# ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
Ilya Sutskever
Geoffrey Hinton

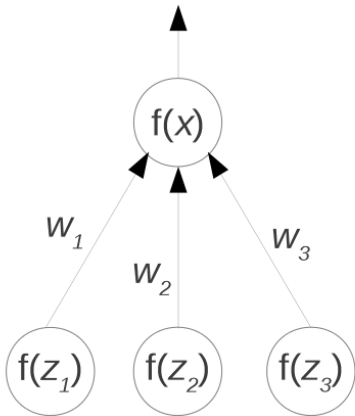University of Toronto
Canada

**Main idea**

Architecture
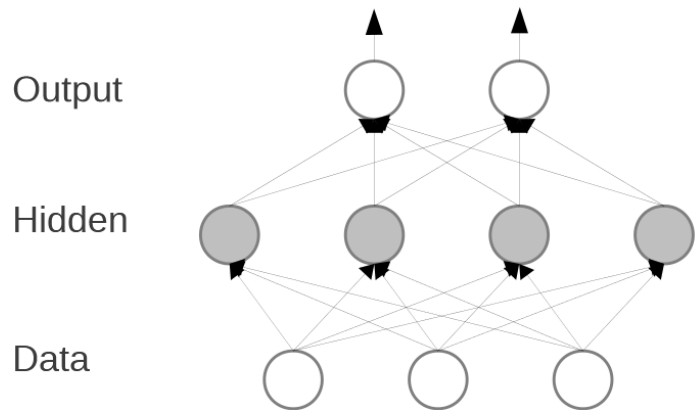Technical details

# Neural networks

- A neuron



$$x = w_1 f(z_1) + w_2 f(z_2) + w_3 f(z_3)$$

$x$ is called the total input
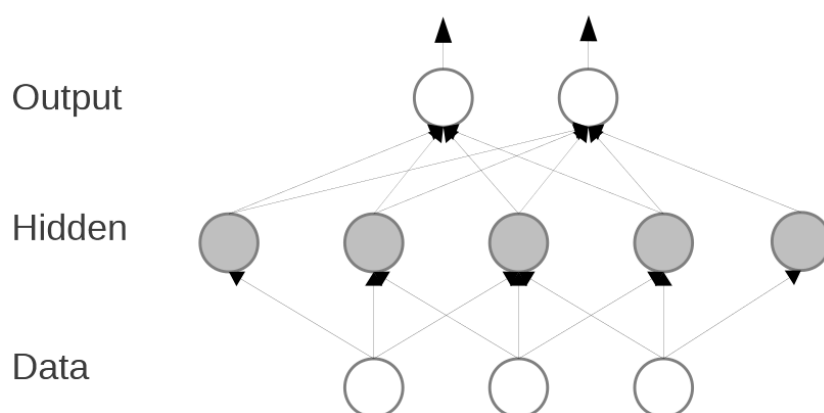to the neuron, and f($x$)
is its output

- A neural network



A neural network computes a differentiable
function of its input. For example, ours computes:
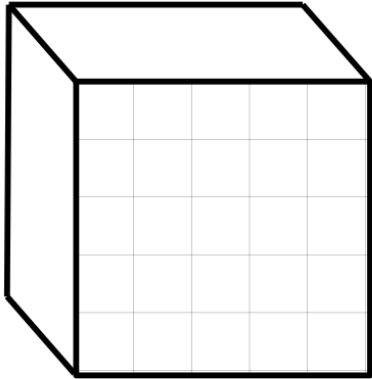$p$(label | an input image)

# Convolutional neural networks

- Here's a one-dimensional convolutional neural network

- Each hidden neuron applies **the same localized, linear filter** to the input
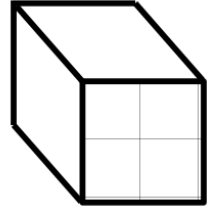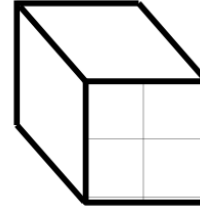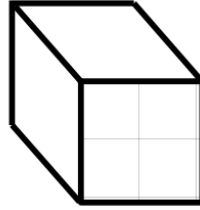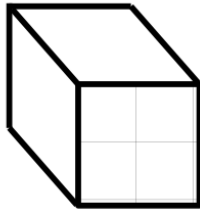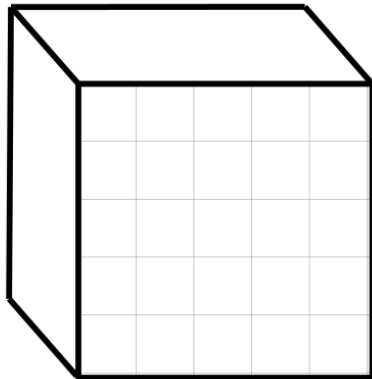
# Convolution in 2D

