

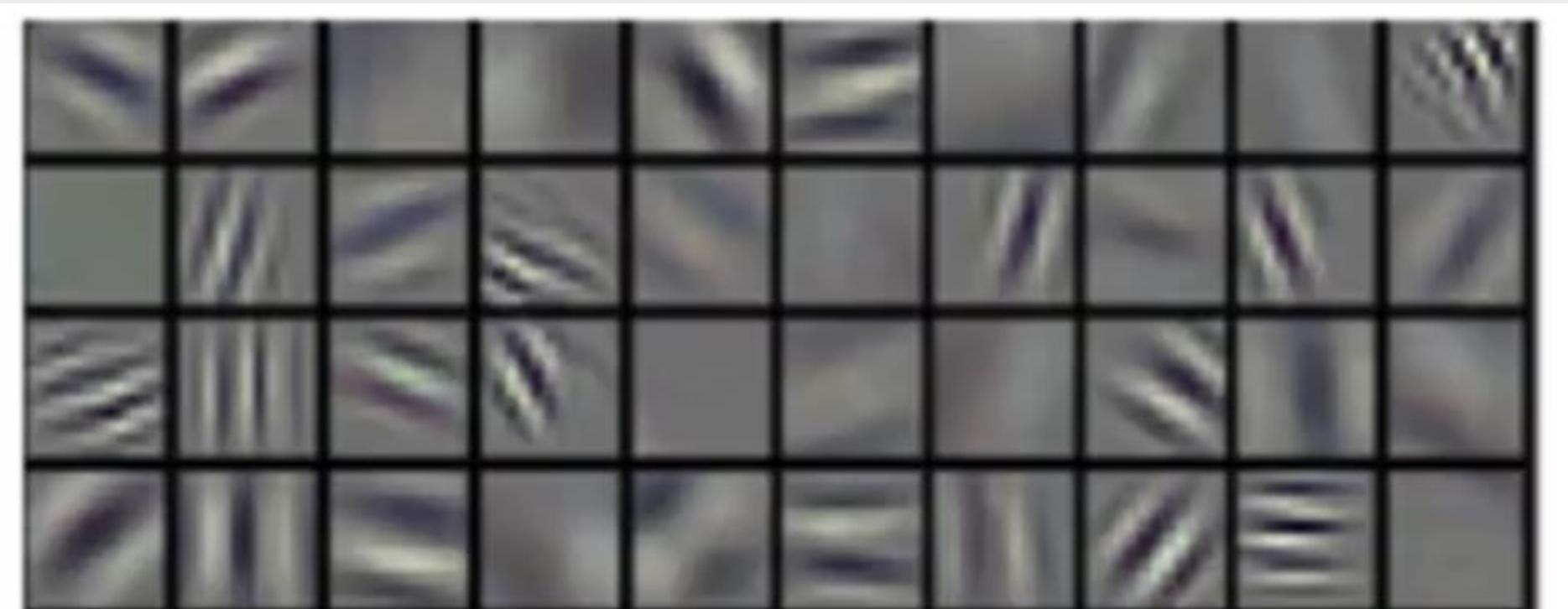
TRANSFER LEARNING

¿Qué es?

Reutilizar un modelo preentrenado para resolver una nueva tarea relacionada, con poca o ninguna necesidad de entrenar desde cero.

Motivación

- Las primeras capas de una red neuronal son las más difíciles (es decir, las más lentas) de entrenar. Debido al problema de Vanishing Gradient
- Sin embargo, estas características "primitivas" deberían ser comunes en muchas tareas de clasificación de imágenes. Las capas finales son más específicas al problema



Mantén congeladas las capas iniciales y re-
entrena las últimas

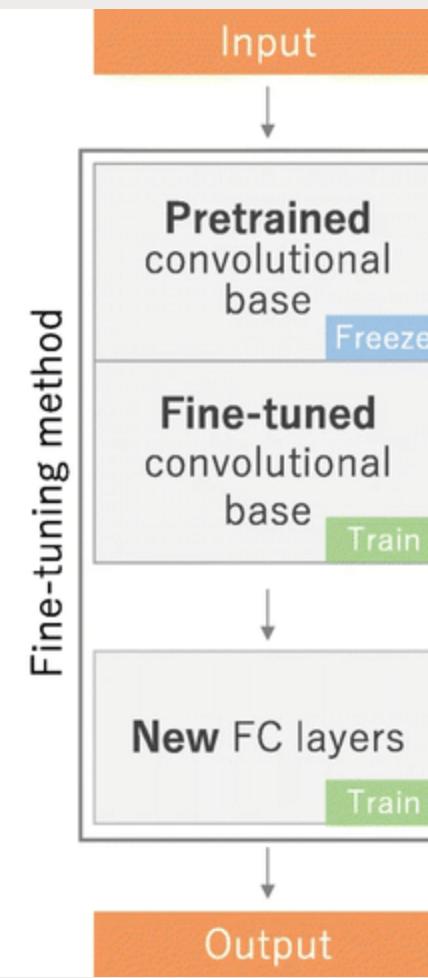
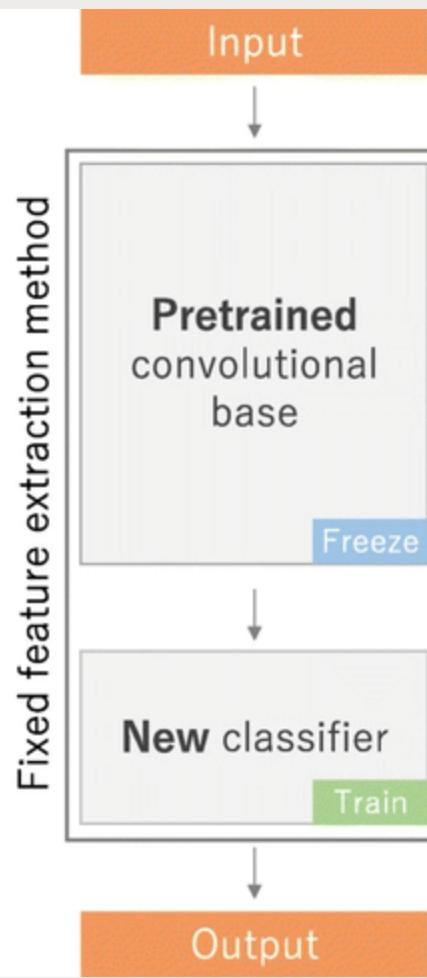
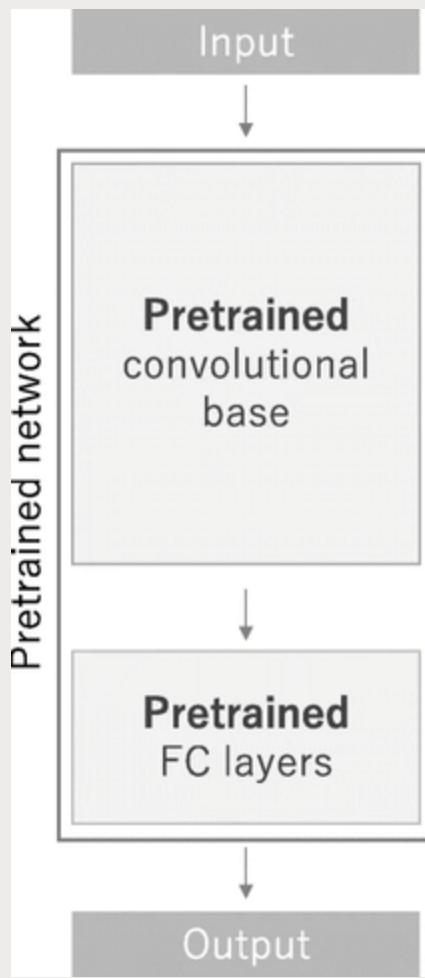
¿Cuándo usar Transfer Learning?

- Cuando no tienes suficiente datos para entrenar desde cero.
- Cuando tu tarea está relacionada visualmente con la tarea original.
- Cuando quieres ahorrar tiempo computacional.

Estrategias

1. Feature Extraction (Extracción de Características)
2. Fine-Tuning (Ajuste fino del modelo)
3. Full Model Fine-Tuning (Reentrenamiento total)

*Depende de



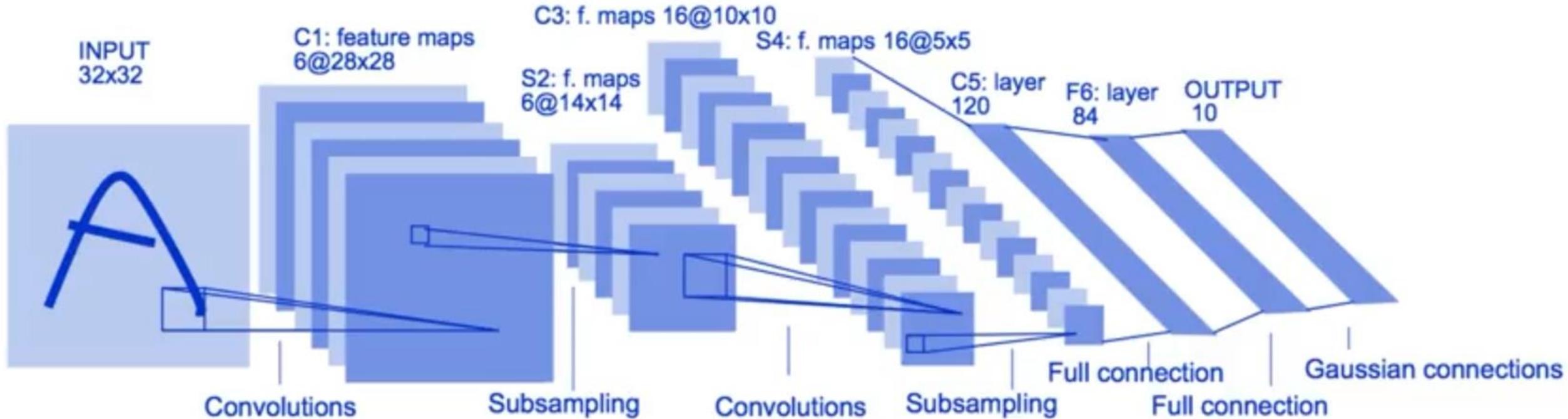
Estrategia	¿Congela capas?	¿Riesgo de overfitting?	¿Necesita muchos datos?	¿Cuándo usarla?
Feature Extraction	Sí	Bajo	No	Dataset similar, pocos datos
Fine-Tuning Parcial	Parcial	Medio	Moderado	Dataset algo distinto
Fine-Tuning Total	No	Alto	Sí	Dataset grande, tarea diferente

¿Qué modelos hay pre-entrenados?

- LeNet
- AlexNet (y el nacimiento de los embeddings)
- VGG
- Inception
- ResNet

LENET

Creada por Yann LeCun en 1998,
para el dataset de MNIST



Kernel 5x5 stride 1

ALEXNET (2012)

El modelo que
mató al Machine
Learning Clásico

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.toronto.edu

Ilya Sutskever
University of Toronto
lysotsk@cs.toronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.toronto.ca

Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest using the 1000 different classes. On the test set, we achieved top-1 and top-5 error rates of 15.3% and 7.5%, which are 17.0% and 15.8% better than the second-best submission. The network, which has 60 million parameters and 650,000 neurons, consists of three convolutional layers, some of which are followed by max-pooling layers, and two fully-connected layers, with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operations. To reduce memory in the fully-connected layers we employed a recently-developed regularization method called “drop-out” that is able to be efficiently implemented in our model. In the model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

1 Introduction

Current approaches to object recognition make essential use of machine learning methods. To improve their performance, we can collect larger datasets, learn more powerful models, and use better techniques for preventing overfitting. Until recently, datasets of labeled images were relatively small and the best performing models had less than 100,000 parameters. Now, with ImageNet [16], CIFAR-10/100 [12], Simple recognition tasks can be solved quite well with datasets of this size, especially if they are augmented with label-preserving transformations. For example, the current best error rate on the MNIST digit-recognition task (0.3%) approaches human performance [4]. But as datasets grow, so do the challenges. We find that it is increasingly difficult to learn features necessary to use much larger training sets. And indeed, the shortcomings of small image datasets have been well recognized (e.g., Pinto et al. [21]), but it has only recently become possible to collect large datasets of images at high resolution. The new largest datasets include LabelMe [42], which consists of over 15 million labeled images of thousands of fully-segmented objects, and ImageNet [48], which consists of over 15 million labeled high-resolution images in over 22,000 categories.

To learn about thousands of objects from millions of images, we need a model with a large learning capacity. However, the inherent complexity of the object recognition task means that this problem cannot be solved even by a model as large as ImageNet. As our model should also have lots of prior knowledge to correctly identify objects, we must add layers of pre-training. Our network consists of four convolutional layers, each learned from scratch, and two fully-connected layers, each learned from a combination of pre-training and fine-tuning. These components to standard feed-forward neural networks with fully-connected layers, CNNs have much fewer connections and parameters and so they are easier to train, while their theoretically-best performance is likely to be only slightly worse.

