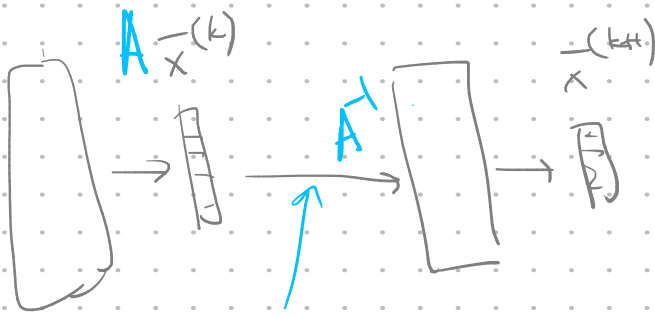


$$p(D|\theta) \rightarrow \max$$

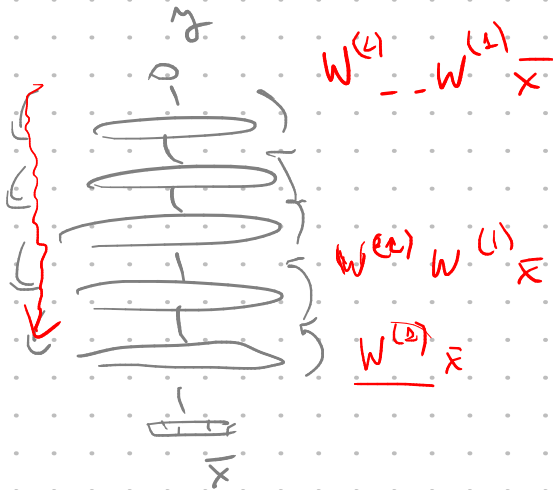
$$p(\theta|D) \rightarrow \max$$

$$\log p(\theta|D) = \underbrace{\log p(D|\theta)}_{(-)} + \underbrace{\log p(\theta)}_{(-)} + \text{const}$$

$L + L_{\text{reg}} \rightarrow \min$



Unsupervised pretraining



exploding gradients
vanishing gradients

$$y = \sum_i y_i \quad y_i = w_i x_i$$

$$\hat{y} = \sum_i \hat{y}_i = \hat{w}^T \hat{x} = h(\hat{w}^T \hat{x})$$

$$\text{Var}[y_i] = \text{Var}[w_i x_i] = \boxed{\mathbb{E}[x_i]^2 \text{Var}[w_i]} + \cancel{\mathbb{E}[w_i]^2 \text{Var}[x_i]} + \text{Var}[x_i] \text{Var}[w_i]$$



1) $E[x_i] = 0$

$Var[y_i] = Var[x_i] Var[w_i]$

$Var[y] = n \cdot Var[w_i] Var[x_i]$

$w_i \sim Unif\left(-\frac{\sqrt{3}}{\sqrt{n}}, \frac{\sqrt{3}}{\sqrt{n}}\right)$

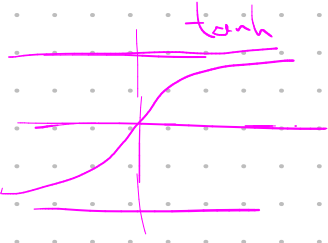
≈ 1

$x \frac{1}{3}$

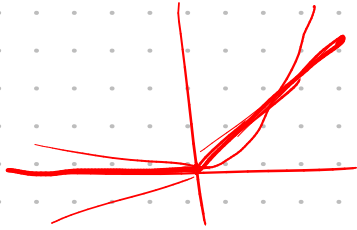
$Var[w_i] = \frac{\left(\frac{2\sqrt{3}}{\sqrt{n}}\right)^2}{12} = \frac{1}{3n}$

