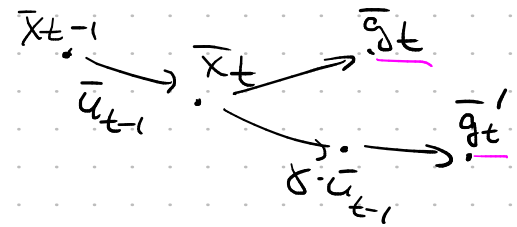


SGD / momentum / NAG

① Adaptive SGD

$F(\bar{x})$

$\bar{x}_t$

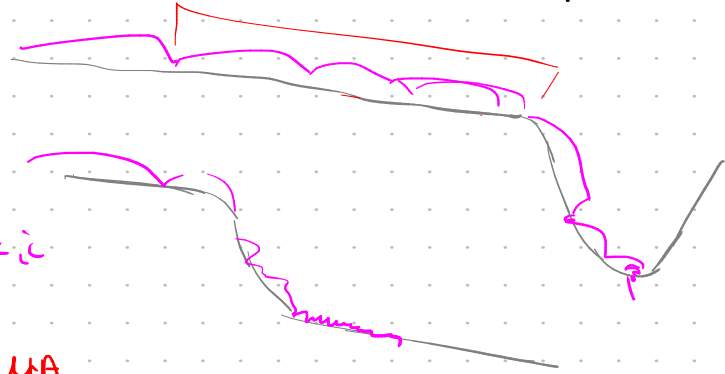


- Adagrad

$$\bar{G}_0 = 0$$

$$G_{t,i} = G_{t-1,i} + g_{t,i}^2$$

$$x_{t+1,i} = x_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} \cdot g_{t,i}$$



- Exponential moving average (EMA)

$$G_{t,i} = G_{t-1,i} \cdot \gamma + (1-\gamma) \cdot g_{t,i}^2 =$$

$$= (1-\gamma) g_{t,i}^2 + (1-\gamma)\gamma g_{t-1,i}^2 + (1-\gamma)\gamma^2 g_{t-2,i}^2 + \dots$$

RMSprop

- AdaDelta

$\bar{x}$  - current  
 $F(\bar{x})$  - next

$\bar{g} = m/c$

$$\bar{x} := \bar{x} - \alpha \cdot \bar{g}_{m/c}$$

$$\bar{x}' := \bar{x} - \frac{\alpha}{\sqrt{\bar{g}^2 + 1}} \cdot g$$

✓

$$\bar{x} := \bar{x} - \frac{1}{2} H^{-1} \bar{g}_{m/c}$$

$$x_{t+1,i} = x_{t,i} - \alpha \cdot \frac{\sqrt{R_{t,i} + \epsilon}}{\sqrt{G_{t,i} + \epsilon}} \cdot g_{t,i}$$

$$R_{t,i} = \gamma \cdot R_{t-1,i} + (1-\gamma) g_{t,i}^2$$

- Adam / Nadam

$$x_{t+1,i} = x_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} \cdot m_{t,i}$$

$$G_{t,i} = \beta_2 G_{t-1,i} + (1-\beta_2) g_{t,i}^2$$

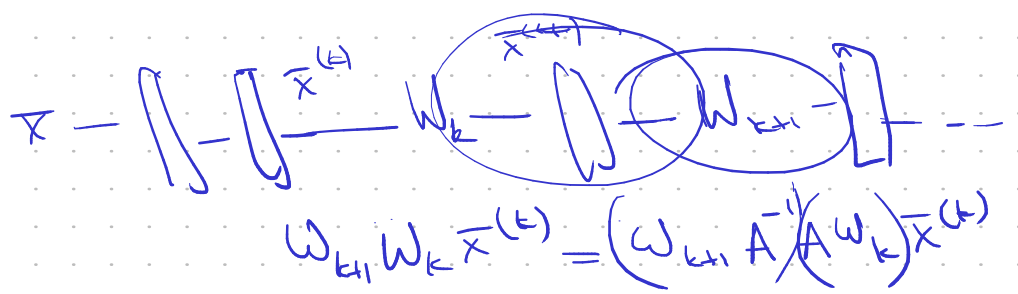
$$m_{t,i} = \beta_1 m_{t-1,i} + (1-\beta_1) g_{t,i}$$

$$\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-6,7}$$

- Weight decay  $L(\bar{w}) + \alpha \cdot \|\bar{w}\|^2$

AdamW  $\hookrightarrow \bar{w} := (1-\eta) \bar{w} - \alpha \nabla_{\bar{w}} L(\bar{w}) = \bar{w} - \alpha \nabla_{\bar{w}} \left( L(\bar{w}) + \frac{\eta}{2\alpha} \|\bar{w}\|^2 \right)$

