

(1)

## Zhang - Modern CNNs

AlexNet (Krizhevsky, 2012)

- ↳ first large scale network to beat conventional CV methods on a large scale

VGG network (Simonyan & Zisserman, 2014)

- ↳ make use of a number of repeating blocks of elements

network-in-network (Lin et al 2013)

- ↳ convolves whole NNs patchwise over inputs

GoogleLeNet (Szegedy et al 2015) multi-branch convos

ResNet (He et al. 2016) one of most pop off their architectures in CV

ResNext blocks (Xie et al 2017) for sparser connects

Densenet (Huang et al 2017) generalization of the residual architecture

(2)

## 8.1. Deep Convolutional Neural Networks (Alexnet)

CNNs did not dominate immediately after LeCun (1995) LeNet

LeNet was only established on small datasets

Pre-2012 neural networks were generally overshadowed by other methods

- Historically, CV practitioners would craft feature extraction & pass this as inputs into the models

- rather than learning

- feature extraction was the main goal  
→ the learning Algo was an after thought

- Historical pipeline:

- (1) Data collect

- (2) Preprocess data

- (3) Feed into feature extractors

- (4) Input into classifier model

## 8.1.1. Representation learning

Up until 2012 the most important part of the pipeline was representation calculated / crafted by hand

LeCun, Hinton, Bengio, Ng, Amari, Schmidhuber

believed features ought to be learned

First modern CNN is AlexNet

evolutionary improvement over LeNet

Main diff between LeNet & AlexNet is increase in computational power

### 8.1.1.1 Missing ingredient: Data

Trad methods are based on convex optim (linear & kernel)

Deep models need vast amounts of complex data to outperform this

ImageNet (Deng et al 2009) = 1 mil examples  
1000 from 1000 examples

### 8.1.1.2 Missing ingredient: Hardware

## 8.1.2 AlexNet

8 layer CNN

Showed that features obtained by learning can be manual features

① Image

② Conv w/ 96 filters

③ Max Pool

④ Conv

⑤ Max Pool

⑥ Conv  $\rightarrow$  Conv  $\rightarrow$  Conv

⑦ Max Pool

⑧ FC  $\rightarrow$  FC  $\rightarrow$  FC

AlexNet changed sigmoid activation to a simpler ReLU activation

### 8.1.4 Discussion

AlexNets improve on LeNets using dropout for accuracy & for ReLU for ease of training

AlexNets downfall is its MCP layers which are very comp heavy

### 8.2. Networks using Blocks (VGG)

Alex = evidence that Deep CNNs can achieve good results

but did not provide a general template to guide researchers in designing new networks

Following section introduces common concepts used in designing deep networks

Researchers used to think in terms of individual neurons

then whole layers

now to blocks

↳ repeating patterns of layers

Super modern approach is to repurpose pre-trained models called foundational models

VGG = Visual geometry group

Idea of using blocks first emerged from  
ox uni group

easy to implement repeated code structures  
using loops & subroutines

### 8.2.11. ~~blocks~~ VGG Blocks

Historically, the basic CNN sequence is:

- ① conv w/ padding
- ② ReLU for non-lin
- ③ Max pool to reduce reso

issue here is that the reso reduce  
has a limit so ~~res~~ before all params  
are used up

Slim & Zisserman (2014) idea was to use  
nilli convs between pooling

they wanted to know whether deep or  
wide networks were better

Deep & narrow = better

Stacking  $3 \times 3$  convs has been gold standard

Liu (2022) conv net for the 2020s revises  
this

## 8.2.2 VGG Network

VGG can be thought of in two ways

- Blocks of Convolutional & Pooling Sequences
- A FC layer

Conv layers do not change dimensions

example VGG:

conv → conv → conv → ... → FC

VGG is more computationally expensive than AlexNet

### Summary

VGG might be the first modern CNN

introduced

- (1) Blocks of multiple convolutions
- (2) Preference for deep & narrow

## 8.3 Network-in-Network (Nin)

LeNet, AlexNet & VGG all share common extract features exploiting spatial structure via a sequence of convs & pool layers

↳ finished w/ fully connected layers

these FC layers are a problem & require a lot of computation

NiN offer

(1) use  $1 \times 1$  convs to add local nonlinearities across all channel activations

