

<https://introml.mit.edu/>

6.390 Intro to Machine Learning

Spring 25 Final Review

Shen Shen

May 13, 2025

5pm, Room 34-101

Outline

- Semester Overview
- Weekly Review
- Exam Prep Resources, live Q&A
- Future courses

- Intro to ML
 - Linear Regression
 - Gradient Descent
 - Linear Classification
-
- Features, Neural Networks I
 - Neural Networks II
 - Convolutional Neural Networks
 - Representation Learning
 - Transformers
-
- Markov Decision Processes
 - Reinforcement Learning
-
- Non-parametric Models

Many other ways to dissect

Model class

- linear models
- linear model on non-linear features
- fully connected feed-forward
- convolutional
- transformers
- Q-table
- tree, k-nearest neighbor, k-means

Learning process:

- training/validation/testing
- overfitting/underfitting
- regularization
- hyper parameters

Modeling choices:

- Supervised:
 - regression
 - classification
- Unsupervised/self-supervised
- Reinforcement/sequential

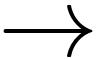
Optimization:

- analytical solutions
- gradient descent
- back propagation
- value iteration, Q-learning
- non-parametric methods

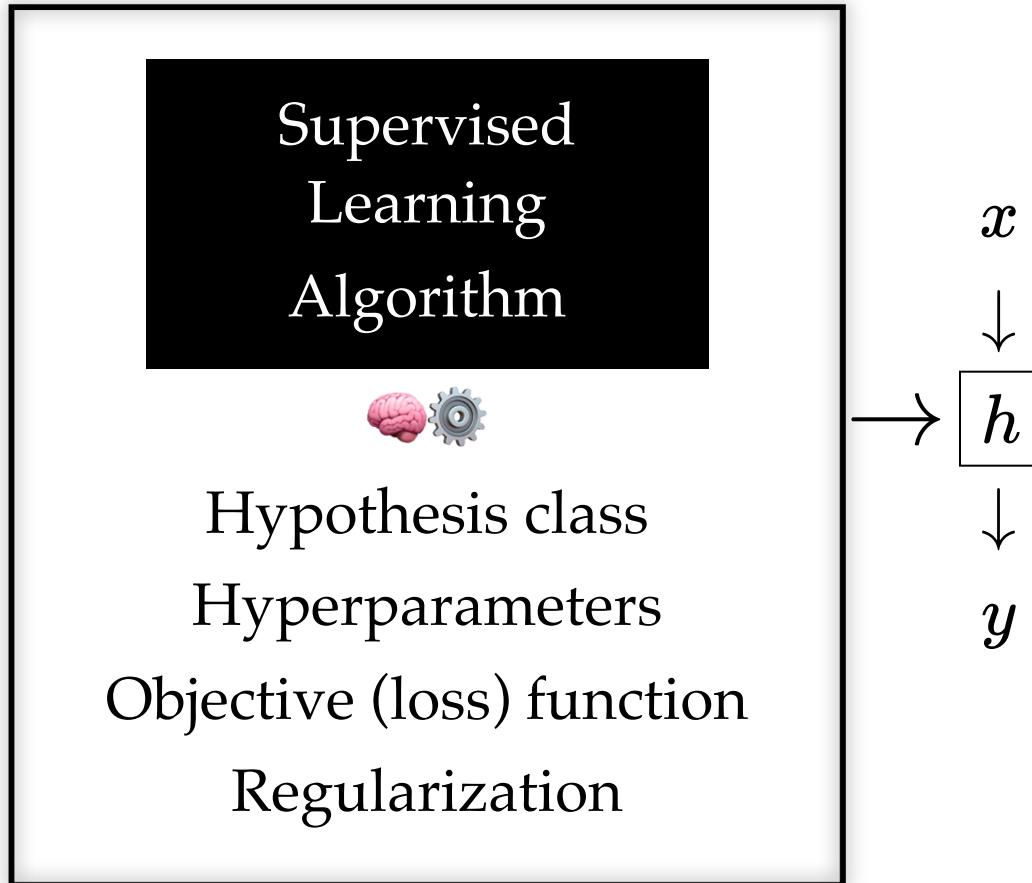
Week 1 - IntroML

- Terminologies
 - Training, validation, testing
 - Identifying overfitting and underfitting
- Concrete processes
 - Learning algorithm
 - Validation and Cross-validation
 - Concept of hyperparameter

"Learn" a model



$\mathcal{D}_{\text{train}} \rightarrow$



x



h



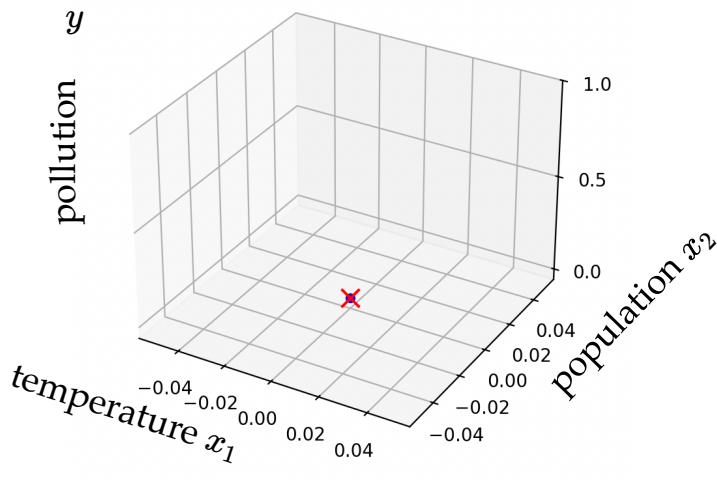
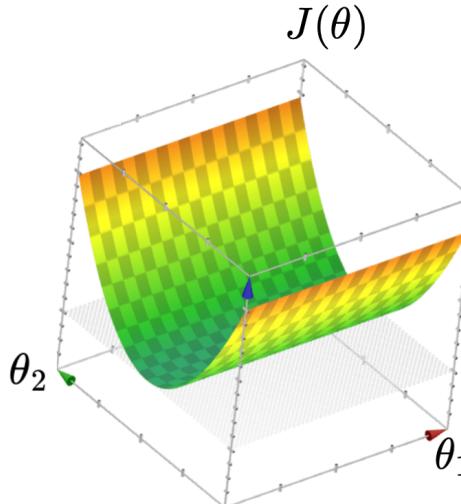
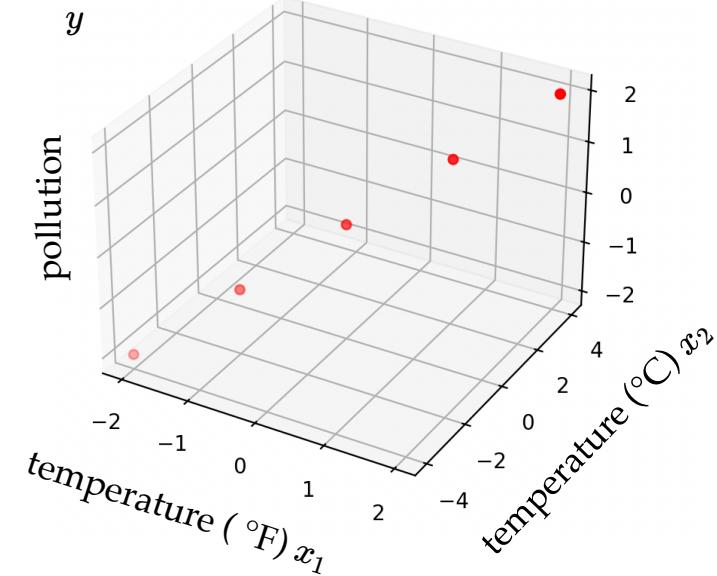
y



"Use" a model

Week 2 - Regression

- Problem Setup
- Analytical solution formula $\theta^* = (X^\top X)^{-1} X^\top Y$ (and what's X)
- When $X^\top X$ not invertible (optimal solutions still exist; just not via the "formula")
 - Practically (two scenarios)
 - Visually (obj fun no longer of "bowl" shape, instead has "half-pipe" shape)
 - Mathematically (loss of solution uniqueness)
- Regularization
 - Motivation, how to, when to

Case	Example	Objective Function Looks Like	Optimal Parameters
2a. less data than features			infinitely many optimal parameters (that define optimal hyperplanes)
2b. linearly dependent features			

Week 3 - Gradient Descent

- The gradient vector (both analytically and conceptually)
- The gradient-descent algorithm and the key update formula
- (Convex + small-enough step-size + gradient descent + global min exists + run long enough) guarantee convergence to a global min
 - What happens when any of these conditions is violated
- How does the stochastic variant differ (Set up, run-time behavior, and conclusion)

For $f : \mathbb{R}^m \rightarrow \mathbb{R}$, its *gradient* $\nabla f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is defined at the point $p = (x_1, \dots, x_m)$ as:

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_m}(p) \end{bmatrix}$$

1. The gradient generalizes the concept of a derivative to multiple dimensions.
2. By construction, the gradient's dimensionality always matches the function input.