## ② Dropout

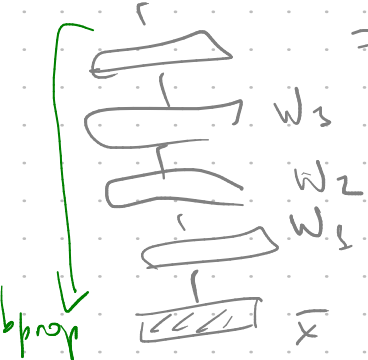


## ③ Weight initialization

2006-2007  
unsupervised  
pretraining

$$= \underbrace{W_k W_{k-1} W_{k-2} \dots W_3 x}_{\text{exploding gradients}}$$

- exploding gradients
- vanishing gradients



$$x \Rightarrow \left( \sum \right) - \bar{w}^T x = y \Rightarrow h - h(\bar{w}^T x)$$

Var[y] ≈ Var[x<sub>i</sub>]

$$y = \sum w_i x_i = \sum y_i$$

Var[y] = n · Var[y<sub>i</sub>]

$$\text{Var}[y_i] = \text{Var}[w_i x_i] = E[w_i^2 x_i^2] - (E[w_i x_i])^2 =$$

$$= E[x_i^2] - \text{Var}[w_i] + E[w_i]^2 \text{Var}[x_i] + \text{Var}[w_i] \cdot \text{Var}[x_i]$$

$$\text{Var}[y] = n \cdot \text{Var}[w_i] \cdot \text{Var}[x_i]$$

$$\text{Var}(\text{Unif}(a,b)) = \frac{(b-a)^2}{12}$$

$$w_i \sim \text{Unif}\left(-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right)$$

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

(Glorot)

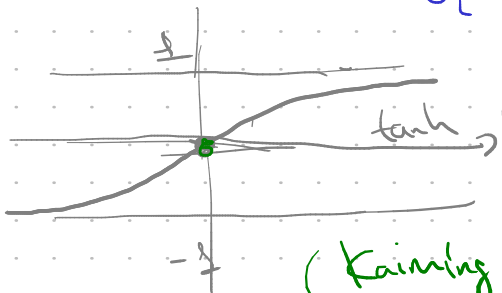
$$O(1/n) \approx 1$$

Xavier init

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

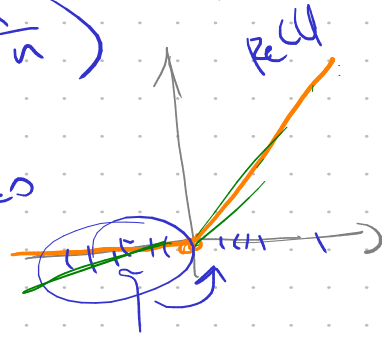
$$w_i \sim \mathcal{N}\left(0, \frac{1}{n}\right)$$

$$E[\text{ReLU}(\cdot)] > 0$$



(Kaiming He)

