

Evolving of Deep Neural Networks: Algorithm and Applications



Ming Yang
Co-founder & VP, Software



**Horizon
Robotics**

Define the brain of things



Co-founder &
VP of Software



AI Research



Multimedia
Analytics



PhD.,
Dept. of EECS



BE. And ME.,
Dept. of EE

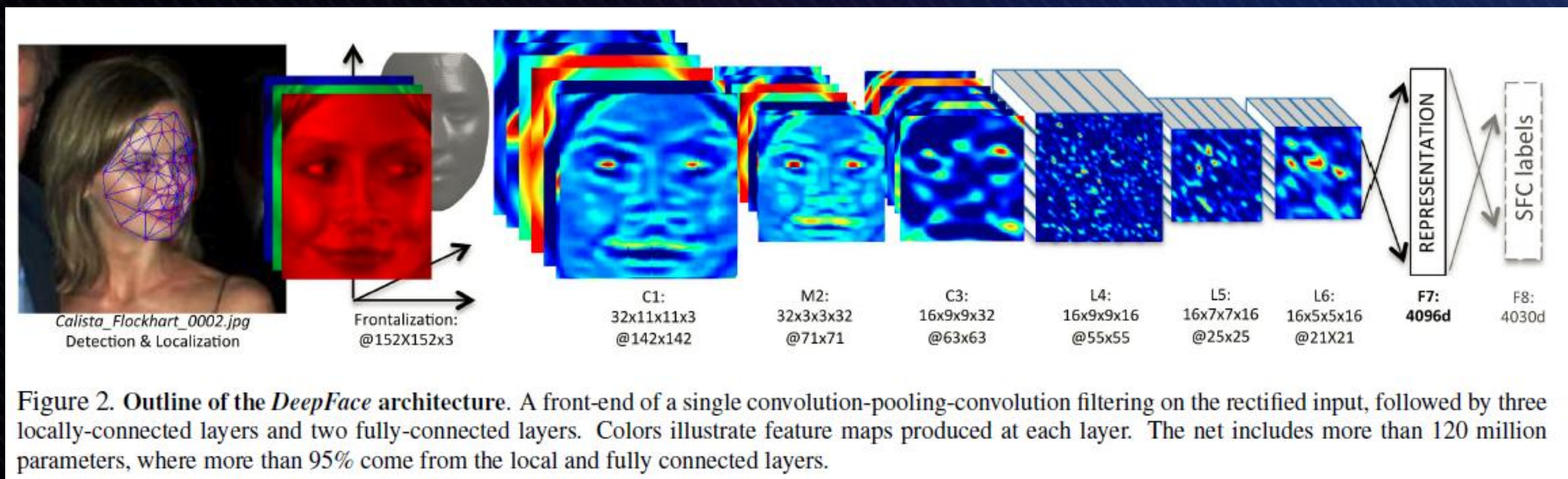
Dr. Ming Yang is one of the founding member of the Facebook Artificial Intelligence Research (FAIR) and a former senior researcher at NEC Labs America. Dr. Yang is a well-recognized researcher in computer vision and machine learning. He co-authored 14 US patents, and over 20 publications in top conferences like CVPR and ICCV and 8 publications in the top international journal T-PAMI with more than **4750** citations.

During his tenure at Facebook, Dr. Yang led the deep learning research project “**DeepFace**”, which had a significant impact in the deep learning research community and got widely reported by various media including Science Magazine, MIT Tech Review and Forbes.

He received his B.Eng. and M.Eng. Degree from the Dept. of Electrical Engineering at Tsinghua University and Ph.D. degree from the Dept. of Electrical Engineering and Computer Science at Northwestern University.



Science Magazine, Jan. 2015

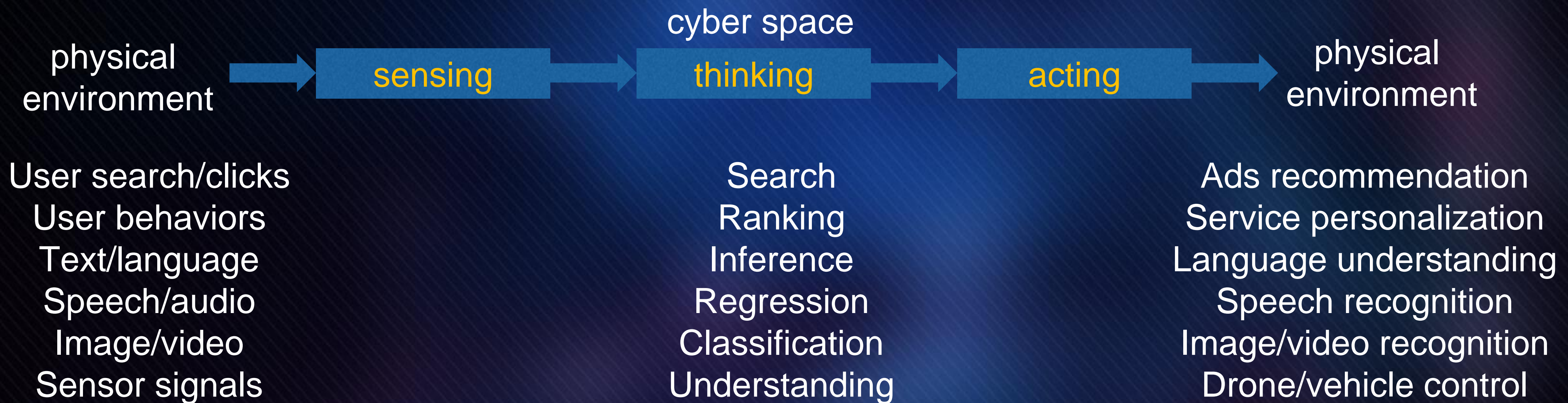


DeepFace: closing the gap to human level performance in face verification, Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, Lior Wolf, IEEE Conference on Computer Vision and Pattern Recognition, CVPR, June 2014.

The most cited paper (1200+ citations) on face recognition using deep neural networks.



Artificial intelligence (AI) is defined as “the study and design of **intelligent agents**, in which an intelligent agent perceives its environment and takes actions that maximize its chance of success.





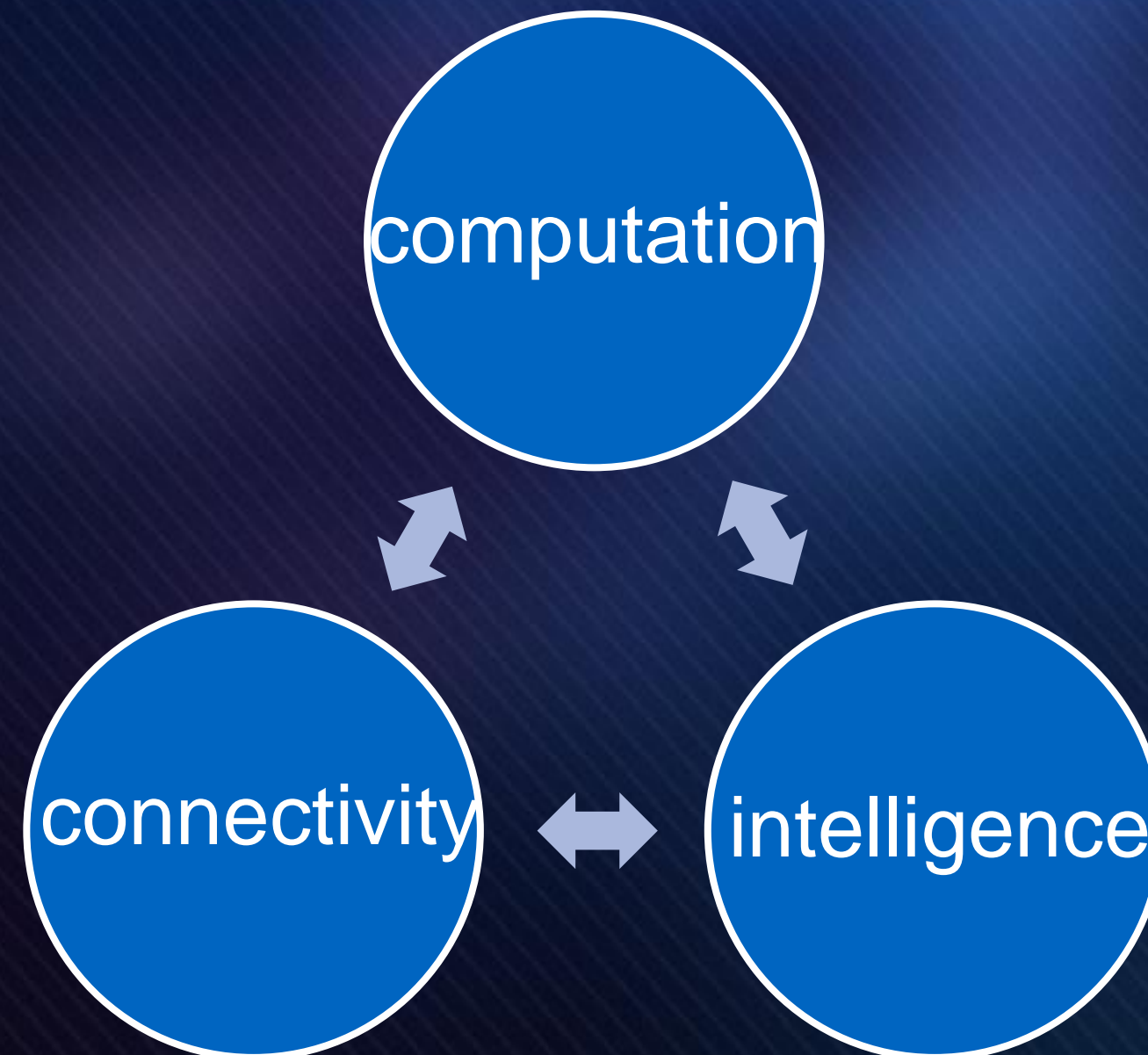
PC + Connection



Mobile computing



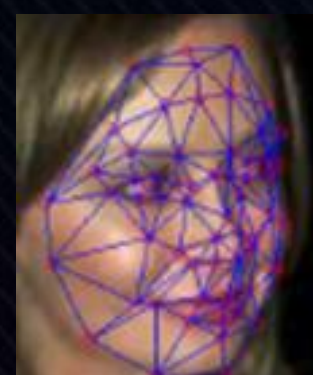
IoT + intelligence



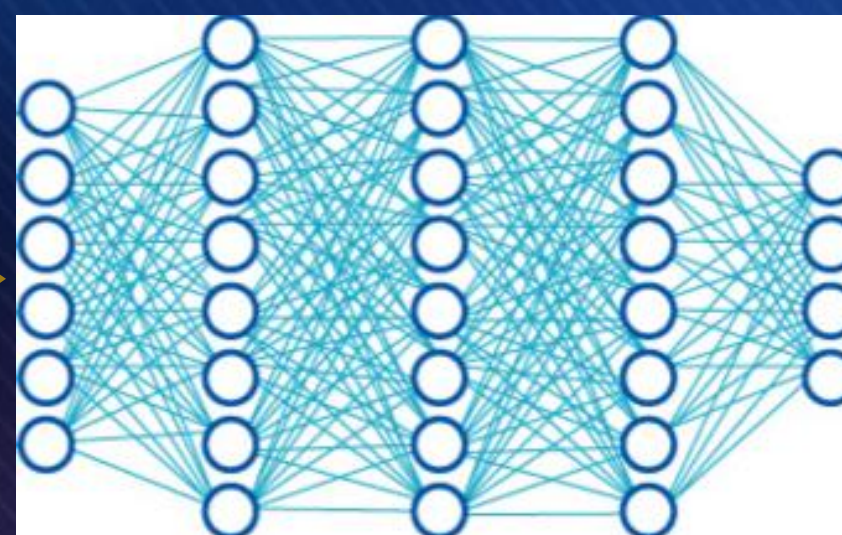
raw data



representations



Joe went to the kitchen. Fred went to the kitchen. Joe picked up the milk.
Joe travelled to the office. Joe left the milk. Joe went to the bathroom.
Where is the milk now? A: office
Where is Joe? A: bathroom
Where was Joe before the office? A: kitchen



end-to-end learning



ANN/CNN/RNN/LSTM < Neural networks < Deep learning < AI

GAN, 2016

RL, AlphaGo, 2016.3.9



Autonomous driving, 2015

RNN/LSTM 2014-2015

Face recognition, 2014

Facebook AI Research, 2013

ImageNet+GPU

A. Krizhevsky, etc., 2012

Google brain / DistBelief

J. Dean, A. Ng, etc., 2011

DNN for speech recognition

G. Hinton, etc., 2010

Layer-by-layer training

G. Hinton, etc., 2006

Convolutional neural networks, 1989-1998

Back propagation (BP)

P. Werbos, Y. LeCun, etc., 1980s

Artificial neural networks (ANN)

K. Fukushima, etc., 1980

Perceptrons vs. XOR

M. Minsky, 1969

Perceptron, artificial neuron

F. Rosenblatt, 1957

On every smart device / thing!

1950

2000

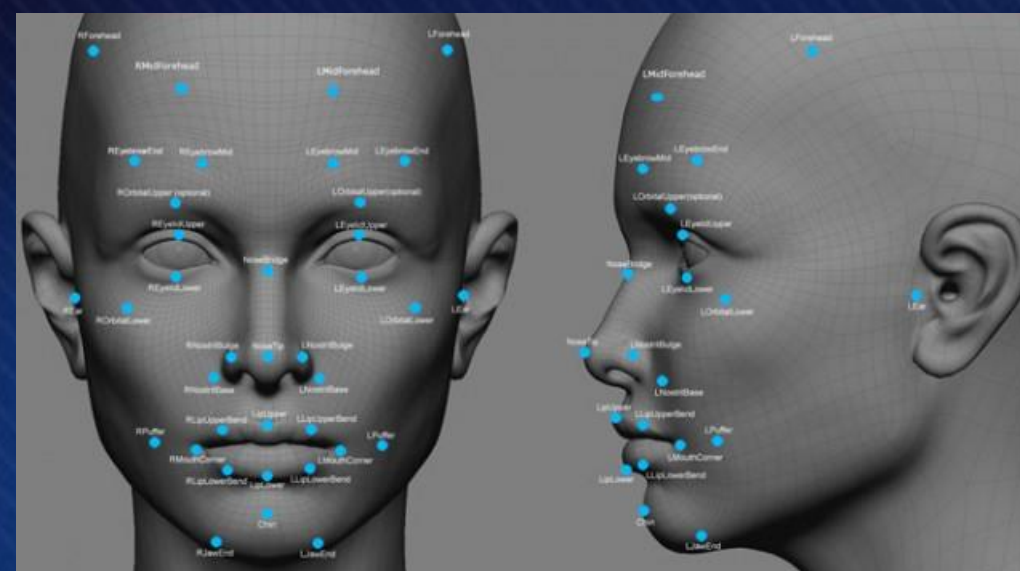
2017

2020

Speech



Vision



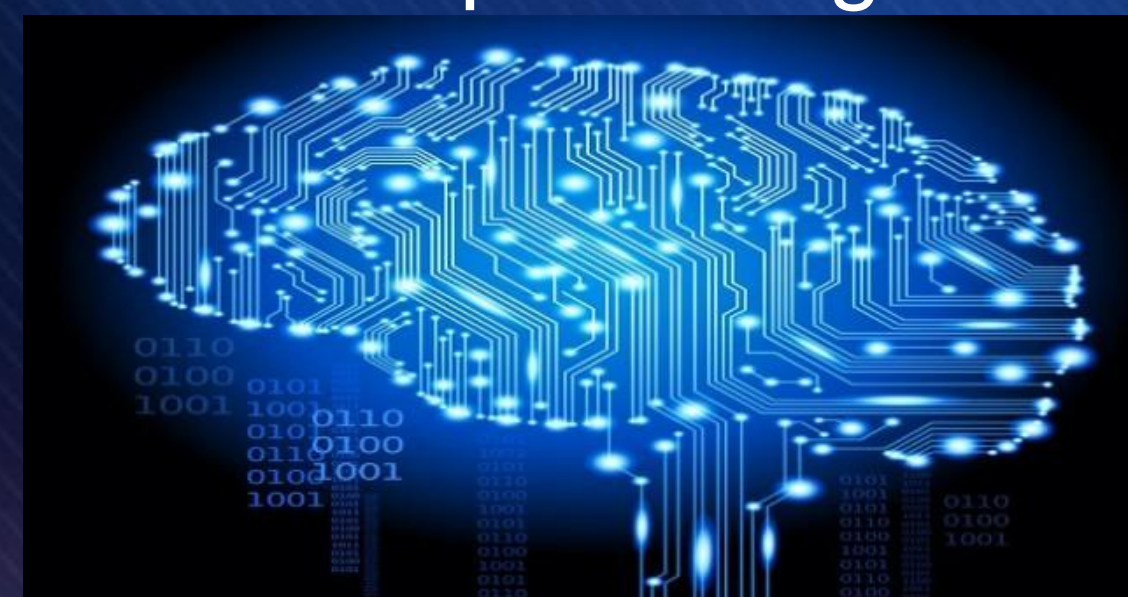
Automotive



Ads/Search



Deep learning



Hardware



Fintech



Industry

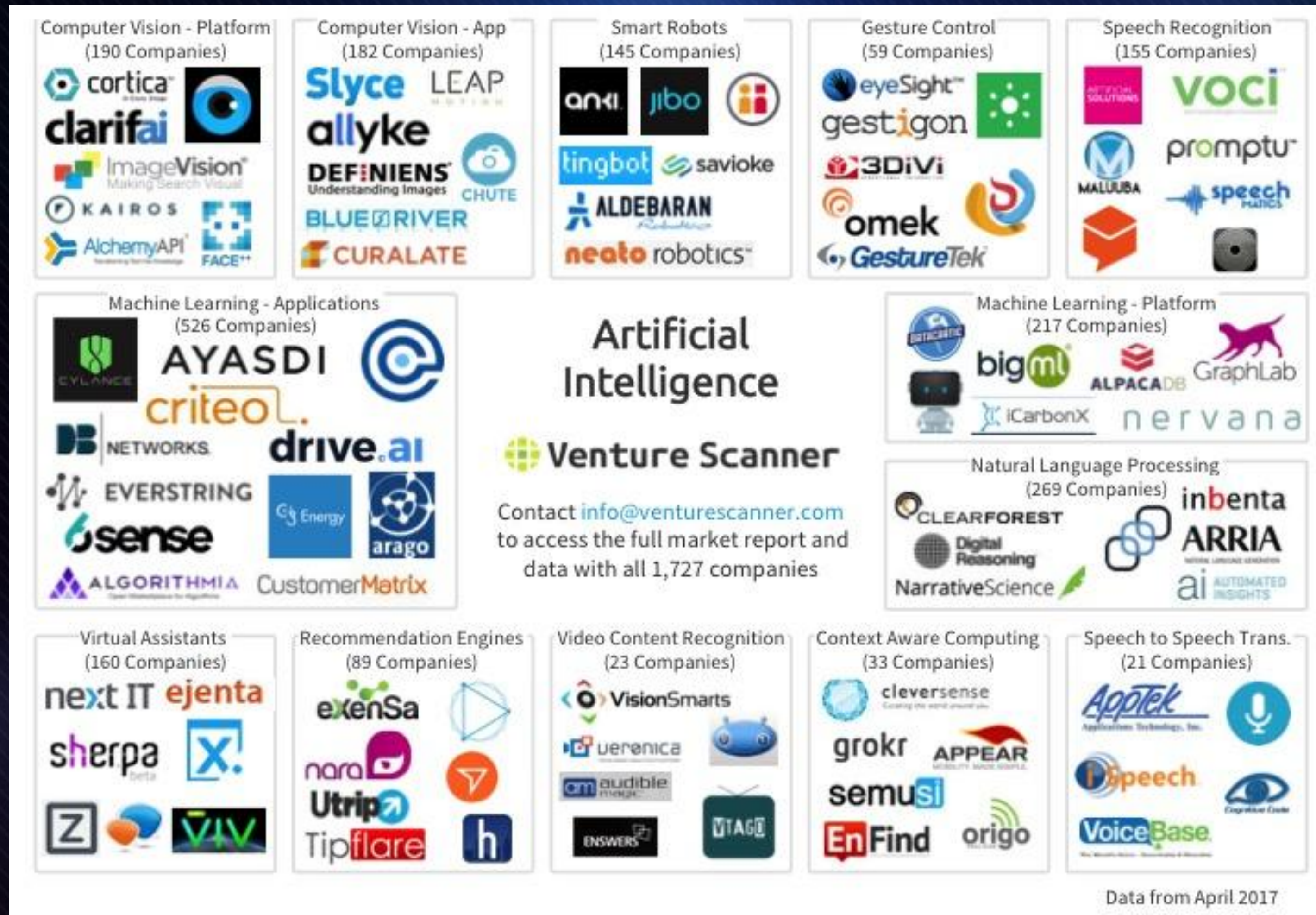


Healthcare



A thousand of AI startups (incomplete)

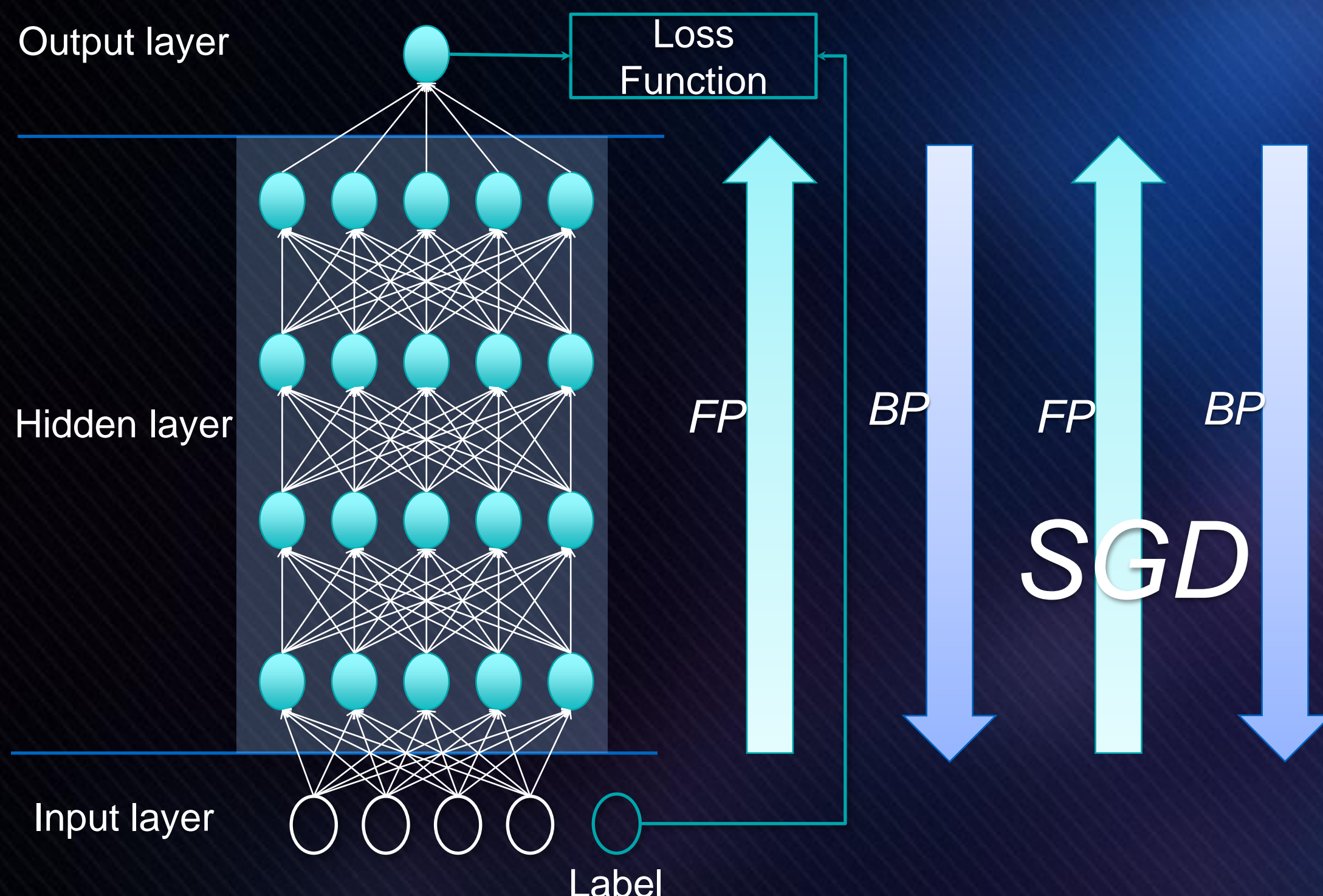
1,727 AI startups across 13 categories, with a combined funding amount of 14.5B\$ (by VentureScanner, 2017/4)



957 startups with 4.8B\$ investment, in 2016/3

(Deep) Multi-layer neural networks

What's new?