Sometimes the gradient is undefined or ill-behaved, but today it is well-behaved unless stated otherwise.

3. The gradient can be symbolic or numerical.

example: $f(x, y, z) = x^2 + y^3 + z$

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_m}(p) \end{bmatrix}$$

its *symbolic* gradient:

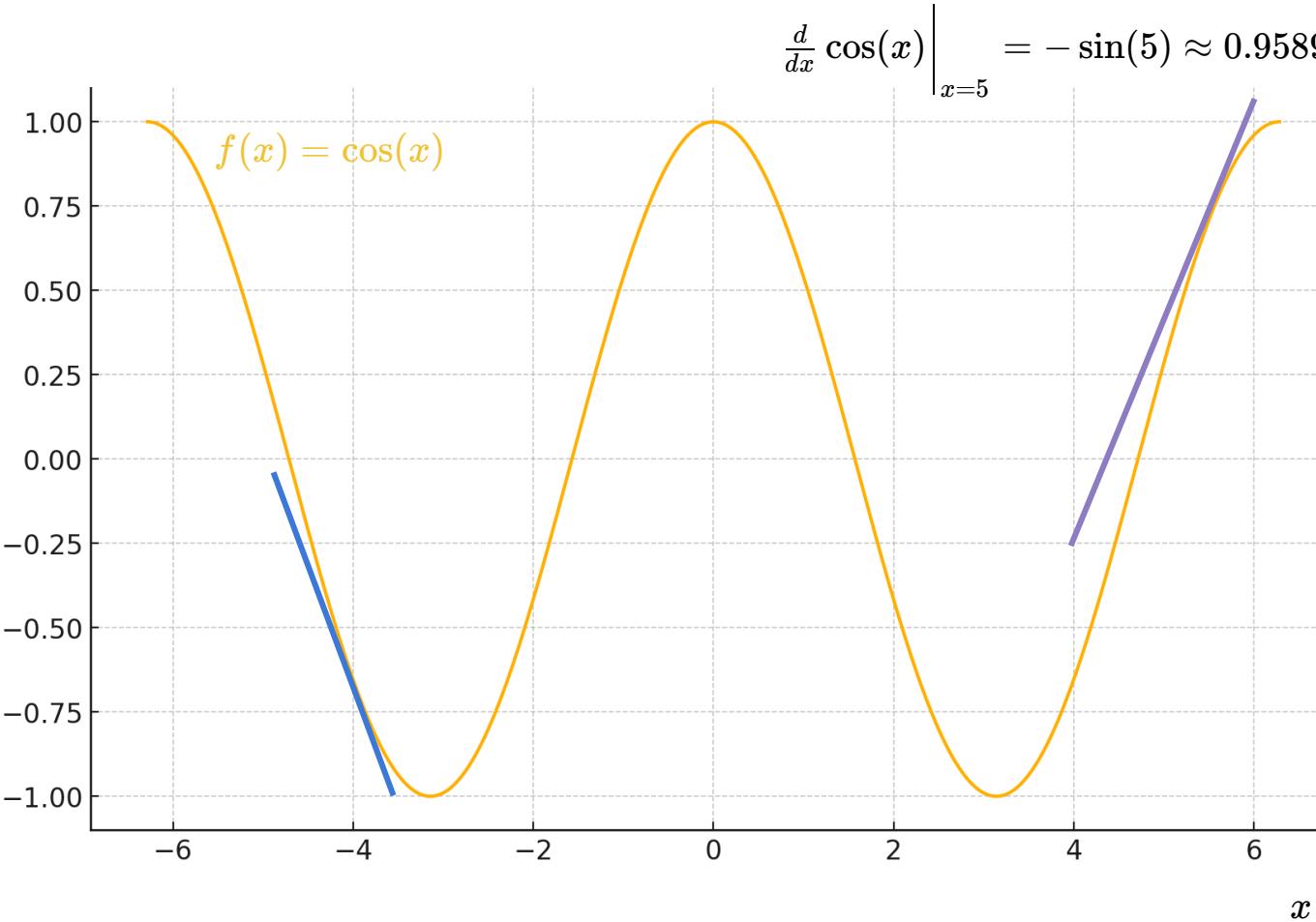
$$\nabla f(x, y, z) = \begin{bmatrix} 2x \\ 3y^2 \\ 1 \end{bmatrix}$$

evaluating the symbolic gradient at a point gives a *numerical* gradient:

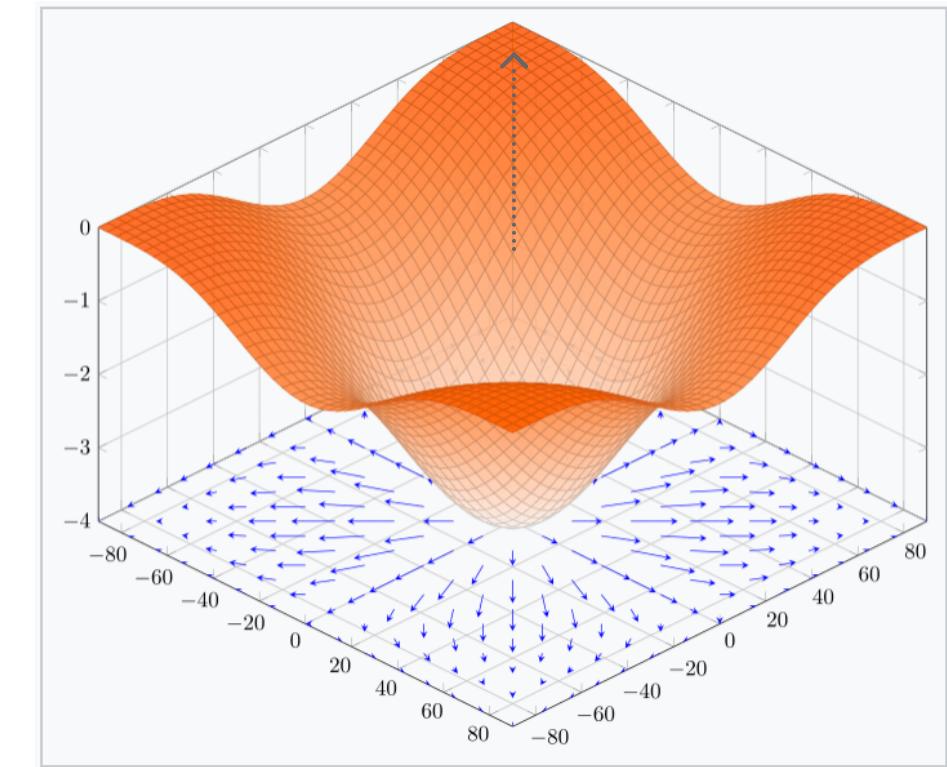
$$\nabla f(3, 2, 1) = \nabla f(x, y, z) \Big|_{(x,y,z)=(3,2,1)} = \begin{bmatrix} 6 \\ 12 \\ 1 \end{bmatrix}$$

just like a derivative can be a function or a number¹¹.

4. The gradient points in the direction of the (steepest) *increase* in the function value.



$$\frac{d}{dx} \cos(x) \Big|_{x=-4} = -\sin(-4) \approx -0.7568$$



The gradient of the function $f(x,y) = -(\cos^2 x + \cos^2 y)^2$ □

	initial guess of parameters	learning rate, aka, step size	precision
hyperparameters			

