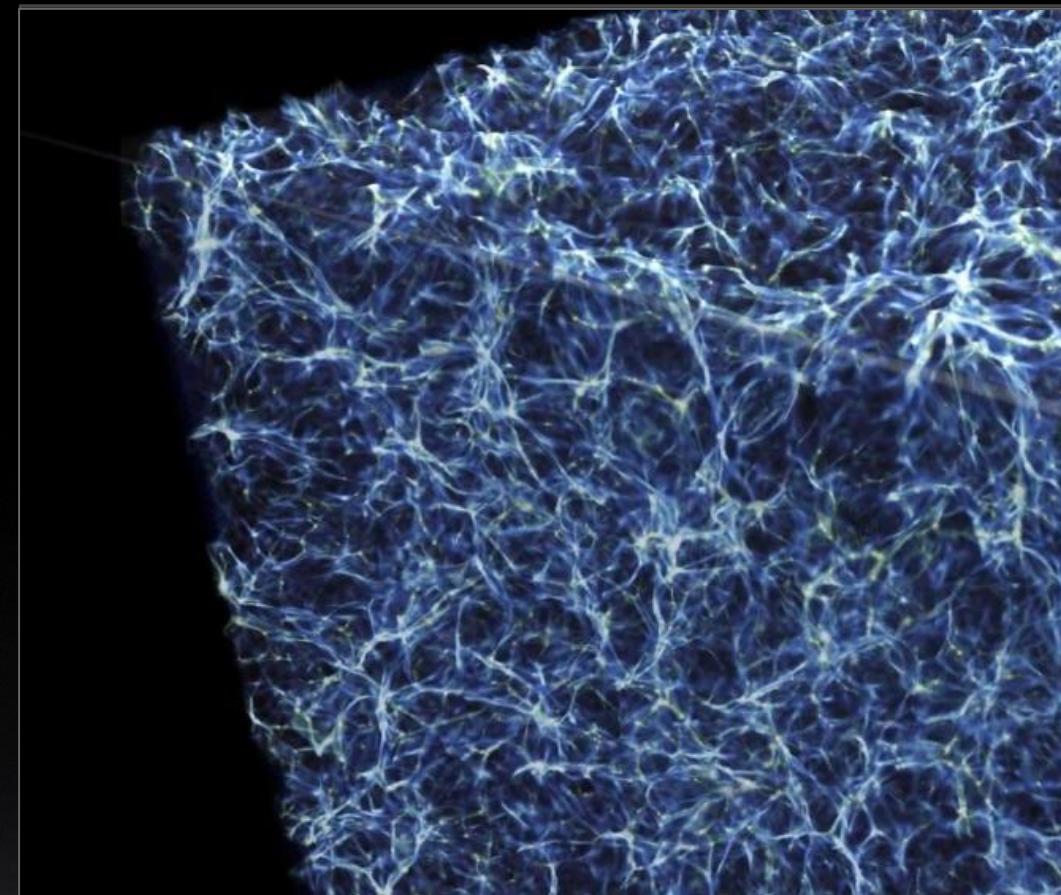
CONVOLUTIONAL NEURAL NETWORKS

WHAT ARE CNNs USED FOR?

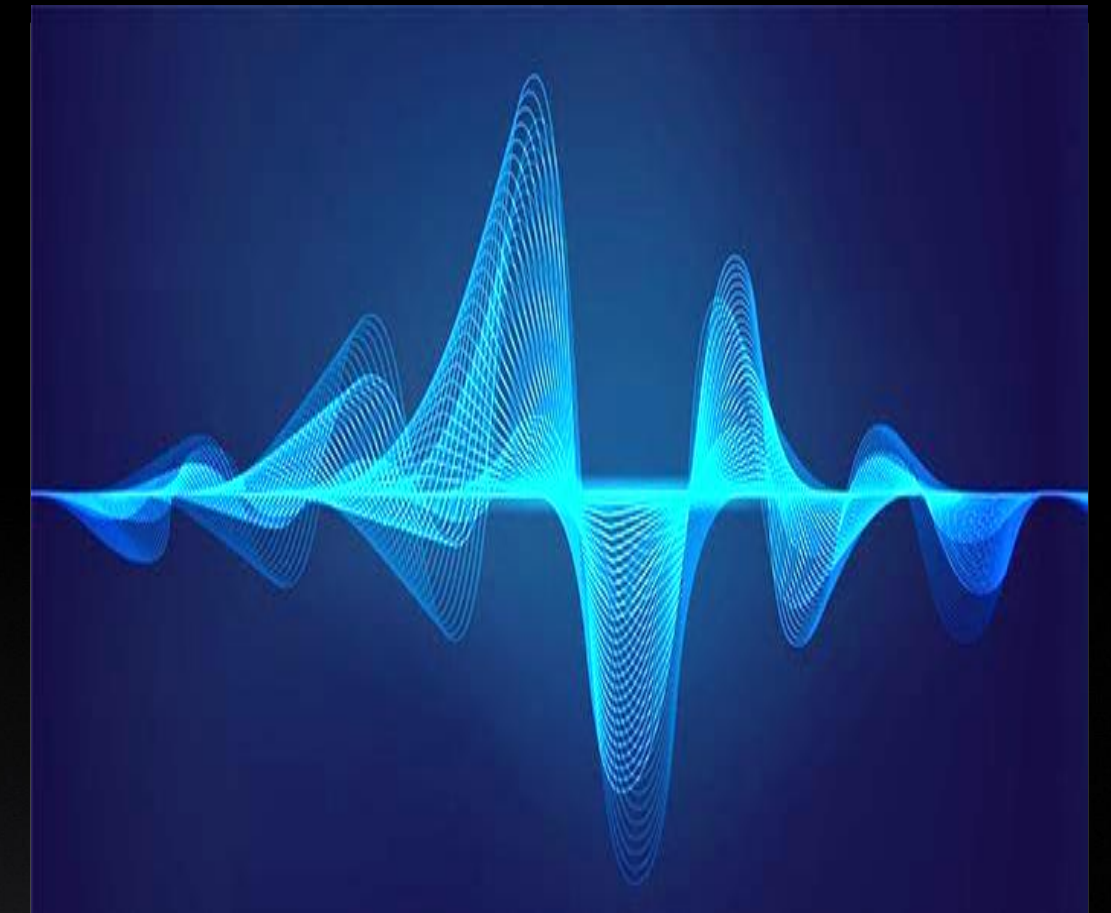
Problems with translational invariance



Computer Vision
Invariance in 2d space



Computational Physics
Invariance in 3d space



Audio and Time Series
Invariance in time

COMPUTER VISION TASKS

Each task requires a different model and data setup



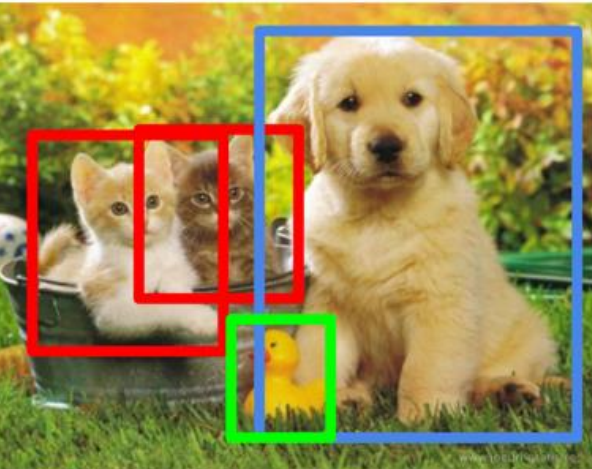

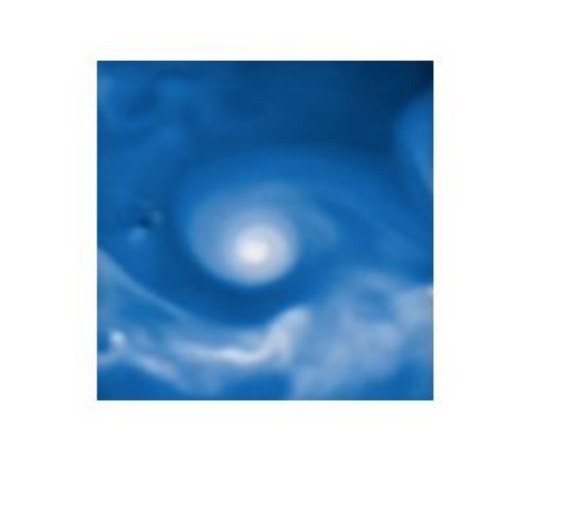
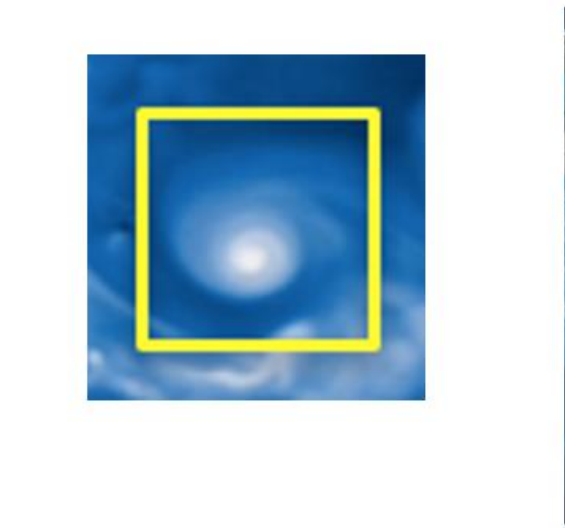
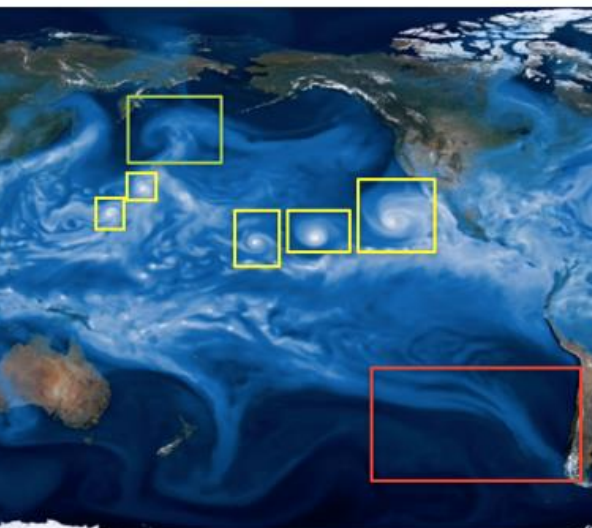
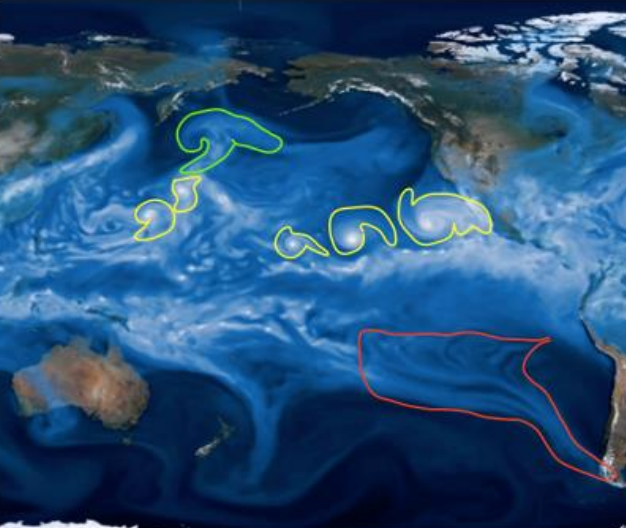
Classification	Classification + Localization	Object Detection	Instance Segmentation
			
			