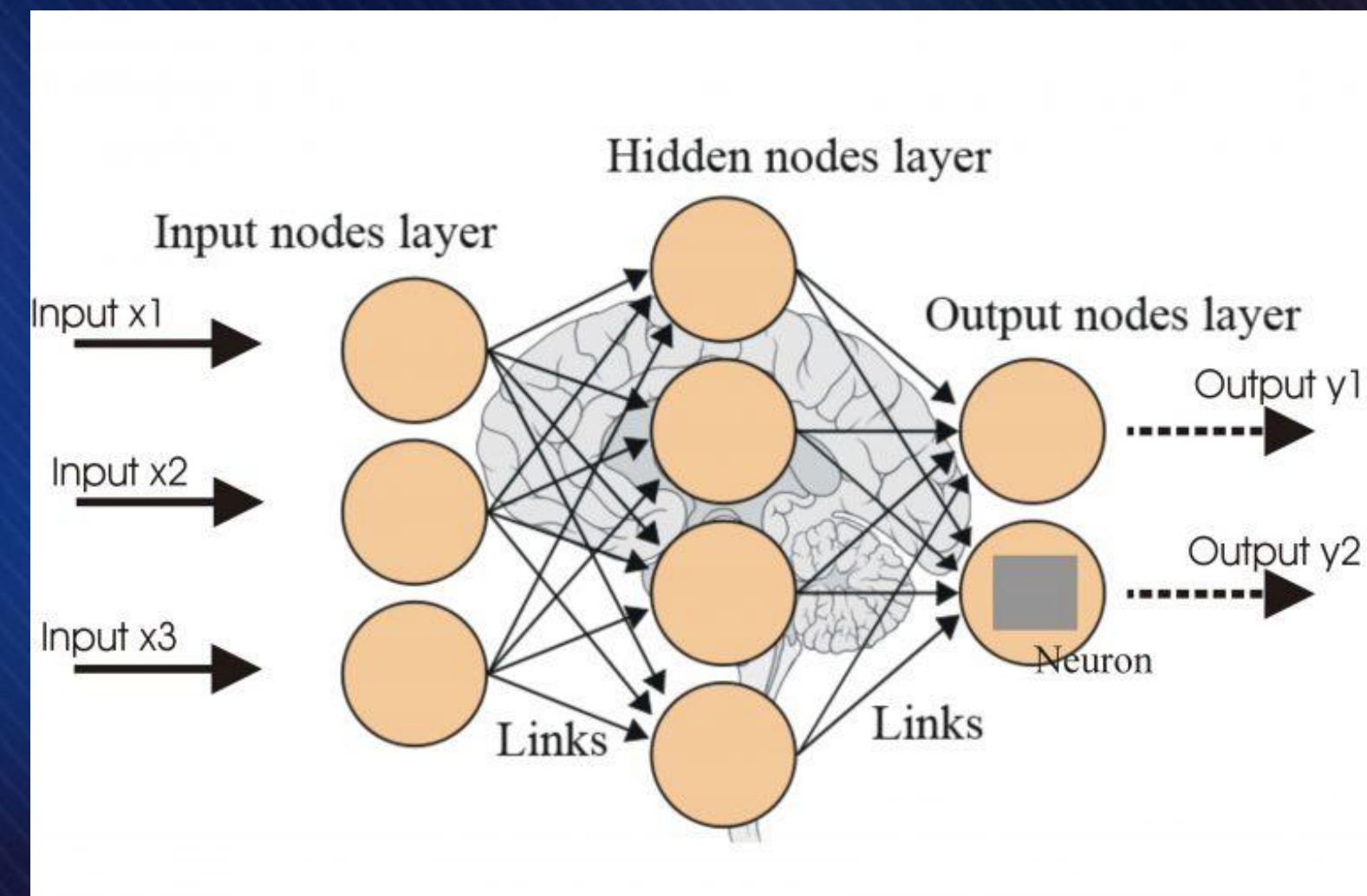
Challenges

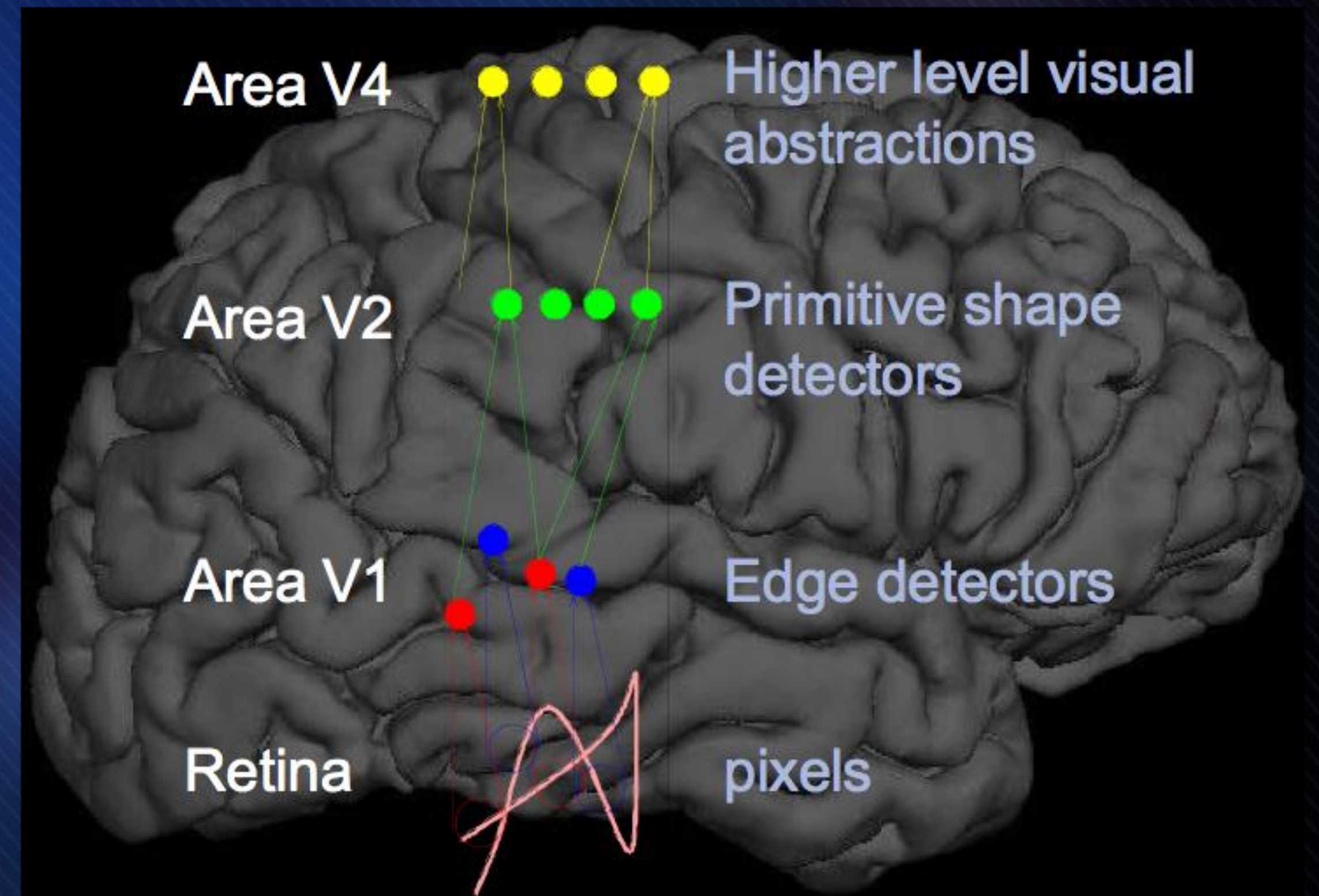
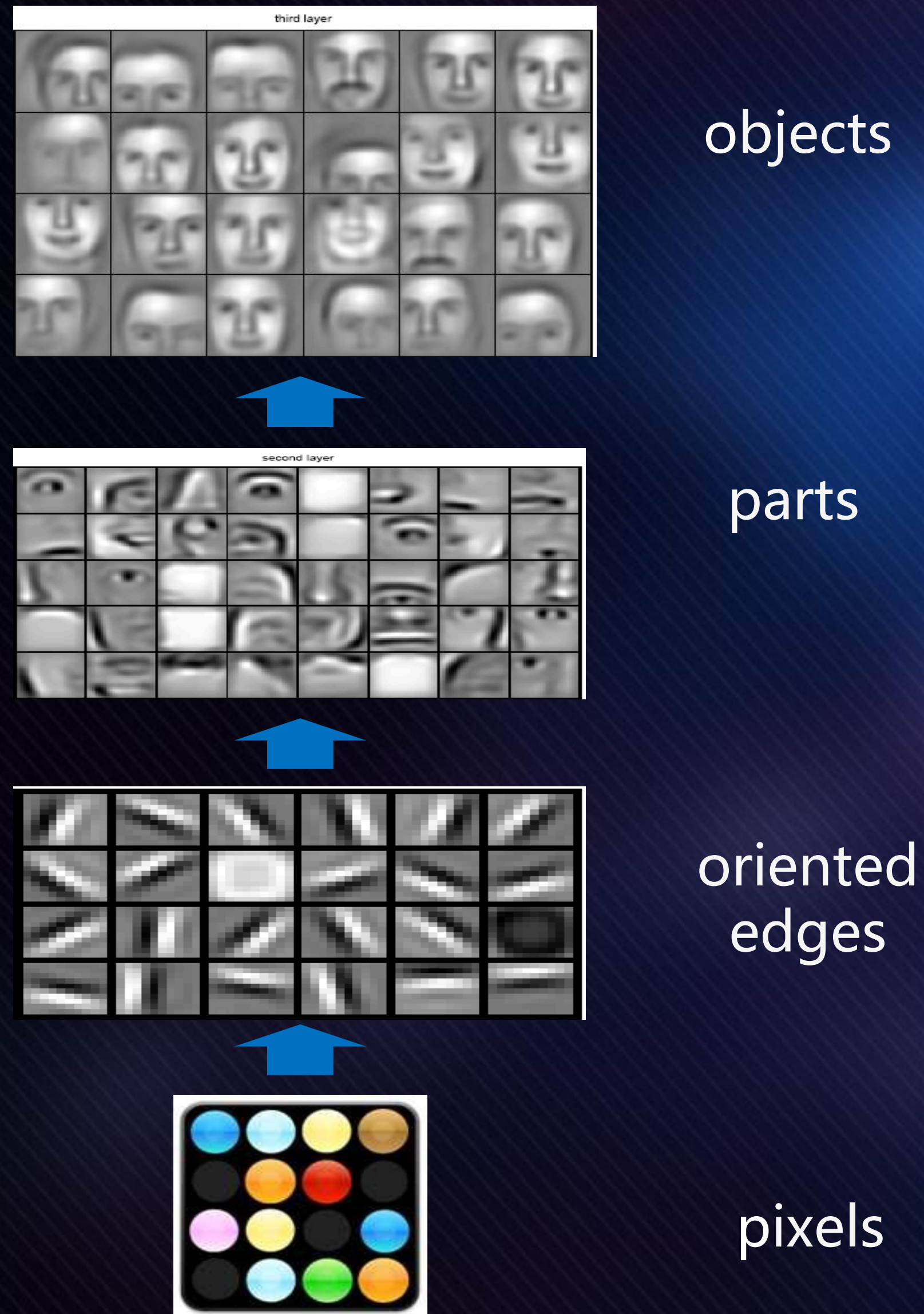
- ✓ over-fitting
 - ✓ noise sensitive
 - ✓ a black box
 - ✓ black magic
- an art to tune the network architecture

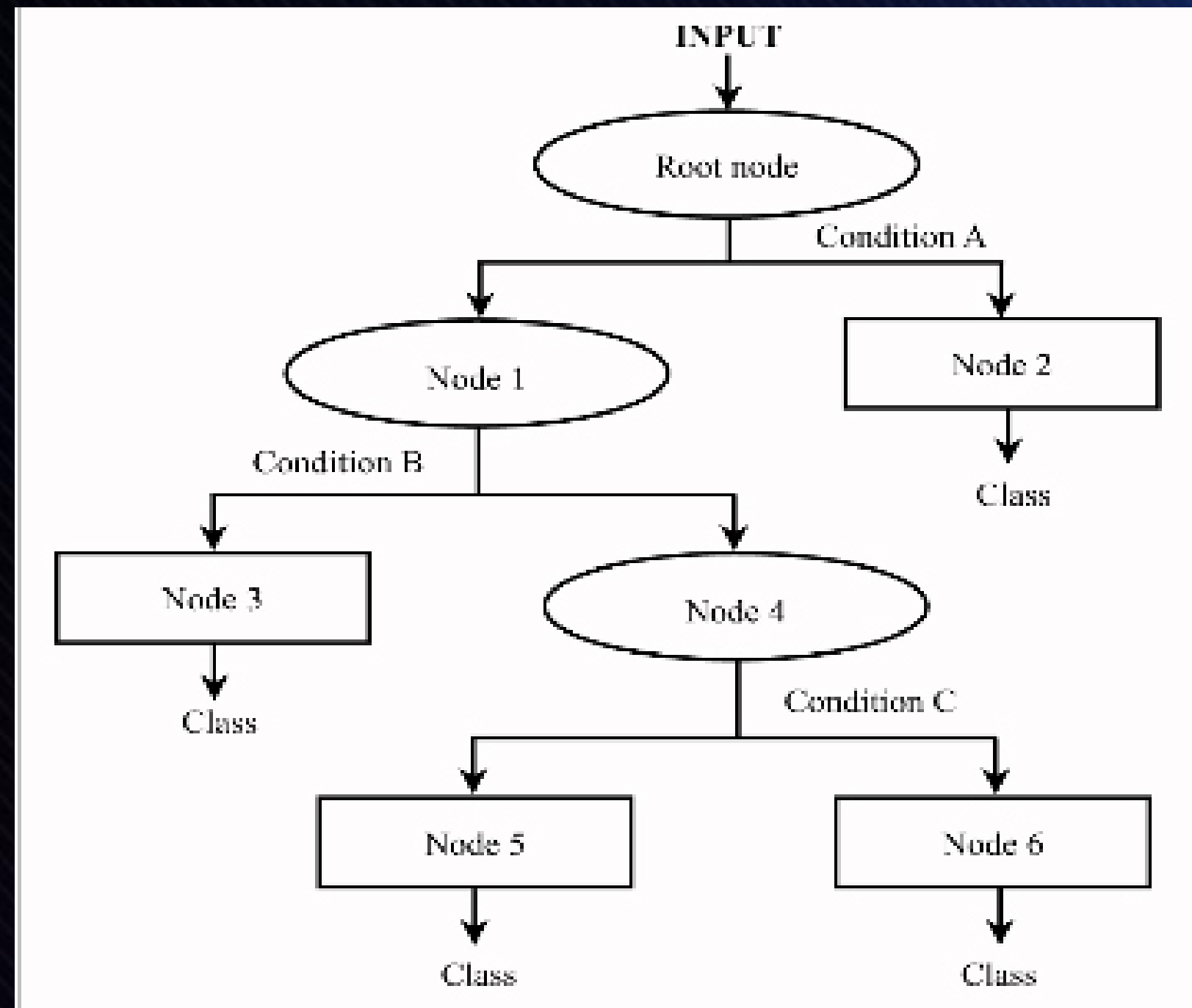
Now

- ✓ big big big data
 - ✓ GPU + parallelism
 - ✓ training tricks
 - ✓ **black magic**
- an art to tune the network architecture**

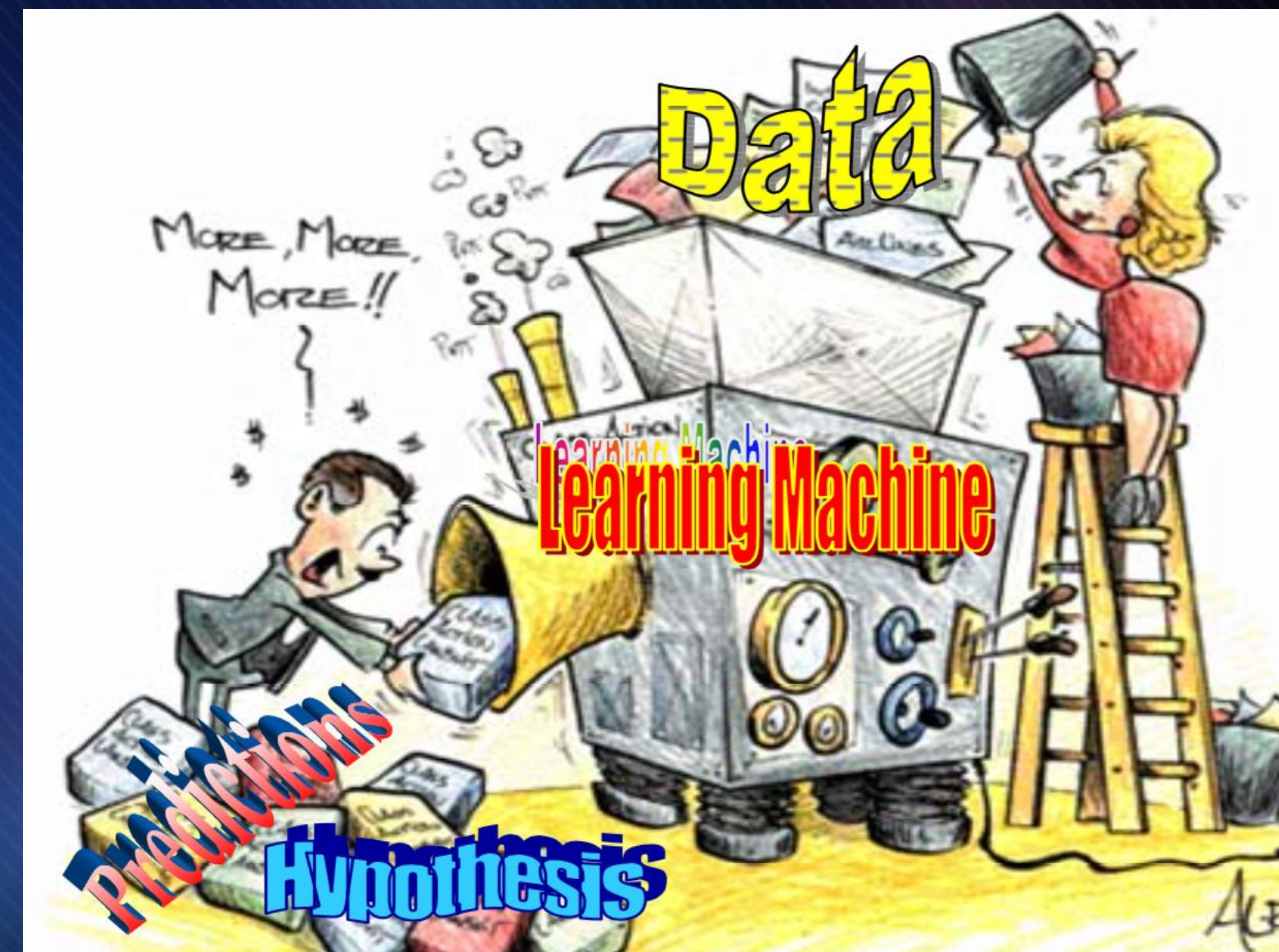
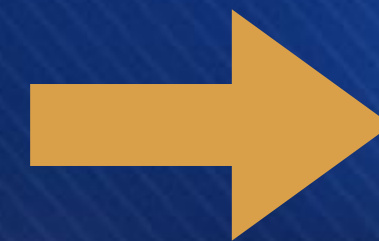
- Biological Inspired
- Big big big data
- End-to-end learning
- A flexible modeling language





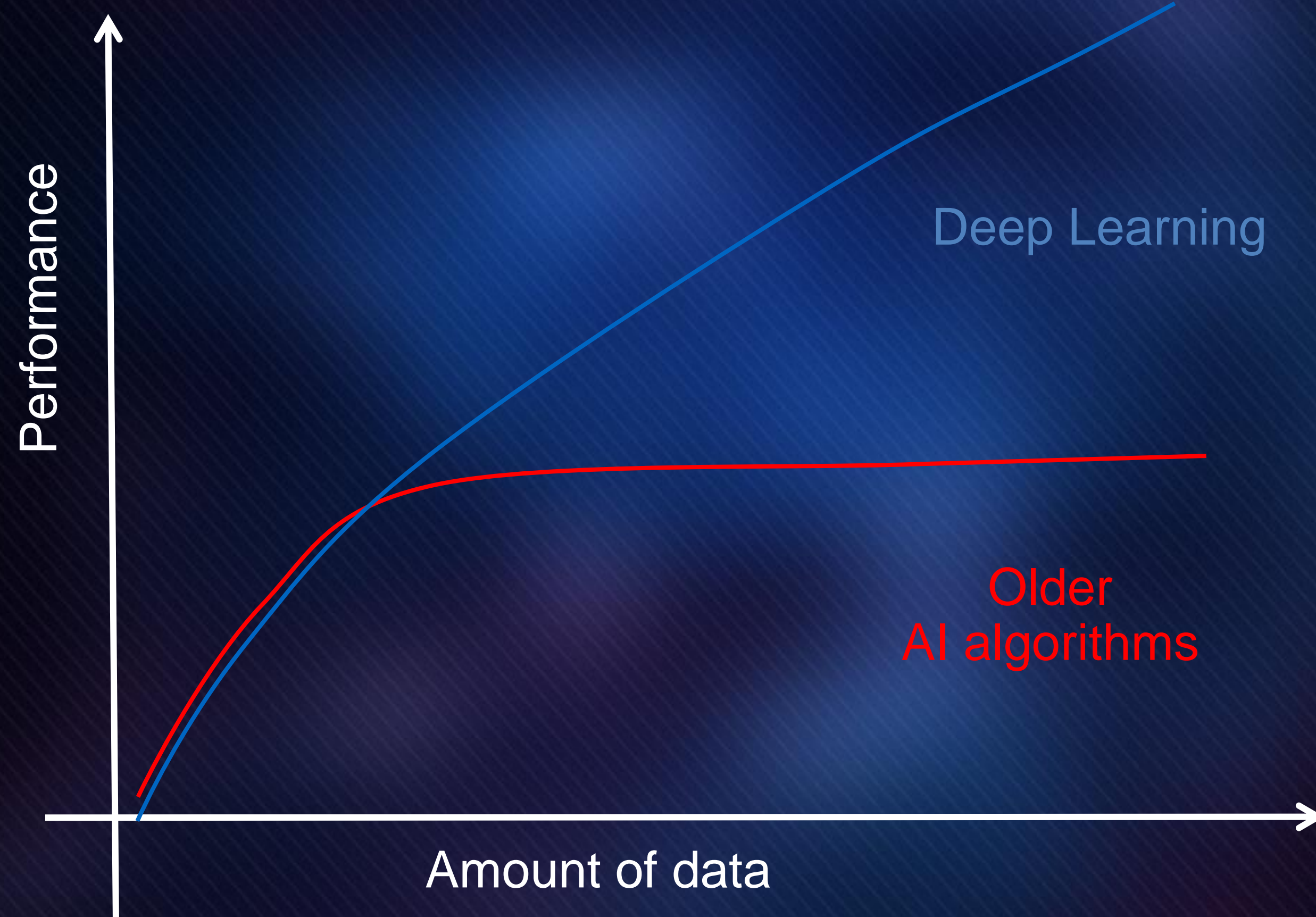


Rule-based AI

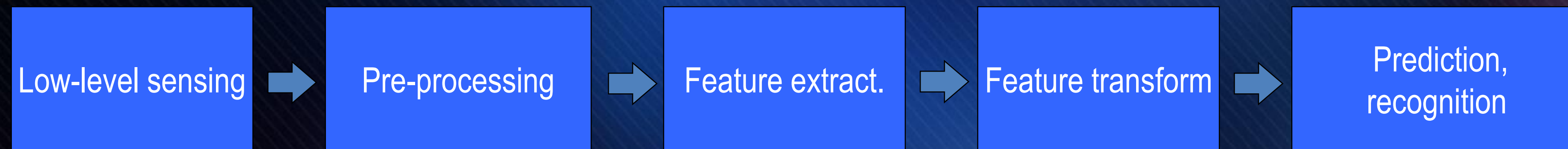


Data-driven AI

Big data, the more the better?

