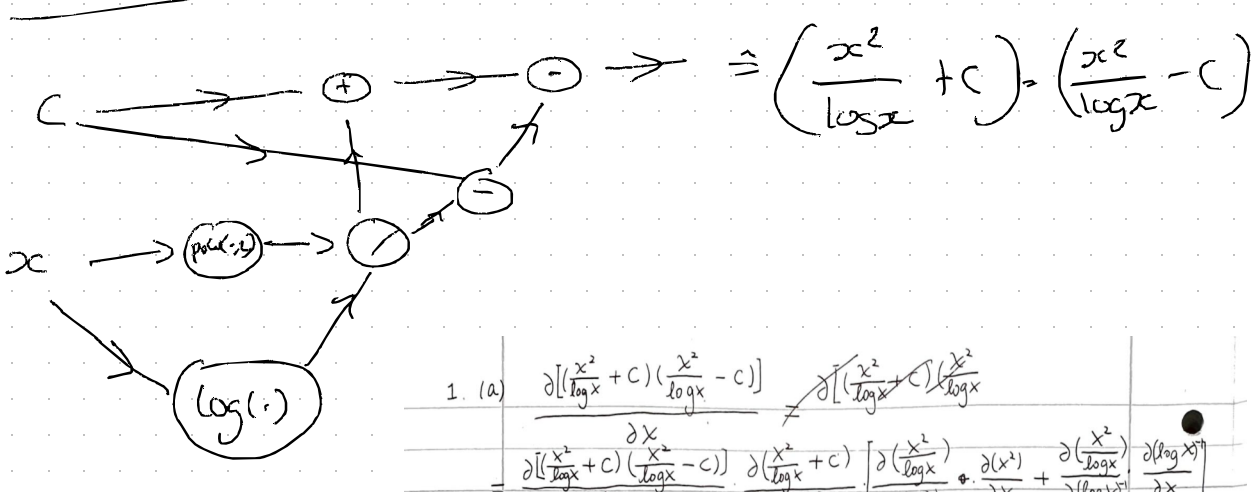


Sheet 07



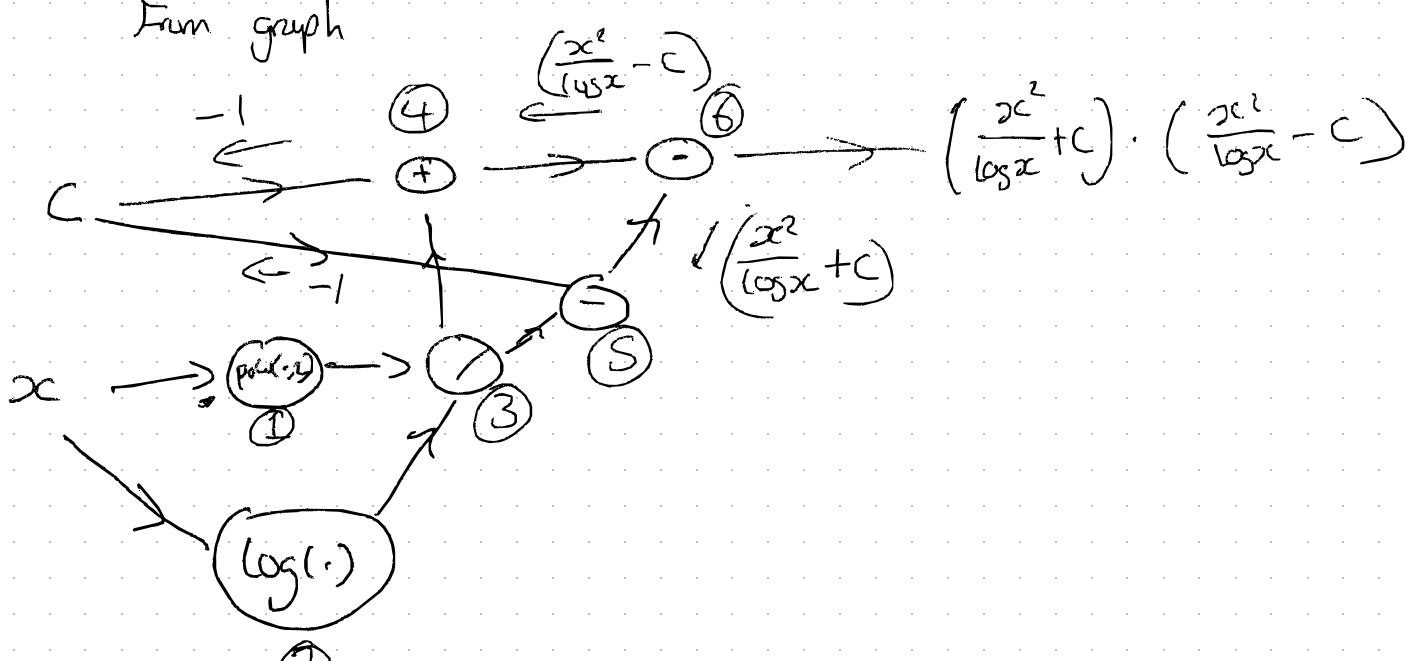
$$1. (a) \frac{\partial}{\partial x} \left[ \left( \frac{x^2}{\log x} + c \right) \left( \frac{x^2}{\log x} - c \right) \right] = \frac{\partial}{\partial x} \left[ \frac{x^2}{\log x} + c \right] \cdot \frac{\partial}{\partial x} \left[ \frac{x^2}{\log x} - c \right] + \left( \frac{x^2}{\log x} + c \right) \cdot \frac{\partial}{\partial x} \left[ \frac{x^2}{\log x} - c \right]$$