Xavier init (float) 2010

$E[y] = 0 \Rightarrow E[h(y)] = 0$



2) ReLU



$E[ReLU(\cdot)] \neq 0$

ReLU $\rightarrow Var[w_i] = \frac{2}{n}$

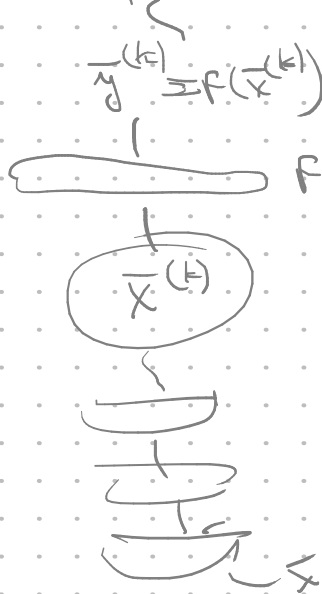
He init Kaiming He

Batch normalization

Covariate shift



Internal covariate shift



$\bar{x}^{(k+1)} = Norm(y^{(k)})$

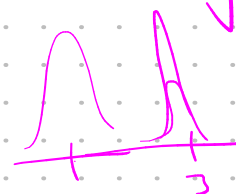
Norm

$\bar{y}^{(k)} = F(x^{(k)})$

F

$\bar{x}^{(k)}$

$x^{(k+1)} = \frac{y^{(k)} - E_D y^{(k)}}{\sqrt{Var_D y^{(k)}}}$



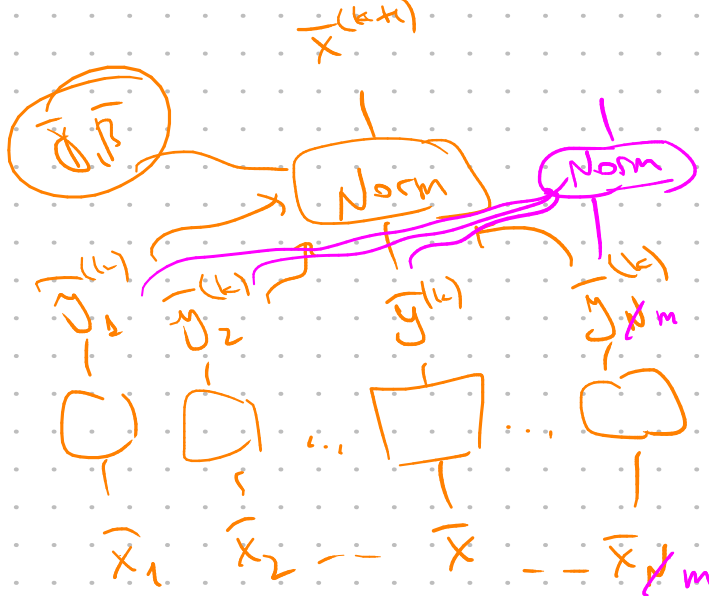
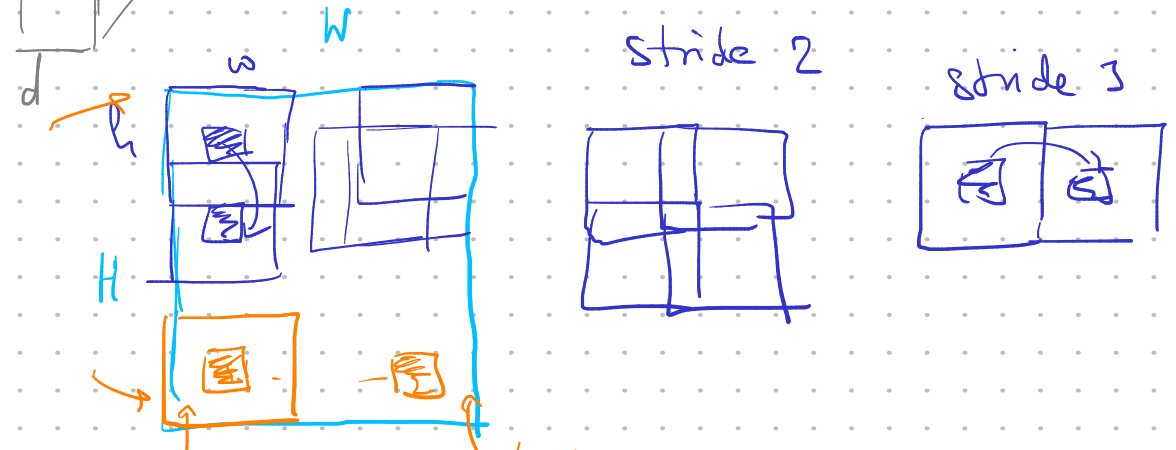
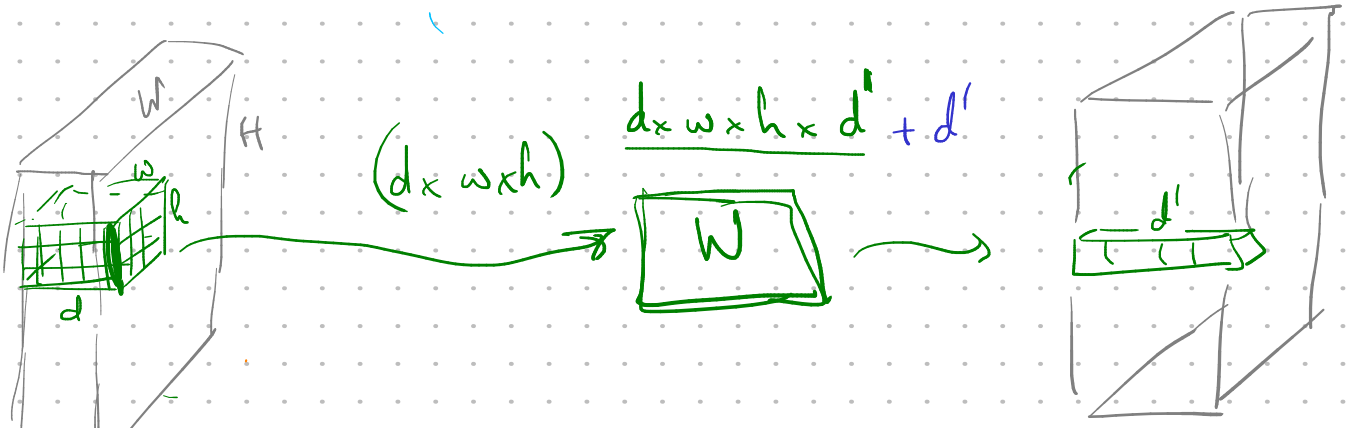
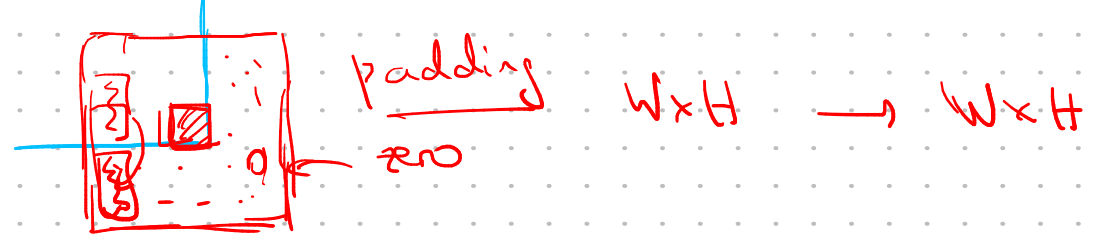


