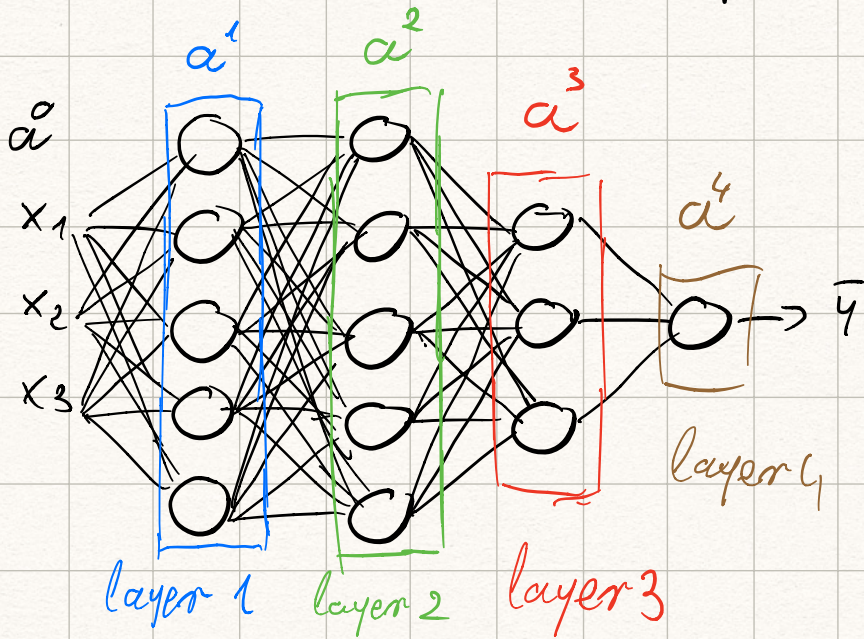


# Forward propagation in deep networks



$$1) \quad z^1 = w^1 a^0 + b^1$$

$$a^1 = f^1(z^1)$$

$$2) \quad z^2 = w^2 a^1 + b^2$$

$$a^2 = f^2(z^2)$$

$$3) \quad z^3 = w^3 a^2 + b^3$$

$$a^3 = f^3(z^3)$$

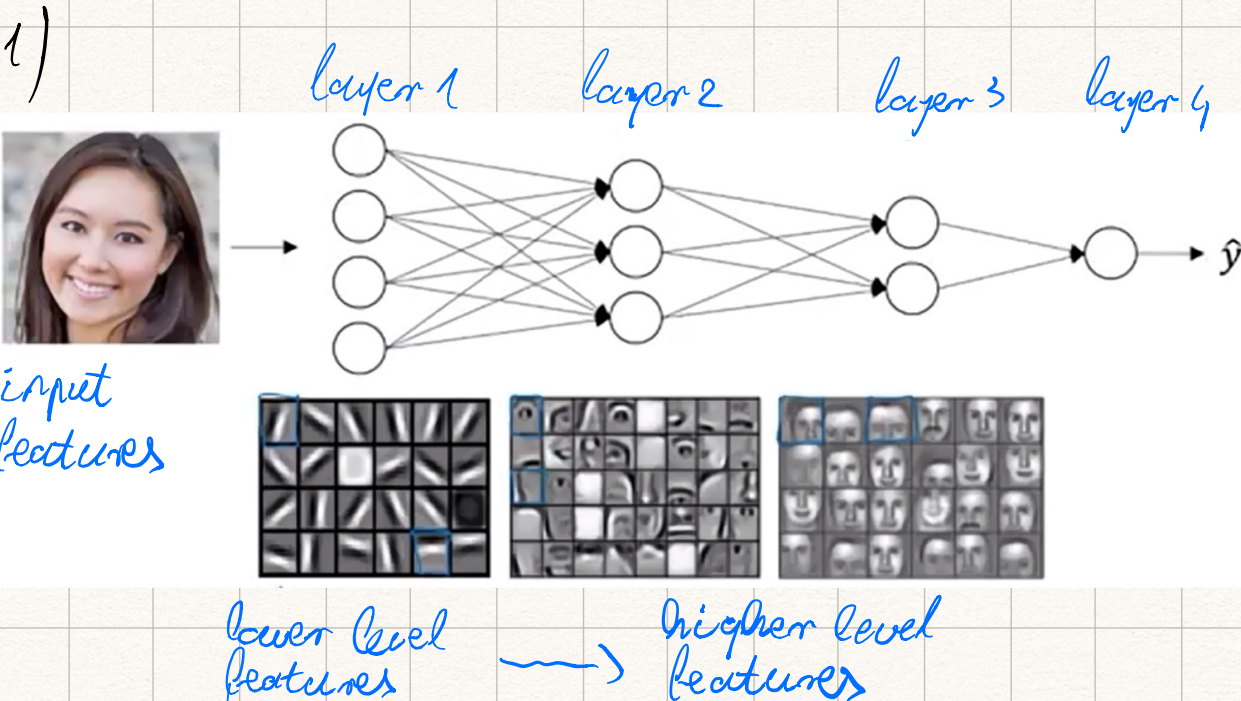
$$4) \quad z^4 = w^4 a^3 + b^4$$

$$a^4 = f^4(z^4)$$

$$\Rightarrow \quad z^l = w^l \cdot a^{l-1} + b^l$$

$$a^l = f^l(z^l)$$

# Why deep representation



2) shallower networks require exponentially more hidden units

# Building blocks of deep networks

**Forward:** Input  $\Rightarrow a^{l-1}$ , output  $\Rightarrow a^l$ , cache  $\Rightarrow z^l$

$$z^l = w^l a^{l-1} + b^l \quad \text{for layer } l$$