(2) use global average pool to integrate across layers

$1 \times 1$  can be considered a fully connected layer @ each pixel

NiN has no FC layer @ the end

reduces params but may increase run time

## 8.4 Multi-Branch Networks (googlenet)

Mix of NIN, blocks (vgg) & conv kernels

lst to have clear distinct between:

① Stem (ingest)

② Body (data process)

③ head (Pred)

Stem = 2-3 convs to learn low level

body = conv blocks

head = maps obtained features to task  
(Classif, Seg, Det etc)

### key Contribution

Problem = Selecting convolutional kernels in

Prob is to select correct kernel:  $1 \times 1$  or  $3 \times 3$

Solution: to concat multibranch

Block in google net = Inception Block

4 (n) parallel branches each w/ diff spatial sizes

Also implement  $1 \times 1$  layers to reduce channels

If the branches are coordinated in the padding to give the same dimension output

Finally all Branches are concatenated

Google net is computationally complex

large num of relatively arbitrary hyperparams:

- ① number channels chosen
- ② num blocks prior to Dim redud

Reason for this is that tools for Auto network Ref or Design Explorer had not been implemented yet introduced

## 8.5 Batch normalization

training Deep NNs is difficult

Specifically the time to converge

Batch normalization helps w/ this

Batch combined w/ residual blocks makes it possible to routinely train networks w/ over 100 layers

Also acts as regularization

Generally we preprocess data before training

Standardization plays nicely w/ optimizers as it puts params on a similar scale

Batch is applied to individual layers

in each training iteration

① normalize the inputs by subtracting their mean & dividing by their standard deviation  
(based on current mini batch)

② Apply scale coeff & an offset to recover the lost degrees of freedom

Normalization + Batch statistics =  
Batch Normalization

this makes the size of batch chosen very important

- ↳ Can't work w/ mini batches of size 1

After standardizing, the resulting minibatch has zero mean & unit variance

A small constant is added to avoid ever dividing by 0

throughout Deep learning, noise leads to faster training & less overfitting  
- regularization

Batch norm layers are slightly diff for MCP vs convolutions

### fully connected

before activation & after affine transform

### Convolution

Also after conv but before activation

Diff is application on a per-channel basis across all locations

each channel has its own scale & shift params

(13)

In the context of convolutions, BN is well defined even w/ minibatch = 1

have locations across an image to average

leads to concept of layer normalization

like batch norm but one observation  
@ a time

## 8.6 Residual Networks: Resnet & Resnext

As we build deeper networks it is important to understand how adding layers can increase the complexity

More important

↳ More layers should make the network more expressive not just different

only if larger functions contain the smaller ones are we guaranteed that increasing them will increase the expressive power

At the heart of resnet is idea is that every additional layer should more easily contain the identity function as one of the elements

leads to residual block

### 8.6.2 Residual Blocks

in a normal block, the  $f(x)$  needs to be learned & passed into the activation func

