

Today: Finish CNN  
Graph Neural Nets

Reading: Prince through Ch 11  
Start Ch 13

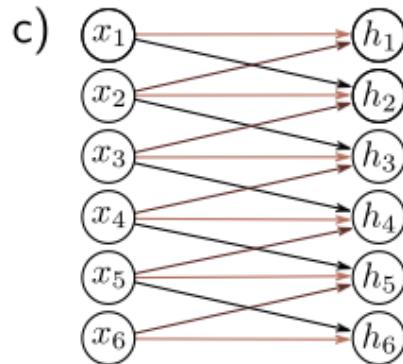
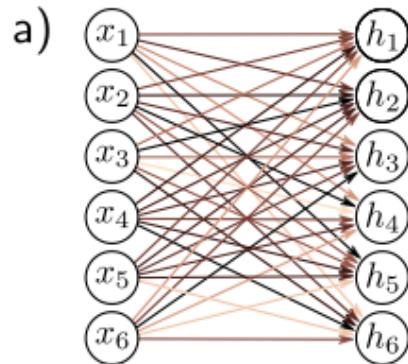
Discussion Tomorrow

Fri TA OH S31Cory 2-3 PM

Architecture Order In Class:

MLPs  $\rightarrow$  [CNNs  $\rightarrow$  Graph NN]  $\rightarrow$  RNN/State-space  $\rightarrow$  Transformers

Goal



b)

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$h_1$	Dark Brown	Medium Brown	Light Brown	Dark Brown	Medium Brown	Light Brown
$h_2$	Medium Brown	Black	Medium Brown	Medium Brown	Medium Brown	Medium Brown
$h_3$	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Black
$h_4$	Light Brown	Black	White	Dark Brown	Medium Brown	Light Brown
$h_5$	Black	Medium Brown	Dark Brown	Medium Brown	Medium Brown	Light Brown
$h_6$	Light Brown	Dark Brown	Medium Brown	Black	Dark Brown	Medium Brown

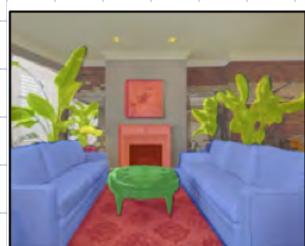
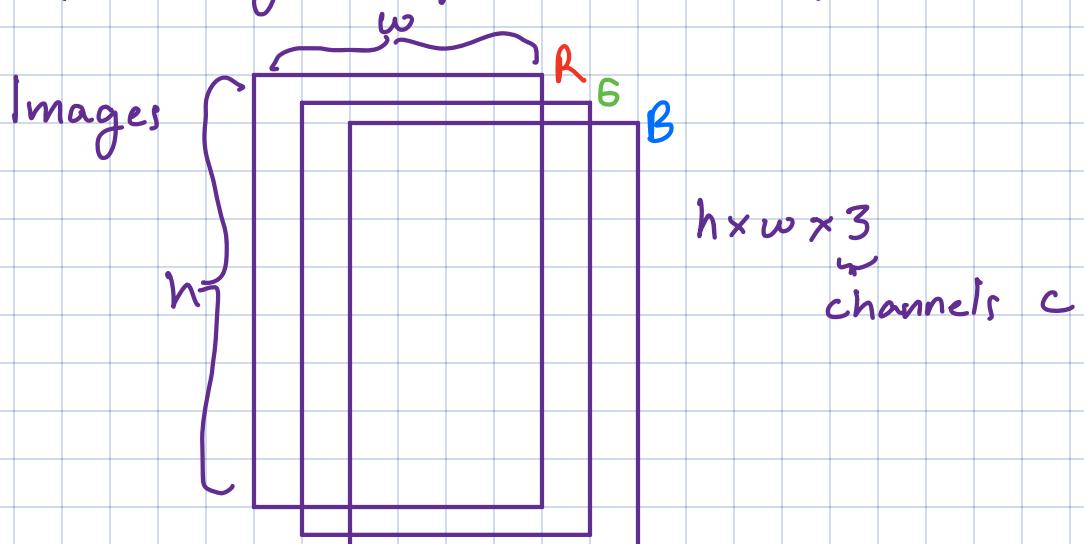
d)

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$h_1$	Dark Brown	Medium Brown	Light Brown	Dark Brown	Medium Brown	Light Brown
$h_2$	Medium Brown	Black	Medium Brown	Medium Brown	Medium Brown	Medium Brown
$h_3$	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Black
$h_4$	Light Brown	Black	White	Dark Brown	Medium Brown	Light Brown
$h_5$	Black	Medium Brown	Dark Brown	Medium Brown	Medium Brown	Light Brown
$h_6$	Light Brown	Dark Brown	Medium Brown	Black	Dark Brown	Medium Brown

MLP

CNN

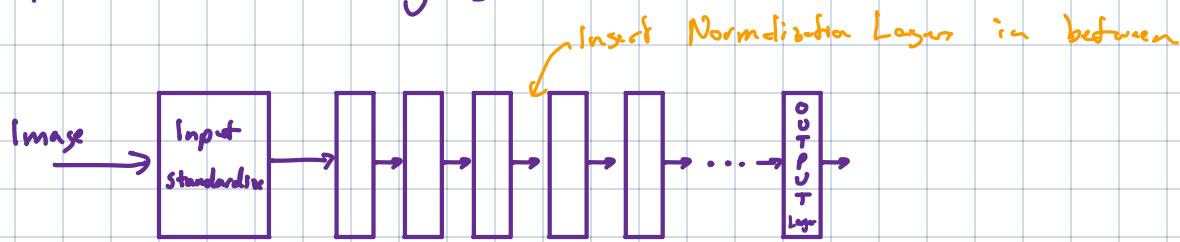
Inspired by computer vision problems: classification,



semantic segmentation,  
etc.