a)  $\frac{x^4}{\log^2 x} - c^2$

derivative:  $\frac{\partial}{\partial x} \left( \frac{x^4}{\log^2 x} \right) = \frac{4x^3}{\log^2 x} - \frac{2x^4}{\log^3 x} \cdot \frac{1}{x}$

$$= 2x^3 \left( \frac{2}{\log^2 x} - \frac{1}{\log^3 x} \right)$$

From graph



labelled nodes

$$u_1 = x^2, \quad u_2 = \log x, \quad u_3 = \frac{u_1}{u_2} = \frac{x^2}{\log x}$$

$$u_4 = u_3 - c, \quad u_5 = u_3 + c, \quad u_6 = u_4 \cdot u_5$$

Derivatives:

1)  $\frac{\partial u_1}{\partial x} = 2x$

2)  $\frac{\partial u_2}{\partial x} = \frac{1}{x}$

3)  $\frac{\partial u_3}{\partial x} = \frac{1}{u_2} \frac{\partial u_1}{\partial x} - \frac{u_1}{u_2^2} \frac{\partial u_2}{\partial x} = \frac{2x \log x - x}{\log^2 x}$

4)  $\frac{\partial u_4}{\partial x} = \frac{\partial u_3}{\partial x}$

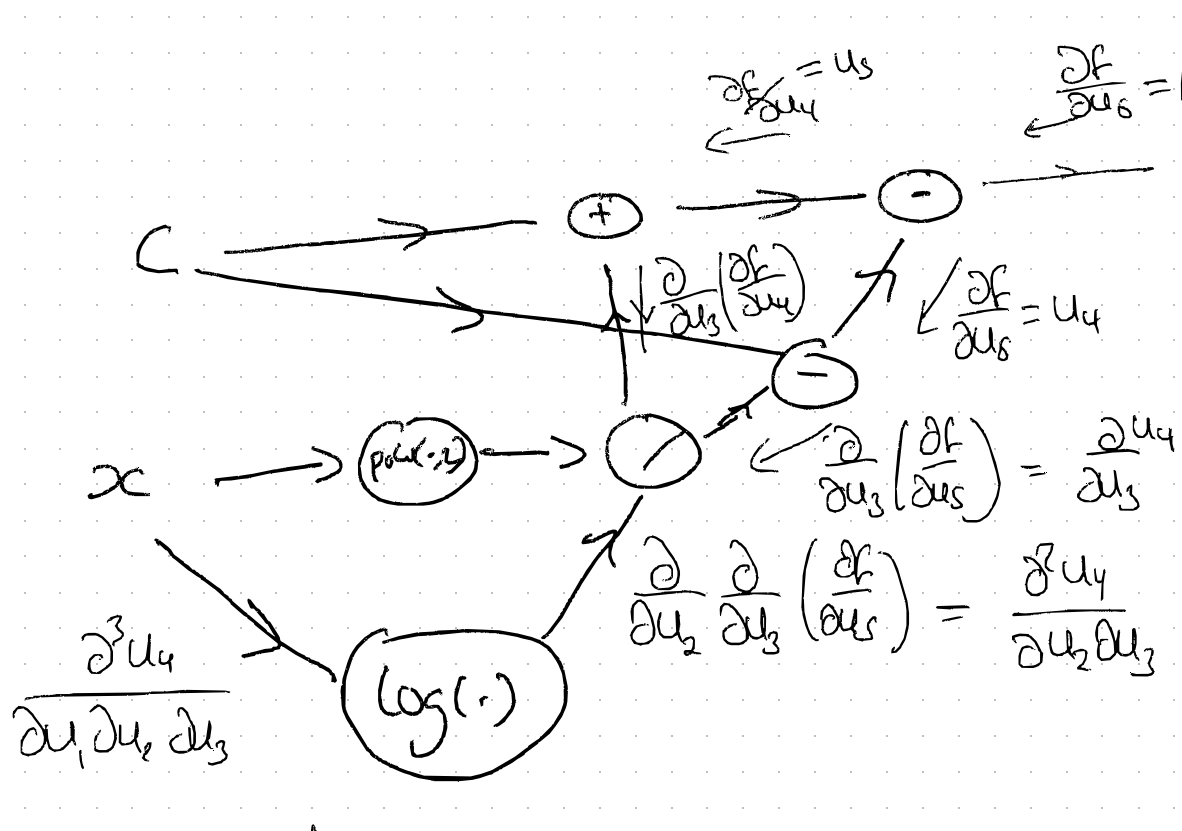
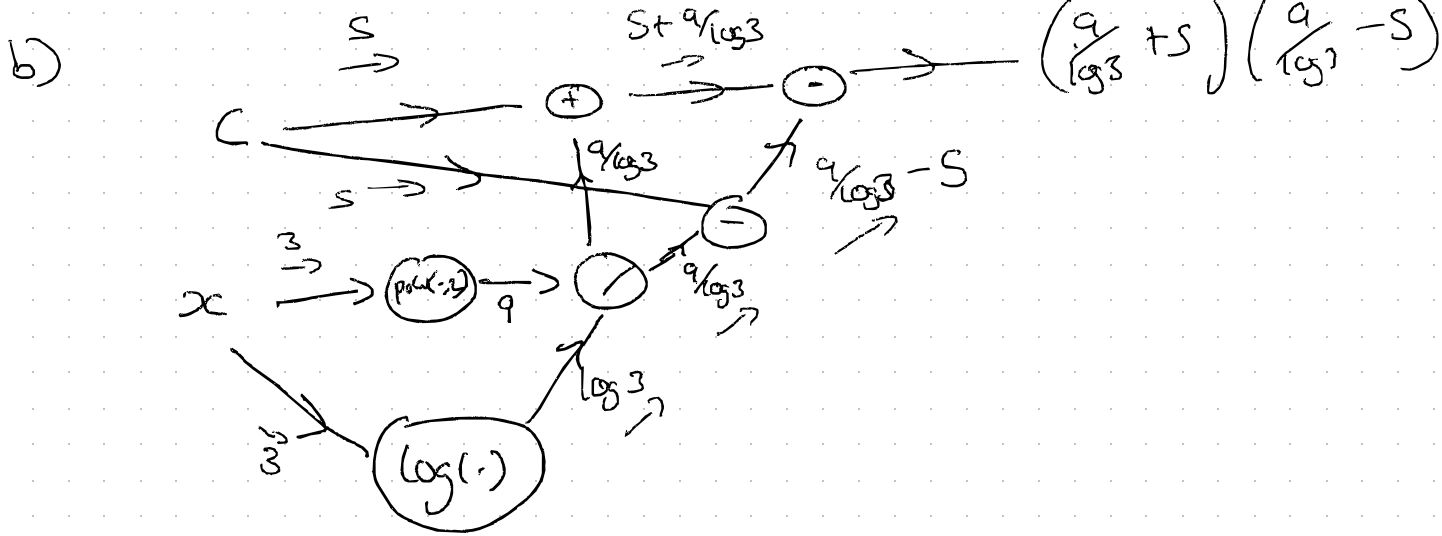
5)  $\frac{\partial u_5}{\partial x} = \frac{\partial u_3}{\partial x}$

6)  $\frac{\partial u_6}{\partial x} = u_4 \frac{\partial u_5}{\partial x} + u_5 \frac{\partial u_4}{\partial x}$

$$\Rightarrow f = u_6 \Rightarrow \frac{\partial f}{\partial x} = u_4 \frac{\partial u_5}{\partial x} + u_5 \frac{\partial u_4}{\partial x}$$

$$\frac{\partial u_5}{\partial x} = \frac{\partial u_4}{\partial x} = \frac{\partial u_3}{\partial x}$$

$$\Rightarrow \frac{2x \log x - x}{\log^2 x} \left[ \left( \frac{x^2}{\log x} - c \right) + \left( \frac{x^2}{\log x} + c \right) \right] = \left[ \frac{2x \log x - x}{\log^2 x} \right] \frac{2x^2}{\log x}$$



etc.