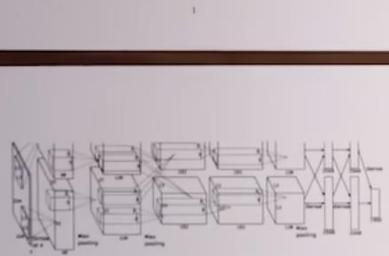


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer parts at the top of the figure while the other runs the layer parts at the bottom. The GPUs communicate only at certain layers. The network’s input is 100Mx256-dimensional, and the total number of neurons in the network (excluding layers) is given by $352,460 + 186,624 + 64,896 + 64,896 + 4096 = 4096 \times 4096 = 16M$.

4 Reducing Overfitting

Our neural network architecture has 60 million parameters. Although the 1000 classes of ILSVRC make training examples impose 100 bits of constraints on the mapping from image to label, this may not be enough to constrain so many parameters without considerably overfitting. Below, we describe the two primary ways in which we combat overfitting.

4.1 Data Augmentation

The easiest and most common method to reduce overfitting on image data is to artificially enlarge the dataset using label-preserving transformations (e.g., [25, 44]). We employ two distinct forms of data augmentation, both of which alter transformed images in different ways. The first form of data augmentation very little, so that the transformed images do not need to be stored on disk. In our implementation, the transformed images are generated in Python code on the CPU while the GPU is training on the preprocessed batch of images. So these data augmentation schemes are, in effect, computed on the fly.

The first form of data augmentation consists of generating translations and horizontal reflections. We do this by extracting random 224 × 224 patches (and their horizontal reflections) from the 256 × 256 images and training our network on these extracted patches. This increases the size of our training set by a factor of 2048, though the resulting training examples are, of course, highly degenerate. The second form of data augmentation is to store the transformed images and have them loaded on to use much smaller networks. At test time, the network makes a prediction by extracting five 224 × 224 patches (the four corner patches and the center patch), as well as their horizontal reflections, and averaging the predictions made by the network’s softmax layer on the test patches.

The second form of data augmentation consists of altering the intensities of the RGB channels in training images. Specifically, we perform PCA on the set of RGB pixel values to extract the ImageNet principal component. This was to do with the fact that the validation error rate improved with the current learning rate. The learning rate was initialized at 0.01 and

Despite the attractive qualities of CNNs, and despite the relative efficiency of their local architecture, they have still been prohibitively expensive to apply in large scale to high-resolution images. Luckily, current GPUs, paired with a highly-optimized implementation of 2D convolution, are powerful enough to handle such large networks. In fact, the ImageNet competition uses such architectures as ImageNet to contain enough labeled examples to train such a model without severe overfitting.

The specific contributions of this paper are as follows: we trained one of the largest convolutional neural networks to date on the subsets of ImageNet used in the ILSVRC-2010 and ILSVRC-2012 competitions [2] and achieved by far the best results ever reported on these datasets. We wrote a highly-optimized GPU implementation of 2D convolution and all the other components required to train convolutional neural networks, where we make extensive use of ReLU units. Our network contains a number of new and unusual features which improve its performance and reduce its training time, which are detailed in Section 3. The size of our network made overfitting a significant problem, even with 2048 training examples. To combat this, we used data augmentation, label-preserving cropping, overfitting, which are described in Section 4. Our final network contains five convolutional and three fully-connected layers, and this depth seems to be optimal: we found that removing any successive layer of reach of which contain no more than 1% of the model’s parameters) will result in inferior performance.

In the end, the network’s size is limited mainly by the amount of memory available on current GPUs and by the amount of training time that we are willing to tolerate. Our network takes between five and six days to train on two GTX 580 4GB GPUs. All of our experiments suggest that our results can be improved simply by waiting for faster GPUs and bigger datasets to become available.

2 The Dataset

ImageNet is a dataset of over 15 million labeled high-resolution images belonging to roughly 22,000 categories. The images were collected from the web and labeled by human labelers using Amazon’s Mechanical Turk crowd-sourcing tool. Starting in 2010, as part of the PASCAL Visual Object Challenge (VOC) and the Microsoft Research ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) has been held. ILSVRC uses a subset of ImageNet containing roughly 1000 categories and 1000 images. In all, there are roughly 1.2 million training images, 50,000 validation images, and 150,000 testing images.

ILSVRC-2010 is the only version of ILSVRC for which the test set labels are available, so this is the version on which we performed most of our experiments. Since we trained our model in the ILSVRC-2010 contest, in this section we report our results on this test set. The validation set, for which test set labels are unavailable, on ImageNet, it is customary to report two error rates: top-1 and top-5, where the top-5 error rate is the fraction of test images for which the predicted label is not one of the five labels corresponding to the true label.

ImageNet contains no absolute-resolution images, which our network requires a constant input dimensionality. Now we downsampled the images to a fixed resolution of 256 × 256. Given a rectangular image, we first rescaled the image so that the shorter side was of length 256, and then cropped out the central 256 × 256 patch from the resulting image. We did not pre-process the images in any other way, except for subtracting the mean activity over the training set from each pixel. So we trained our network on the (centered) raw RGB values of the pixels.

3 The Architecture

The architecture of our network is summarized in Figure 2. It contains eight learned layers — five convolutional and three fully-connected. Below, we describe some of the novel or unusual features of our network’s architecture. Sections 3.1–3.4 are sorted according to our estimation of their importance, with the most important first.

<http://code.google.com/p/cuda-convnet/>

1

with magnitudes proportional to the corresponding eigenvalues times a random variable drawn from a Gaussian with zero mean and standard deviation 0.1. Therefore to each RGB image pixel $I_{ijk} = [I_{ijk}, I_{ijk}^2, I_{ijk}^3]^T$ we add the following quantity:

$$[\delta p_i, \delta q_i, \delta u_i, \delta v_i, \delta z_i]^T$$

where p_i and q_i are 6th eigenvector and eigenvector of the 3×3 covariance matrix of RGB pixel values, respectively, and δu_i is the aforementioned random variable. Each δz_i is drawn only once for each pixel of a particular training image until that image is used for training again, at which point it is drawn again. This technique captures an important property of natural images, namely, that object identity is invariant to changes in the intensity and color of the illumination. This scheme reduces the top-1 error rate by over 1%.

4.2 Dropout

Combining the predictions of many different models is a very successful way to reduce test errors [11, 13]. It appears to be less effective for big neural networks that already take several days to train. However, it is very effective for the type of convolutional layers that we use, which shows a factor of two during training. The recently-introduced technique, called “drop-out” [10], consists of setting to zero the output of each hidden neuron with probability 0.5. The neurons which are dropped out are chosen independently at random. This is done before every step of backpropagation. So every time an input is presented, the neural network samples a different architecture, but all these architectures share weights. This technique reduces complex co-adaptations of neurons, since a neuron cannot rely on the presence of particular other neurons. It is therefore, forced to learn to identify objects based on different features, and to learn different parts of the other neurons. At test time, we use all the neurons but multiply their outputs by 0.5, which is a reasonable approximation to taking the geometric mean of the predictive distribution produced by the exponentially-many dropout networks.

We use dropout in the first two fully-connected layers of Figure 2. Without dropout, our network exhibits substantial overfitting. Dropout roughly doubles the number of iterations required to converge.

5 Details of Learning

We trained our models using stochastic gradient descent with a batch size of 128 examples, momentum of 0.9, and weight decay of 0.0005. We found that this small amount of weight decay was important for the model to learn. In other words, the learned weights in the network do not need to be stored on disk. In our implementation, the transformed images are generated in Python code on the CPU while the GPU is training on the preprocessed batch of images. So these data augmentation schemes are, in effect, computed on the fly.

$$\nabla_{\theta+1} := (1 - \eta_k) \nabla_{\theta} - \eta_k \nabla_{\theta} - \epsilon \left(\frac{\partial L}{\partial \theta} \right)_{\theta=\theta_k}$$

where ϵ is the iteration index, v is the momentum variable, ϵ is the learning rate, and $\left(\frac{\partial L}{\partial \theta} \right)_{\theta=\theta_k}$ is the average over the k th batch D_k of the derivative of the objective with respect to v , evaluated at θ_k .

We initialized the weights in each layer from a zero-mean Gaussian distribution with standard deviation of 0.01. We initialized the biases in each layer from a uniform distribution with standard deviation of 0.01. We initialized the weights in the first two convolutional layers, as well as in the fully-connected layers, with the constant 1. This initialization accelerates the early stages of learning by providing the ReLU with positive inputs. We initialized the neurons in the remaining layers to the constant 0.

The first form of data augmentation consists of altering the intensities of the RGB channels in training images. Specifically, we perform PCA on the set of RGB pixel values to extract the ImageNet principal component. This was to do with the fact that the validation error rate improved with the current learning rate. The learning rate was initialized at 0.01 and

3.1 ReLU Nonlinearity

The standard way to model a neuron’s output f as a function of its input x is with $f(x) = \text{tanh}(x)$. Or $f(x) = x$ for $x > 0$ and $f(x) = 0$ for $x \leq 0$. Using gradients, these saturating nonlinearities are much slower than the non-saturating nonlinearity $f(x) = \max(0, x)$. In fact, we refer to neurons with this nonlinearity as Rectified Linear Units (ReLU). Deep convolutional neural networks with ReLU train several times faster than their counterparts with tanh [11]. Our network contains a number of new and unusual features which improve its performance and reduce its training time, which are detailed in Section 3. The size of our network made overfitting a significant problem, even with 2048 training examples. To combat this, we used data augmentation, label-preserving cropping, overfitting, which are described in Section 4. Our final network contains five convolutional and three fully-connected layers, and this depth seems to be optimal: we found that removing any successive layer of which contain no more than 1% of the model’s parameters) will result in inferior performance.

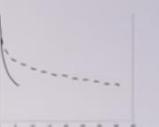


Figure 1: A four-layer convolutional neural network diagram showing the flow of data from input to output. The input is a 256x256 image. It passes through four layers: Layer 1 (5x5 kernel, stride 2), Layer 2 (5x5 kernel, stride 2), Layer 3 (3x3 kernel, stride 1), and Layer 4 (1x1 kernel, stride 1). The final output is a 1000-class softmax layer. The diagram illustrates the spatial dimensions of the feature maps at each layer.

3.2 Training on Multiple GPUs

A single GTX 580 GPU has only 4GB of memory, which limits the maximum size of the networks that can be trained on it. It turns out that 1.2 million training images are enough to train networks which are too big to fit on one GPU. Therefore we spread the net across two GPUs. Current GPUs are particularly well-suited to cross-GPU parallelization, as they are able to read from and write to one another’s memory directly, without going through host machine memory. The parallelization works as follows: we split the training set into two sets of 600,000 images each, and we use additional GPU’s communicate only in certain layers. This means that, for example, the kernels of layer 3 take input from all kernel maps in layer 2. However, kernels in layer 4 take input from only one kernel map in layer 3, which is why the patterns of connectivity in a kernel map are preserved. This is important for maintaining the patterns of connectivity in a kernel map after normalization, but this allows us to precisely tune the amount of computation used in it to an acceptable fraction of the amount of computation.

The resultant architecture is somewhat similar to that of the “columnar” CNN employed by Corrigan et al. [5], except that our columns are not independent (see Figure 2). This scheme reduces our top-1 error rate by 1.1% and our top-5 error rate by 1.2%, respectively, as compared with a net with half as many kernels in each fully-connected layer trained on one GPU. The two-GPU net takes slightly less time to train than the one-GPU net.

¹The one-GPU net actually has the same number of kernels as the two-GPU net in the final convolutional layer. This is because most of the net’s parameters are in the first fully-connected layer, which takes the last convolutional layer as input. To make the two-GPU net work, we had to move the first fully-connected layer to the middle of the network. We did this by splitting the two-GPU net into two separate networks, which do not have access to the final convolutional layer (nor the fully-connected layers which follow it). Therefore this computation is biased in favor of the one-GPU net, since it is bigger than “half the size” of the two-GPU net.

The first convolutional layer filters the 224 × 224 × 3 input image with 96 kernels of size 11 × 11 × 3.

3

3.3 Local Response Normalization

ReLU has the desirable property that they do not require input normalization to prevent them from exploding. If at least one training example produces a positive input to a ReLU, learning will happen in that neuron. However, we still find that the following local normalization scheme aids generalization. Denoting by $\alpha_{x,y}$ the activity of a neuron composed by applying kernel i at position (x,y) and applying the ReLU nonlinearity, the response-normalized activity $\tilde{\alpha}_{x,y}$ is given by the expression

$$\tilde{\alpha}_{x,y} = \alpha_{x,y} \left(k + \alpha \frac{\sum_{n=1}^{N-1} \alpha_{x,y+n}}{\sum_{n=1}^{N-1} \alpha_{x,y+n}} \right)^{\beta}$$

where the sum ratio over N “adjacent” kernel maps at the same spatial position, and N is the total number of kernels in the layer. The ordering of the kernel maps is of course arbitrary and does not matter here. This scheme is inspired by the type found in real neurons, creating competition for big activities amongst neurons outputting using different kernels. The constants k , α , and β are hyperparameters. We chose values $\alpha = 0.001$ and $\beta = 0.75$. We also tried the normalization scheme proposed by Srivastava et al. [35], but ours would be more correctly termed “brightness normalizations”, since we do not subtract the mean activity. Response normalization reduces our top-1 and top-5 error rates by 1.4% and 1.2%, respectively. We also verified the effectiveness of this scheme on the CIFAR-10 dataset, a four-layer neural network trained with a different normalization scheme [35].

