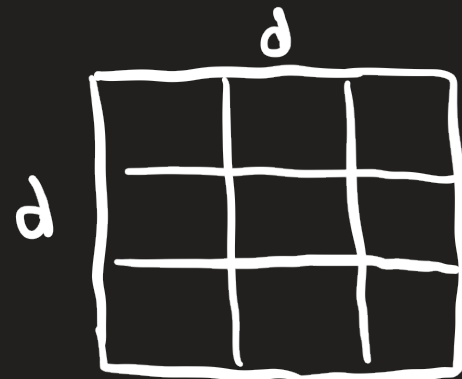


ML and Opt lecture 17

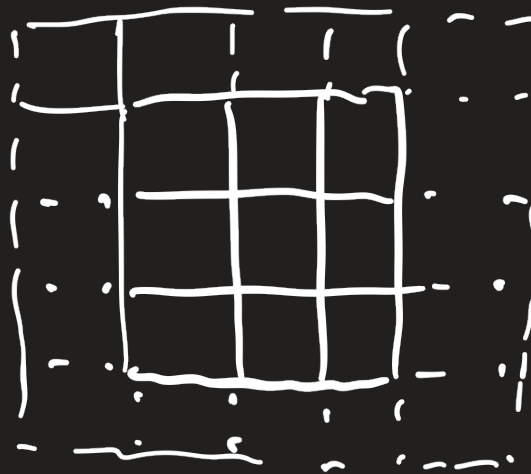
- Transposed Convolutions
- Backprop for Convolutional Layers (sketch)
- Vanishing and Exploding Gradients
- Google Inception architecture (v1); auxiliary losses

Transposed Convolutions

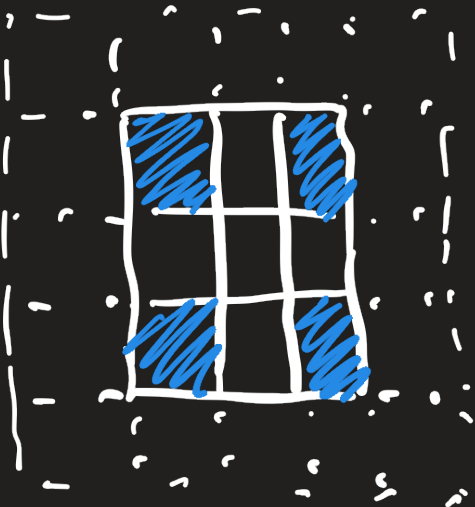
Convolution



pad P



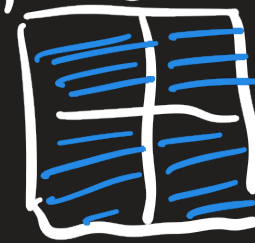
convolve
w/ stride S
filter of size K



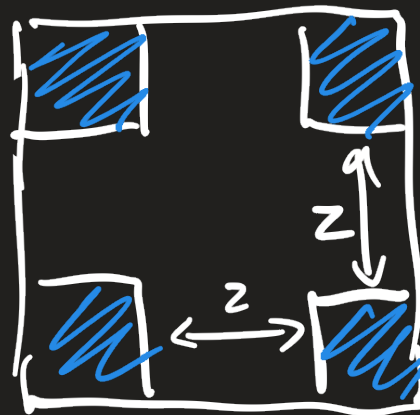
Transposed Convolution

insert Z

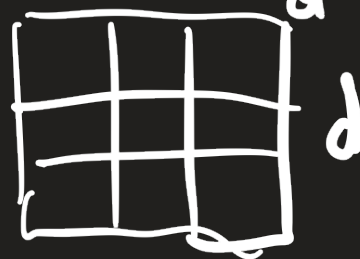
zeros b/w
each row & column



pad P'



convolve w/ stride 1
filter of size K



To "invert" a convolutional filter K , we learn via BP, a filter \bar{K} . We need choose Z, P' so that the images from K when passed through \bar{K} have dimension d .