$\Rightarrow$  Backprop:

①  $\frac{\partial f}{\partial u_6} = 1$

②  $\frac{\partial f}{\partial u_4} = u_5, \quad \frac{\partial f}{\partial u_5} = u_4$

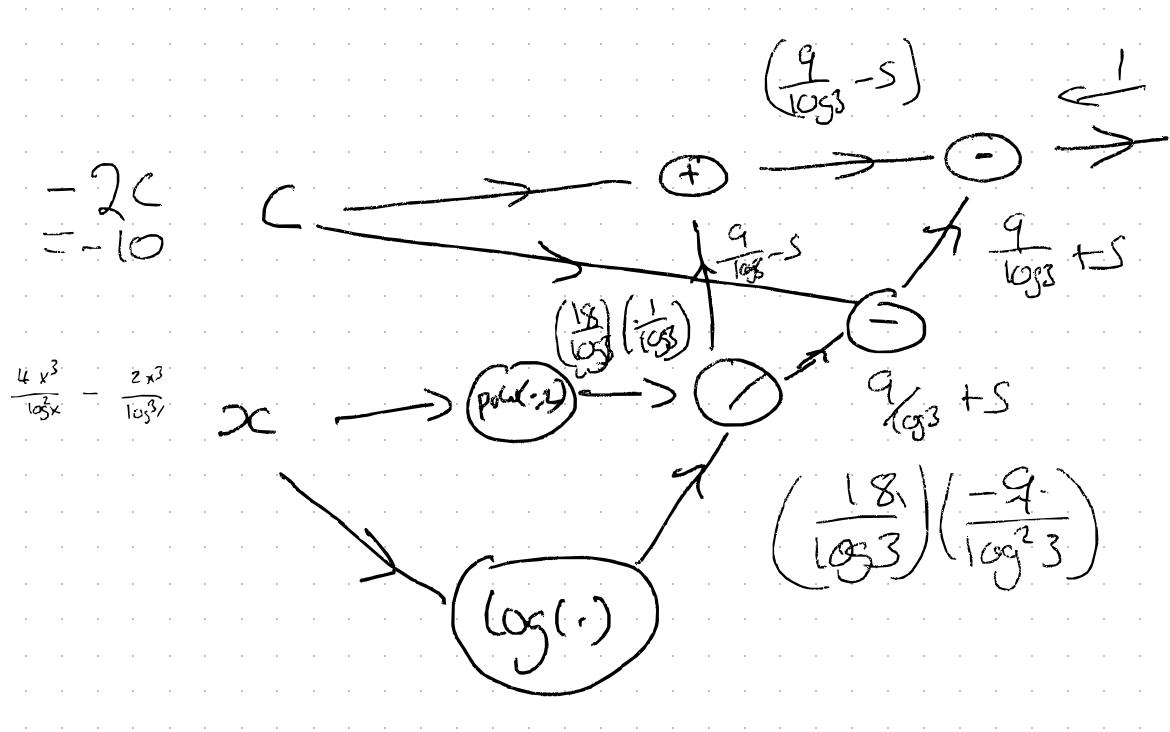
③  $\frac{\partial f}{\partial c} = -1$   
 $\frac{\partial f}{\partial c} = -\frac{\partial u_4}{\partial c} + 1 = 0$

④  $\frac{\partial f}{\partial u_3} = \frac{\partial f}{\partial u_4} \frac{\partial u_4}{\partial u_3} = u_5 \frac{\partial u_4}{\partial u_3}$   
 $\frac{\partial f}{\partial u_3} = \frac{\partial f}{\partial u_5} \frac{\partial u_5}{\partial u_3} = u_4 \frac{\partial u_5}{\partial u_3}$   
 $\left( \frac{x^2}{\log x} - c \right) \cdot \frac{\partial}{\partial \left( \frac{x^2}{\log x} \right)} \left( \frac{x^2}{\log x} + c \right) + \dots = \frac{2x^2}{\log^2 x}$

⑤  $\frac{\partial f}{\partial u_2} = \frac{\partial f}{\partial u_3} \frac{\partial u_3}{\partial u_2} = \left( \frac{2x^2}{\log^2 x} \right) \frac{\partial}{\partial \log x} \left( \frac{x^2}{\log x} \right)$

