Why is there a smaller number of parameters?

— 5×5 filter: 5×5+1 (bias) = 26 parameters

2 layers of 3×3 → 2(3×3+1) = 2×10 = 20 parameters

— 7×7 filter: 7×7+1 (bias) = 50 parameters

3 layers of 3×3 → 3(3×3+1) = 3×10 = 30 parameters

↳ enough to cover the same pixel area / lattice

— 11×11 filter: 11×11+1 (bias) = 122 parameters

5 layers of 3×3 → 5(3×3+1) = 5×10 = 50 parameters

↳ enough to cover the same pixel area / lattice

However, cannot go deeper than 16/19 layers, as issues cropped up with gradients, which would only be solved later, with residual connections in the ResNet family (ResNet: Microsoft Research)

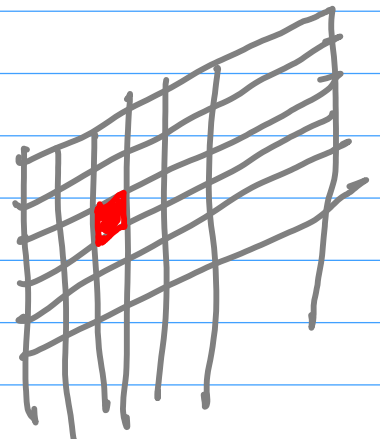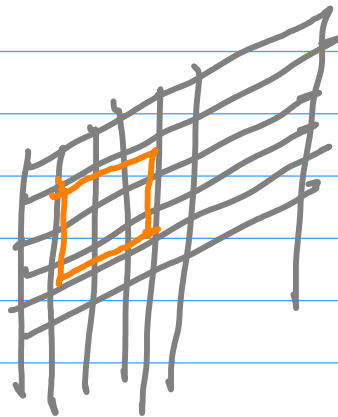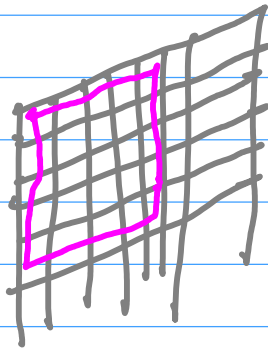Physical Significance: cost effectiveness: is in terms of the number of parameters of the convolution (# of weights to be tuned, to avoid overfitting issues) The number of multiplications/additions/ operations is not an important parameter for this application (we had also considered computing this as well for all layers of LeNet for instance: operation count)

Next up: ResNet (we have already covered the main points)

# Inception Architecture (Google): 2015 variants

GoogleNet , Batch Normalisation,
Factorising Convolutions

we will see what this is

we have seen this before

→ we have seen the basic philosophy in VGG-16/19

## Factorising Convolutions:



Stage I

Stage II
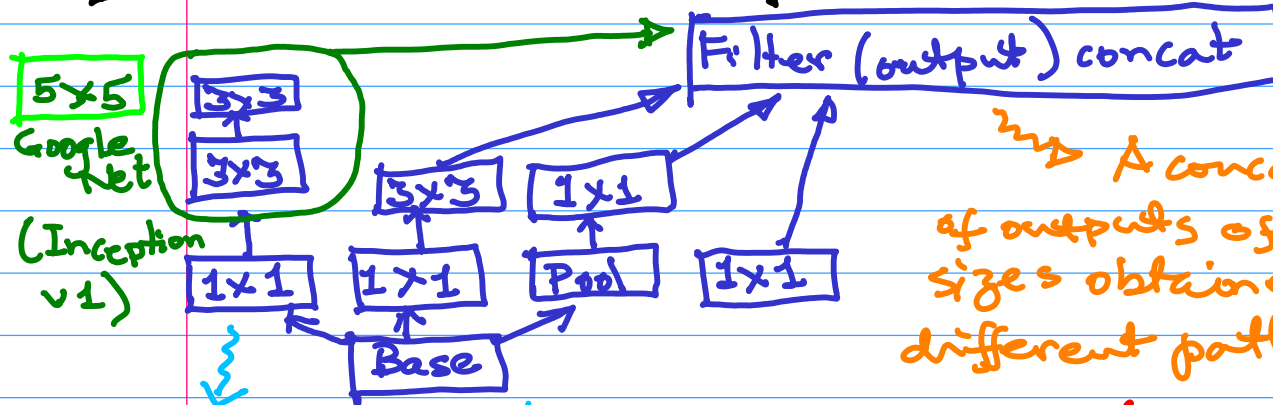
Hollywood Movie "Inception":
"dream within a dream"