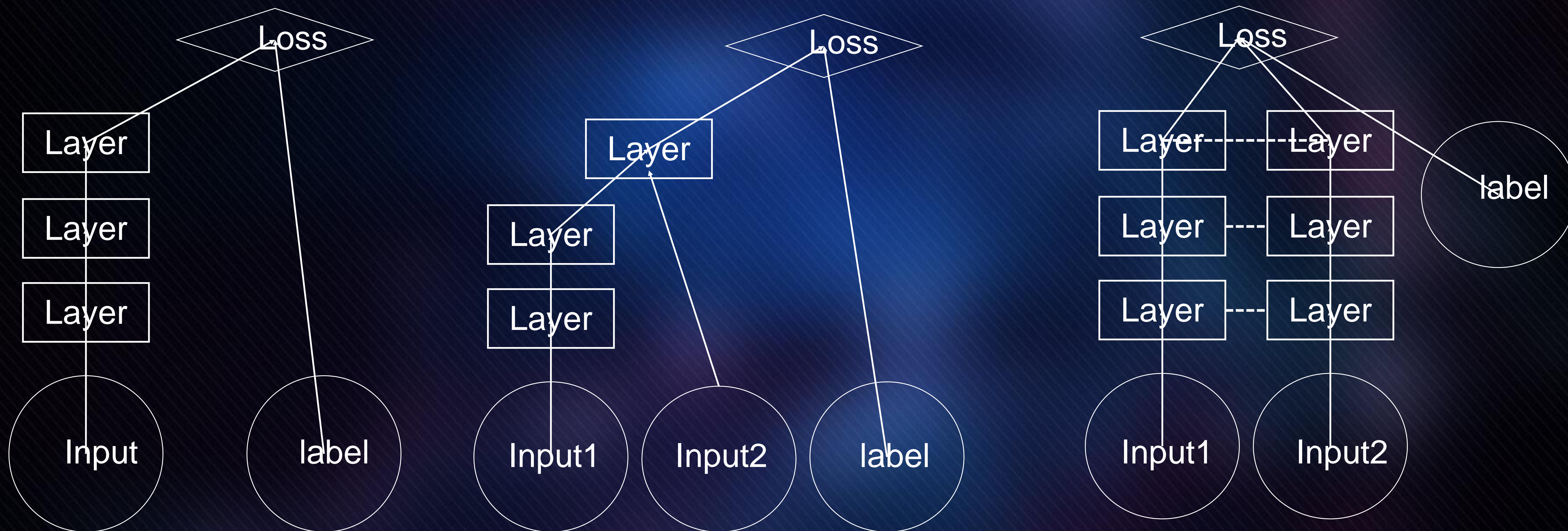


Most Efforts in Machine Learning



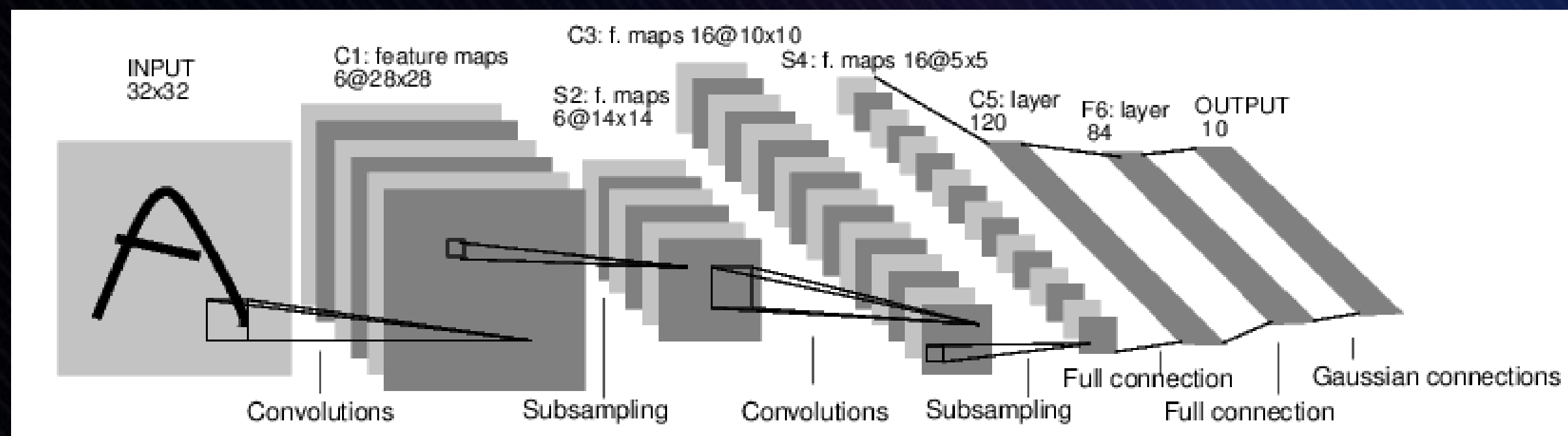
- Most critical for accuracy
- Account for most of the computation for testing
- Most time-consuming in development cycle
- Often hand-craft in practice

Deep Learning

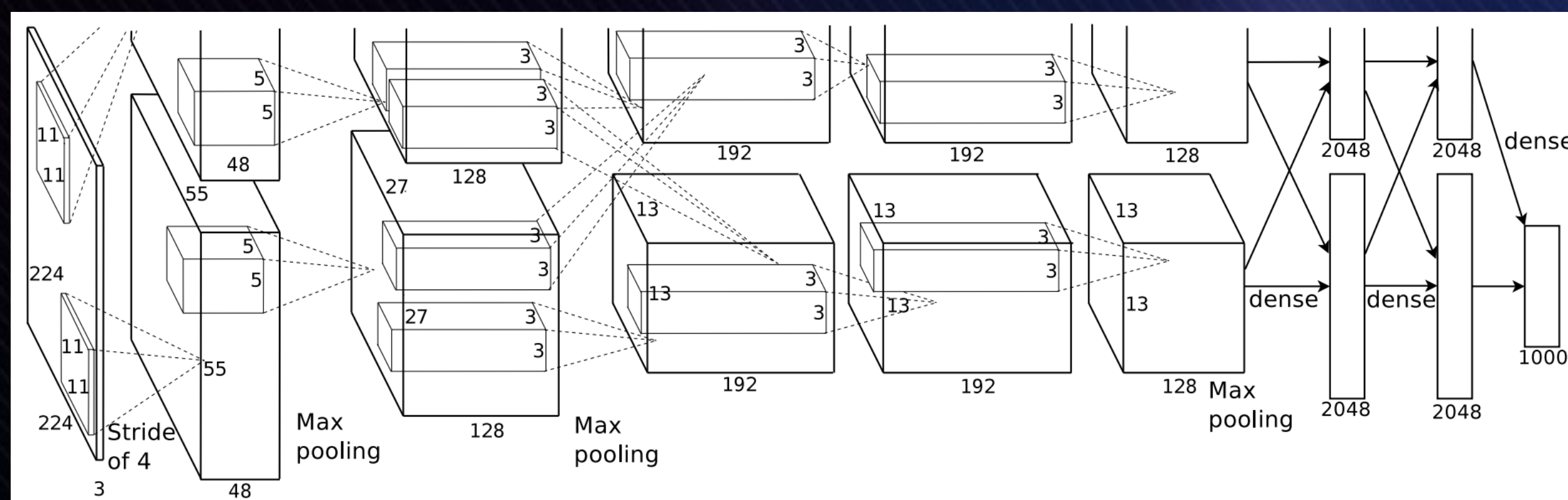


The network architecture is very flexible, so long as it is differentiable and learnable!

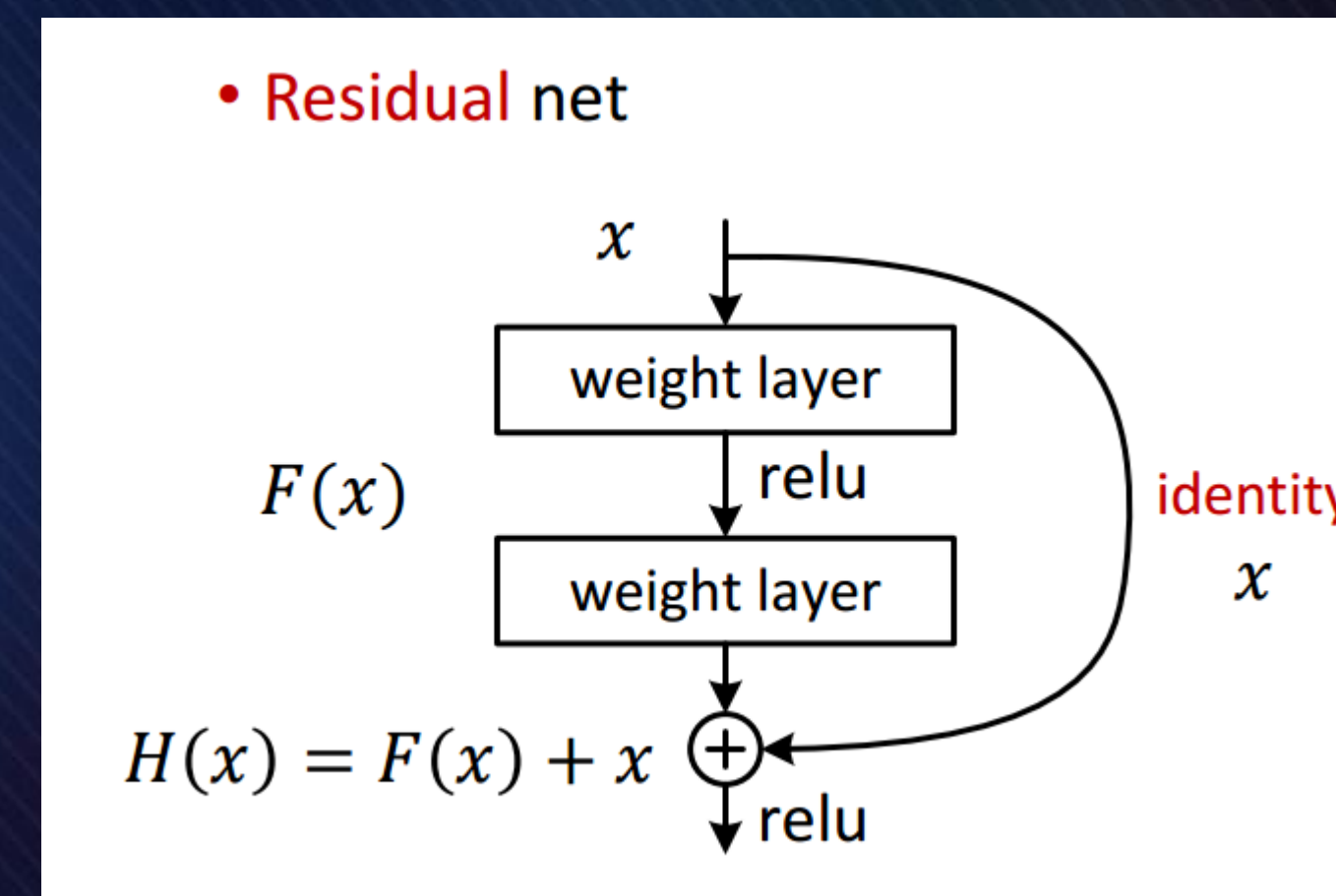
- Network **performance, size, operations**
- Network architecture : the deeper the better? spatially or temporally
- Smaller storage: model compression, or with low precision operations
- Faster computation : less operations, or hardware acceleration



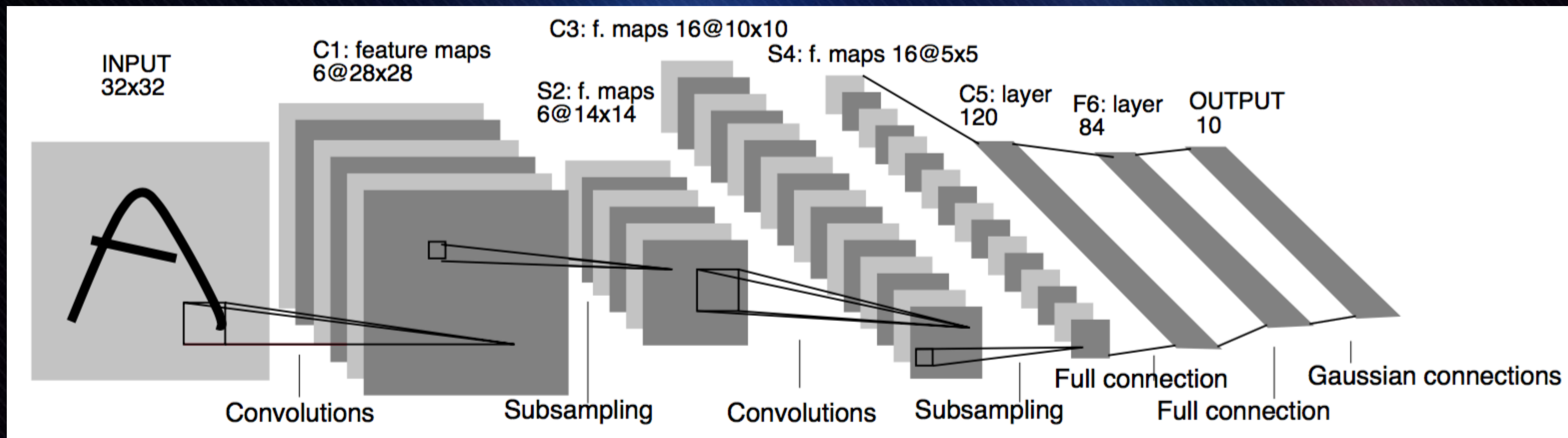
Gradient-based learning applied to document recognition, Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Proc. of IEEE, 1998



ImageNet classification with deep convolutional neural networks, A. Krizhevsky, L. Sutskever, G. Hinton, NIPS 2012

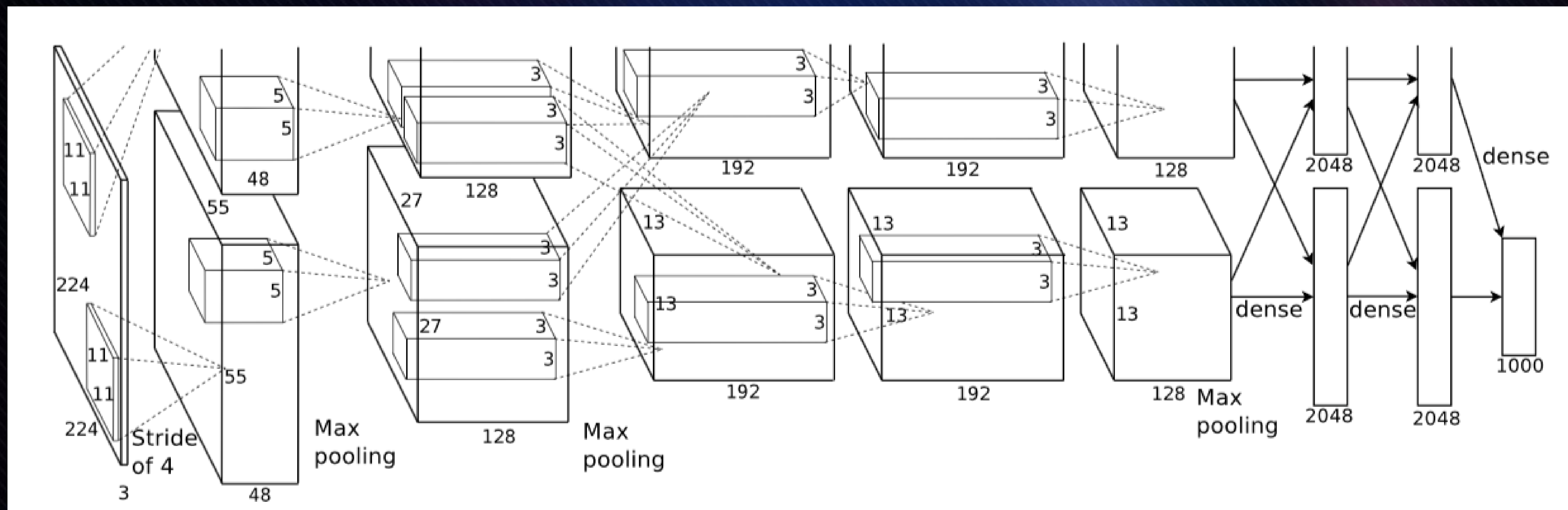


LeNet-5



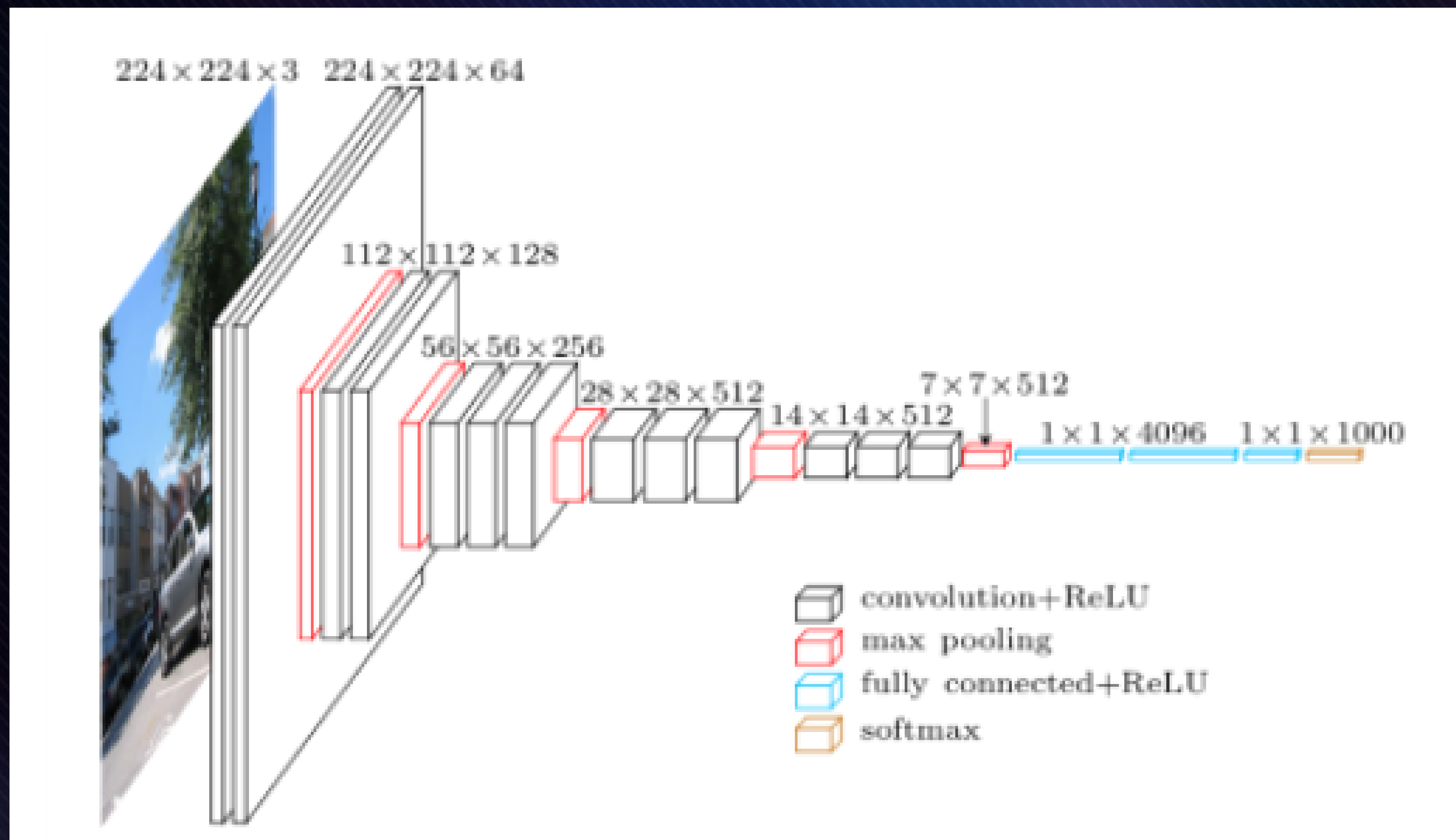
Gradient-based learning applied to document recognition, Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Proc. of IEEE, 1998

AlexNet



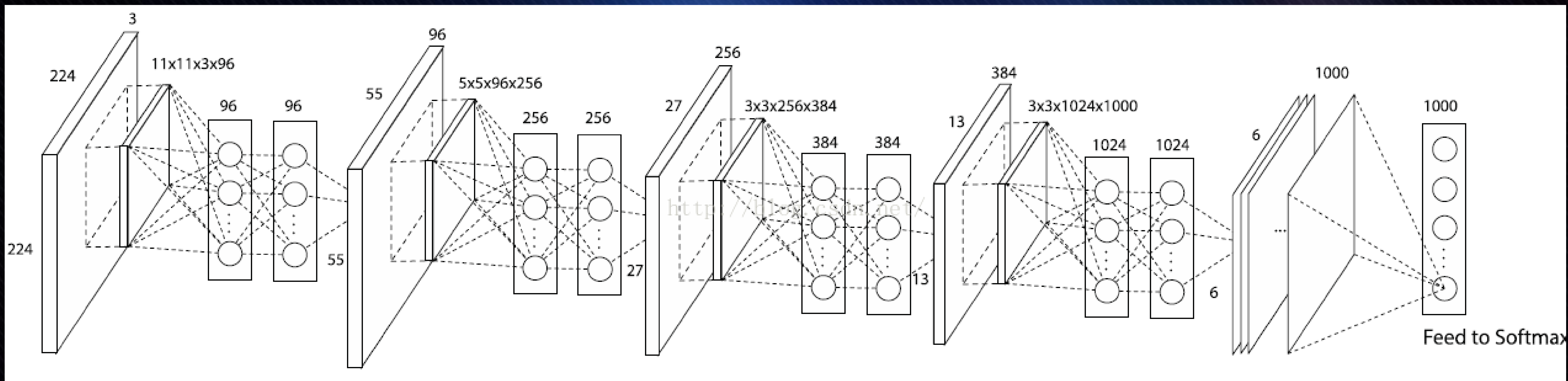
ImageNet classification with deep convolutional neural networks, A. Krizhevsky, L. Sutskever, G. Hinton, NIPS 2012

VGG



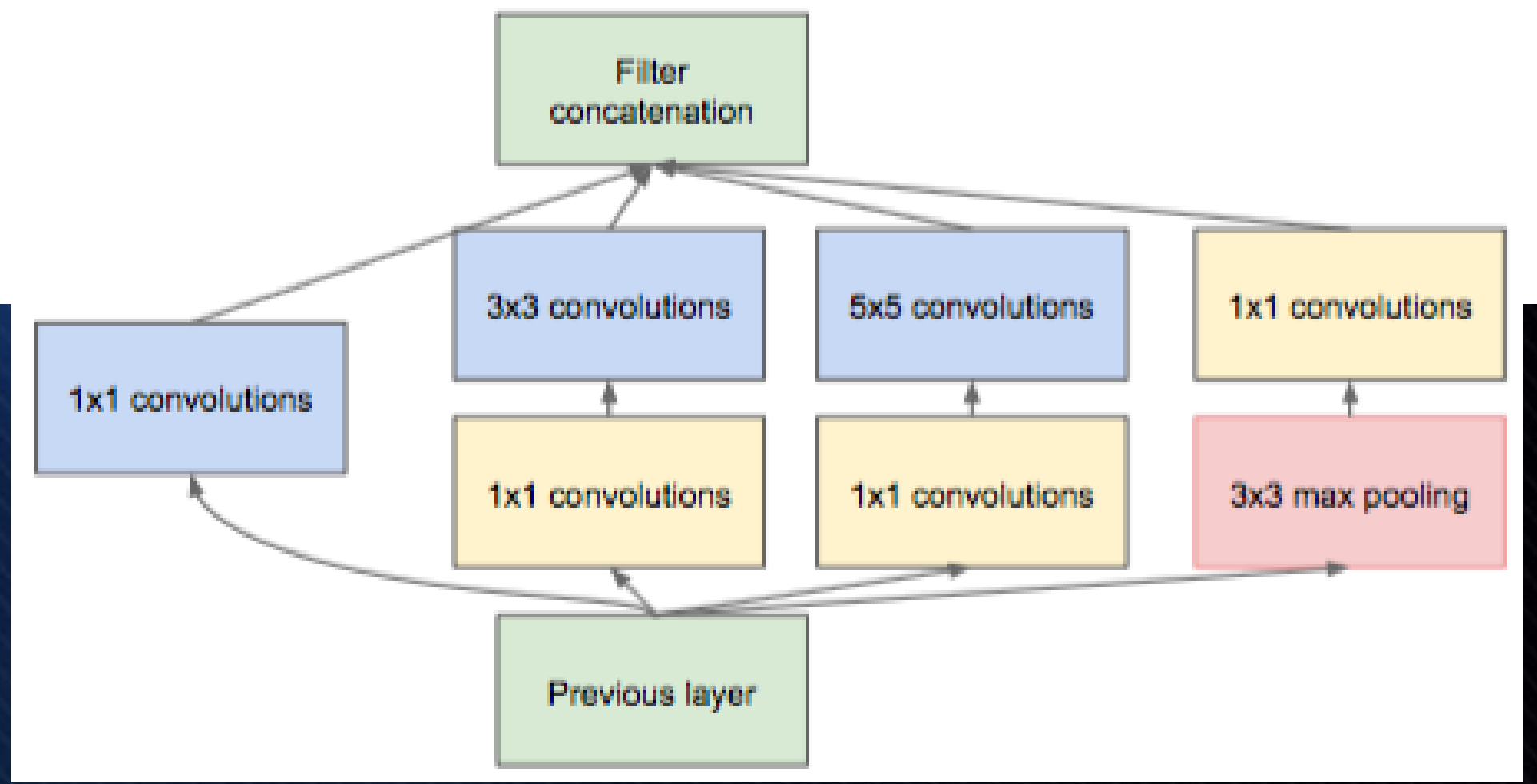
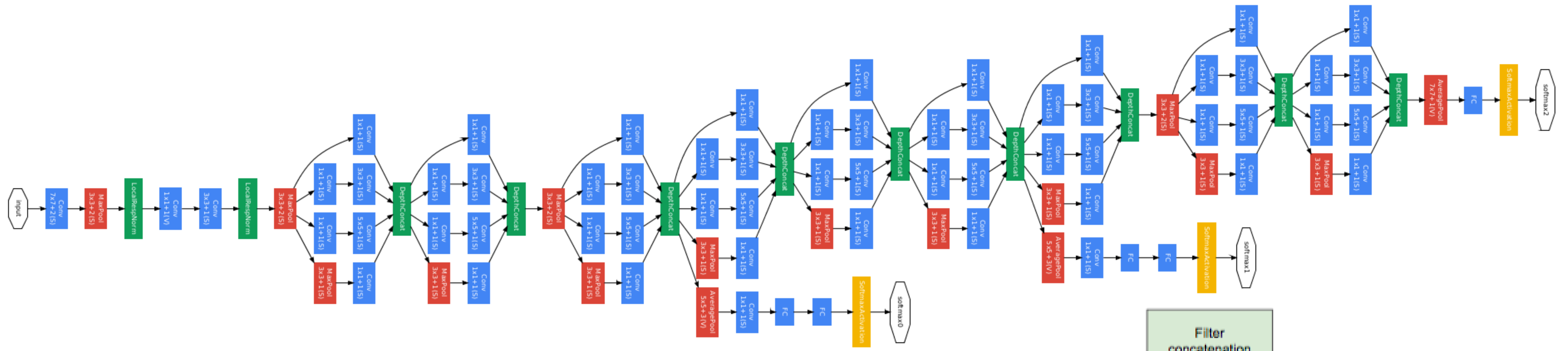
Very deep convolutional networks for large-scale image recognition, K. Simonyan, A. Zisserman, ICLR 2015