1 Gradient-Descent (Θ_{init} , η , f , $\nabla_{\Theta}f$, ϵ)

2 Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

3 Initialize $t = 0$

4 **repeat**

5 $t = t + 1$

6 $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$

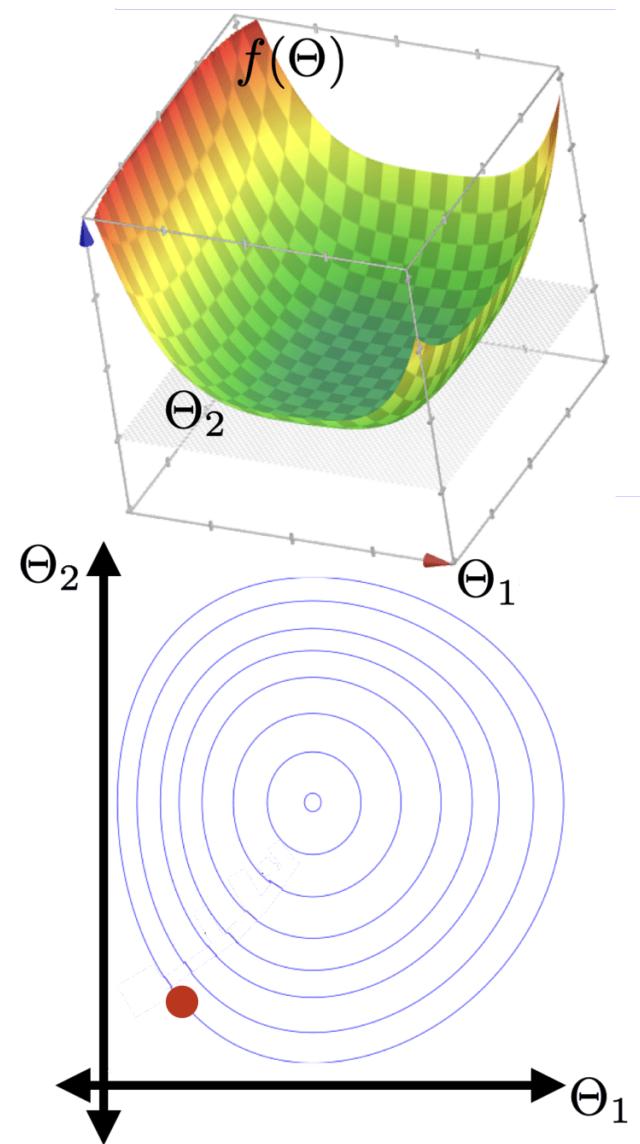
7 **until** $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$