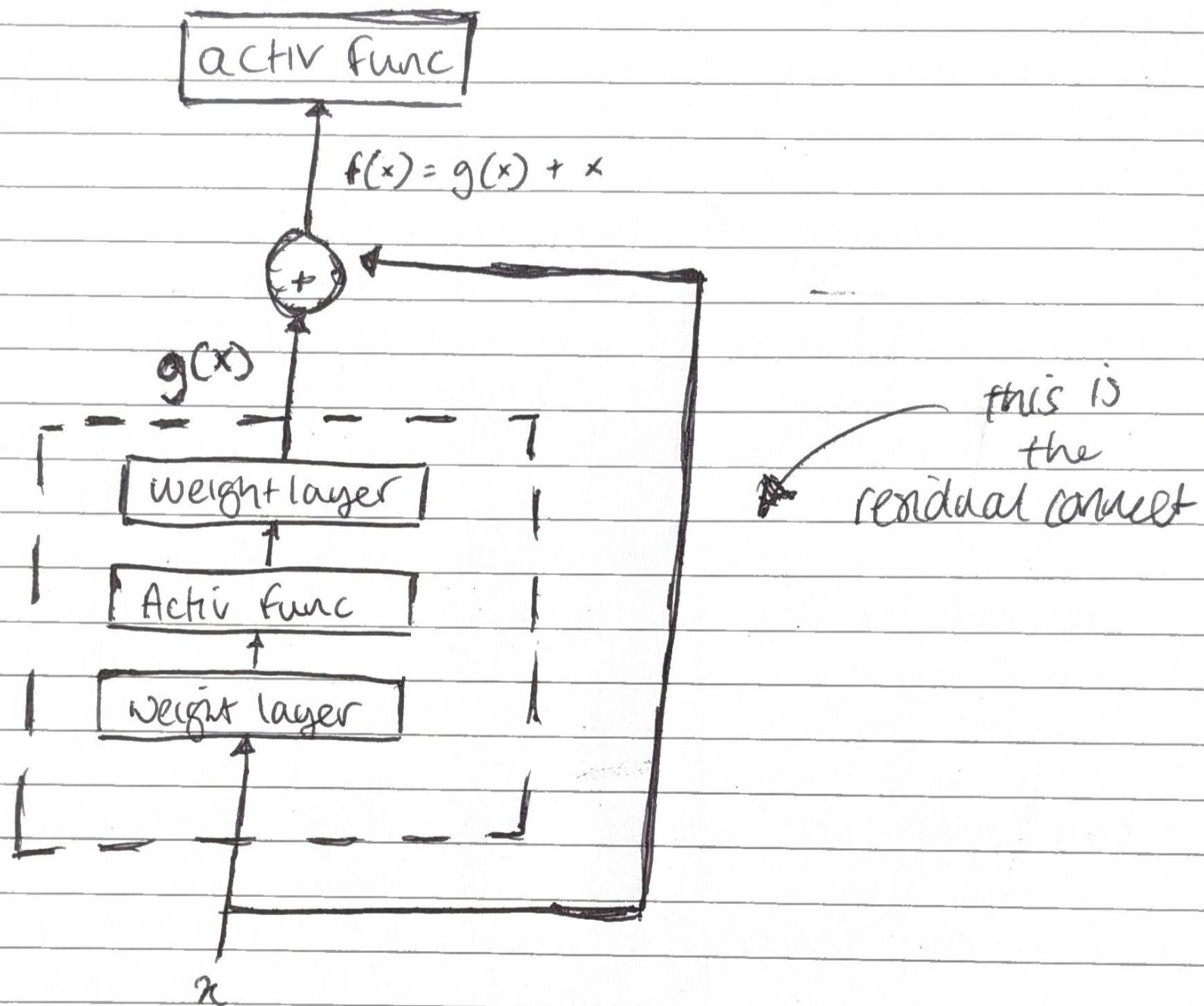
Residual block learns the residual map  $y(x) = f(x) - x$

$$y(x) = f(x) - x$$

(15)

If the ident mapping  $f(x) = x$  is the desire mapping

then the residual map =  $g(x) = 0$



### 8.6.3 ResNet Model

first two layers of resnet are the same as googlenet

After:

googlenet uses four modules made up of inception blocks

ResNet uses four modules made up of residual blocks  
 ↪ of which uses several residual blocks

it has 18 layers in total - ResNet-18

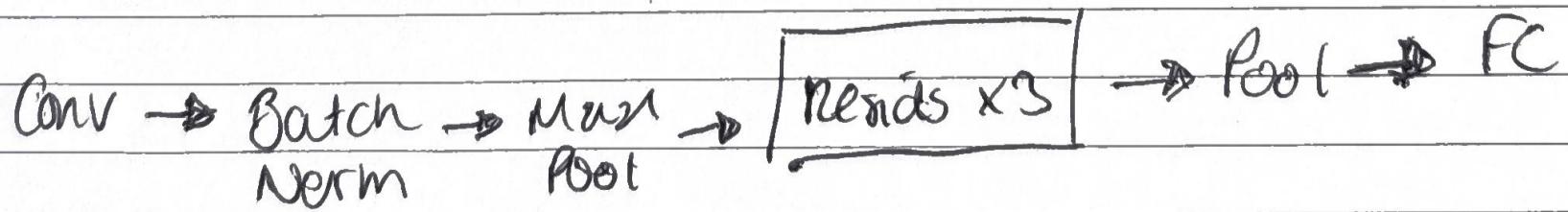
By config diff num of channels & residual blocks you get diff Resnet models