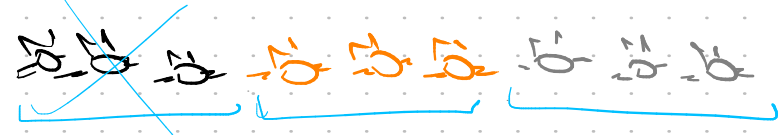
He init:  $\text{Var}[w_i] = 2/n$



## ④ Batch normalization

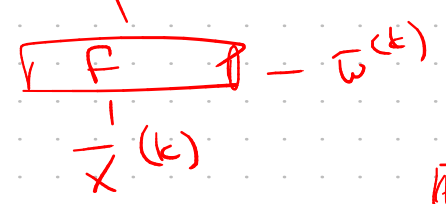
intercept covariate shift

covariate shift:  $\bar{w}^{(k+1)} - \bar{w}^{(k)}$

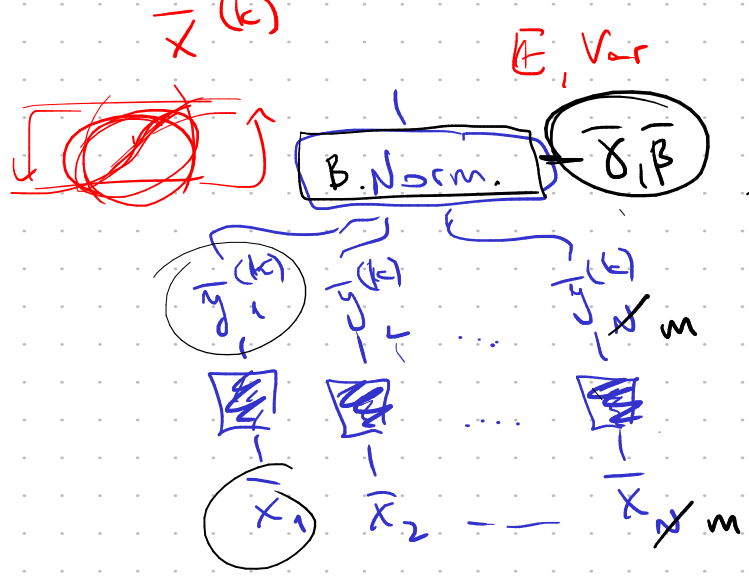


**Norm**  $\cdot \bar{x}^{(k+1)} = \text{Norm}(\bar{y}^{(k)}) = \frac{\bar{y}^{(k)} - \mathbb{E}[\bar{y}^{(k)}]}{\sqrt{\text{Var}[\bar{y}^{(k)}]}}$

