# IMAGE RECOGNITION AND IMAGENET

Janu Verma

Hike, New Delhi

@januverma

# About Me

- Sr. Machine Learning Scientist at Hike, New Delhi — deep learning, recommendation systems, NLP, network analysis.

- Previously :

    - IBM Research, New York

    - Cornell University

    - Kansas State University

    - Cambridge University

- Researcher in machine learning , data visualization, and mathematics.

- Public Speaking, Teaching ML and data science.

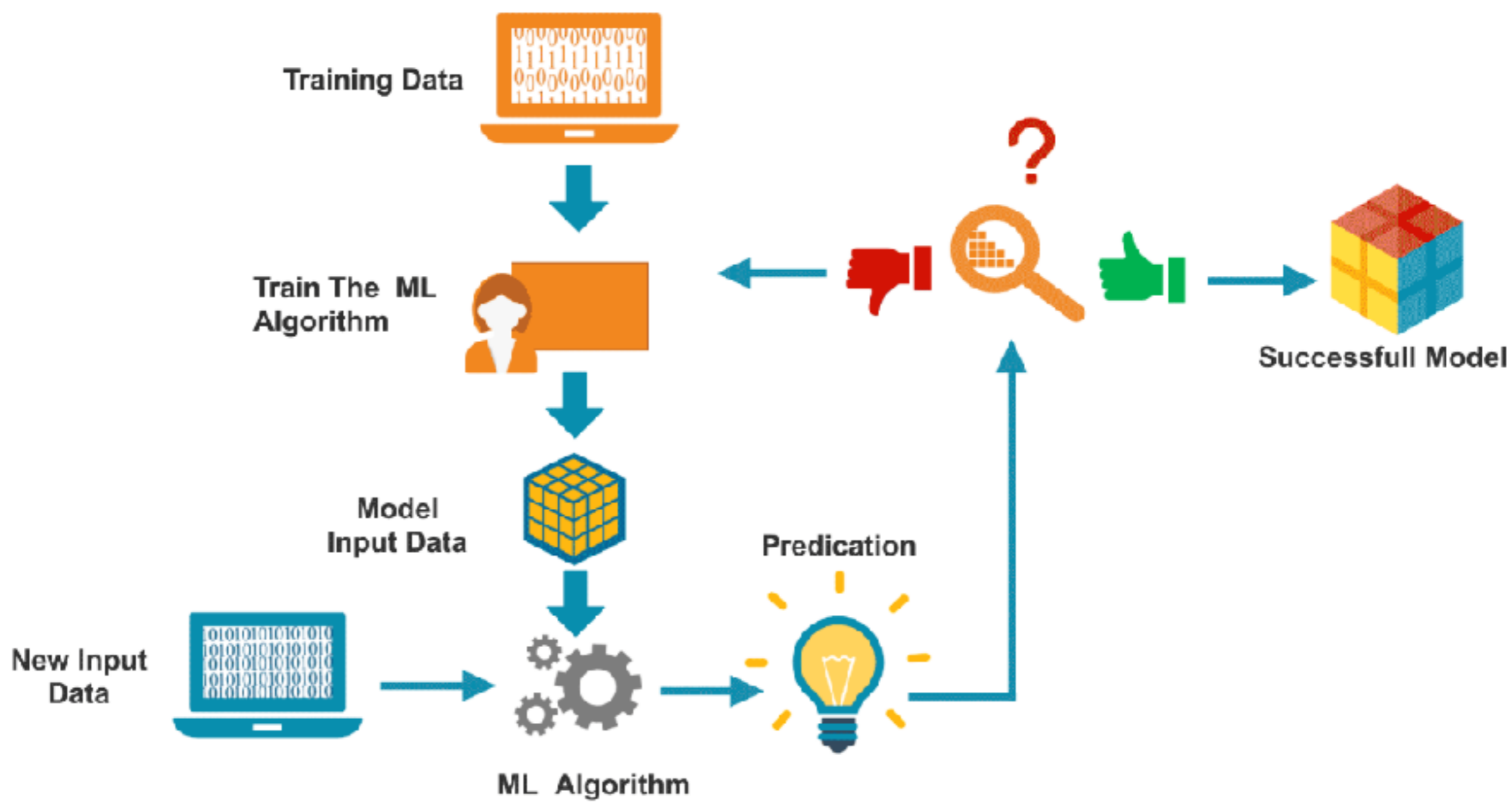- Advise Startups on ML, data science, and hiring & team building.