8 **Return** $\Theta^{(t)}$

```

1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

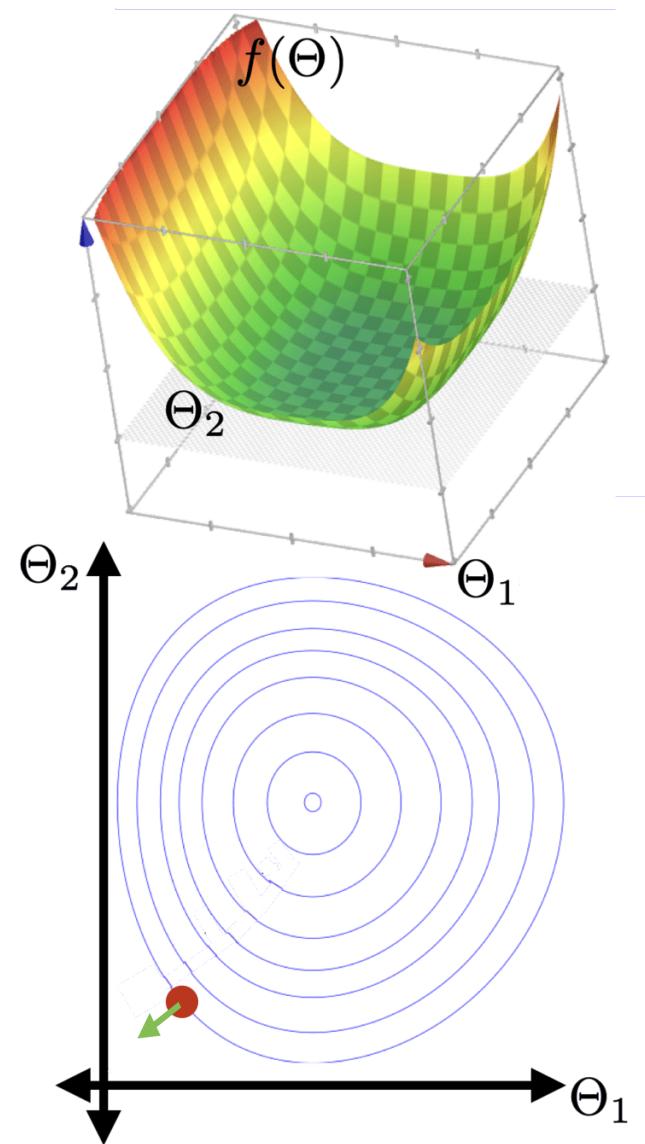
```



```

1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

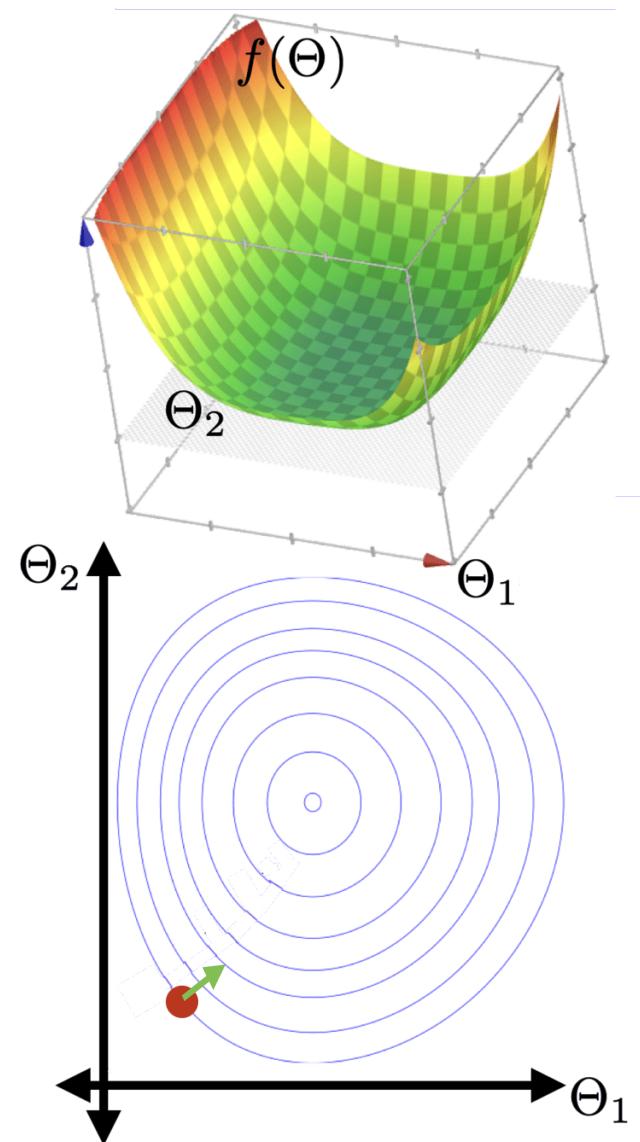
```



```

1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

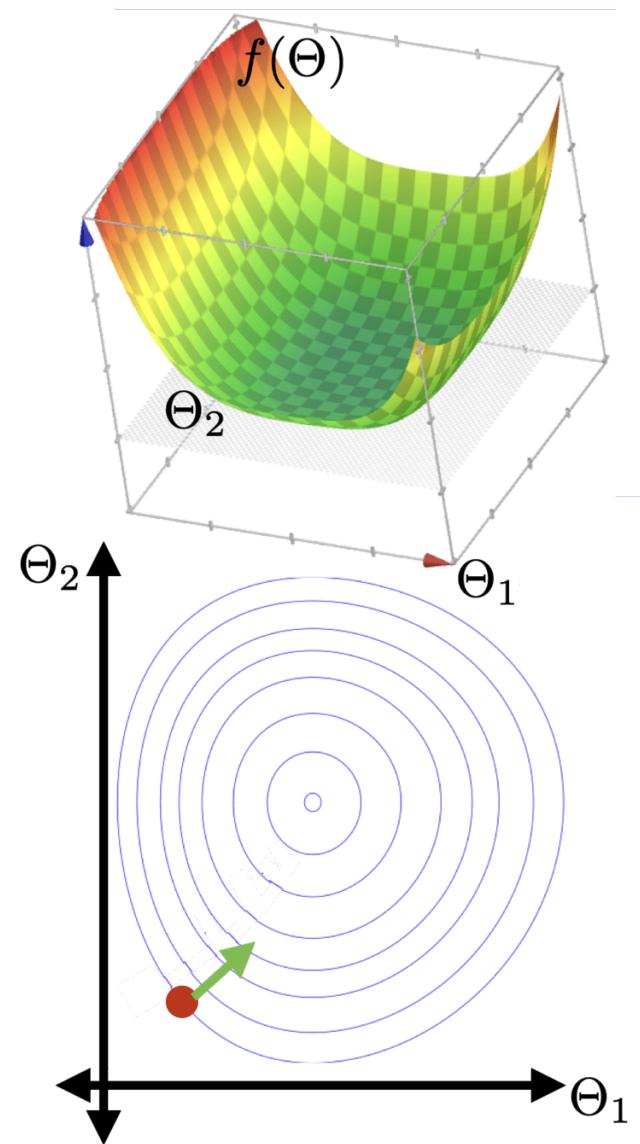
```



```

1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

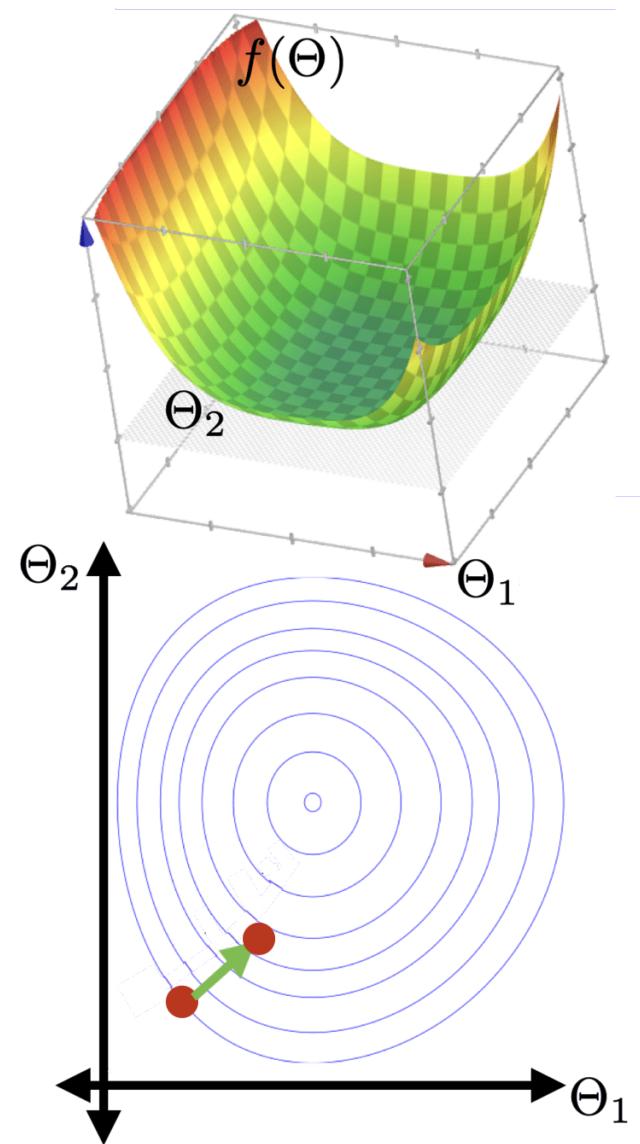
```



```

1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

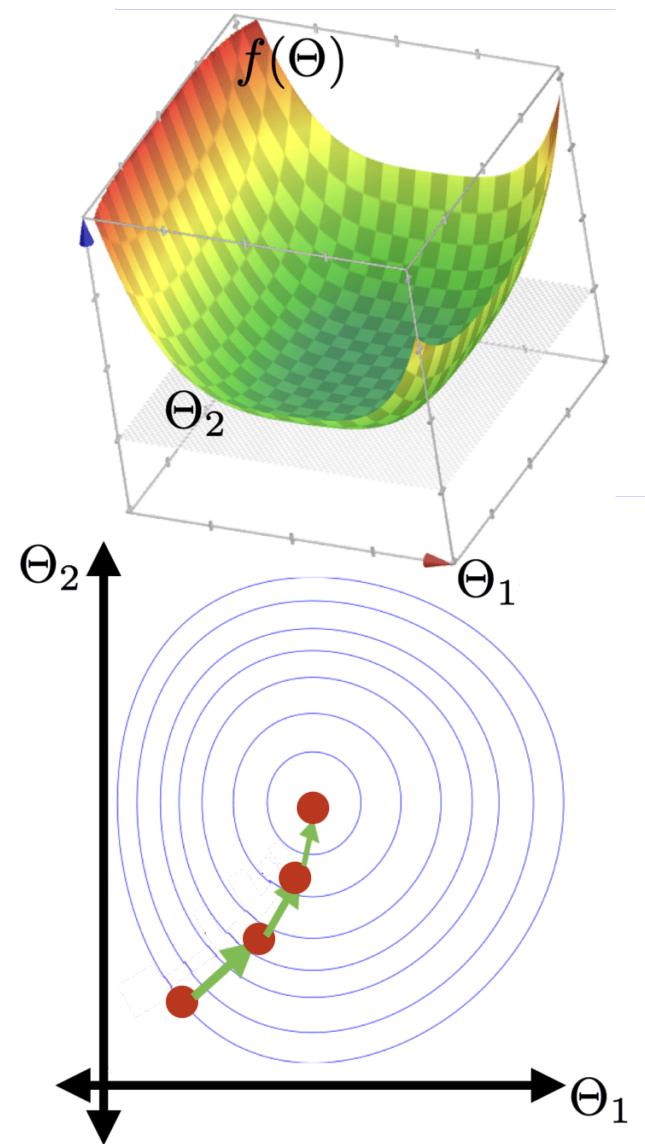
```



```

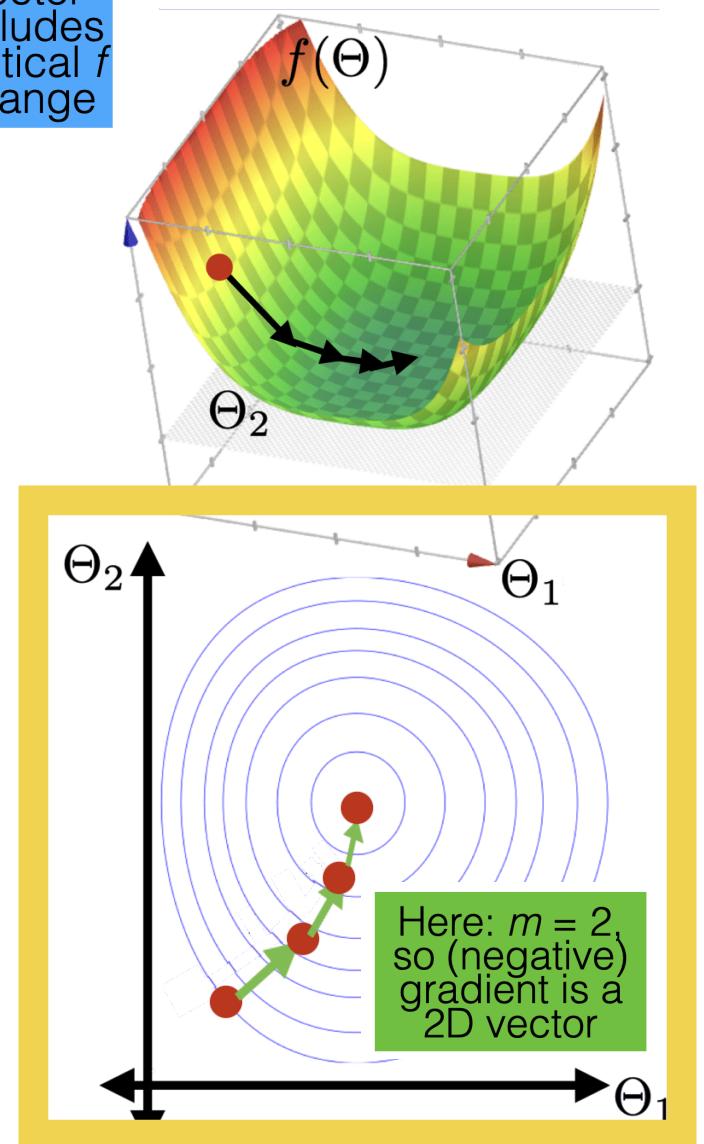
1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$ 
3 Initialize  $t = 0$ 
4 repeat
5      $t = t + 1$ 
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$ 
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$ 
8 Return  $\Theta^{(t)}$ 

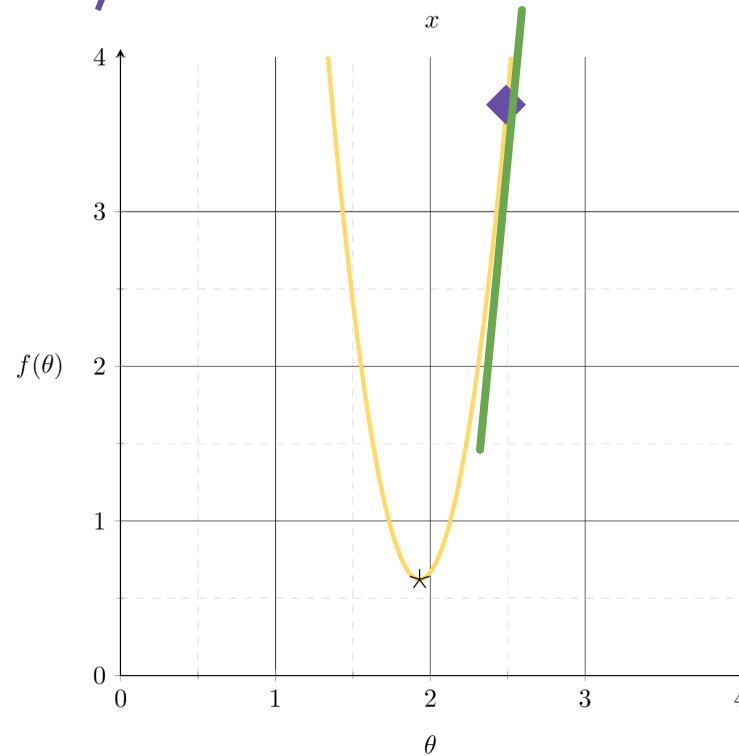
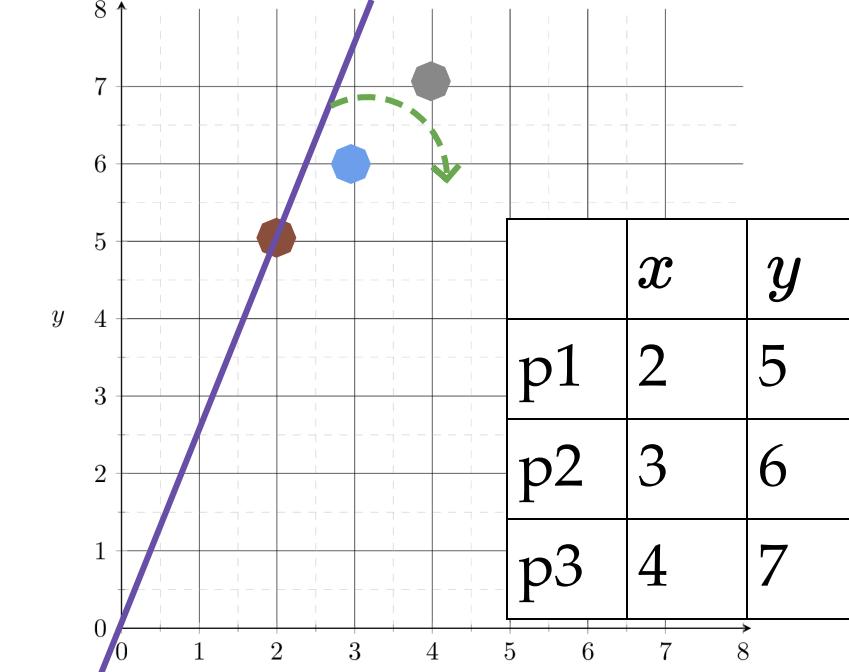
```