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

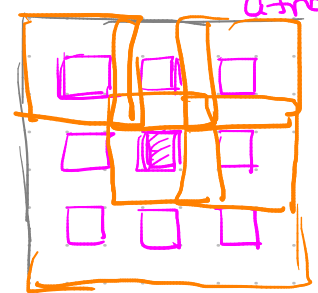
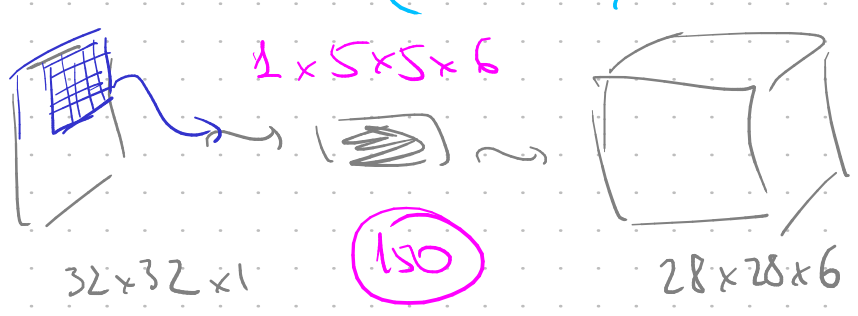
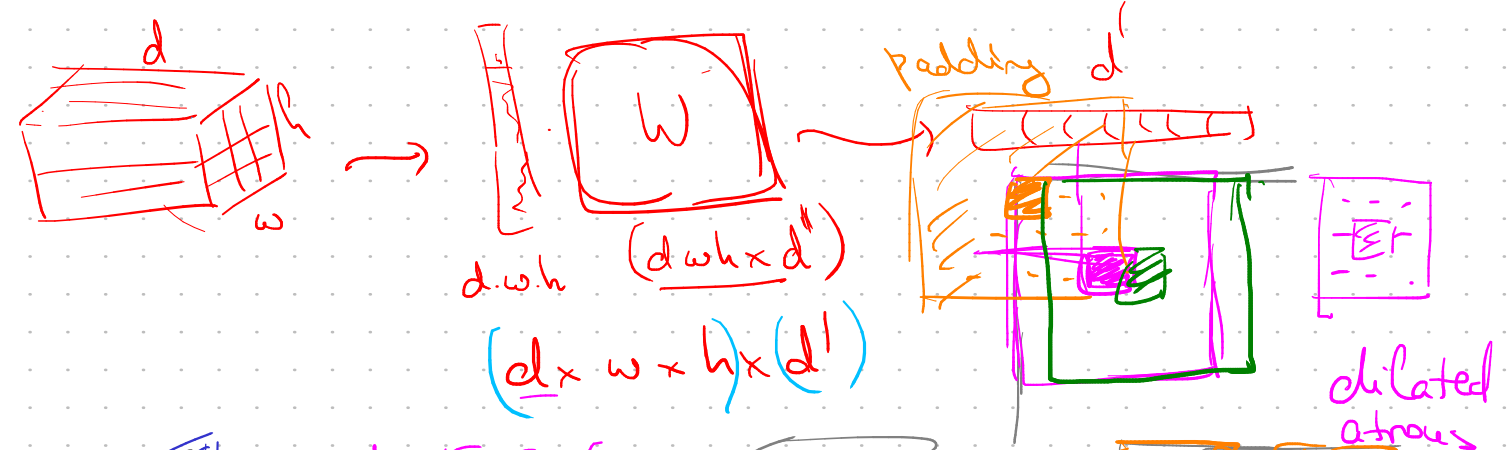
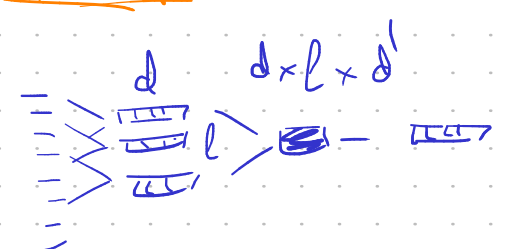
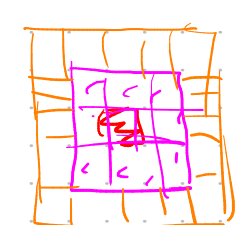
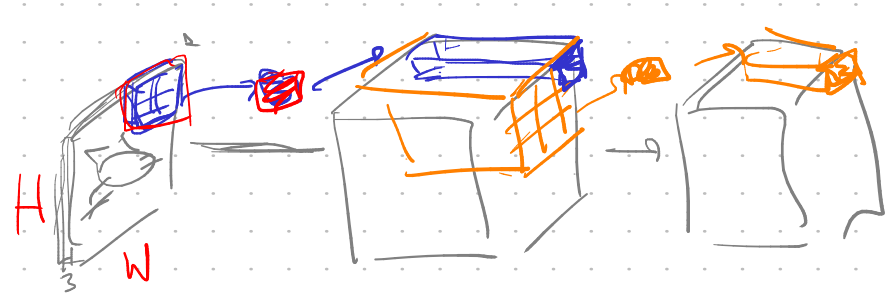


No boxogram cross - layer normalization



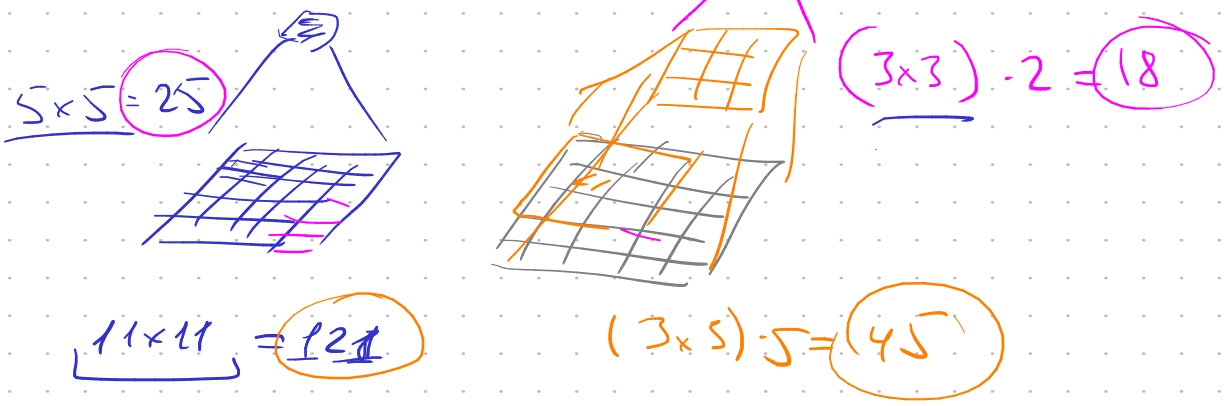
$\text{BN}(\bar{y}^{(k)}) = \frac{\bar{y}^{(k)} - \mathbb{E}_m[\bar{y}^{(k)}]}{\sqrt{\text{Var}_m[\bar{y}^{(k)}]}} \odot \bar{\gamma} + \bar{\beta}$