i.e. ResNet-152

Resnet architecture is similar to googlenet

↳ but simpler structure & easier to modify

Resnet 18



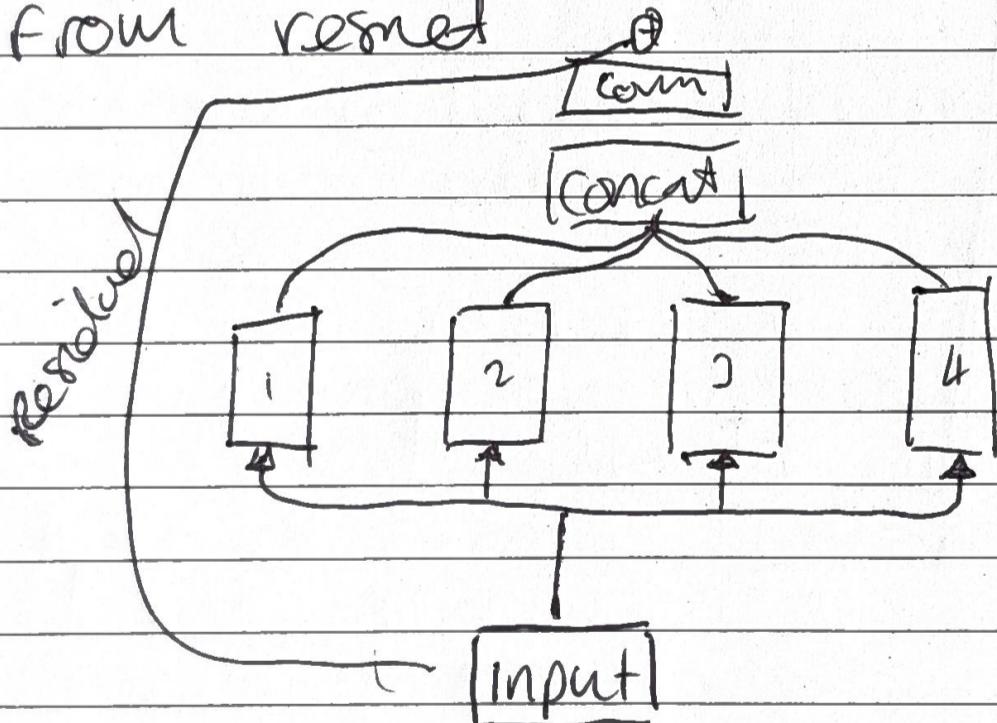
### 8.6.5 ResNet()

One of the challenges Resnet faces is the trade-off between non-linearity & dimensionality within a given block

Could add more non-linearity by increasing the number of layers

We can take inspo from Resnets which has information flowing through the block in separate groups

This idea is what lead to ResNext from Resnet



This structure reduces the number of params needed

initially the perceived challenge is that no info is shared between channels

ResNext amends this in two ways:

① grouped convolution

②

### Summary & Discussion

Nested functions classes are desirable as they allow us to obtain significantly more powerful classes when adding capacity

↳ rather than just being diff

Achieve by allowing layers to flow from input to output  
(residual connections)

this changes the inductive bias from

$$f(x) = 0 \quad \text{to} \quad f(x) = x$$

Residual Mapping can learn the identity function more easily

such as pushing params to zero

NNs w/ resids → inputs can forward prop faster through the next work

(14)

Allows for training deeper networks

original Resnet paper allowed for upto  
152 layers

## 8.7 Densely Connected Networks (DenseNet)

(20)

ResNet significantly changed the view of how to parametrize functions in Deep Nets

DenseNets is a logical extension

Skipped

## 8.8 Designing Convolutional Network Arch

Ever since AlexNet it has become popular to construct very deep networks by stacking blocks of convolutions

$3 \times 3$  convs were popularised by Vgg

NiN showed  $1 \times 1$  convs help add local non-linearities

• NiN also solved prob of aggregating info @ the head of a network by agg across all locations

Googlenet added multiple branches of diff conv widths

Resnets changes the inductive bias towards the identity mapping ( $\text{f}(\text{f}(\text{x})) = \text{x}$ ) allowing for very deep networks

upto now we have emitted networks obtained via neural architecture search (NAS)

Cost of doing so is very large