# Training: Stability and Effectiveness

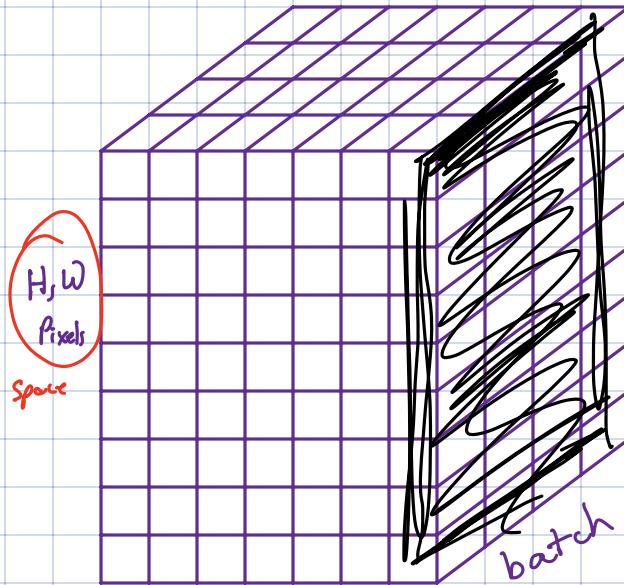
## Normalization Layers



Idea: RMS Norm Layer

$$d\text{-dim} \left\{ \vec{h}_e \right\} \xrightarrow{\text{RMS Norm}} \tilde{h}_e \quad \text{has } \|\tilde{h}_e\|_{\text{RMS}} = 1$$

$$\begin{aligned} \tilde{h}_e &= \frac{\vec{h}_e}{s} \quad \text{where } s = \sqrt{d} \sqrt{\sum_i h_e[i]^2} \\ &= \frac{\vec{h}_e}{\|\vec{h}_e\|_{\text{RMS}}} + \epsilon \\ &= \frac{\vec{h}_e}{\|\vec{h}_e\|_{\text{RMS}}} + \epsilon \end{aligned}$$



Batch Norm

Average over space & batch

Let  $B$  be what is averaged over

$$m = \frac{1}{|B|} \sum_{i \in B} h_{in}$$

$$s = \sqrt{\frac{1}{|B|} \sum_{i \in B} (h_{in} - m)^2}$$

$$h_{out}[i,j] = \gamma \left( \frac{h_{in}[i,j] - m}{s + \epsilon} \right) + \delta$$

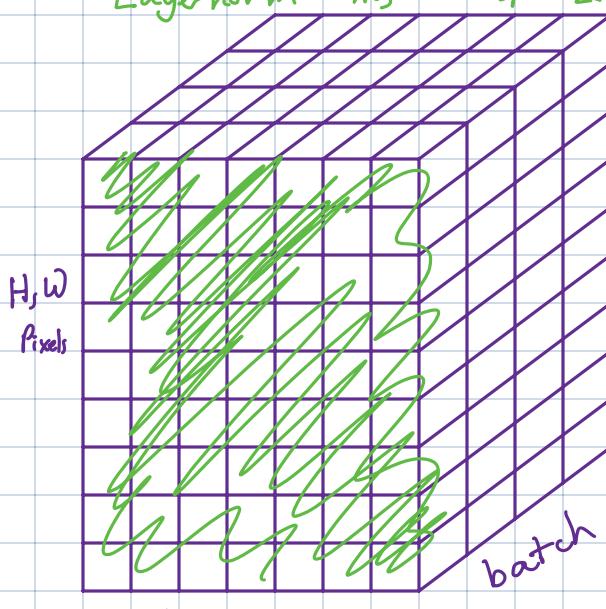
↑ learnable

Learnable center  $\delta$   
scale  $\gamma$

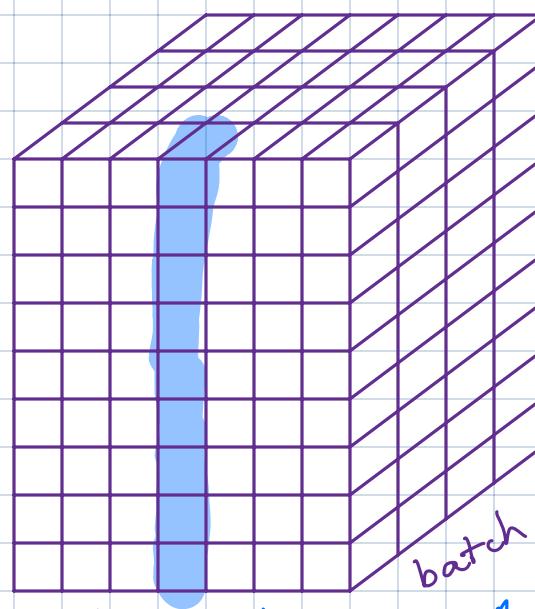
Initialize  $\Rightarrow \delta = 0$   
 $\gamma = 1$

Typically: Don't use weighting on  $\delta, \gamma$ .

channels  
LayerNorm: Avg over space & channels



$H, W$   
Pixels



InstanceNorm: Avg across space only

In LayerNorm  
Can have different  
 $\gamma$  and  $\delta$   
per channel.