3D
vector
includes
vertical f
change

```
1 Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )  
2 Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$   
3 Initialize  $t = 0$   
4 repeat  
5      $t = t + 1$   
6      $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$   
7     until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$   
8 Return  $\Theta^{(t)}$ 
```





- the MSE of a linear hypothesis:

$$f(\theta) = \frac{1}{3} [(2\theta - 5)^2 + (3\theta - 6)^2 + (4\theta - 7)^2]$$

- and its gradient w.r.t. θ :

$$\nabla_{\theta} f = \frac{2}{3} [2(2\theta - 5) + 3(3\theta - 6) + 4(4\theta - 7)]$$

Gradient of an ML objective

Using our example data set,

- the MSE of a linear hypothesis:

$$f(\theta) = \frac{1}{3} [(2\theta - 5)^2 + (3\theta - 6)^2 + (4\theta - 7)^2]$$

- and its gradient w.r.t. θ :

$$\nabla_{\theta} f = \frac{2}{3}[2(2\theta - 5) + 3(3\theta - 6) + 4(4\theta - 7)]$$

Using any dataset,

- the MSE of a linear hypothesis:

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

- and its gradient w.r.t. θ :

$$\nabla f(\theta) = \frac{2}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)}) x^{(i)}$$

Gradient of an ML objective

- the MSE of a linear hypothesis:

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\theta^\top x^{(i)} - y^{(i)} \right)^2$$

- and its gradient w.r.t. θ :

$$\nabla f(\theta) = \frac{2}{n} \sum_{i=1}^n \left(\theta^\top x^{(i)} - y^{(i)} \right) x^{(i)}$$

In general,

- An ML objective function is a finite sum

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

- and its gradient w.r.t. θ :

$$\nabla f(\theta) = \nabla \left(\frac{1}{n} \sum_{i=1}^n f_i(\theta) \right) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta)$$



👉 (gradient of the sum) = (sum of the gradient)

Gradient of an ML objective

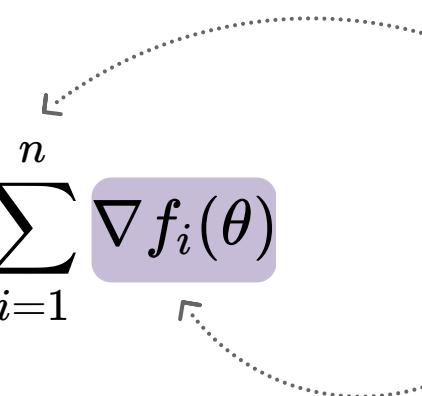
In general,

- An ML objective function is a finite sum

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

- and its gradient w.r.t. θ :

$$\nabla f(\theta) = \nabla \left(\frac{1}{n} \sum_{i=1}^n f_i(\theta) \right) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta)$$



need to add n of them

each of these $\nabla f_i(\theta) \in \mathbb{R}^d$

Costly!

Let's do stochastic gradient descent (on the board).

Stochastic gradient descent

```
Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )  
    Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$   
    Initialize  $t = 0$   
    repeat  
         $t = t + 1$   
         $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta}f(\Theta^{(t-1)})$   
    until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$   
Return  $\Theta^{(t)}$ 
```

Stochastic

```
Gradient-Descent (  $\Theta_{\text{init}}, \eta, f, \nabla_{\Theta}f, \epsilon$  )  
    Initialize  $\Theta^{(0)} = \Theta_{\text{init}}$   
    Initialize  $t = 0$   
    repeat  
         $t = t + 1$   
        randomly select  $i$  from  $\{1, \dots, n\}$   
         $\Theta^{(t)} = \Theta^{(t-1)} - \eta(t) \nabla_{\Theta}f_i(\Theta^{(t-1)})$   
    until  $|f(\Theta^{(t)}) - f(\Theta^{(t-1)})| < \epsilon$   
Return  $\Theta^{(t)}$ 
```

$$\nabla f(\Theta) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\Theta)$$

$$\approx \nabla f_i(\Theta)$$

for a randomly picked data point i

Week 4 - Classification

- (Binary) linear classifier (sign based)
- (Binary) Logistic classifiers (sigmoid, NLL loss)
- Linear separator (the equation form \Leftrightarrow the visual form with normal vector)
- Linear separability (its interplay with features)
- How to handle multiple classes
 - Softmax generalization (Softmax, cross-entropy)
 - Multiple sigmoids

<https://shenshen.mit.edu/demos/separator.html>

linear
binary classifier

linear *logistic*
binary classifier

features

$$x \in \mathbb{R}^d$$

parameters

$$\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$$

linear
combination

$$\theta^T x + \theta_0 = z$$

predict

$$\begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{cases} 1 & \text{if } \sigma(z) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

	linear regressor	linear binary classifier	linear logistic binary classifier
features		$x \in \mathbb{R}^d$	
parameters		$\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$	
linear combo		$\theta^T x + \theta_0 = z$	
predict	z	$\begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$	$\begin{cases} 1 & \text{if } g = \sigma(z) > 0.5 \\ 0 & \text{otherwise} \end{cases}$
loss	$(g - y)^2$	$\begin{cases} 0 & \text{if } g = a \\ 1 & \text{otherwise} \end{cases}$	$-[y \log g + (1 - y) \log (1 - g)]$
optimize via	closed-form or gradient descent	NP-hard to learn	<ul style="list-style-type: none"> gradient descent only need regularization to not overfit

	linear logistic <i>binary</i> classifier	one-out-of- K classifier
training data	$x \in \mathbb{R}^d, y \in \{0, 1\}$	$x \in \mathbb{R}^d, y : K$ -dimensional one-hot
parameters	$\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$	$\theta \in \mathbb{R}^{d \times K}, \theta_0 \in \mathbb{R}^K$
linear combo	$\theta^T x + \theta_0 = z \in \mathbb{R}$	$\theta^T x + \theta_0 = z \in \mathbb{R}^K$
predict	$\sigma(z) = \frac{\exp(z)}{\exp(0) + \exp(z)}$ positive if $\sigma(z) > 0.5$	$\text{softmax}(z) = \begin{bmatrix} \exp(z_1) / \sum_i \exp(z_i) \\ \vdots \\ \exp(z_K) / \sum_i \exp(z_i) \end{bmatrix}$ category corresponding to the largest entry in $\text{softmax}(z)$

Negative log-likelihood K – classes loss (aka, cross-entropy)

g : softmax output

g_k : probability or confidence in class k

y : one-hot encoding label

$$\mathcal{L}_{\text{nllm}}(g, y) = - \sum_{k=1}^K y_k \cdot \log(g_k)$$

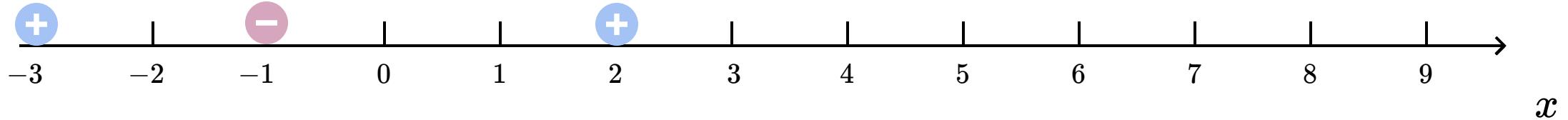
y_k : either 0 or 1

- Generalizes negative log likelihood loss $\mathcal{L}_{\text{nll}}(g, y) = - [y \log g + (1 - y) \log (1 - g)]$
- Appears as summing K terms, but
- for a given data point, only the term corresponding to its true class label matters.