²Overlapping Pooling

Pooling layers in CNNs summarize the outputs of neighboring groups of neurons in the same spatial neighborhood. Traditionally, the neighborhoods summarized by adjacent pooling units do not overlap (e.g., [17, 11, 44]). To be more precise, a pooling layer can be thought of as computing a grid of pooling units, where the receptive field of a unit is the union of the receptive fields of the units in the same row and column of the pooling layer. If we set $s < r$, we obtain traditional local pooling as commonly employed in CNNs. If we set $s = r$ and $r = 3$, this scheme reduces the number of pooling units to $1/9$ of what it would be without overlapping pooling. This is what we use throughout our work, with $s = 2$ and $r = 3$. This scheme reduces the number of neurons and layers to 1/8 of what it would be without overlapping pooling. We generally observe during training that models with overlapping pooling find it slightly more difficult to overfit.

3.4 Overall Architecture

Now we are ready to describe the overall architecture of our CNN. As depicted in Figure 2, the net consists of eight layers which comprises the first five are convolutional and the remaining three are fully-connected. The output of the last fully-connected layer is fed to a 1000-way softmax which produces a distribution over the 1000 class labels. Our network maximizes the multinomial logistic regression objective, which is equivalent to minimizing the average across training cases of the log-probability of the correct class.

The kernels of the second, fourth, and fifth convolutional layers are connected to all kernel maps in the second layer. The neurons in the fully-connected layers are connected to all neurons in the previous layer. Max-pooling layers, of the kind described in Section 3.4, follow both response-normalization layers as well as the fifth convolutional layer. The ReLU non-linearity is applied to the output of every convolutional and fully-connected layer.

The first convolutional layer filters the 224 × 224 × 3 input image with 96 kernels of size 11 × 11 × 3. The stride of a 4 pixels is provided here. d is the distance between the receptive field centers of neighboring neurons. With a stride of 4 pixels, d is the distance between the receptive field centers of neighboring neurons.

³We can describe this network in detail due to space constraints, but it is specified precisely by the code and parameters file provided here: <http://code.google.com/p/cuda-convnet/>

4



Figure 4: (Left) Eight ILSVRC-2010 test images and the first label (cat). The first label is a dog, but it is the top 3. (Right) Five ILSVRC-2012 test images in the first position. The first label is a dog, but it is the top 3. (Right) Five ILSVRC-2012 test images in the first position. The remaining columns show the six training images that produce feature vectors in the last hidden layer with the smallest Euclidean distance from the feature vector for the test image.

In the left panel of Figure 4 we qualitatively assess what the network has learned by comparing its top-5 errors on four test sets. Nearest neighbors of the test set images, such as the ones in the first column, can be measured. The size of the test set’s nearest neighbor is, for example, only other types of cat are considered plausible labels for the test set. In some cases (grille, sheep) there is genuine ambiguity about the intended focus of the photograph.

Another way to probe the network’s visual knowledge is to consider the feature activations induced by an image at the last, 4096-dimensional hidden layer. If two images produce feature vectors which are similar, we can measure the size of the test set’s nearest neighbor. For example, the first two images in the first column are more similar to each other than to the rest of the test set. Notice that at the global level, the network is able to distinguish between different types of dogs and cats. For example, the retrieved dogs and cats appear in a variety of poses. We present the results for many more test images in the supplementary material.

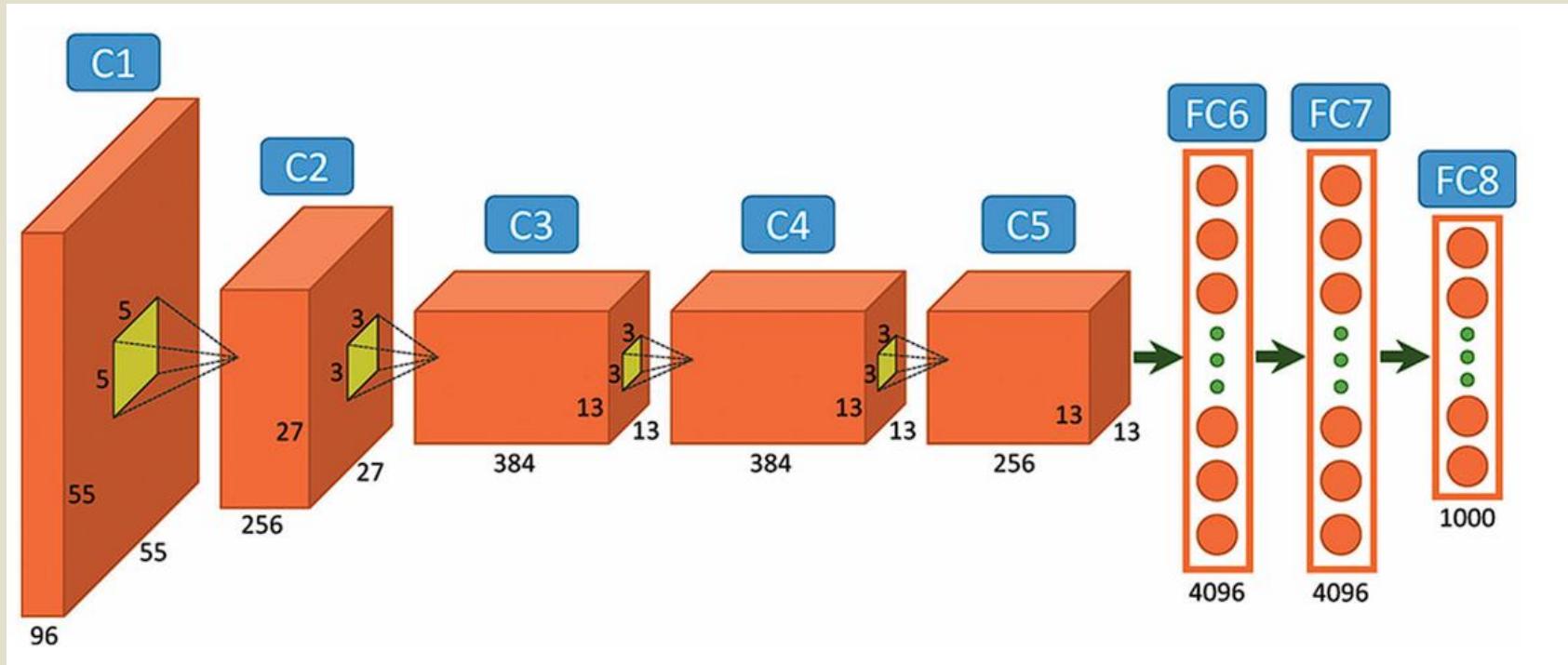
Comparing similarity by using Euclidean distance between raw 4096-dimensional real-valued vectors is inefficient, but it could be made efficient by using an auto-encoder to compute latent vectors. For example, we can use a neural network to map the raw 4096-dimensional vectors to a 100-dimensional latent space. This kind of representation makes it easier to compare two images. For example, the network can be trained to map images to a latent space where images with similar patterns of edges, whether or not they are semantically similar, are close together.

⁴Discussion

Our results show that a large, deep convolutional neural network is capable of achieving record-breaking results on a highly challenging task using purely supervised learning. In particular, removing any of the middle layers results in a loss of about 2% for the top-1 performance of the network. So the depth really is important for achieving our results.

To simplify our experiments, we did not use any unsupervised pre-training even though we expect that it will help, especially if we obtain enough unlabeled training data. In fact, we did not even consider unsupervised pre-training in our new paper [4], which does not make use of image labels and hence has a tendency to re-use images with similar patterns of edges, whether or not they are semantically similar.

Finally, we would like to use very large and deep convolutional nets on video sequences. We believe that the temporal structure provides very helpful information that is missing in the image domain.



Input: 224x224x3

Conv1: 96 filters of 11x11, stride 4, padding 2 \rightarrow ReLU \rightarrow MaxPool 3x3, stride 2

Conv2: 256 filters of 5x5, padding 2 \rightarrow ReLU \rightarrow MaxPool 3x3, stride 2

Conv3: 384 filters of 3x3, padding 1 \rightarrow ReLU

Conv4: 384 filters of 3x3, padding 1 \rightarrow ReLU

Conv5: 256 filters of 3x3, padding 1 \rightarrow ReLU \rightarrow MaxPool 3x3, stride 2

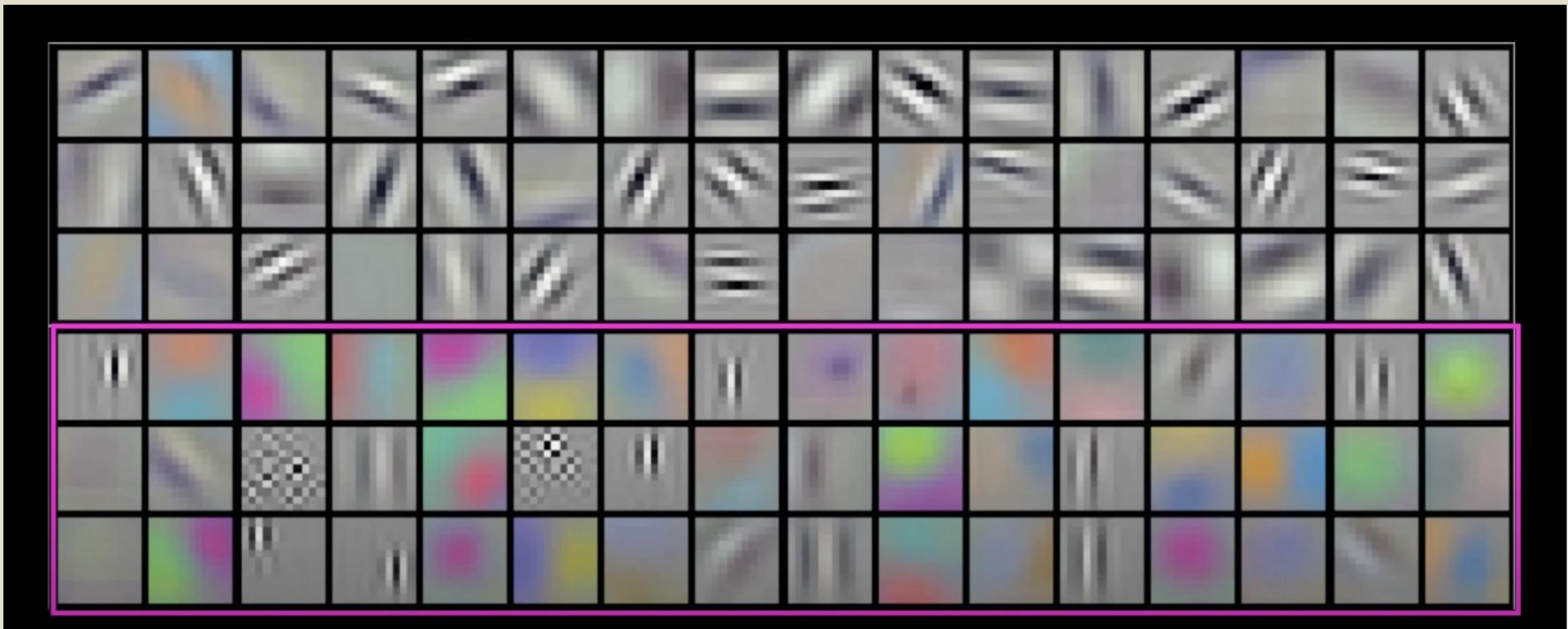
Flatten \rightarrow FC1: 4096 units \rightarrow ReLU \rightarrow Dropout