Network in network



Network in network, M. Lin, Q. Chen, S. Yan, 2014

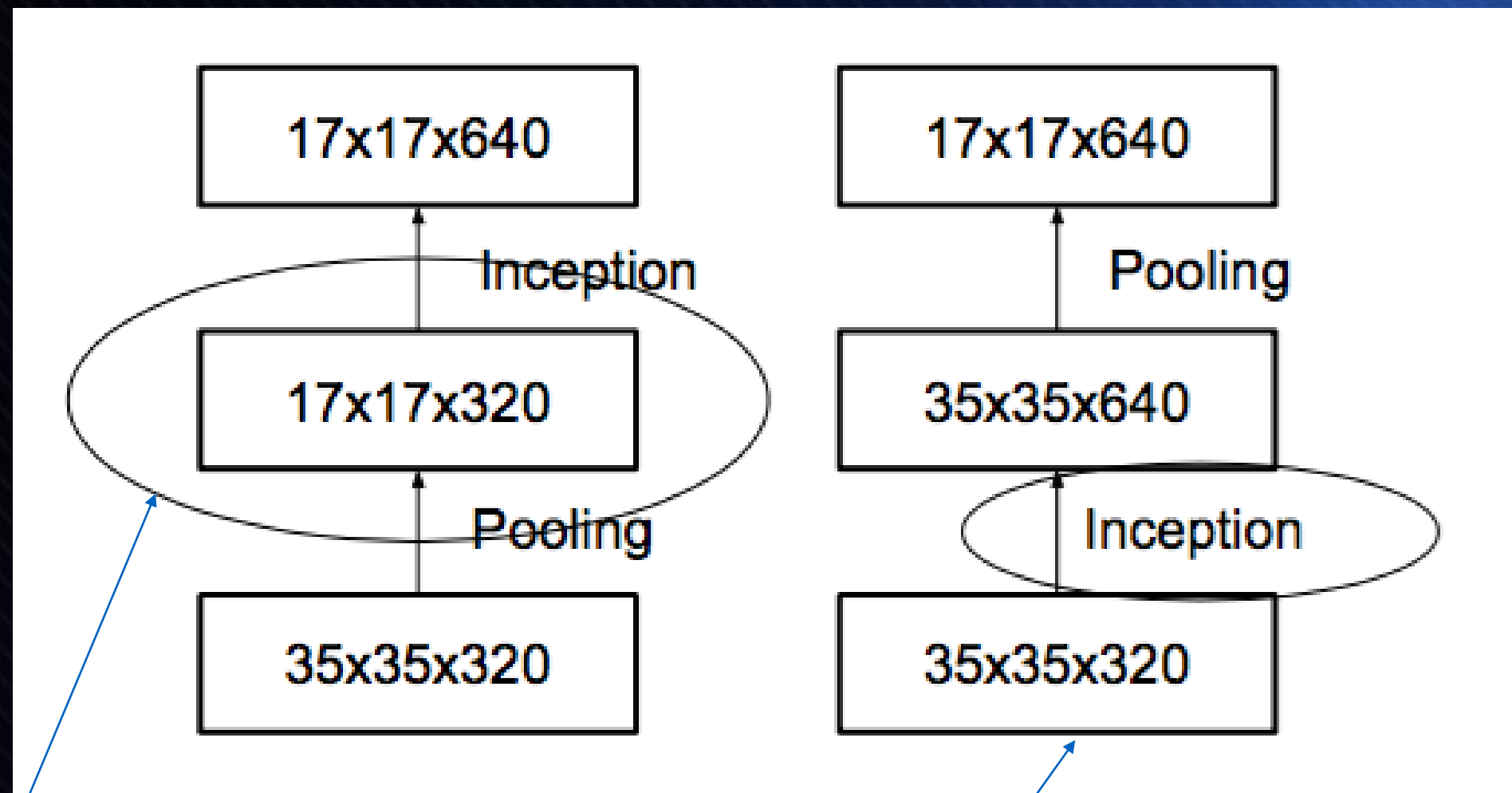
GoogLeNet



Going deeper with convolutions, C. Szegedy, W. Liu, Y. Jia, etc., CVPR 2015

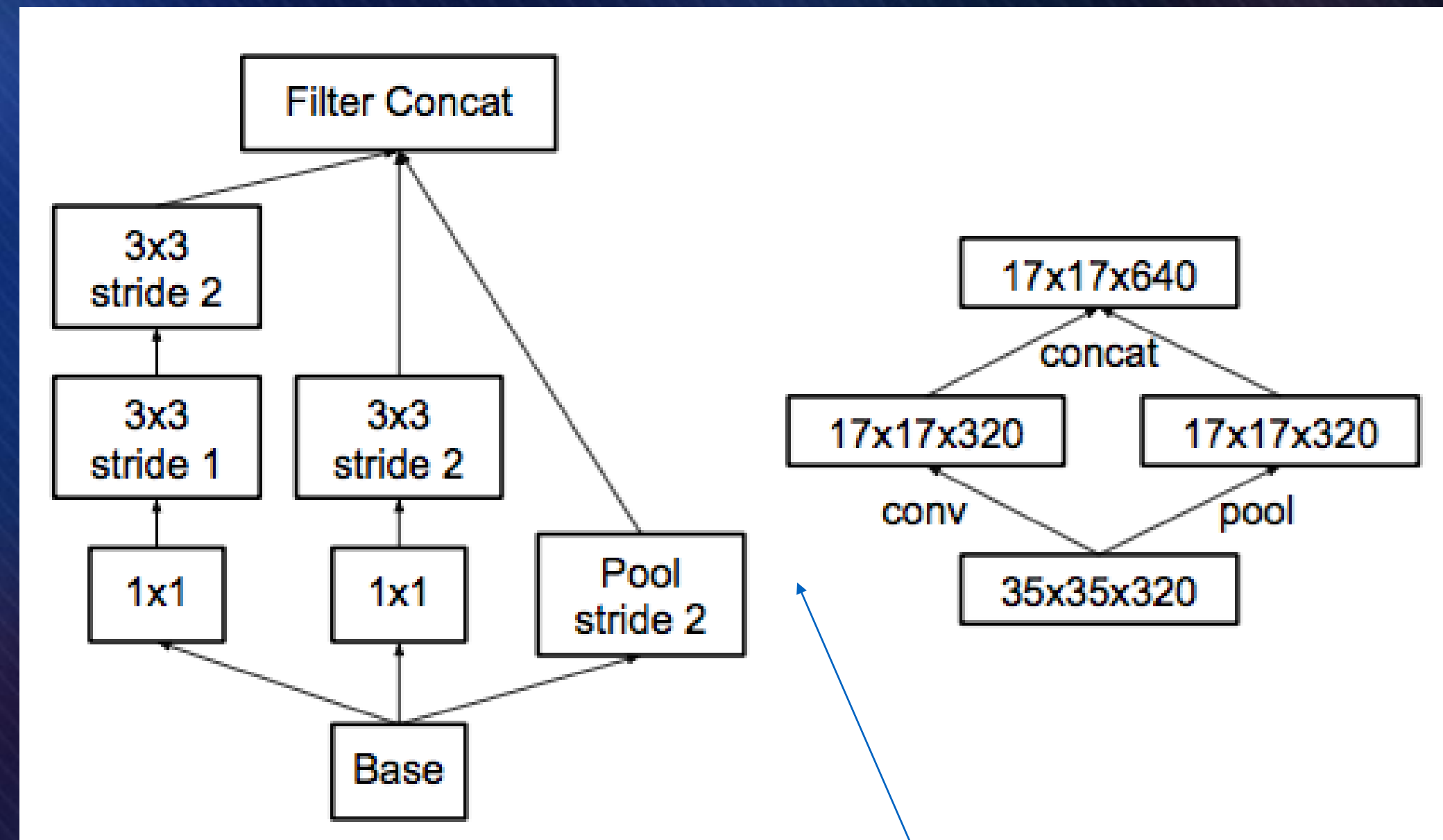
Design Principles of Inception v3

- Avoid representational bottlenecks



A bottleneck

Expensive

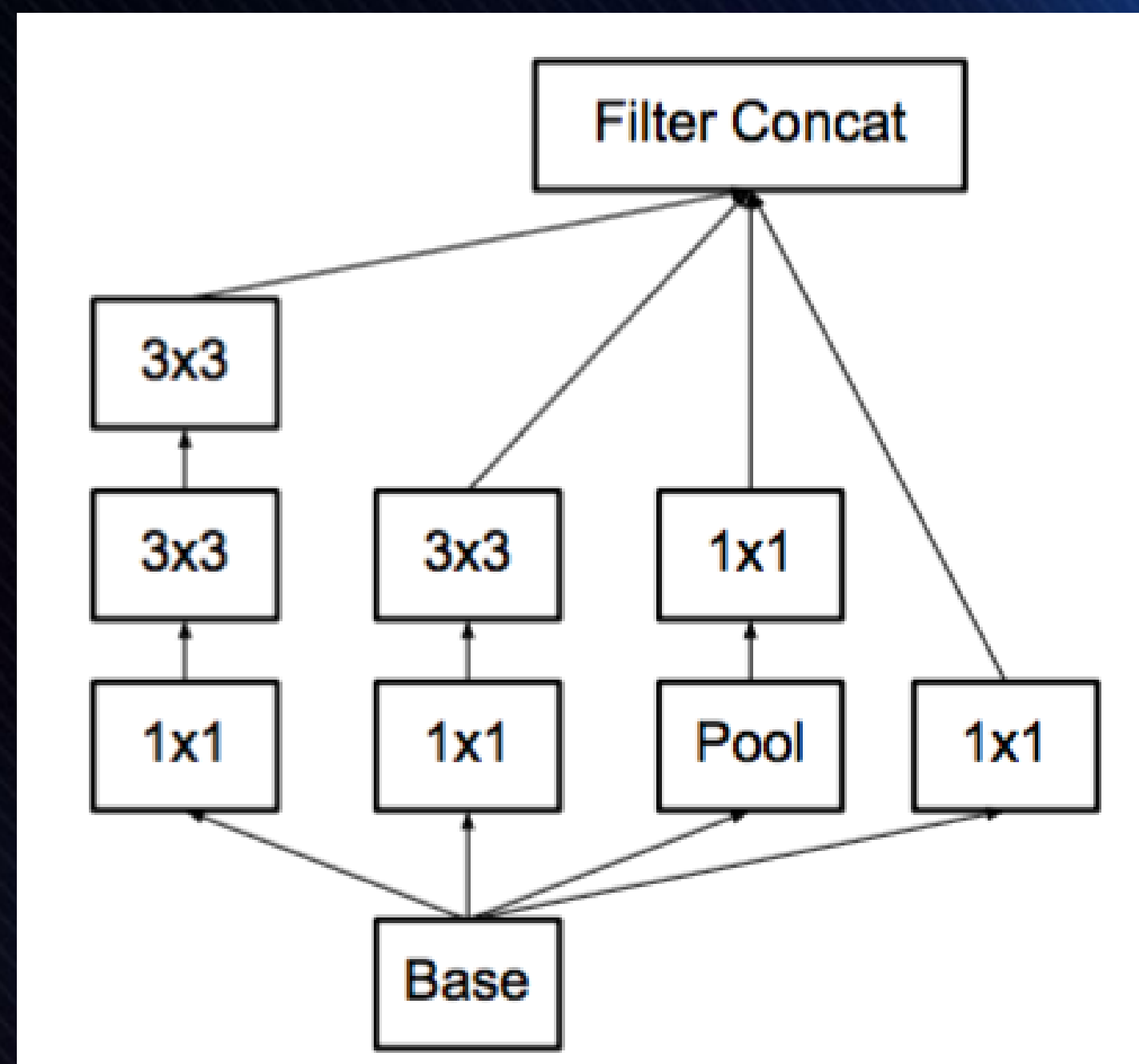


Proposed Solution

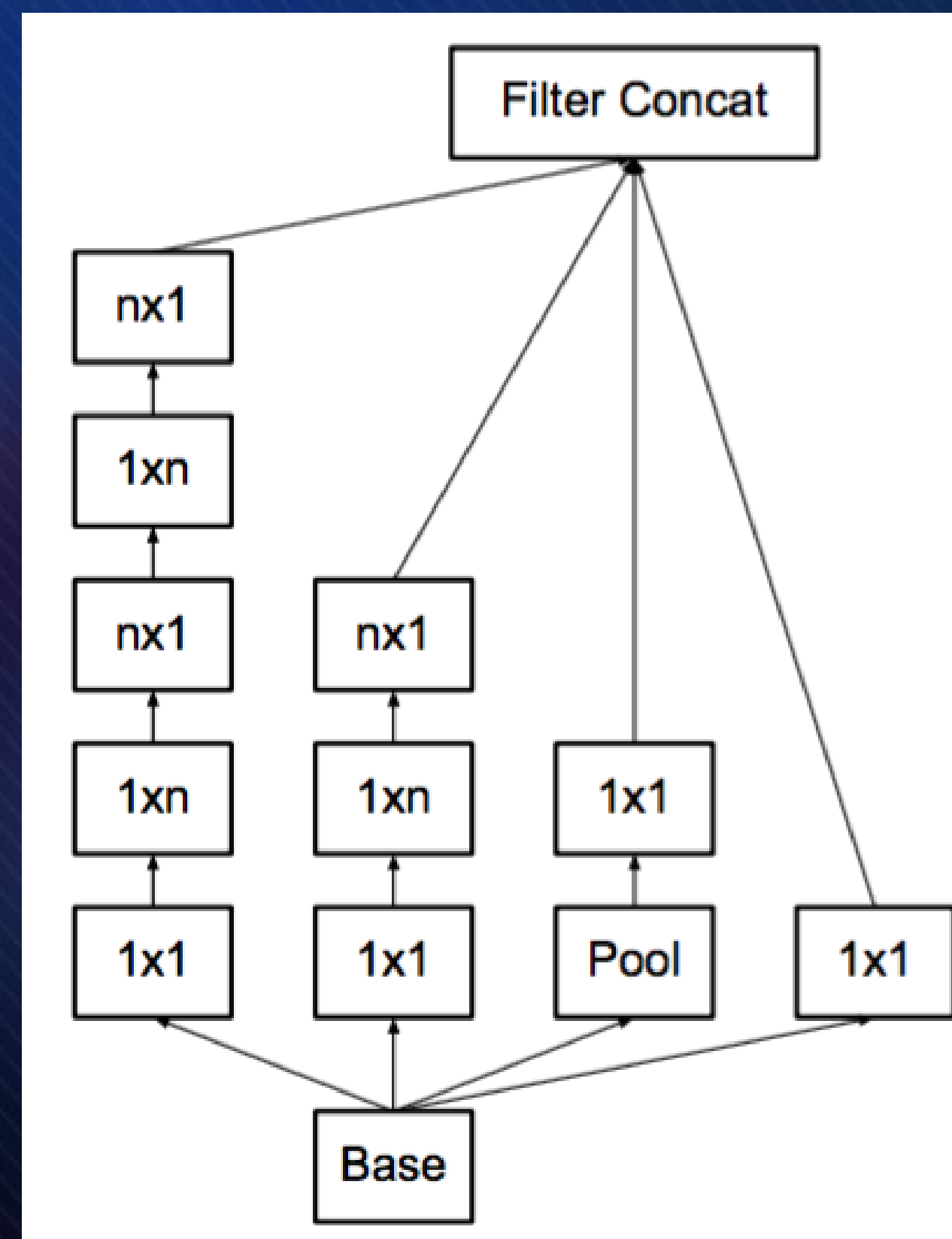
Rethinking the inception architecture for computer vision, C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, 2015

Design Principles of Inception v3

- Avoid representational bottlenecks
- Spatial aggregation can be done over lower dimensional embeddings



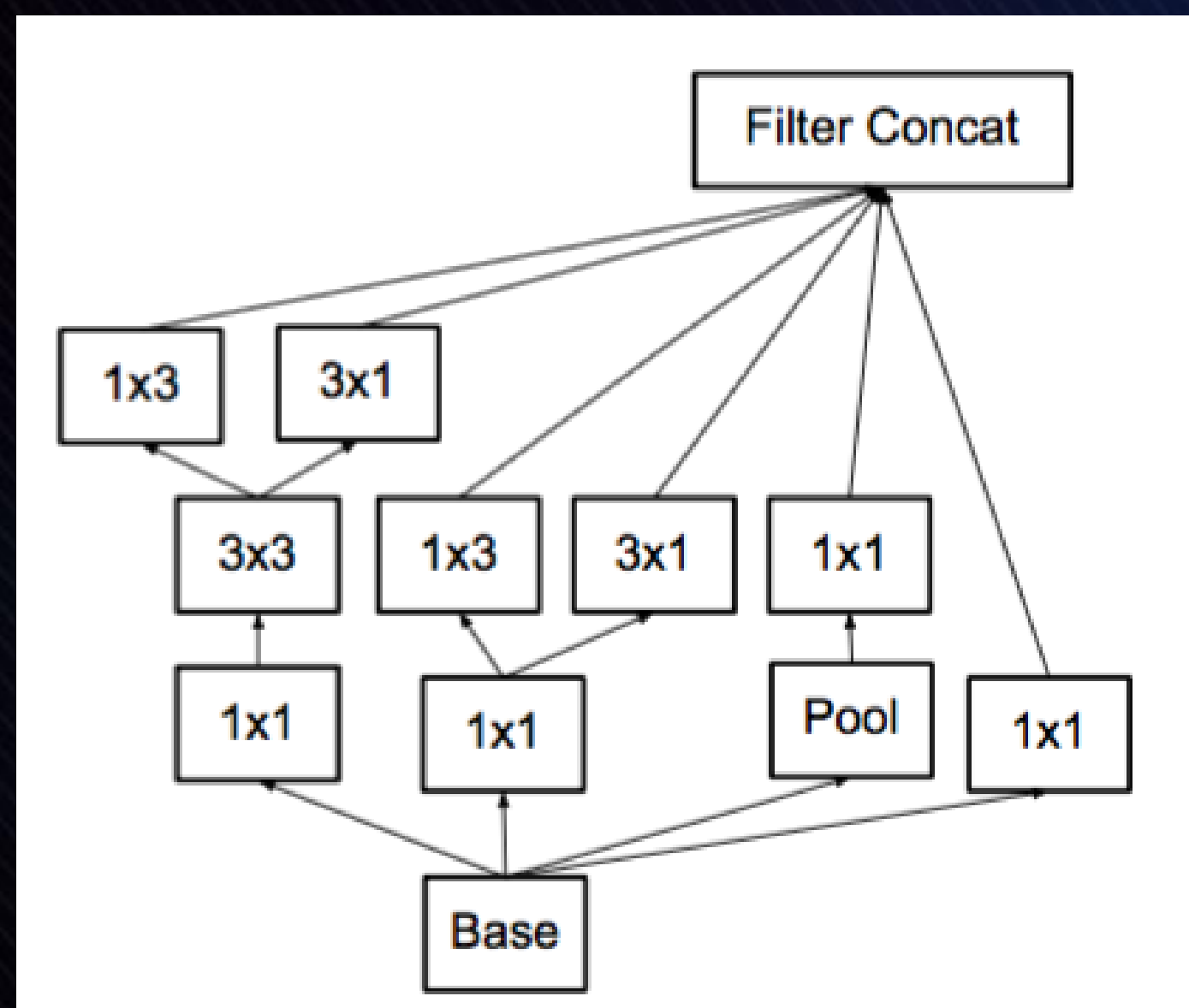
Used on 35 x 35 grids



Used on 17 x 17 grids
n=7

Design Principles of Inception v3

- Avoid representational bottlenecks
- Spatial aggregation can be done over lower dimensional embeddings
- Higher dimensional representations are easier to process locally



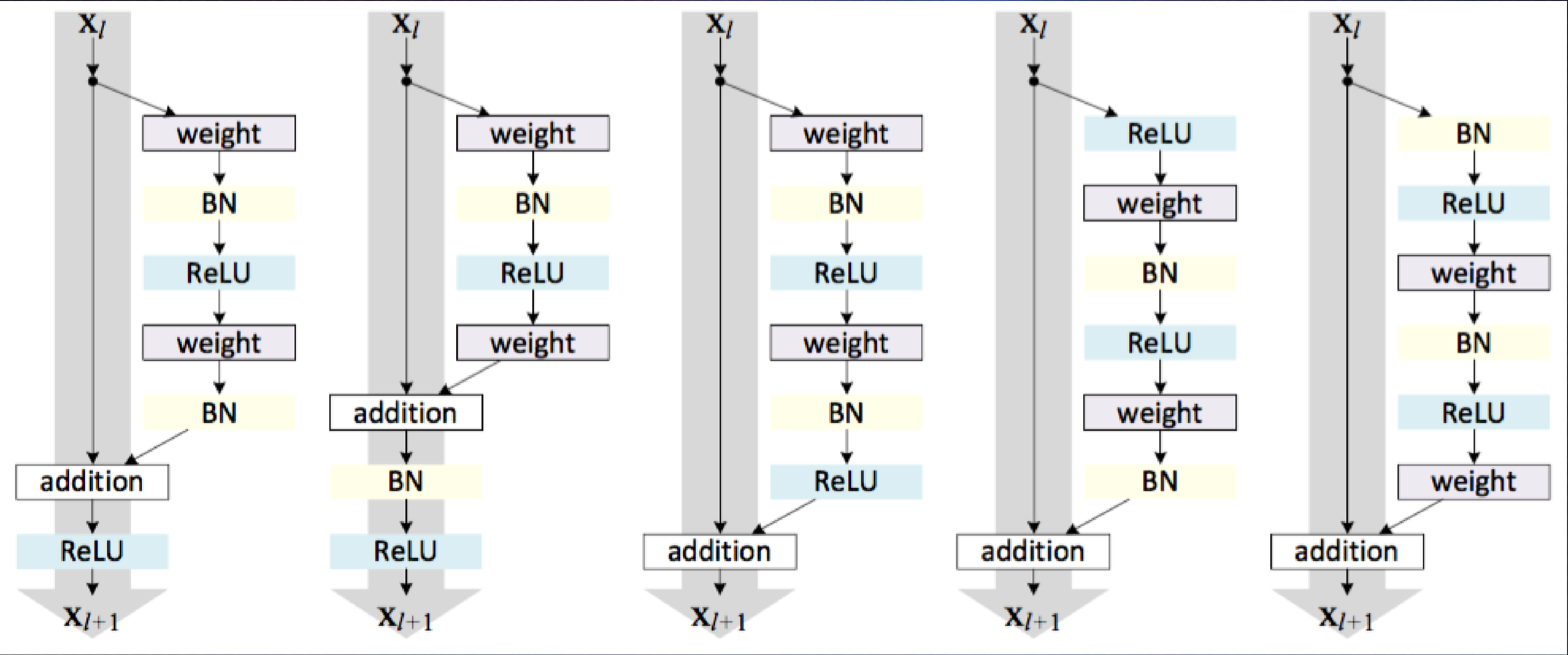
Used on the coarsest (8 x 8) grids

Design Principles of Inception v3

- Avoid representational bottlenecks
- Spatial aggregation can be done over lower dimensional embeddings
- Higher dimensional representations are easier to process locally
- Balance the width and depth of the network