# Residual / Skip - connections

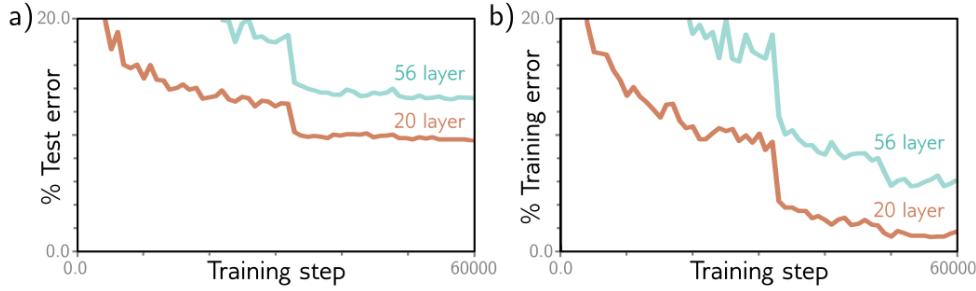
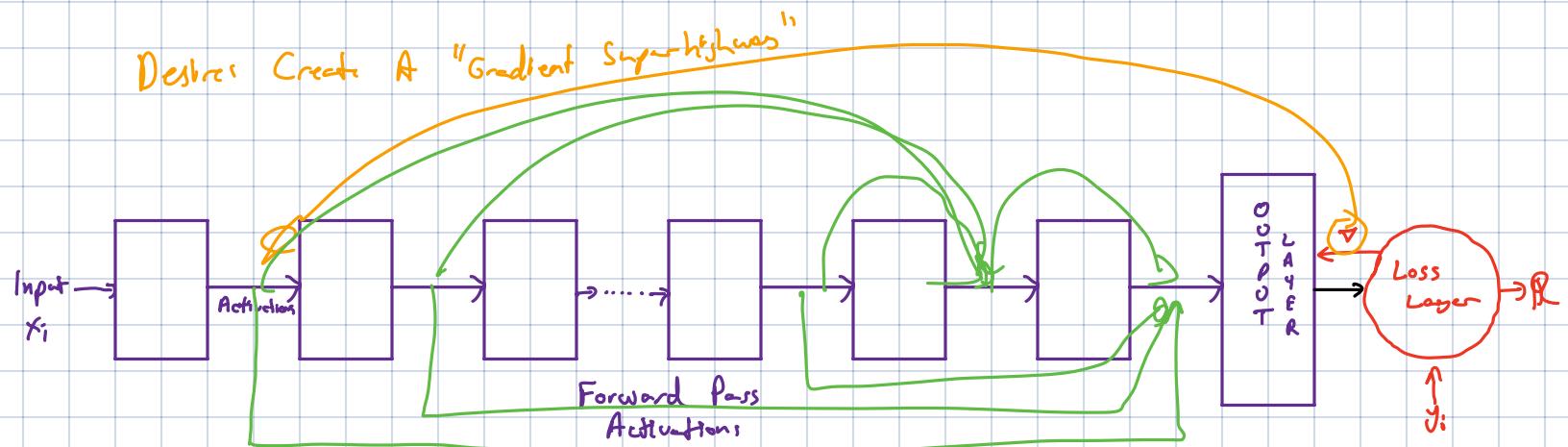
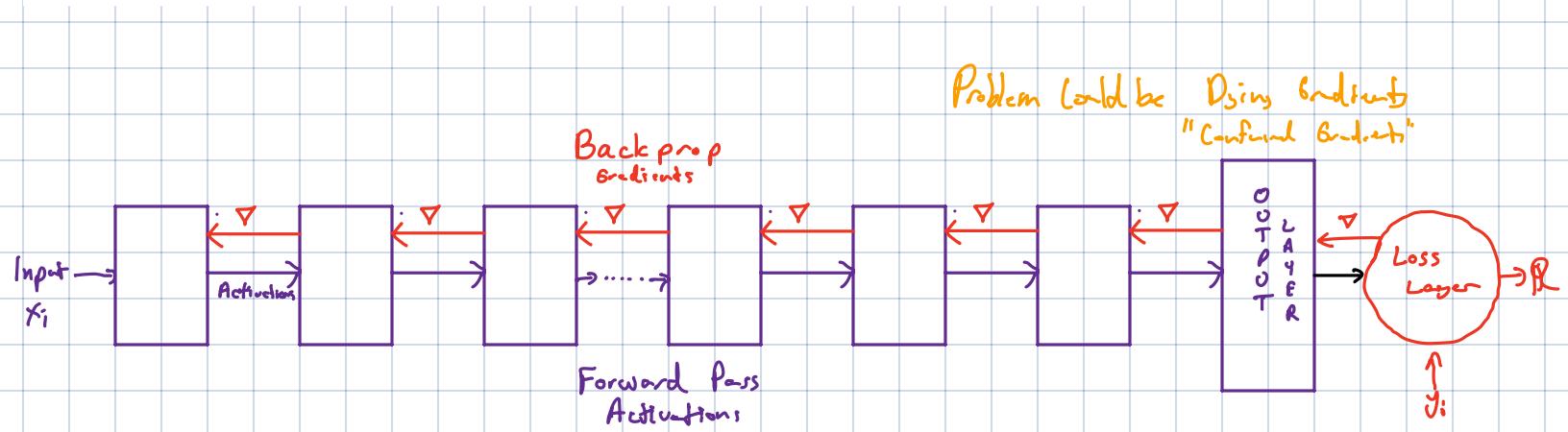