## "Inception Module A"



5×5
GoogleNet
(Inception v1)

Filter (output) concat

⟿ A concatenation of outputs of different sizes obtained along different pathways

(To be discussed later)

→ concept of scales

this is a bit like looking at a signal across scales.
(gaussian pyramid): a concatenation of the outputs of the features in a signal at all possible scales → this information can be put to use later.
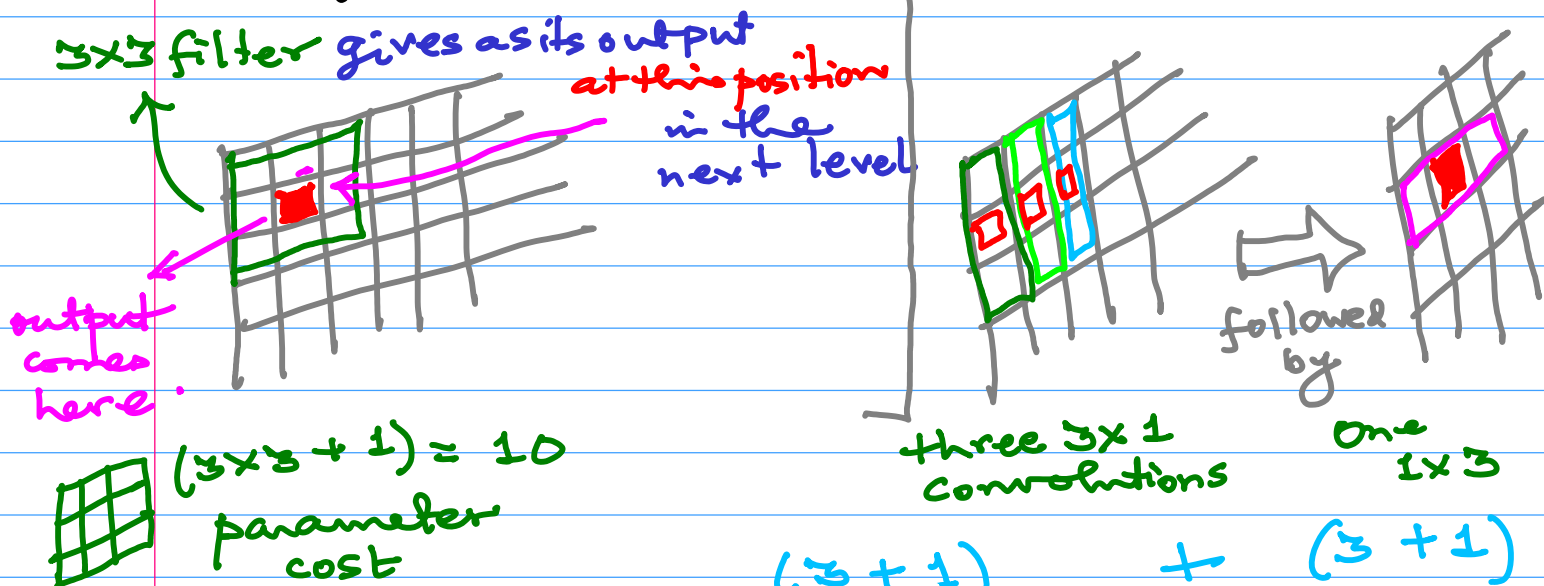
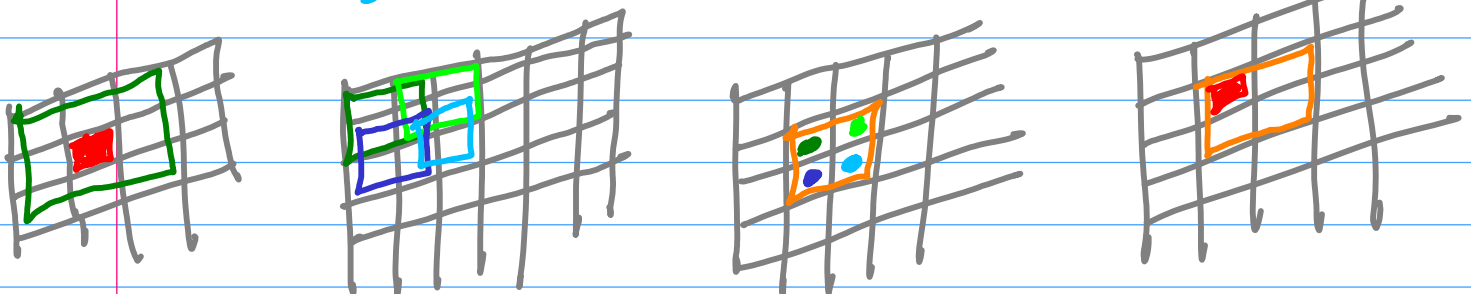Handcrafted feature detectors: SIFT/SURF