Image Credit: NERSC

CLASSIFICATION

Example: Classifying Land Use



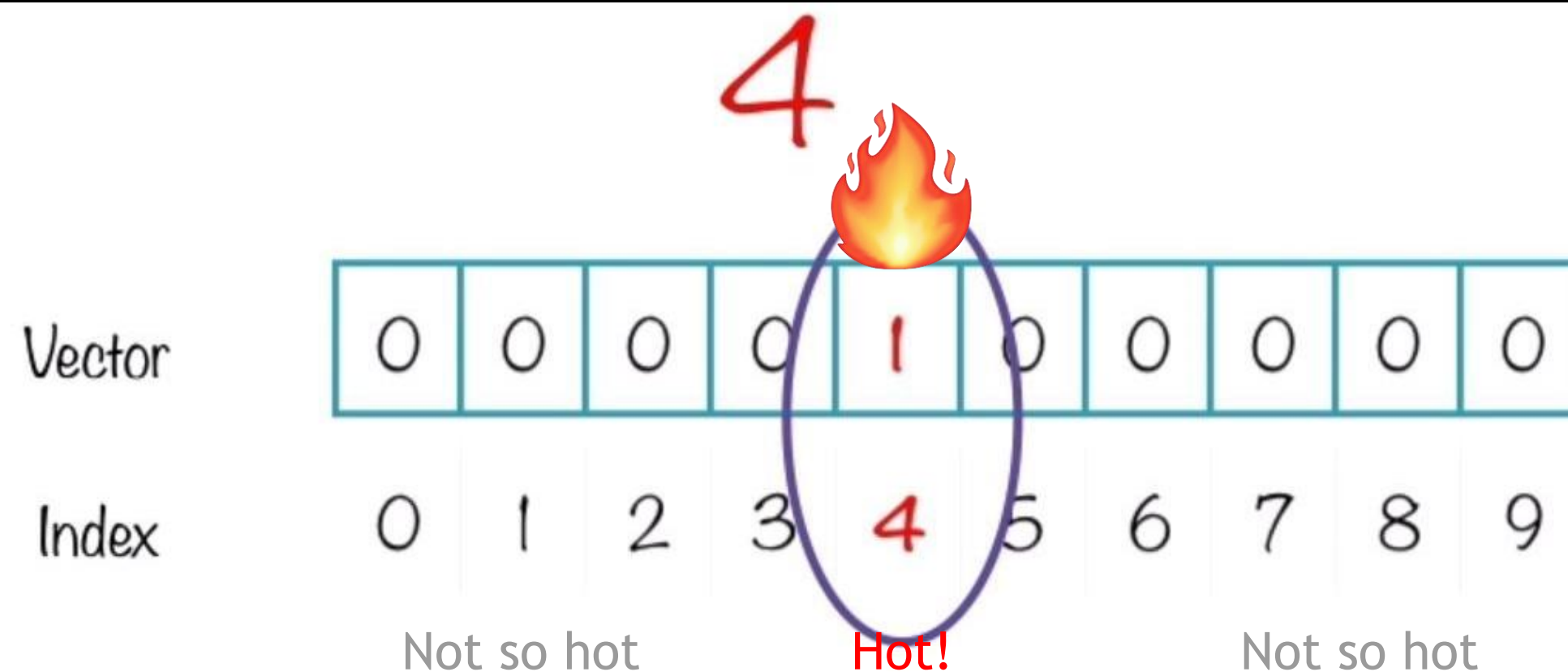
ONE-HOT ENCODING

Input: Pixels, Output: One-hot encoding

INPUT: PIXEL VALUES



OUTPUT: ONE-HOT VECTOR

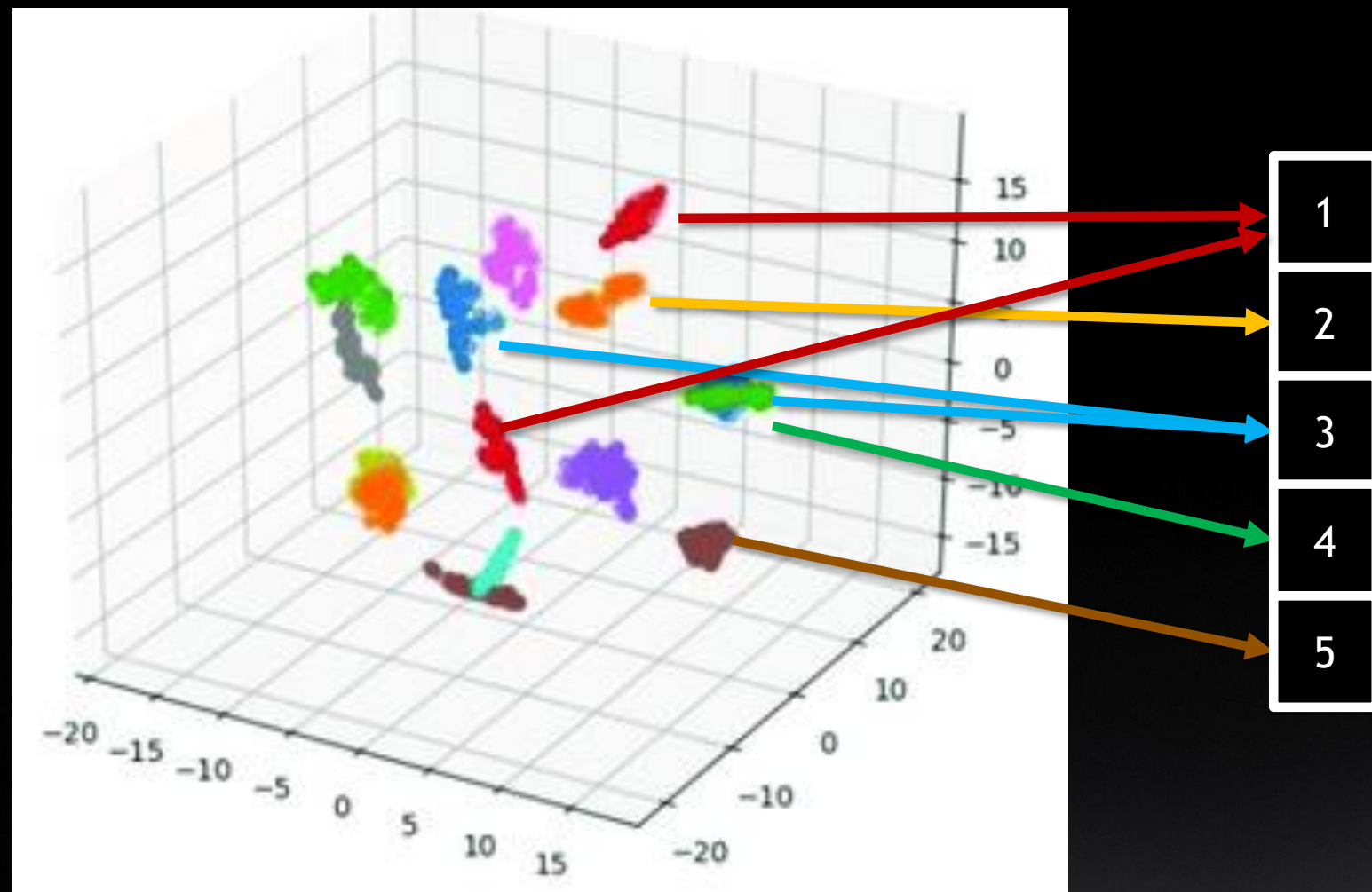


<https://blog.carbonteq.com/practical-image-recognition-with-tensorflow/>

IMAGES ARE POINTS, WITH MANY DIMENSIONS

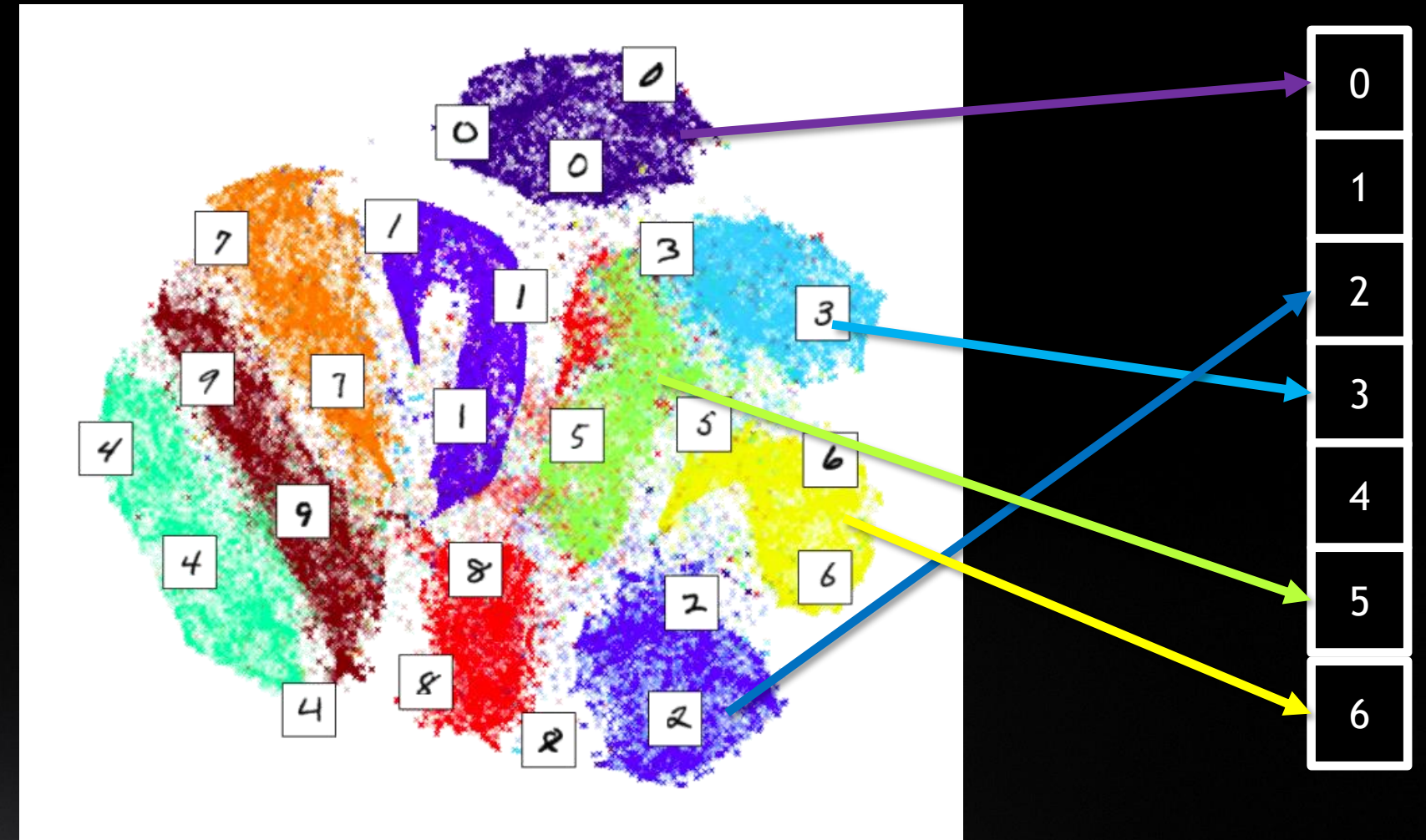
IN: 3-D Vector

OUT: 1 hot vector

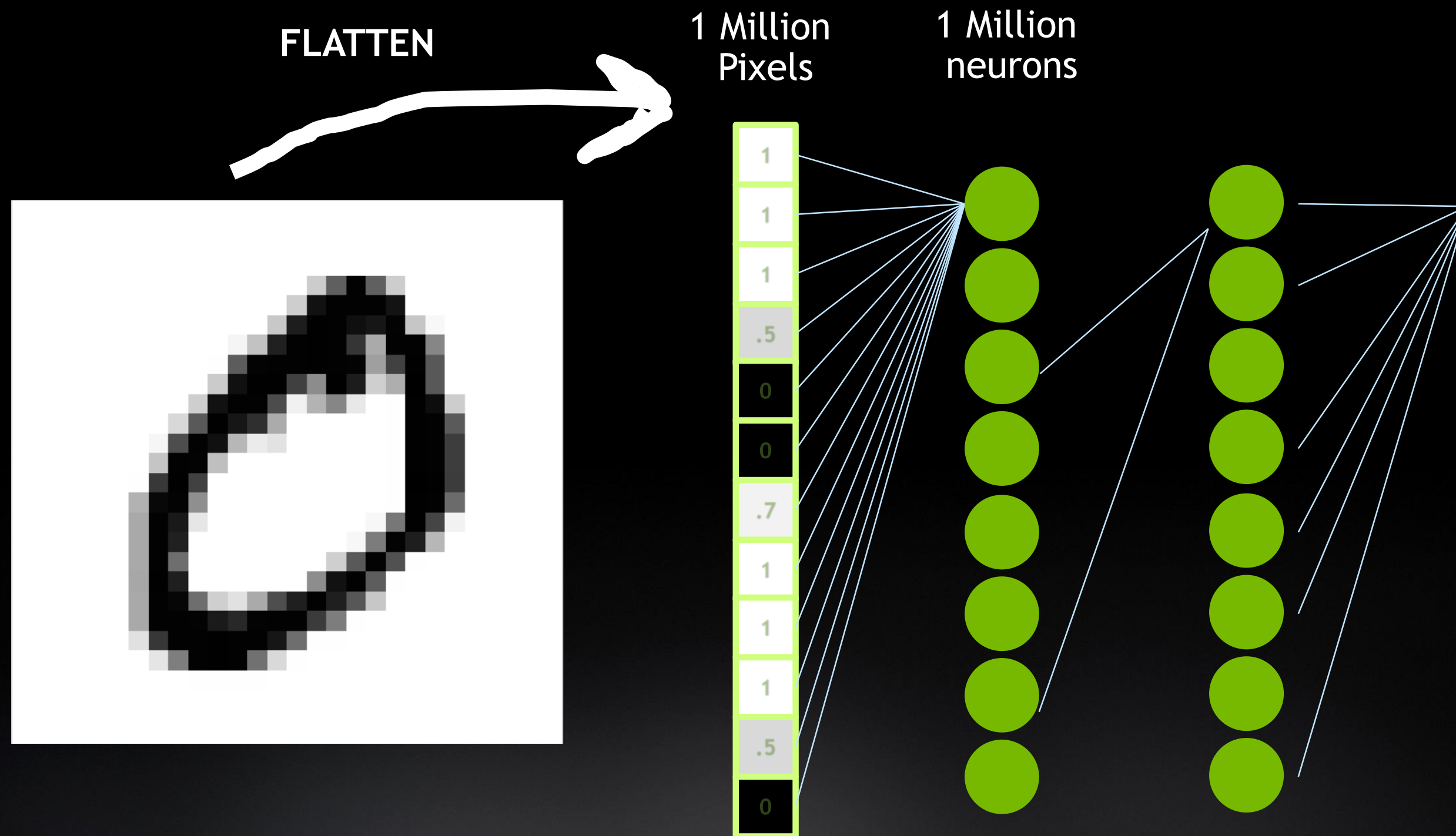


IN: 784-D Vector

OUT: 1-hot vector



FULLY CONNECTED NETWORKS AND IMAGES DON'T MIX



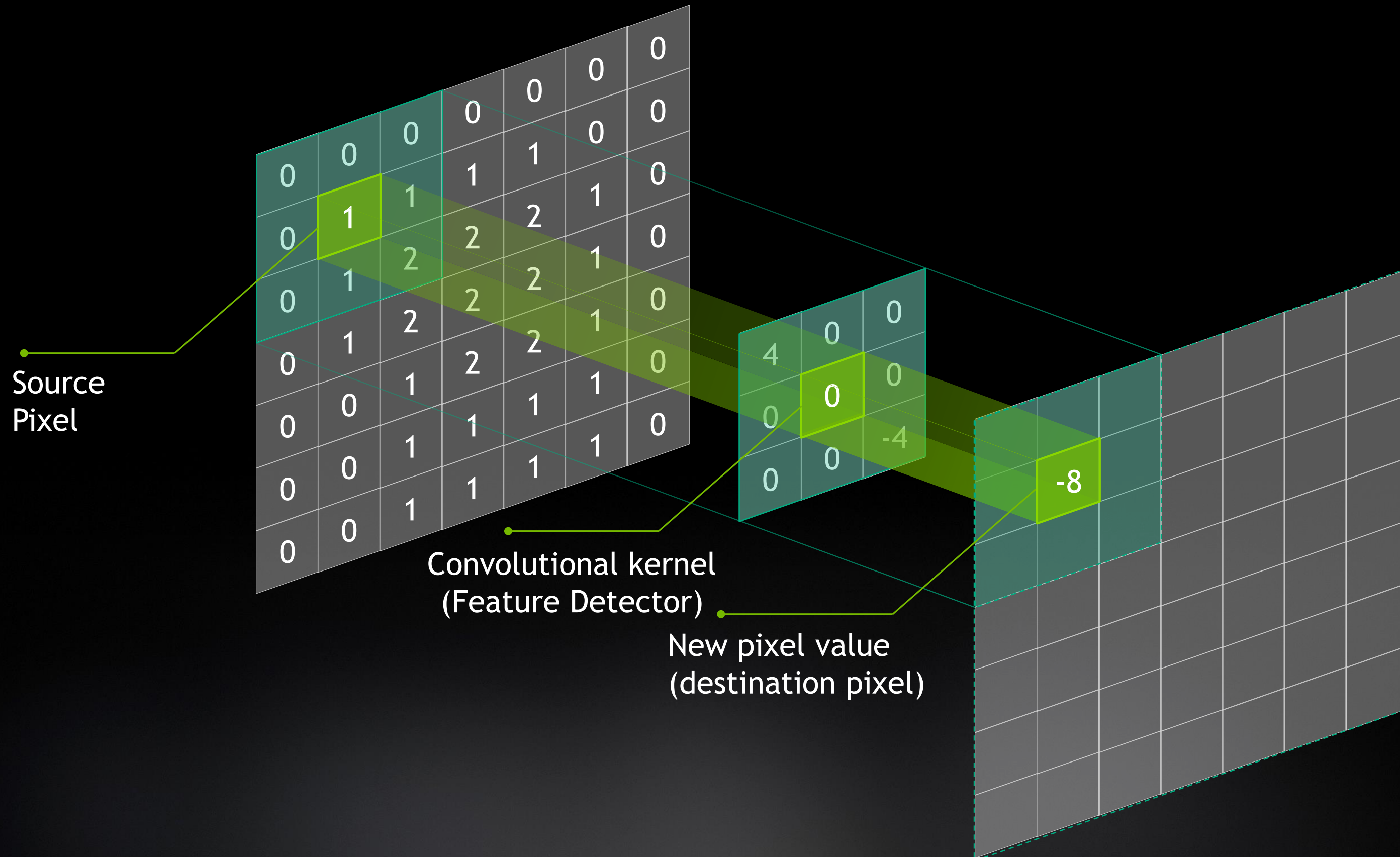
TRANSLATIONAL EQUIVARIANCE

Objects in nature look the same from place to place

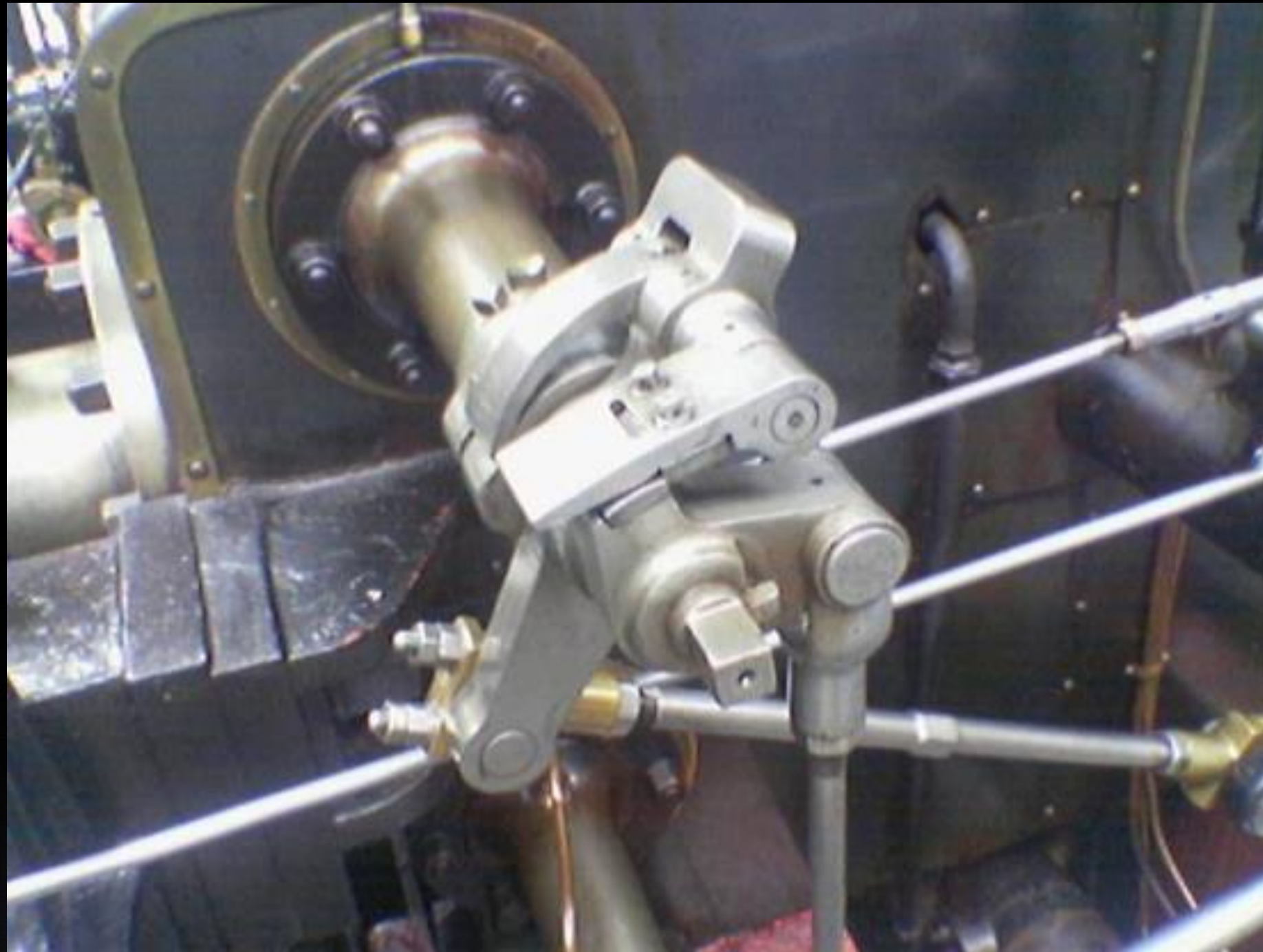


WHAT IS A CONVOLUTION?

A small matrix transformation, applied at each point of the image



CONVOLUTION EXAMPLE: SOBEL FILTER

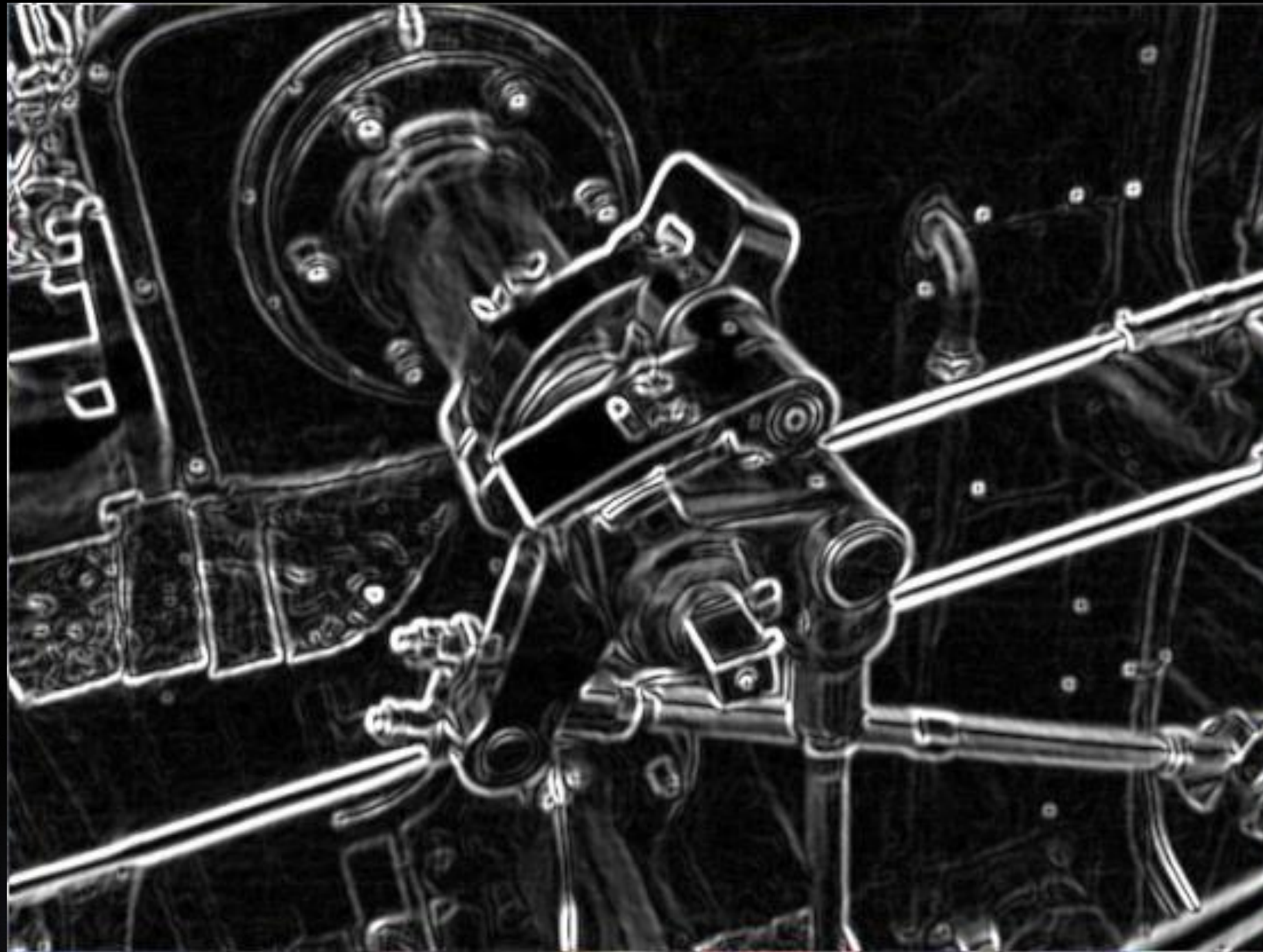


$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

$$G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

$$G = \sqrt{G_x^2 + G_y^2}$$

CONVOLUTION EXAMPLE: SOBEL FILTER



$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

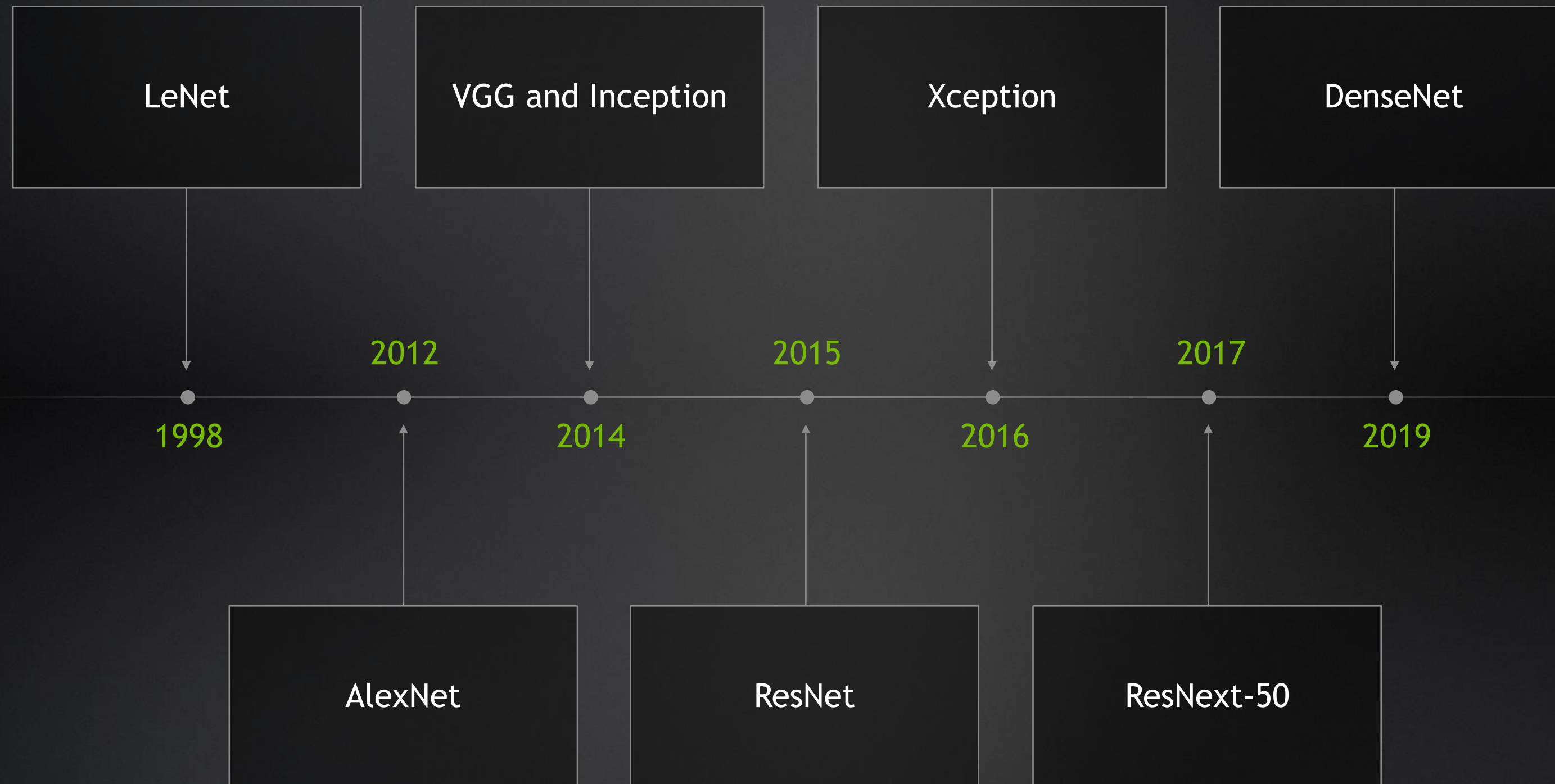
$$G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

$$G = \sqrt{G_x^2 + G_y^2}$$



CLASSIFICATION

CLASSIFIER EVOLUTION OVER TIME



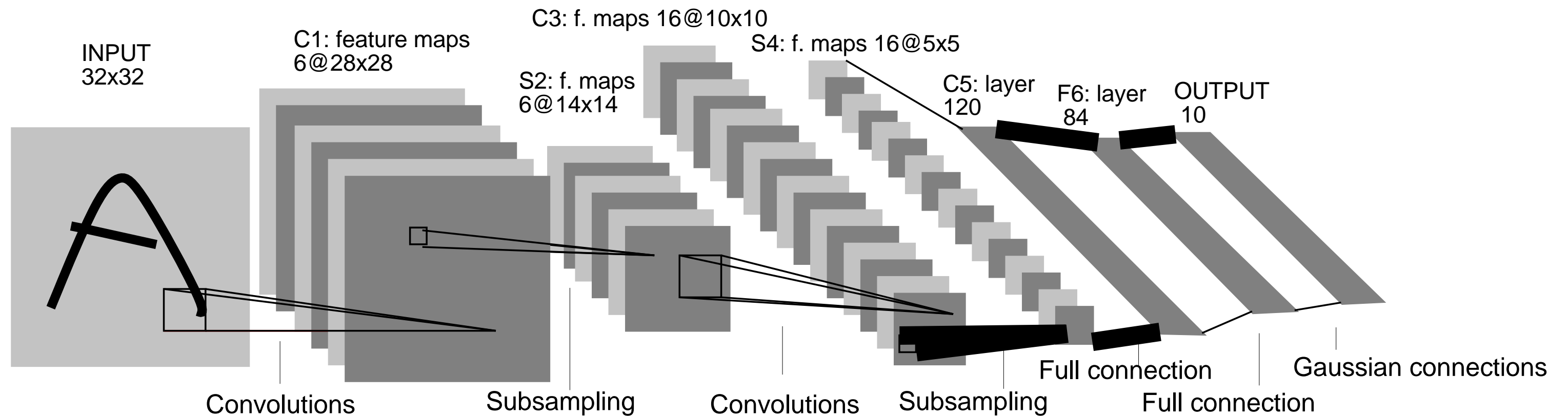
LENET-5

(1988) Yann LeCun. Hand written recognition. 60k parameters.

PROC. OF THE IEEE, NOVEMBER 1998

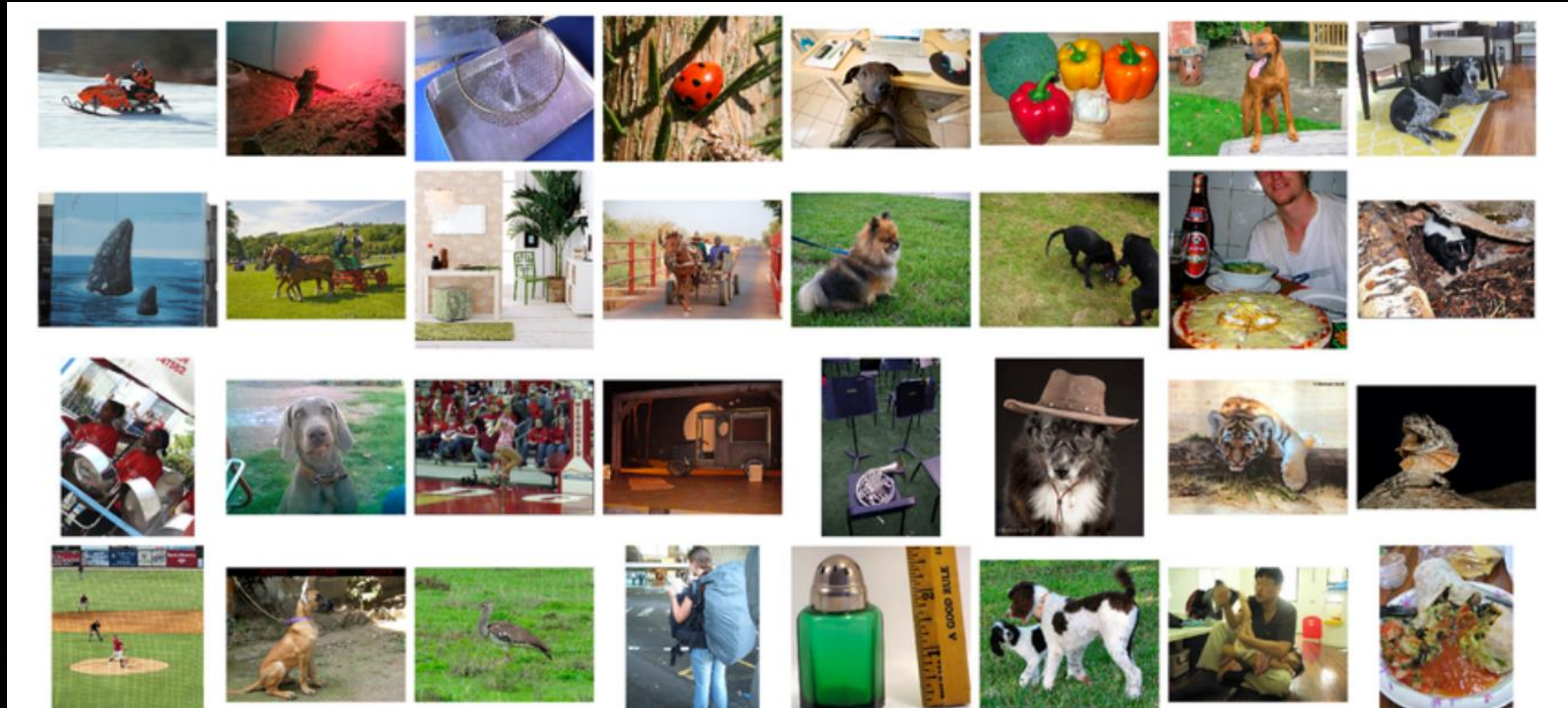
Gradient-Based Learning Applied to Document Recognition

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner



IMAGENET ILSVR COMPETITION

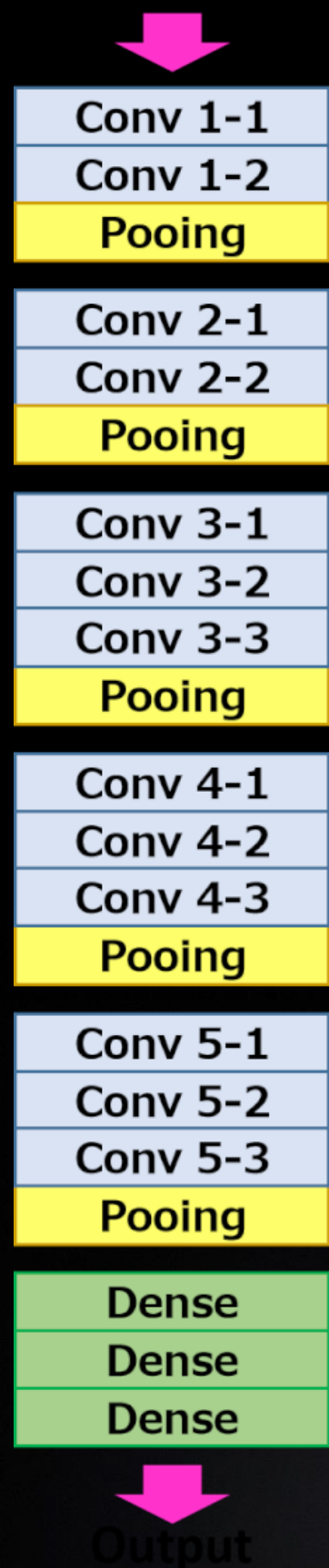
Large Scale Visual Recognition Competition (2010-2017)



<https://en.wikipedia.org/wiki/ImageNet>

VGG-16

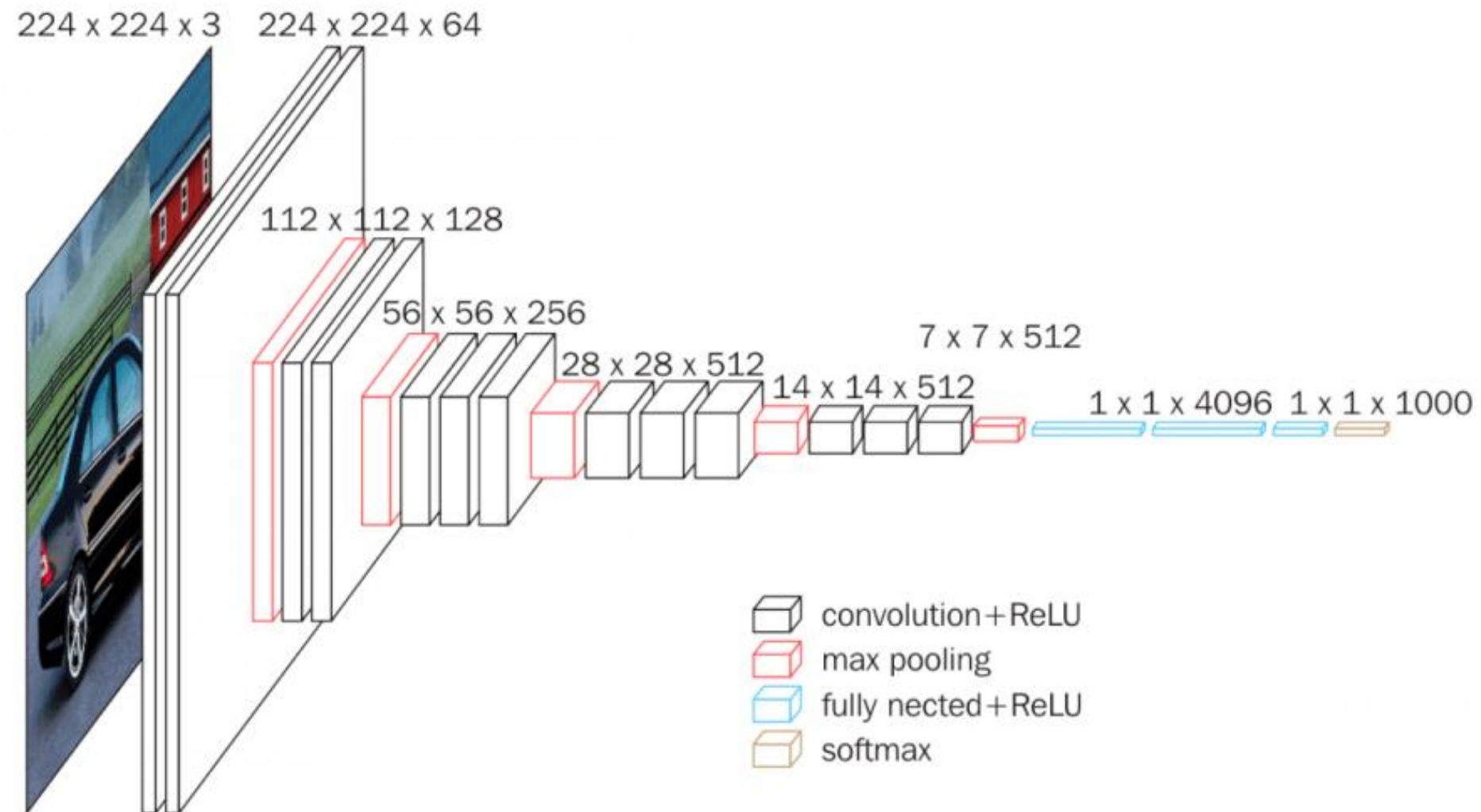
2014. ImageNet runner up. Simple, clean architecture.



VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

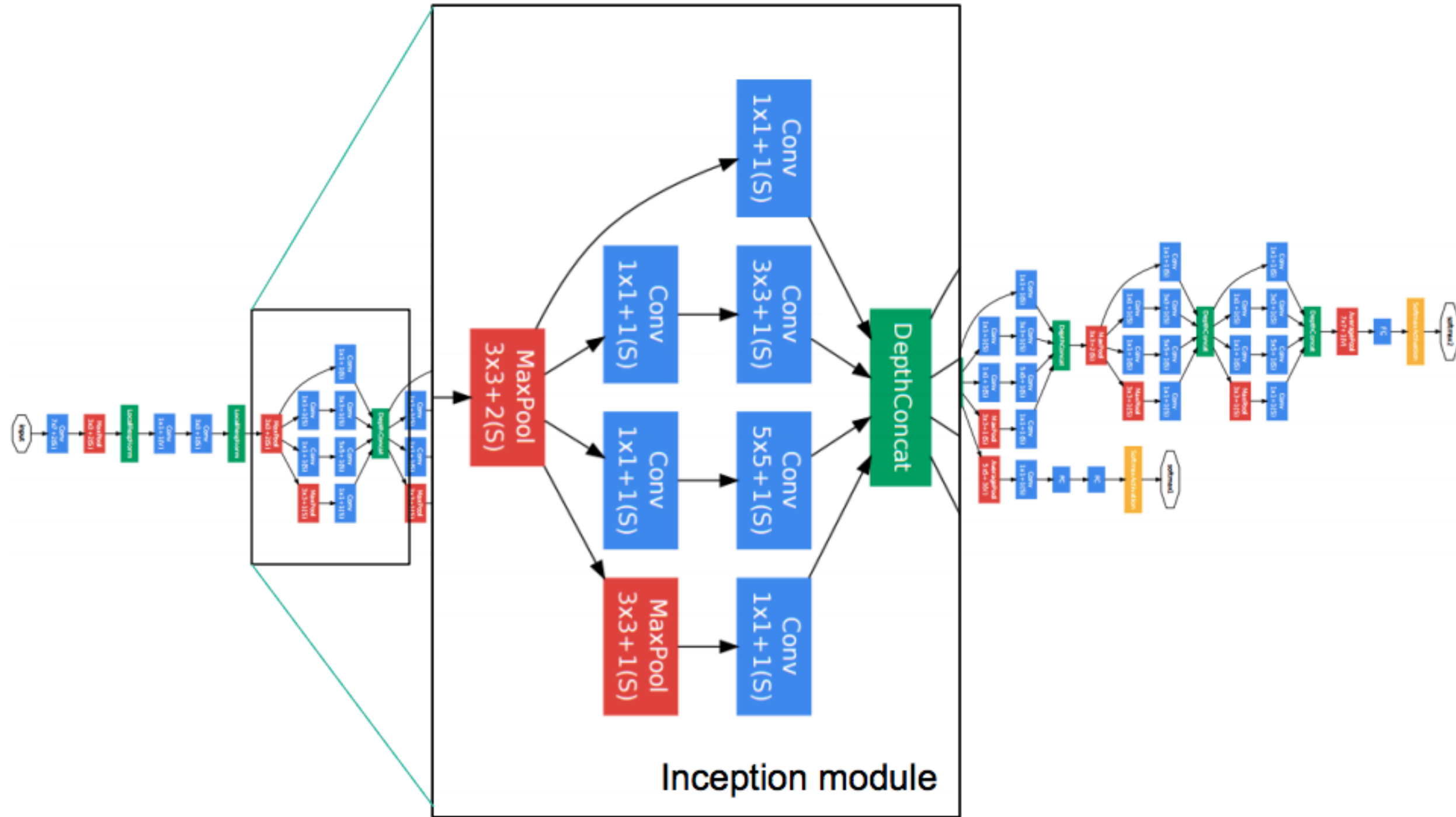
Karen Simonyan* & Andrew Zisserman⁺

Visual Geometry Group, Department of Engineering Science, University of Oxford



INCEPTION-V1 (GOOGLENET)

2014. Train different size convolutions in parallel



FULLY CONVOLUTIONAL NETWORKS (FCN)

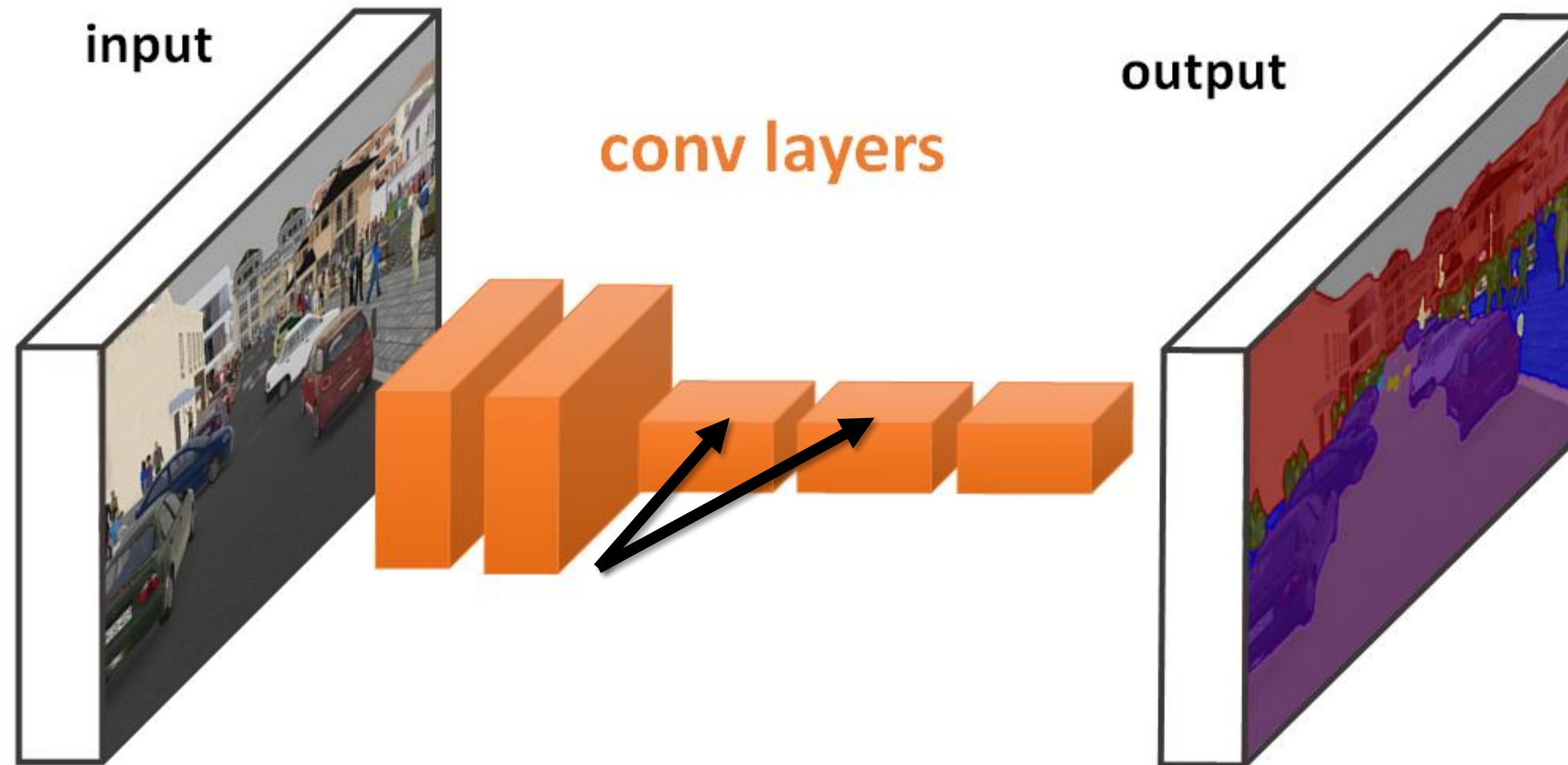
2015: Convert fully connect layers into convolutions of the same size

Fully Convolutional Networks for Semantic Segmentation

Jonathan Long*

Evan Shelhamer*
UC Berkeley

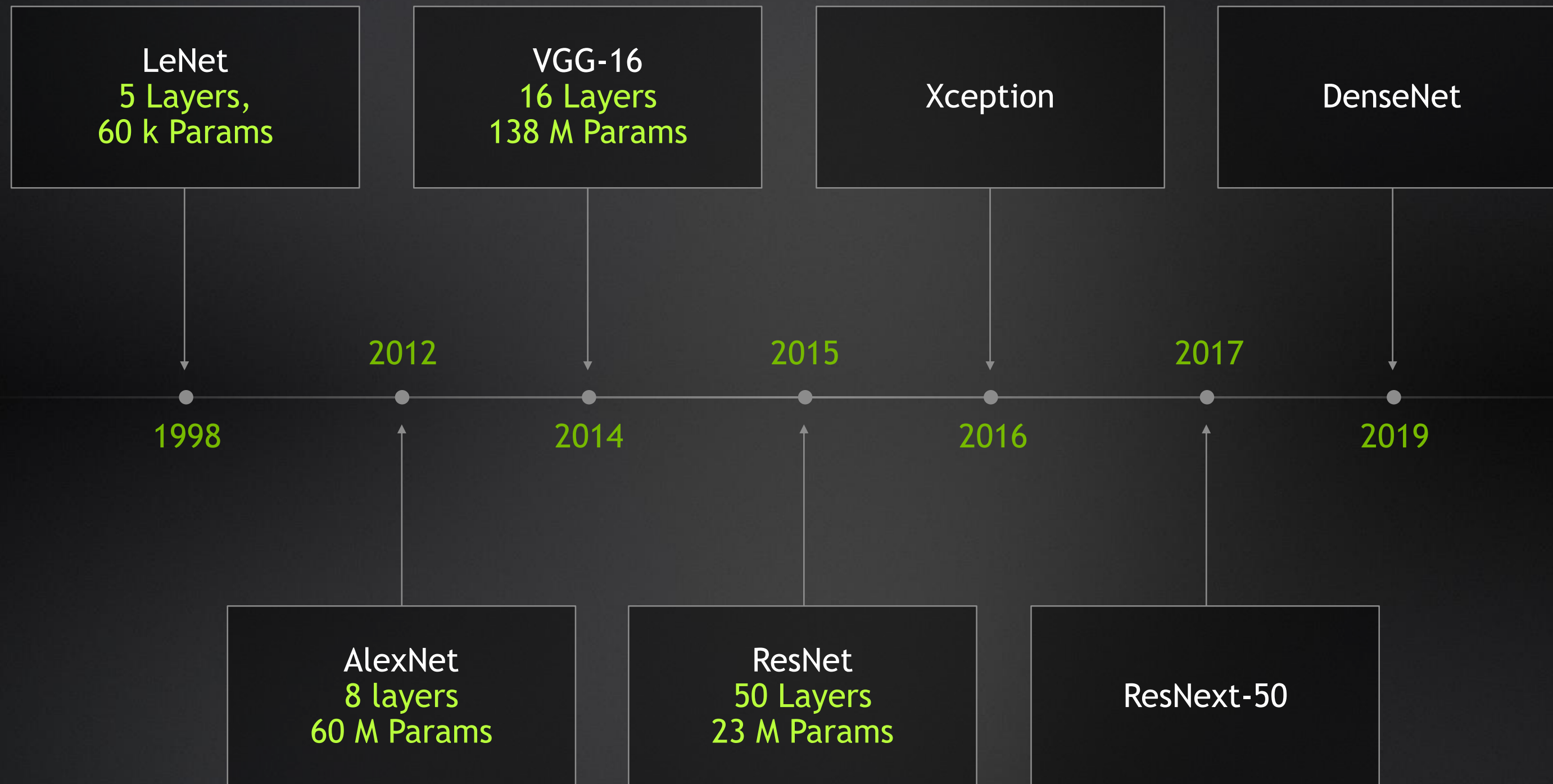
Trevor Darrell





RESNETS

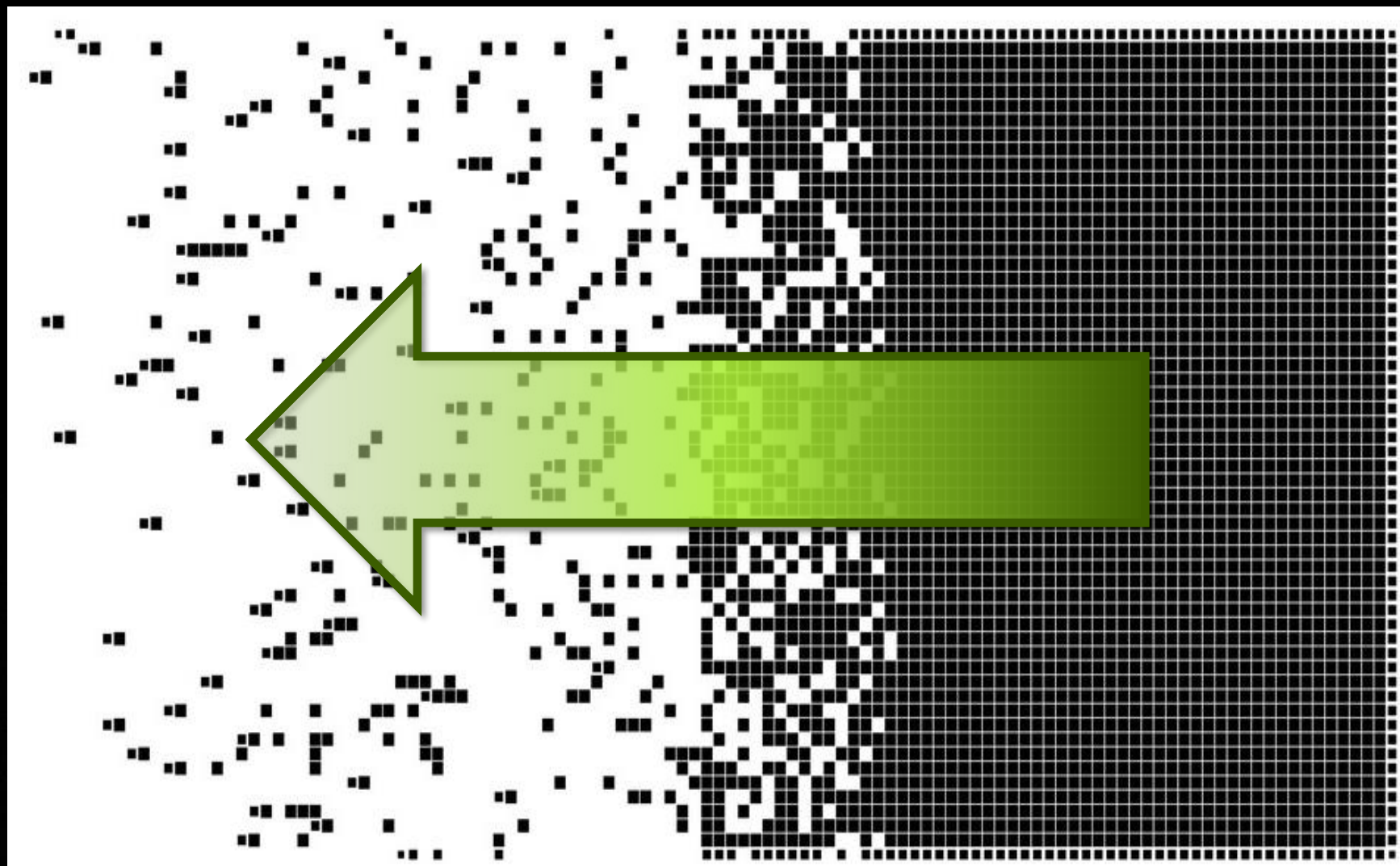
MODELING TRENDS: DEEPER AND LARGER



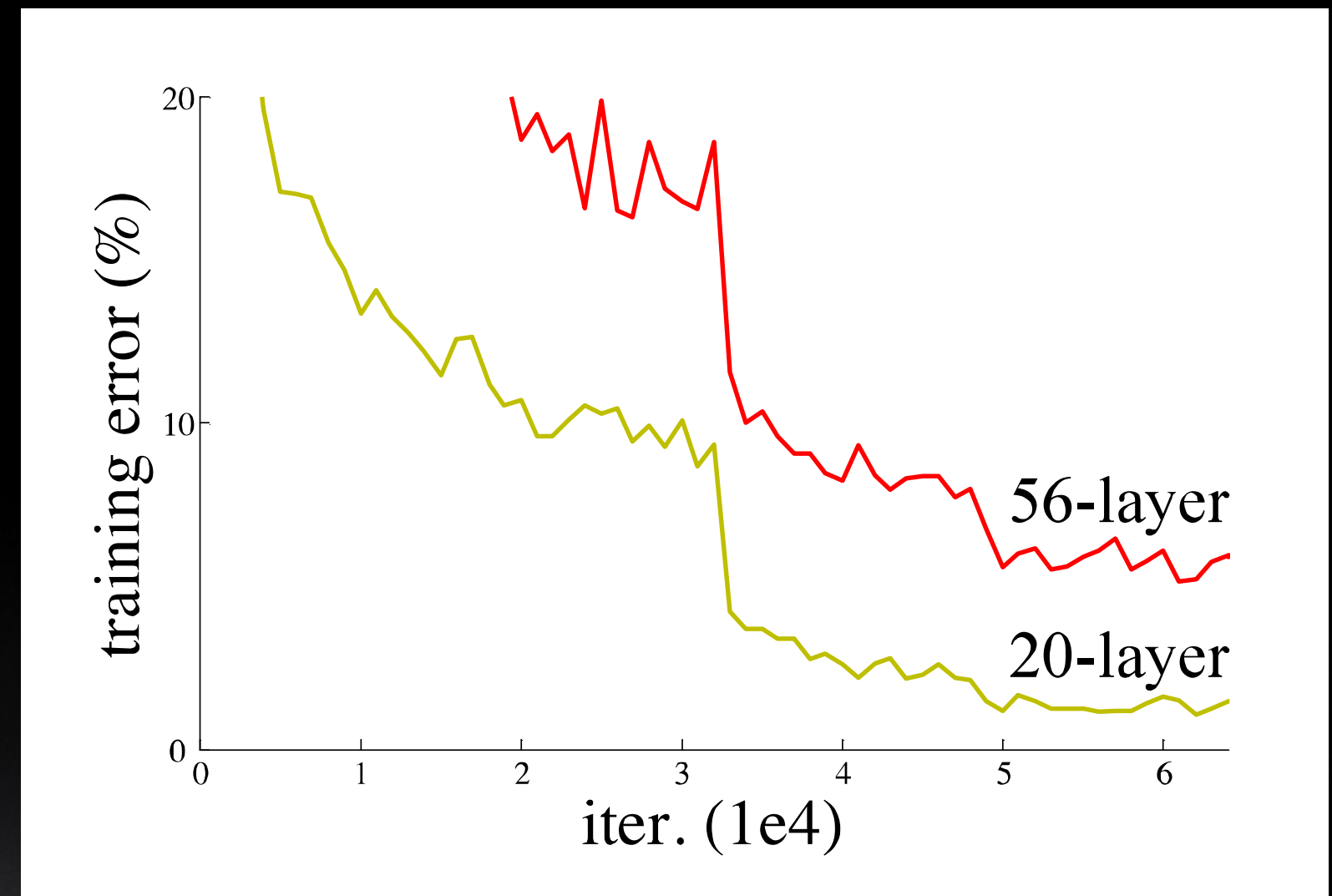
PROBLEM: VANISHING GRADIENTS

Error signal decays exponentially as it propagates backward through the network

ERROR SIGNAL VANISHES DURING BACKPROP



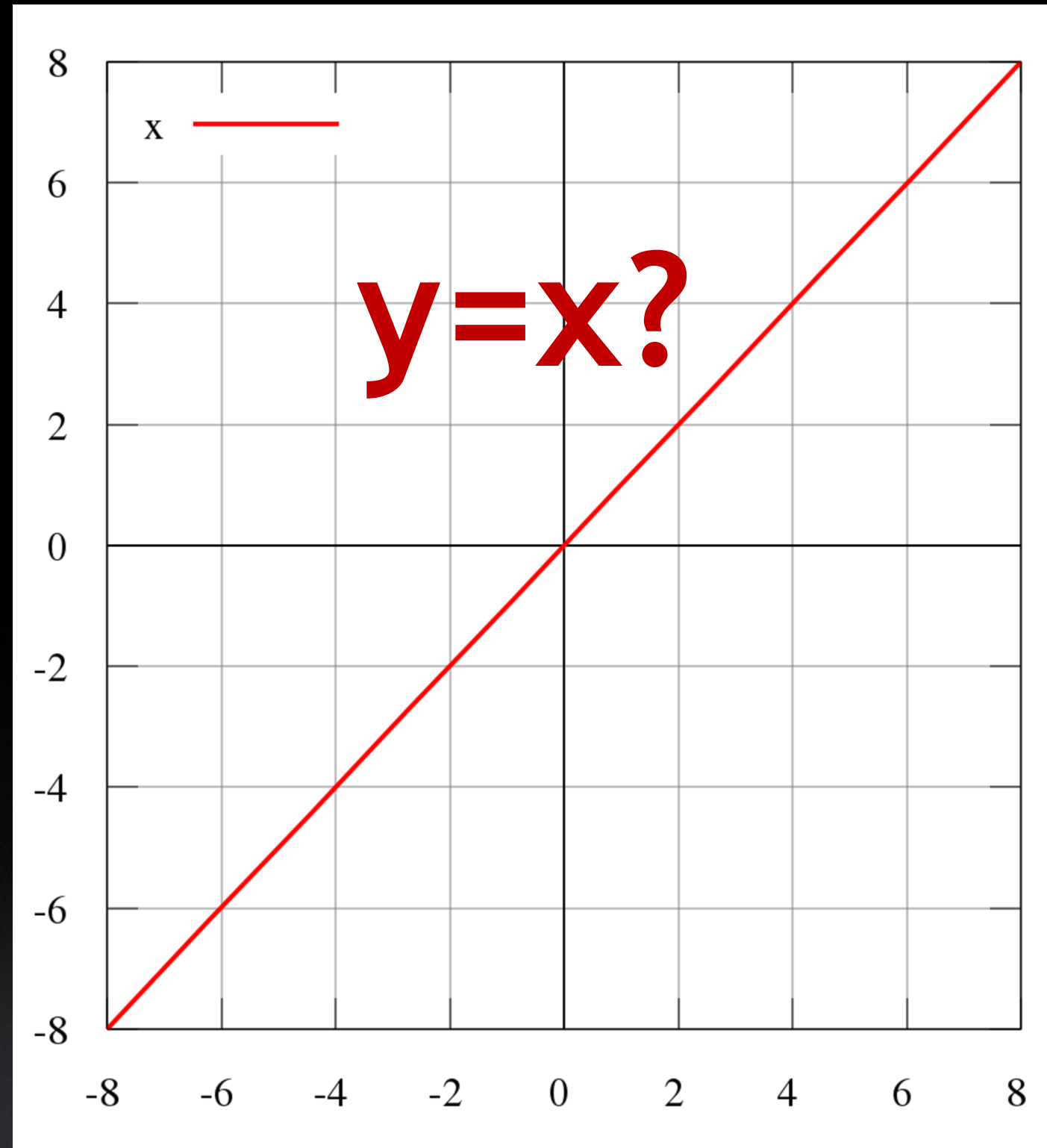
DEEPER NETWORKS WERE HARDER TO TRAIN



<https://www.arxiv-vanity.com/papers/1512.03385/>

PROBLEM: THE MISSING IDENTITY

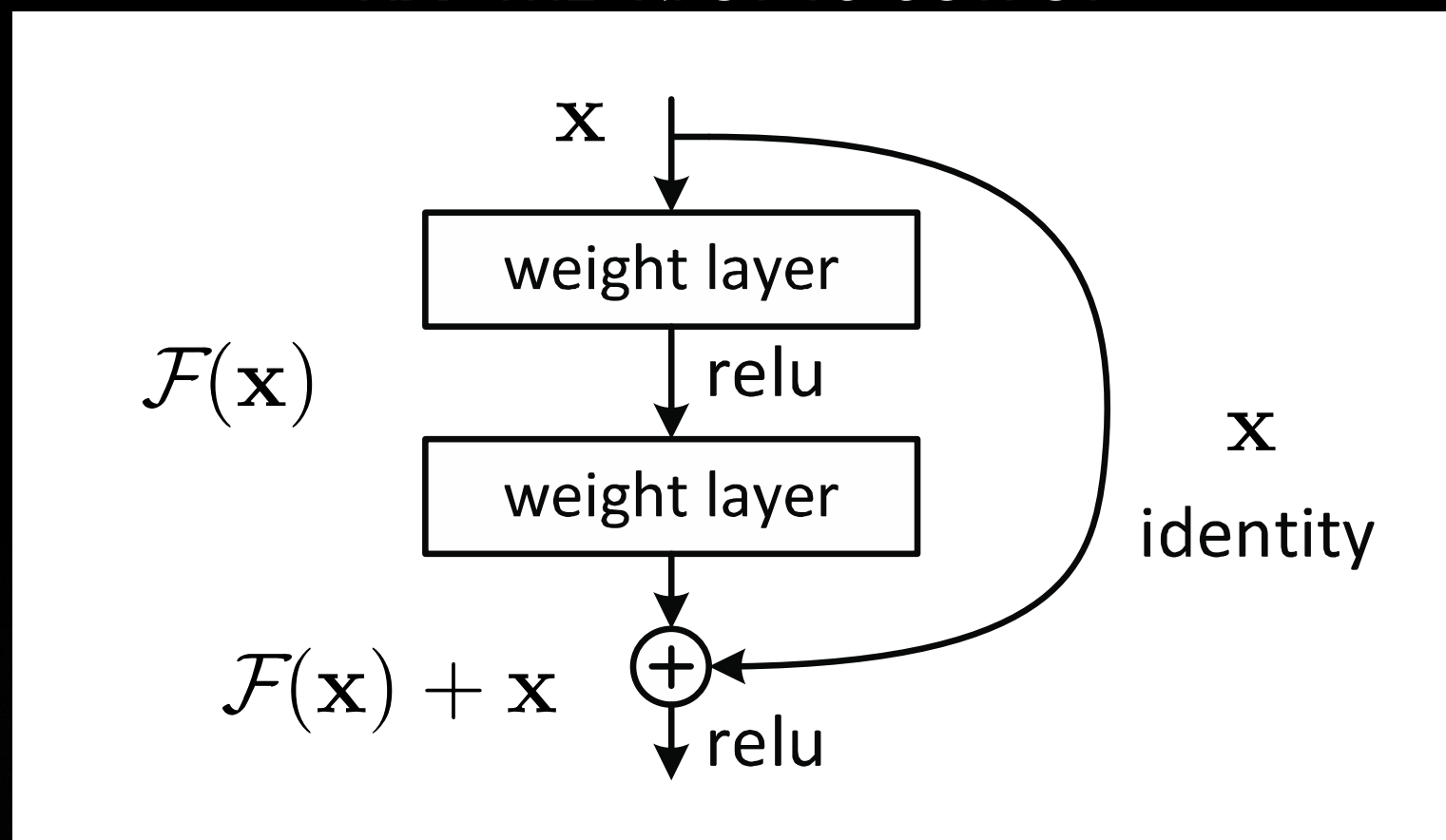
2015: Neural nets had a hard time learning the identity function!