Diagram illustrating the Batch Normalization formula: $BN(x) = \frac{x - E_m \bar{x}}{\sqrt{Var_m \bar{x}}} \odot \gamma + \beta$. The mean and variance are labeled as 'mean' and 'variance' respectively. The gamma and beta parameters are also shown.

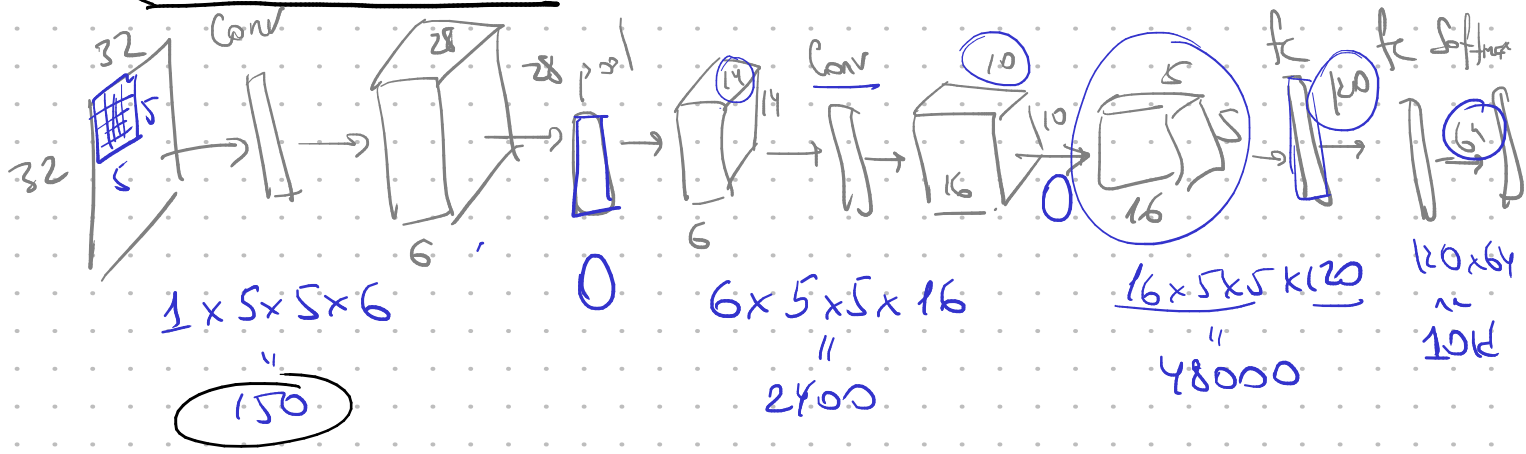
CNN



Stride 1 $W \times H$ $w \times h$ $\rightarrow (W - (w - 1)) \times (H - (h - 1))$