Fig 11.2 in Prince

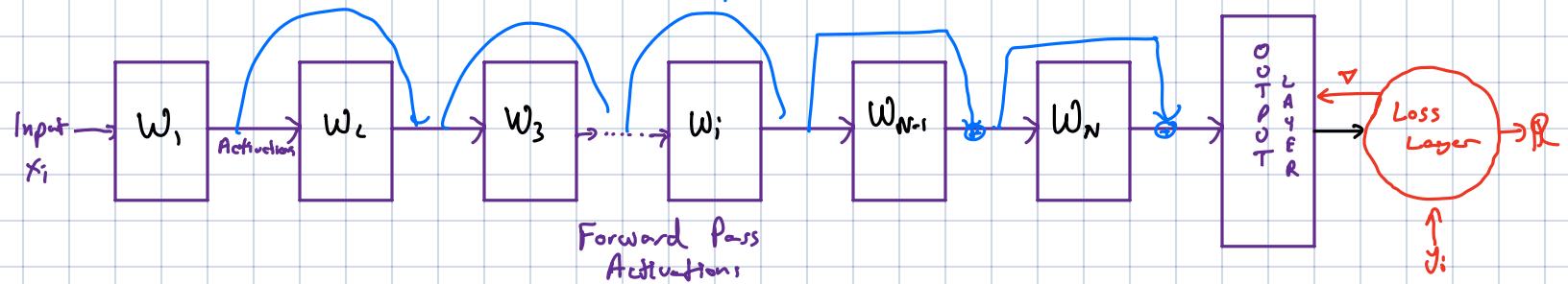
← Feels like Underfitting.



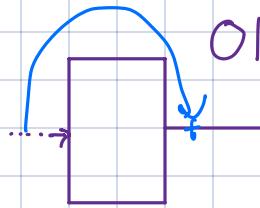
Idea 1: Add links from outputs of all earlier layers to the input of a given layer.

Blows up activations 😢

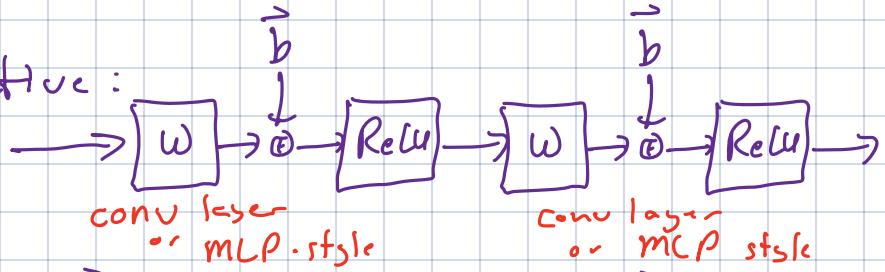
skip Connection: Add input to the output of a layer



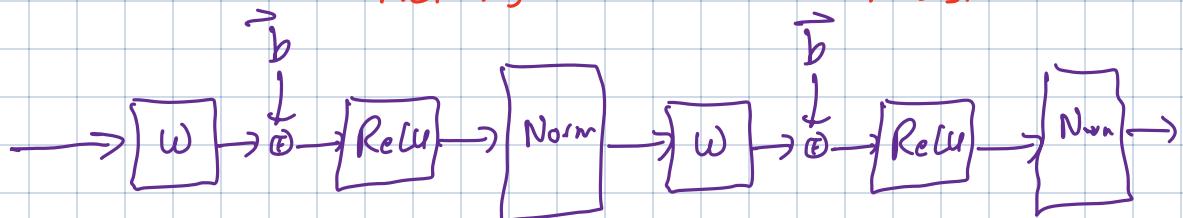
What goes inside a block that's skipped over?



Old Perspective:

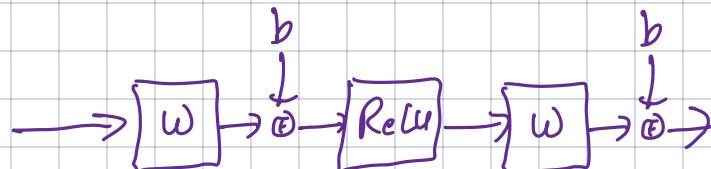


or



What matters most?