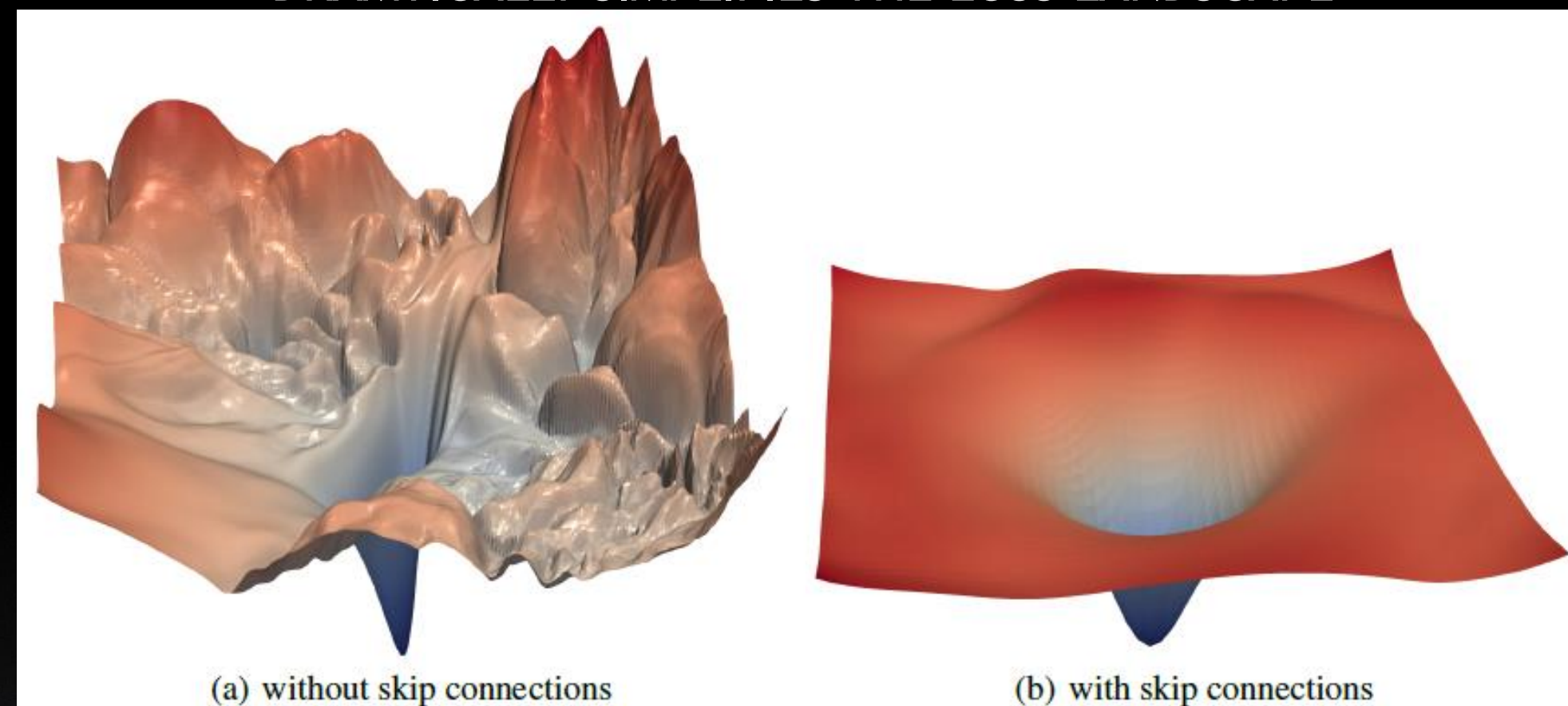
RESNETS AND SKIP CONNECTIONS

(aka Highway Networks)

ADD THE INPUT TO OUTPUT



DRAMTICALLY SIMPLIFIES THE LOSS LANDSCAPE



<https://arxiv.org/pdf/1512.03385.pdf>

<https://arxiv.org/abs/1712.09913>

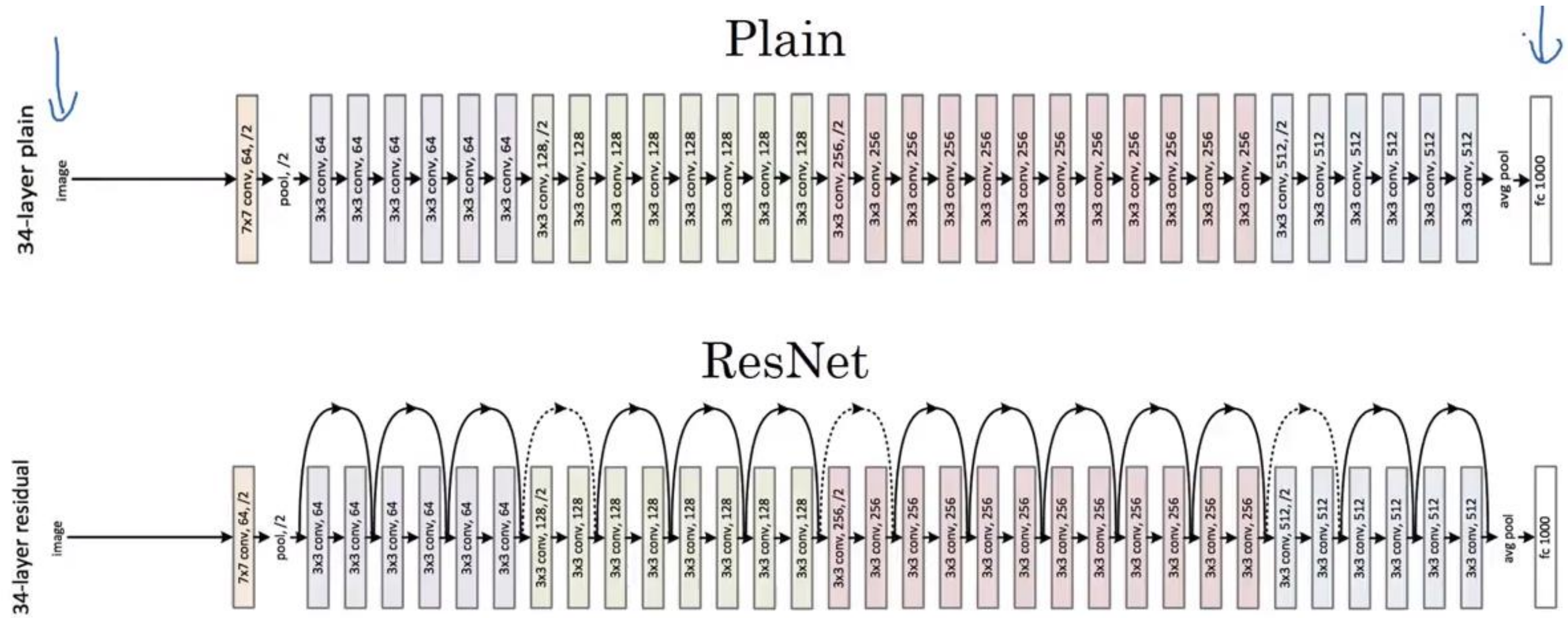
https://jithinjk.github.io/blog/nn_loss_visualized.md.html

RESNET-50

2015 Microsoft Research. 50 Layers, 23M params.

Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun
Microsoft Research



DENSENET

2017

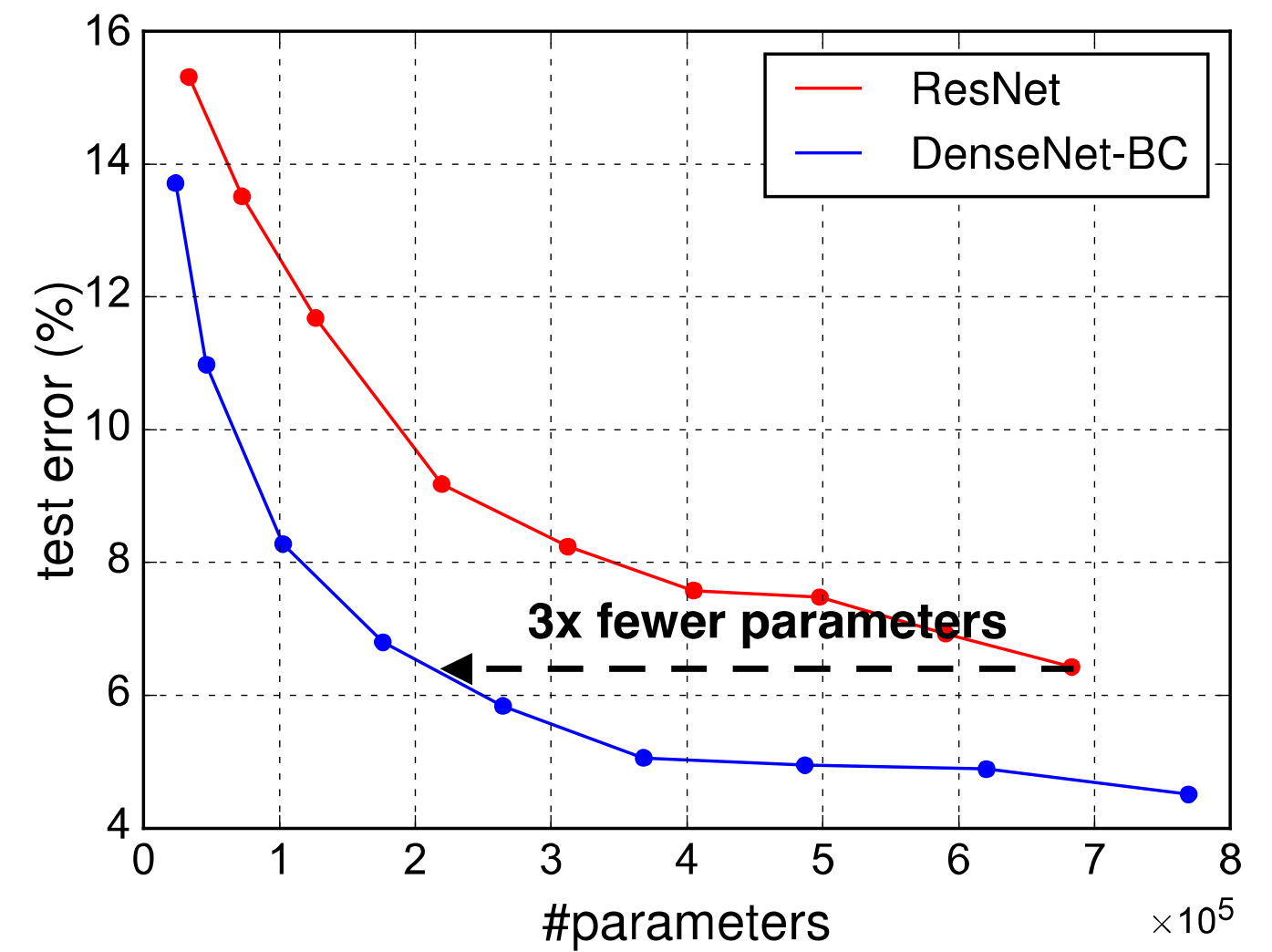
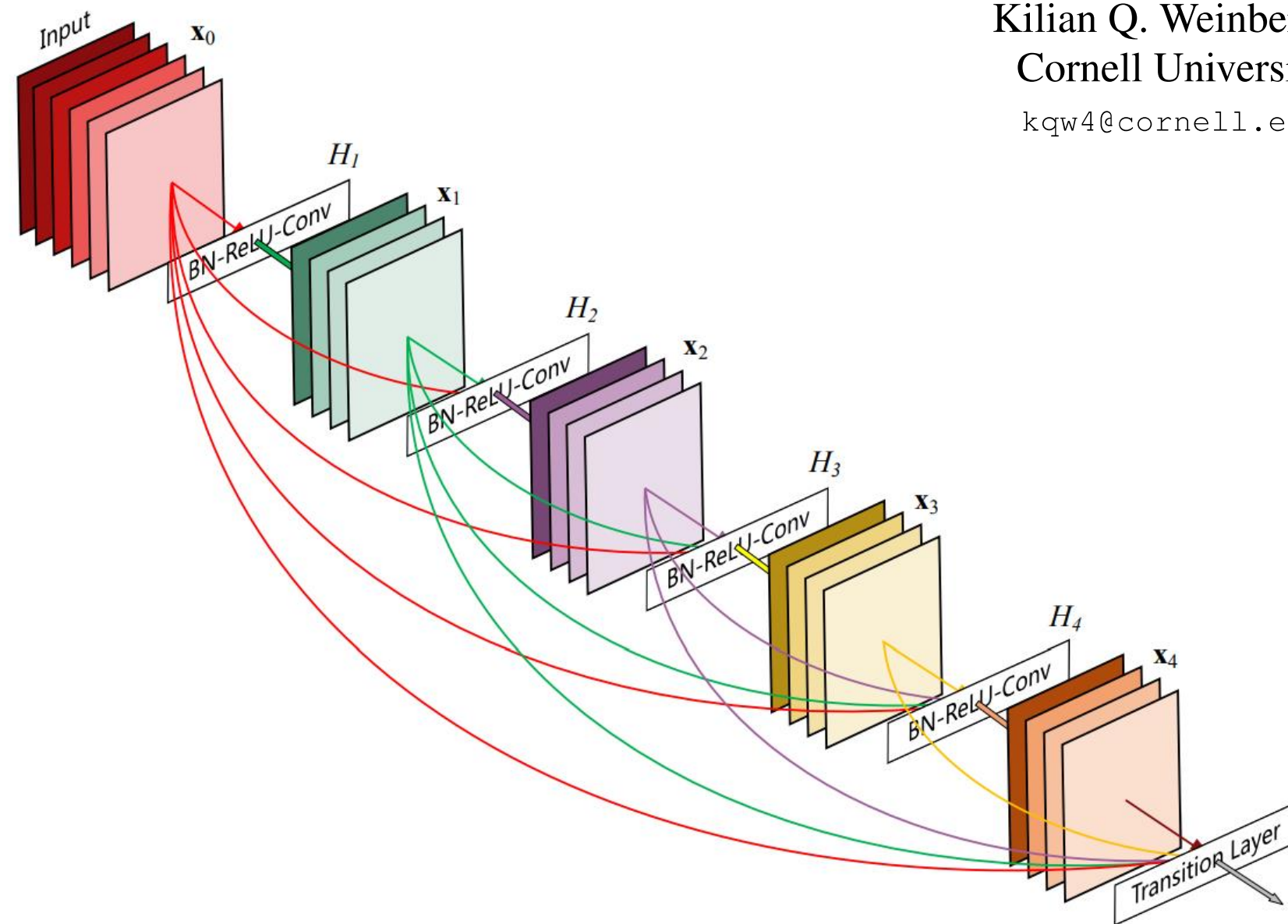
Densely Connected Convolutional Networks

Gao Huang*
Cornell University
qh349@cornell.edu

Zhuang Liu*
Tsinghua University
liuzhuang13@mails.tsinghua.edu.cn

Laurens van der Maaten
Facebook AI Research
lvdmaaten@fb.com

Kilian Q. Weinberger
Cornell University
kqw4@cornell.edu



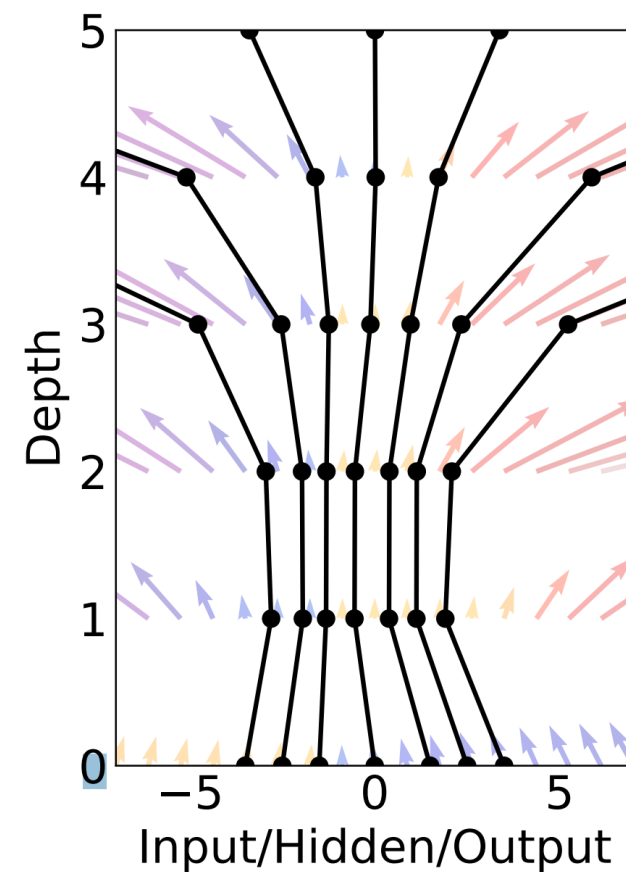
NEURAL-ODES

2018

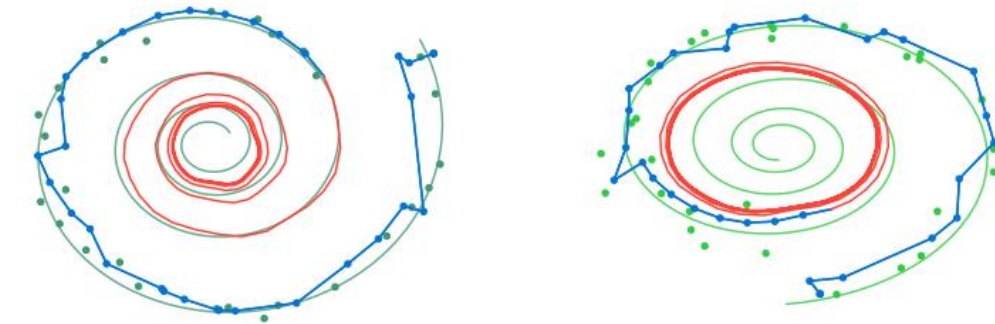
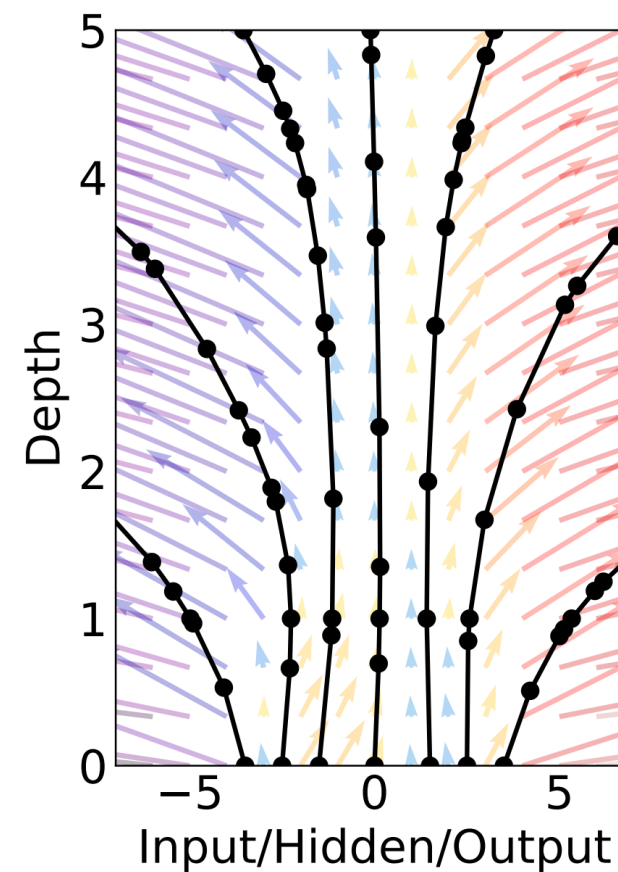
Neural Ordinary Differential Equations

Ricky T. Q. Chen*, Yulia Rubanova*, Jesse Bettencourt*, David Duvenaud
University of Toronto, Vector Institute
{rtqichen, rubanova, jessebett, duvenaud}@cs.toronto.edu

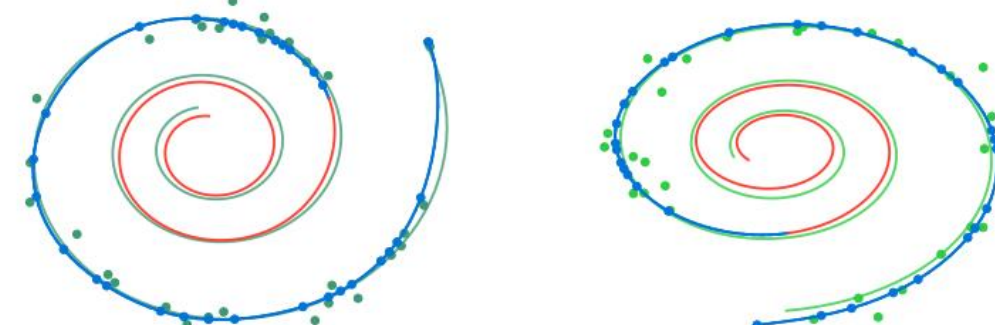
Residual Network



ODE Network

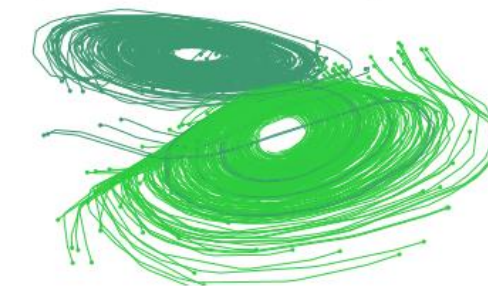


(a) Recurrent Neural Network



(b) Latent Neural Ordinary Differential Equation

- Ground Truth
- Observation
- Prediction
- Extrapolation



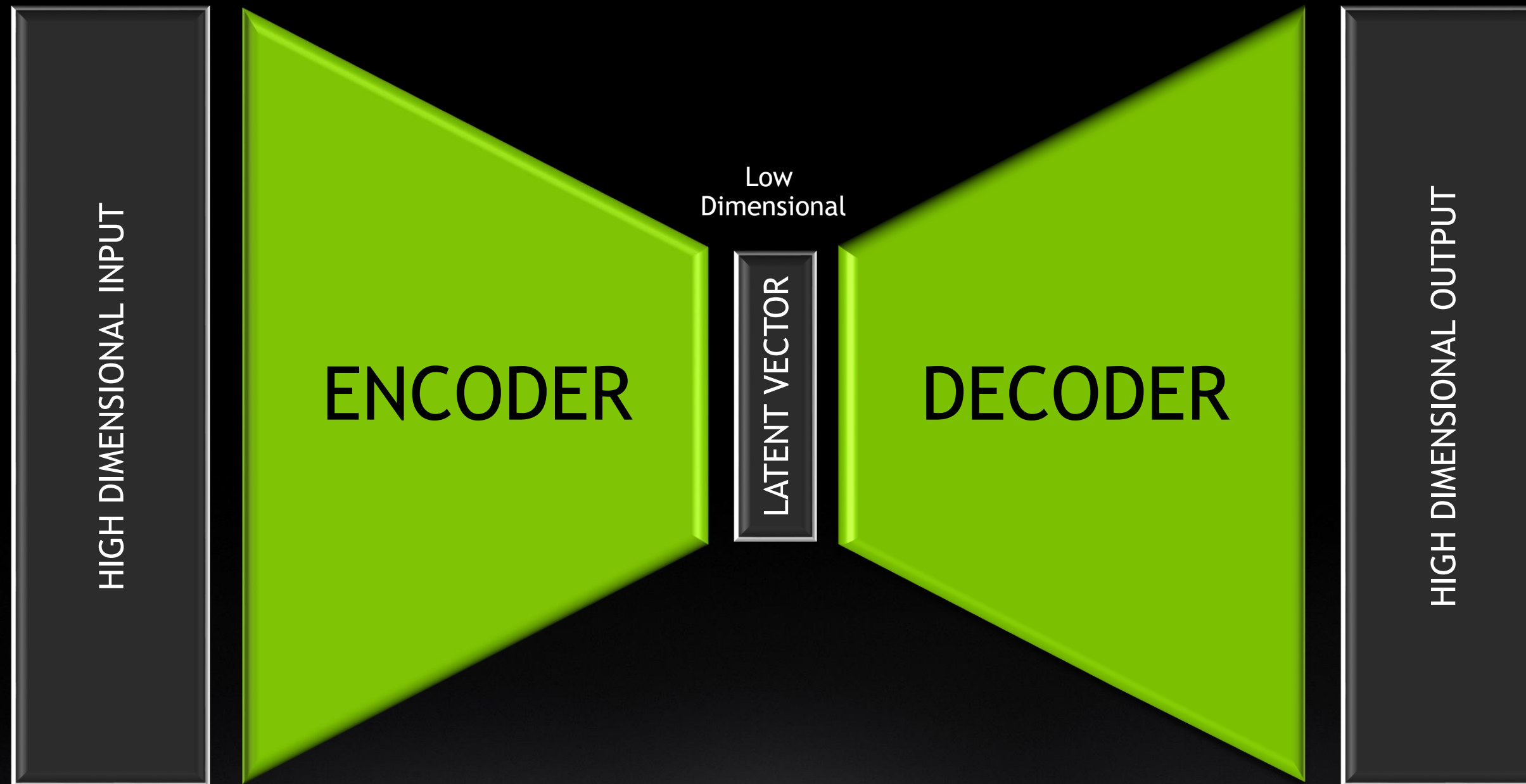
(c) Latent Trajectories



ENCODER-DECODERS

ENCODERS AND DECODERS

Networks connecting high and low dimensional spaces



CLASSIFIER: IMAGE → CLASS ENCODER

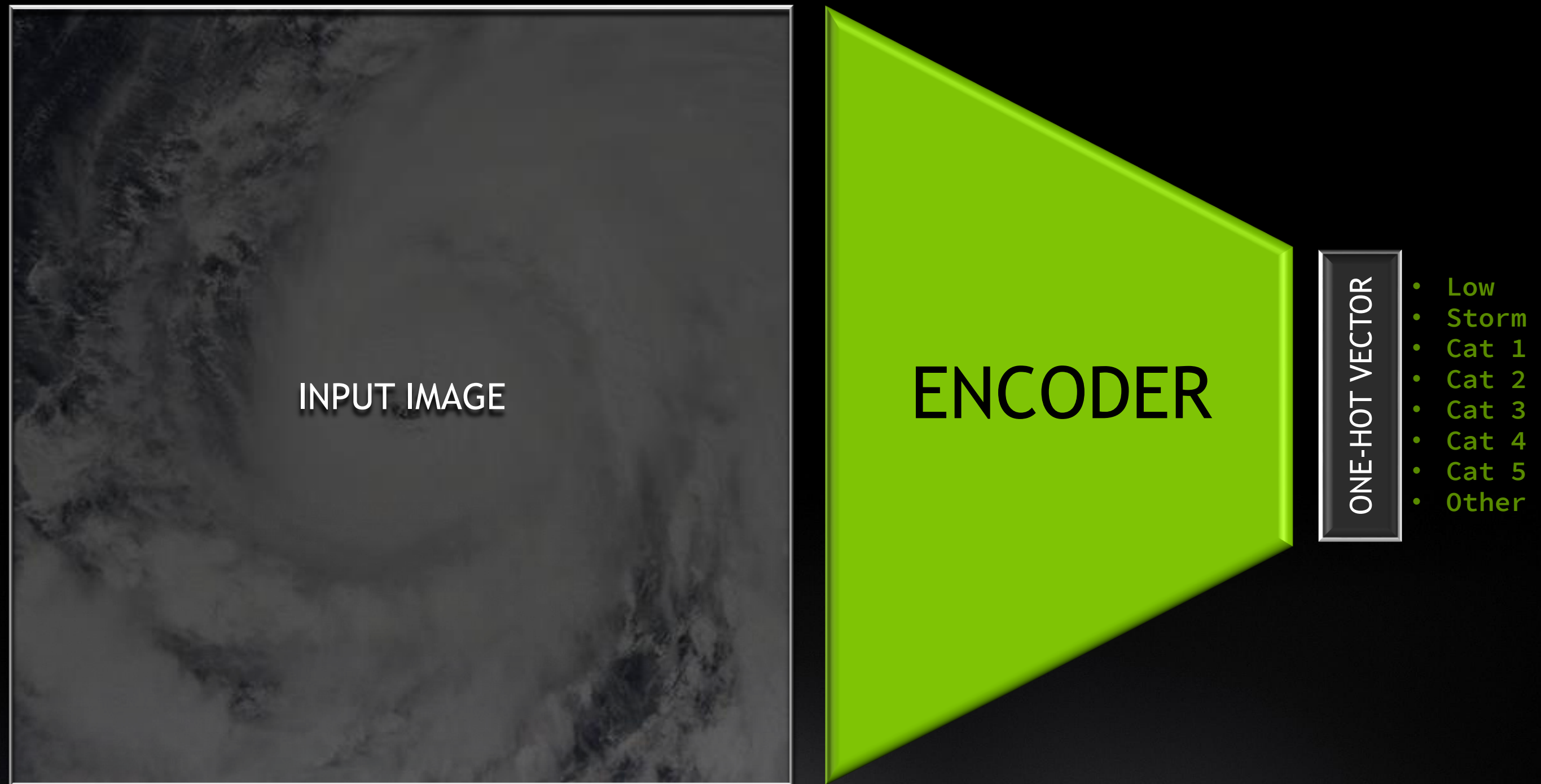
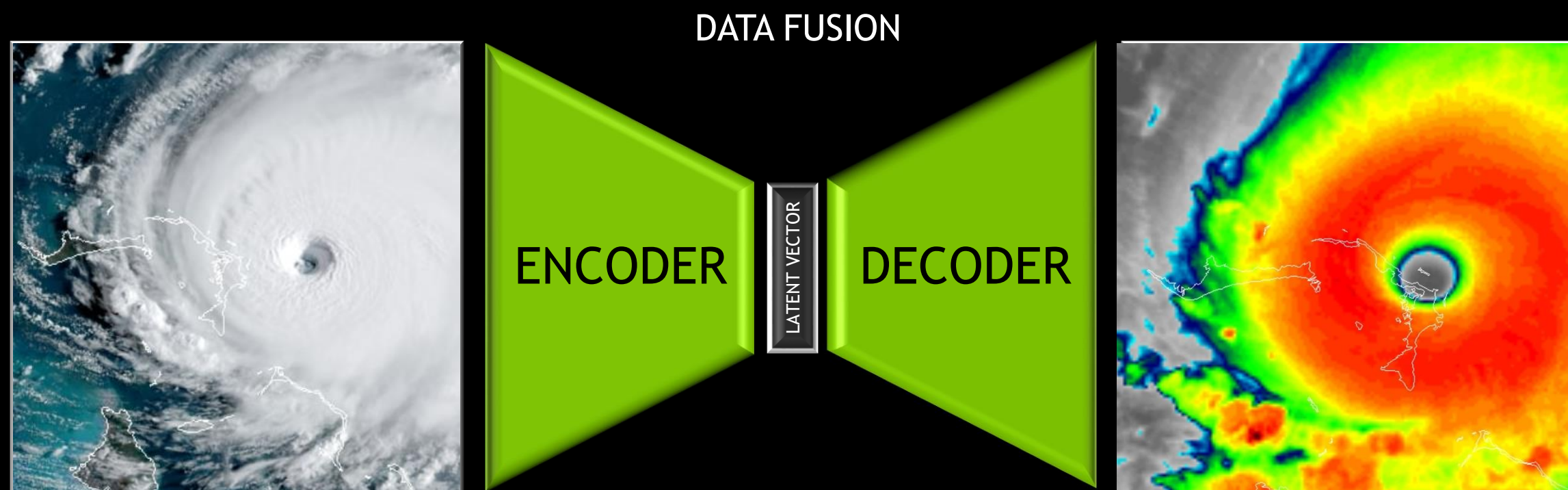
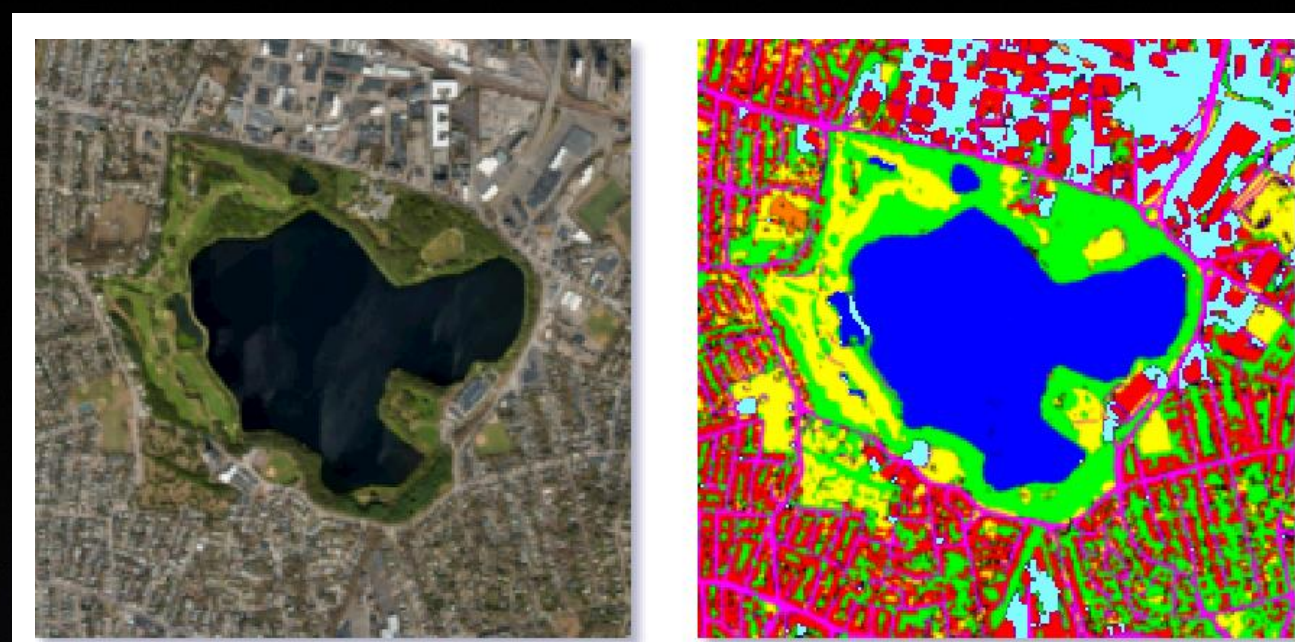