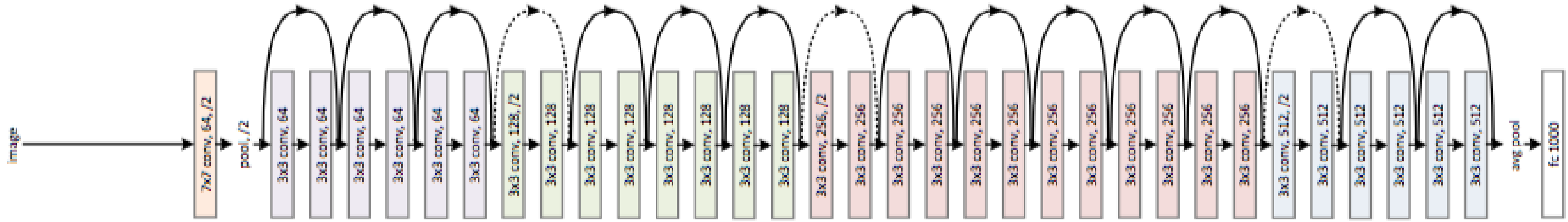
ResNet

Residual block v1



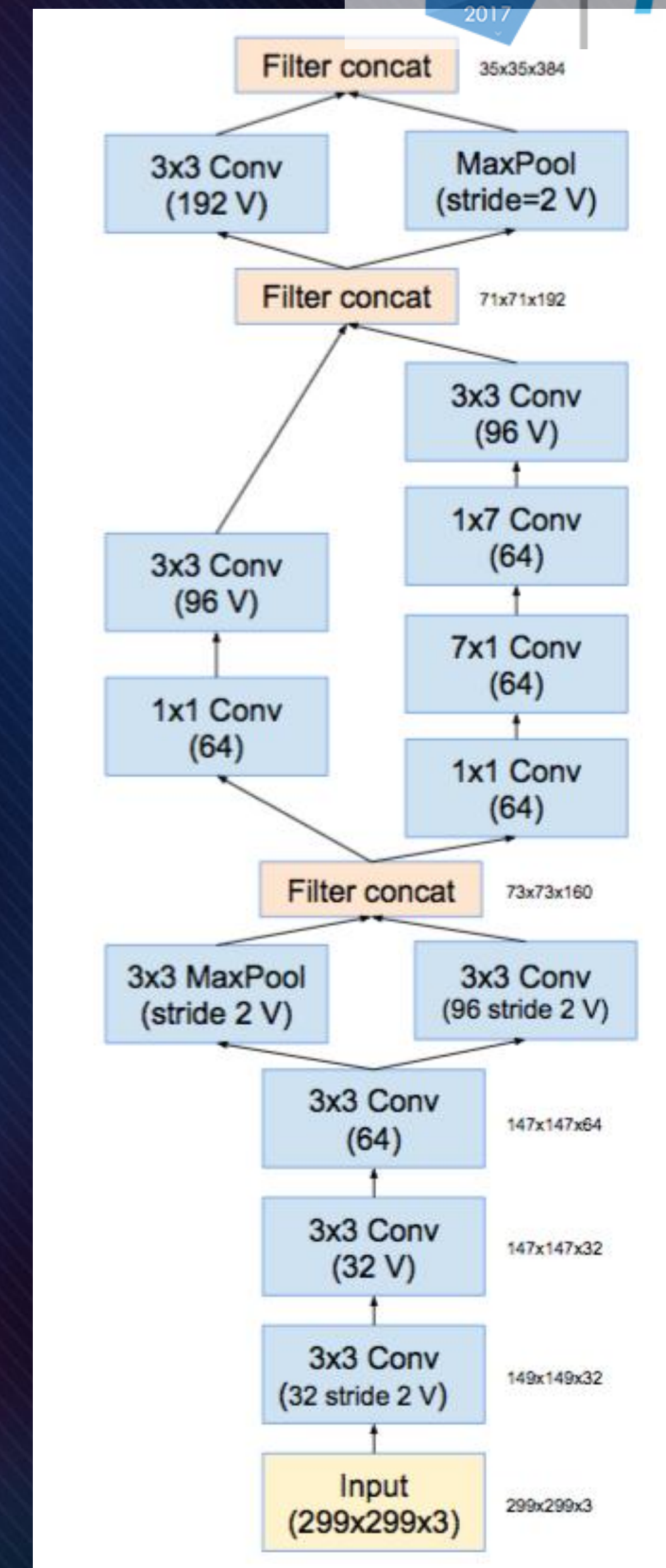
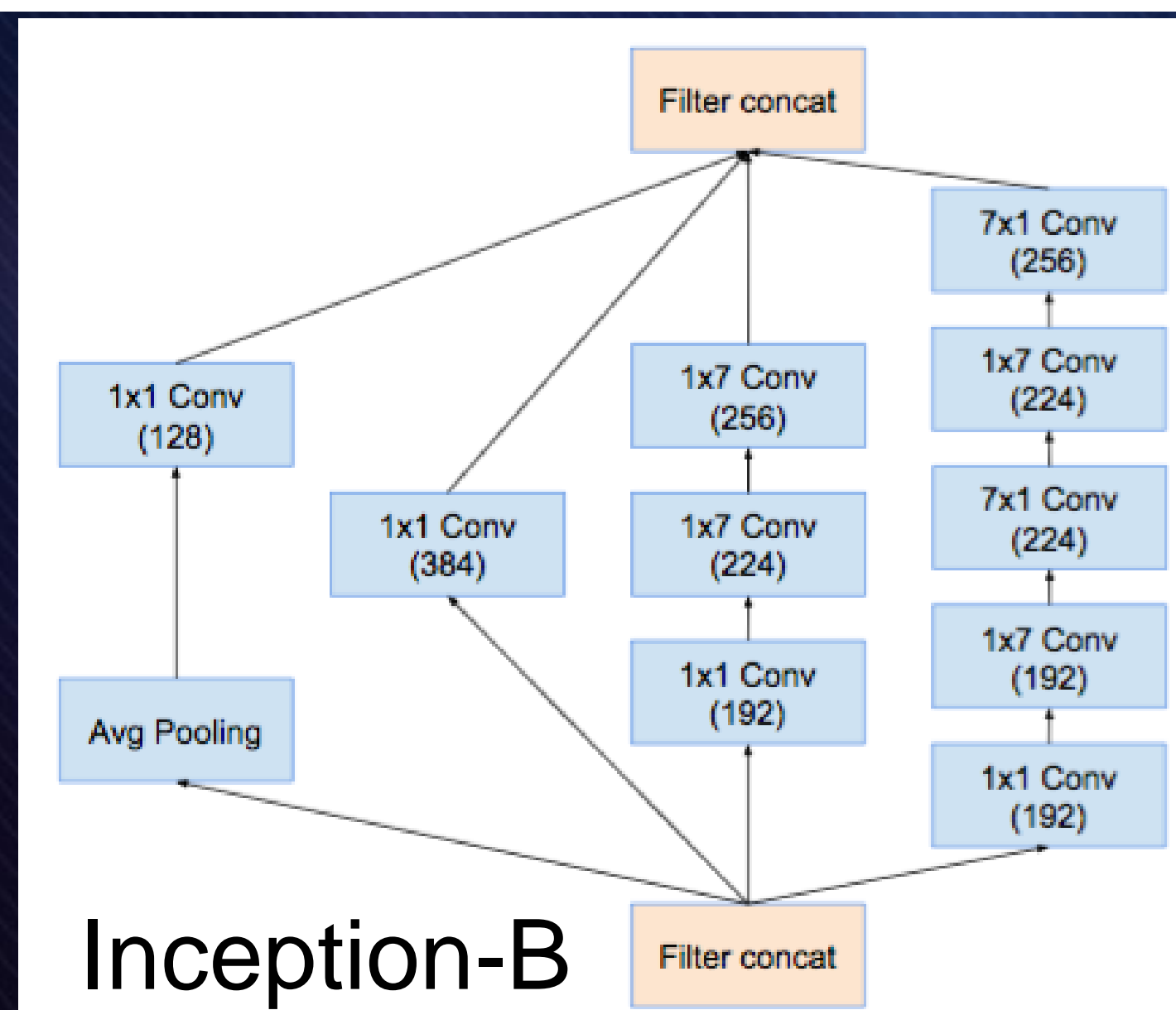
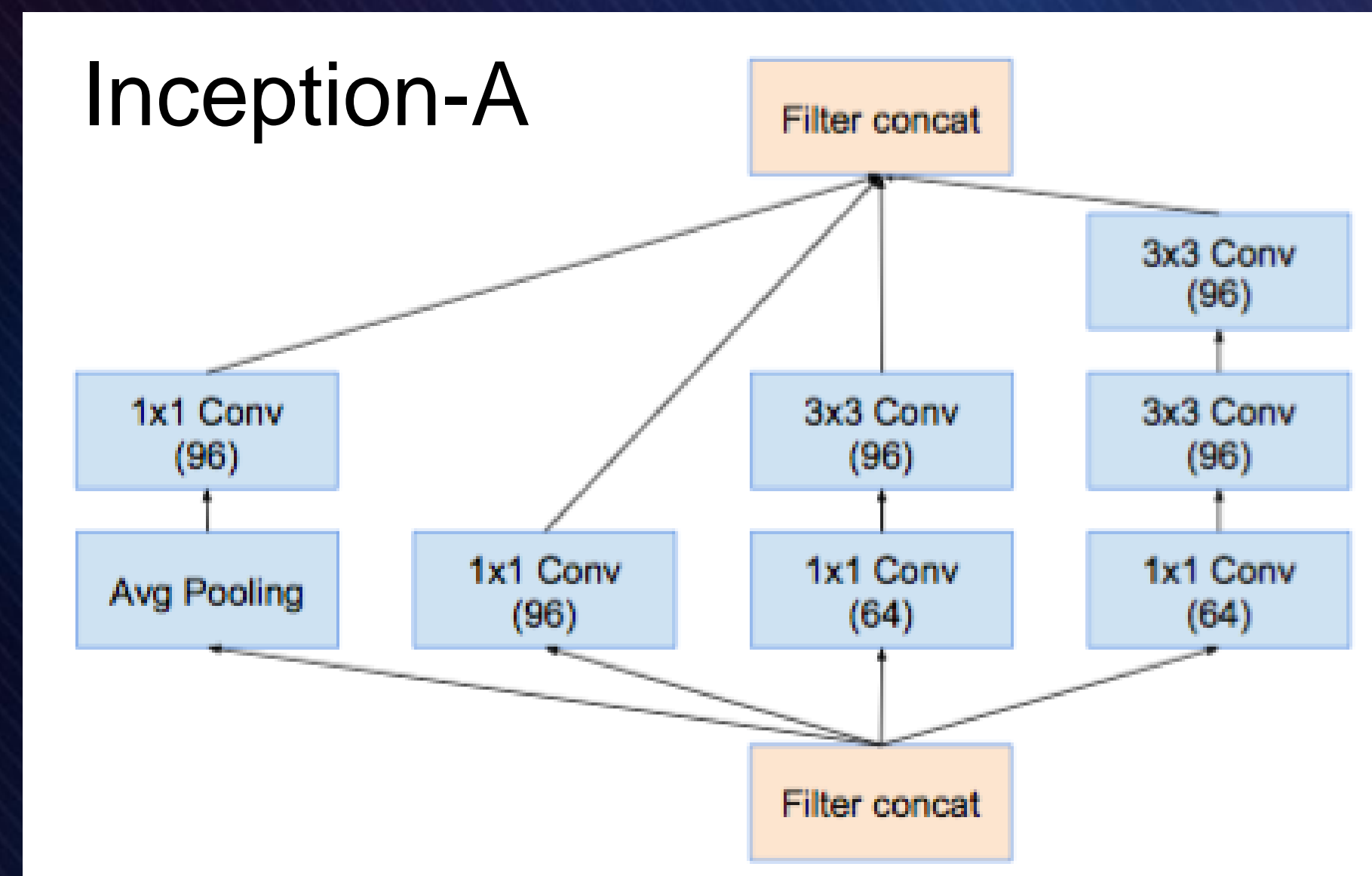
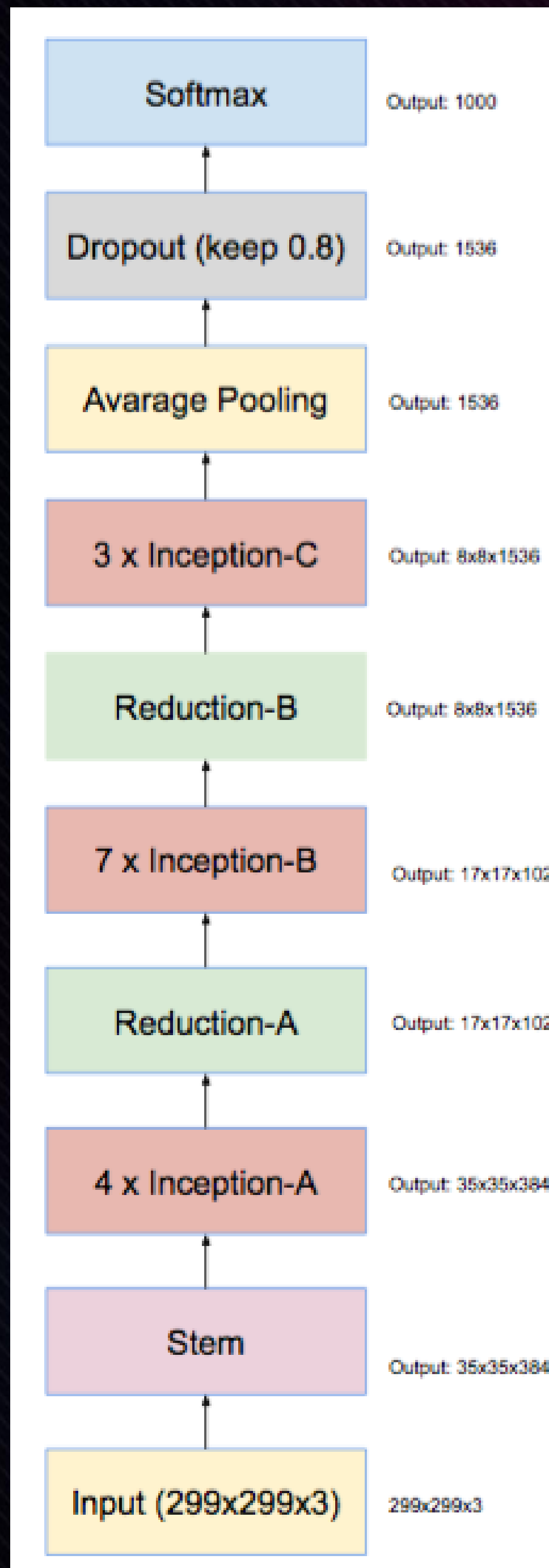
Residual block v2

34-layer residual



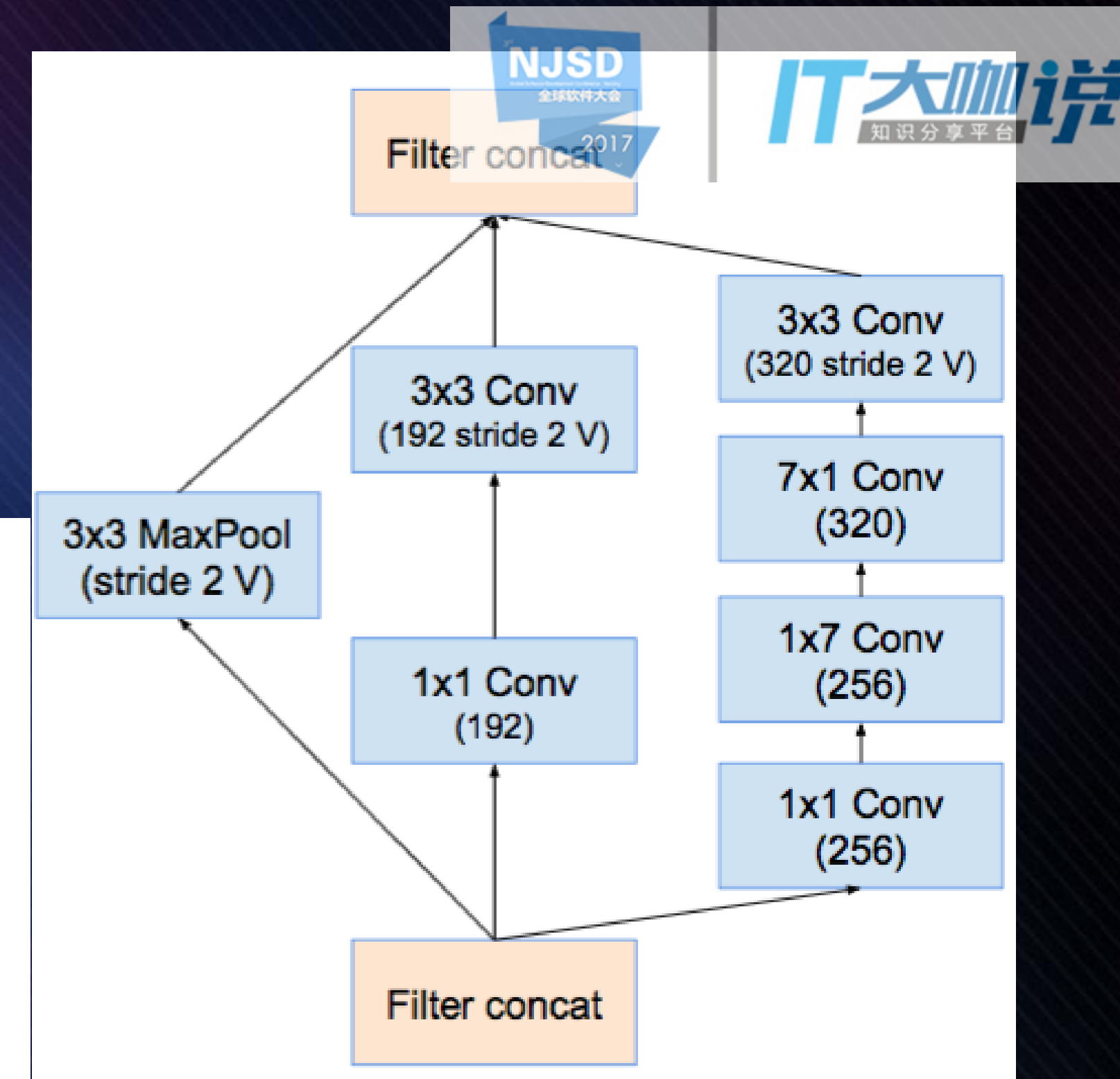
Deep residual learning for image recognition, K. He, X. Zhang, S. Ren, J. Sun

Inception v4

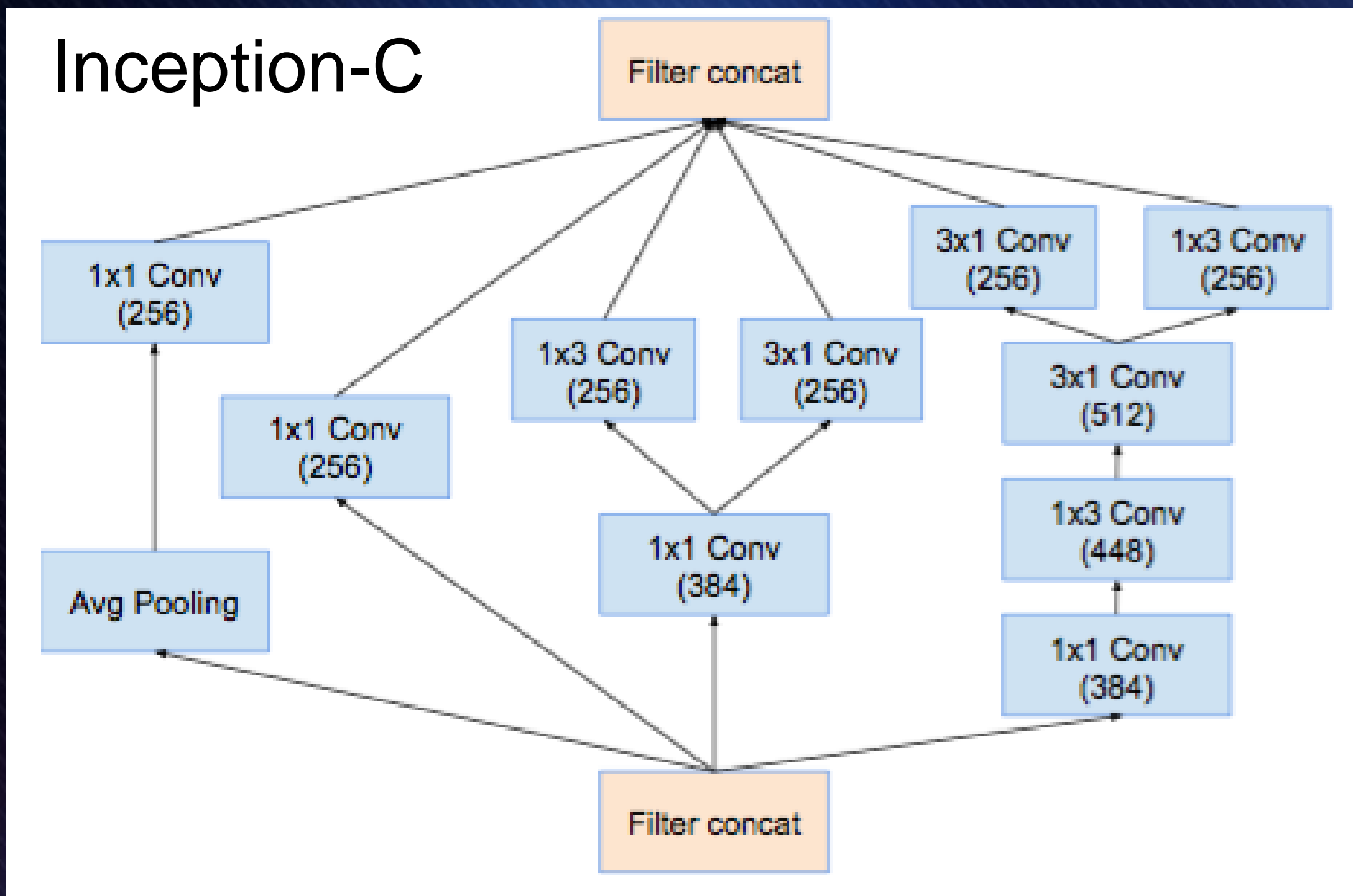


Inception v4

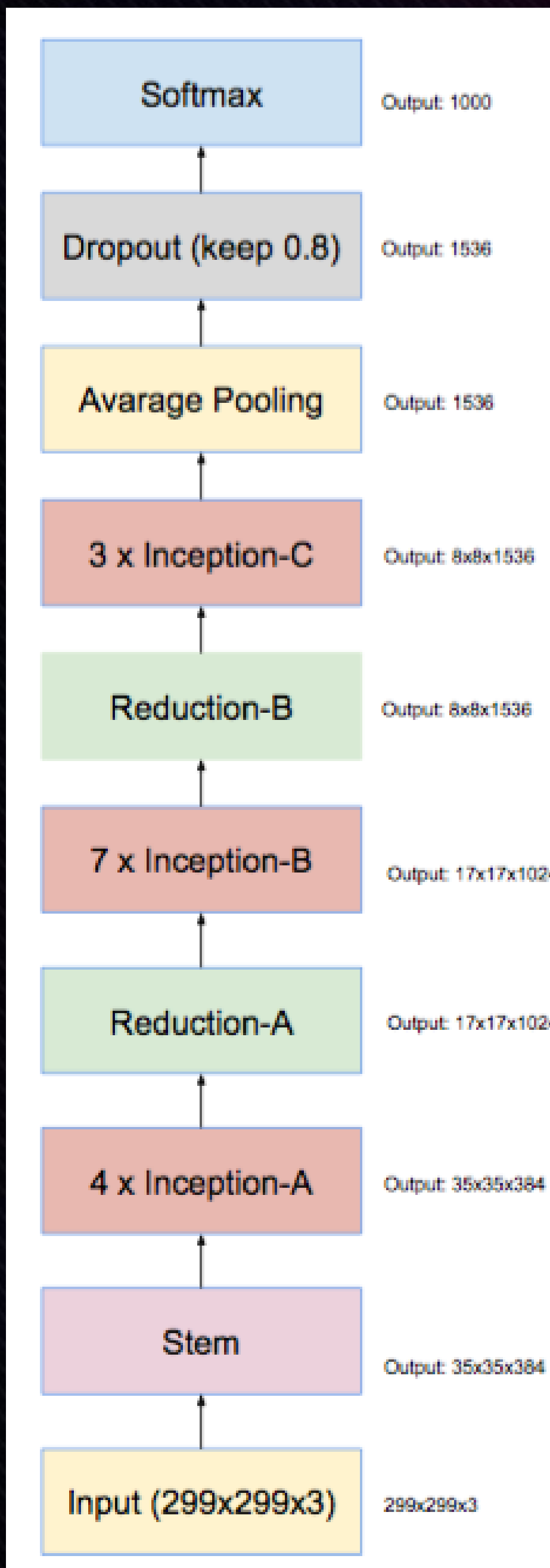
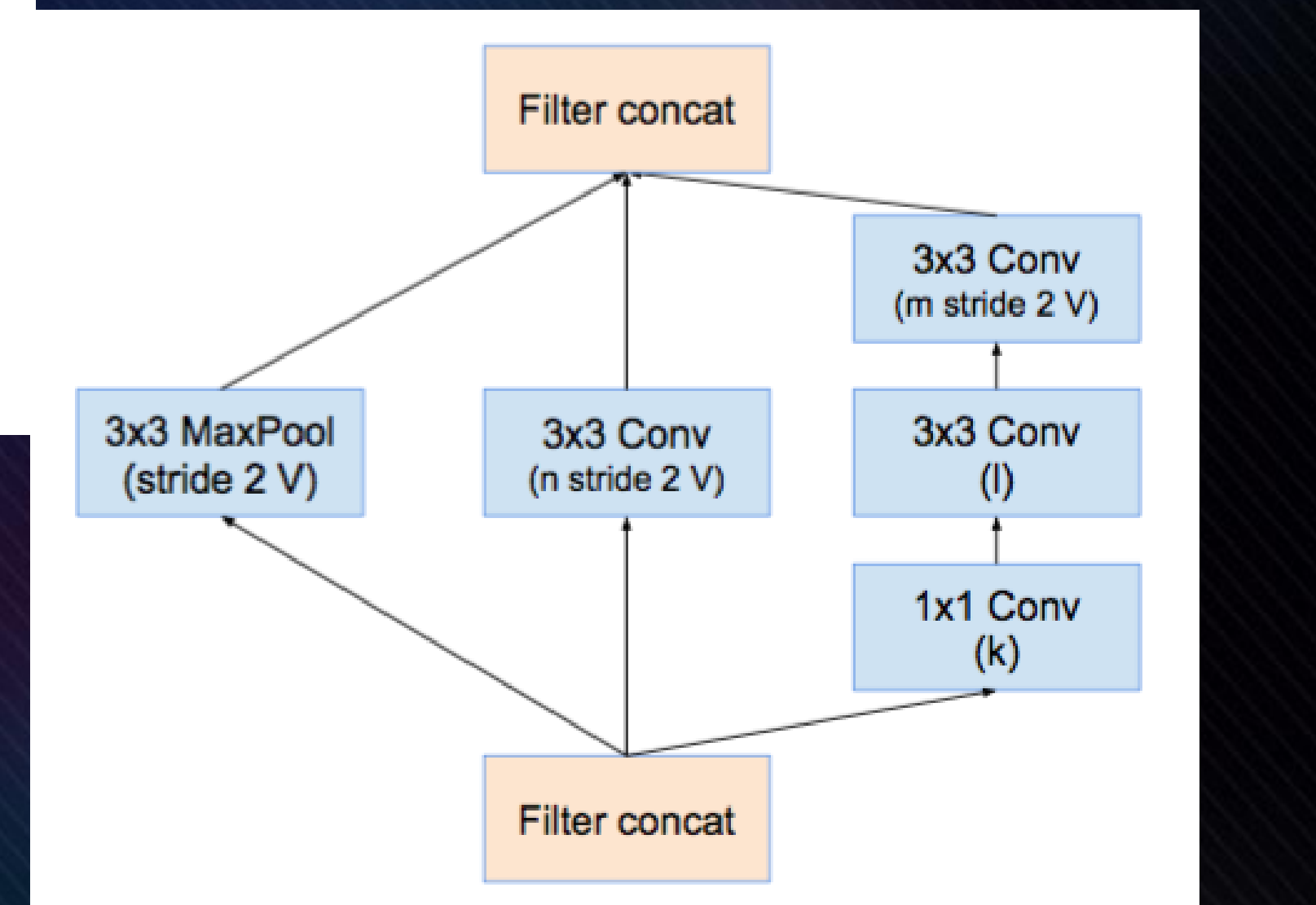
Reduction-A



