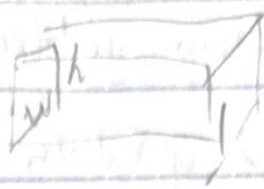
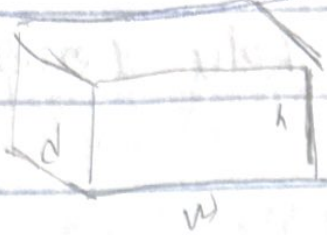


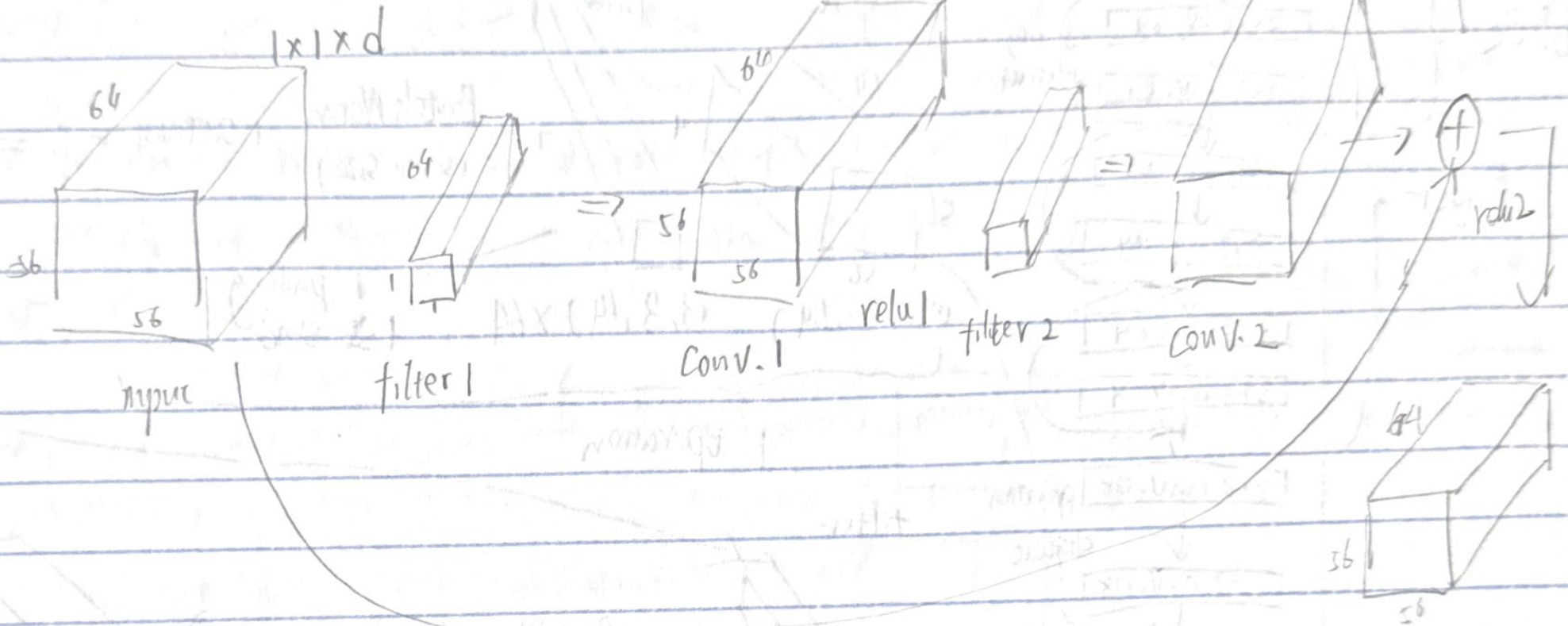
# Resnet

more parameters ( $w$ ) should more likely to fit, better performance  
 address vanishing gradient and degradation



vanishing gradient =  
 for deep network, the gradient from where the loss function calculated easily shrink to 0 after several application of chainrule

through skip connection, gradient can flow directly backward from later layers to initial filters

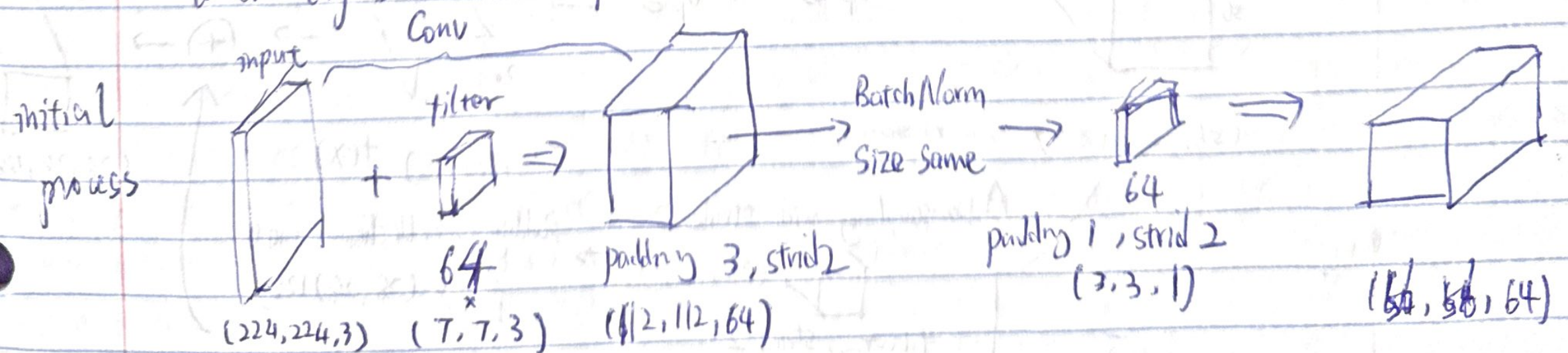


★ we just skip 1, since second didn't add to  $x$  for relu  
 ★ what we try to learn:  $f(x)$

Resnet makes Learning Identity easier

learning  $f(x) + x \Rightarrow 0 + x$  is easier than learning  $H(x) \Rightarrow x$

Gradients can flow directly through the skip connections backward from later layers to initial filters





Only 1 max pooling and 1 average pooling layer at the end of ResNet  
 or the beginning first Conv layer  
 Downsampling of the volume is achieved by increasing stride x pooling

## Resnet Layer

1 operation = a convolution + a Batch Norm + a Relu (Activation)  
 the last operation in block doesn't have Relu (exception)

Basic Block = 2 operation in block

Bottleneck Block = 3 operation in block