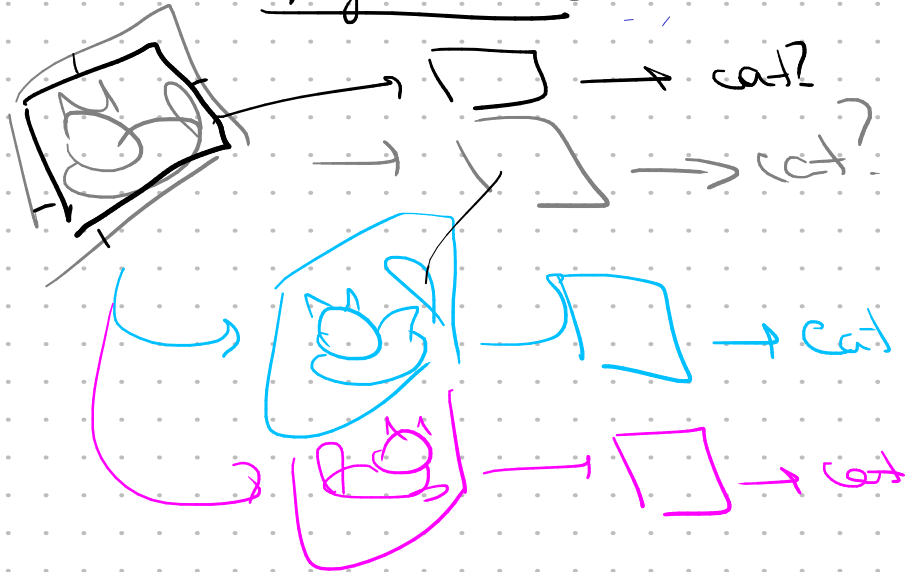


$$(32 \times 32) \times (28 \times 28 \times 6)$$

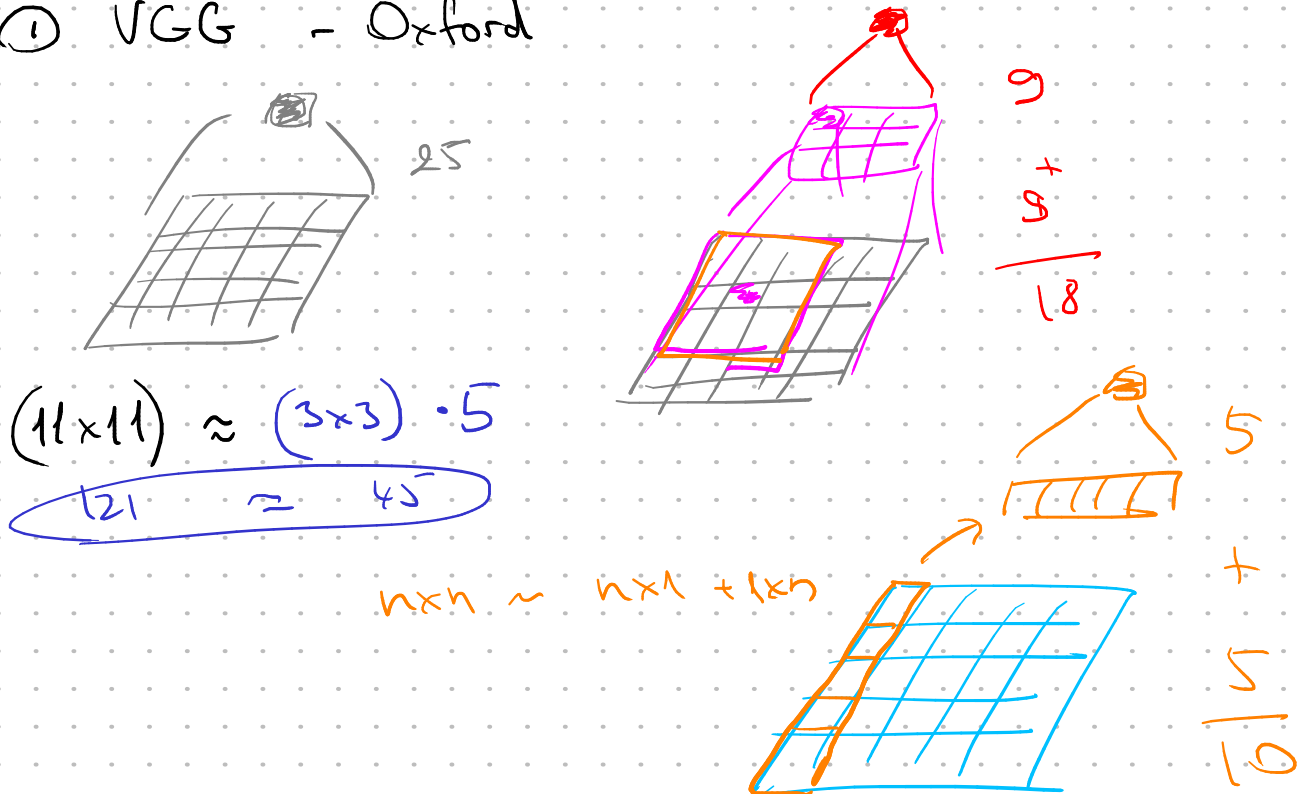


Augmentations

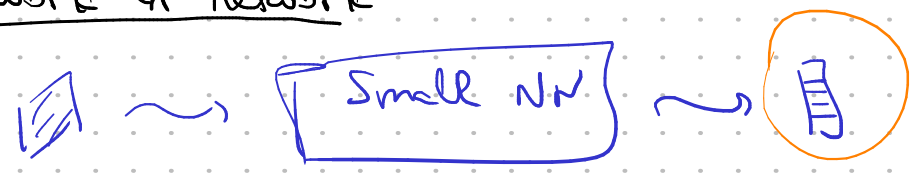
augmentations



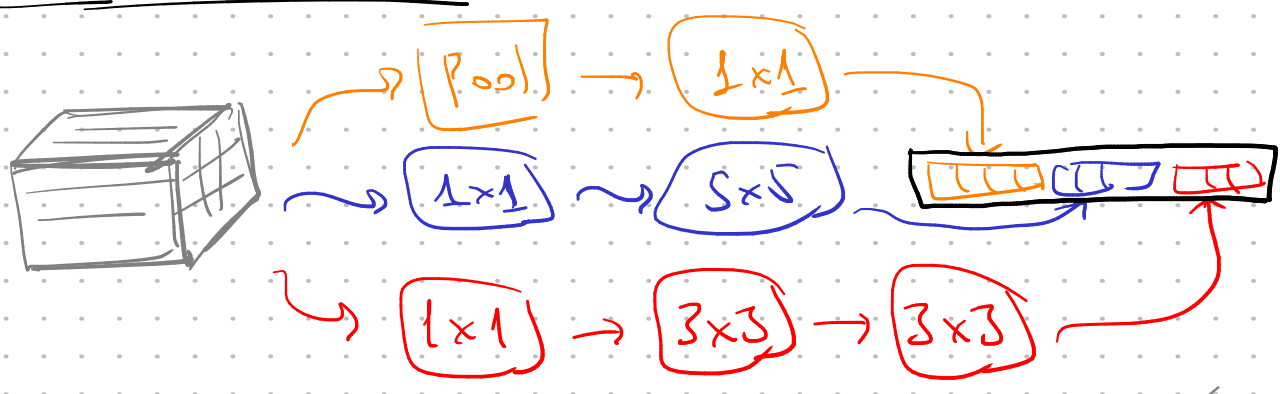
① VGG - Oxford



2 Network in network

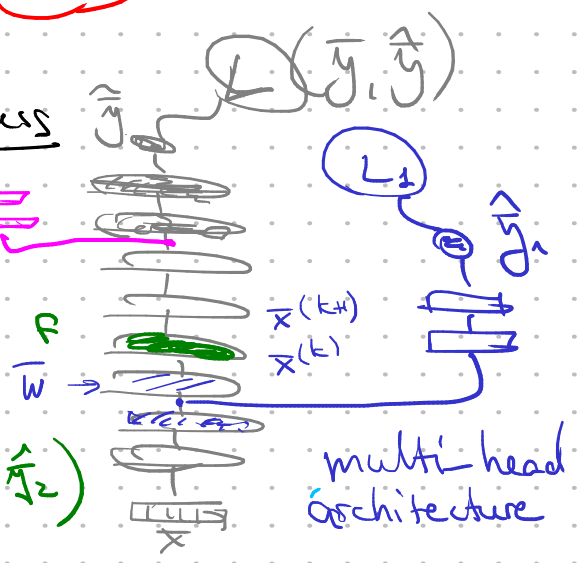


Inception module



3 GoogLeNet - Auxiliary classifiers

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x^{(k+1)}} \frac{\partial x^{(k+1)}}{\partial x^{(k)}} \frac{\partial x^{(k)}}{\partial w}$$