IMAGE TO IMAGE

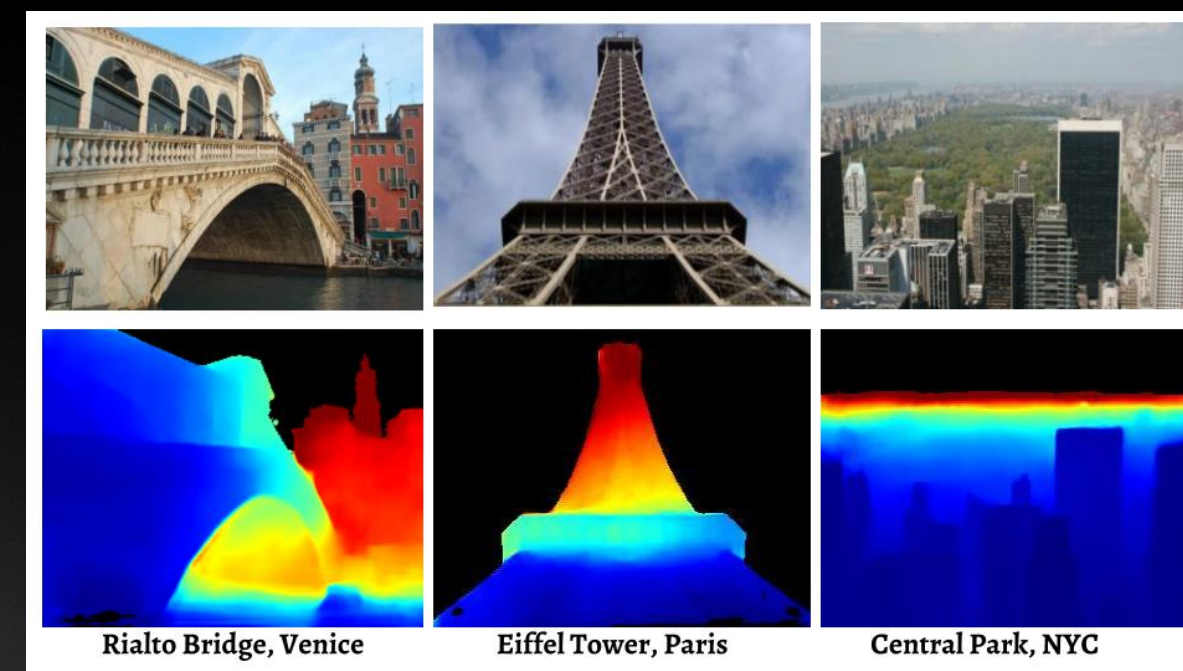
Encoder-Decoder network with Images at both ends



SEGMENTATION



DEPTH PREDICTION

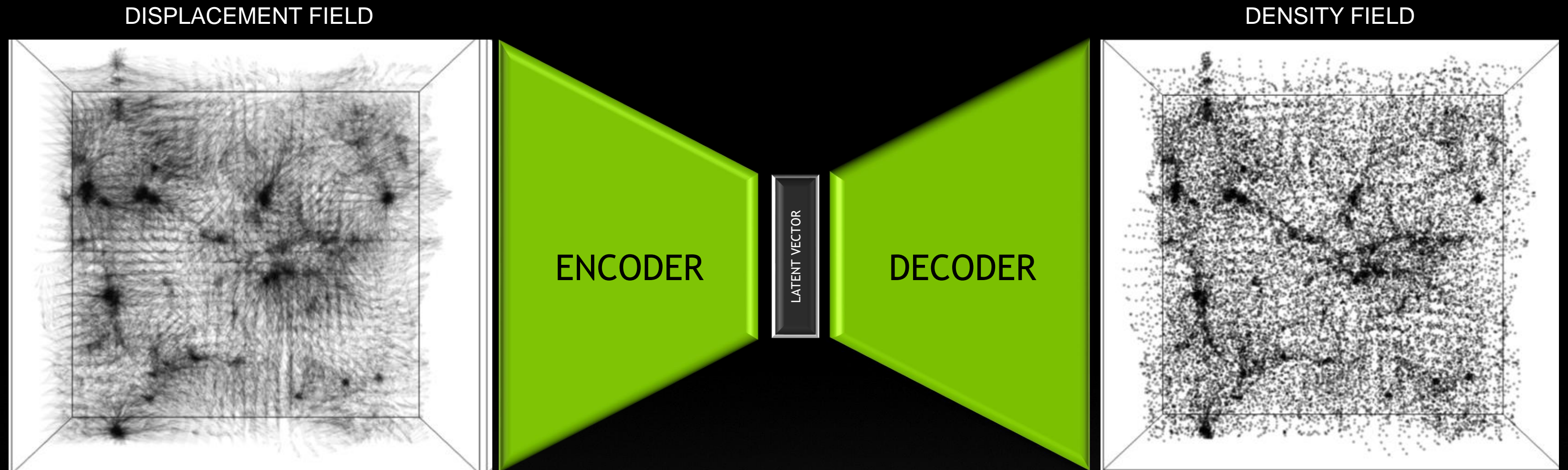


<http://liu.diva-portal.org/smash/get/diva2:1182913/FULLTEXT01.pdf>

<https://research.cs.cornell.edu/megadepth/>

VOLUME TO VOLUME

Input and Output can have 1,2,3 spatial dimensions or more

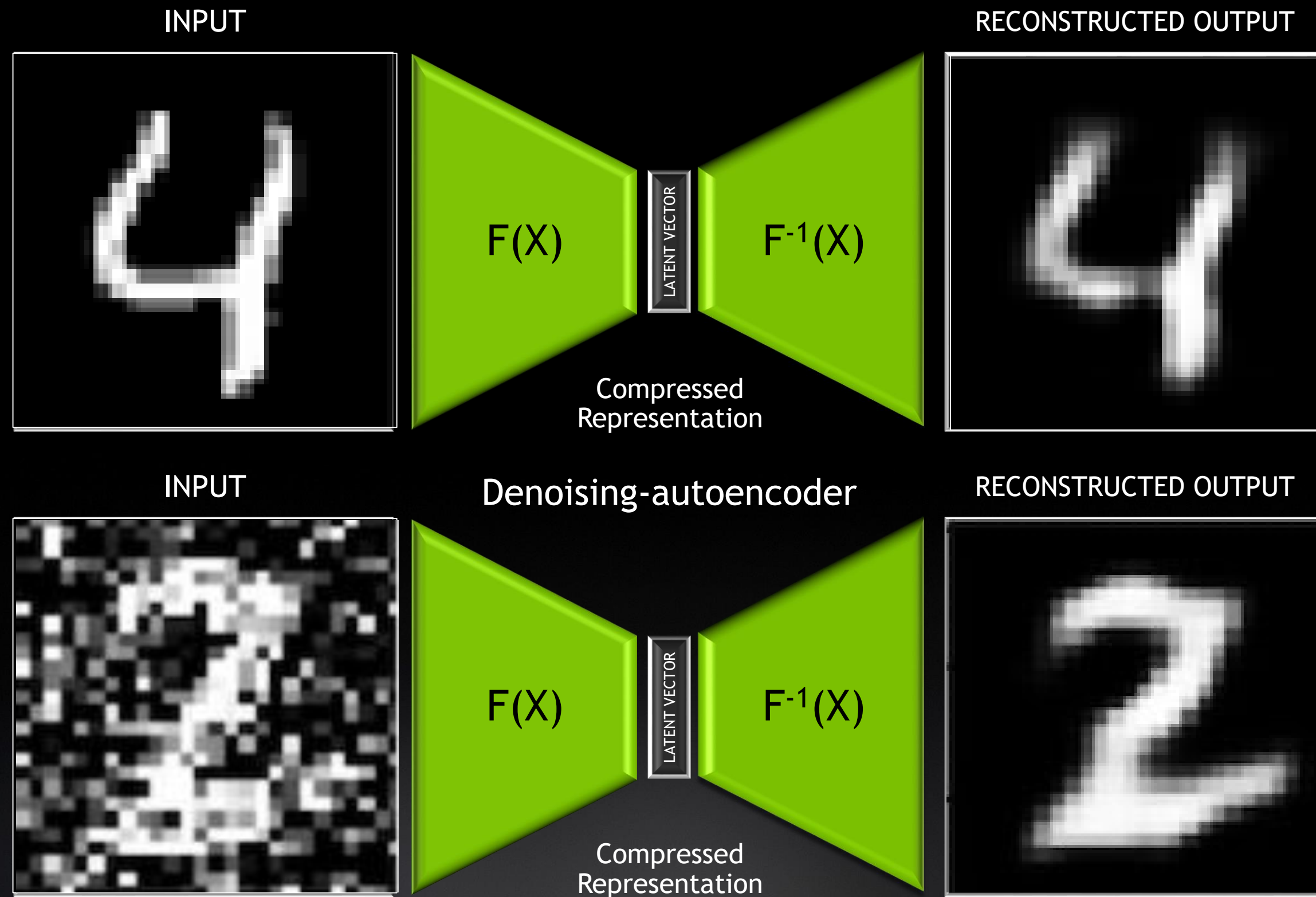


D3M: Learning to Predict the Cosmological Structure Formation

<https://arxiv.org/pdf/1811.06533.pdf>

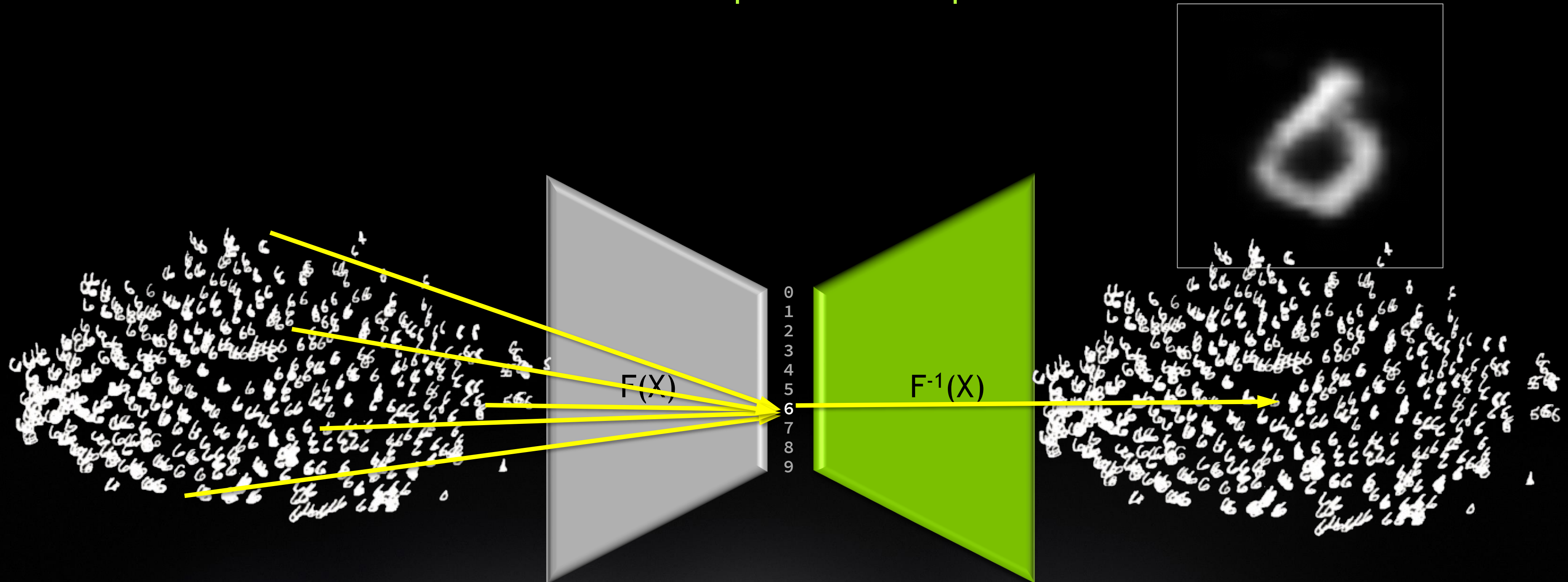
AUTOENCODER

Adaptive Data Compression and Noise Removal



WHY IS MY DECODER OUTPUT FUZZY?

Decoded Output is Not Unique

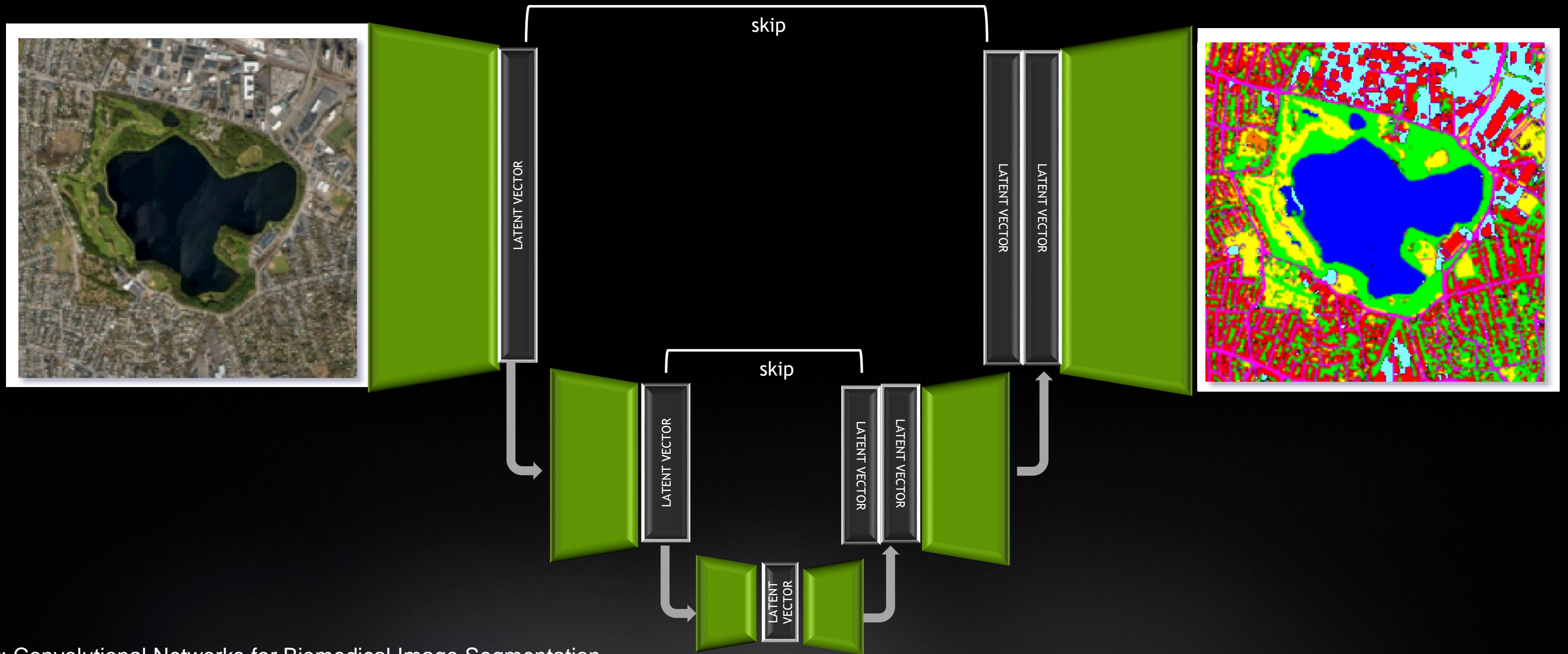


Many instances map
to one class

One class represents
many instances

UNET (2015)

Nested encoder-decoders at multiple spatial scales



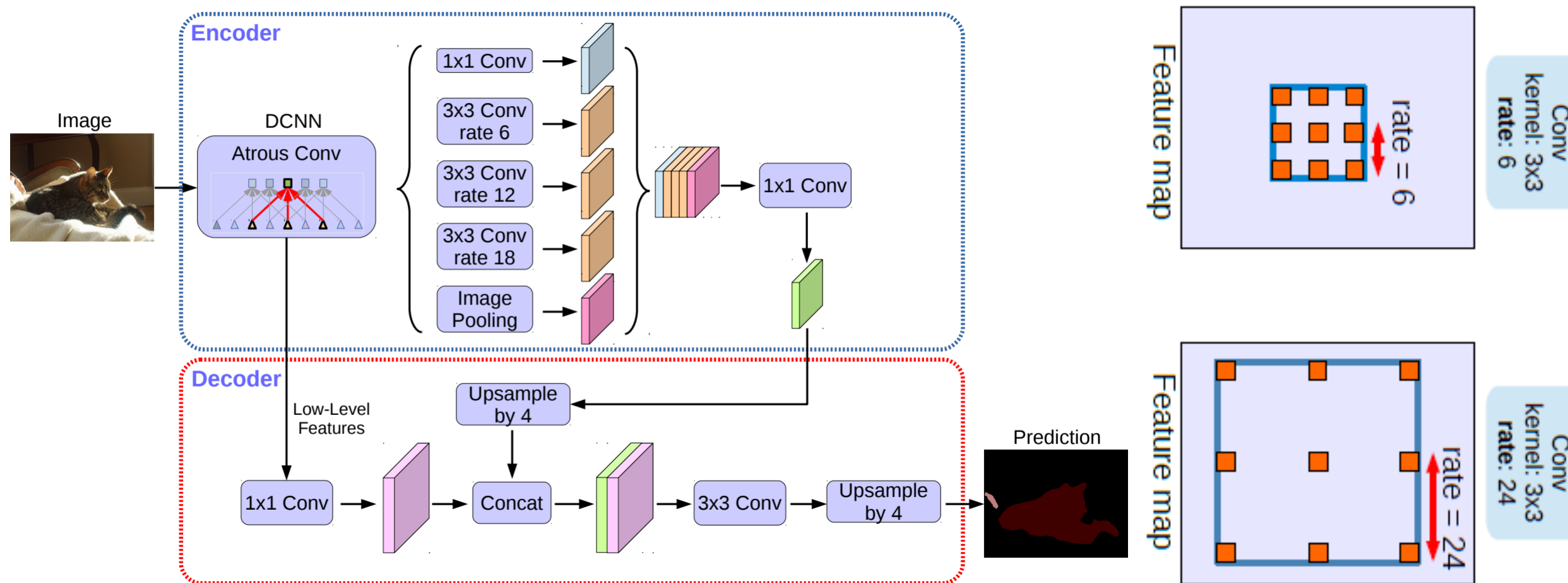
U-Net: Convolutional Networks for Biomedical Image Segmentation
Olaf Ronneberger, Philipp Fischer, and Thomas Brox
<https://arxiv.org/pdf/1505.04597.pdf>

DEEPLAB V3+

Another encoder decoder design for accurate segmentation

Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation

Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam

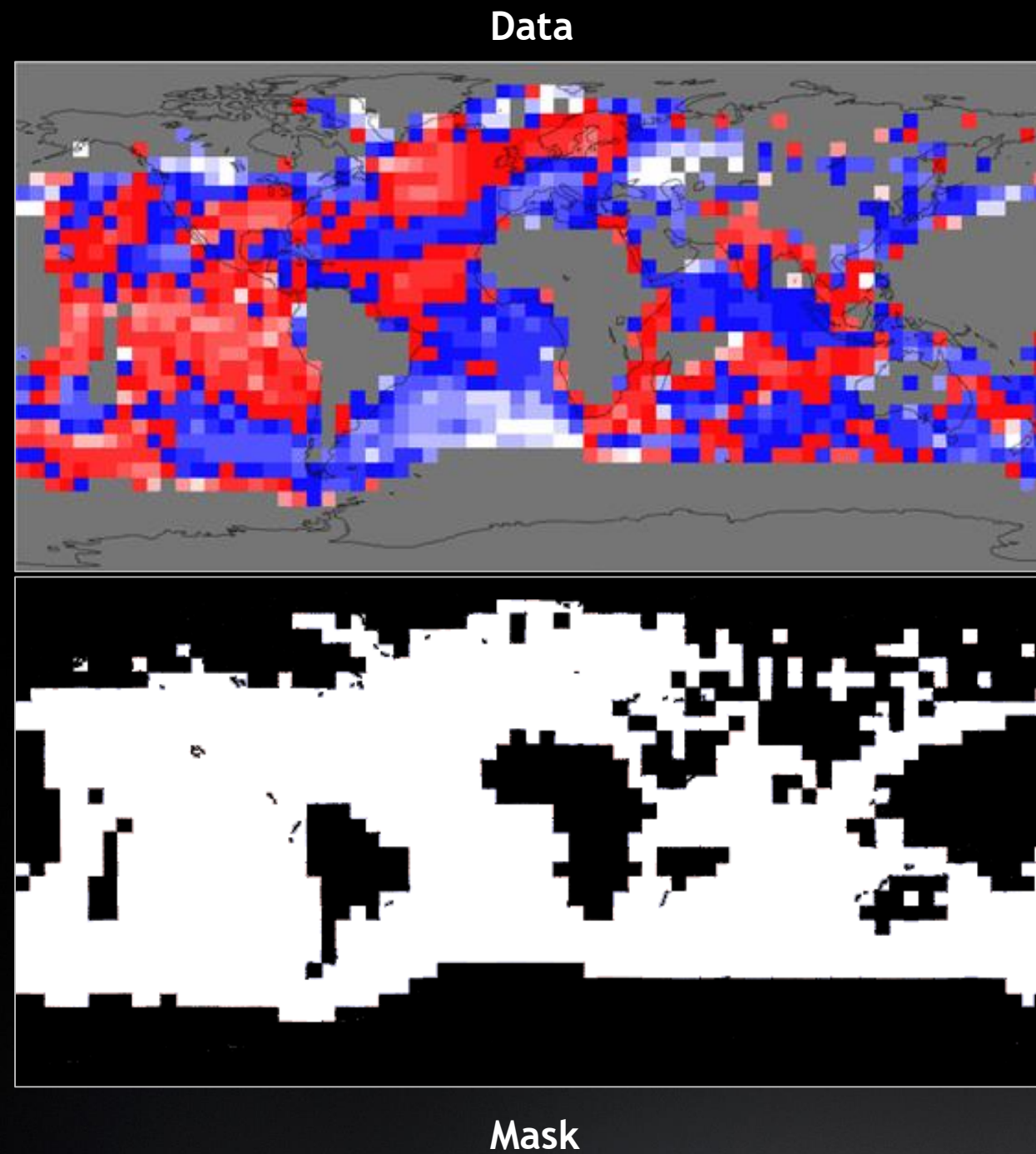




MASKED CONVOLUTIONS AND INPAINTING

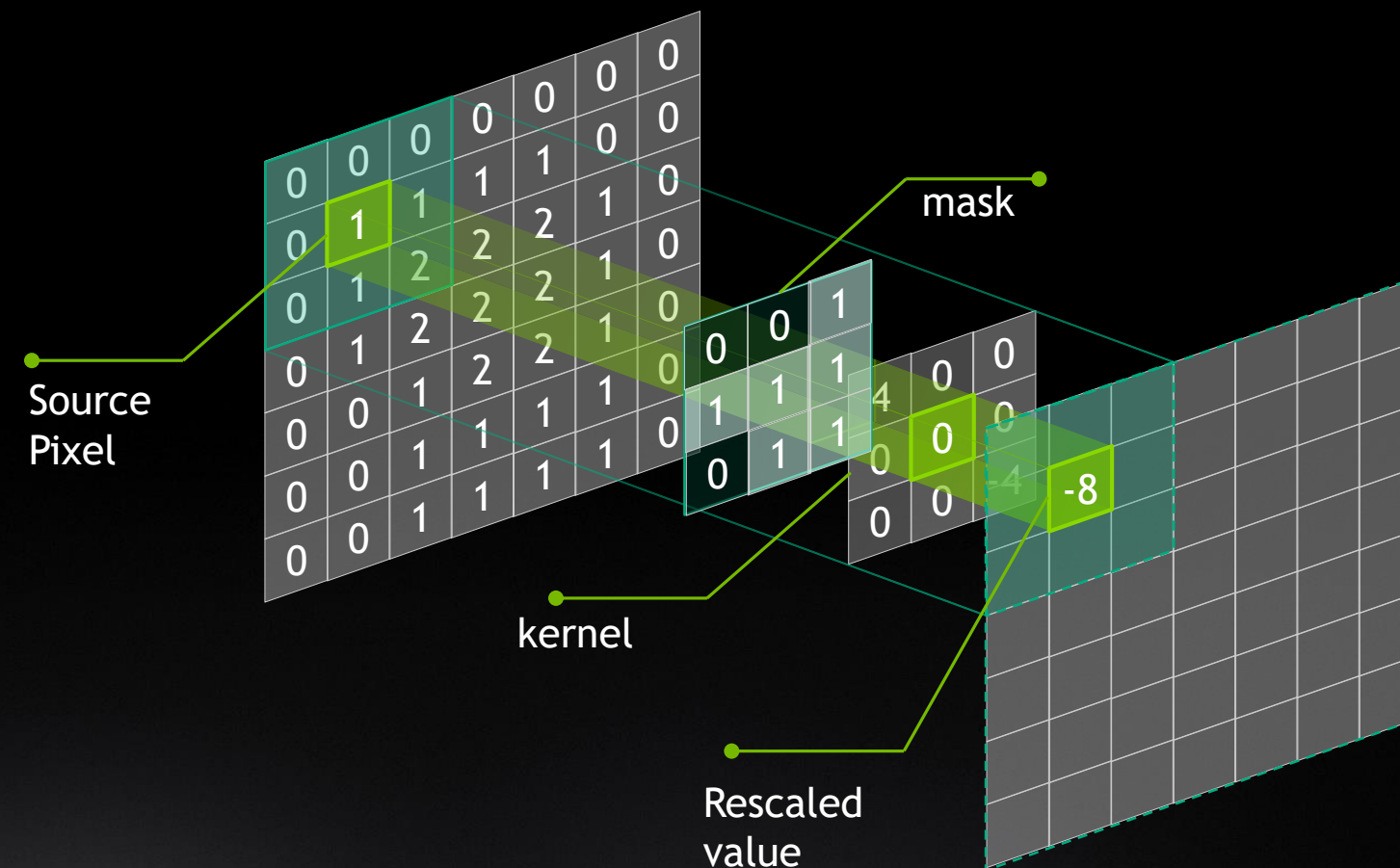
PARTIAL CONVOLUTIONS

Convolutions that ignore missing or invalid data



Masked Convolution Operation

$$x' = \begin{cases} \mathbf{W}^T (\mathbf{X} \odot \mathbf{M}) \frac{\text{sum}(\mathbf{1})}{\text{sum}(\mathbf{M})} + b, & \text{if } \text{sum}(\mathbf{M}) > 0 \\ 0, & \text{otherwise} \end{cases}$$