Week 5 - Features, Neural Networks I

- Feature transformations
 - Apply a fixed feature transformation
 - Hand-design feature transformation (e.g. towards getting linear separability)
 - Interplay between the number of features, the quality of features, and the quality of learning algorithms
- Forward-pass (for evaluation)

Not linearly separable in x space

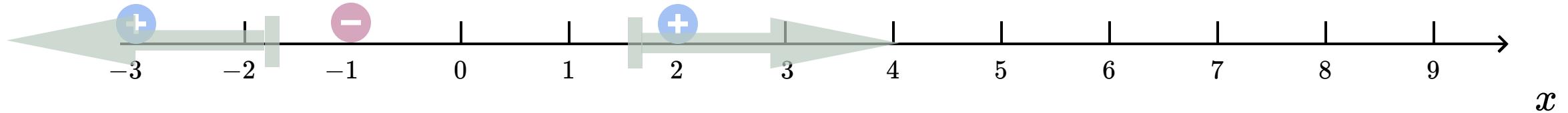


\Downarrow transform via $\phi(x) = x^2$

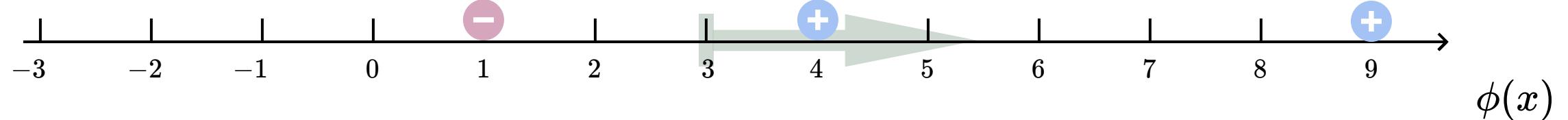


Linearly separable in $\phi(x) = x^2$ space

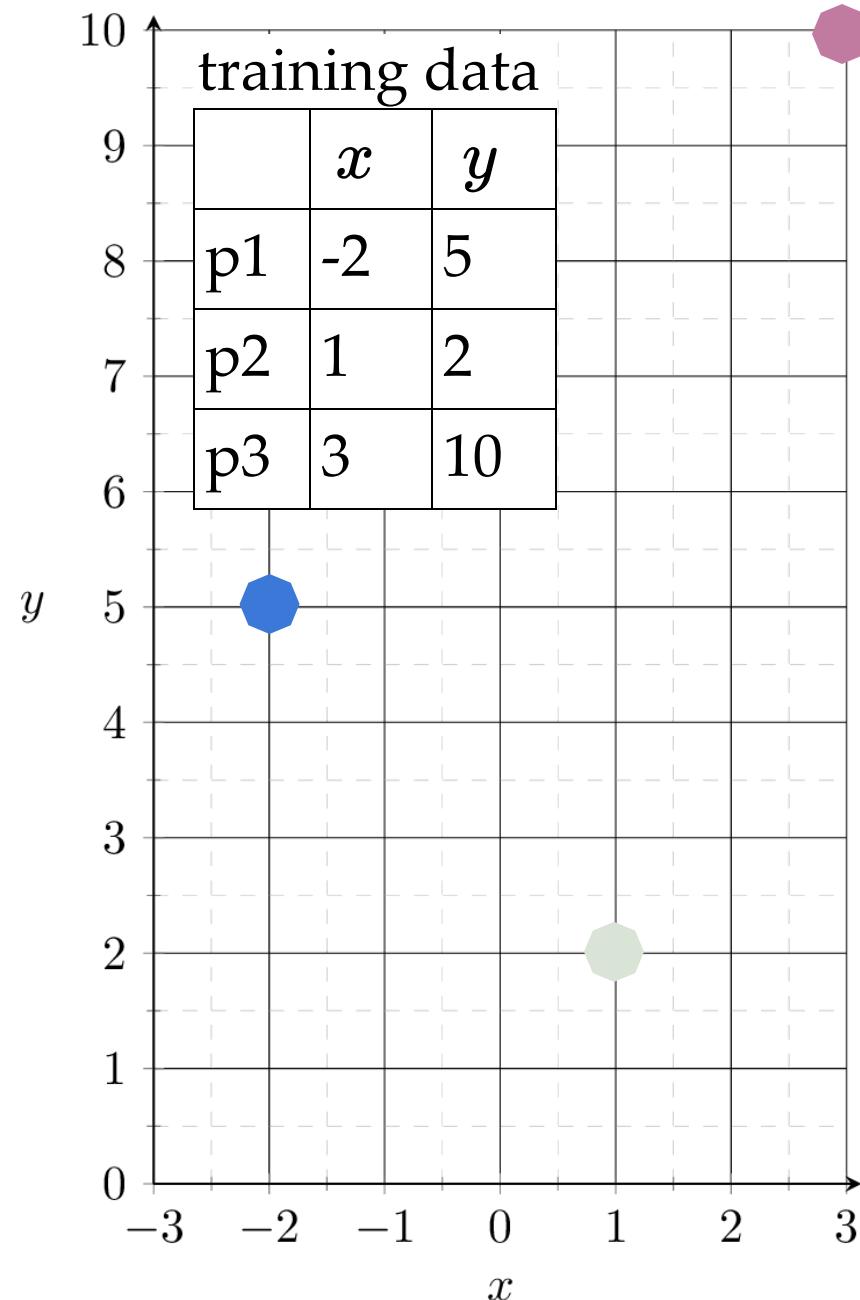
Non-linearly separated in x space, e.g. predict positive if $x^2 \geq 3$