Scale Invariant Feature Transform (has information at various scales, which in turn leads to scale-invariance)

**Inception Module B**    Asymmetric Convolutions

One 3×3 convolution is replaced with two stages three 3×1 convolutions followed by one 1×3, or vice versa

3×3 filter gives as its output at this position in the next level

output comes here.

$(3\times3+1) = 10$ parameter cost

three 3×1 convolutions        followed by         One 1×3

$(3+1)$            +            $(3+1)$

$= 8$ (less)

Question: Why don't we use two 2×2 filters filter layers to replace one 3×3?
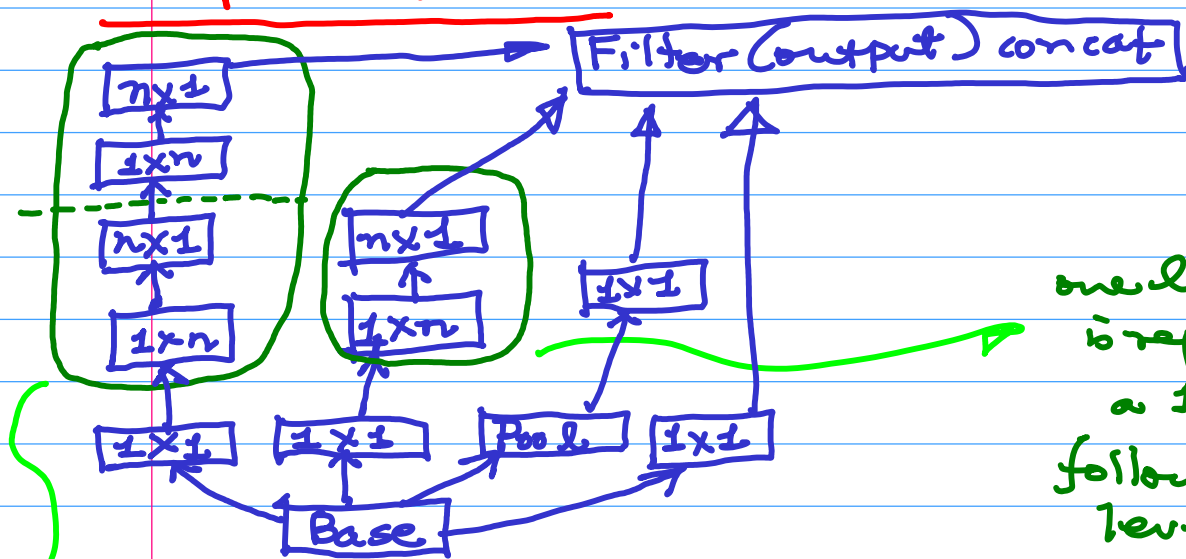
(after all, this was the basic VGG idea!)