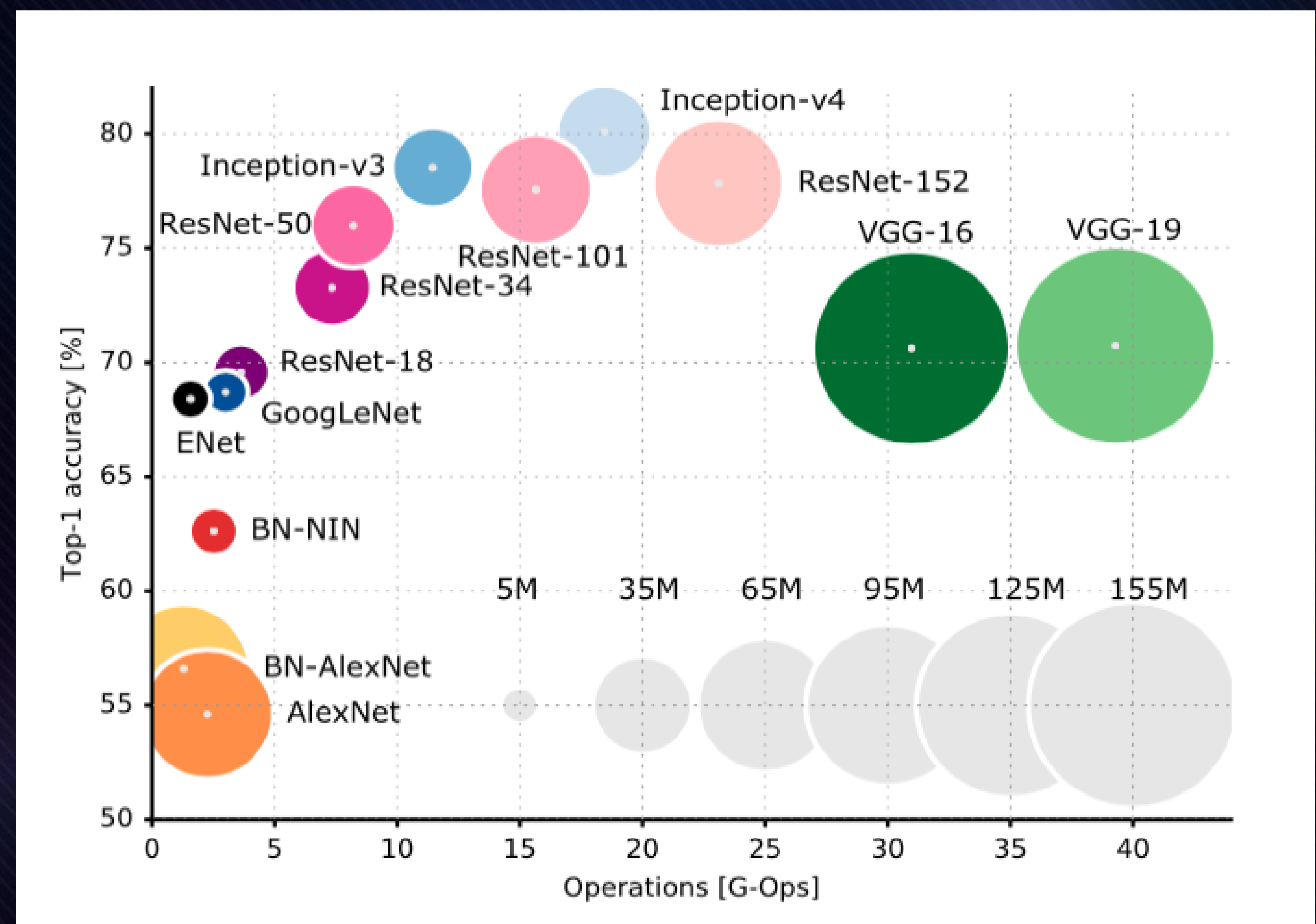
Inception-C



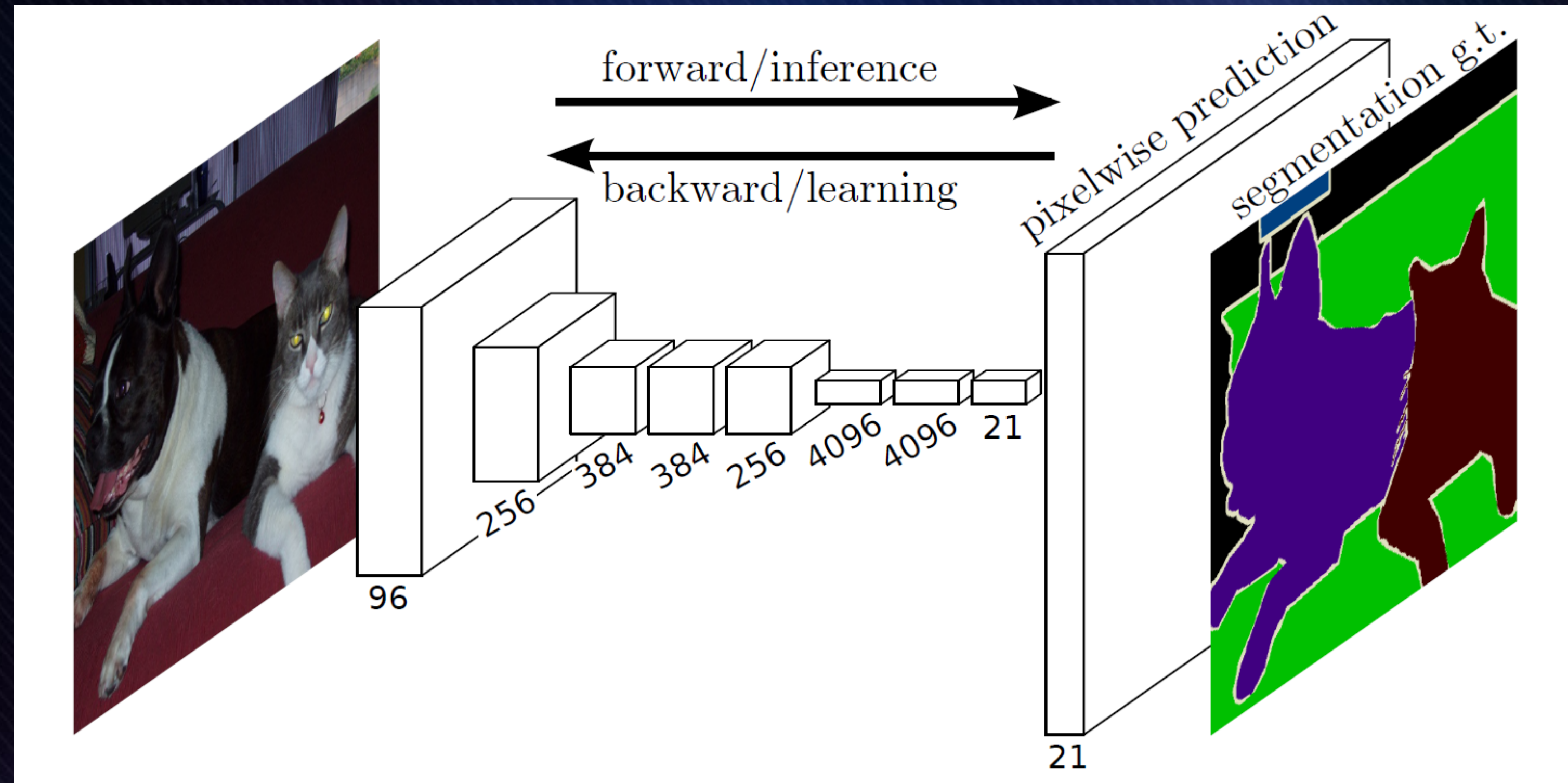
Reduction-B



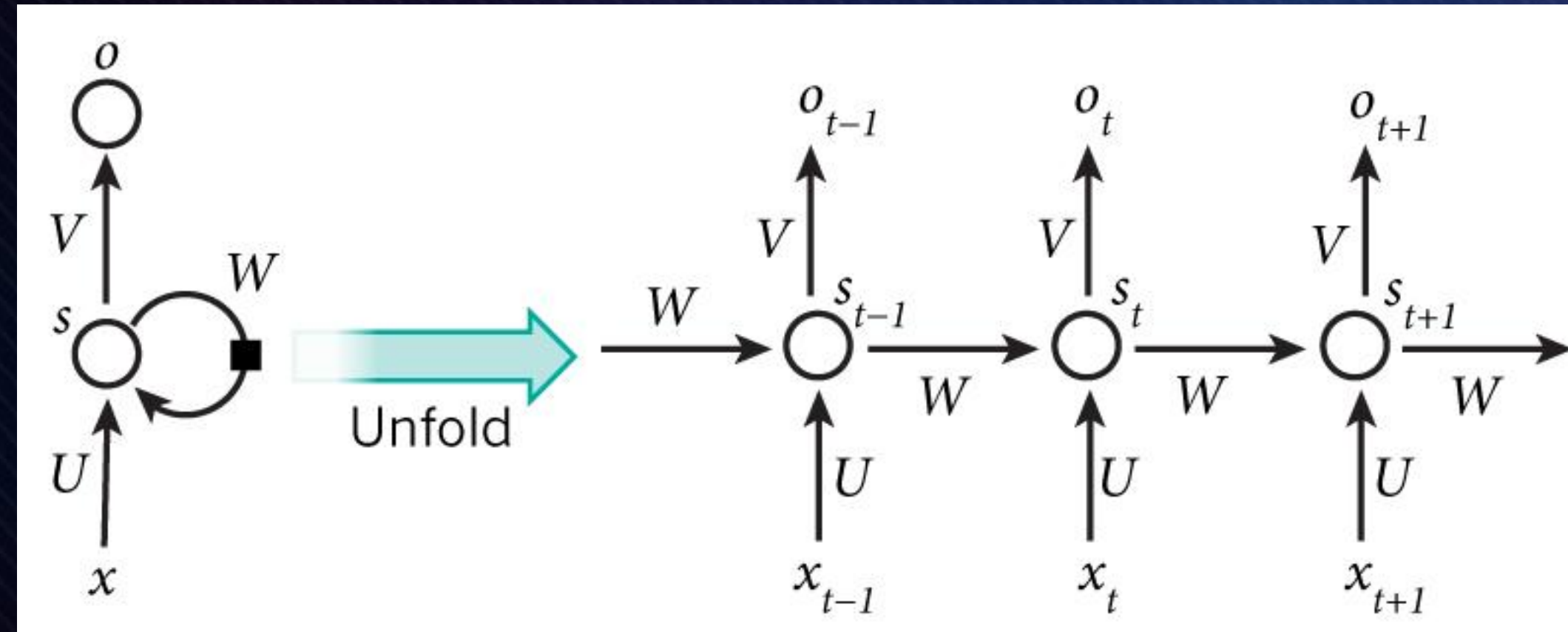
Performances, Size and Operations



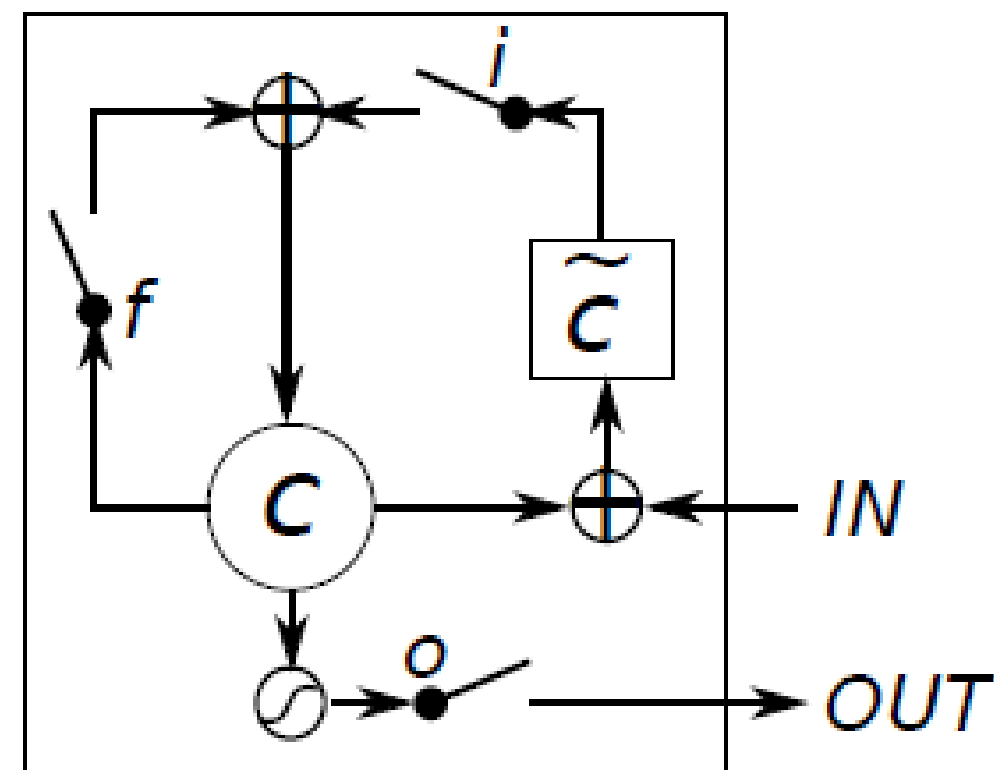
Fully convolutional network



Fully convolutional networks for semantic segmentation, E. Shelhamer, J. Long, and T. Darrell, 2016

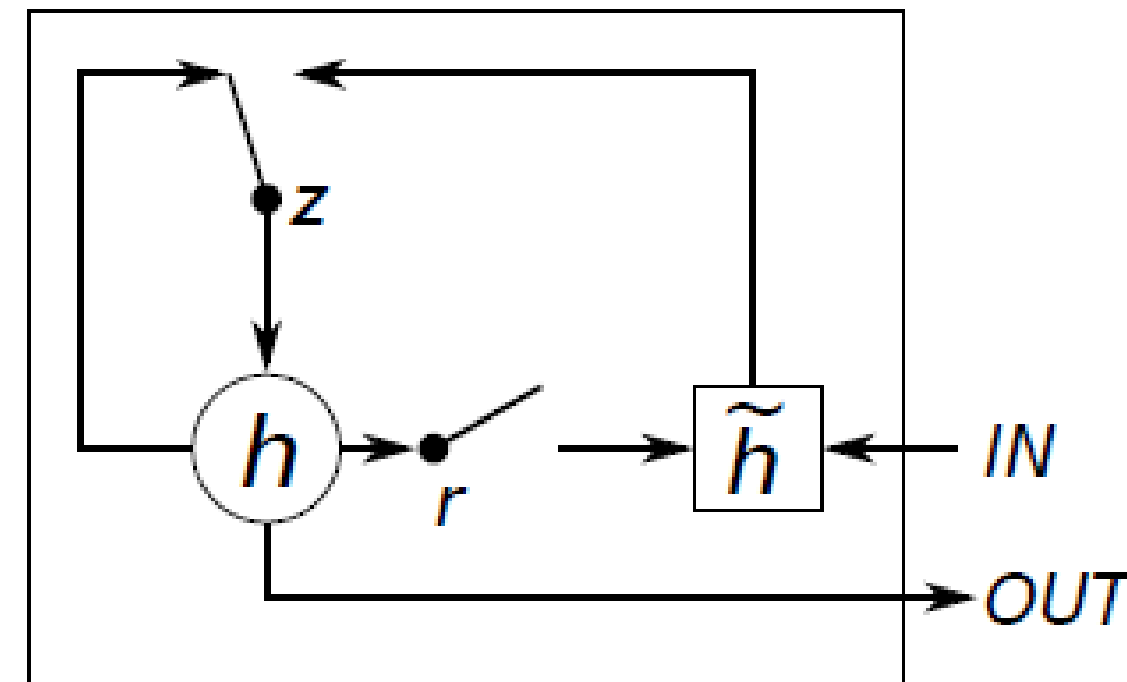


LSTM



(a) Long Short-Term Memory

Gated Recurrent Unit



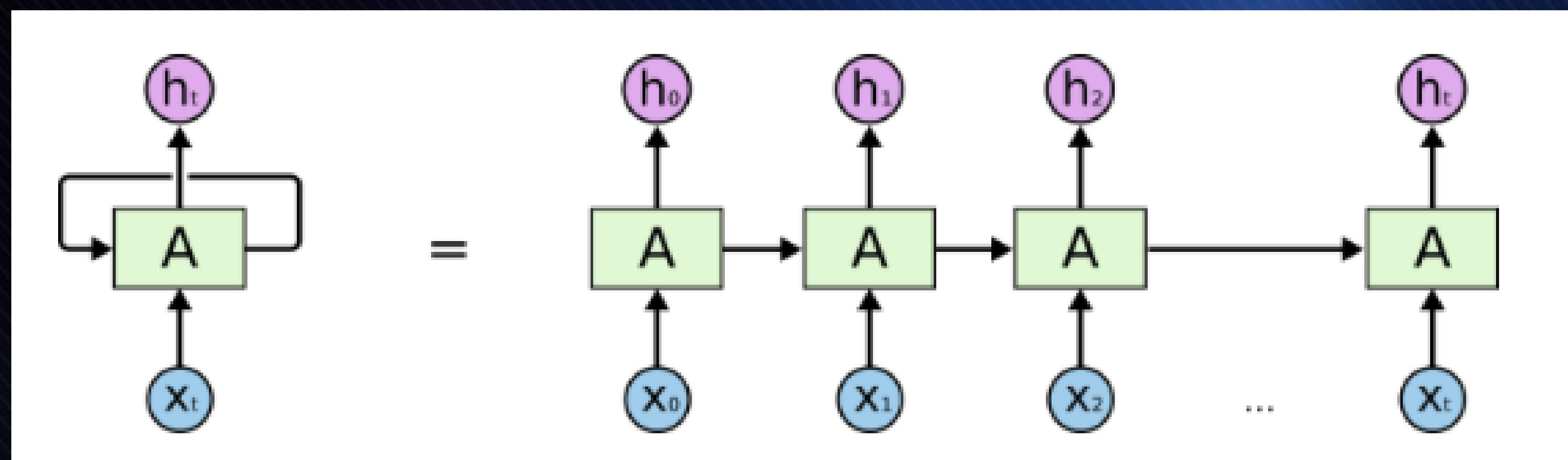
(b) Gated Recurrent Unit

Long short-term memory, S. Hochreiter and J. Schmidhuber, 1999

Learning phrase representations using RNN encoder-decoder for statistical machine translation, K.Cho, etc