Let o be the size of the image from the forward convolution

$$o = \frac{d + 2p - k + 1}{s}$$

pytorch:
Conv2D

Then d' be the size of the image from the transposed convolution when the input is size o

$$d' = o + (o-1) \cdot z + 2p' - k + 1$$

Want to choose z and p' so $d' = d$

Claim is: $z = s - 1$

$$p' = k - p - 1$$

pytorch:
Conv2DTranspose

$$\begin{aligned} d' &= \left(\frac{d + 2p - k + 1}{s} + 1 \right) + \left(\frac{d + 2p - k}{s} \right) \cdot (s - 1) + 2(k - p - 1) - k + 1 \\ &= 1 + d + 2p - k + 2k - 2p - 2 - k + 1 \\ &= d \quad \checkmark \end{aligned}$$

Conv2d

Input: $d \times d$

Kernel: $K \times K$

Padding: P

Stride: S

Output: $o \times o$

$$\text{where } o = \frac{d + 2p - K}{S} + 1$$

Pads, then convolves
with stride

Conv2DTranspose

Input: $o \times o$

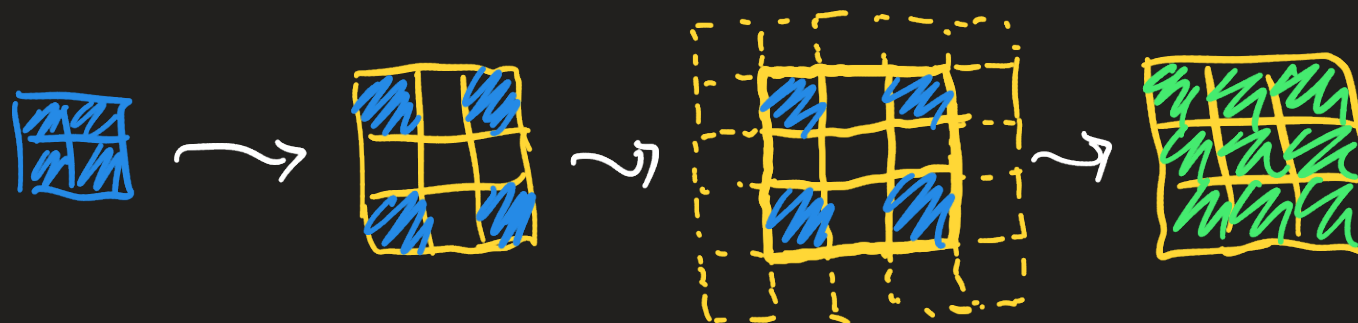
Kernel: $K \times K$

Padding: $K - P - 1$

Spacing: $S - 1$

Output: $d \times d$

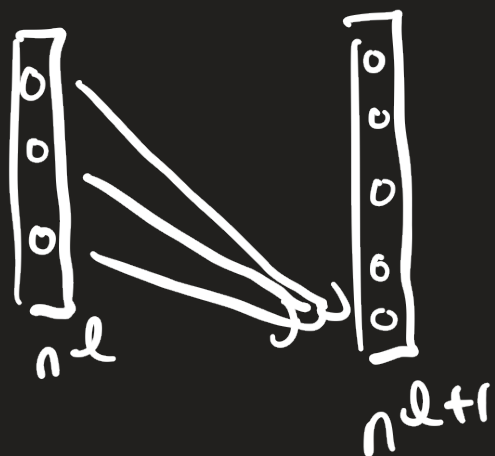
Inserts spaces, pads, then
convolves with stride 1