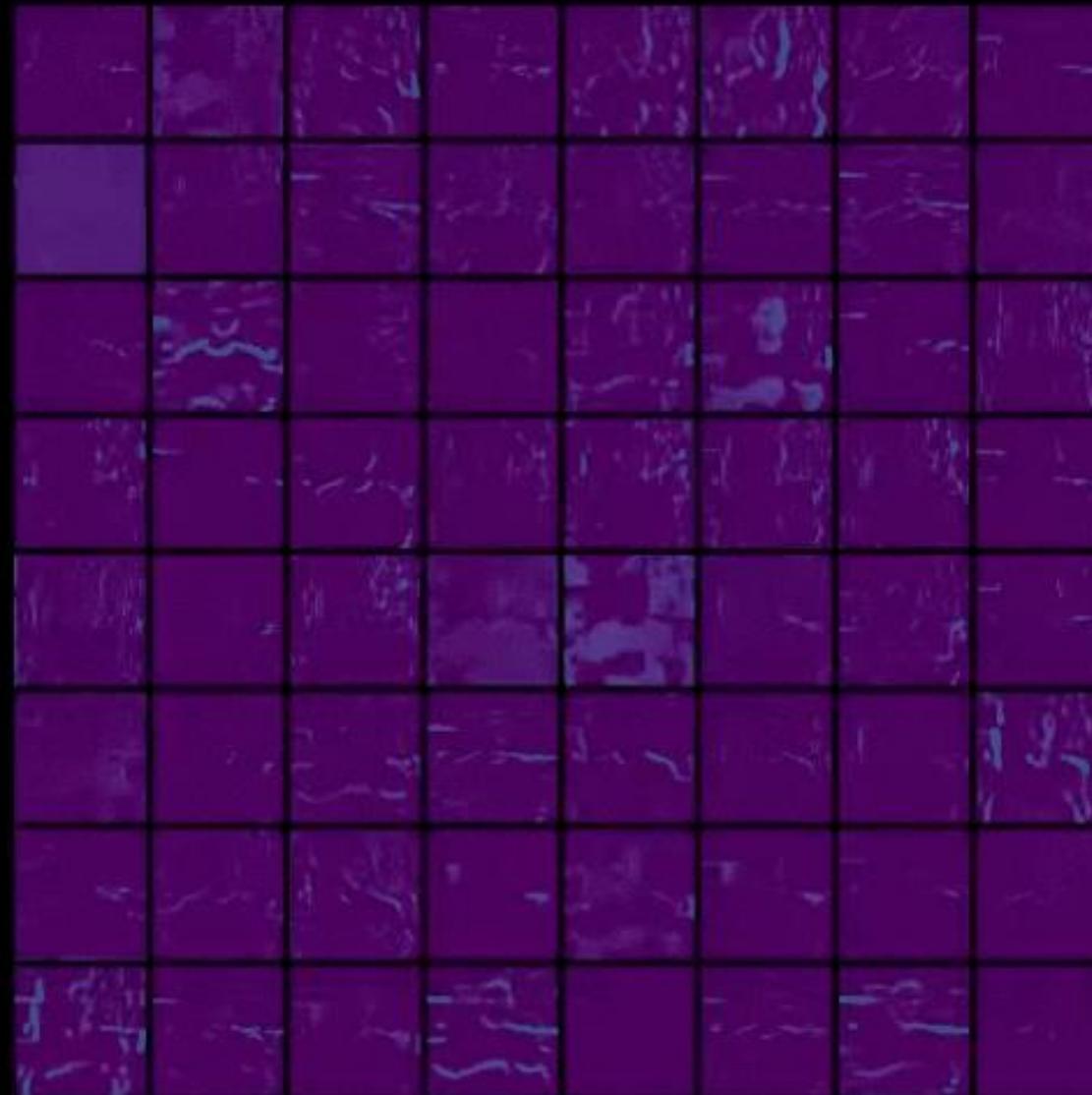
FC2: 4096 units \rightarrow ReLU \rightarrow Dropout

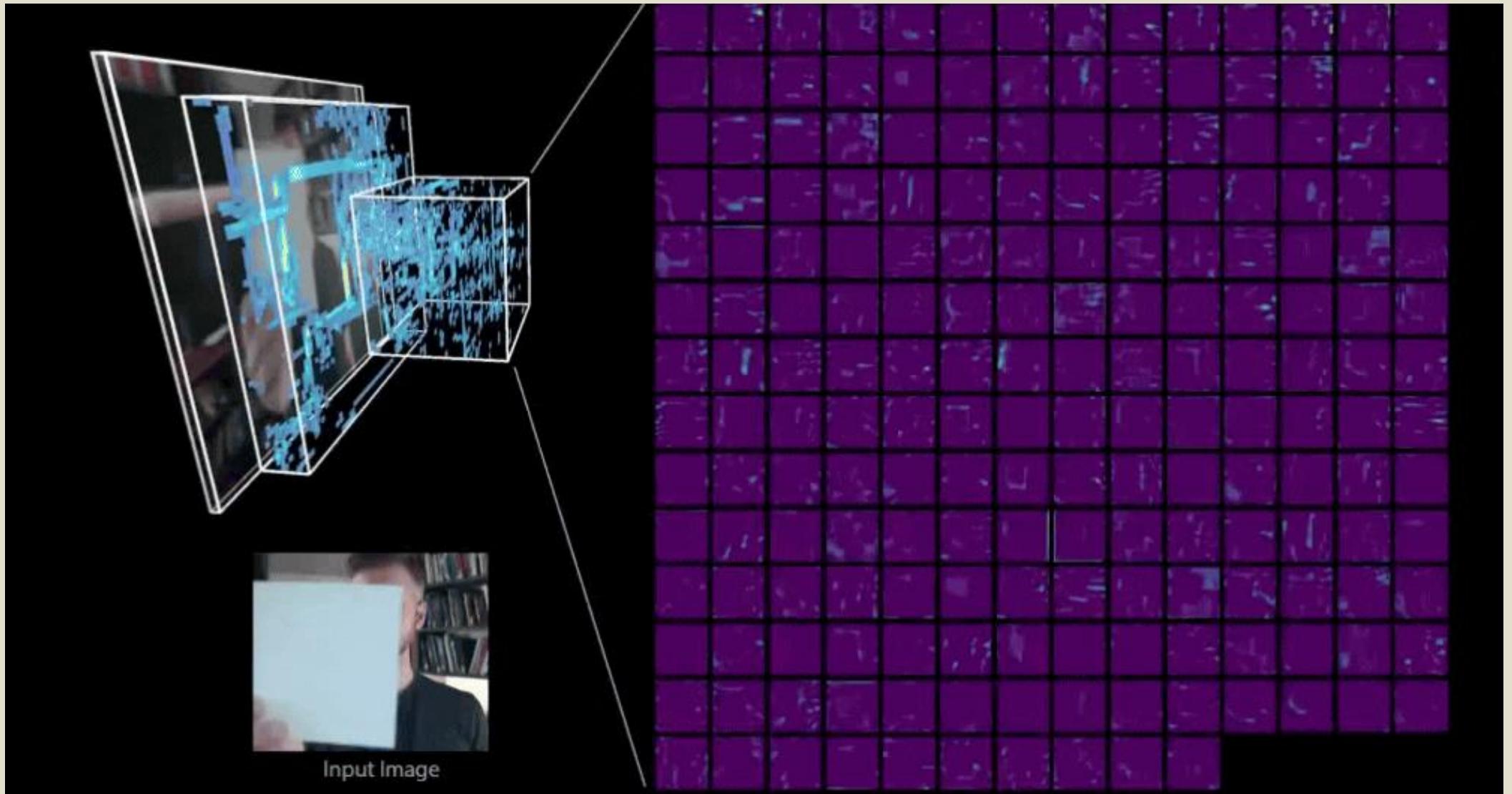
FC3: 1000 units (softmax)

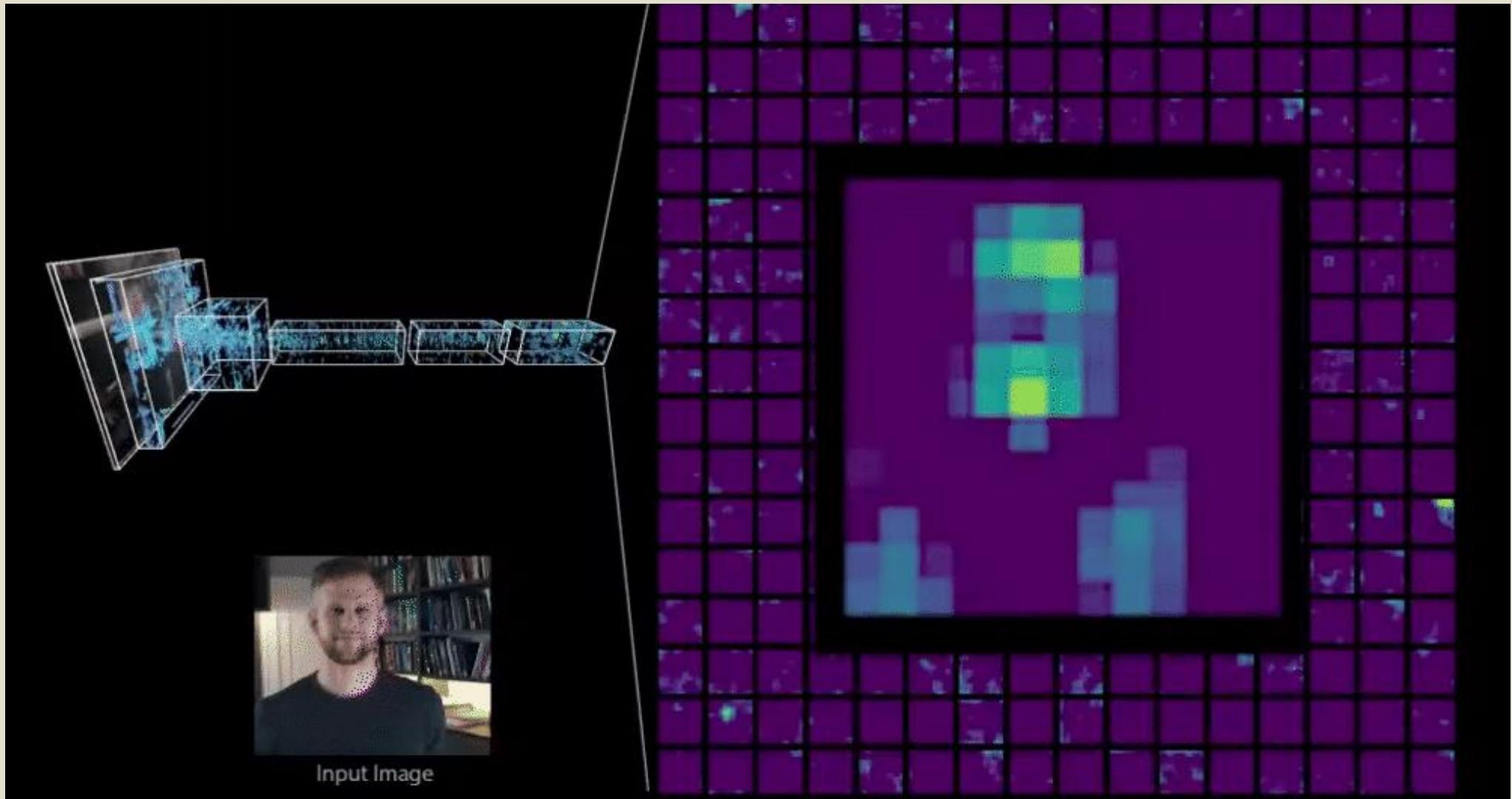
Componente	LeNet-5	AlexNet	Cambio clave / Motivación
Tamaño del input	32×32 (MNIST)	$224 \times 224 \times 3$ (ImageNet)	Adaptado a imágenes naturales RGB
Profundidad de la red	7 capas (2 conv + 2 pool + 2 FC + softmax)	8 capas (5 conv + 3 FC)	Mucho más profunda y compleja
Función de activación	tanh	ReLU (Rectified Linear Unit)	Evita saturación y acelera convergencia
Uso de GPU	CPU	Entrenado con 2 GPUs (NVIDIA GTX 580)	Aceleración masiva en entrenamiento
Tamaño del dataset	~60K imágenes (MNIST)	1.2 millones (ImageNet)	Requiere modelos más potentes

Componente	LeNet-5	AlexNet	Cambio clave / Motivación
Regularización	Ninguna	Dropout en capas fully connected	Combate overfitting
Data Augmentation	No	Sí (cropping, flipping, jittering de color)	Mejora generalización
Stride y padding	Manual y reducido	Padding explícito, strides agresivos	Control fino sobre tamaño espacial
Normalization (pre-BN era)	Parcial	Local Response Normalization (LRN)	Imita inhibición lateral biológica (hoy obsoleto)
Paralelismo en GPU	No	División del modelo entre 2 GPUs	Distribución de parámetros para computar eficiencia