# Machine Learning
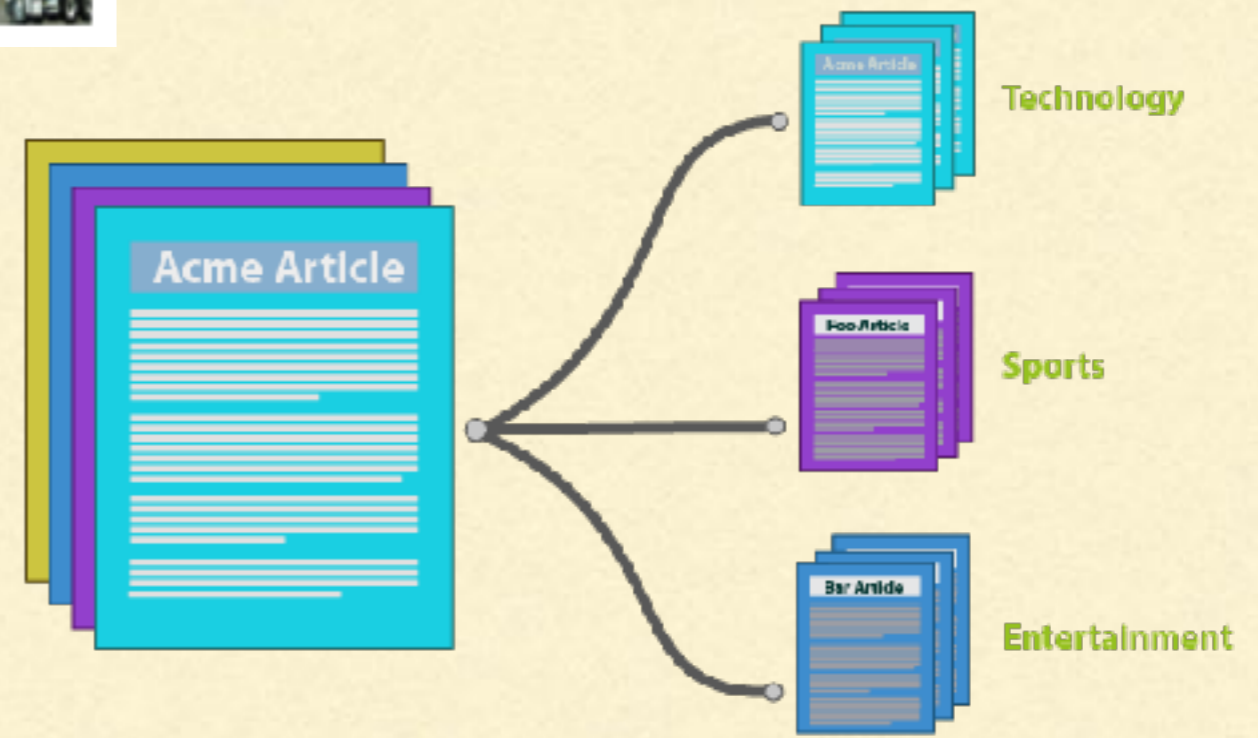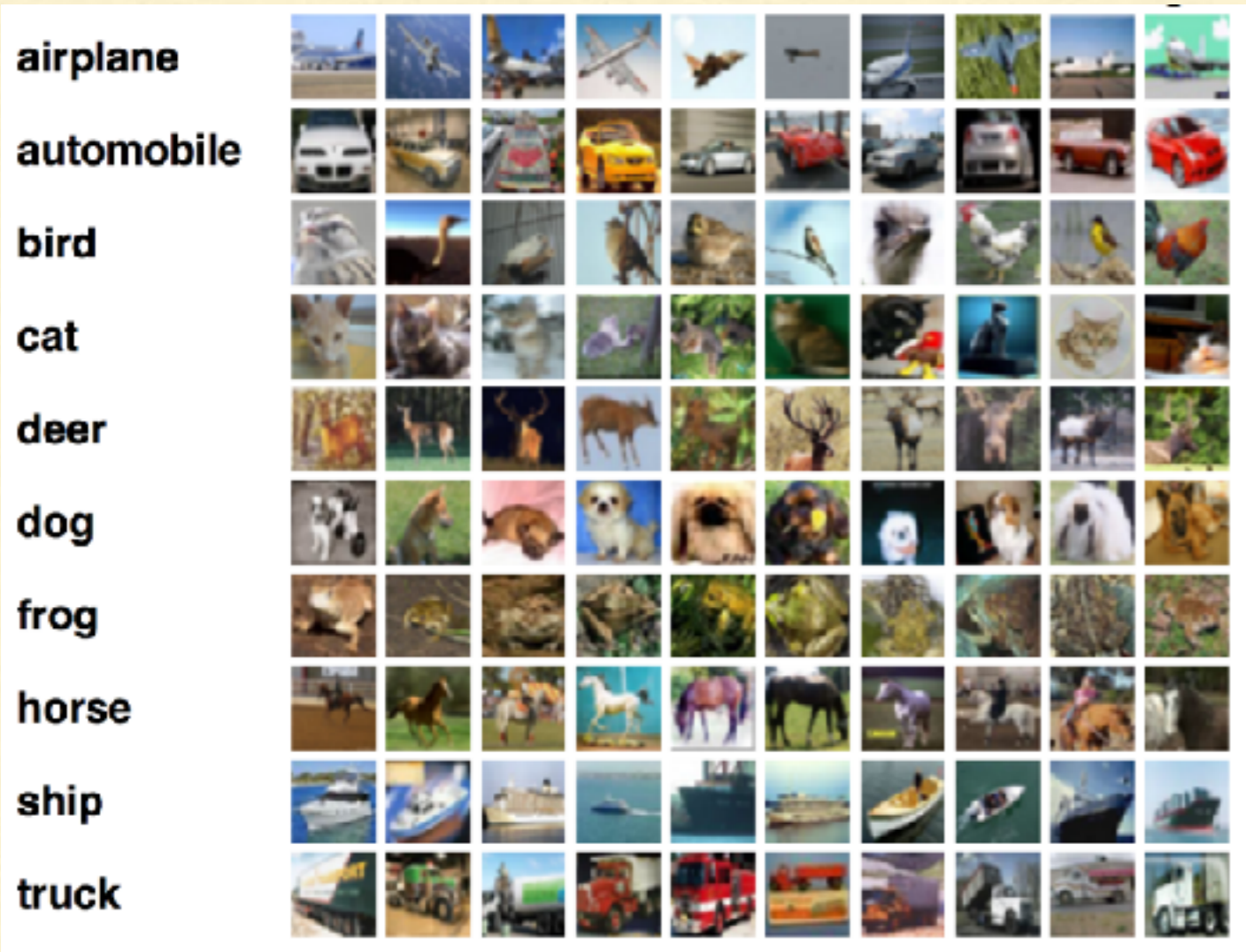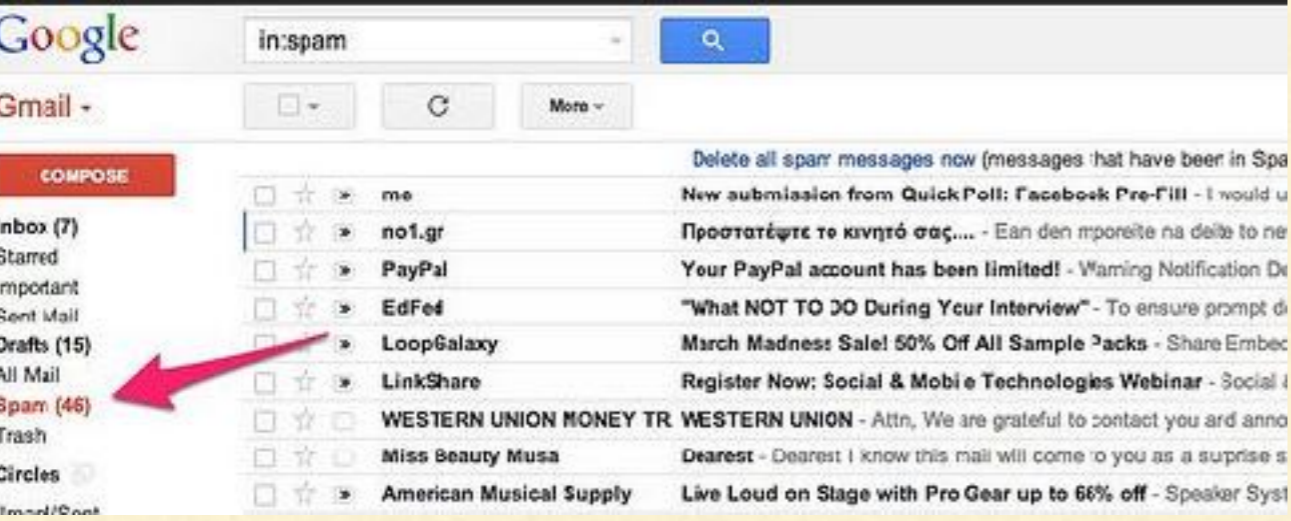
- **Machine learning** is a field of computer science that designs systems with the ability to automatically learn and improve from experience without being explicitly programmed.