CORE



INTUITION: Pre-ResNet Perspective:  $f_i(F_i(F_2(F_3(\dots f_L(x))\dots))$

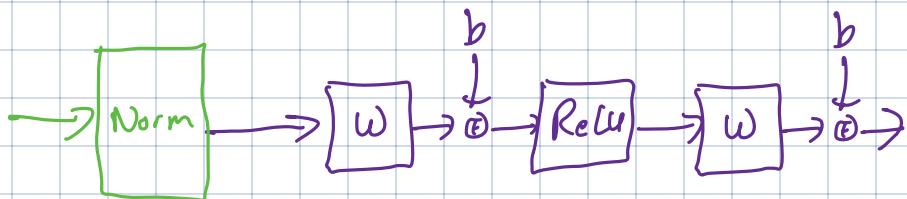
ResNet Perspective  $f_{i+1}(x) = f_i(x) + g_i(f_i(x))$

Feels like  $\frac{d}{dt} x(t) = g(x(t))$

Where should normalizations go?

Modern Default: Inside the block near the start

e.g.



Called "Pre-Norm" convention

Associated with more stable training

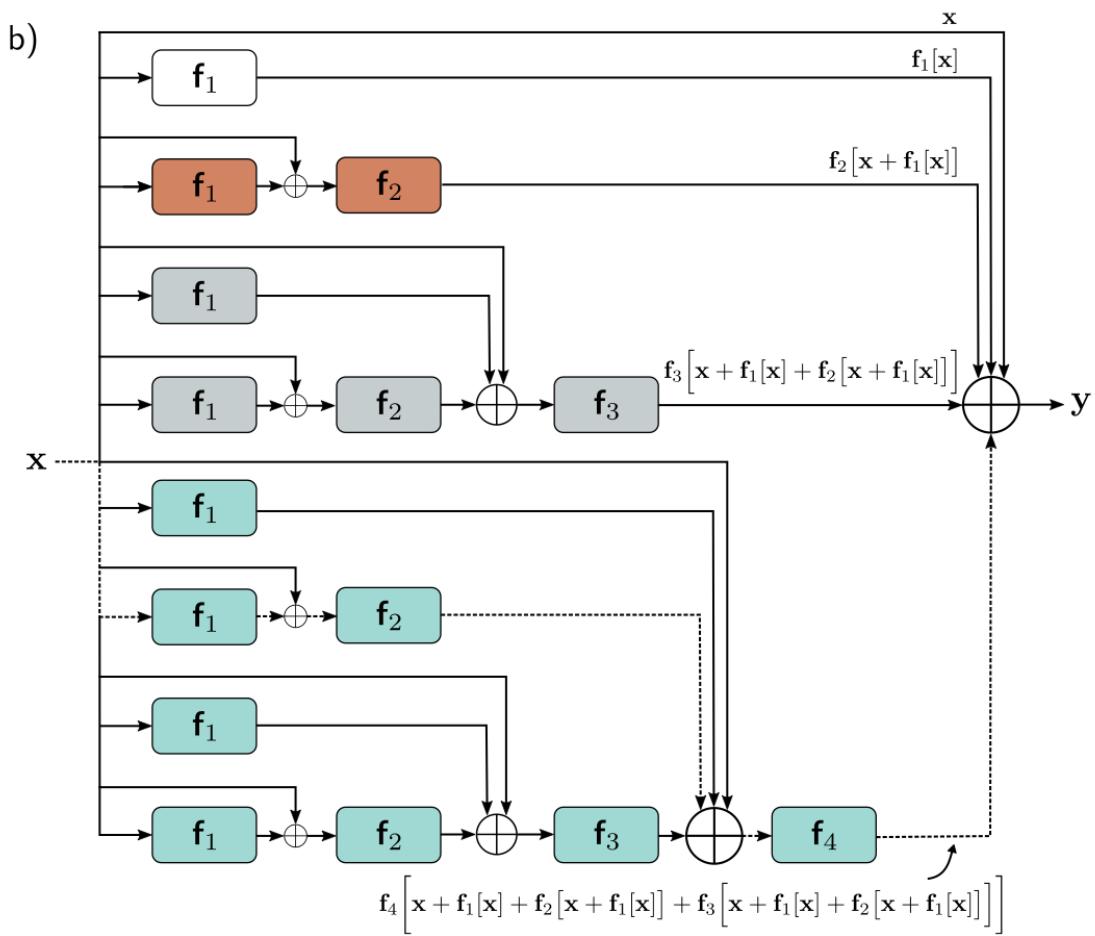
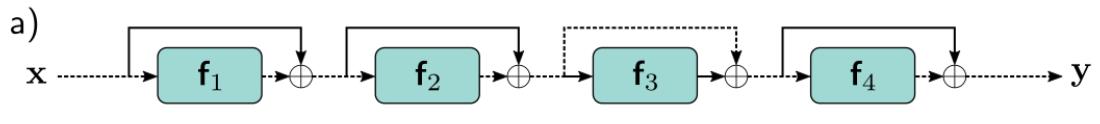