Input "image"

Filter bank

Output map

# Local pooling

# Overview of our model

- **Deep**: 7 hidden "weight" layers

- **Learned**: all feature extractors initialized at white Gaussian noise and learned from the data

- Entirely supervised

- **More data = good**

○ **Convolutional layer:** convolves its input with a bank of 3D filters, then applies point-wise non-linearity

□ **Fully-connected layer:** applies linear filters to its input, then applies point-wise non-linearity

Image

# Overview of our model

- Trained with stochastic gradient descent on two NVIDIA GPUs for about a week
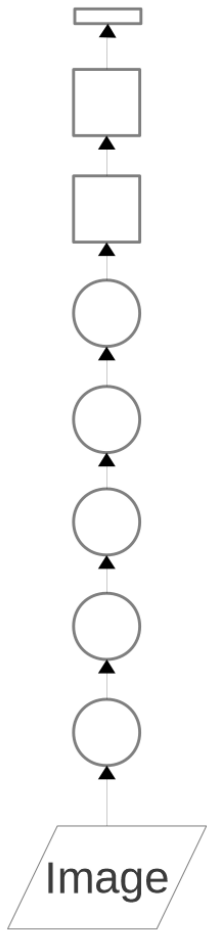
- 650,000 neurons

- 60,000,000 parameters

- 630,000,000 connections

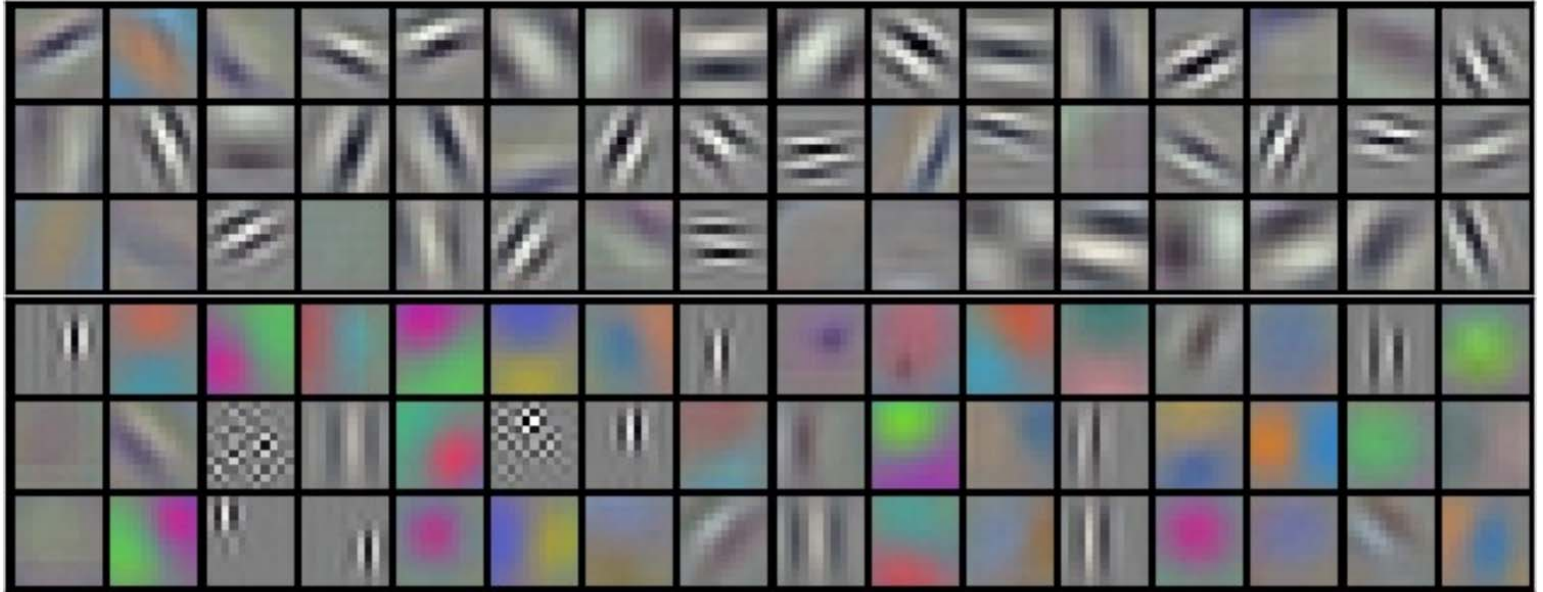- **Final feature layer:** 4096-dimensional