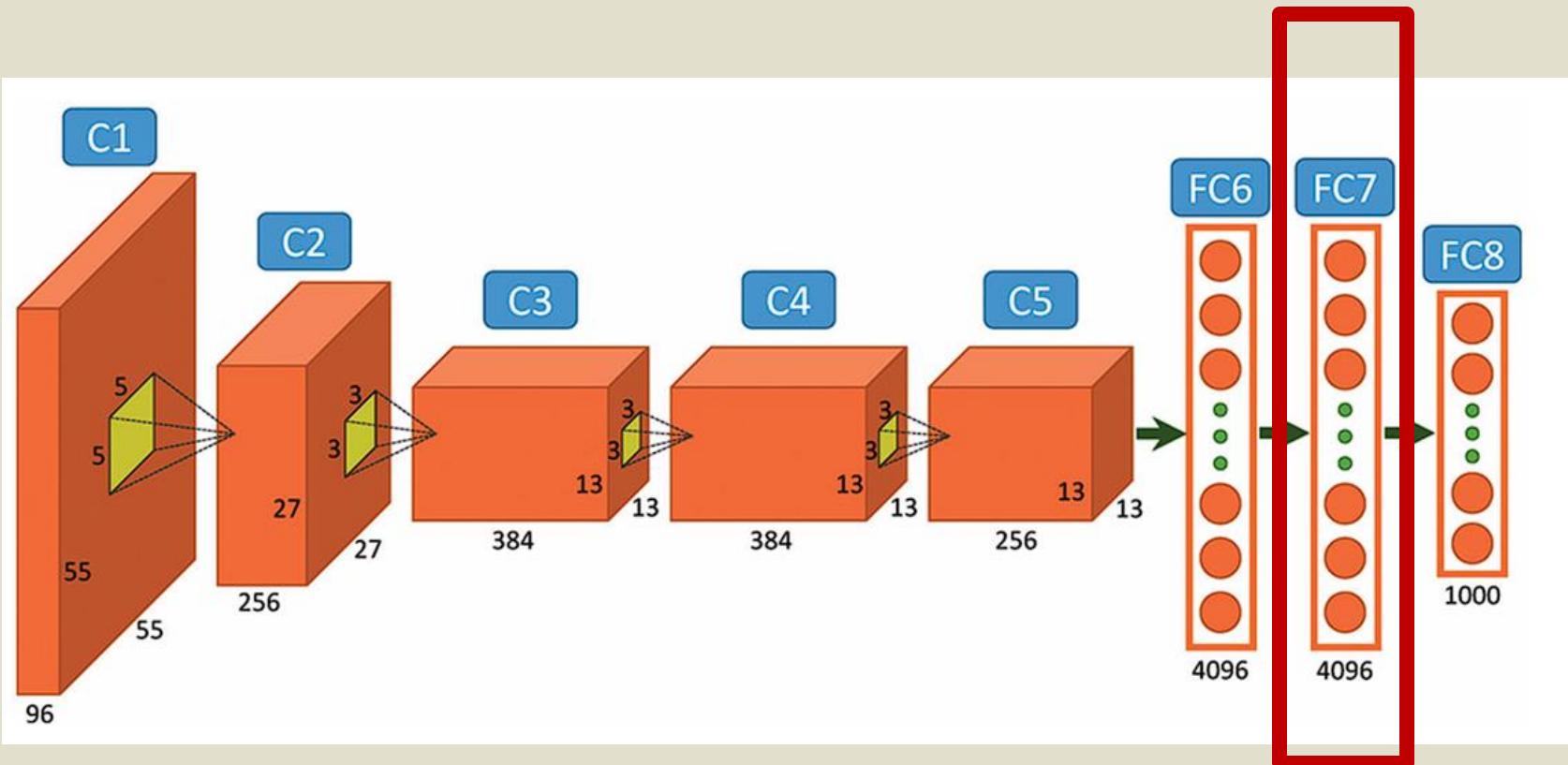


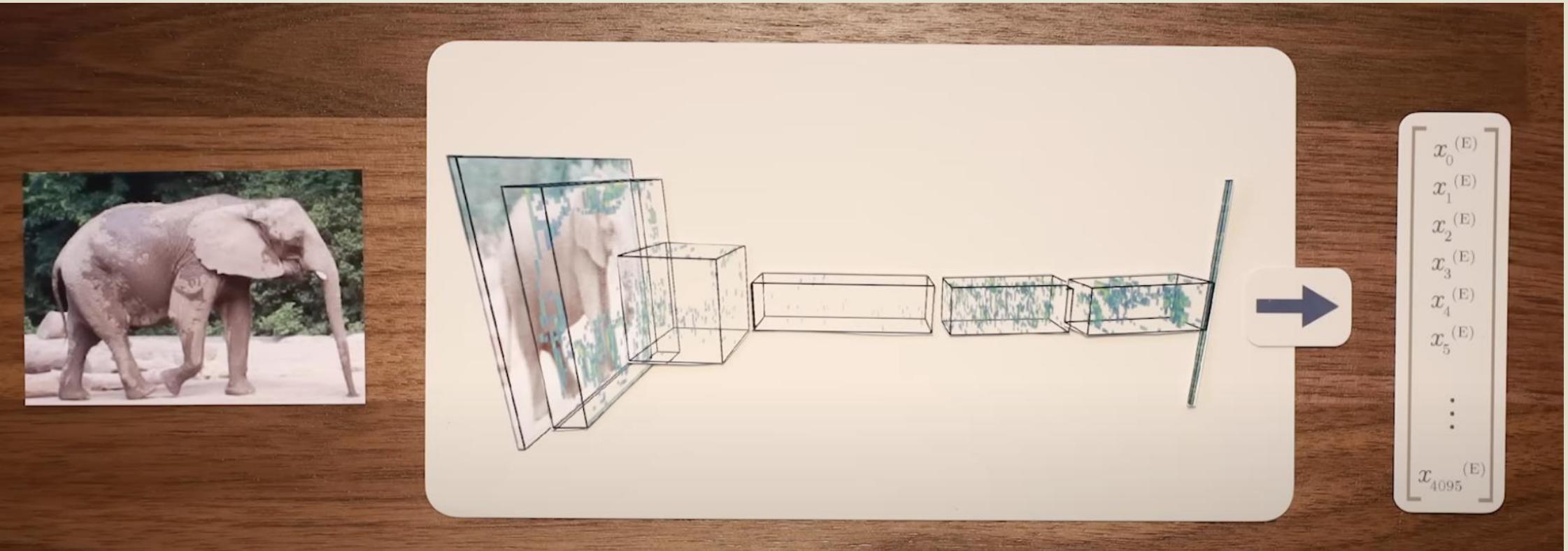


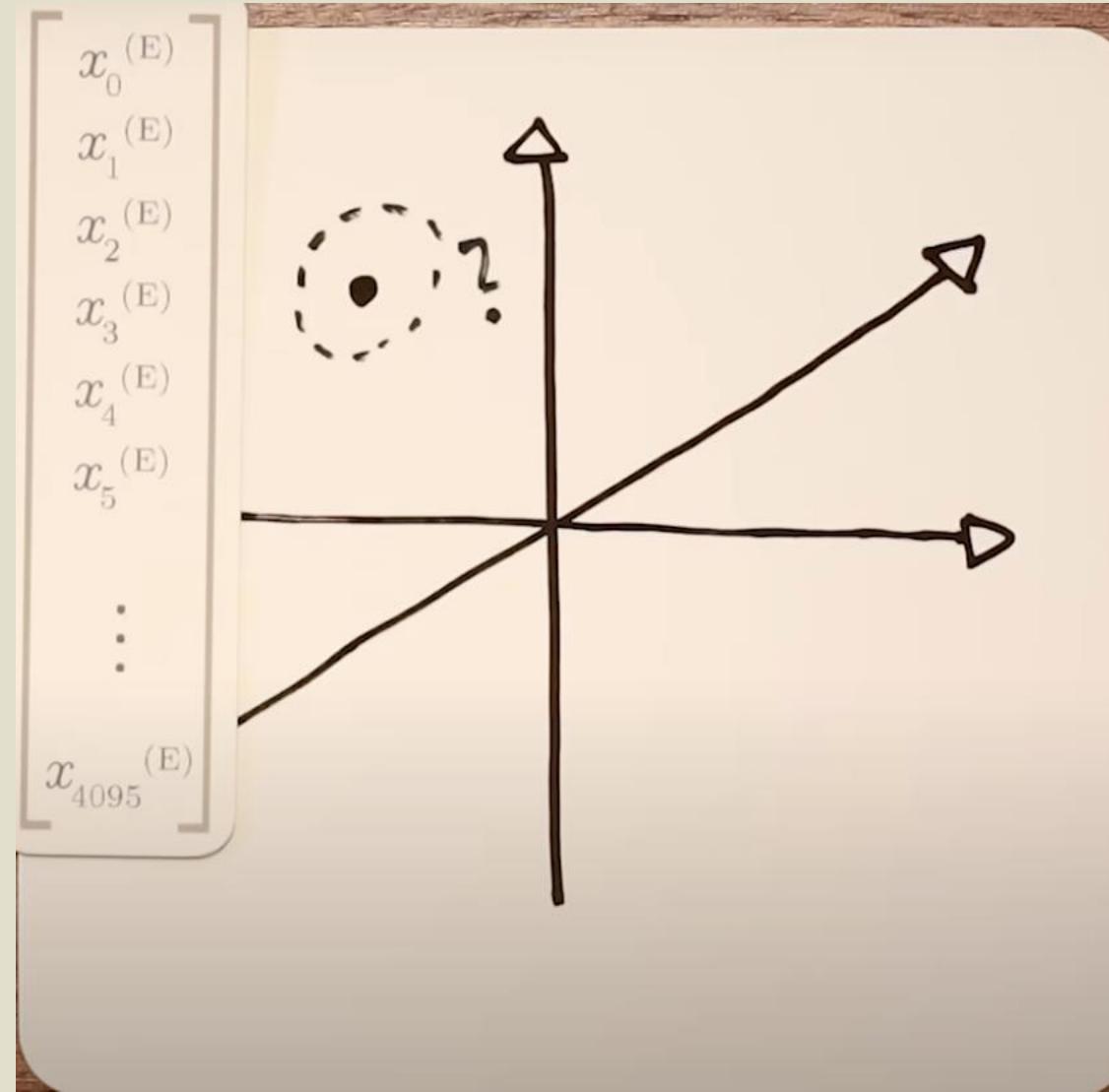


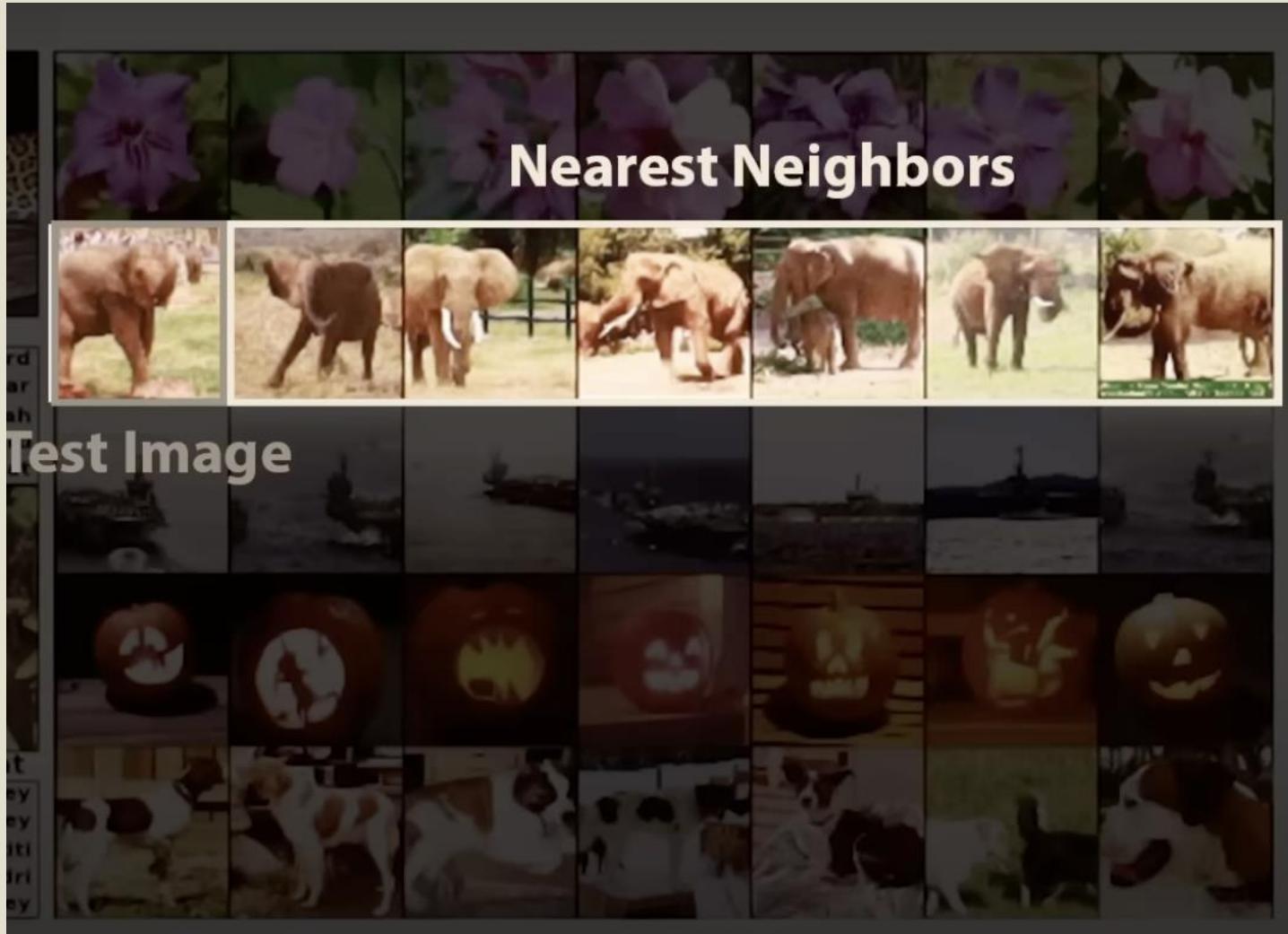


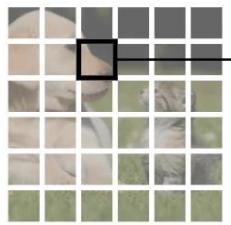
No siendo suficiente, descubren los embeddings