parameter
count

$3 \times 3 + 1 = 10$

parameter count

$= 2 \times 2 + 1 \qquad 2 \times 2 + 1$

the total parameter count hasn't
changed!

# Inception Module B



one level of 7×7 is replaced with a 1×7 level followed by a 7×1 level

Take n=7
Two levels of 7×7:
each level of 7×7 is replaced
with a 1×7 level followed by
a 7×1 level

· what is the physical significance of the concatenated filter output? The output size is often larger (in cases, double as well)

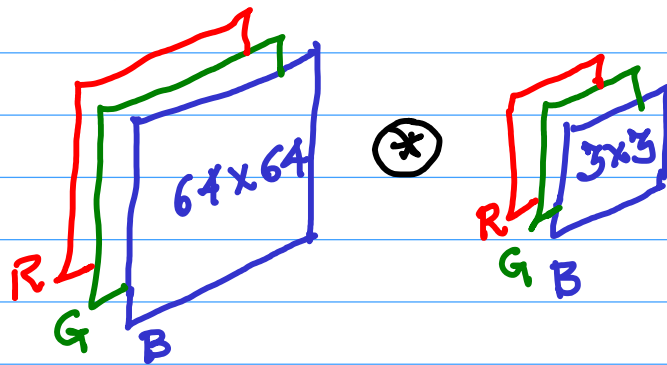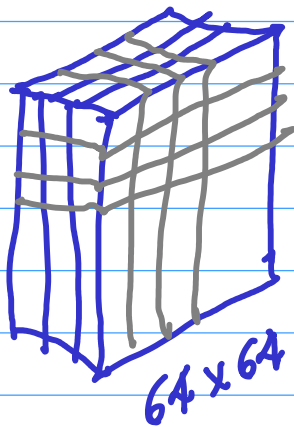→ features from all scales

e.g.,  | 17 × 17 × 640 |

↑ Inception ↑

| 17 × 17 × 320 |

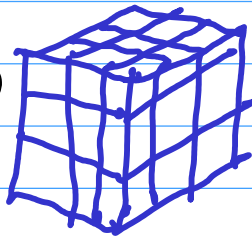The output image size is the same as the input image size but has a richer representation in terms of the number of channels.

what is a 1x1 convolution? seems to be pointless
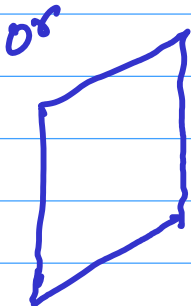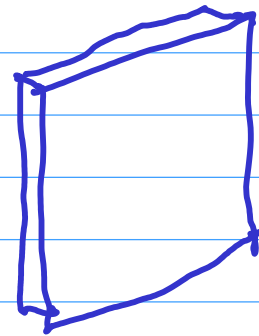(an image convolved with one value: plus a
bias parameter)

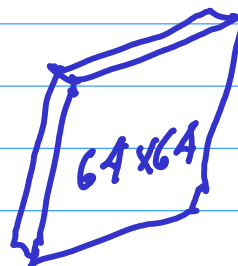if we want to create
1 pooled output

(in place of

$64 \times 64$ ⊛ $3 \times 3$ →

R
G B

R
G B

(for each pixel
position)

"Rubik's
cube"

or

$64 \times 64$

$64 \times 64 \times 3$
(RGB)