**Convolutional layer:** convolves its input with a bank of 3D filters, then applies point-wise non-linearity

**Fully-connected layer:** applies linear filters to its input, then applies point-wise non-linearity

Image

# 96 learned low-level filters

Main idea

→ **Architecture**

Technical details

# Training

Forward pass

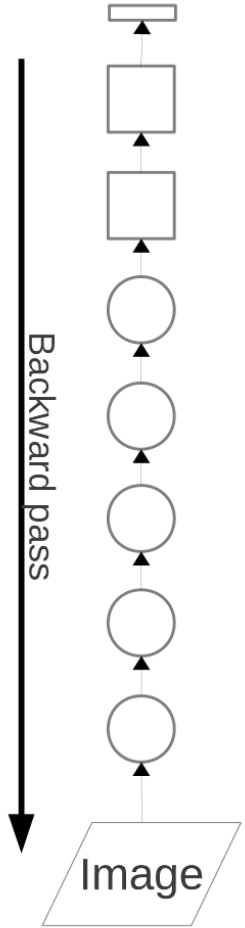Backward pass

Using stochastic gradient descent and the *backpropagation algorithm* (just repeated application of the chain rule)

Image

Image

One output unit per class

$x_i$ = total input to output unit $i$

$$f(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^{1000} \exp(x_j)}$$

We maximize the log-probability of the correct label, $\log f(x_t)$

# Our model

- Max-pooling layers follow first, second, and fifth convolutional layers

- The number of neurons in each layer is given by 253440, 186624, 64896, 64896, 43264, 4096, 4096, 1000

Main idea
Architecture
**Technical details**

# Input representation

- Centered (0-mean) RGB values.



An input image (256x256)          Minus sign          The mean input image

# Neurons

$f(x) = \tanh(x)$



$f(x) = \max(0, x)$



$f(x)$

$w_1$

$w_2$

$w_3$

$f(z_1)$

$f(z_2)$

$f(z_3)$

$x = w_1 f(z_1) + w_2 f(z_2) + w_3 f(z_3)$

$x$ is called the total input to the neuron, and $f(x)$ is its output

Very bad (slow to train)

Very good (quick to train)

# Data augmentation

- Our neural net has 60M real-valued parameters and 650,000 neurons

- It overfits a lot. Therefore we train on 224x224 patches extracted randomly from 256x256 images, and also their horizontal reflections.
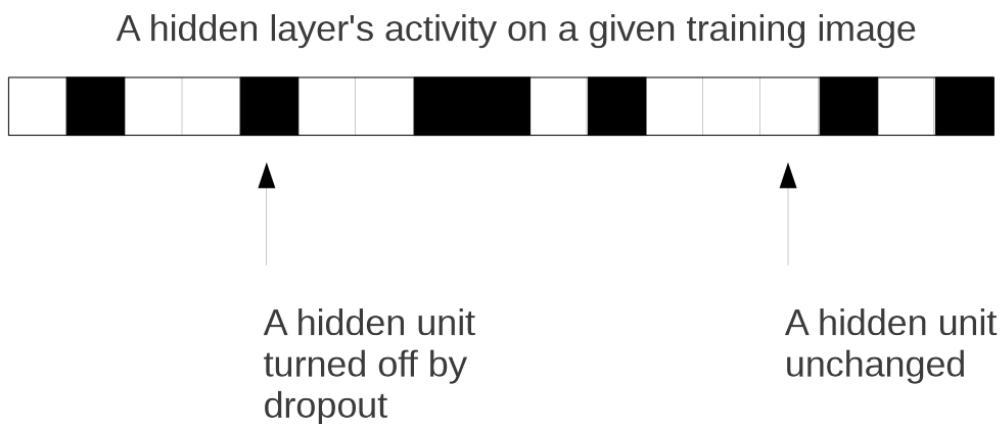
# Testing

- Average predictions made at five 224x224 patches and their horizontal reflections (four corner patches and center patch)

- Logistic regression has the nice property that it outputs a probability distribution over the class labels

- Therefore no score normalization or calibration is necessary to combine the predictions of different models (or the same model on different patches), as would be necessary with an SVM.