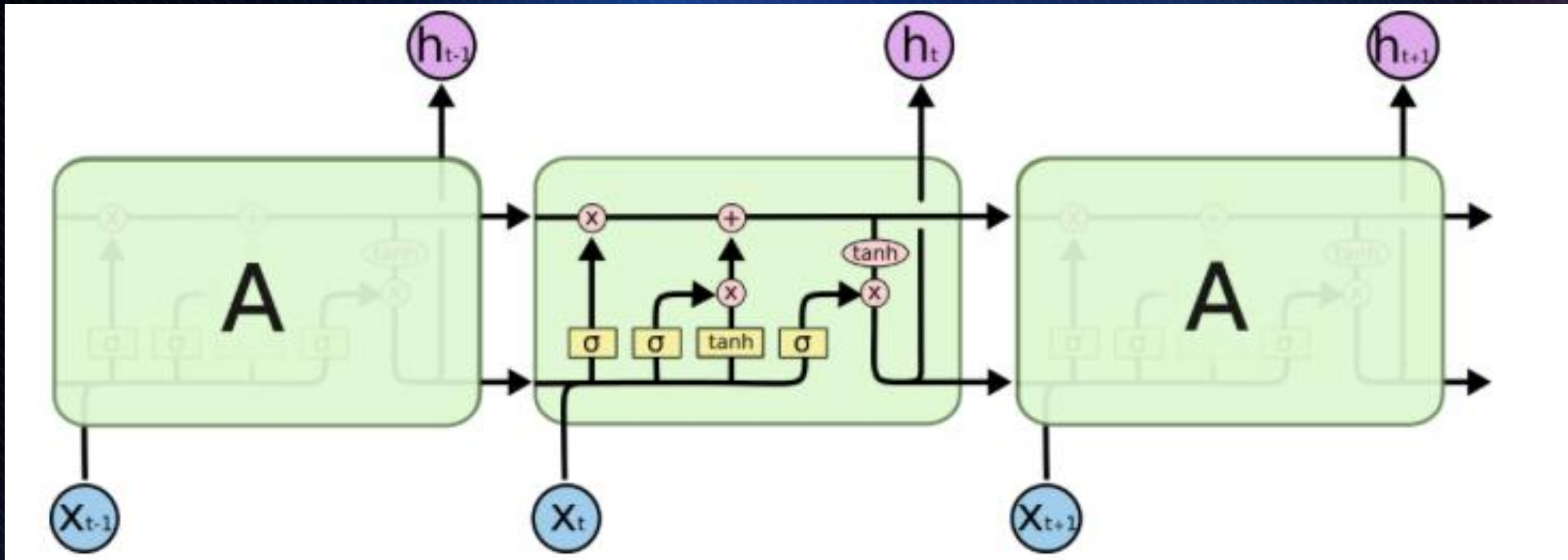
Recurrent Neural Networks

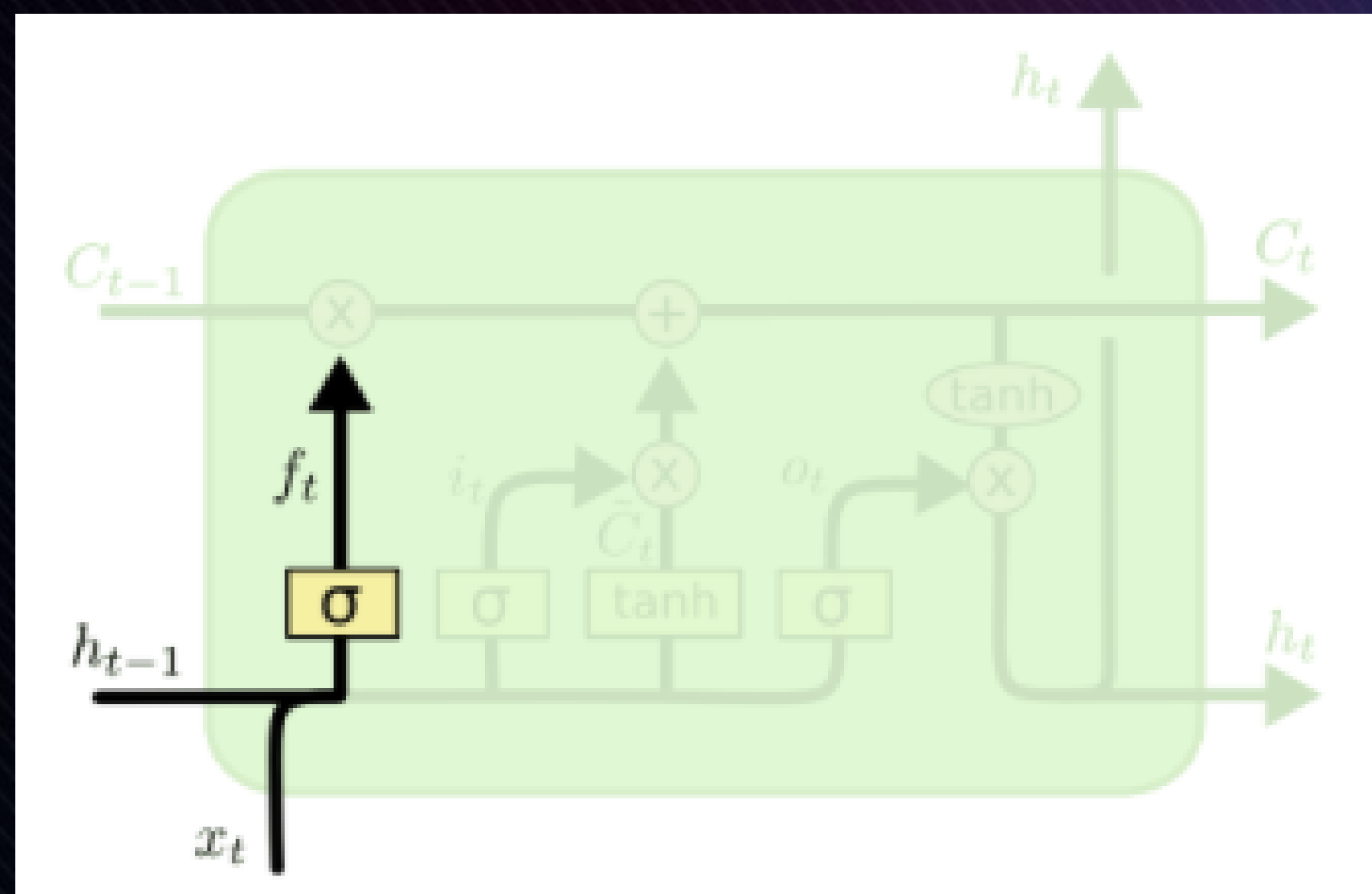
- Training RNNs used to be extremely difficult
 - An Unrolled RNN is equivalent to a very deep net with tiled weights
 - derivatives are susceptible to vanishing or exploding



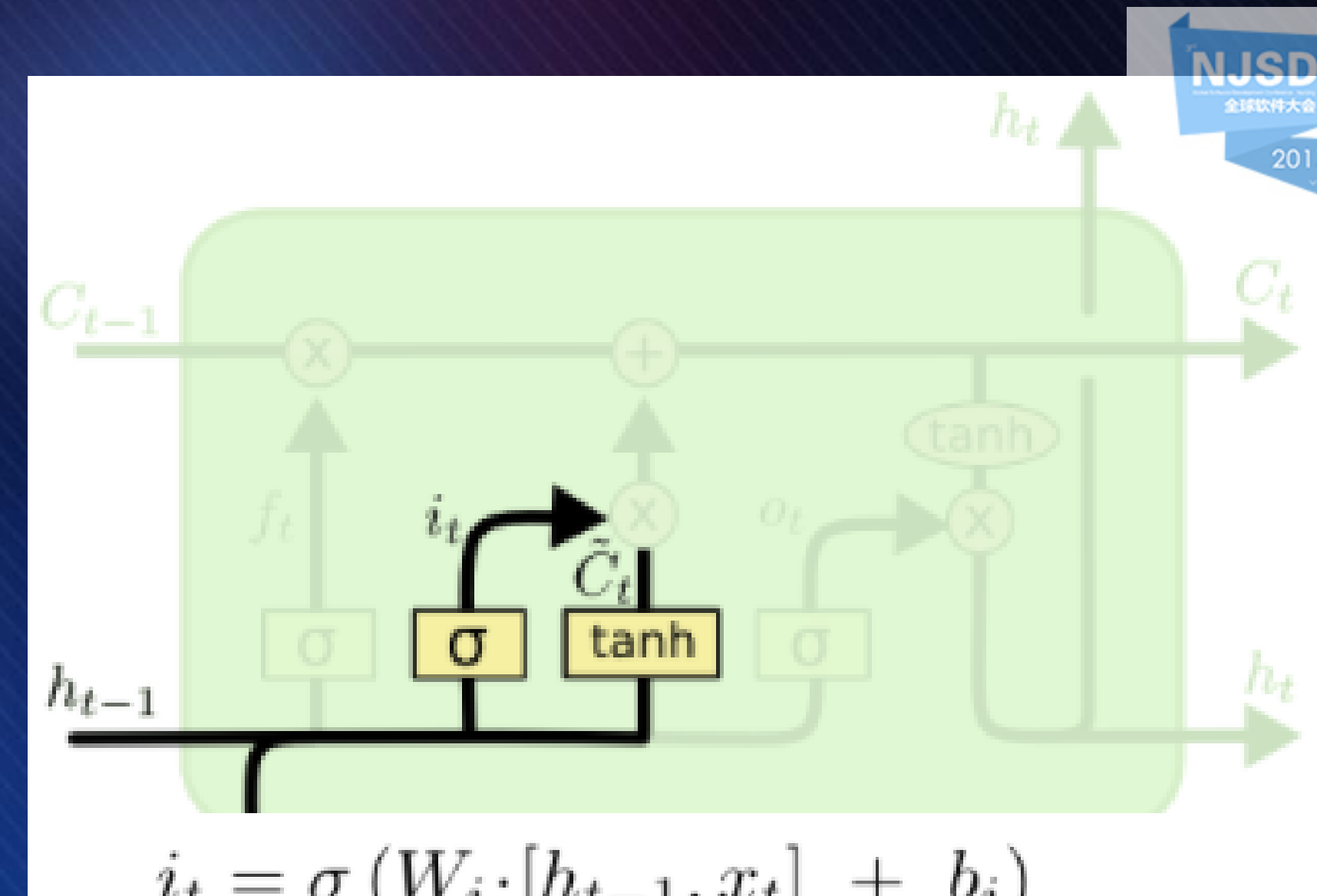
Long Short-Term Memory Units (LSTMs)

- A special kind of RNN have a better memory



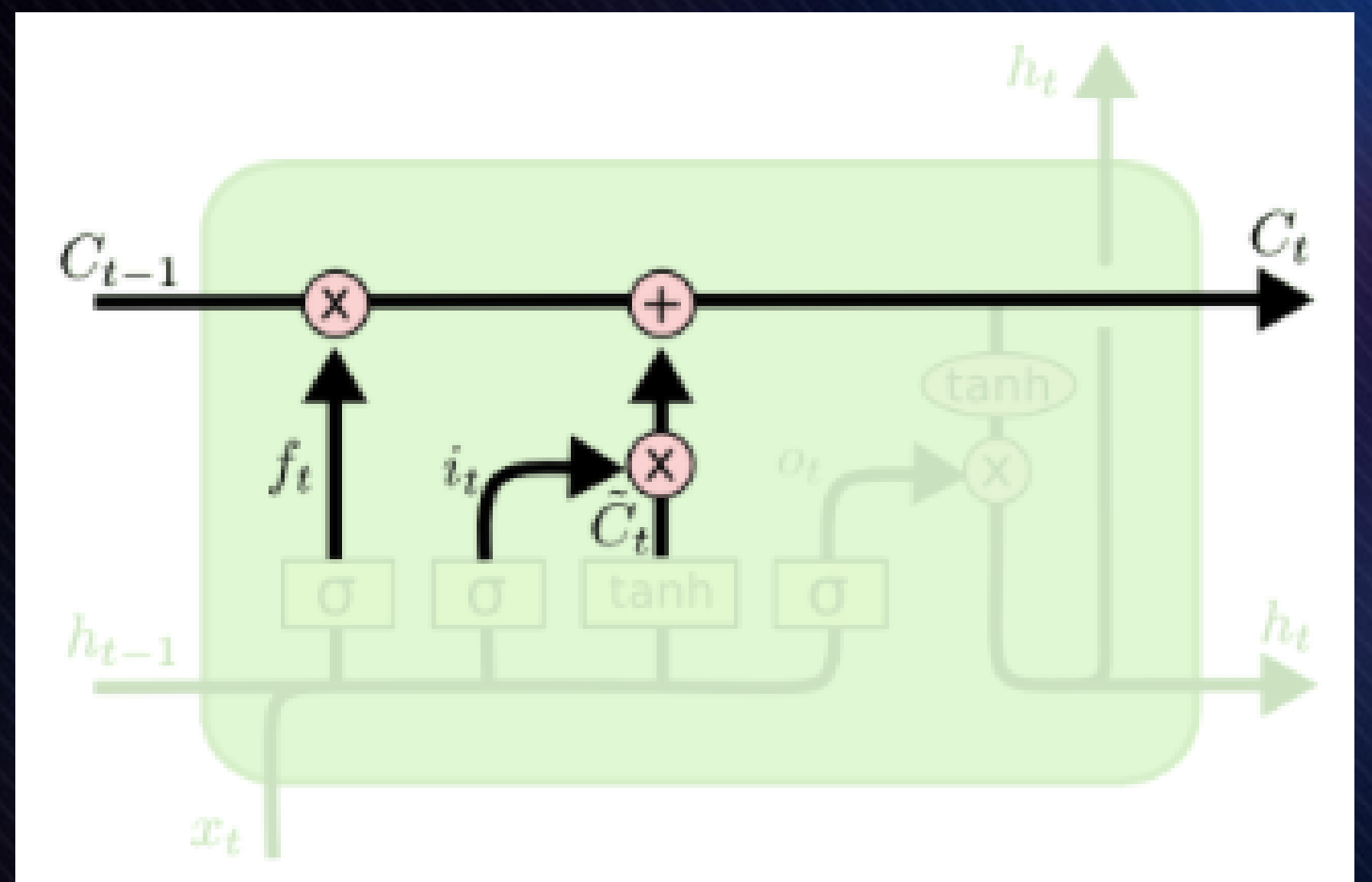


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

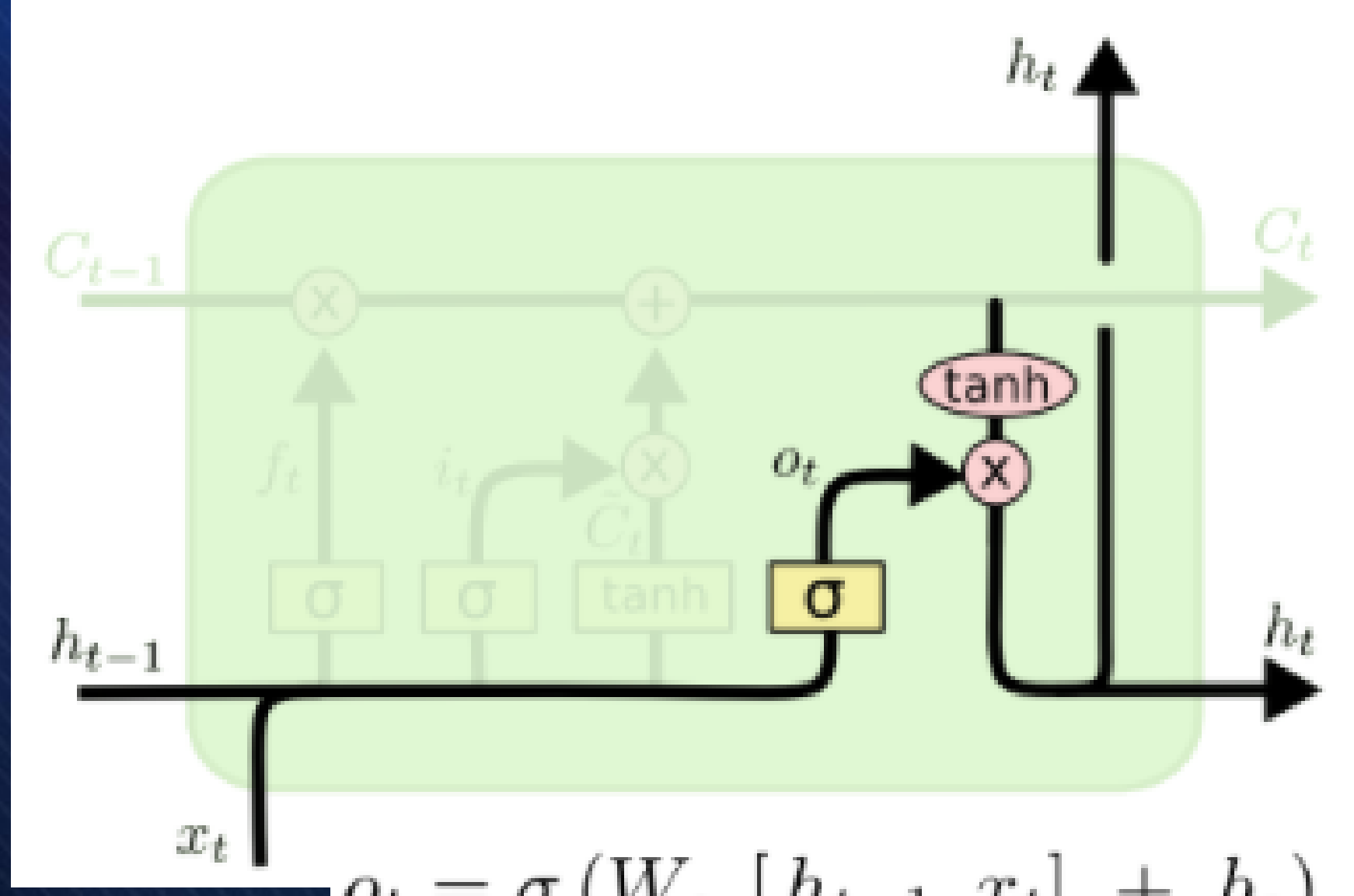


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



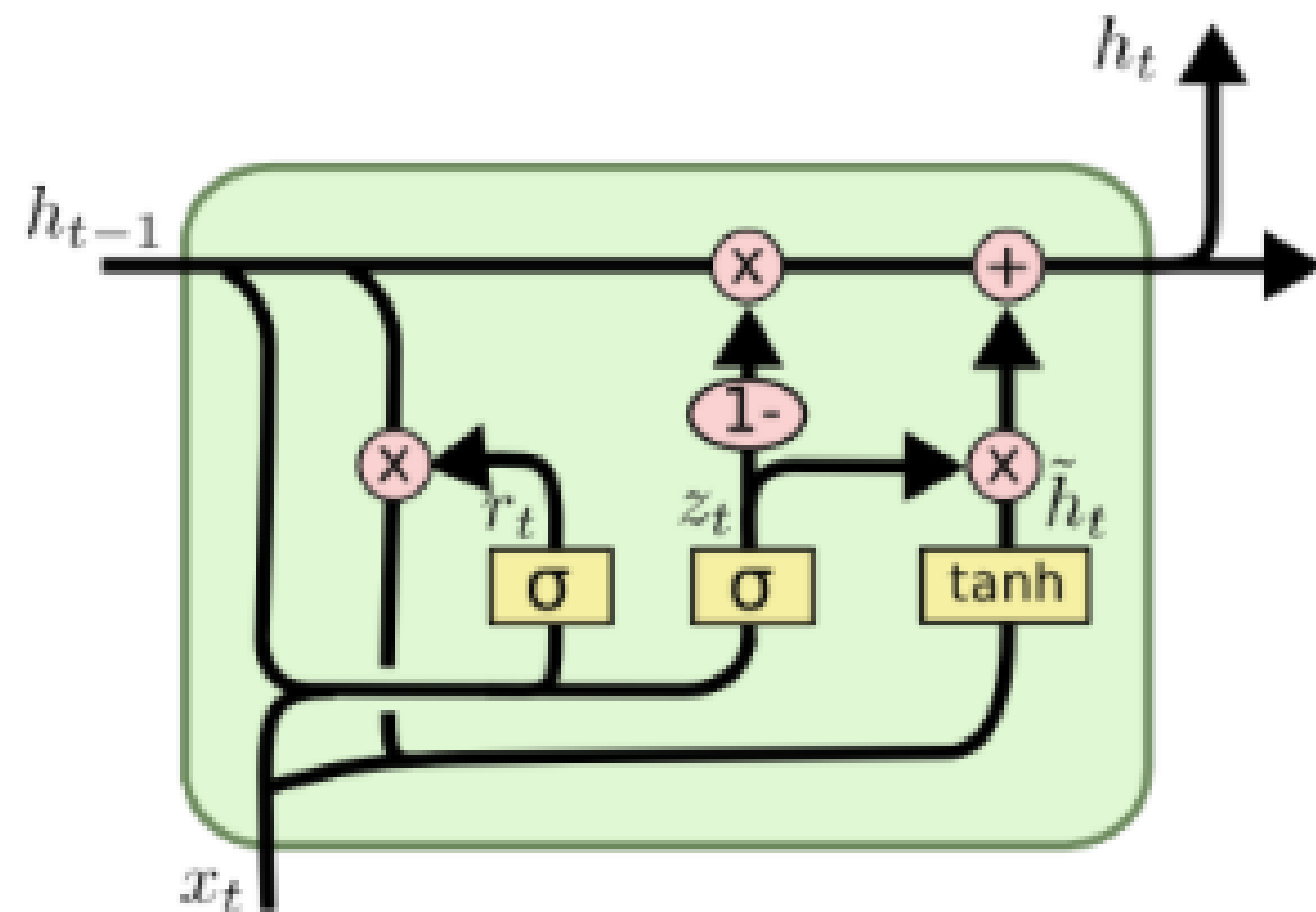
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Gated Recurrent Unit



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

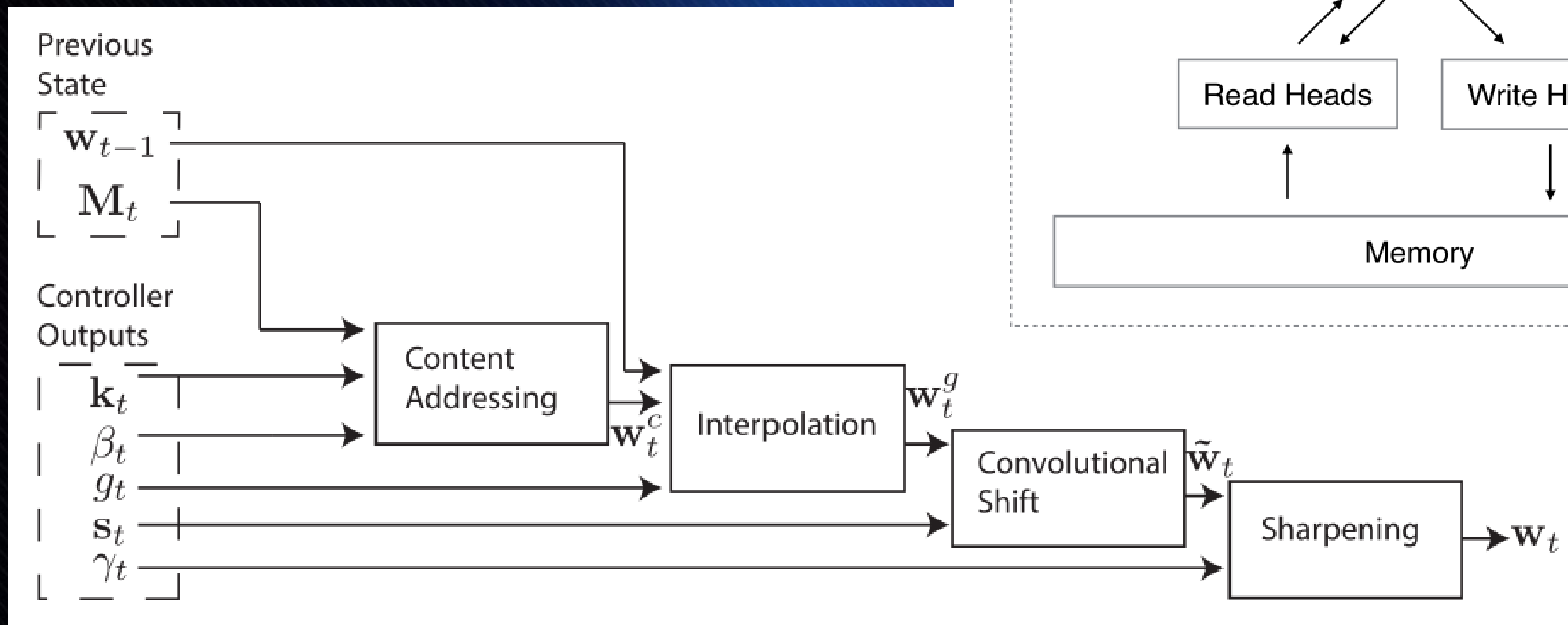
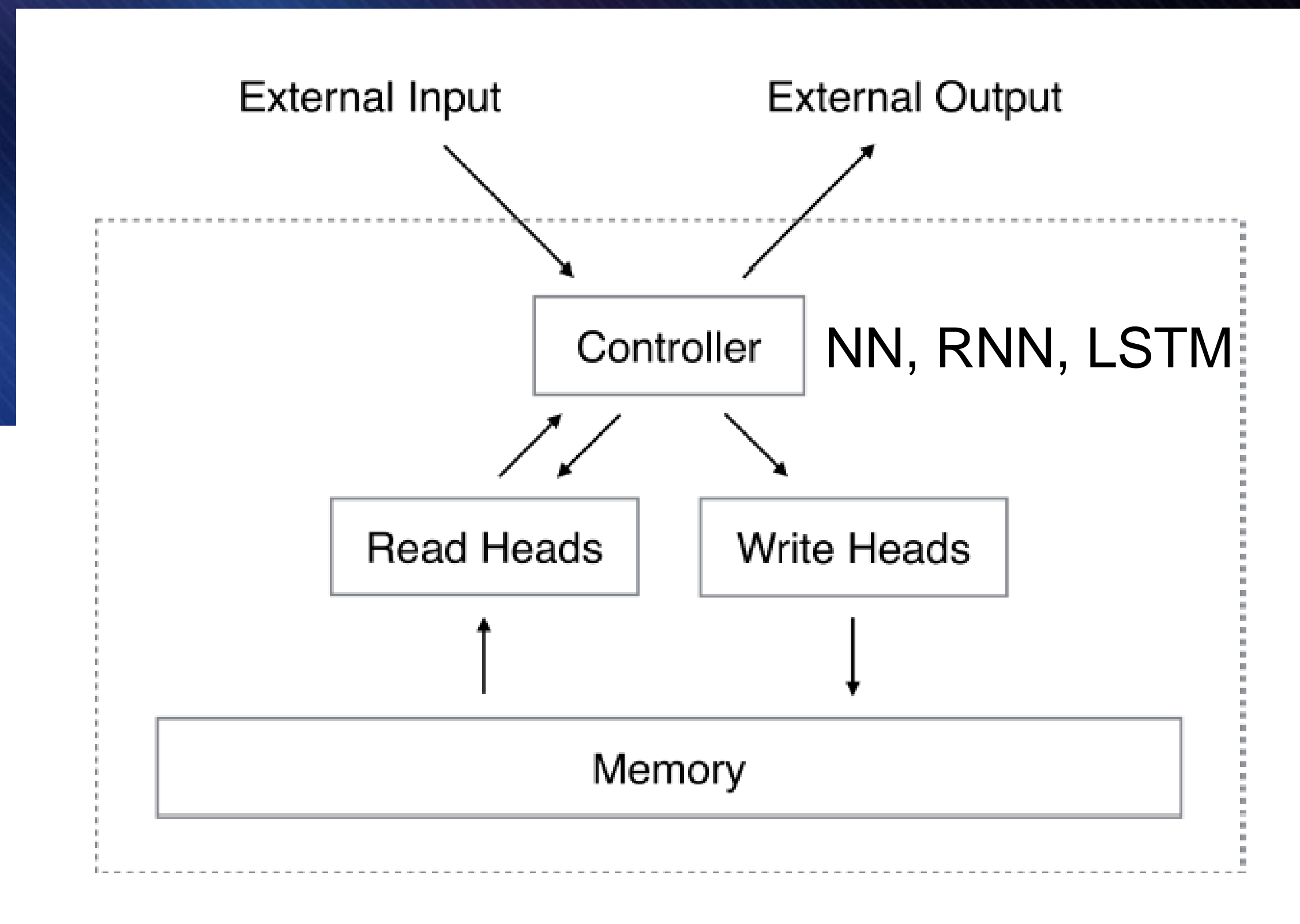
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