\Downarrow transform via $\phi(x) = x^2$

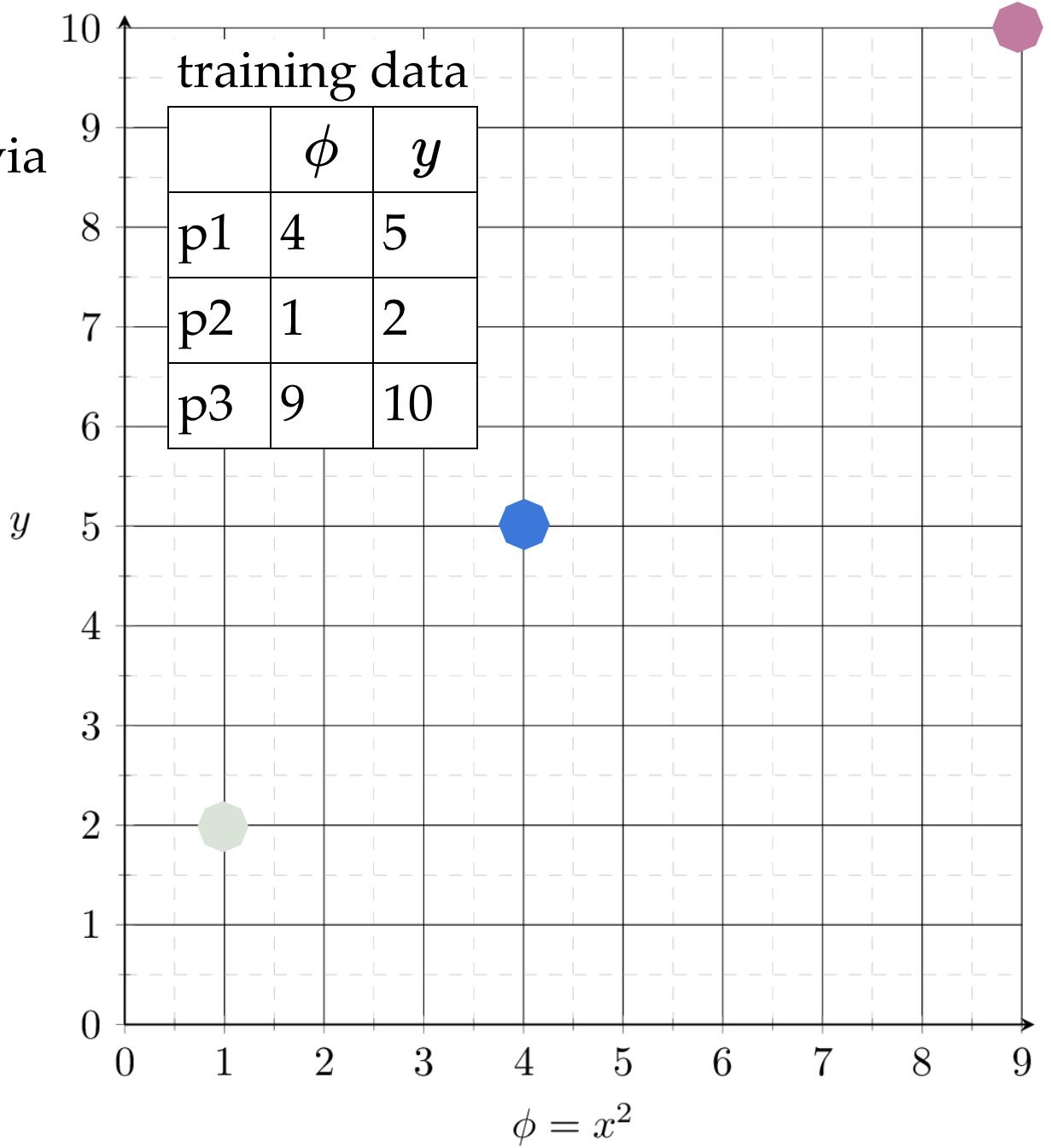


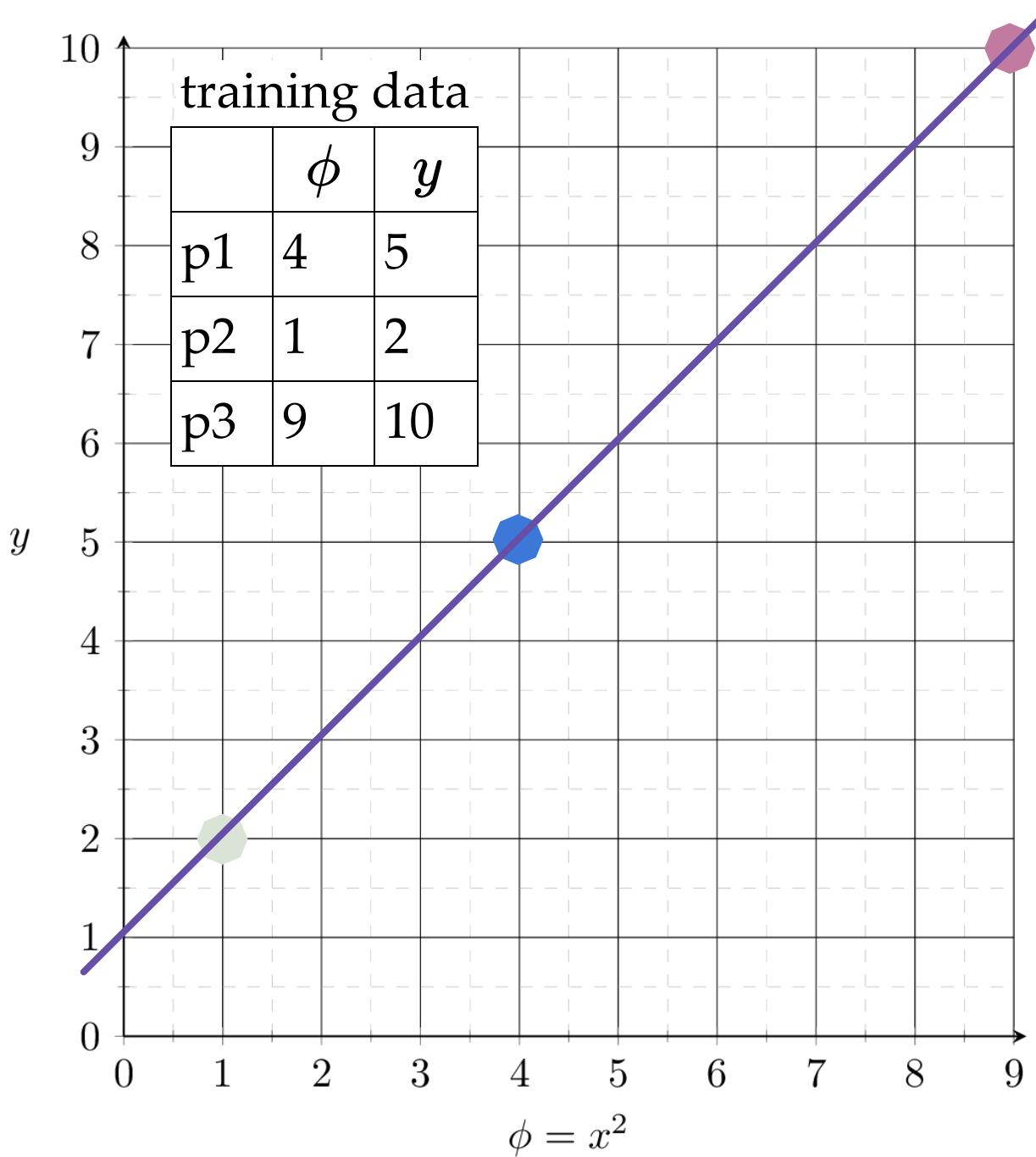
Linearly separable in $\phi(x) = x^2$ space, e.g. predict positive if $\phi \geq 3$