- Computer systems that access data and use statistical/mathematical techniques to learn patterns and make inference based on probabilistic responses.

- *Supervised learning* involves providing example inputs and respective outputs to the the program, which 'learns' to make predictions on the outputs for new, unseen inputs.

- e.g. a classification system to categorize the images as cats or dogs.

Machine Learning Pipeline

airplane

automobile

bird

cat

deer

dog

frog

horse

ship

truck



Acme Article

Technology

Sports

Entertainment

**People You May Know**

Add people you know as friends and connect with public profiles you like.

| Lala Lalabs | Brian Crecente | Brian Ashcraft |
| Add as friend | Add as friend | Add as friend |
| justin bieber | Camile Gozon | Karla Danielle Beger |
| Add as friend | Add as friend | Add as friend |
| Taylor-Alison Swifty | Adam Rifkin | Luke Plunkett |
| Add as friend | Add as friend | Add as friend |

Who to follow · Refresh · View all

Andrew Robb  @AndrewRob...
Followed by EssentialView and ot...
Follow

ABC South East SA  @abcsouthe...
Followed by LGA South Australia ...
Follow

Ian Shuttleworth  @Shutts10
Followed by Jason McConnell an...
Follow

Browse categories · Find friends

# Deep Learning

- Deep learning has seen tremendous success in past few years with State-of-The-Art performance in many real-world applications.

- All examples mentioned earlier are dominated by deep learning models.

- But training a decent deep learning model requires large amount of data.

- Success of deep learning can be attributed to development of great computation resources (e.g. nvidia, distributed systems) and availability of huge amount of data (Google).

- Deep learning tools existed, but data was missing to actualise the full potential.

- Only big companies have data large enough for deep learning.

- Again:  *Algorithms open, data proprietary.*

# IMAGE RECOGNITION

- Identify objects in the images.

# CONVOLUTIONAL NEURAL NETWORKS

- CNNs were responsible for major breakthroughs in Image Classification.

- Core of most Computer Vision systems today, from Facebook's automated photo tagging to self-driving cars.

- Also being applied to NLP problems like text classification.

- A CNN comprises of many *convolutional layers* which perform a '*convolution*' operation on the previous layer instead of a '*feed-forward*'.

# A convolution can be though of as a sliding window function applied to a matrix



Image

Convolved Feature

Averaging each pixel with its neighboring values blurs an image:



| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 |

cf: Denny Britz, WildML

# Edge detection by convolutions

# CNN

- CNNs are basically just several layers of convolutions with *nonlinear activation functions* like **ReLU** or **tanh** applied to the results.

- Traditional feed-forward networks - each input neuron to each output neuron in the next layer.

- CNNs - Convolutions over the input layer to compute the output.

- Local connections - each region of the input is connected to a neuron in the output.

- Each layer applies different filters, typically hundreds or thousands like the ones showed above, and combines their results.

- Also pooling layers.

Let's say you want to classify whether or not there's an elephant in an image. Because you are sliding your filters over the whole image you don't really care *where* the elephant occurs. In practice, *pooling* also gives you invariance to translation, rotation and scaling, but more on that later. The second key aspect is (local) *compositionality*. Each filter *composes* a local patch of lower-level features into higher-level representation. That's why CNNs are so powerful in Computer Vision. It makes intuitive sense that you build edges from pixels, shapes from edges, and more complex objects from shapes.

A typical convolutional neural network

# MNIST

# MNIST

- One of the most popular dataset for image analysis. Contains hand-written digits for recognition.

- 60k training examples and 10k in test set.

- Early success of neural networks.

- Established efficacy of Convolution neural networks (ConvNets) in image recognition.

- LeNet-5 : Deep ConvNet by Yan LeCun in 1998. Used by several banks to recongnize hand-written numbers of checks.

- Many new deep ConvNet architectures have been proposed to improve performance on this dataset.

- SOTA : Dynamic Routing Between Capsules (Sabour, Frost, Hinton) Nov 2017

INPUT 32x32

C1: feature maps 6@28x28

C3: f. maps 16@10x10

S2: f. maps 6@14x14

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions

Subsampling

Convolutions

Subsampling

Full connection

Full connection

Gaussian connections

LeNet-5 architecture as published in the original paper.

# LENET - 5

- First major application of Convolutional Neural Networks for image recoginition.