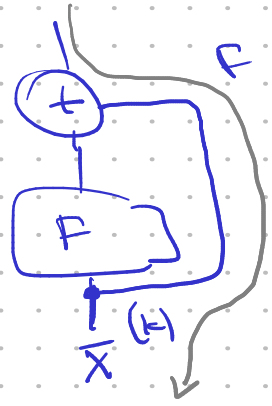
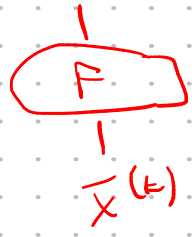
$$L = L(y, \hat{y}) + \lambda_1 L(y, \hat{y}_1) + \lambda_2 L(y, \hat{y}_2)$$

4 Residual connections

Kaiming He

$$\bar{x}^{(k+1)} = F(x^{(k)})$$

$$\bar{x}^{(k+1)} = x^{(k)} + F(x^{(k)})$$

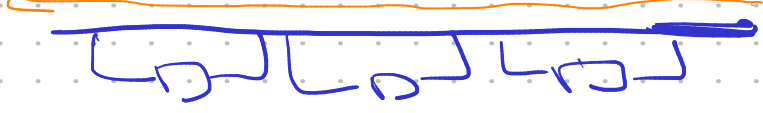


$$F \approx x^{(k)} \rightarrow (\bar{x}^{(k+1)} - x^{(k)})$$

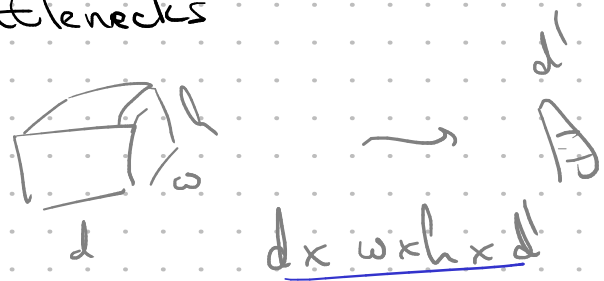
residue

$$\frac{\partial \bar{x}^{(k+1)}}{\partial x^{(k)}} = \frac{\partial F}{\partial x^{(k)}} + I$$