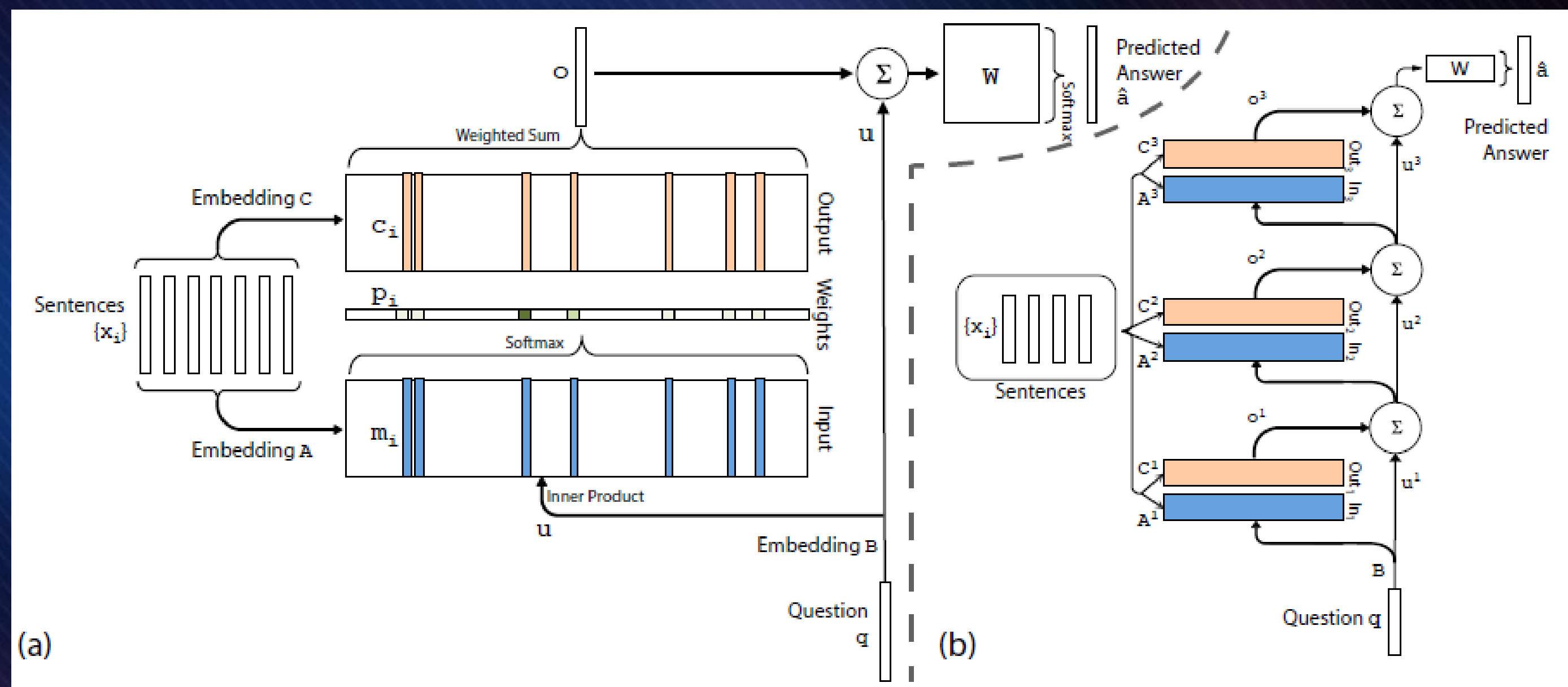
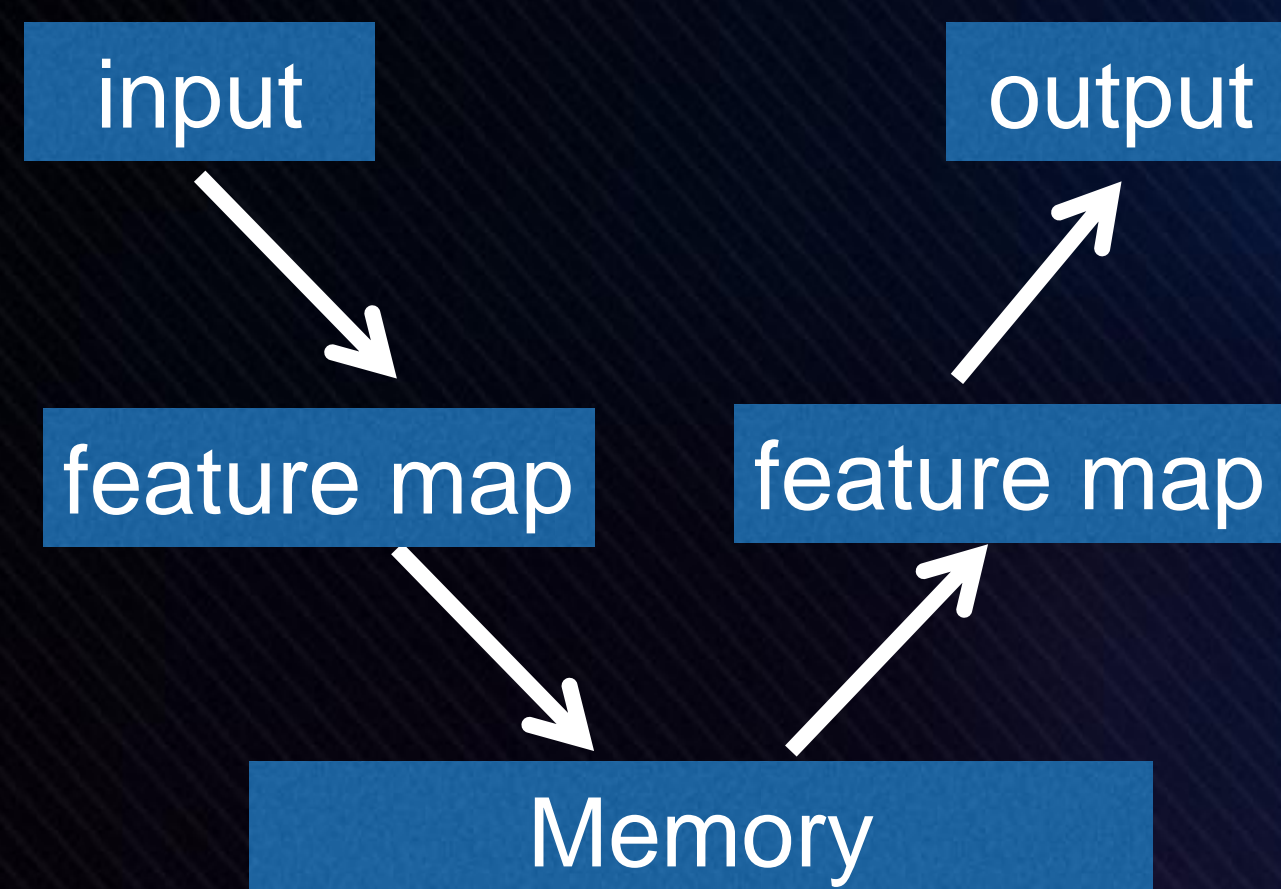
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Tape-like memory

Learn to copy, sort, etc.



Associative long-term memory with reasoning



Joe went to the kitchen. Fred went to the kitchen. Joe picked up the milk.
 Joe travelled to the office. Joe left the milk. Joe went to the bathroom.
 Where is the milk now? **A: office**
 Where is Joe? **A: bathroom**
 Where was Joe before the office? **A: kitchen**

Memory networks, J. Weston, S. Chopra, A. Bordes, 2014

End-to-End memory networks, S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus, 2015

Fixed point neural networks

- Inference of NNs works well on fixed point number
 - 8-bit networks converted from pre-trained float nets without sacrificing accuracy
 - Training from scratch will enable extreme low bit (1-bit) precision networks
- Fixed point are hardware friendly
 - Reduce model size and increase power efficiency
 - Speed-up inference

Binary neural networks

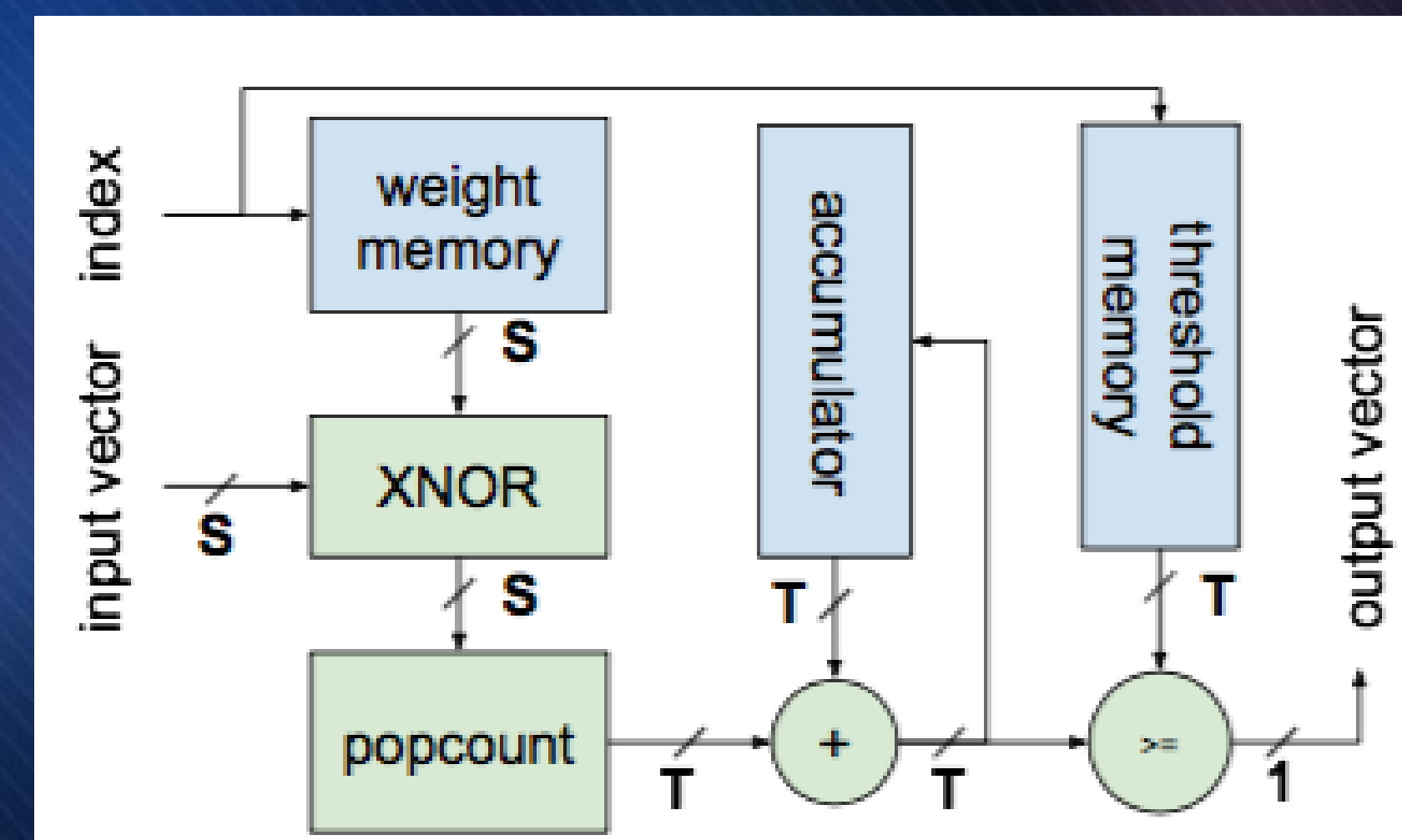
- Binary weights and activations
 - Enable working on large images and keep power efficiency
 - Replace MACC with XNOR and bit-count
- Less unique filters allow further optimization
- BNNs achieved promising results on small datasets

	Main arithmetic operation	2D Unique filters	Memory Saving (inference)	Accuracy on Cifar10
Float Neural Nets	float MAC	100%	1x	89.1%
Binary Neural Nets	XNOR-bitcount	~40%	~32x	88.8%

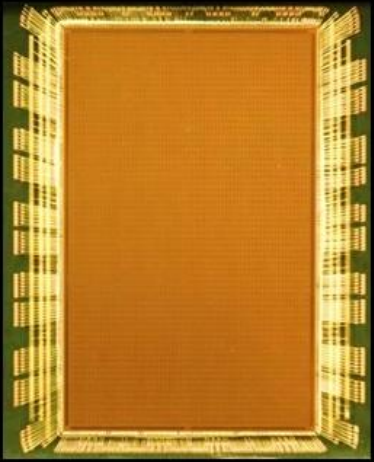
Binary neural networks

- Some tricks for training:
 - Put batch norm everywhere
 - Avoid network bottle-neck
 - Be patient with the slow converge
- Improve BNNs performances by
 - using more parameters
 - allowing extra high bit layer
 - working with other high bit networks

Typical convolutional layer in BNNs
 MACC -> XNOR-popcount
 Batch norm and binarization merged into a threshold operation

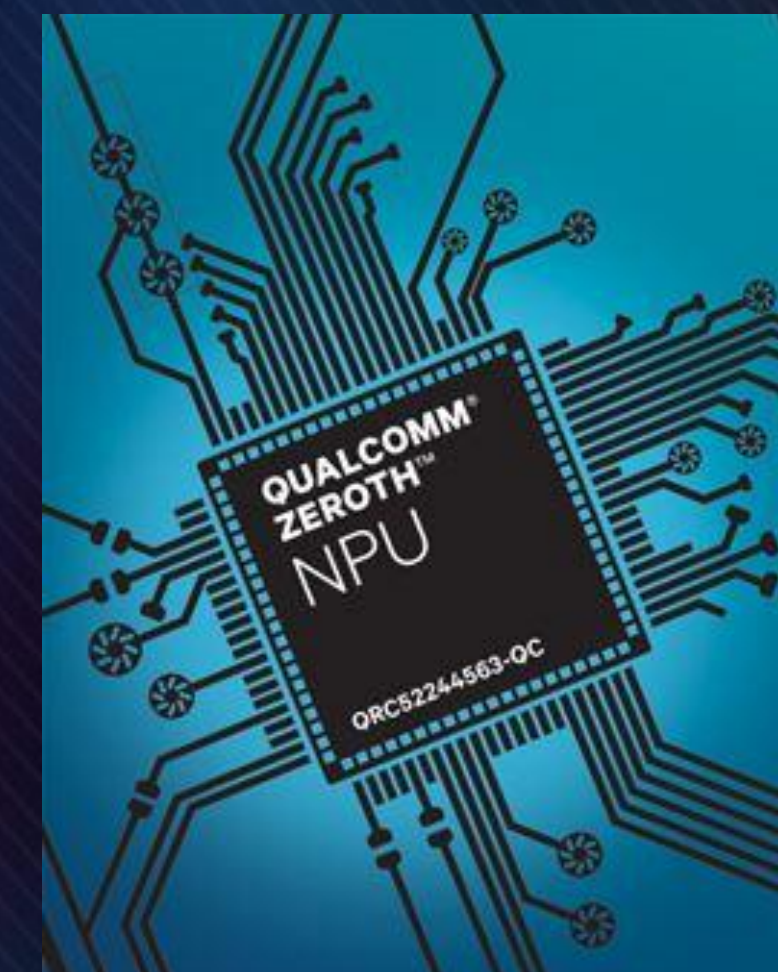
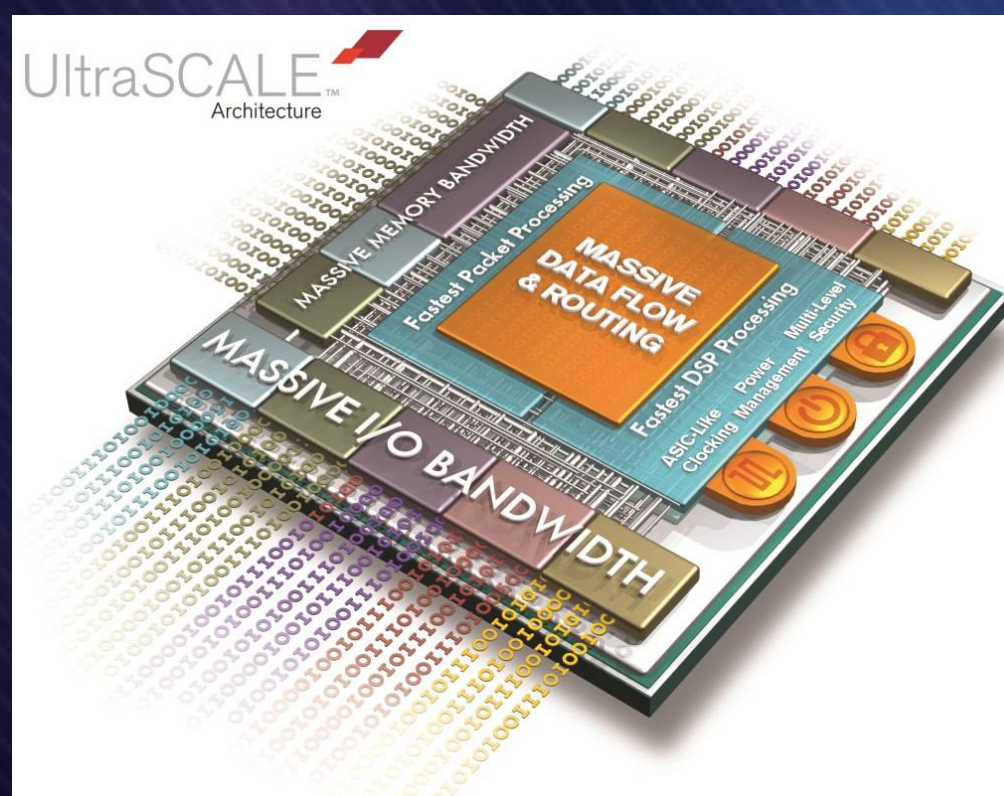


Umuroglu et. al. "FINN: A Framework for Fast, Scalable Binarized Neural Network Inference"



1 M Neurons
256 M Synapses
5.4 B Transistors
Realtime
73 mW

TrueNorth



The Mission of Horizon Robotics

Create the “brain” platform of smart things, to make human life more convenient, safer, and more fun.

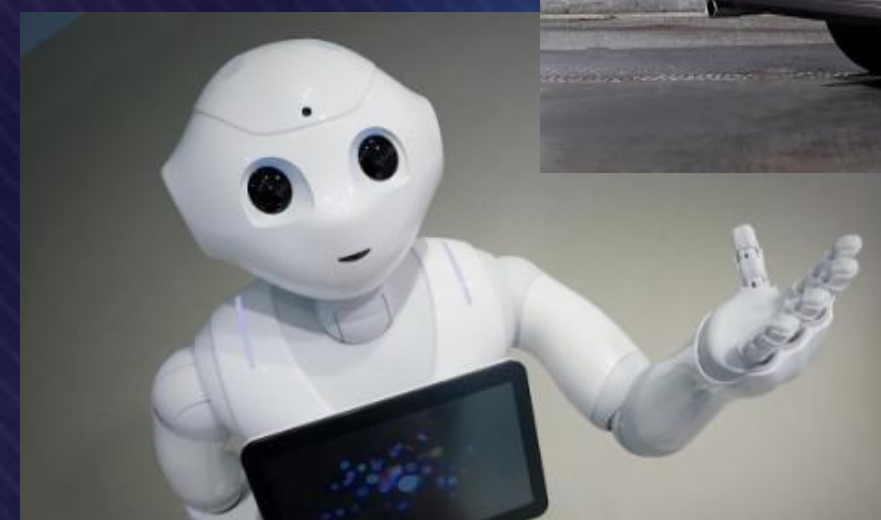
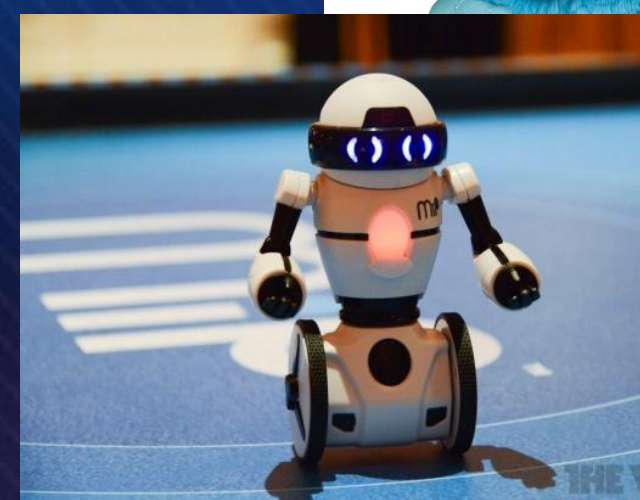
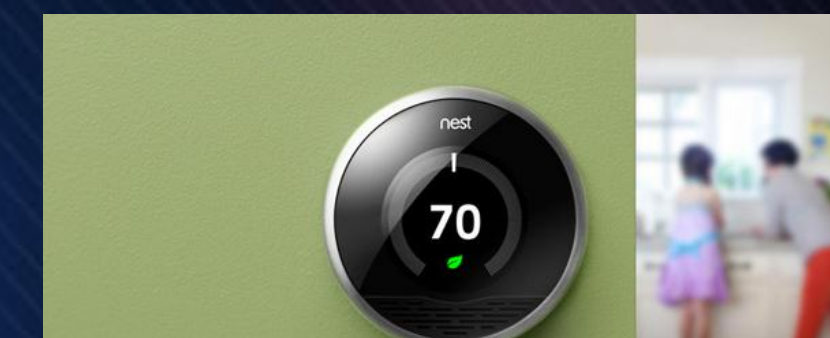
- Founded in July 14th, 2015, HQ in **Beijing**, and R&D site in **Nanjing**, an office in **Shenzhen** and **Shanghai**
- Experienced engineers from Baidu, Facebook, Google, Huawei, Nokia, Microsoft, TI, NVIDIA
- 40% have oversea experiences, 14% have PhD degrees, 100% are seasoned engineers
- Pioneers of many accomplishments in AI



Nanjing
office



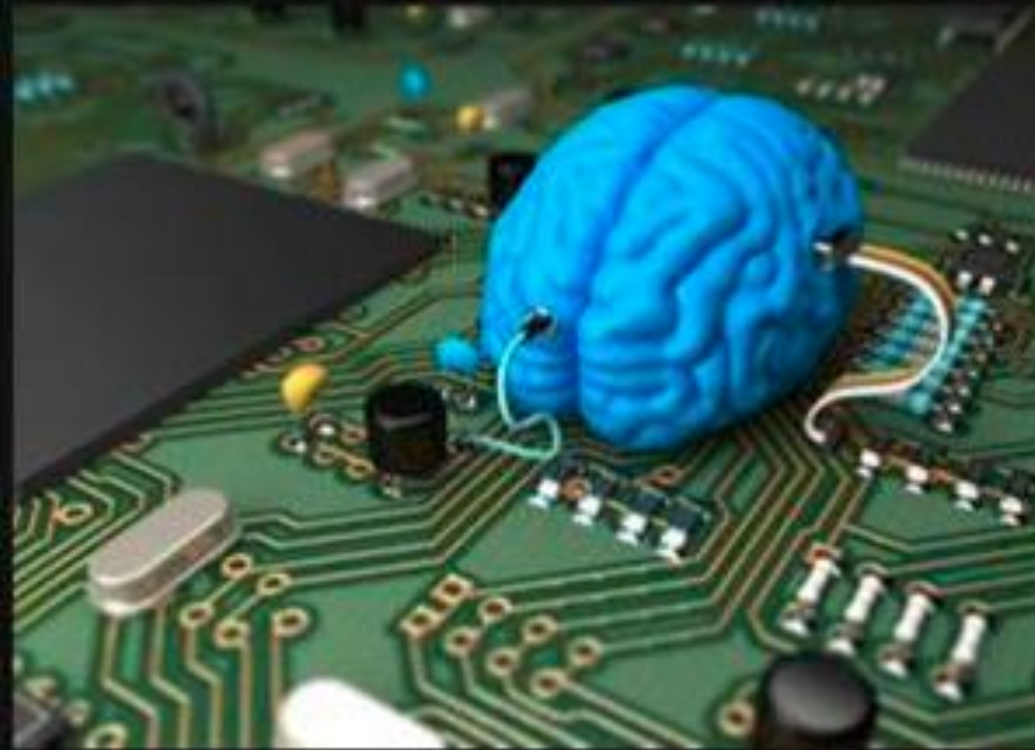
- In the future, all the devices are not only connected, but more importantly, “AI inside”.
- **sense the environment**
- **interact with people**
- **make control decisions**
- A local brain on the device
- **Perception & HCI**
- **Low-latency & real-time**
- **Low-power & low cost**
- **Privacy protection**



Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
from theano import *
```

What I actually do

Copied from Dr. Jiang Wang