- Provided foundation to all modern state-of-the-art computer vision models.

- tanh activations used.

- Averge pooling scheme.

- Not all convolutional filters were used by the whole image. For saving compute time.

- Achieves error rate of < 1%.

# IMAGENET



**The data that transformed AI research—and possibly the world**

# IMAGENET

- Data that changed everything. Big time!!

- Fei Fei Li and team in 2009.

- Previous datasets didn't capture the variability of the real world.

- Even identifying pictures of cats was infinitely complex.

- **WordNet:** a hierarchal structure for the English language. Like a dictionary, but words would be shown in relation to other words rather than alphabetical order.

- **Idea:** WordNet could have an image associated with each of the words, more as a reference rather than a computer vision dataset.

mammal ⟶ placental ⟶ carnivore ⟶ canine ⟶ dog ⟶ working dog ⟶ husky

- S: (n) Eskimo dog, **husky** (breed of heavy-coated Arctic sled dog)
    - *direct hypernym* / *inherited hypernym* / *sister term*
        - S: (n) working dog (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
            - S: (n) dog, domestic dog, Canis familiaris (a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) *"the dog barked all night"*
                - S: (n) canine, canid (any of various fissiped mammals with nonretractile claws and typically long muzzles)
                    - S: (n) carnivore (a terrestrial or aquatic flesh-eating mammal) *"terrestrial carnivores have four or five clawed digits on each limb"*
                        - S: (n) placental, placental mammal, eutherian, eutherian mammal (mammals having a placenta; all mammals except monotremes and marsupials)
                            - S: (n) mammal, mammalian (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
                                - S: (n) vertebrate, craniate (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
                                    - S: (n) chordate (any animal of the phylum Chordata having a notochord or spinal column)
                                        - S: (n) animal, animate being, beast, brute, creature, fauna (a living organism characterized by voluntary movement)
                                            - S: (n) organism, being (a living thing that has (or can develop) the ability to act or function independently)
                                                - S: (n) living thing, animate thing (a living (or once living) entity)
                                                    - S: (n) whole, unit (an assemblage of parts that is regarded as a single entity) *"how big is that part compared to the whole?"; "the team is a unit"*
                                                        - S: (n) object, physical object (a tangible and visible entity; an entity that can cast a shadow) *"it was full of rackets, balls and other objects"*
                                                            - S: (n) physical entity (an entity that has physical existence)
                                                                - S: (n) entity (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

IMAGENET

The ImageNet hierarchy derived from WordNet.

cf: *http://www.image-net.org/papers/imagenet_cvpr09.pdf*

# IMAGENET

- WordNet contains approximately 100,000 phrases and ImageNet has provided around 1000 images on average to illustrate each phrase.

- It consisted of 3.2 million labelled images, separated into 5,247 categories, sorted into 12 subtrees like "mammal," "vehicle," and "furniture."

- Originally published as a poster in CVPR, without attracting much fanfare.

- "There were comments like 'If you can't even do one object well, why would you do thousands, or tens of thousands of objects?" - Jia Deng, co-creator of Imagenet.

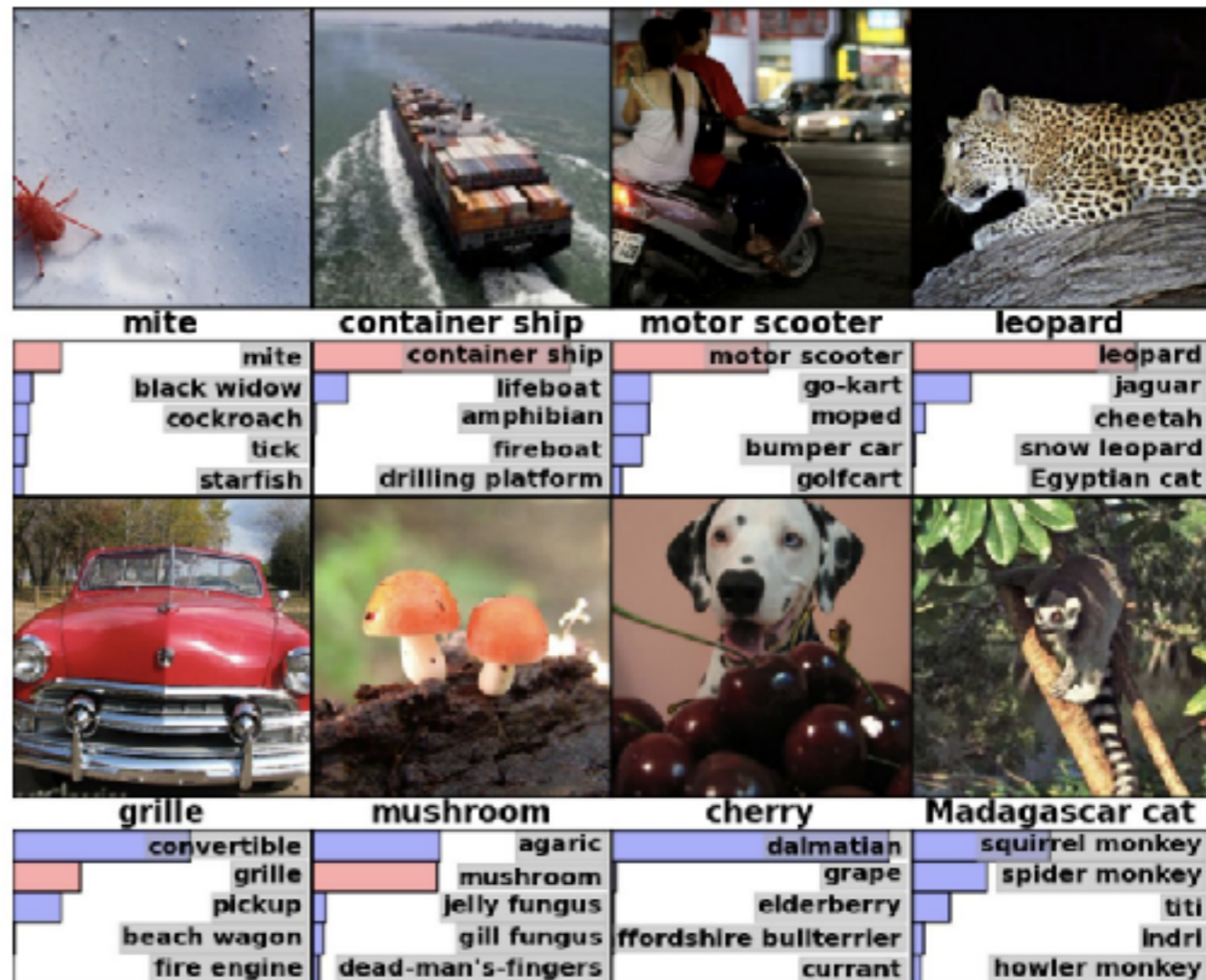*If data is the new oil, it was still dinosaur bones in 2009.*