⑥  $\frac{\partial f}{\partial u_1} = \frac{\partial f}{\partial u_3} \frac{\partial u_3}{\partial u_1} = \left( \frac{2x^2}{\log^2 x} \right) \frac{\partial}{\partial x^2} \left( \frac{x^2}{\log x} \right)$

New calculating values



very tedious by hand

$$\frac{2x^2}{\log x} \left( -\frac{x^2}{\log^2 x} \frac{\partial}{\partial x} \log x + \frac{1}{\log x} \frac{\partial}{\partial x} x^2 \right)$$

$$= \frac{2x^2}{\log x} \left( -\frac{2x}{\log^2 x} + \frac{2x}{\log x} \right) = \boxed{\frac{4x^3}{\log^2 x} - \frac{2x^3}{\log^3 x}}$$

$$\frac{\partial f}{\partial c} = -2c$$

$$\frac{\partial f}{\partial x}$$

d) using python

$$\frac{\partial f}{\partial x} = 48.757 \checkmark$$

$$\frac{4 \times 27}{\log^2 3} - \frac{2 \times 27}{\log^3 3} = 48.757 \checkmark$$

$$\frac{\partial f}{\partial c} = -10.0 = -2c \checkmark$$

$\rightarrow$  as predicted

```
import torch
```

✓ 8.4s

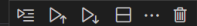
Python

```
# Define x and c as PyTorch tensors with requires_grad=True to enable gradient computation
x = torch.tensor(3.0, requires_grad=True) # x = 3
c = torch.tensor(5.0, requires_grad=True) # c = 5

# Define the function f(x, c)
u1 = x**2
u2 = torch.log(x)
u3 = u1 / u2
u4 = u3 + c
u5 = u3 - c
f = u4 * u5 # f = (u3 + c) * (u3 - c)
```

✓ 0.1s

Python



```
# Compute the gradients using backward()
f.backward()

# Output the computed derivatives
df_dx = x.grad # ∂f/∂x
df_dc = c.grad # ∂f/∂c

print(f"Derivative with respect to x: {df_dx.item()}")
print(f"Derivative with respect to c: {df_dc.item()}")
```

✓ 0.2s

Python

```
Derivative with respect to x: 48.756893157958984
Derivative with respect to c: -10.0
```

## 2) ADAM optimiser

① 
$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

momentum (mean)      ← gradient

decays rate, often  $\beta = 0.9$

↳ moving average of gradients

② 
$$v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2$$

Second moment (variance)      ← grad<sup>2</sup>

↳ moving average of squared gradients.

③ Bias correction:

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t} \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

↳ Adjust for initialisation bias

④ Update parameters

$$w_{t+1} = w_t - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t^2 + \epsilon}}$$

learning rate

prevent poles

updates weights with scaling  $\sim \frac{\hat{m}_t}{\sqrt{v_t}}$

b) Initialise to zero

$$\Rightarrow m_t = (1-\beta_1) g_t \Rightarrow \hat{m}_t = g_t$$

$$v_t = (1-\beta_2) g_t^2 \Rightarrow \hat{v}_t = g_t^2$$

$$w_{t+1} = w_t - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

$$\frac{\hat{m}_t}{\sqrt{\hat{v}_t}} = \frac{g}{\sqrt{g^2}} = \boxed{\text{Sign}(g)}$$

c) Second iteration:

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

$$v_{t+1} = \beta_2 (1-\beta_2) g_t^2 + (1-\beta_2) g_{t+1}^2$$

$$\hat{m}_{t+1} = \frac{\beta_1 g_t + g_{t+1}}{1+\beta_1}$$

$$\hat{v}_{t+1} = \frac{\beta_2 g_t^2 + g_{t+1}^2}{1+\beta_2}$$

$$\Rightarrow \frac{\hat{m}_{t+1}}{\sqrt{\hat{v}_{t+1}}} = \frac{g_{t+1}(1 + g_t/g_{t+1}) \frac{\sqrt{1+\beta_2}}{1+\beta_1} \text{Sign}(g) \times \frac{1+\epsilon}{\sqrt{1+\epsilon^2}} \frac{\sqrt{1+\beta_1}}{1+\beta_1}}{\sqrt{g_{t+1}^2 (1 + (g_t/g_{t+1})^2)} \frac{\sqrt{1+\beta_2}}{1+\beta_1}}$$