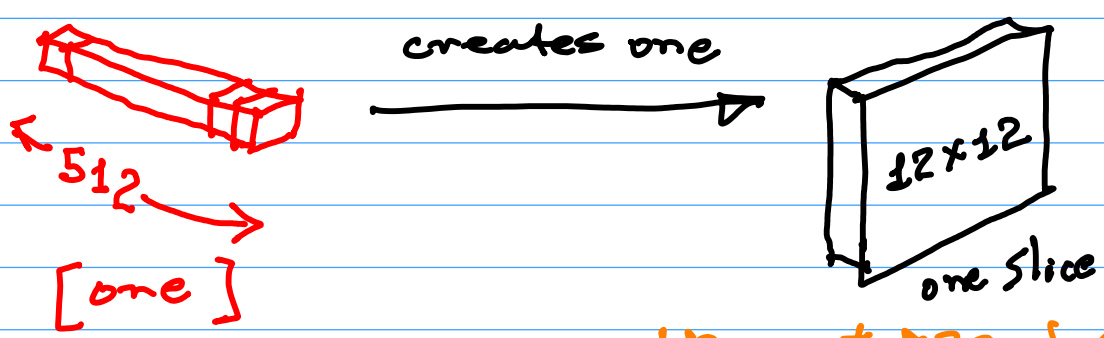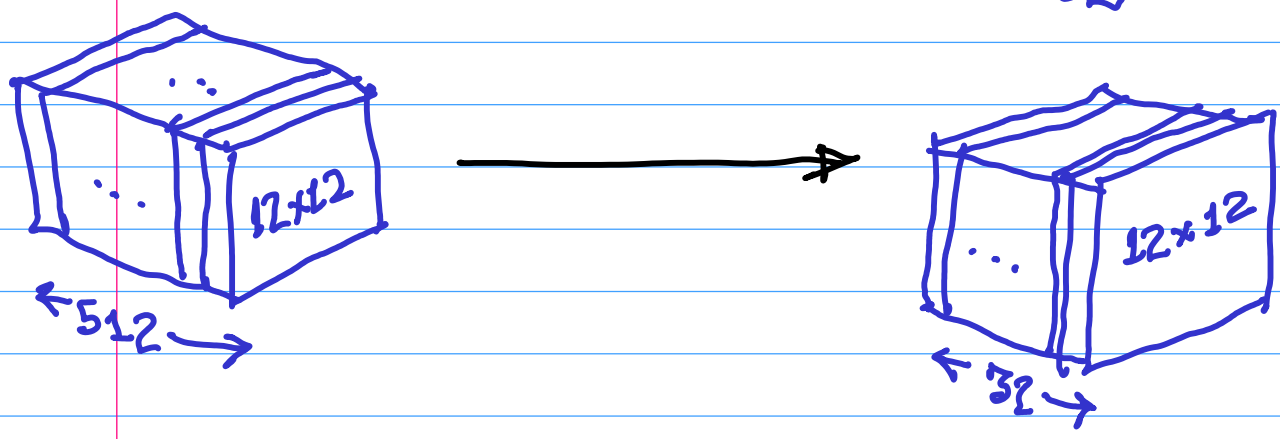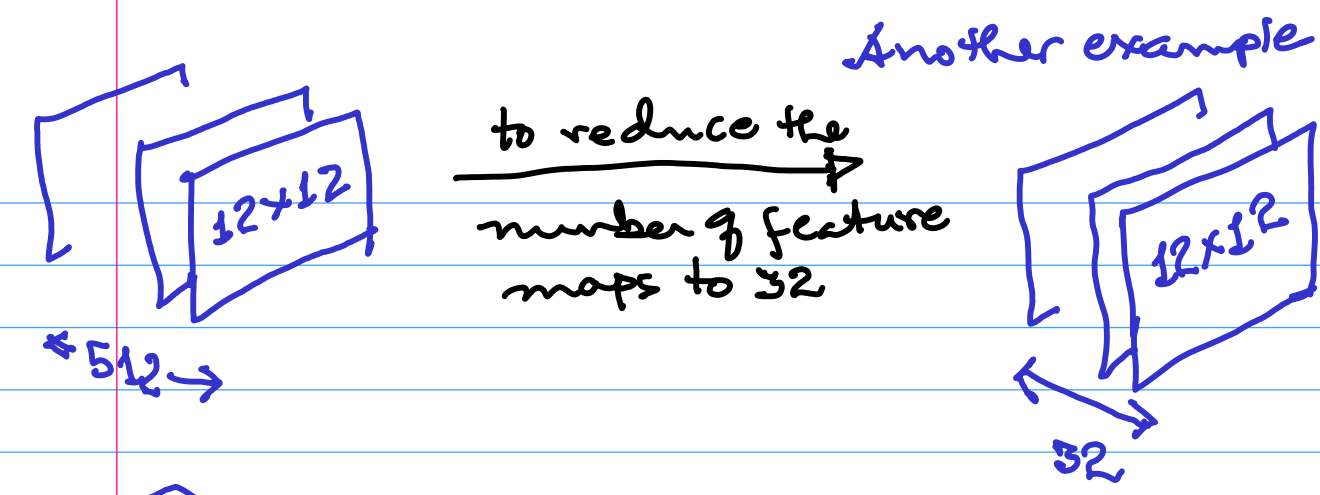$64 \times 64$        $1 \times 1 \times 3$ ⟹ $64 \times 64$  or  $64 \times 64$