# IMAGENET Competition

# IMAGENET Competition

- The ImageNet Large Scale Visual Recognition Challenge (ILSVRC) based on the data in Imagenet opened in 2010.

- Soon became a benchmark for how well image classification algorithms fared against the most complex visual dataset assembled at the time.

- Algorithms performed better when trained on Imagenet.

- Competition ran for 8 years.

- In 2012, the deep neural network submitted by Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton performed 41% better than the next best competitor, demonstrating that deep learning was a viable strategy for machine learning

- Accuracy went from 25% to < 5% through the course of ILSVRC.

ImageNet Large Scale Visual Recognition Challenge results

# IMAGENET

- If the deep learning boom we see today could be attributed to a single event, it would be the announcement of the 2012 ImageNet challenge results.

- Deep learning research exploded.

- Imagenet went from a poster on CVPR to benchmark of most of the presented papers today.

- "It was so clear that if you do a really good on ImageNet, you could solve image recognition," - Ilya Sutskever

- Without Imagenet, the deep learning revolution would have been delayed.

- After LeNet-5 for reading handwritten cheques, deep ConvNets (and Hinton?) needed a much bigger data to be useful in the real world.

# Progress in Image Recognition

# AlexNet

- ILSVRC 2012.

- First successful use of deep convnet for large scale image classification. Possible because of large amounts of labelled data from ImageNet as well as computations on 2 GPUs.

- **ReLU** non-linearity activation functions, finding that they performed better and decreased training time relative to the *tanh* function. The ReLU non-linearity now tends to be the default activation function for deep networks.

- **Data augmentation** techniques that consisted of image translations, horizontal reflections, and mean subtraction. They techniques are very widely used today for many computer vision tasks.

- **Dropout** layers in order to combat the problem of overfitting to the training data.

- Proposed style of having **successive convolution** and **poolinglayers**, followed by **fully-connected** layers at the end is still the basis of many state-of-the-art networks today.
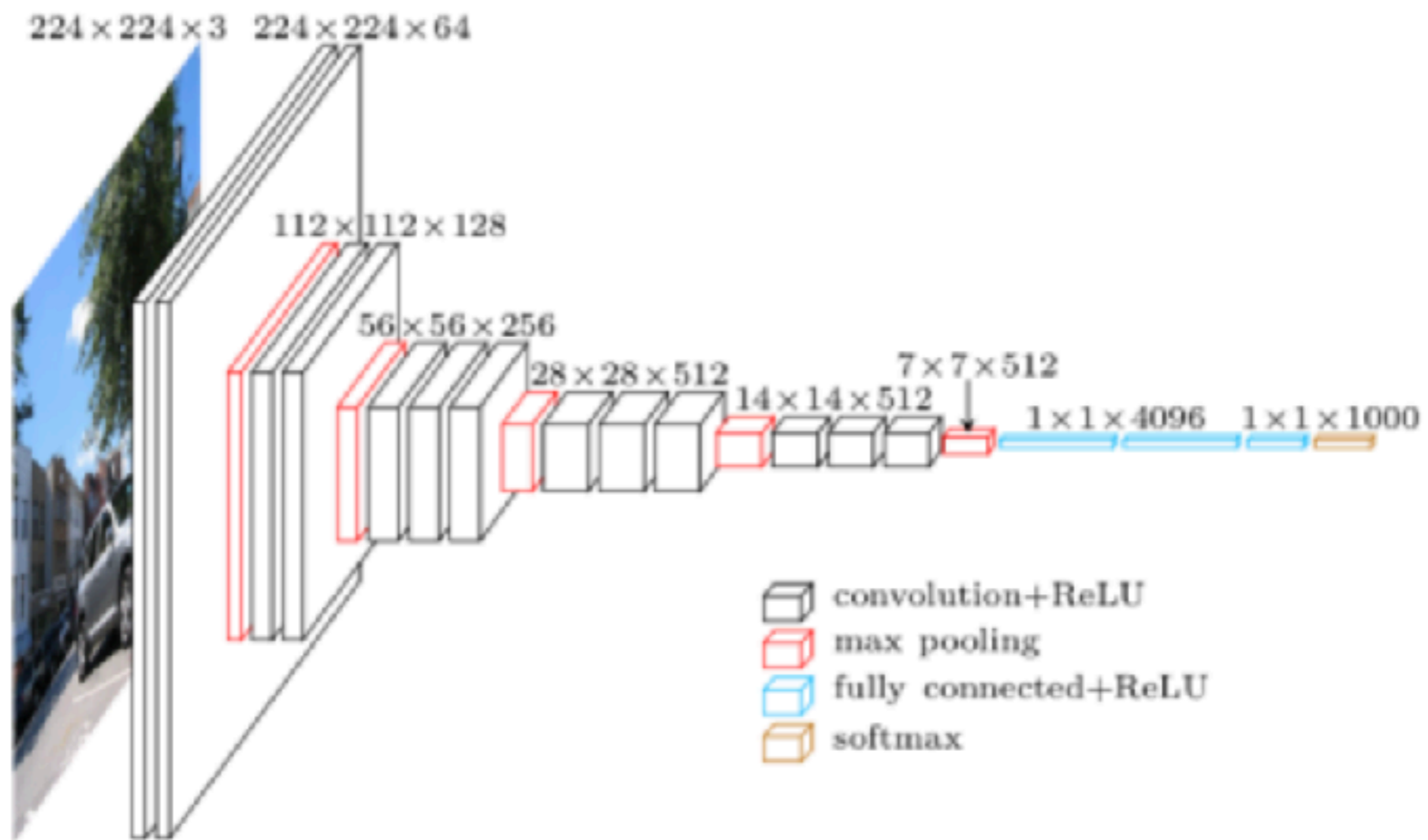
# Clarifai

- ILSVRC 2013.