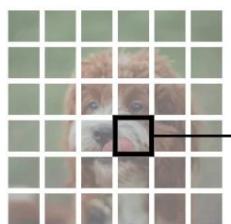
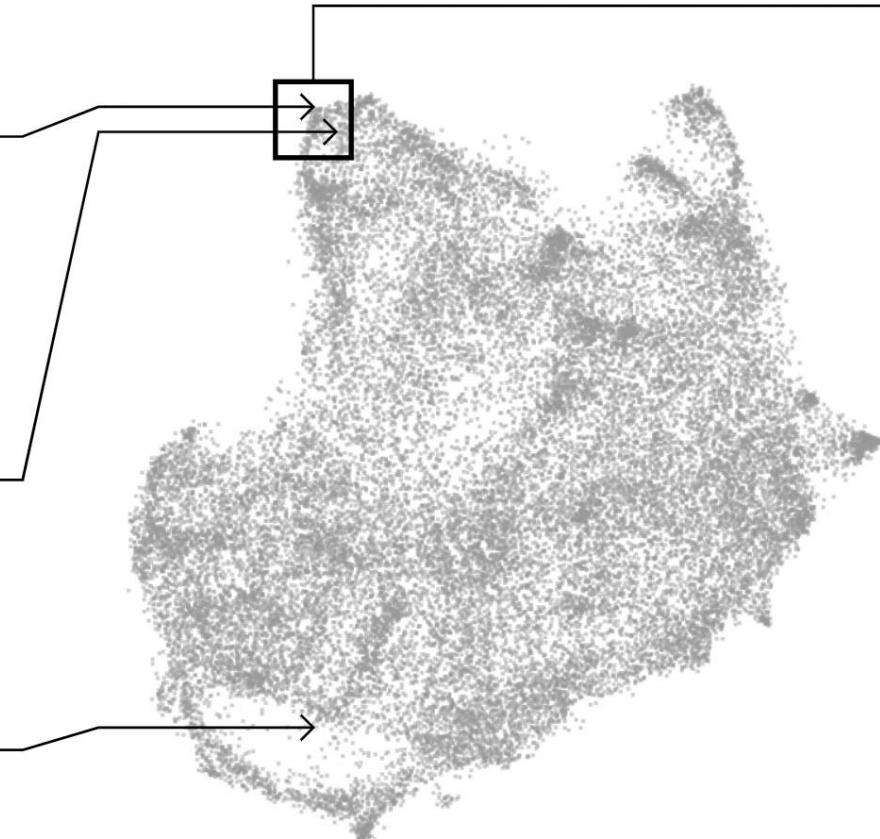




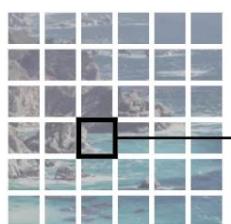




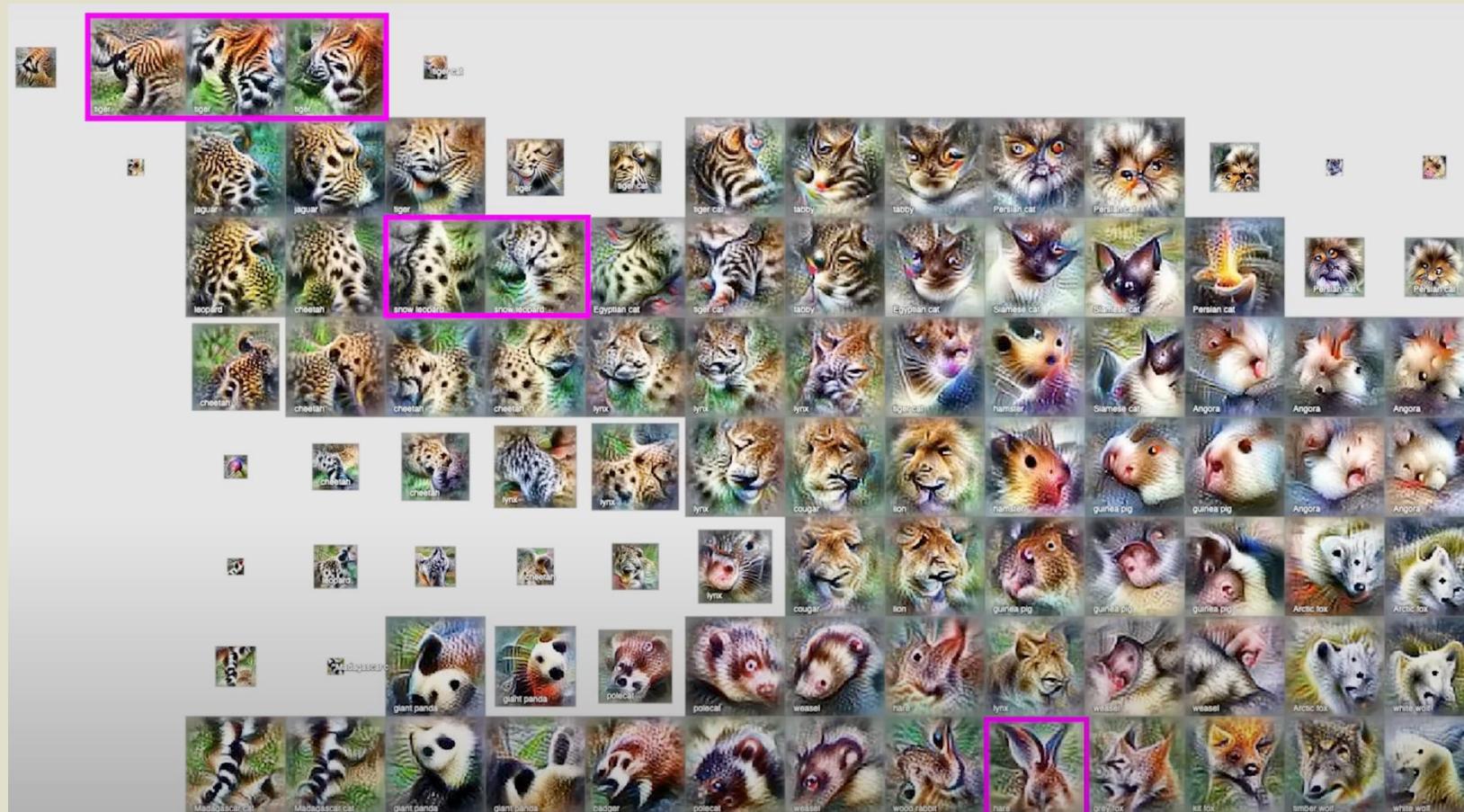
[0.370, 0.698, ...]



[0.012, 0.540, ...]



[0.034, 0.678, ...]