Fig 11.4  
in Prince

# One challenge...

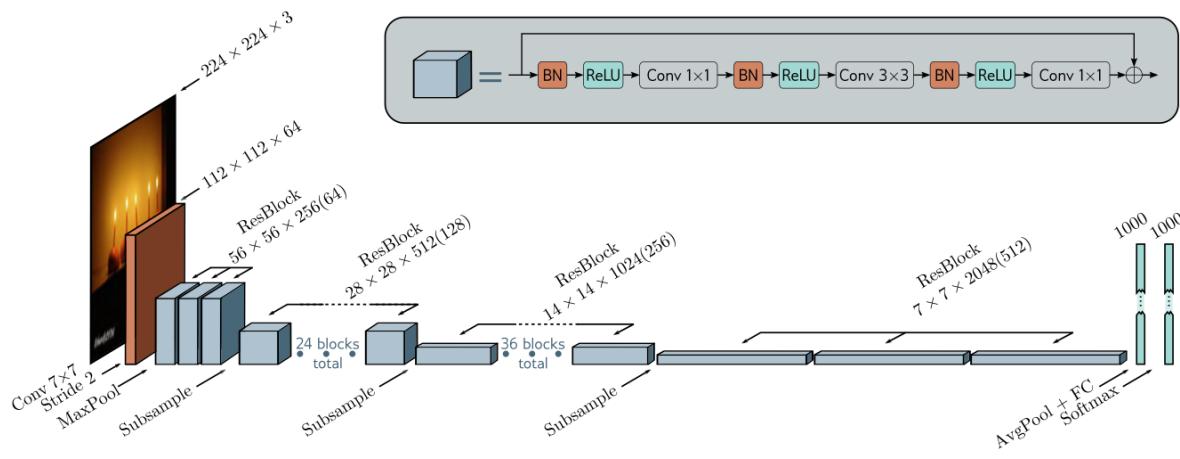


Fig 11.8  
in Prince

ResNet-200

What do we do for subsampling? (e.g. a block that has a maxpool, stride, etc)  
along with other stuff



How to skip??

Need something that changes shape.

e.g. Linear Layer of some kind. (Matrix with right shape)  
 $P$

Note: IF the entire block is just a maxpool downsampling or a subsampling, then that particular block doesn't need a skip connection because gradients pass through cleanly.



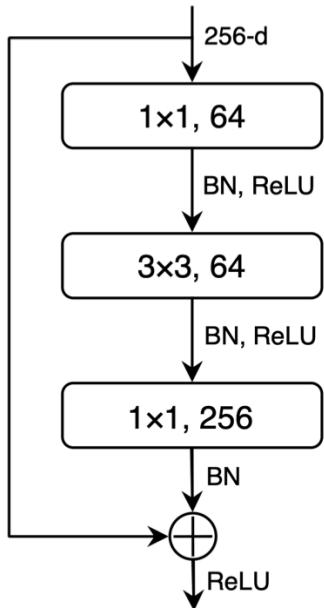
Added in Office Hours.

A more modern perspective:

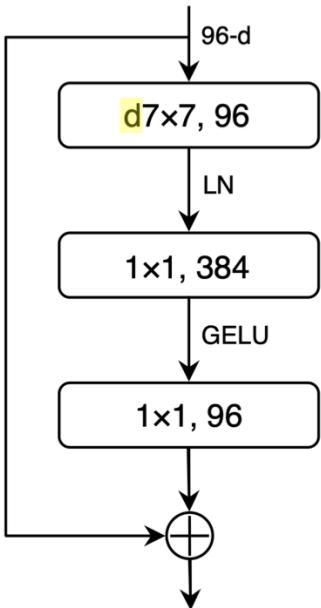
"A ConvNet for the 2020s"  
by Liu, et.al. in 2022

(Meta FAIR & UC Berkeley)

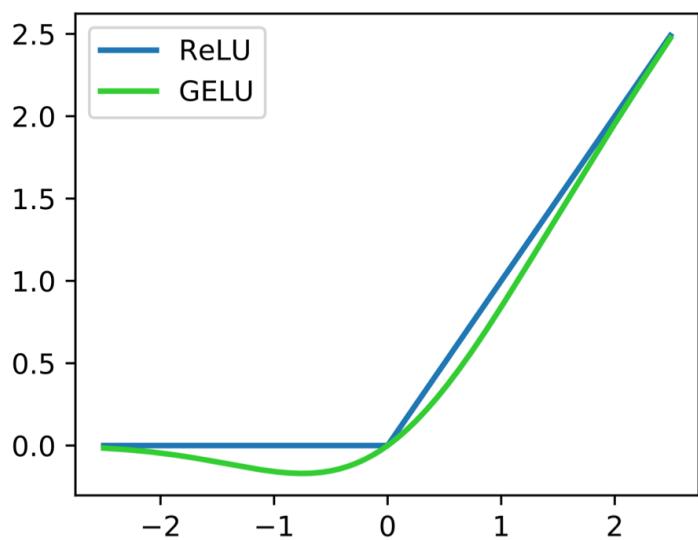
### ResNet Block



### ConvNeXt Block



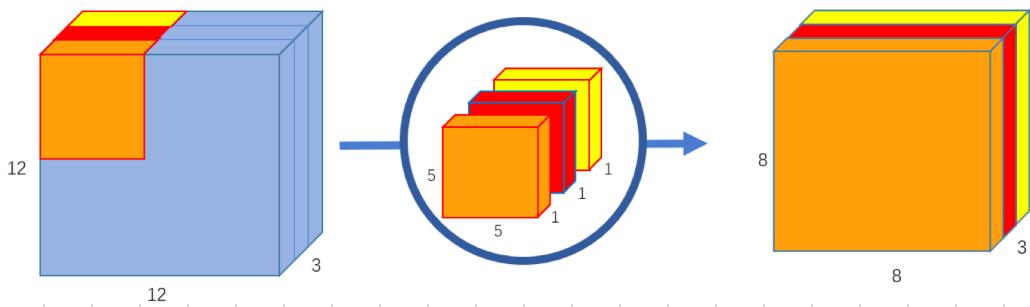
Nonlinearities



Pic from Wikipedia

$\times \Phi(x)$

## Depthwise - convolution



Just treat each channel separately as though it was a one-input-channel & one-output-channel conv.