Assuming  $|g_t/g_{t+1}|$  small

$$\approx \text{Sign}(g) (1+\epsilon) (1 - \frac{1}{2} \epsilon^2) \sqrt{\frac{1+\beta_2}{(1+\beta_1)^2}}$$

$$= \text{Sign}(g) (1+\epsilon) \sqrt{\frac{1+\beta_2}{(1+\beta_1)^2}} + O(\epsilon^2)$$

↳ Small perturbation to gradient  
extra weighting

d) ~~Gradient rescaling is possible solution~~

use ~~$$\hat{g}_t = \frac{g_t}{1-\beta_1^t}$$~~

Alternatively: change learning rate  $\rightarrow$  Start small, and then grow.

↳ Adaptive learning rate

e) MLP: weights  $w$  trained with Adam and L2 Reg.

Is there difference with L2 penalty  $\|w\|_2^2$  in loss or weight decay directly onto weights

L2 reg:

$$L' = L + \lambda \|w\|_2^2$$

For Adam:  $g_t = \nabla L$

$$\Rightarrow g_t \rightarrow \nabla L' = \boxed{\nabla L + 2\lambda w}$$

↳  $w_{t+1}$  affected implicitly

↳ Additional term in grad.

↳ Effective strength of Reg varies

weight decay:

$$w_{t+1} = w_t - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} - \underbrace{\alpha \lambda w_t}_{\text{AdamW}}$$

↳ additional decoupled weight decay

↳ does not interfere with adaptive gradient descent

$\Rightarrow$  often favored

$\Rightarrow$  const. effective strength of grad.

3)

work backwards from final  
 $7 \times 7 \times 512$



$\Rightarrow$  Field of view is  $7 \times 7$

Max pooling  $2 \times 2 \rightarrow$  Receptive field doubles (going backward)

Conv  $3 \times 3 \Rightarrow$  New field = old +  $2 \times$  Dilation

Here Dilation = 1

$\Rightarrow$  new = old + 2

working backwards:

$\square \rightarrow$

Block 5: Max pool,  $3 \times$  conv

Backwards  $3 \text{ conv} \Rightarrow$  Max pool

1 pixel in  $7 \times 7$  new  $1 \rightarrow 2$  Max pool

$3 \rightarrow 5$

$5 \rightarrow 7$

$7 \rightarrow 14$  Conv

Block 4:

Rev:  $3 \text{ conv}$ , max pool.

$8 \rightarrow 16$

$16 \rightarrow 18$

$18 \rightarrow 20$

$20 \rightarrow 22$

Block 3: Same again

$22 \rightarrow 44$

$44 \rightarrow 46$

$46 \rightarrow 48$

$48 \rightarrow 50$

Block 2: Max pool,  $2 \text{ conv}$

$50 \rightarrow 100$

$100 \rightarrow 102$

$102 \rightarrow 104$

Block 1:

$104 \rightarrow 208$

$208 \rightarrow 210$

$210 \rightarrow 212$

Each pixel in  $7 \times 7$  has info from

$212 \times 212$  subset of  
 initial image

$$(b) \text{ \# parameters} = \cancel{3 \times 3 \times 64}$$

(with bias)  
in convolutional  
layers

$$\begin{aligned} & (3 \times 3 \times 3 + 1) \times 64 + (3 \times 3 \times 64 + 1) \times 64 \\ & + (3 \times 3 \times 64 + 1) \times 128 + (3 \times 3 \times 128 + 1) \times 128 \\ & + (3 \times 3 \times 128 + 1) \times 256 + (3 \times 3 \times 256 + 1) \times 256 \times 2 \\ & + (3 \times 3 \times 256 + 1) \times 512 + (3 \times 3 \times 512 + 1) \times 512 \times 2 \\ & + \cancel{(3 \times 3 \times 512 + 1) \times 512} + (3 \times 3 \times 512 + 1) \times 512 \times 3 \end{aligned}$$

# parameters in  
fully connected  
layers

$$\begin{aligned} & = 7 \times 7 \times 512 \times 4096 + 4096 \\ & + 4096 \times 4096 + 4096 \\ & + 4096 \times 1000 + 1000 \\ & = 123642856 \approx 1.24 \times 10^8 \end{aligned}$$

$$\text{\# parameters in total} = 138357544 \approx 1.38 \times 10^8$$

$$\frac{\cancel{\text{\# p of so}} \text{\# p of convolutional}}{\text{\# p of fully connected}} = \cancel{12} 11.90\%$$