TASK: INPAINTING

Repair an image that has missing data

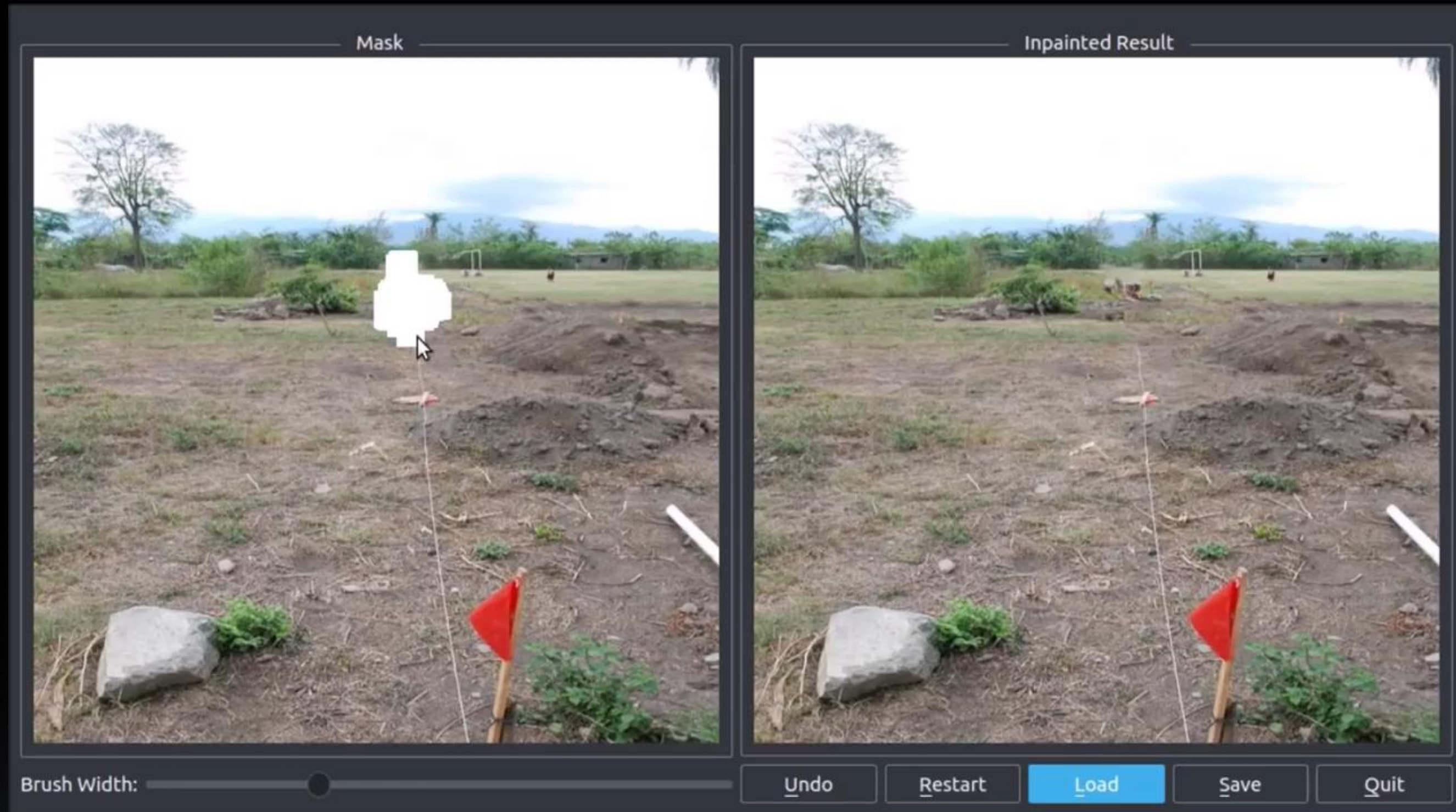


Image Inpainting for Irregular Holes Using Partial Convolutions

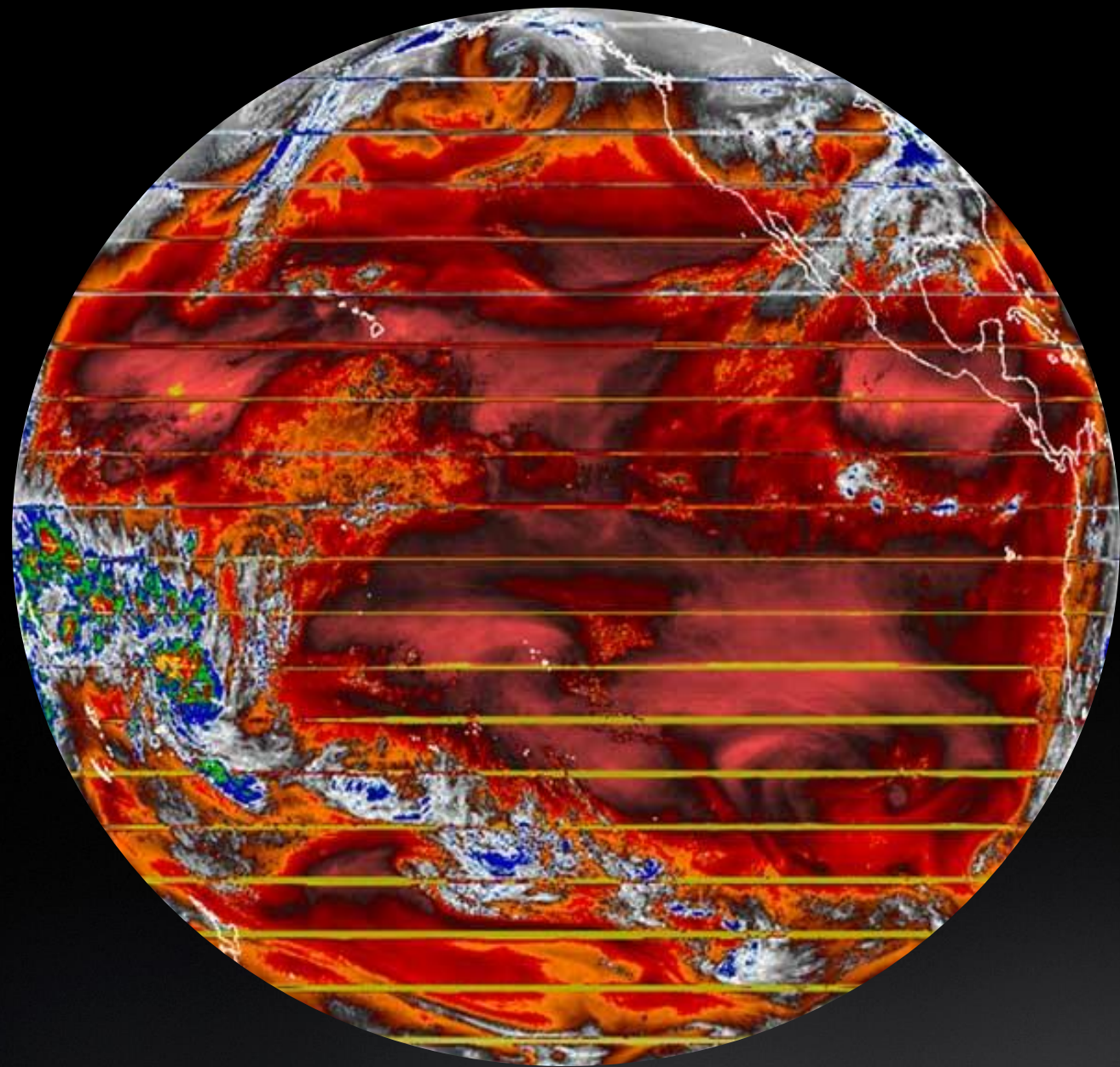
Guilin Liu Fitsum A. Reda Kevin J. Shih Ting-Chun Wang Andrew Tao Bryan Catanzaro

NVIDIA Corporation

INPAINTING APPLICATIONS

Fill in missing observations, or remove unwanted objects

GOES-17: Repair Missing Data



Remove Clouds, Haze, Shadows

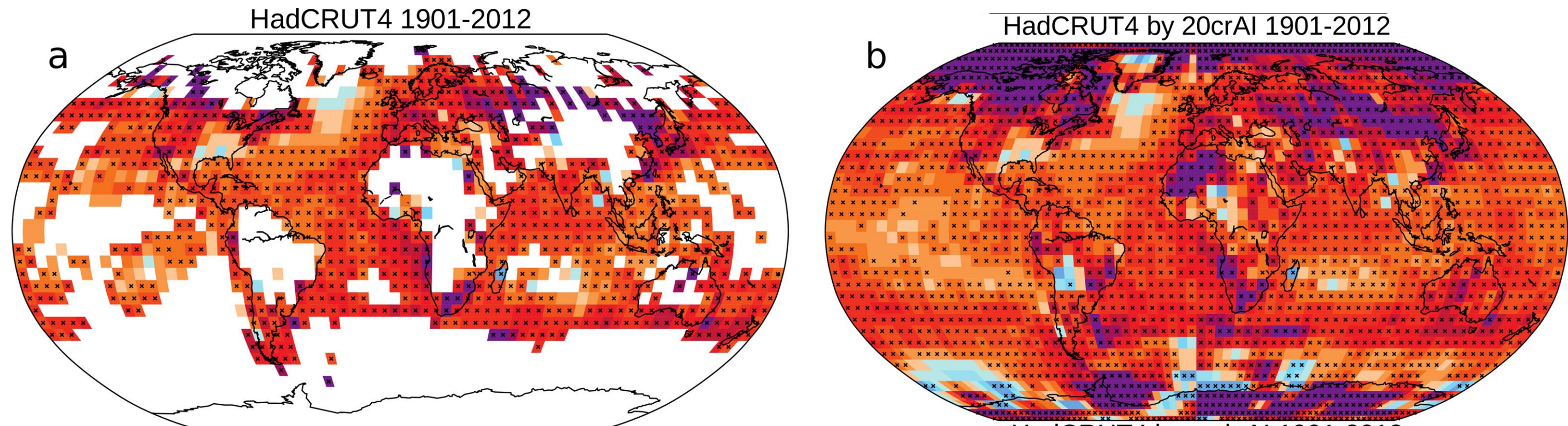


TRANSFER LEARNING FOR INPAINTING

Train on model data to repair observational data

Artificial intelligence reconstructs missing climate information

Christopher Kadow^{1,2}, David Matthew Hall³ and Uwe Ulbrich²



<https://www.nature.com/articles/s41561-020-0582-5>



GENERATIVE MODELS

GENERATIVE MODELS

Generate Specific Examples from a Learned Distribution



A Style-Based Generator Architecture for Generative Adversarial Networks

<https://thispersondoesnotexist.com>

<https://arxiv.org/pdf/1812.04948.pdf>

GENERATE NEW EXAMPLES

You can generate a new example of nearly any type of data



Neural Rhapsody
maia (2019)
inspired by Wolfgang Mozart (1756-1791)

$\text{♩} = 120$

Context (human-written): In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

GPT-2: The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

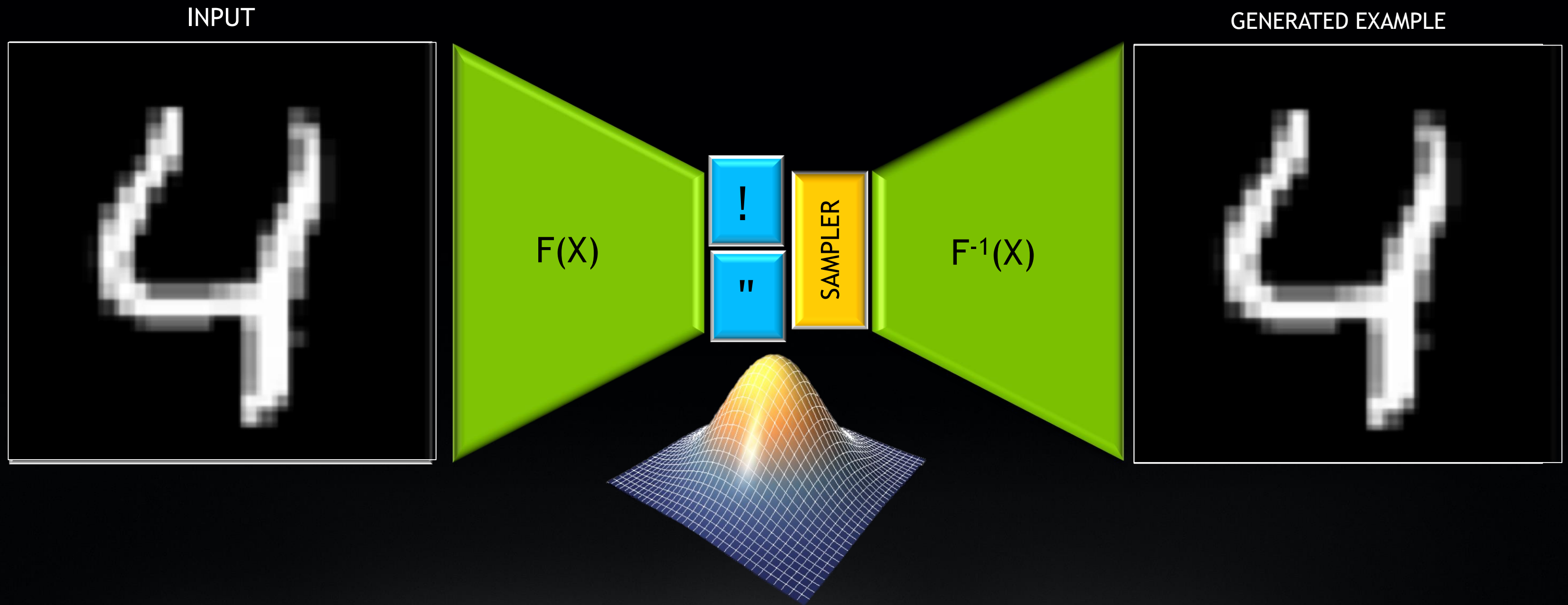
Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top

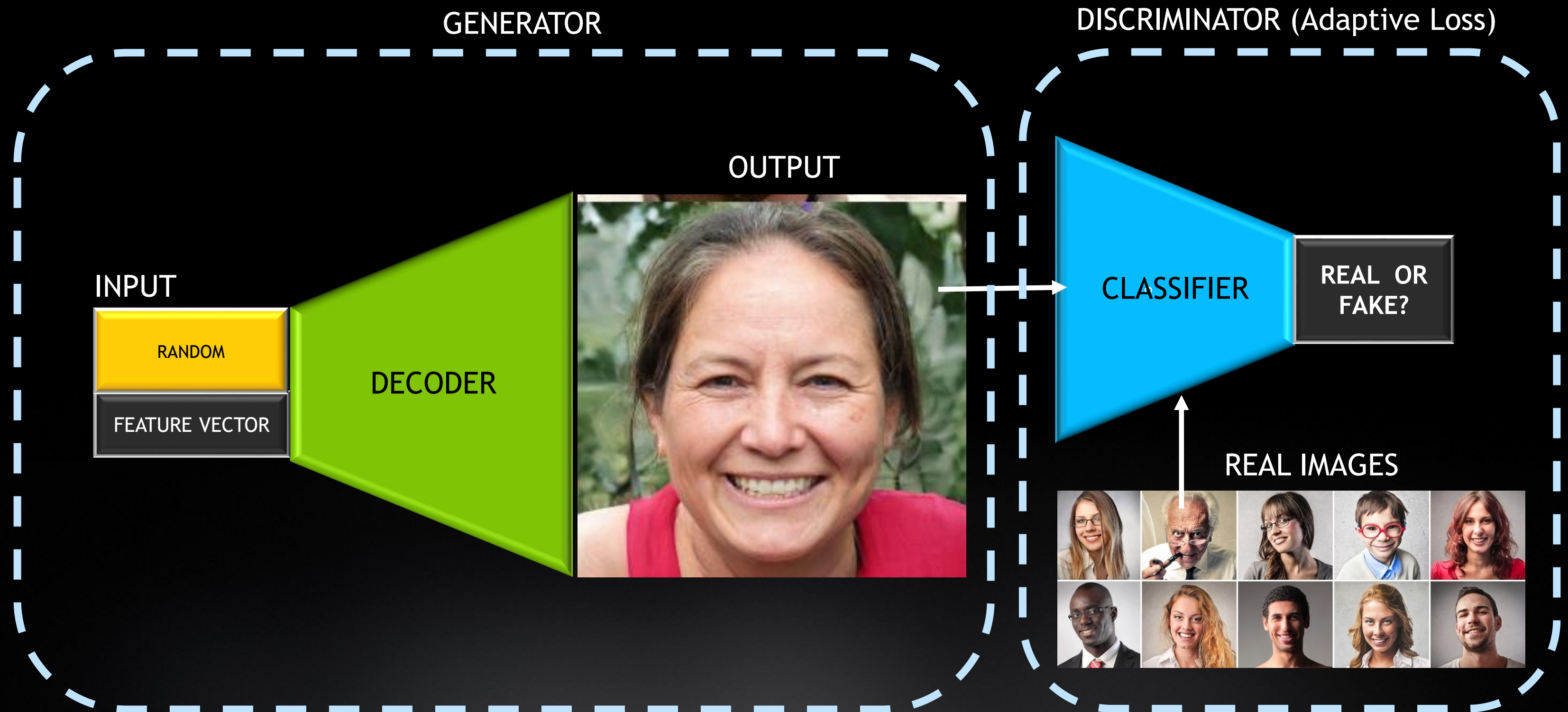
VAE: VARIATIONAL AUTOENCODER

An autoencoder that learns Gaussian Distributions



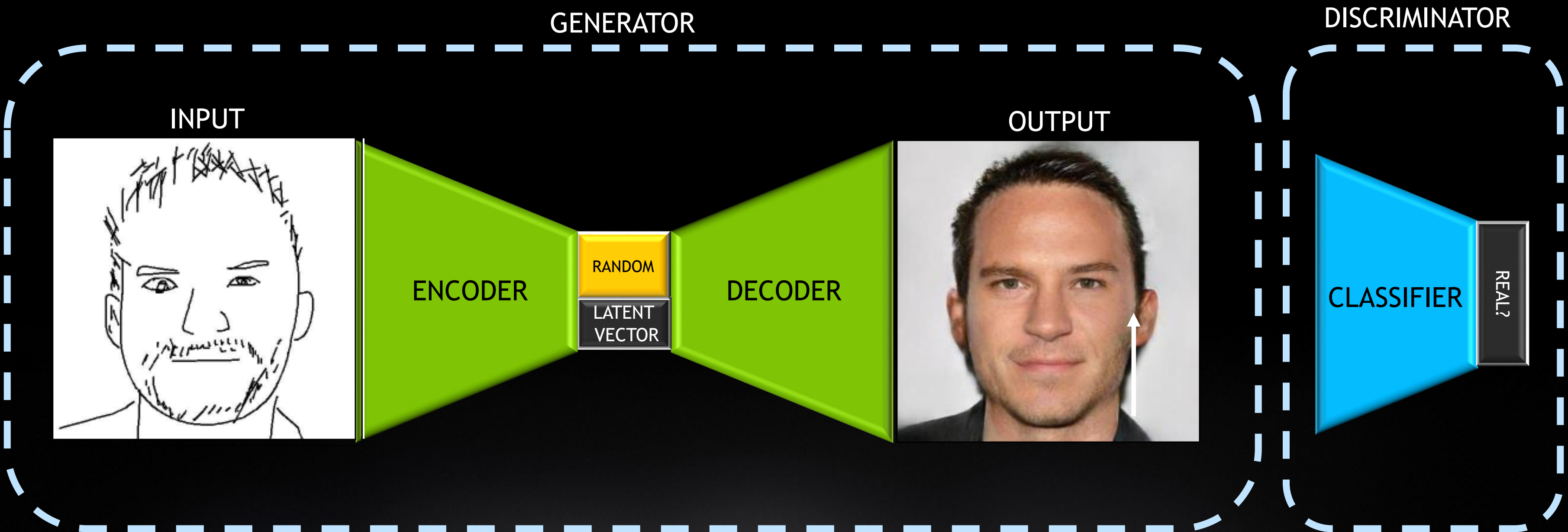
GAN: GENERATIVE ADVERSARIAL NETWORK

A trick for training a decoder to produce samples indistinguishable from real ones



CONDITIONAL GAN

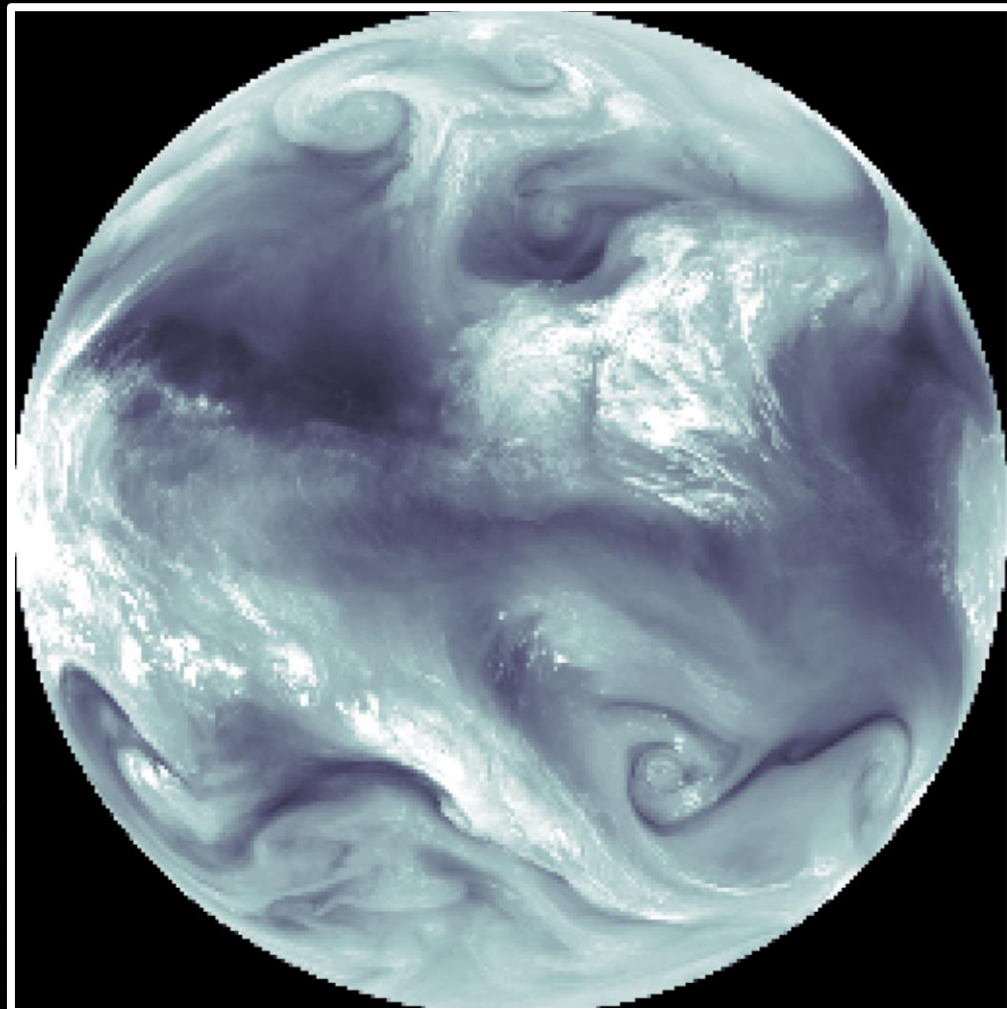
Generate Synthetic Images, conditioned upon the input