$1 \times 1 \times 3$

"cross-channel downsampling"
"cross-channel pooling"

12×12

to reduce the number of feature maps to 32

*Another example*

12×12

← 512 →

32

12×12

← 512 →

12×12

← 32 →

creates one

512

[ one ]

12×12

one slice

**Physical significance:**
Basically have 512 weights for each pixel position in the 12×12 image

now, we take 31 other such "1×1" entities, each with 512 layers
→ each of these has different weight combinations for the 512 constituent channels

We want 32 such slices
(→ 32 such filters)
The resultant 12×12 image has the weighted sum at each pixel position, of the constituent 512 pixel values at each pixel position

stack the 12×12
slices: all 32 of them,
together

12×12

32

or

32

512

512

32
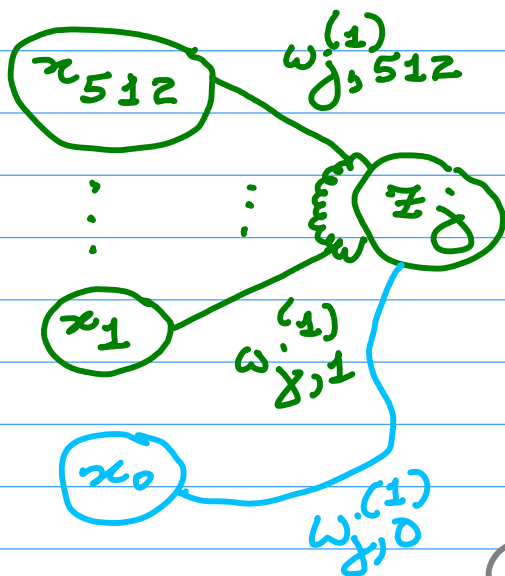
parameters?

$$32 \times (512 \times 1 \times 1 + \underbrace{1}_{\text{bias}})$$

other option: convolutionally
reduce it → more
parameters.

How to introduce non-linearities?

$x_{512}$   $\omega_{j,512}^{(1)}$

$z_j$

$x_1$   $\omega_{j,1}^{(1)}$

$x_0$   $\omega_{j,0}^{(1)}$

By default, activation

$$a_j = \sum_{i=1}^{512} x_i \cdot \omega_{j,i}^{(1)} + \omega_{j,0}^{(1)}$$

Now, we can apply

$$h(a_j)$$

→ activation
function

(possibly non-linear)

e.g., ReLU or tanh or sigmoid