- Matthew Zeiler, a PhD student at NYU in 2013 and Rob Fergus won the 2013 competition.

- Matthew Zeiler built Clarifai based off his 2013 ImageNet win, and is now backed by $40 million in VC funding.

- "widely seen as one of the most promising [startups] in the crowded, buzzy field of machine learning." (Forbes)

- Also, OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks by Yan LeCun and team at NYU.
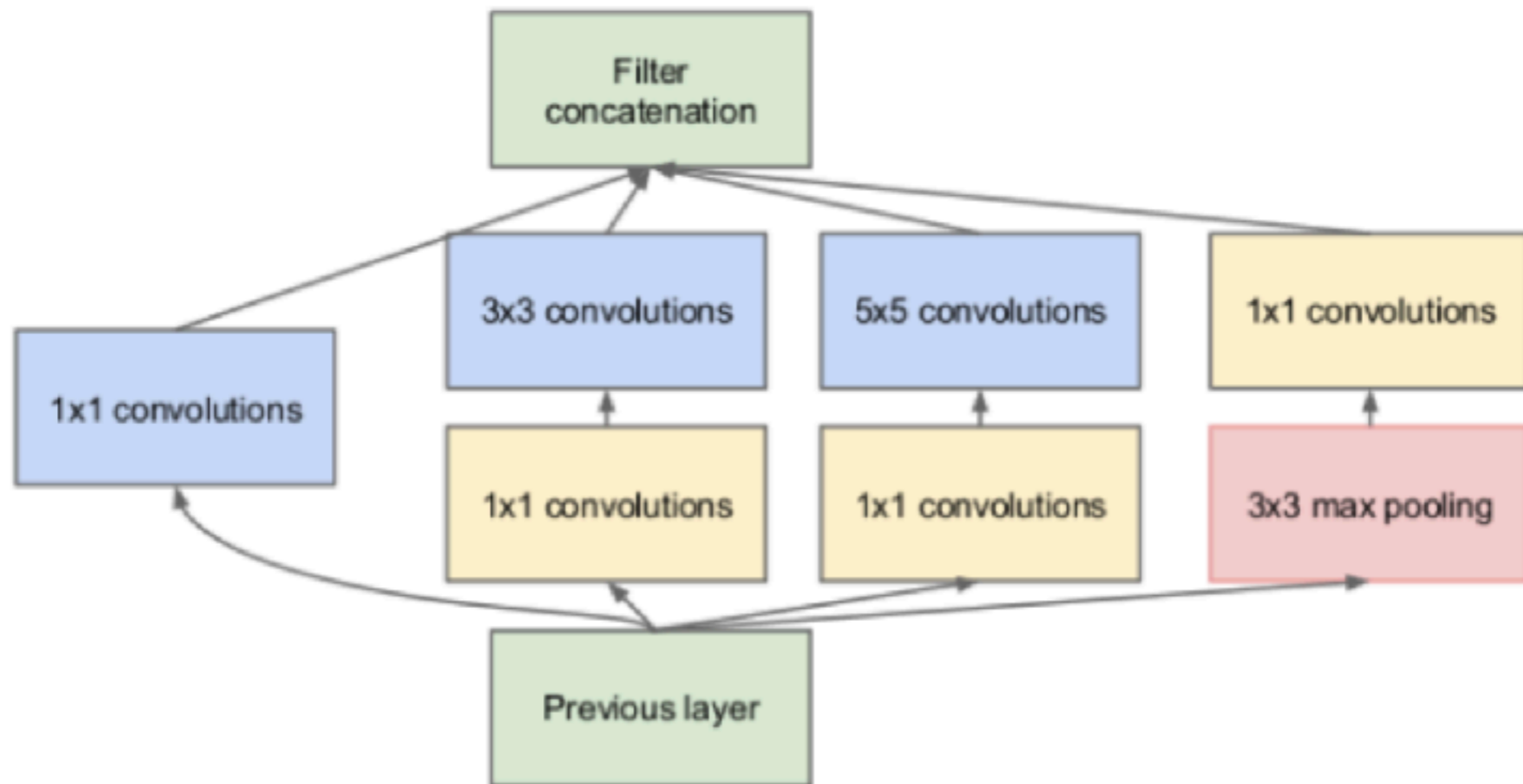
VGGNet Architecture

# VGG

- Visual Geometry Group @ Oxford, UK.

- **Idea:** You didn't really need any fancy tricks to get high accuracy. Just a deep network with lots of small 3x3 convolutions and non-linearities will do the trick!

- Two successive 3x3 convolutions has the equivalent receptive field i.e. the pixels it sees as a single 5x5 filter, and 3 3x3 filter ~ a 7x7 filter. Same stimulation of pixels with added benefits of smaller filters.

    - Decrease in the number of parameters.

    - Using a ReLU function in-between each convolution introduces more non-linearity into the network which makes modeling decision function better.

- As the spatial size of the input volumes at each layer decrease (as a result of the pooling layers), the depth of the volumes increase since need more discriminative features to use for accurate classification.