Efficient Convolutions & Backprop

Parallels b/w MLPs and ConvNets

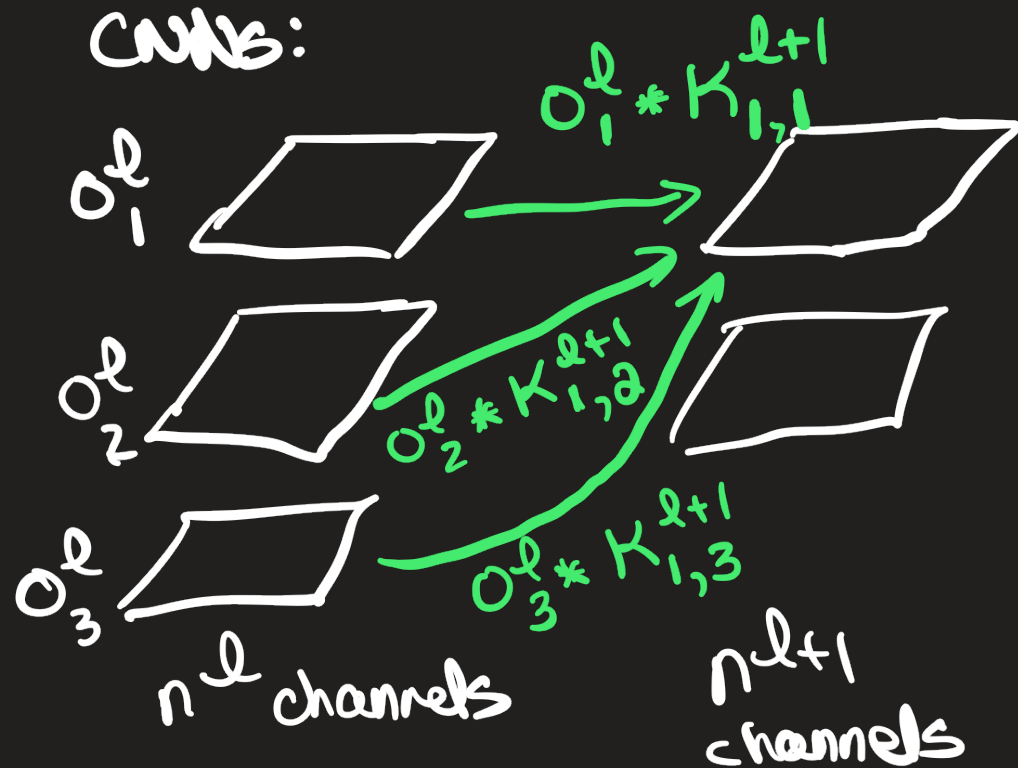
MLP:



$$a^{l+1} = W^{l+1} o^l + b^{l+1}$$

$$o^{l+1} = \sigma(a^{l+1})$$

ConvNets:



$$A_1^{l+1} = o_1^l * K_{1,1}^{l+1} + o_2^l * K_{1,2}^{l+1} + o_3^l * K_{1,3}^{l+1} + b_1^{l+1}$$

a scalar
that is the
same for each
pixel

$$o_1^{l+1} = \sigma(A_1^{l+1})$$

If layer l has n_l channels, then

$$A_i^{l+1} = \left(\sum_{j=1}^{n_l} O_j^l * K_{i,j}^{l+1} \right) + b_i^{l+1}$$

$$O_i^{l+1} = \sigma(A_i^{l+1})$$

The number of parameters for layer $l+1$ is

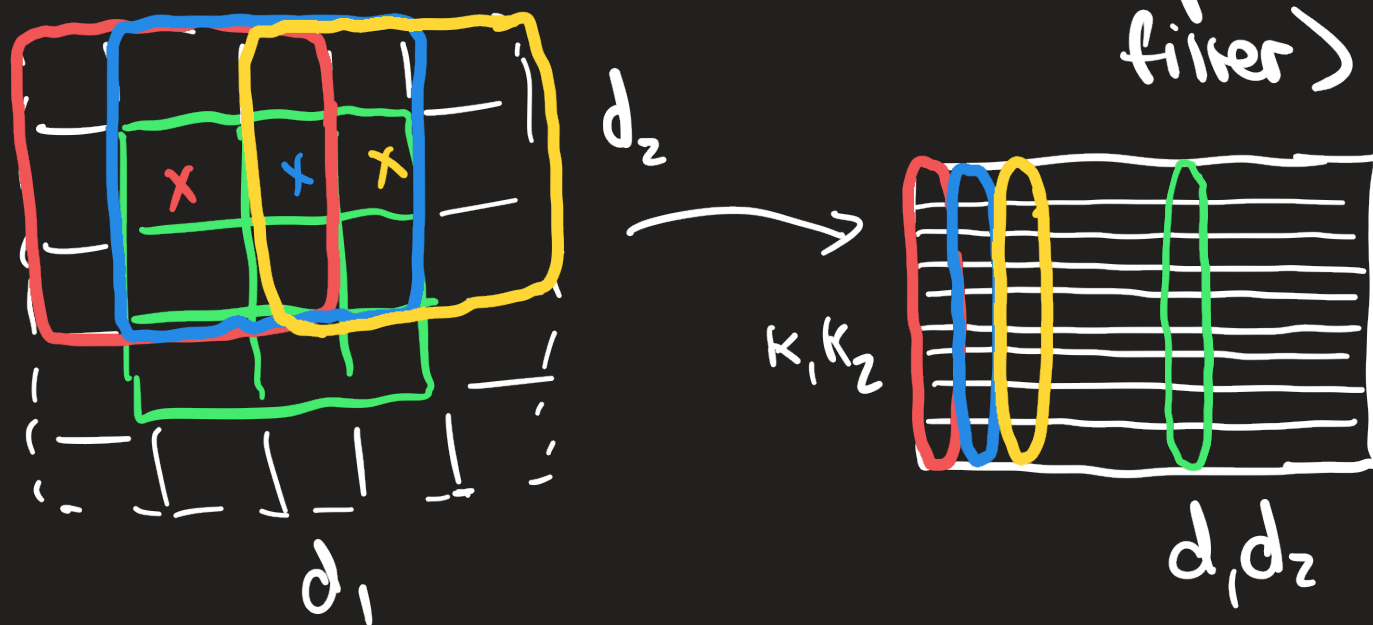
$$n_l n_{l+1} k_1 k_2 + n_{l+1} \quad (\text{assuming all the filters connecting layer } l \text{ to layer } l+1 \text{ are } k_1 \times k_2)$$