More Fruit

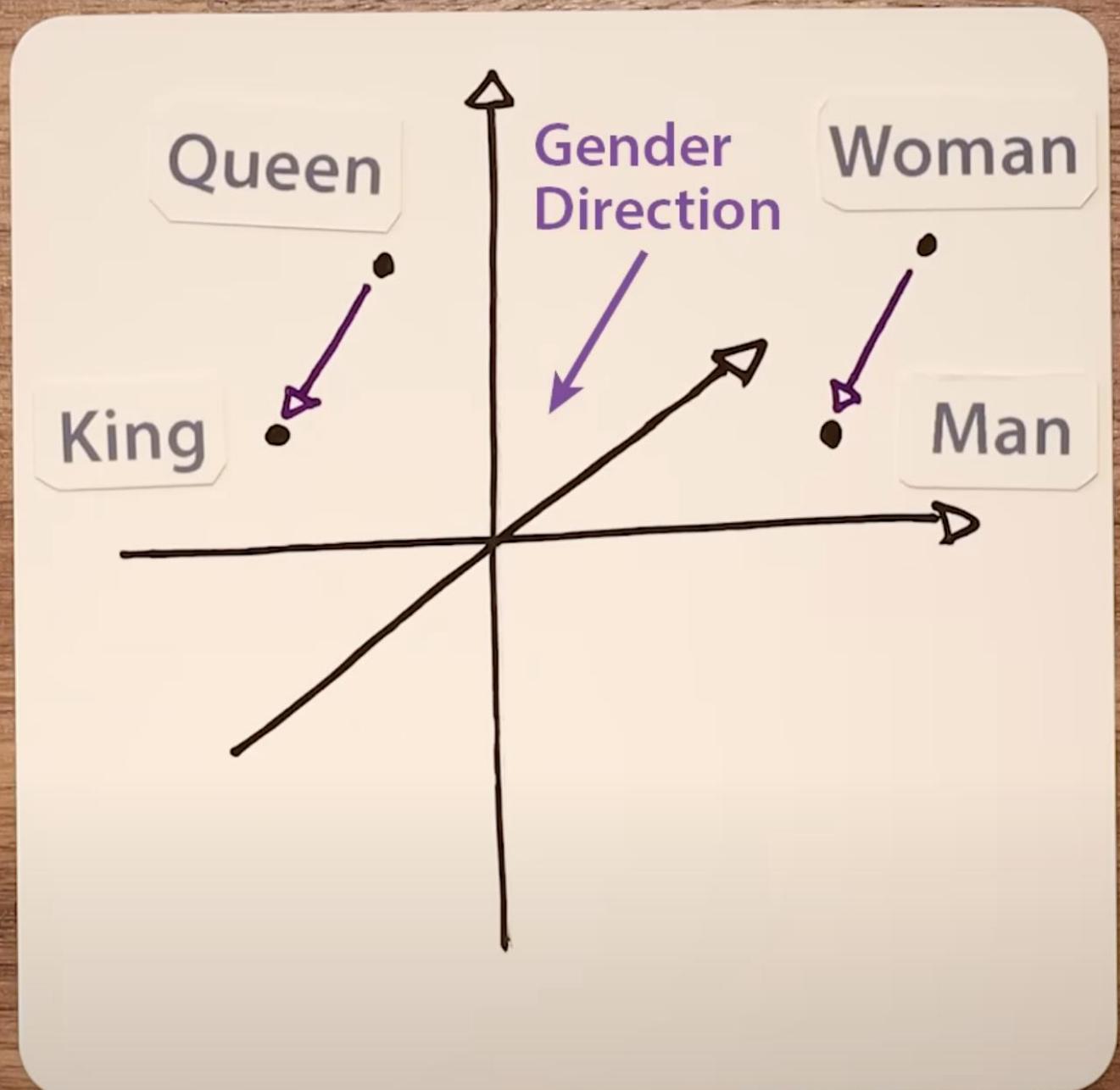


Less Fruit

Andrés Daniel Godoy Ortiz - @adgodoyo

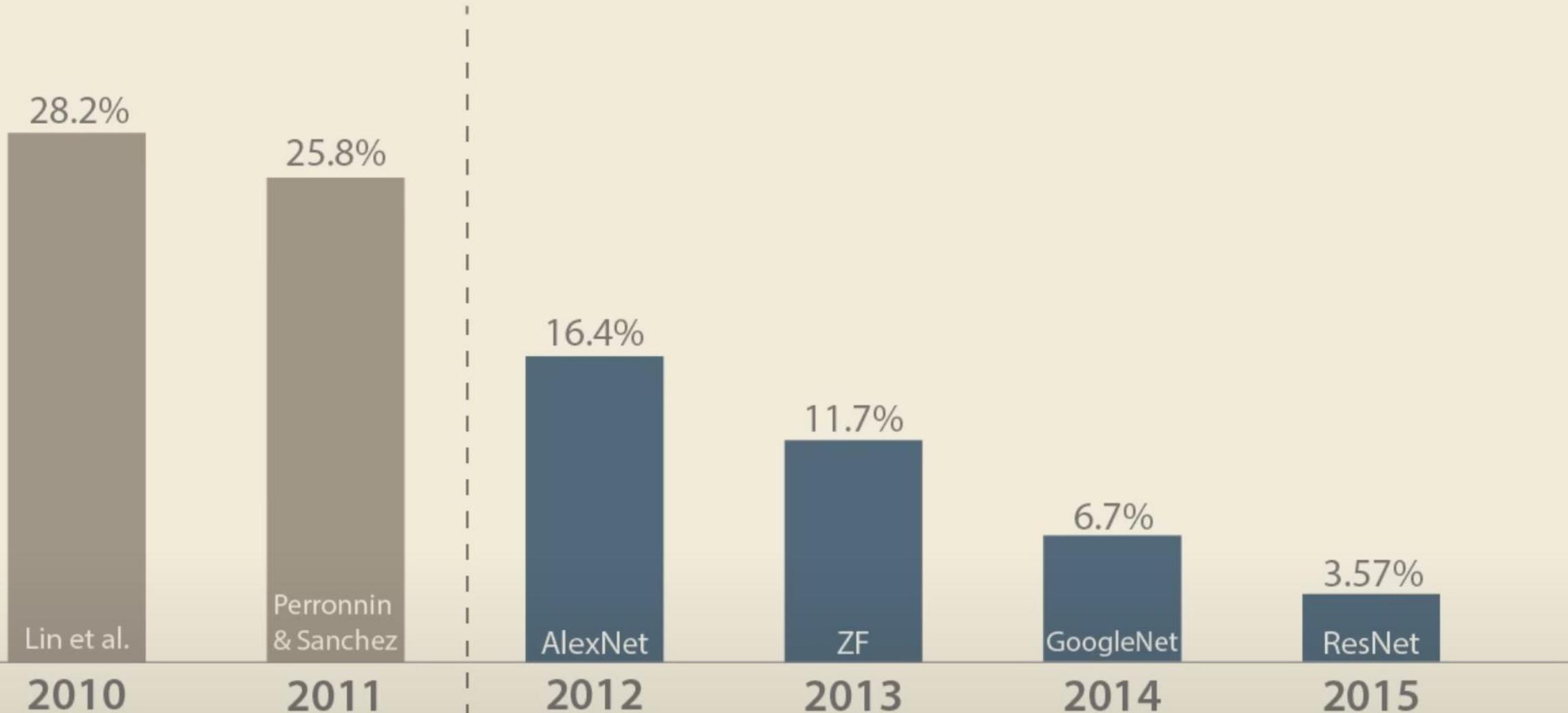
29

LLM



IMAGENET IMAGE CLASSIFICATION TOP-5 ERROR RATE

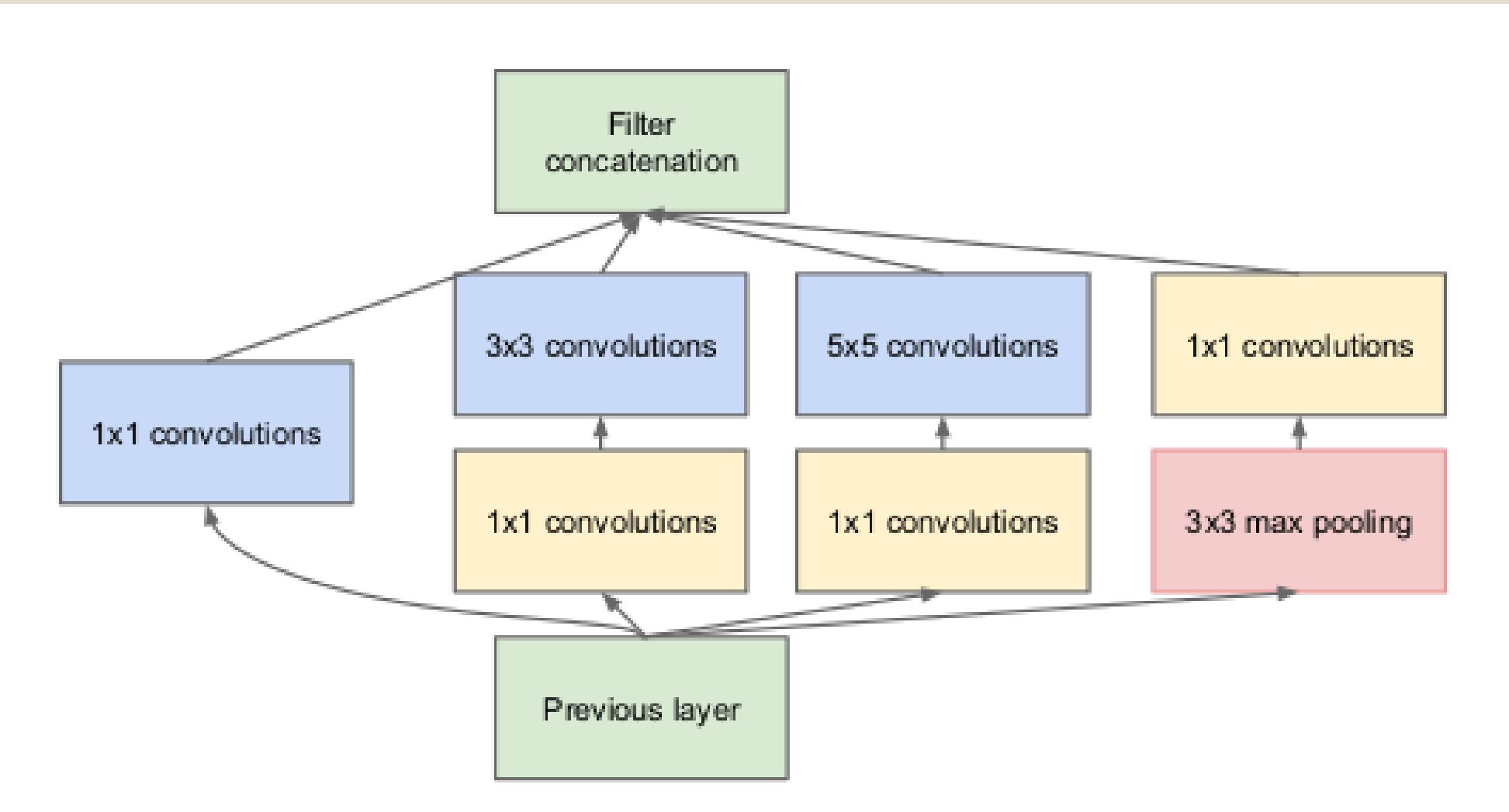
(ILSVRC, lower is better)



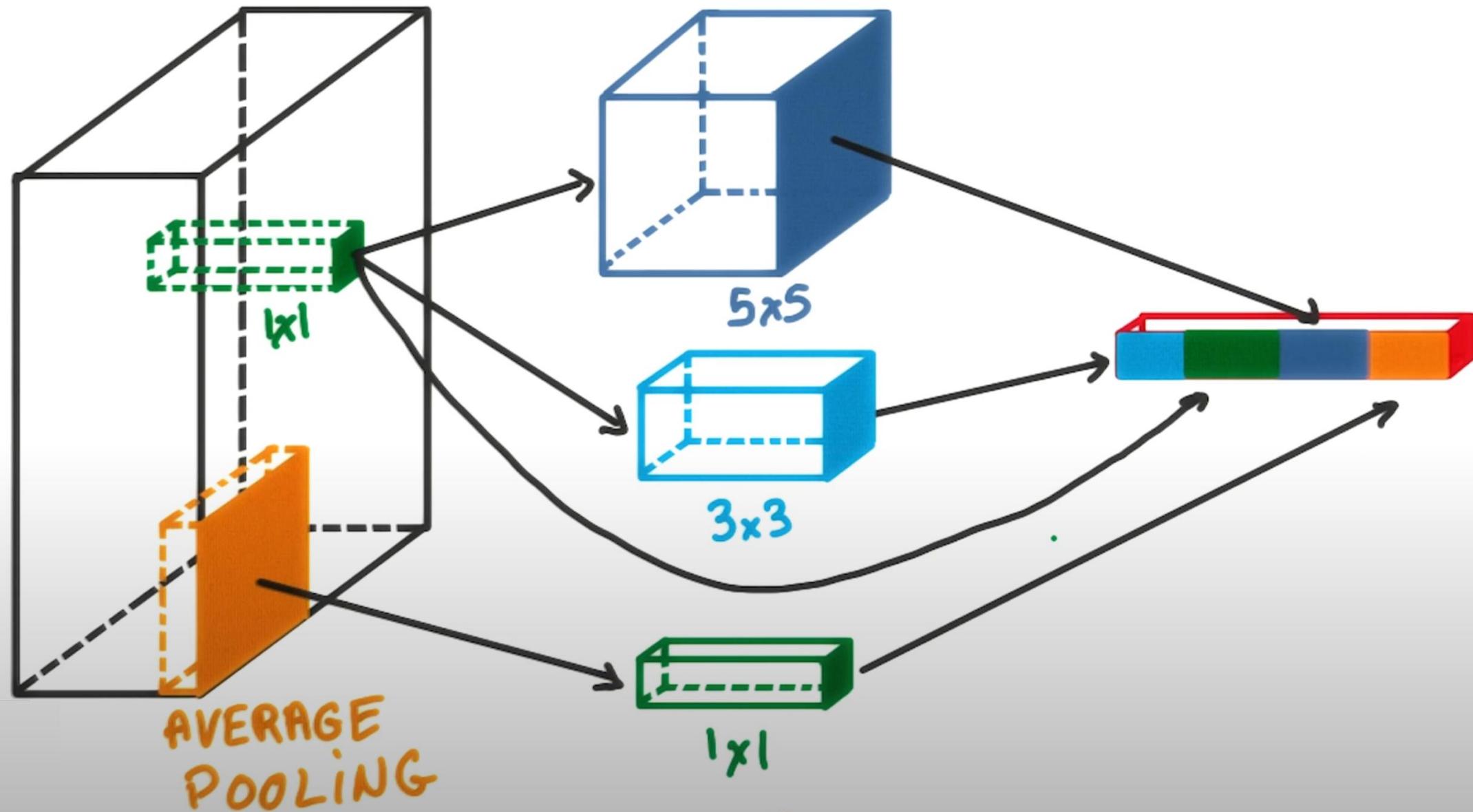
GOOGLENET (2014)

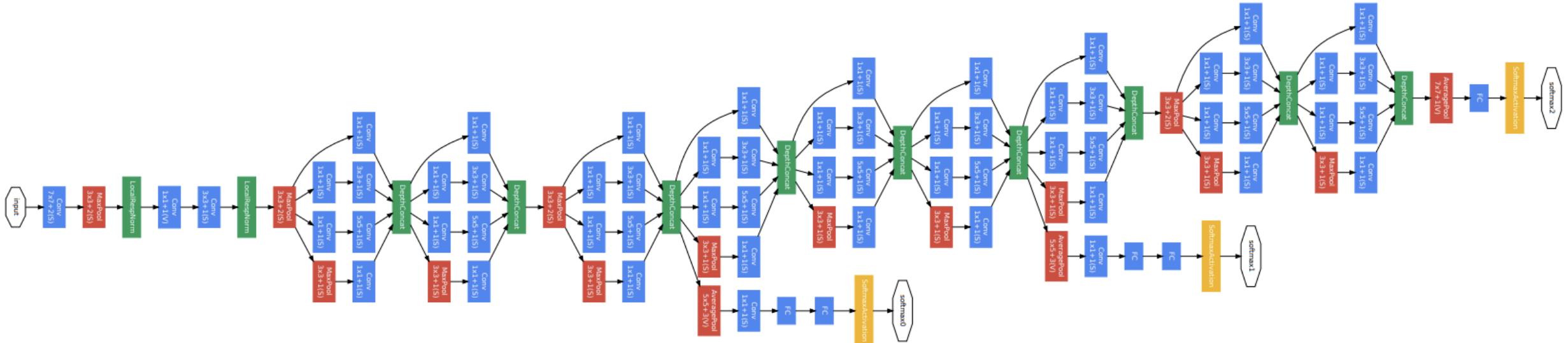
Eficiencia
computacional

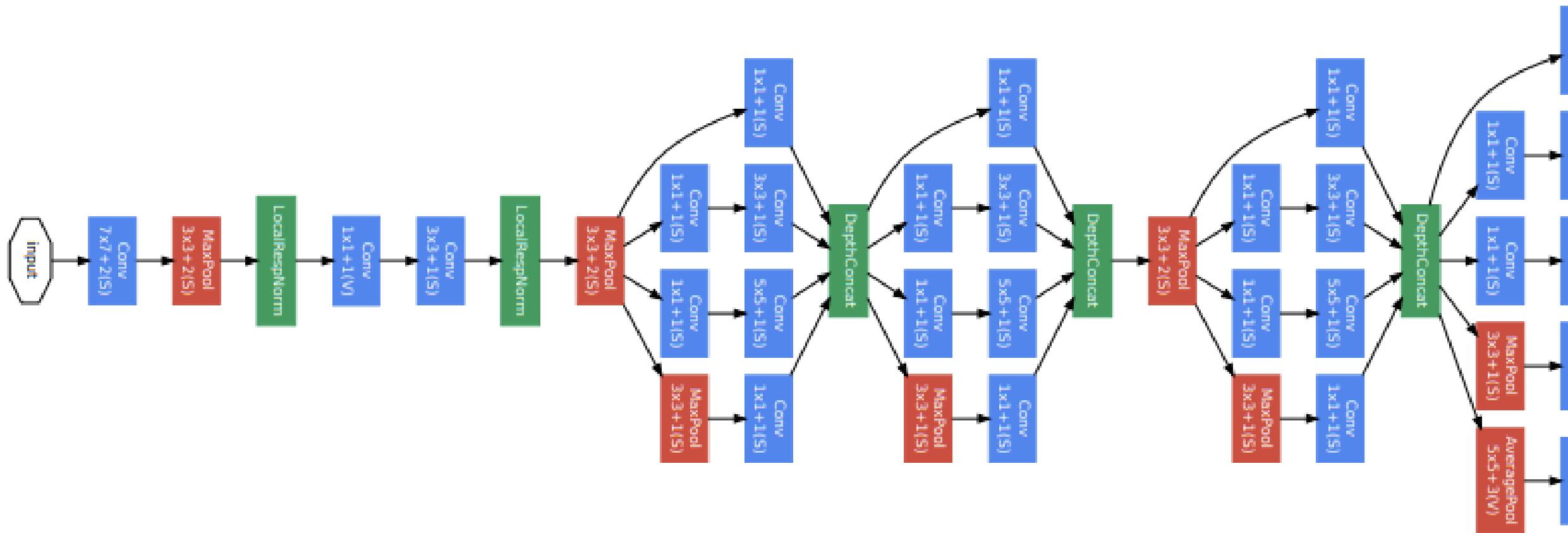
Bloque/módulo Inception

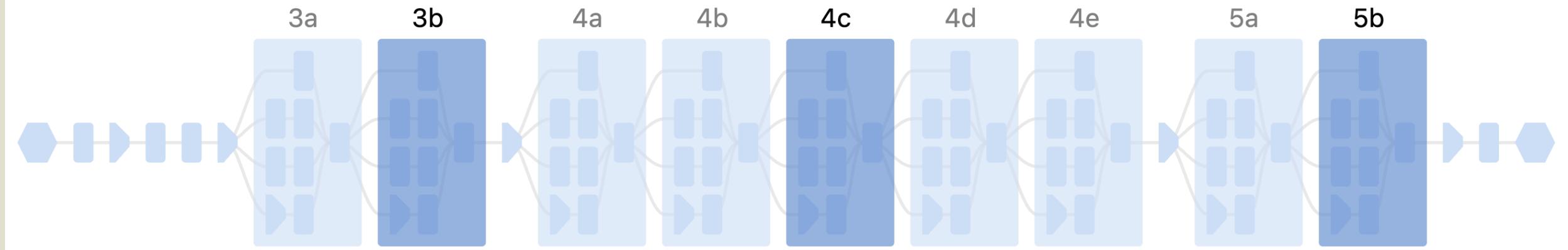


INCEPTION MODULES









distill.pub/2019/activation-atlas/

- Demostró que profundidad no implica explosión de parámetros.
- Introdujo el **parallelismo multi-escala** ahora común en FPN, U-Net, MobileNet, EfficientNet y NAS-Nets.
- Popularizó el uso extensivo de **1×1 conv** como “mezclador de canales” (antecedente de *self-attention* bajo perspectiva de mezcla lineal)

RESNET (2015)

Anti-Vanishing
gradient

- A medida que hacemos redes más profundas, **esperaríamos mejor rendimiento**, pero ocurre lo contrario: muchas veces el error de entrenamiento **aumenta o se estanca**.
- Esto se llama el **problema de degradación**.
-  No es *overfitting*, ¡es que la red más profunda rinde peor incluso en entrenamiento!