5 CNN



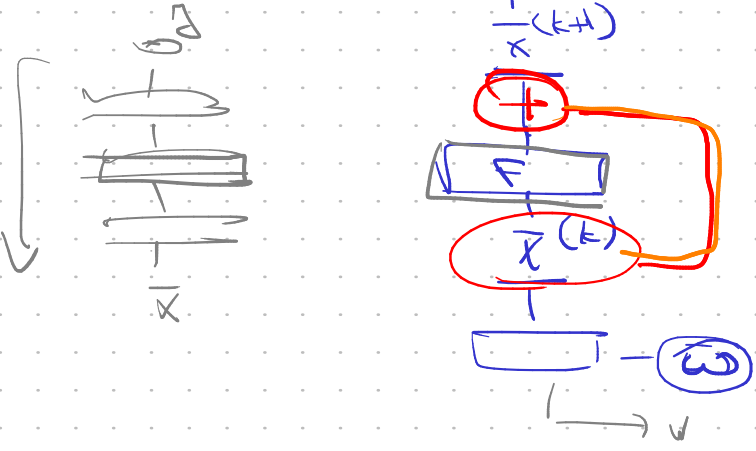
$(32 \times 32) \times (28 \times 28 \times 6)$

6 VGG

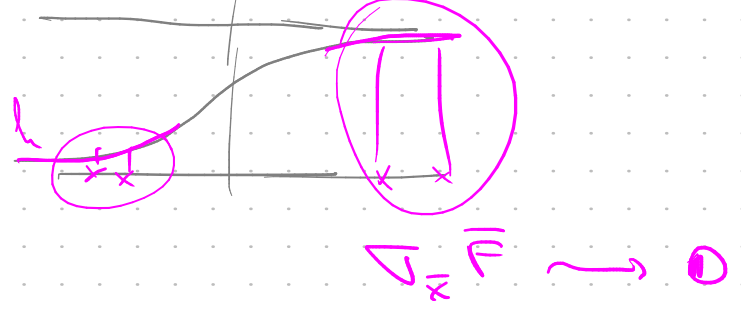


7 Network in Network / Inception

8 Residual connections



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



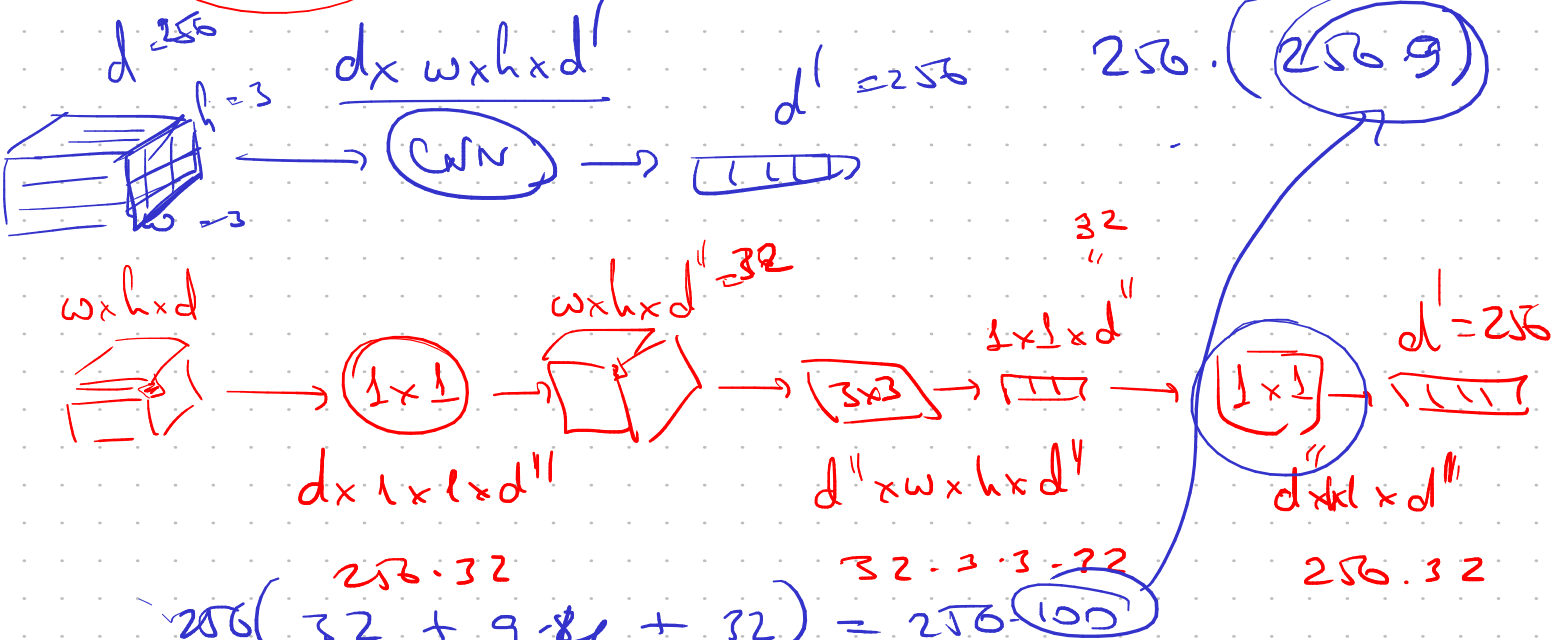
$$\bar{x}^{(k+1)} = f(\bar{x}^{(k)}) + \bar{x}^{(k)}$$