(Sketch) Efficient Computations

cf Manas Sahni "Anatomy of a High-Speed Convolution"

idea: reduce convolution to the **indcol** operation and GEMMs

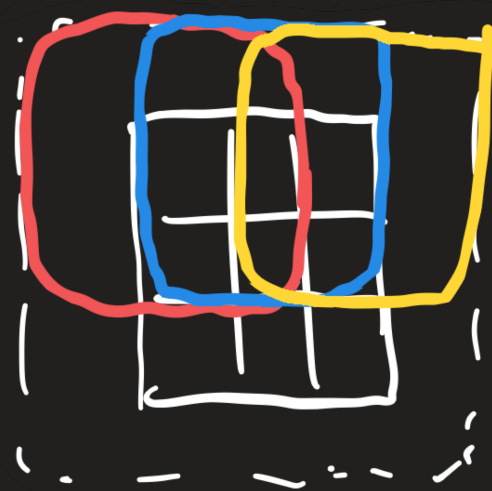
- **indcol** takes an image in $\mathbb{R}^{d_1 \times d_2}$ and maps to a matrix $\mathbb{R}^{k_1 \times k_2 \times d_1 \times d_2}$ (corresponding to zero-padding the input and convolving by a $k_1 \times k_2$ filter)



note that we want to compute
and this implies

$$O_j^l * K_{i,j}^{l+1}$$

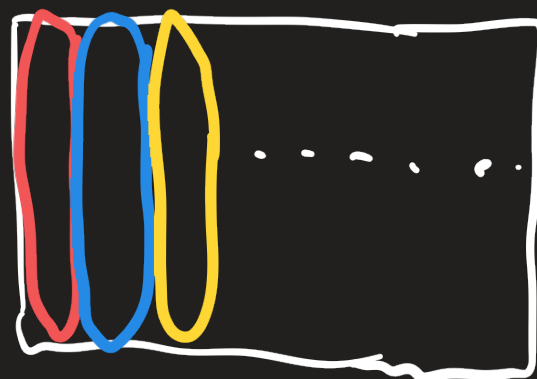
$$\text{vec}(O_j^l * K_{i,j}^{l+1}) = \text{vec}(K_{i,j}^{l+1}) \cdot \text{indcol}(O_j^l)$$



*

1	2	3
4	5	6
7	8	9

\Rightarrow



This means we can write our CNN layers in the form

$\tilde{O}_i^l = \text{vec}(O_i^l)$ — the vectorized channel i on convolutional layer l

$\tilde{a}_i^l = \text{vec}(A_i^l)$ — the vectorized preactivation of channel i on layer l

$\tilde{K}_{i,j}^l = \text{vec}(K_{i,j}^l)$ — the vectorized kernel connecting channel j from layer $l-1$ to channel i from layer l