No weight-sharing across channels.

## Recall Pooling (Downsampling)

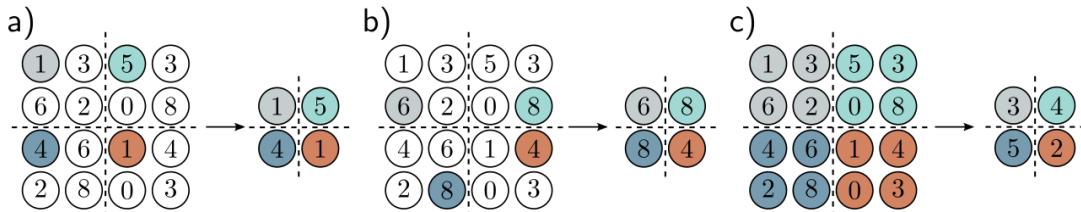


Fig 10.11 in Princ

stride/subsample

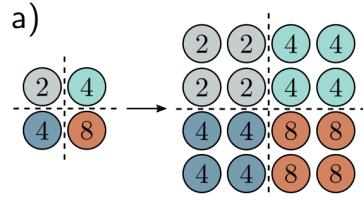
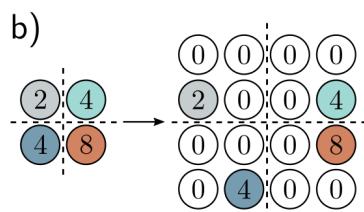
Max-pool

Mean-Pool

$$\text{Output} = \max_{i \in \text{Pooling domain}} h_i$$

$$\text{Output} = \frac{1}{|\text{Pooling domain}|} \sum_{i \in \text{Pooling domain}} h_i$$

What if we wanted the reverse? Go from low-res to high res? Upsampling



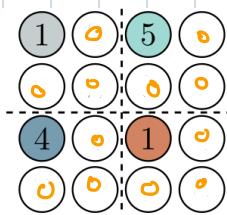
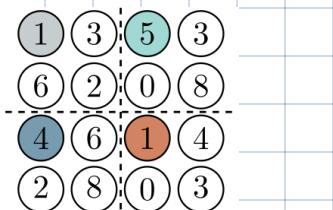
Max-unpool

Must remember which position was the max.

Duplication — like opposite of mean



What is counterpart upsampling for subsampling?



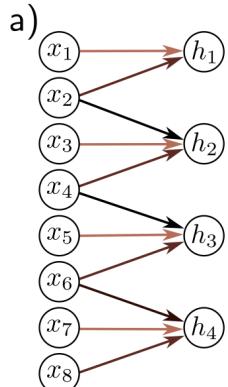
- Ideas:
- Interpolate in a smooth way
  - Randomly sample w.r.t. a distribution
  - Make it learnable  
(Version of a)
  - Fill with zeros

Idea:

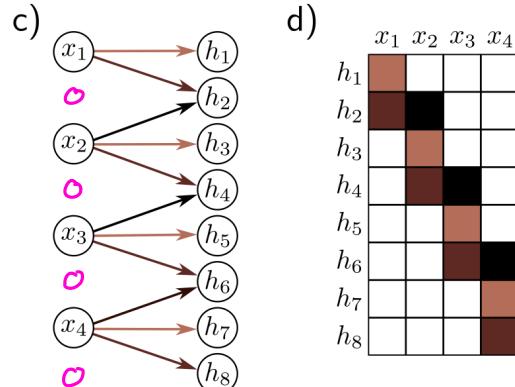
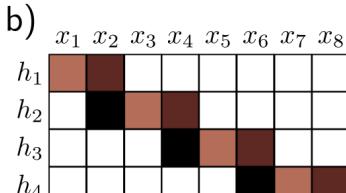
Use (d) followed by a learnable conv.

Behaves like learnable interpolation

Transpose Convolution: Downsample has: Conv + Subsample  $\Rightarrow$  Strided Conv  
 Upsample has: Zero-Fill Upsample + Conv  $\Rightarrow$  Transpose Conv



Strided Conv



Transpose Conv

3x3 examples:

$w_{nw}$	$w_N$	$w_{ne}$
$w_w$	$w_m$	$w_e$
$w_{sw}$	$w_s$	$w_{se}$

circ  $\{ w_n \}$

$\vec{b}$  } Cont

$$\vec{h}_{out} = \vec{b} + \sum_i w_i \vec{h}_{in,i}$$

Standard Conv

R ranges over the  $k \times k$  positions  
in the conv filter.