Recordemos cómo funciona backpropagation

Si tenemos una red clásica con una función compuesta como:

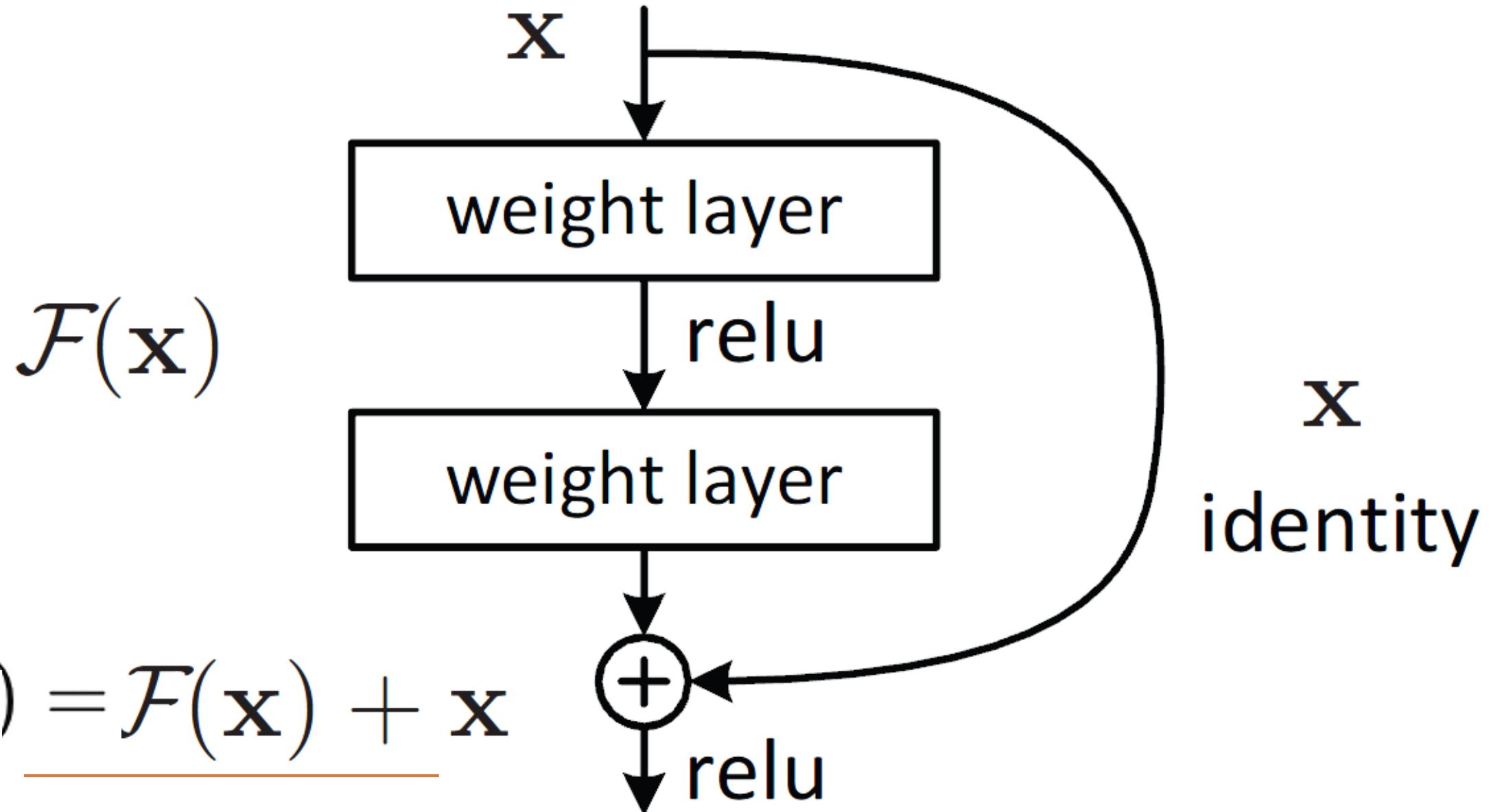
$$y = \text{Salida final} = H(x) = f(g(x))$$

Entonces para calcular la derivada de la pérdida L con respecto a x , usamos la **regla de la cadena**:

$$\frac{dL}{dx} = \frac{dL}{dH} \cdot \frac{dH}{dx} = \frac{dL}{dH} \cdot f'(g(x)) \cdot g'(x)$$

Si tenemos muchas capas profundas (muchas funciones anidadas), esa multiplicación se alarga:

$$\frac{dL}{dx} = \frac{dL}{dz} \cdot f'_n \cdot f'_{n-1} \cdots \cdot f'_1$$



Ejemplo:

Supongamos que $x = 2$, y queremos que el bloque aprenda una función $H(x) = 5$.

En una red clásica, el bloque tendría que aprender directamente esa función:

$$H(x) = 5$$

Pero en ResNet, en lugar de eso, le pedimos al bloque que **aprenda la diferencia** entre lo deseado y lo que ya tenemos:

$$F(x) = H(x) - x = 5 - 2 = 3$$

Entonces el bloque solo tiene que aprender:

$$\text{Salida} = x + F(x) = 2 + 3 = 5$$

Módulo clásico

Observemos lo qué retorna la función

```
import torch
import torch.nn as nn

class ClassicBlock(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear1 = nn.Linear(1, 1)
        self.relu = nn.ReLU()
        self.linear2 = nn.Linear(1, 1)

    def forward(self, x):
        out = self.linear2(self.relu(self.linear1(x)))
        return out
```

Módulo Residual

Observemos lo qué retorna la función

```
import torch
import torch.nn as nn

class ResidualBlock(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear1 = nn.Linear(1, 1)
        self.relu = nn.ReLU()
        self.linear2 = nn.Linear(1, 1)

    def forward(self, x):
        fx = self.linear2(self.relu(self.linear1(x)))
        return fx + x # Esto es  $H(x) = F(x) + x$ 
```

Ruta del gradiente

Red clásica

Si un bloque sin residual produce

$$y = H(x),$$

al hacer backpropagation, el gradiente de la pérdida L con respecto a su entrada x es

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial H(x)}{\partial x}.$$

...

Conexión Residual

En ResNet, cada bloque hace

$$y = F(x) + x.$$

Entonces

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial(F(x) + x)}{\partial x} = \frac{\partial L}{\partial y} \left(\frac{\partial F(x)}{\partial x} + 1 \right).$$

Sin residual: $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot F'(x).$

Con residual: $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot (F'(x) + 1).$

Imagina que en un punto del entrenamiento $\frac{\partial F}{\partial x} = 0.1$ (una transformación pequeña). Y que $\frac{\partial L}{\partial y} = 2$.

- **Sin residual:**

$$\frac{\partial L}{\partial x} = 2 \times 0.1 = 0.2.$$

- **Con residual:**

$$\frac{\partial L}{\partial x} = 2 \times (0.1 + 1) = 2 \times 1.1 = 2.2.$$

El gradiente fluye mejor

Esa estructura **preserva parte del gradiente original**, como si “**una parte de la derivada fuera constante**”, y por eso se dice que “el gradiente fluye mejor”.

$$\frac{\partial \mathcal{L}}{\partial x_0} = \frac{\partial \mathcal{L}}{\partial x_L} \cdot \prod_{i=1}^L \frac{\partial x_i}{\partial x_{i-1}}$$

Pero ahora:

$$\frac{\partial x_{l+1}}{\partial x_l} = \frac{\partial}{\partial x_l} (x_l + F_l(x_l)) = I + \frac{\partial F_l}{\partial x_l}$$

Entonces:

$$\frac{\partial \mathcal{L}}{\partial x_0} = \frac{\partial \mathcal{L}}{\partial x_L} \cdot \prod_{l=0}^{L-1} \left(I + \frac{\partial F_l}{\partial x_l} \right)$$

Comparación con Inception

Inception combina en paralelo distintos tamaños de filtro (1×1 , 3×3 , 5×5) dentro de un módulo, capturando multi-escala. Su flujo es más ancho y ramificado.

ResNet opta por un diseño más sencillo y profundo, usando un único camino principal reforzado por atajos.

Beneficios

Antes de ResNet, aumentar la profundidad de las CNNs más allá de 20-30 capas **empeoraba el entrenamiento**, aunque el modelo tenía más capacidad teórica.

ResNet resolvió esto introduciendo bloques residuales que permiten a la red aprender "**residuos**" en vez de funciones completas.

Arquitectura

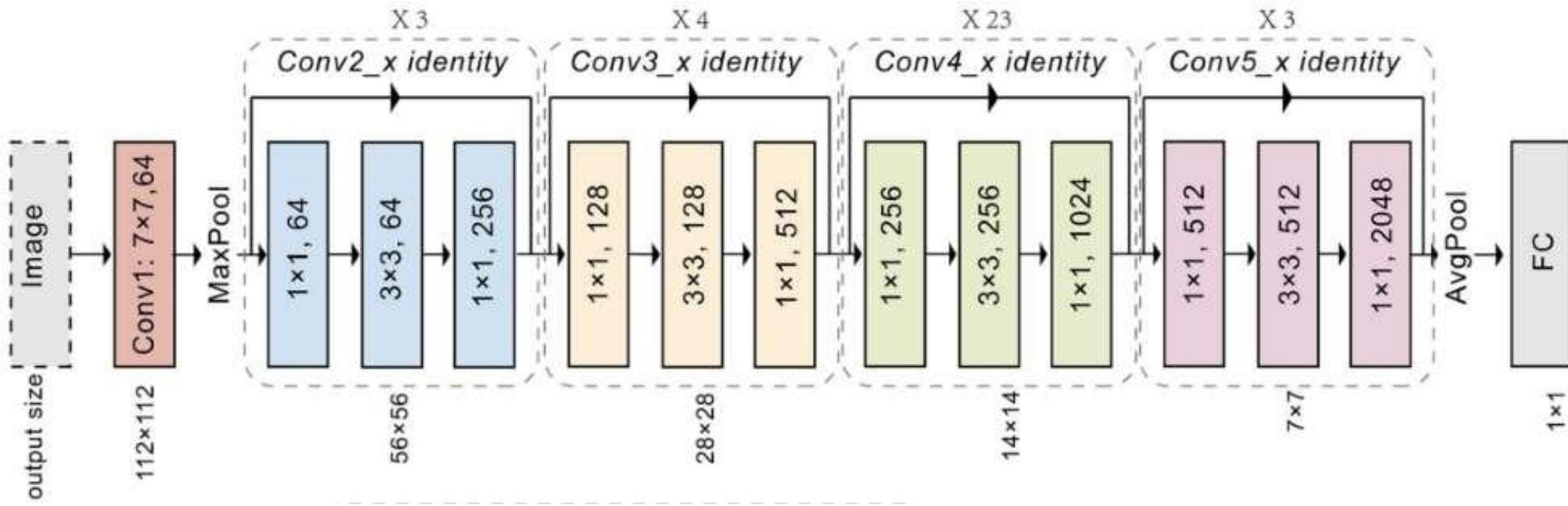
ResNet-18/34: *bloques básicos*

$$F(x) = \text{Conv3x3} \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv3x3} \rightarrow \text{BN}$$

ResNet-50/101/152: *bloques "bottleneck"*

$$F(x) = \text{Conv1x1} \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv3x3} \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{Conv1x1} \rightarrow \text{BN}$$

ResNet-50



Feature Visualization

- <https://distill.pub/2017/feature-visualization/>
- <https://openai.com/index/introducing-activation-atlases/>
- <https://distill.pub/2019/activation-atlas/>