$M_i^l = \text{im2col}(O_i^l)$ — the matrix coming from applying `im2col` to the i th channel on layer l

$$\tilde{a}_i^{l+1} = \sum_{j=1}^{n_{l-1}} \tilde{K}_{i,j}^{l+1} M_j^l + b_i^{l+1}$$

$$\tilde{O}_i^{l+1} = \sigma(\tilde{a}_i^{l+1})$$

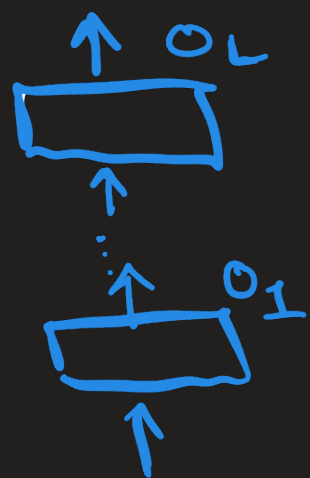
} very similar to MLP forward pass, except uses `im2col`

Issues with deep NNs (not just CNNs):

- overfitting (too much capacity for the amount of training data)
- vanishing & exploding gradients \Rightarrow slow learning!
- hyperparameter selection, e.g.
 - kernel sizes?
 - # channels per layer?
 - # layers?
 - type of pooling & locations?
 - stride & dilation & padding?
 - learning rates? algorithm? minibatch size?
 - weight decay?
 -

Vanishing & Exploding Gradients

Phenomenon that, as the number of layers increases,



as $l \rightarrow 1$ in the backpropagation algorithm,

$$\|\nabla_{w^l} f\|_2 \rightarrow \begin{cases} 0 & \text{vanishing} \\ \infty & \text{exploding} \end{cases}$$

Occurs because of the chain rule.

Recall $o^{l+1} = \sigma(w^{l+1} o^l + b^{l+1})$ so by the multivariate chain rule,

$$\nabla_{o^l} f = J_{o^{l+1}}(o^l)^T \nabla_{o^{l+1}} f$$

and

$$\left[J_{o^{l+1}}(o^l) \right]_{ij} = \left[\frac{\partial o_i^{l+1}}{\partial o_j^l} \right]_{ij} = \sigma'((w^{l+1} o^l + b^{l+1})_i) \cdot (w^{l+1})_{ij}$$

so

$$\begin{aligned} J_{o^{l+1}}(o^l) &= \text{Diag}(\sigma'(\omega^{l+1}o^l + b^{l+1}))\omega^{l+1} \\ &= \text{Diag}(\sigma'(a^{l+1}))\omega^{l+1} \end{aligned}$$

and consequently,

$$\nabla_{o^l} f = (\omega^{l+1})^T \text{Diag}(\sigma'(a^{l+1})) \cdot \nabla_{o^{l+1}} f$$

$\approx \dots$

$$= \prod_{i=l+1}^L [(\omega^i)^T \text{Diag}(\sigma'(a^i))] \cdot \nabla_{o^L} f$$

Two considerations:

1) How $\text{diag}(\sigma'(a^{l+1}))$ behaves



so if a^l far from 0, then this looks like a zero matrix, so

$$\|\nabla_{\theta^l} f\|_2 \ll \|\nabla_{\theta^{l+1}} f\|_2$$

2) If the norm of our weight matrix W^{l+1} is large, then $\|\nabla_{\theta^l} f\|_2 \gg \|\nabla_{\theta^{l+1}} f\|_2$

if it is small, then $\|\nabla_{\theta^l} f\|_2 \ll \|\nabla_{\theta^{l+1}} f\|_2^2$