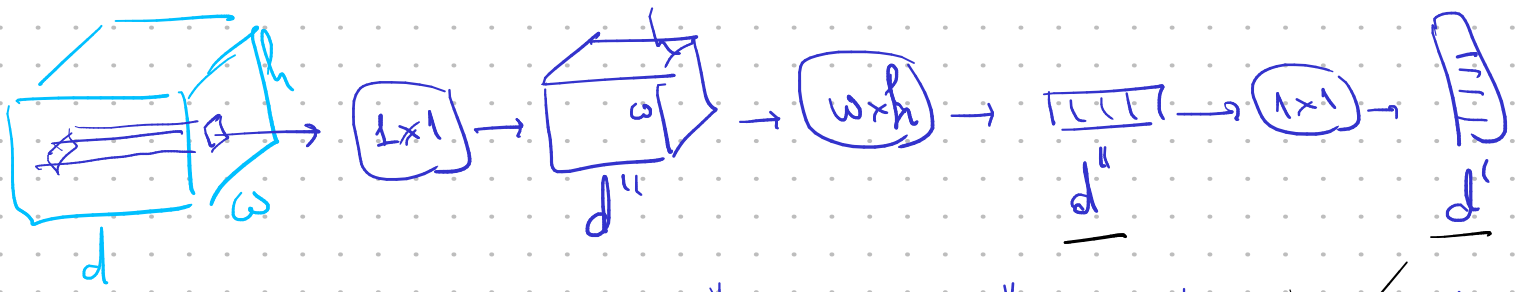
constant error carousel



5 Bottlenecks



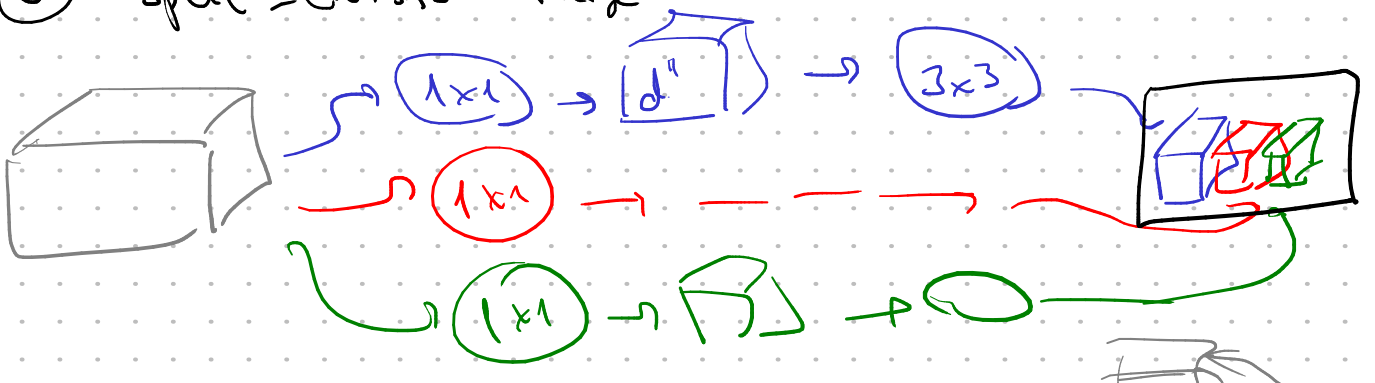
$$256 \times 3 \times 3 \times 256 = 256(256 \cdot 9)$$



$$d \times 1 \times 1 \times d'' + d'' \times w \times h \times d'' + d'' \times 1 \times 1 \times d'$$

$256 \quad 32 \quad 32 \quad 32 \quad 32 \quad 256$
 $256(32 + 4 \cdot 9 + \cancel{32}) \cdot 8$

6 Split-transform-merge



7 DenseNet - skip connections