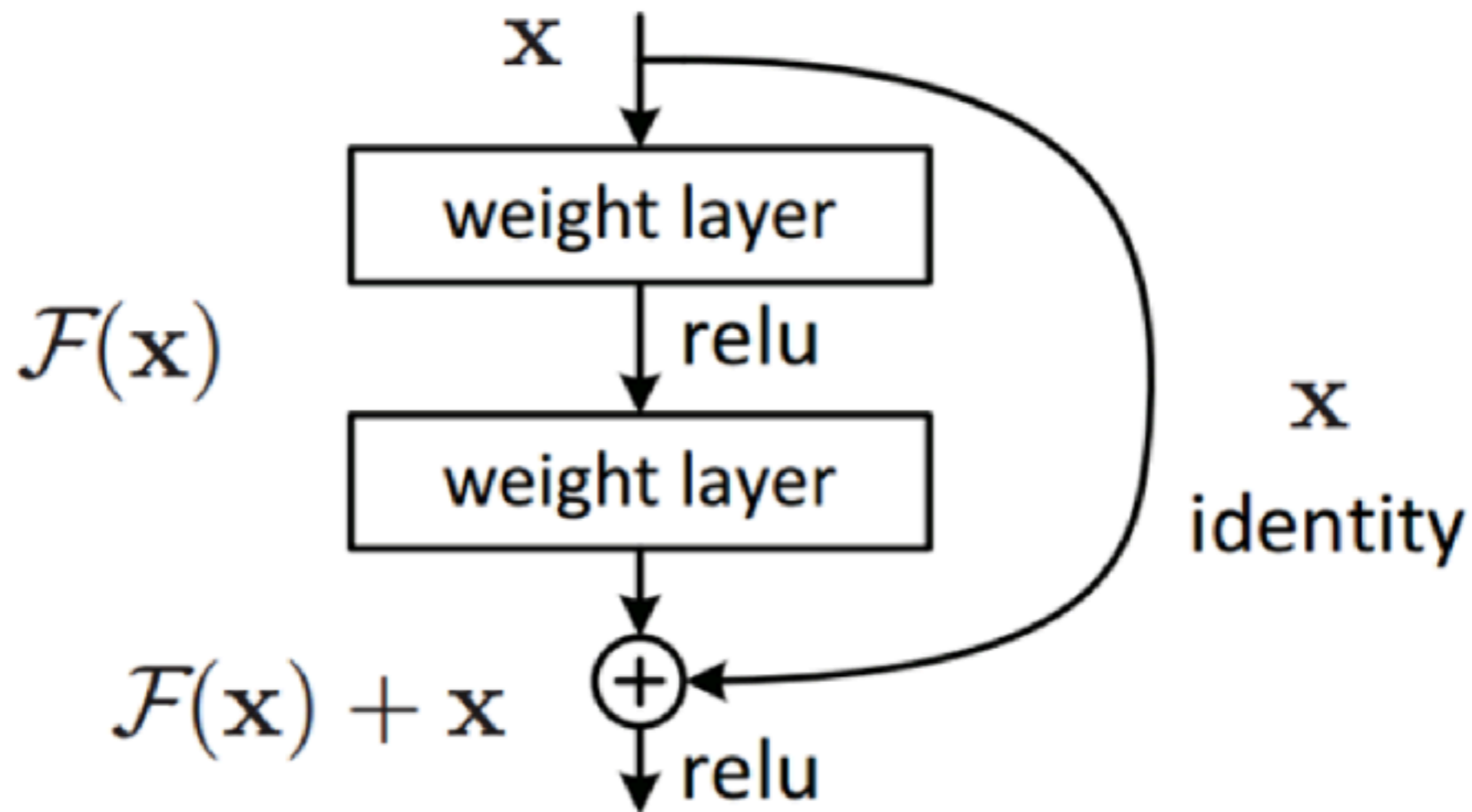- New kind of data augmentation: scale jittering.

Inception Module from GoogLeNet

# GoogleLeNet and Inception

- First to really address the issue of computational resources along with multi-scale processing.

- Through the use of 1x1 convolutions before each 3x3 and 5x5, the inception module reduces the number of feature maps passed through each layer, thus reducing computations and memory consumption.

- The inception module has 1x1, 3x3, and 5x5 convolutions all in parallel. The idea behind this was to let the network decide, through training what information would be learned and used.

- GoogLeNet was one of the first models that introduced the idea that CNN layers didn't always have to be stacked up sequentially. The authors of the paper showed that you can also increase network width for better performance and not just depth.

## Skipping over with a shortcut: ResNet



A residual block

Residual block from ResNet

# ResNet

- The ResNet architecture was the first to pass human level performance on ImageNet.

- Main contribution of *residual learning* is often used by default in many state-of-the-art networks today.

- A naive stacking of layers to make the network very deep won't always help and can actually make things worse.

- The idea is that by using an additive skip connection as a shortcut, deep layers have direct access to features from previous layers. The allows feature information to more easily be propagated through the network. It also helps with training as the gradients can also more efficiently be back-propagated.

- The first "*ultra deep*" network, where it is common to use over 100–200 layers.

# Impacts of IMAGENET

- **Transfer learning:** researchers soon realized that the weights learned in state of the art models for ImageNet could be used to initialize models for completely other datasets and improve performance significantly.

- Achieving good performance with as little as one positive example per category.

- Pre-trained ImageNet models have been used to achieve state-of-the-art results in tasks such as

    - object detection

    - semantic segmentation

    - human pose estimation

    - video recognition.