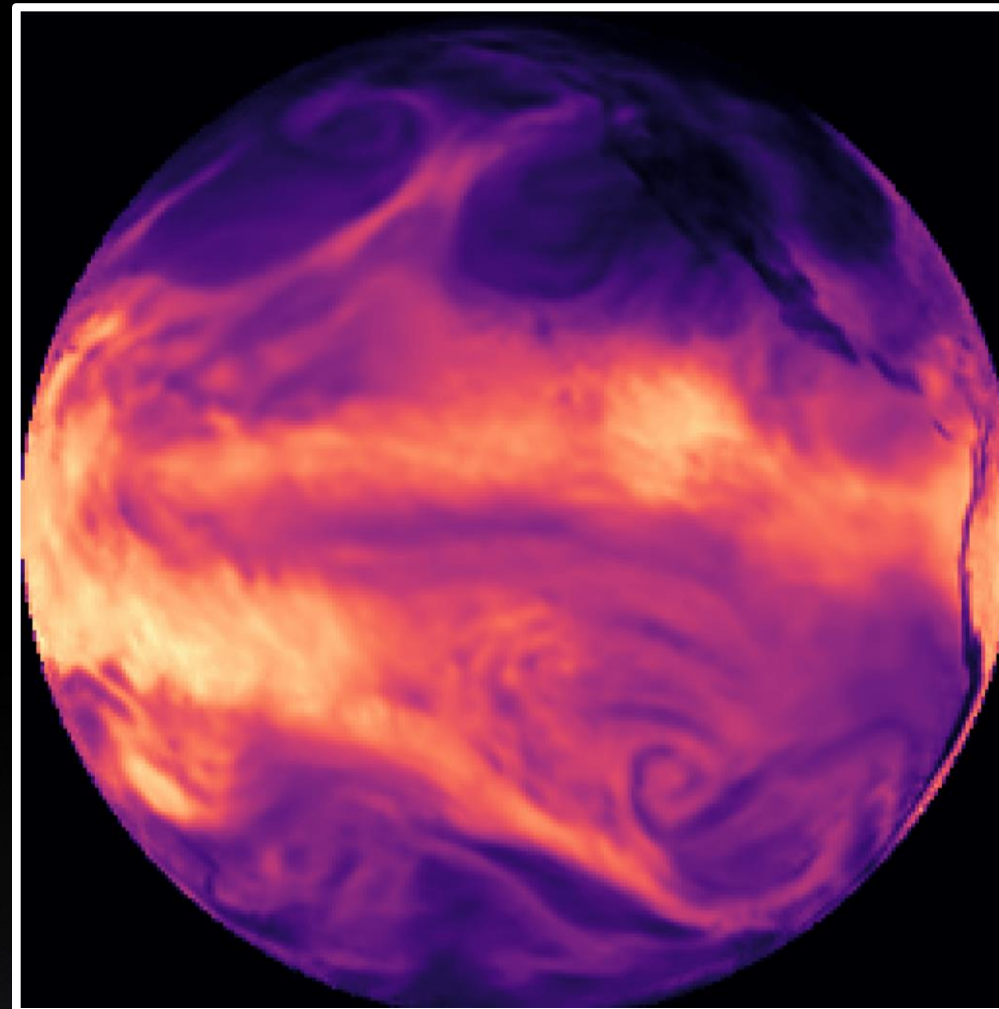
CONDITIONAL GAN EXAMPLE

Map from satellite observations to model variables

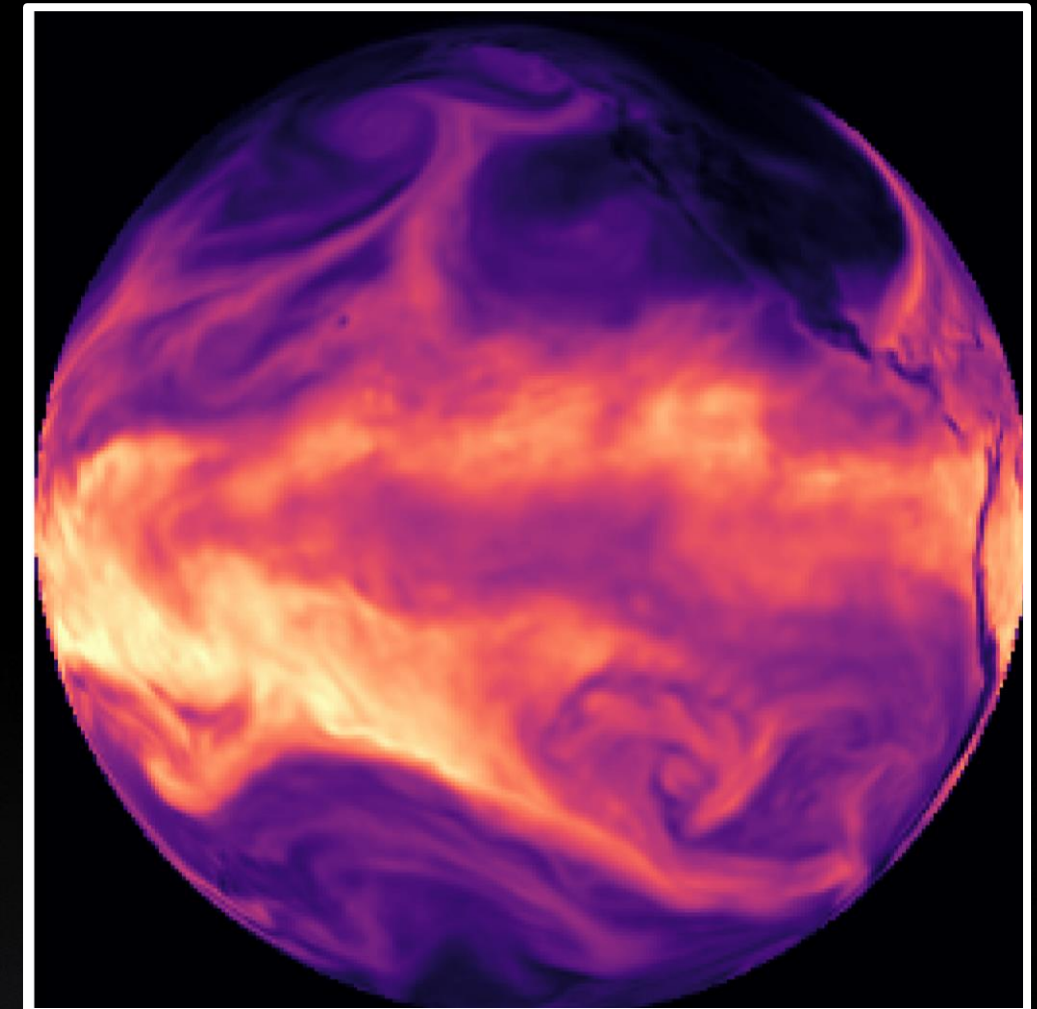
GOES-15 BAND 3



GENERATED WATER VAPOR



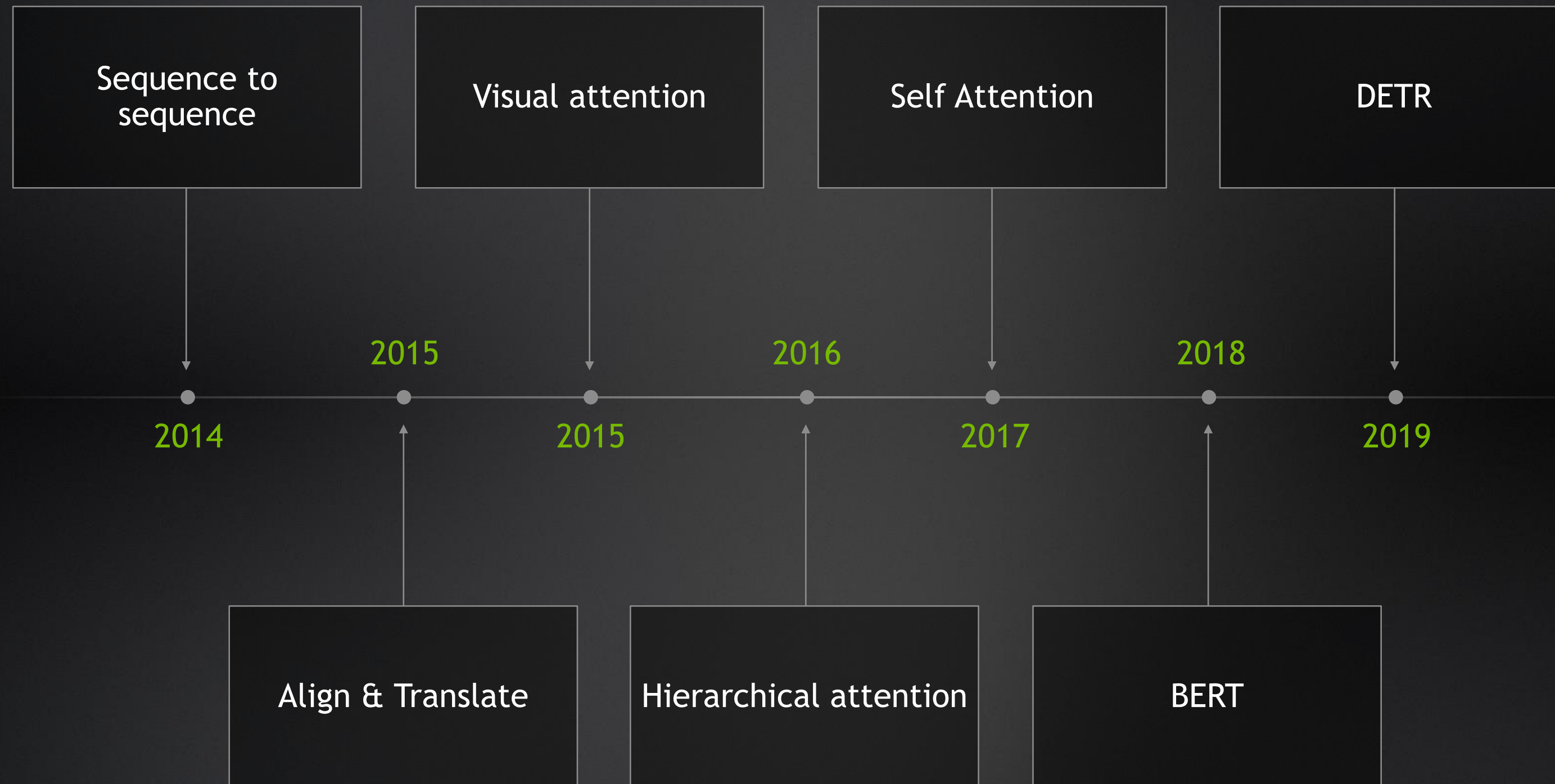
TARGET: GFS WATER VAPOR





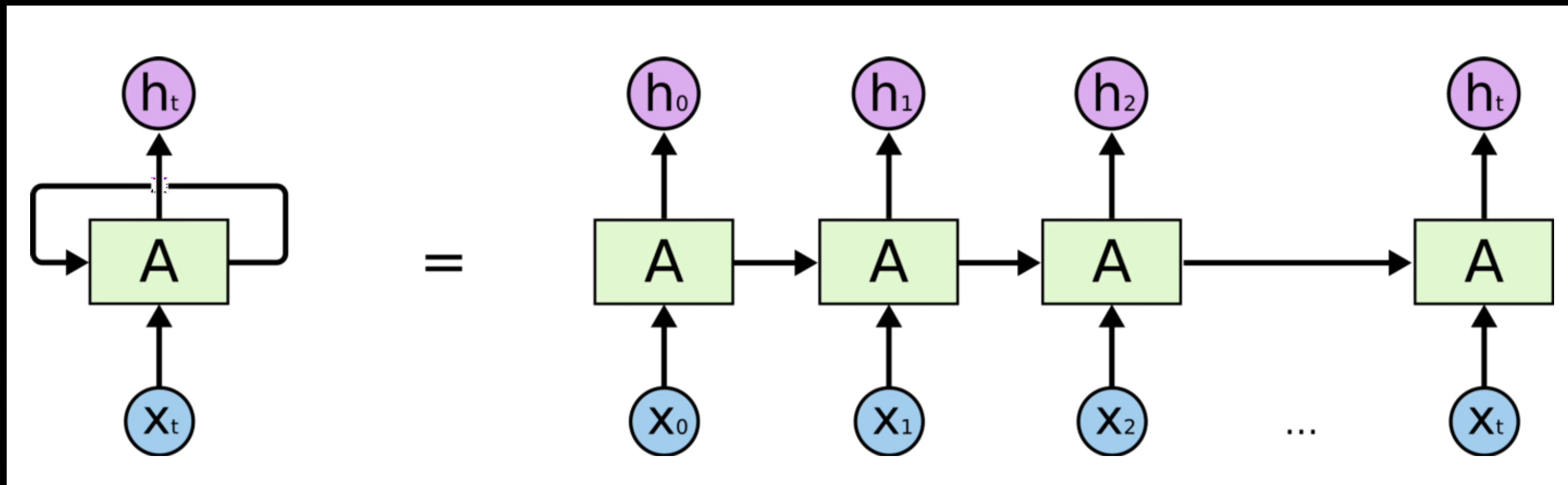
RNNS, ATTENTION, AND TRANSFORMERS

ATTENTION AND TRANSFORMERS



RECURRENT NEURAL NETS

(1986) Neural networks for sequences



RNN, Unrolled over time

LSTM

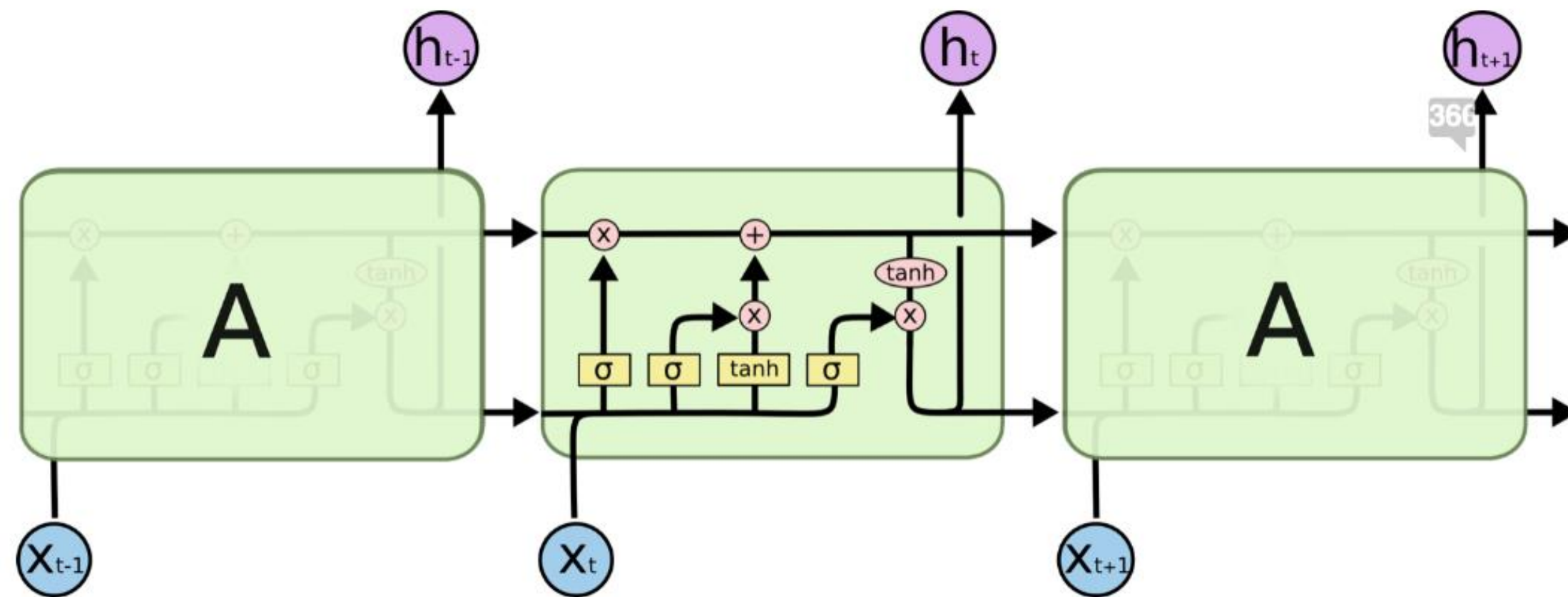
(1997) Long Short-term Memory Units

LONG SHORT-TERM MEMORY

NEURAL COMPUTATION 9(8):1735-1780, 1997

Sepp Hochreiter

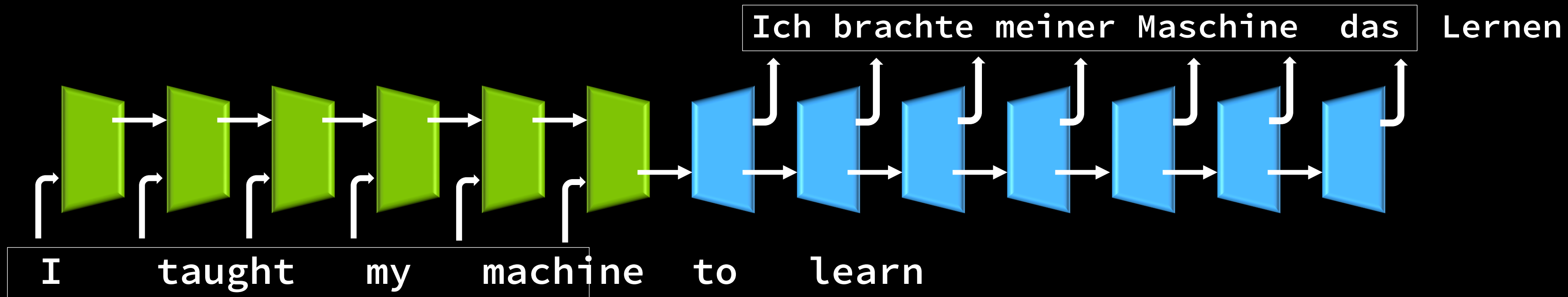
Jürgen Schmidhuber



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

SEQUENCE TO SEQUENCE

(2014) Encoder-decoder pattern for sequences



Sequence to Sequence Learning with Neural Networks

Ilya Sutskever
Google

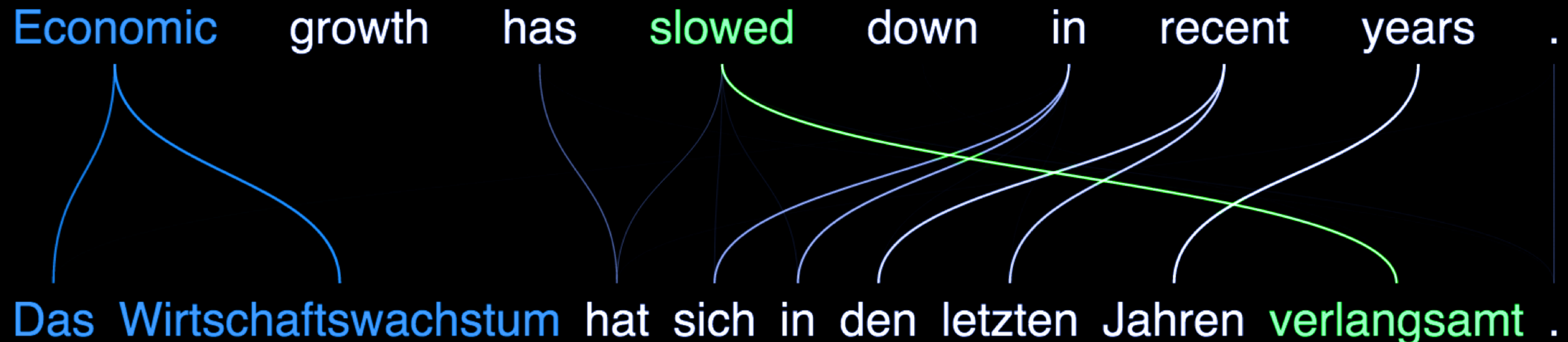
Oriol Vinyals
Google

Quoc V. Le
Google

<https://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf>

ATTENTION

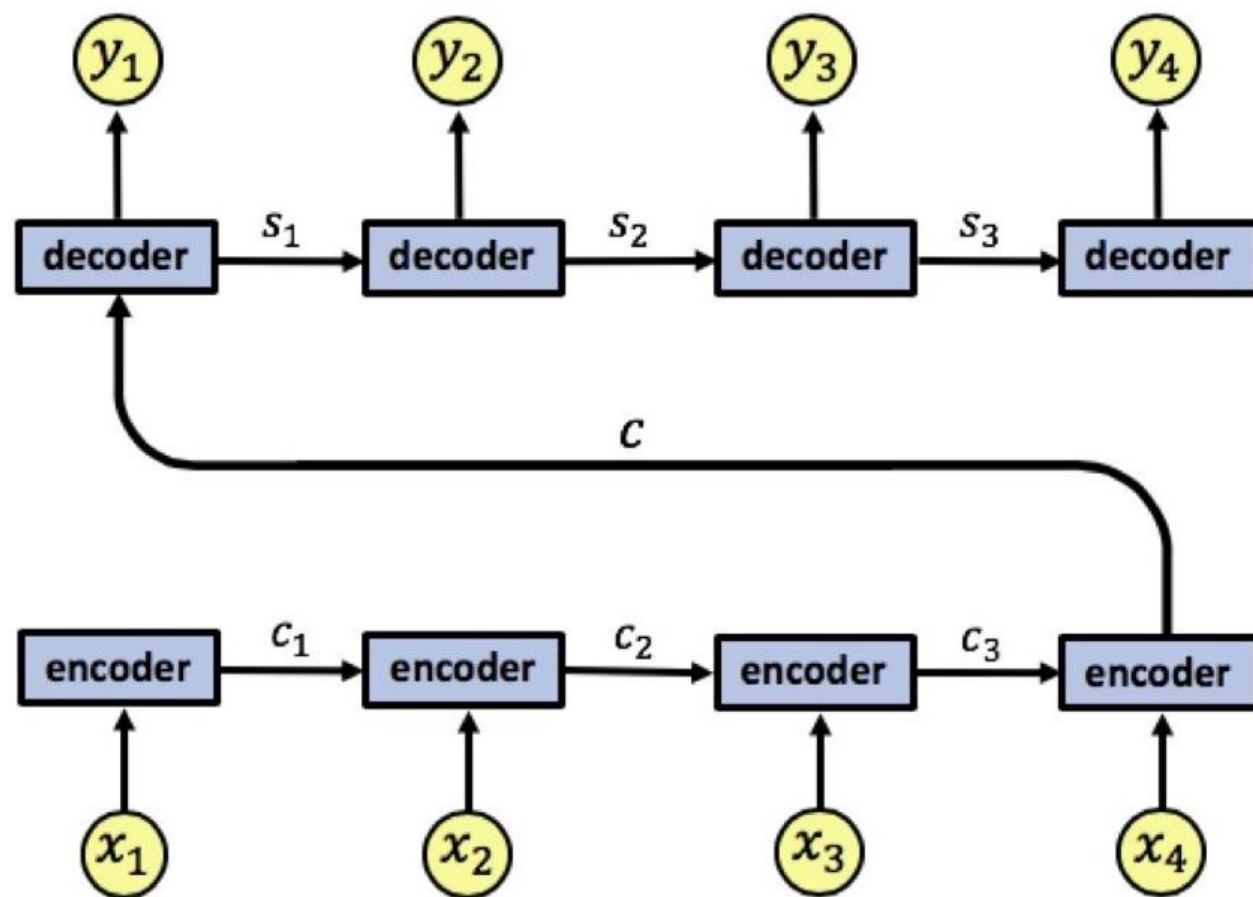
Adjustable weights based on context



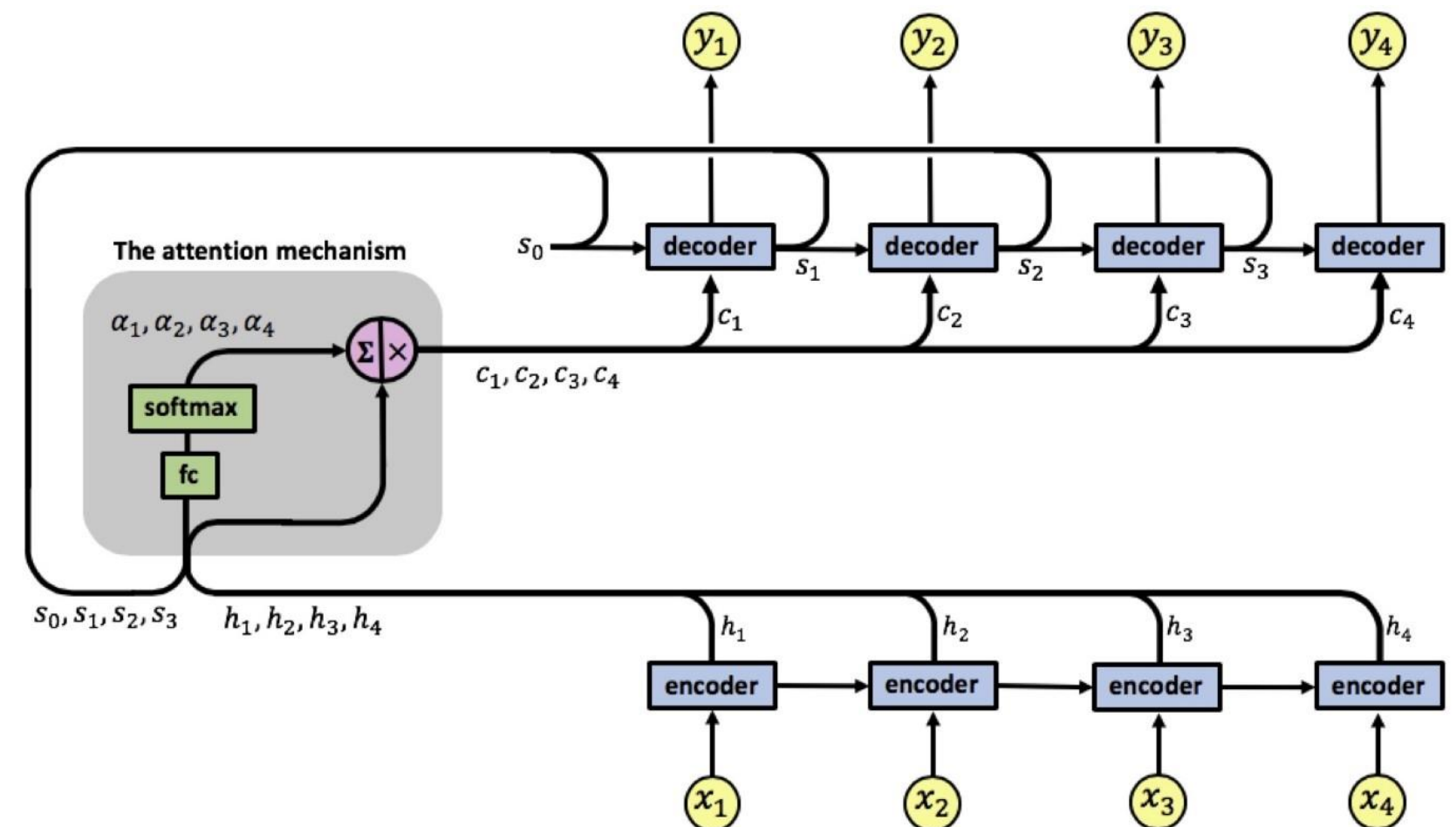
ATTENTION ENCODER-DECODER

A more accurate way to translate long sequences

RNN



RNN + ATTENTION



THE TRANSFORMER

(2017) Attention is all you need!

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

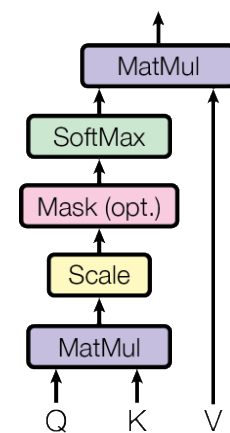
Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Scaled Dot-Product Attention



Multi-Head Attention

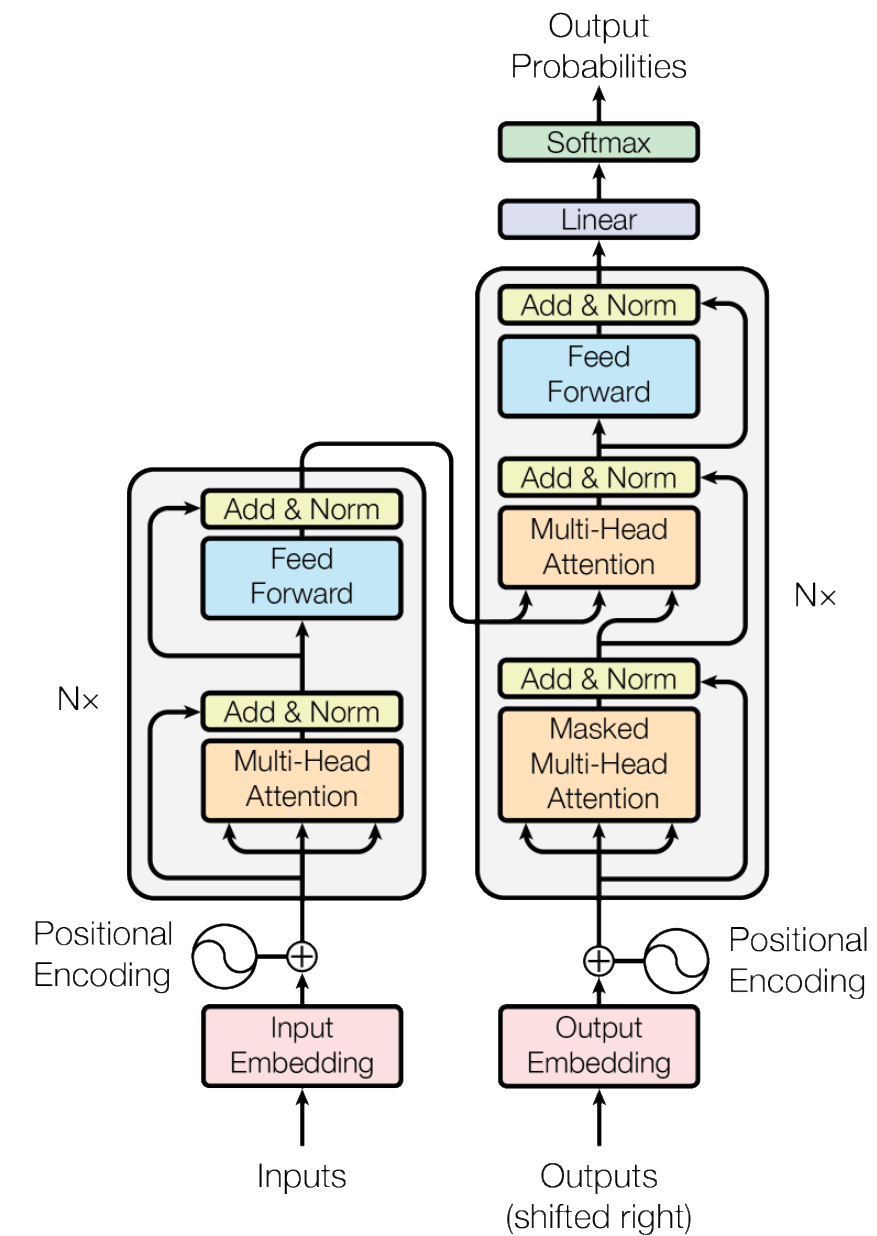
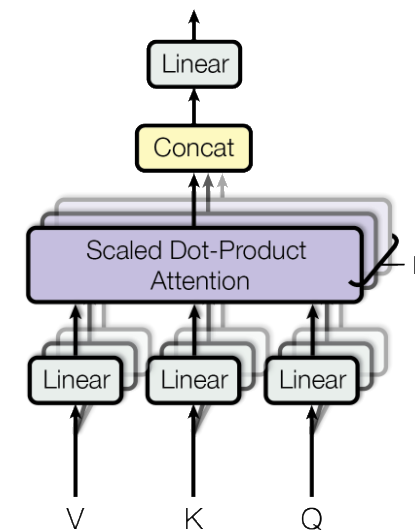
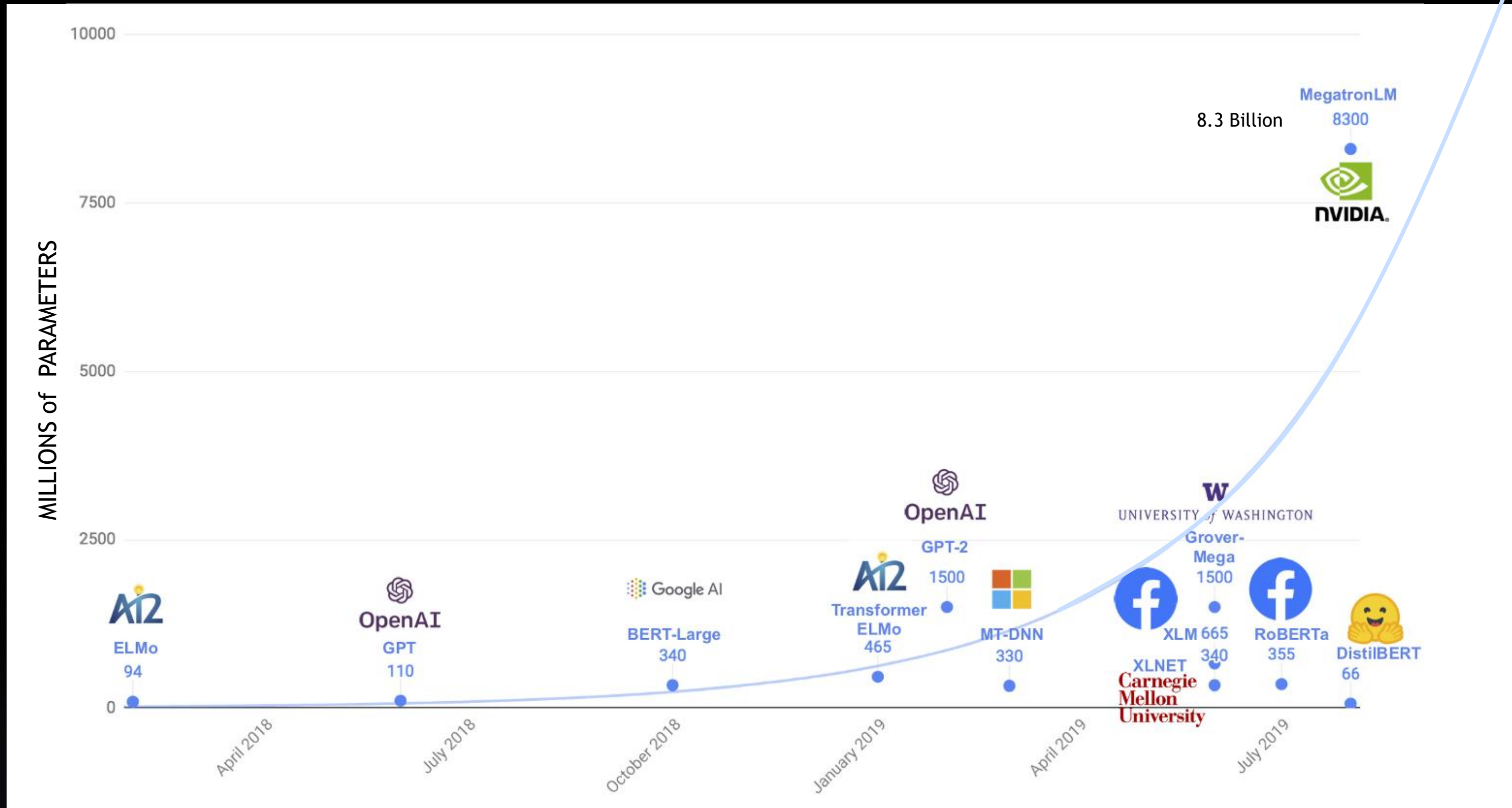


Figure 1: The Transformer - model architecture.