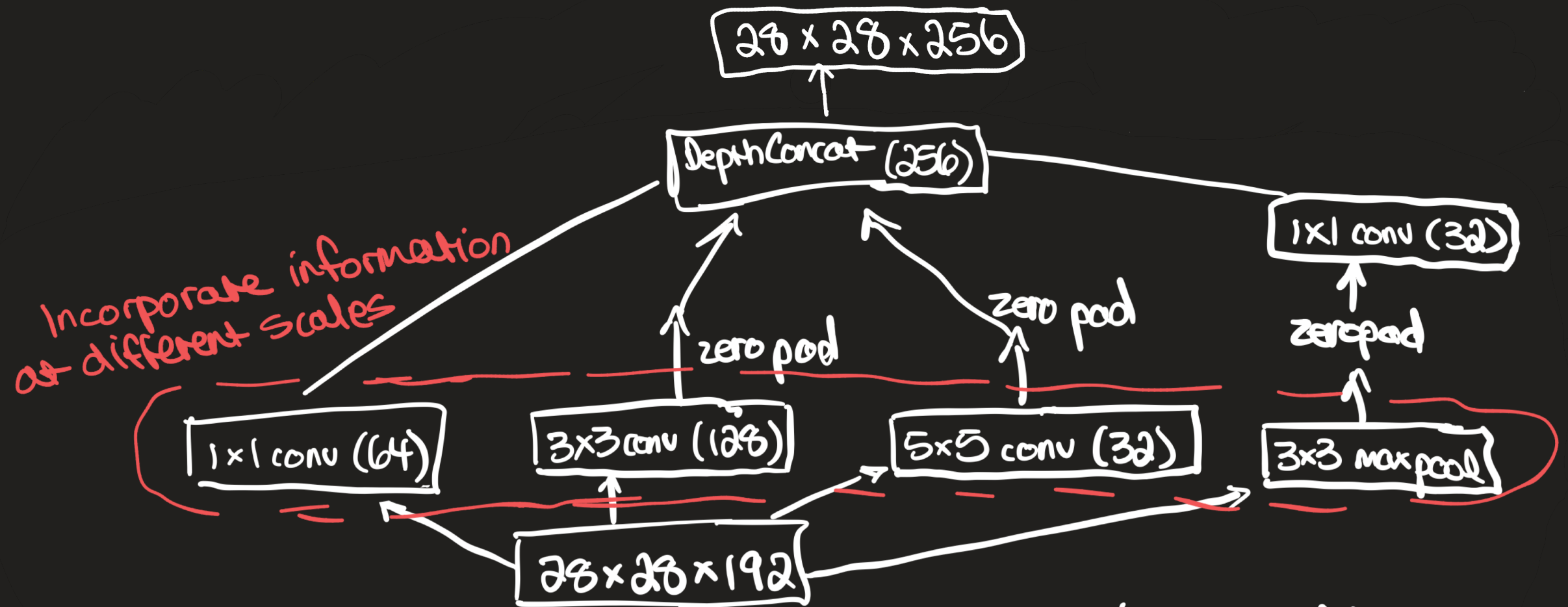
Consequence:

naively training deep NN architectures
either fails or gives poor performance

Takeaways:

- we want to keep activations close to zero so the nonlinearities in the network do not saturate
 - suggests "normalization" layers to keep activations well-behaved.
- we want to keep weight matrices well-behaved:
 - suggests regularizing by norm of weight matrices
- we want to maintain short paths to the output to prevent attenuation/explosion to compound
 - suggests using "auxiliary" losses.