- Applications in domains where the number of training examples is small and annotation is expensive.  (e.g. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition)

*"One thing ImageNet changed in the field of AI is suddenly people realized the thankless work of making a dataset was at the core of AI research. People really recognize the importance the dataset is front and center in the research as much as algorithms."*

*– Fei Fei Li, Creator or ImageNet and Chief Scientist Google Cloud.*

# Beyond IMAGENET
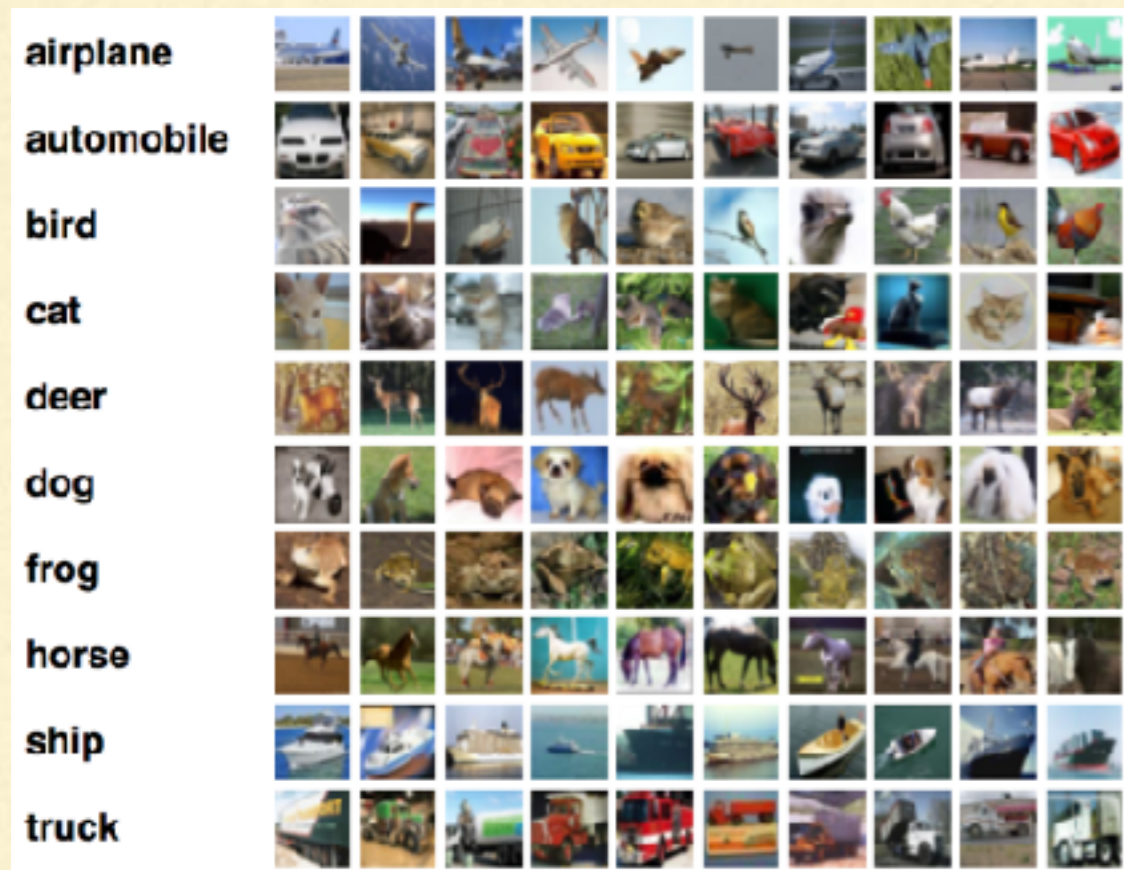
- Google released the Open Images Database, containing 9 million images in 6000 categories.

- YouTube-8M Dataset - to accelerate research on large-scale video understanding, representation learning, noisy data modeling, transfer learning, and domain adaptation approaches for video.

- Visual Genome - A knowledge database to connect structured image concepts to language.

- Fine-tuned ImageNet models are pushing the state of the art in many traditional image datasets like CIFAR, PASCAL etc.

# CIFAR



- For image classification. Beyond digits.

- 60k color images with 10 classes.

- Piggy-banking on the success of IMAGENET and deep convnets.

- SOTA : ShakeDrop regularization 2018

| Method | Reg | Cos | Fil | Depth | #Param | CIFAR -10 (%) | CIFAR -100 (%) |
|---|---|---|---|---|---|---|---|
| Coupled Ensemble (Dutt et al., 2017) | | | | 118 | 25.7M | *2.99 | *16.18 |
| | | | | 106 | 25.1M | *2.99 | *15.68 |
| | | | | 76 | 24.6M | *2.92 | *15.76 |
| | | | | 64 | 24.9M | *3.13 | *15.95 |
| | | | | - | 50M | *2.72 | *15.13 |
| | | | | - | 75M | *2.68 | *15.04 |
| | | | | - | 100M | *2.73 | *15.05 |
| ResNeXt (Xie et al., 2017) | | ✓ | | 26 | 26.2M | +3.58 | - |
| | | | | 29 | 34.4M | - | +16.34 |
| ResNeXt + Shake-Shake (Gastaldi, 2017) | SS | ✓ | | 26 | 26.2M | *2.86 | - |
| | | | | 29 | 34.4M | - | *15.85 |
| ResNeXt + Shake-Shake + Cutout (DeVries & Taylor, 2017b) | SS | ✓ | CO | 26 | 26.2M | *2.56 | - |
| | | | | 29 | 34.4M | - | *15.20 |
| PyramidNet (Han et al., 2017b) | | | | 272 | 26.0M | *3.31 | *16.35 |
| | | ✓ | RE | 272 | 26.0M | 3.42 | 16.66 |
| PyramidDrop (Yamada et al., 2016) | RD | | | 272 | 26.0M | 3.83 | 15.94 |
| | RD | ✓ | RE | 272 | 26.0M | 2.91 | 15.48 |
| PyramdNet + ShakeDrop (Proposed) | SD | | | 272 | 26.0M | 3.41 | **14.90** |
| | SD | | RE | 272 | 26.0M | 2.89 | **13.85** |
| | SD | ✓ | | 272 | 26.0M | 2.67 | **13.99** |
| | SD | ✓ | RE | 272 | 26.0M | **2.31** | **12.19** |

**1** **Upload photo**

The first picture defines the scene you would like to have painted.

**2** **Choose style**

Choose among predefined styles or upload your own style image.

**3** **Submit**

Our servers paint the image for you. You get an email when it's done.

Thanks

@januverma