# Dropout

- Independently set each hidden unit activity to zero with 0.5 probability

- We do this in the two globally-connected hidden layers at the net's output

A hidden layer's activity on a given training image

A hidden unit
turned off by
dropout

A hidden unit
unchanged

# Implementation

- The only thing that needs to be stored on disk is the raw image data

- We stored it in JPEG format. It can be loaded and decoded entirely in parallel with training.

- Therefore only 27GB of disk storage is needed to train this system.

- Uses about 2GB of RAM on each GPU, and around 5GB of system memory during training.

# Implementation

- Written in Python/C++/CUDA

- Sort of like an instruction pipeline, with the following 4 instructions happening in parallel:

  - Train on batch $n$ (on GPUs)

  - Copy batch $n+1$ to GPU memory

  - Transform batch $n+2$ (on CPU)

  - Load batch $n+3$ from disk (on CPU)

# Validation classification



| mite | container ship | motor scooter | leopard |
|------|---------------|---------------|---------|
| mite | container ship | motor scooter | leopard |
| black widow | lifeboat | go-kart | jaguar |
| cockroach | amphibian | moped | cheetah |
| tick | fireboat | bumper car | snow leopard |
| starfish | drilling platform | golfcart | Egyptian cat |

| grille | mushroom | cherry | Madagascar cat |
|--------|----------|--------|----------------|
| convertible | agaric | dalmatian | squirrel monkey |
| grille | mushroom | grape | spider monkey |
| pickup | jelly fungus | elderberry | titi |
| beach wagon | gill fungus | Staffordshire bullterrier | indri |
| fire engine | dead-man's-fingers | currant | howler monkey |

# Validation classification



| lens cap | abacus | slug | hen |
|---|---|---|---|
| reflex camera | abacus | slug | hen |
| Polaroid camera | typewriter keyboard | zucchini | cock |
| pencil sharpener | space bar | ground beetle | cocker spaniel |
| switch | computer keyboard | common newt | partridge |
| combination lock | accordion | water snake | English setter |

| tiger | chambered nautilus | tape player | planetarium |
|---|---|---|---|
| tiger | lampshade | cellular telephone | planetarium |
| tiger cat | throne | slot | dome |
| tabby | goblet | reflex camera | mosque |
| boxer | table lamp | dial telephone | radio telescope |
| Saint Bernard | hamper | iPod | steel arch bridge |

# Validation classification



| koala | tiger | European fire salamander | loggerhead |
|---|---|---|---|
| wombat | tiger | European fire salamander | African crocodile |
| Norwegian elkhound | tiger cat | spotted salamander | Gila monster |
| wild boar | jaguar | common newt | loggerhead |
| wallaby | lynx | long-horned beetle | mud turtle |
| koala | leopard | box turtle | leatherback turtle |

| seat belt | television | sliding door | wallaby |
|---|---|---|---|
| seat belt | television | sliding door | hare |
| ice lolly | microwave | shoji | wallaby |
| hotdog | monitor | window shade | wood rabbit |
| burrito | screen | window screen | Lakeland terrier |
| Band Aid | car mirror | four-poster | kit fox |

# Validation localizations



| bookshop | coyote | cradle | wood rabbit |
|---|---|---|---|
| balance beam | grey fox | cradle | hare |
| cinema | kit fox | bassinet | wood rabbit |
| marimba | red fox | diaper | grey fox |
| parallel bars | coyote | crib | coyote |
| computer keyboard | dhole | bath towel | wallaby |

| bottlecap | harvester | garter snake | Walker hound |
|---|---|---|---|
| bottlecap | harvester | diamondback | beagle |
| magnetic compass | thresher | leatherback turtle | Walker hound |
| puck | plow | sandbar | English foxhound |
| stopwatch | tractor | echidna | muzzle |
| disk brake | tow truck | armadillo | Italian greyhound |

# Validation localizations



| chime | boathouse | Scottish deerhound | electric guitar |
|---|---|---|---|
| wine bottle | apiary | Scottish deerhound | violin |
| digital clock | mobile home | Irish wolfhound | carpenter's kit |
| gar | boathouse | Leonberg | revolver |
| oboe | picket fence | German shepherd | Loafer |
| typewriter keyboard | patio | Tibetan mastiff | corkscrew |

| motor scooter | sturgeon | violin | fire screen |
|---|---|---|---|
| motor scooter | leatherback turtle | violin | fire screen |
| moped | volcano | cello | sundial |
| snowmobile | wreck | acoustic guitar | mailbag |
| police van | alp | drumstick | umbrella |
| moving van | breakwater | electric guitar | purse |

# Retrieval experiments

First column contains query images from ILSVRC-2010 test set, remaining columns contain retrieved images from training set.

# Retrieval experiments