$$a^l = f^l(z^l)$$

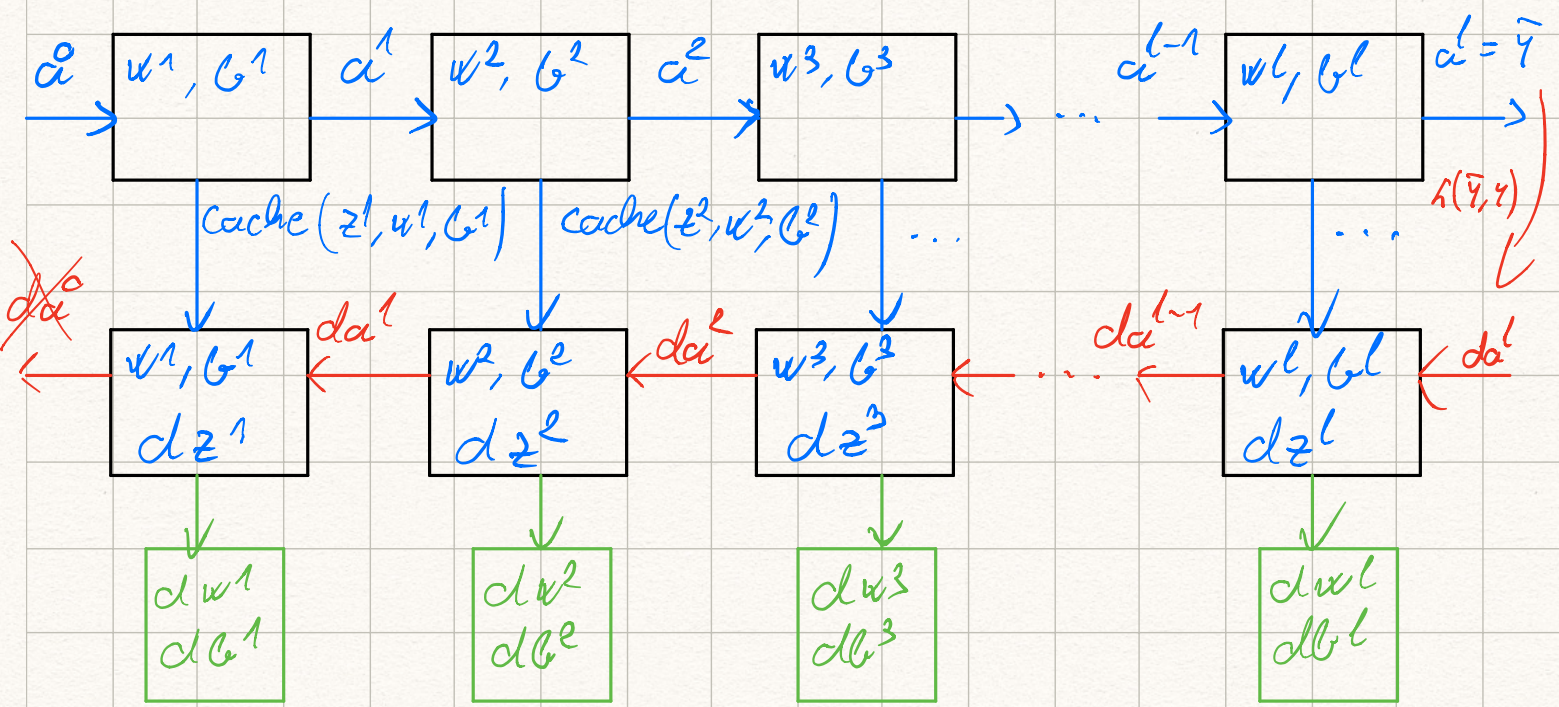
**Backward:** Input  $\Rightarrow da^l$ , output  $\Rightarrow da^{l-1}$   
 cache  $\Rightarrow z^l$   $dw^l$   
 $db^l$

$$dz^l = w^{l+1 T} dz^{l+1} * f^{l+1}(z^l)$$

$$dw^l = dz^l \cdot a^{l-1} \quad \text{for layer } l$$

$$db^l = dz^l$$

$$da^{l-1} = w^{l T} \cdot dz^l$$



$$w^l := w^l - \eta dw^l$$

$$b^l := b^l - \eta db^l$$

# Hyperparameters

Parameters:  $w^l, b^l$

Hyperparameters:

- learning rate  $\eta$
- # iterations
- # hidden layers
- # hidden units
- choice of activation function
- momentum
- minibatch size
- regularization
- 

Hyperparameters control  $w^l$  and  $b^l$ .

Finetuning hyperparameters is an empirical process.