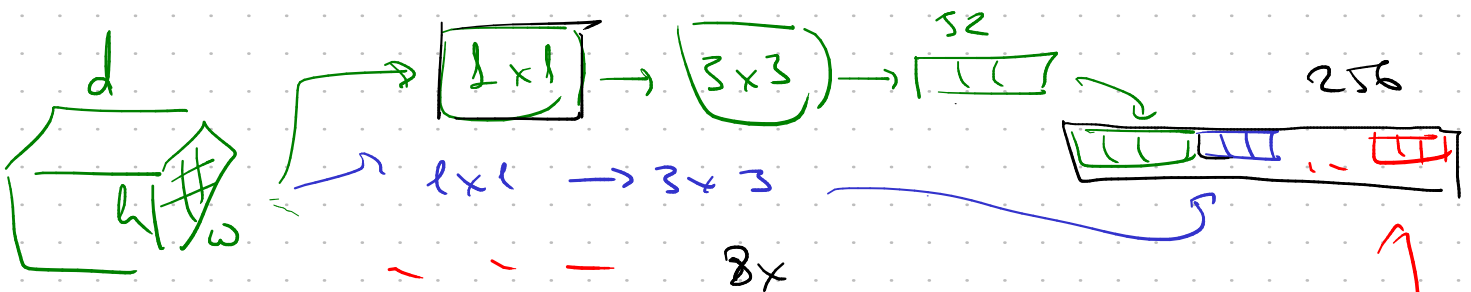
$$F: \bar{x}^{(k)} \mapsto \bar{x}^{(k+1)} - \bar{x}^{(k)}$$

residual

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

9 Bottlenecks / split-transform-merge





$1 \times 1 \rightarrow 3 \times 3$   
 $8 \times$

$1 \times 1 \rightarrow 3 \times 3$

$256 \cdot 800$  vs  $256 \cdot (256 \cdot 9)$