GoogleNet (22 layer) CNN Inception v1 2014



Downside: lots of parameters

E.g. for the 3x3 conv features, we have

$$192 \times 128 \times 3 \times 3 + 128 = 221,312 \text{ parameters}$$

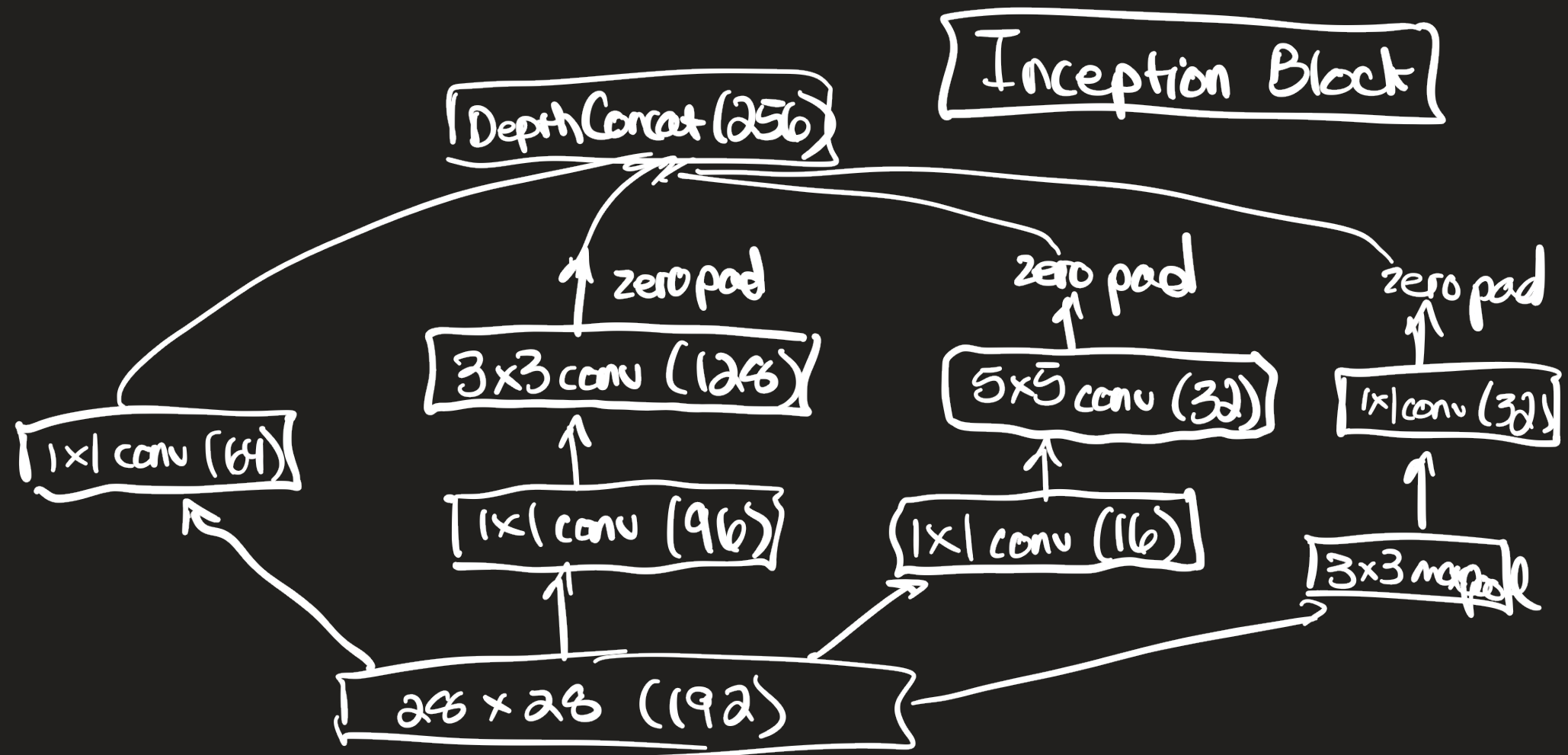
Naive Inception Block:
preserves input dimensions,
uses information at
multiple scales