GPT AND BERT

The Transformer has enabled a series of massive language models



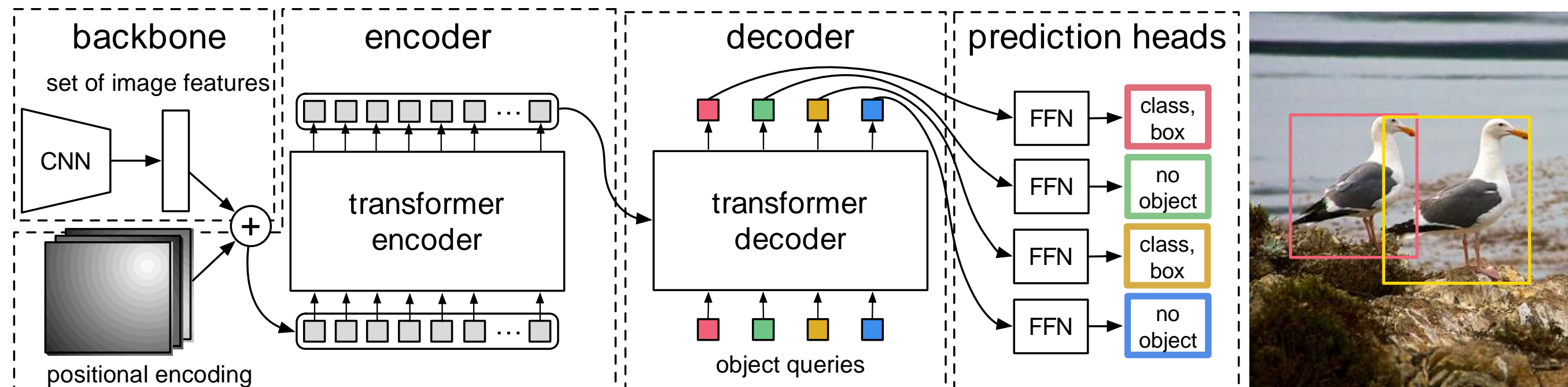
DETNET

(2020) End-to-end object detection with transformers

End-to-End Object Detection with Transformers

Nicolas Carion*, Francisco Massa*, Gabriel Synnaeve, Nicolas Usunier,
Alexander Kirillov, and Sergey Zagoruyko

Facebook AI

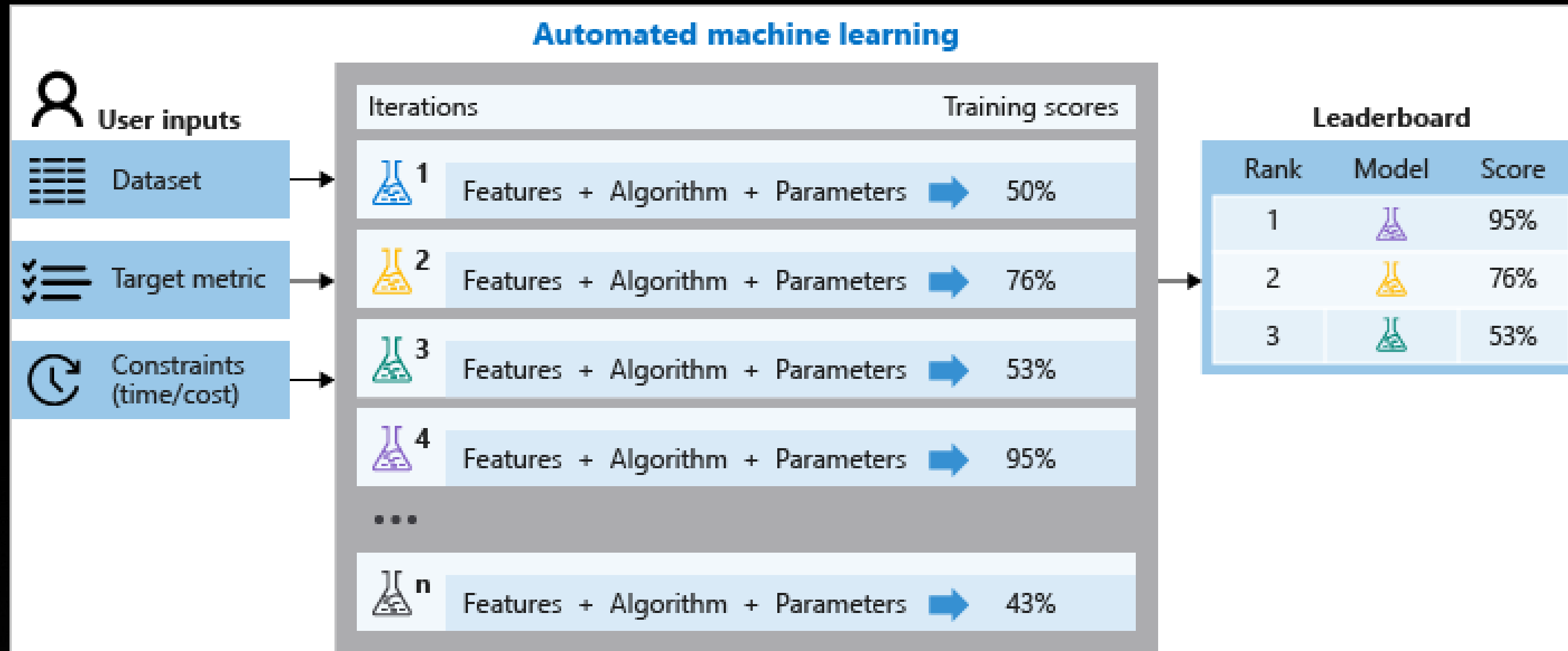




AUTO ML

AUTO ML

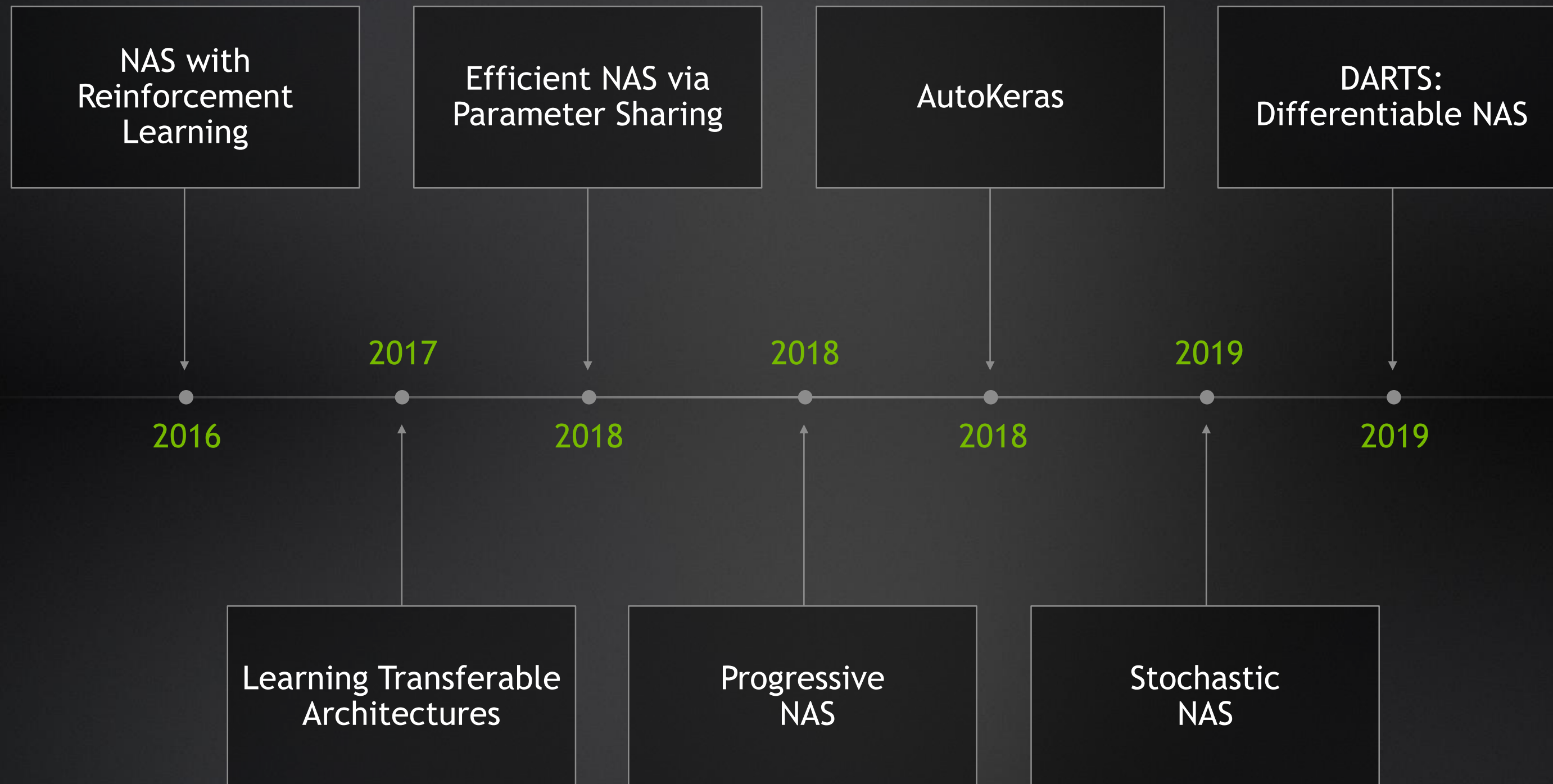
Automating the work of the data scientist



<https://docs.microsoft.com/en-us/azure/machine-learning/concept-automated-ml>

NEURAL ARCHITECTURE SEARCH

Automating model design and selection

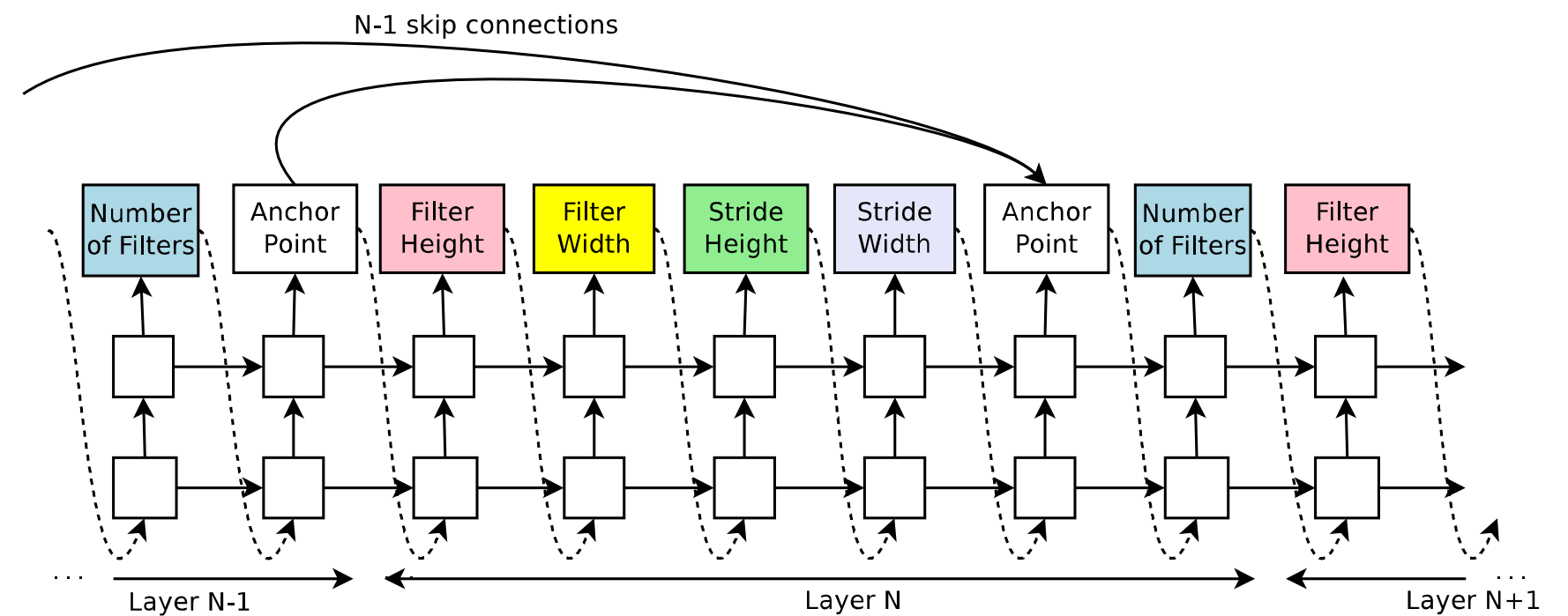
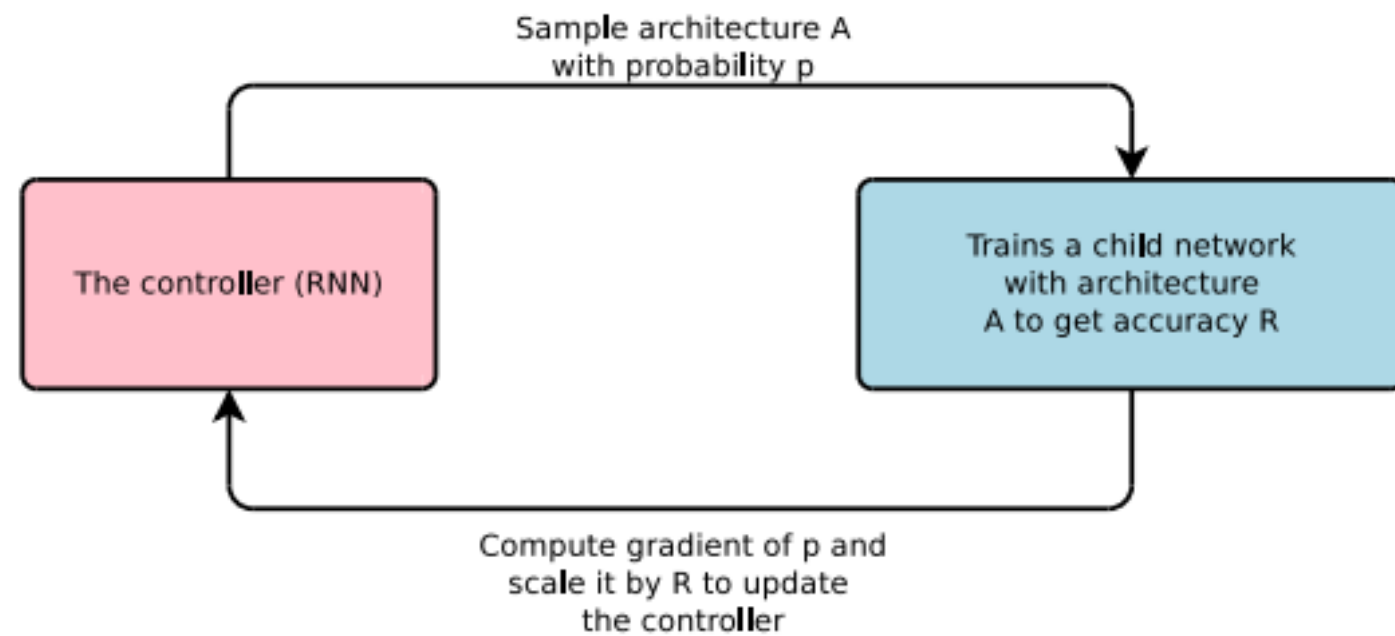


NAS

Google Brain 2016

NEURAL ARCHITECTURE SEARCH WITH REINFORCEMENT LEARNING

Barret Zoph*, Quoc V. Le
Google Brain



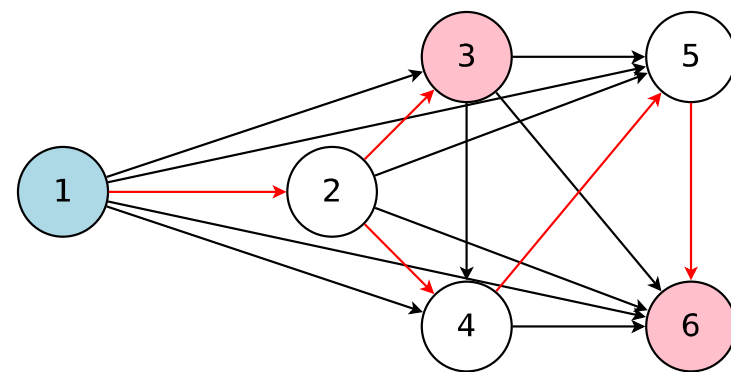
EFFICIENT NAS

2018 Google Brain, CMU, Stanford

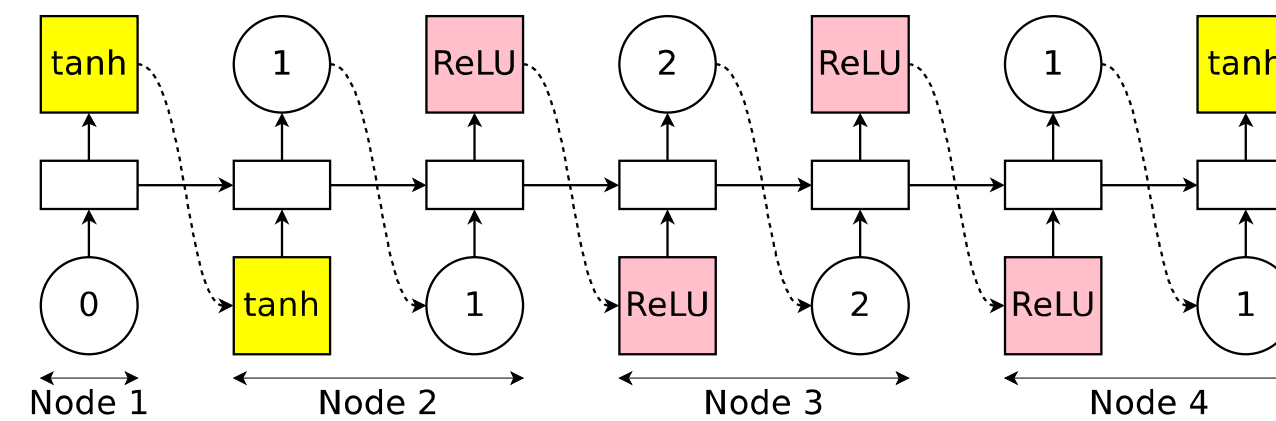
Efficient Neural Architecture Search via Parameter Sharing

Hieu Pham^{*1,2} Melody Y. Guan^{*3} Barret Zoph¹ Quoc V. Le¹ Jeff Dean¹

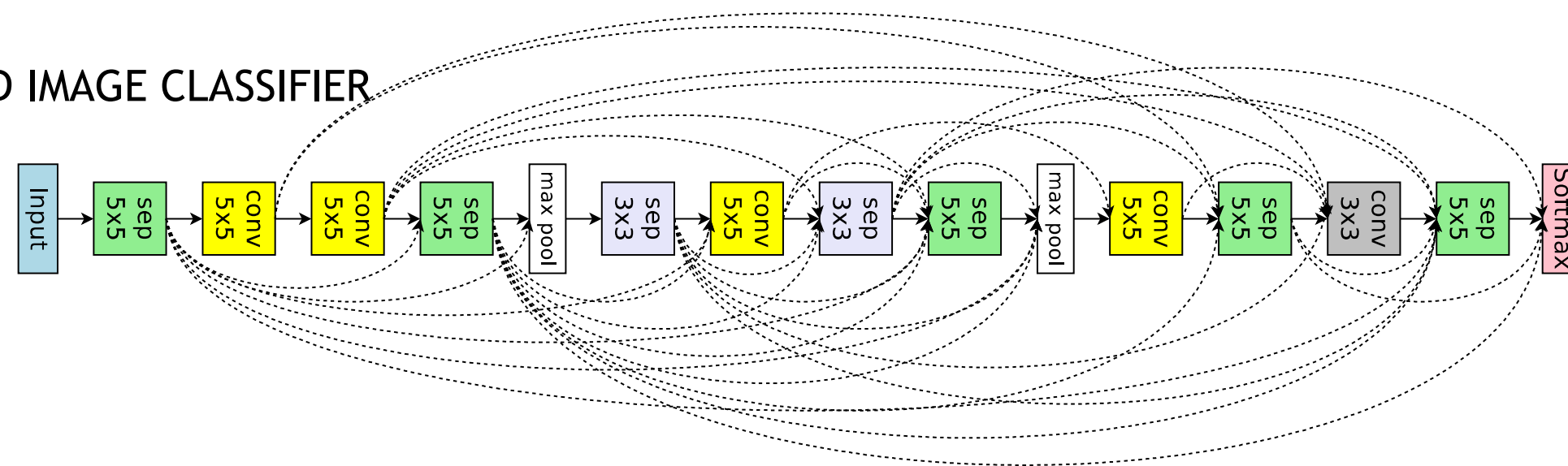
GRAPH SEARCH SPACE



RNN CONTROLLER



DISCOVERED IMAGE CLASSIFIER



DARTS

2019 Google DeepMind, CMU

DARTS: DIFFERENTIABLE ARCHITECTURE SEARCH

Hanxiao Liu*

CMU

hanxiaol@cs.cmu.com

Karen Simonyan

DeepMind

simonyan@google.com

Yiming Yang

CMU

yiming@cs.cmu.edu

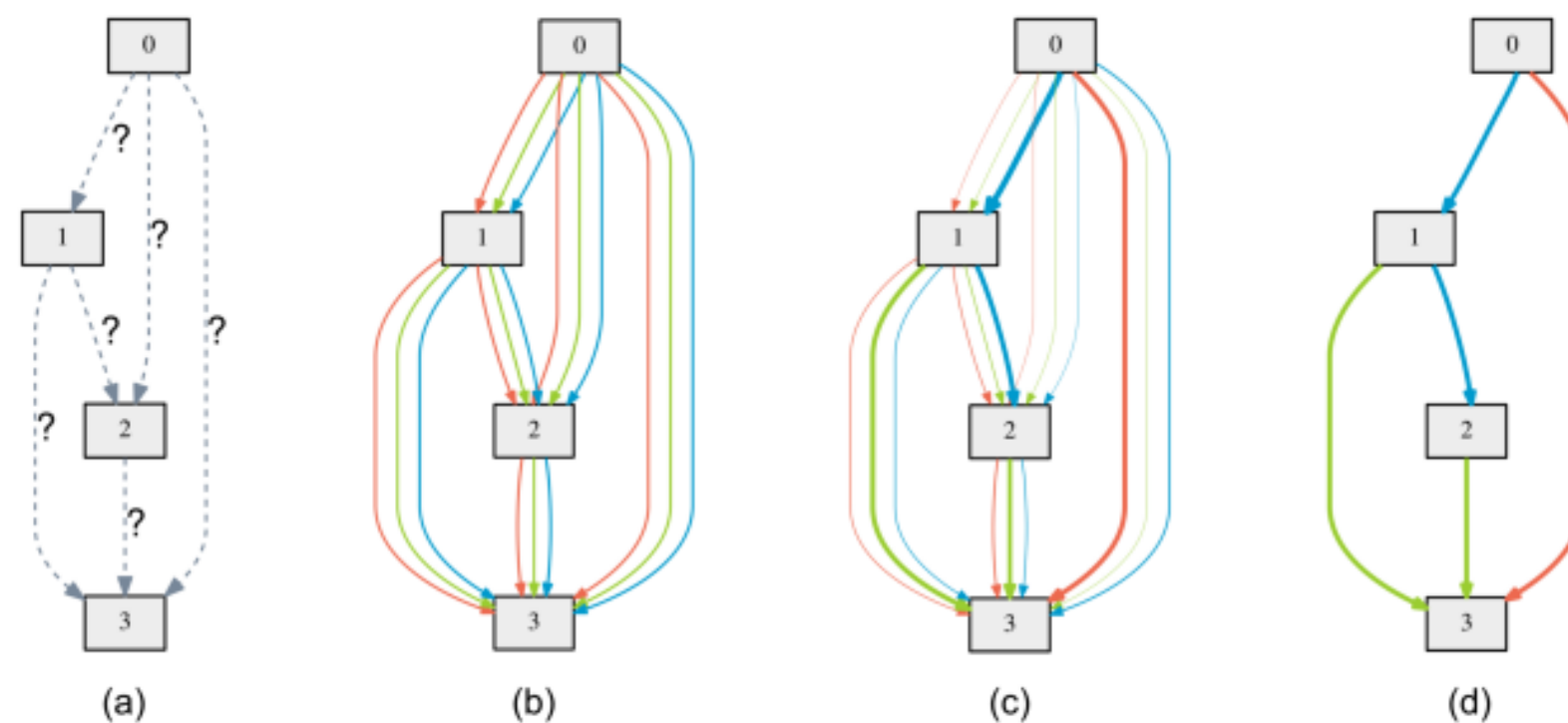


Figure 1: An overview of DARTS: (a) Operations on the edges are initially unknown. (b) Continuous relaxation of the search space by placing a mixture of candidate operations on each edge. (c) Joint optimization of the mixing probabilities and the network weights by solving a bilevel optimization problem. (d) Inducing the final architecture from the learned mixing probabilities.



SUMMARY

- Eventually AutoML and NAS may make model selection obsolete
- Many models in the model zoo
- Model architectures are task specific
- Encode what you know, learn the rest
- Use pretrained models when you can
- Use transfer learning when you can
- AI changes very quickly. Use Arxiv sanity to keep up.

dhall@nvidia.com