Transpose Conv : ranges only over those positions  
that are neighbors given stride

# The code perspective:

## Standard conv

PYTORCH MXNET JAX TENSORFLOW

```
def corr2d(X, K): #@save
    """Compute 2D cross-correlation."""
    h, w = K.shape
    Y = torch.zeros((X.shape[0] - h + 1, X.shape[1] - w + 1))
    for i in range(Y.shape[0]):
        for j in range(Y.shape[1]):
            Y[i, j] = (X[i:i + h, j:j + w] * K).sum()
    return Y
```

From D2L.ai  
Ch 7.2

## Transpose conv

PYTORCH MXNET

```
def trans_conv(X, K):
    h, w = K.shape
    Y = torch.zeros((X.shape[0] + h - 1, X.shape[1] + w - 1))
    for i in range(X.shape[0]):
        for j in range(X.shape[1]):
            Y[i: i + h, j: j + w] += X[i, j] * K
    return Y
```

Ch 14.10

How to do image-level tasks?



## Semantic segmentation

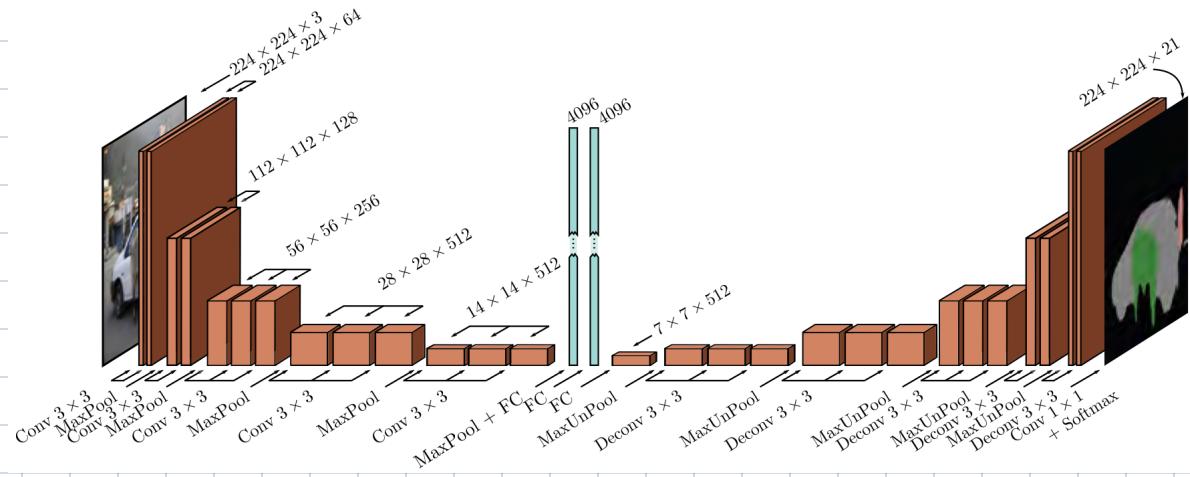


Fig 10.19  
in Prince

Transpose-Conv / Upsampling kind of works, but there's a tension between getting details and big picture. Later layers doing upscaling have no local context.

## U-nets

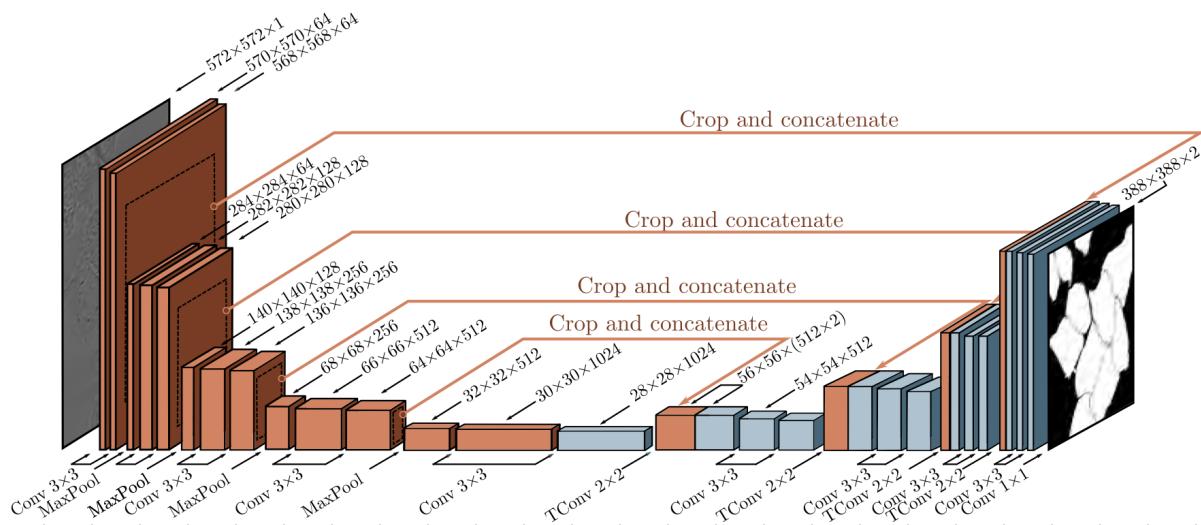


Fig 11.10  
in Prince

Bring information for this pixel from earlier processing.

Like skip connections vis-a-vis gradient flow, but do not replace the block-by-block residual connections.