$$96 \times 192 + 96$$

Issue: too many parameters

$$+ 128 \times 96 \times 3 \times 3 + 128$$

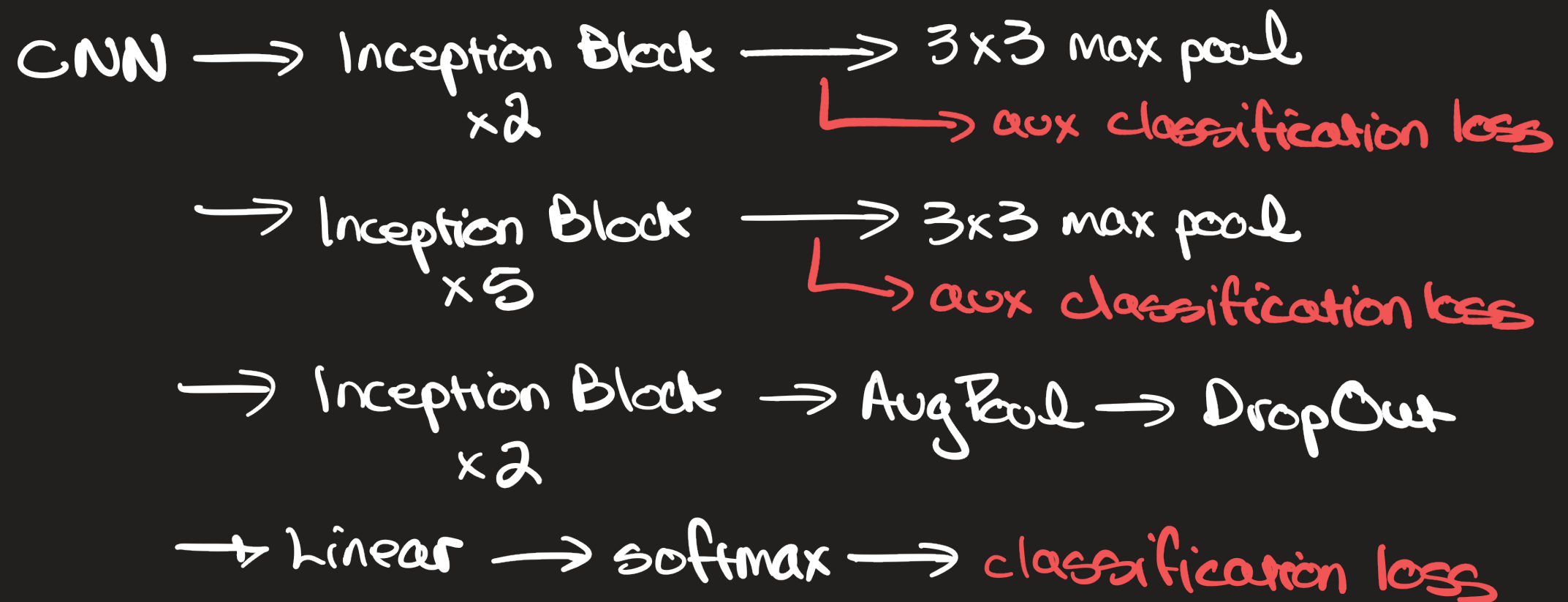
Soln: use 1×1 conv layers for dimensionality reduction



Check: # of 3×3 conv parameters (including the 96 1×1 conv)

$$129,248 = ? \text{ parameters}$$

Original Inception (v1) architecture



Train this architecture to minimize the weighted sum of the three losses. Injects gradient information at intermediate layers to mitigate vanishing/exploding gradients.

Input: $224 \times 224 \times 3$

type	patch size/ stride	output size	depth	$\#1 \times 1$	$\#3 \times 3$ reduce	$\#3 \times 3$	$\#5 \times 5$ reduce	$\#5 \times 5$	pool proj	params	ops
convolution	$7 \times 7 / 2$	$112 \times 112 \times 64$	1							2.7K	34M
max pool	$3 \times 3 / 2$	$56 \times 56 \times 64$	0								
convolution	$3 \times 3 / 1$	$56 \times 56 \times 192$	2		64	192				112K	360M
max pool	$3 \times 3 / 2$	$28 \times 28 \times 192$	0								
inception (3a)		$28 \times 28 \times 256$	2	64	96	128	16	32	32	159K	128M
inception (3b)		$28 \times 28 \times 480$	2	128	128	192	32	96	64	380K	304M
max pool	$3 \times 3 / 2$	$14 \times 14 \times 480$	0								
inception (4a)		$14 \times 14 \times 512$	2	192	96	208	16	48	64	364K	73M
inception (4b)		$14 \times 14 \times 512$	2	160	112	224	24	64	64	437K	88M
inception (4c)		$14 \times 14 \times 512$	2	128	128	256	24	64	64	463K	100M
inception (4d)		$14 \times 14 \times 528$	2	112	144	288	32	64	64	580K	119M
inception (4e)		$14 \times 14 \times 832$	2	256	160	320	32	128	128	840K	170M
max pool	$3 \times 3 / 2$	$7 \times 7 \times 832$	0								
inception (5a)		$7 \times 7 \times 832$	2	256	160	320	32	128	128	1072K	54M
inception (5b)		$7 \times 7 \times 1024$	2	384	192	384	48	128	128	1388K	71M
avg pool	$7 \times 7 / 1$	$1 \times 1 \times 1024$	0								
dropout (40%)		$1 \times 1 \times 1024$	0								
linear		$1 \times 1 \times 1000$	1							1000K	1M
softmax		$1 \times 1 \times 1000$	0								

Table 1: GoogLeNet incarnation of the Inception architecture