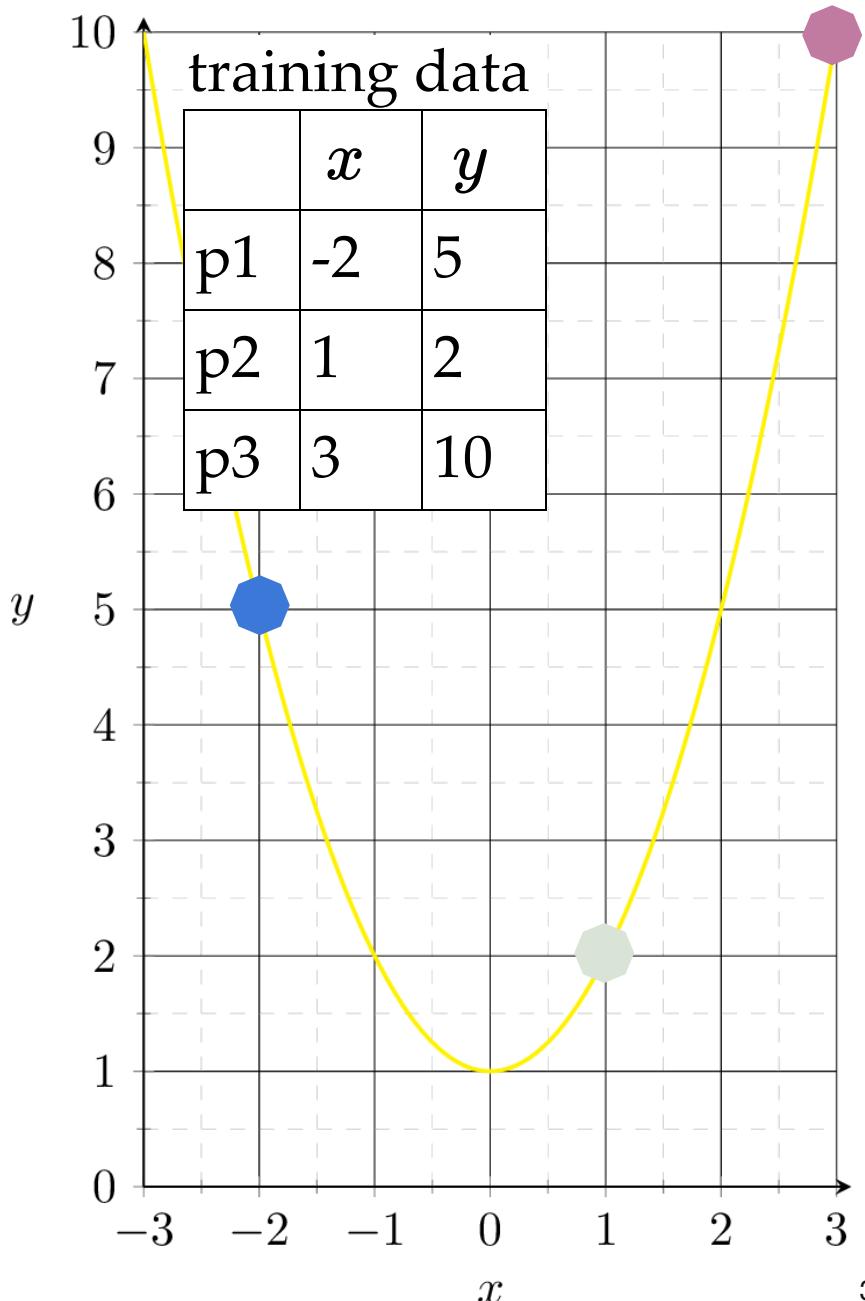
transform via
 $\phi(x) = x^2$

⇒

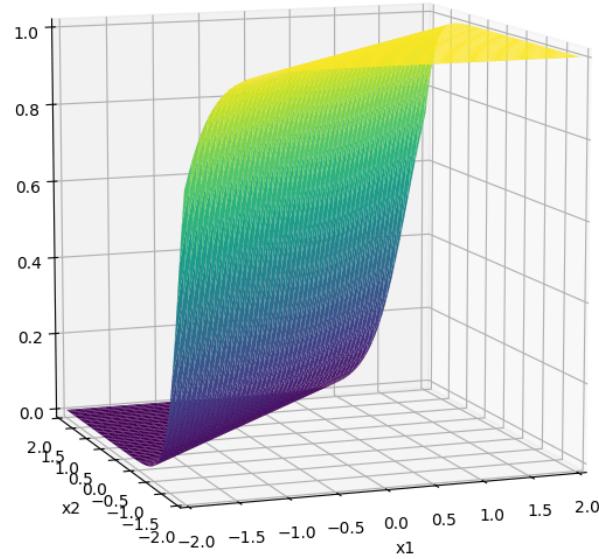




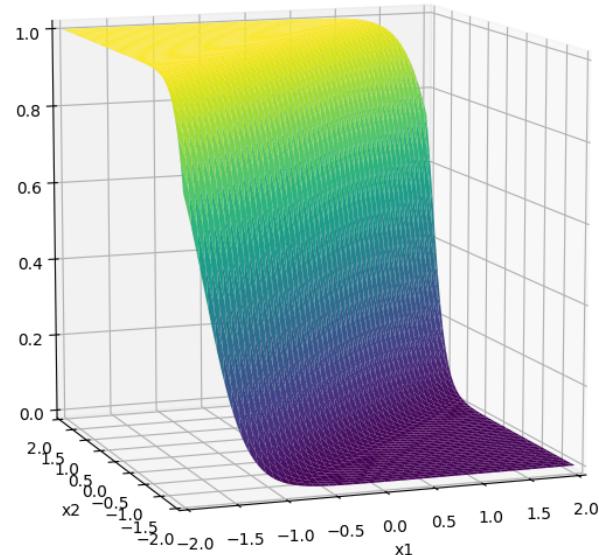
$$\begin{aligned}
 y &= \phi + 1 \\
 &= x^2 + 1 \\
 \Rightarrow
 \end{aligned}$$



$$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$$

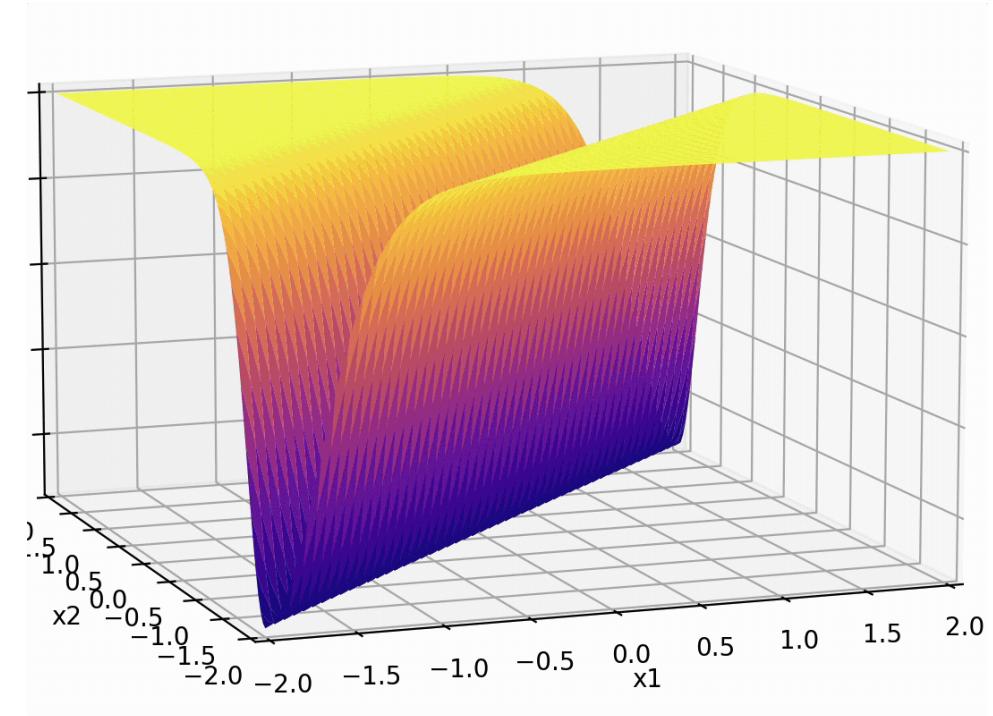


$$\sigma_2 = \sigma(-5x_1 + 5x_2 + 1)$$



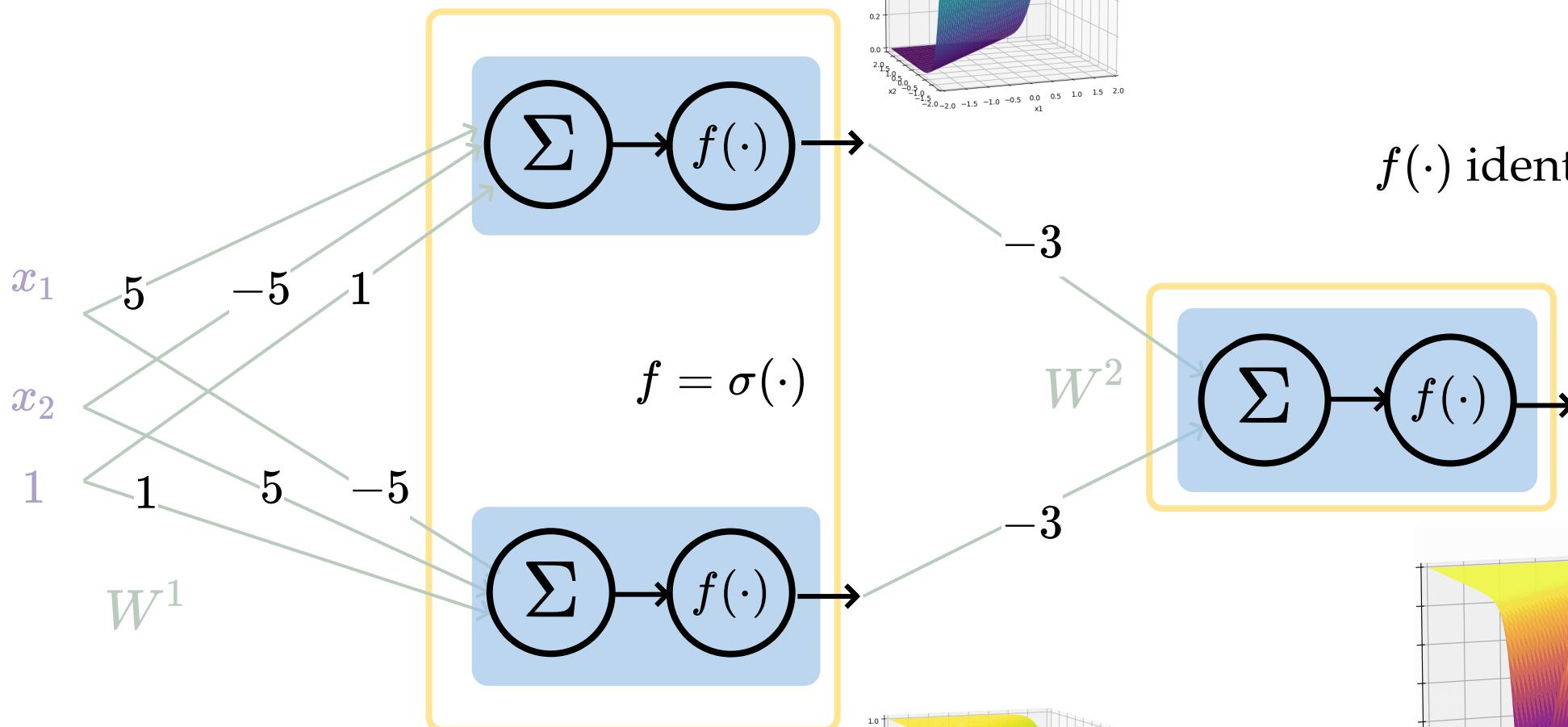
recall this example

$$-3(\sigma_1 + \sigma_2)$$



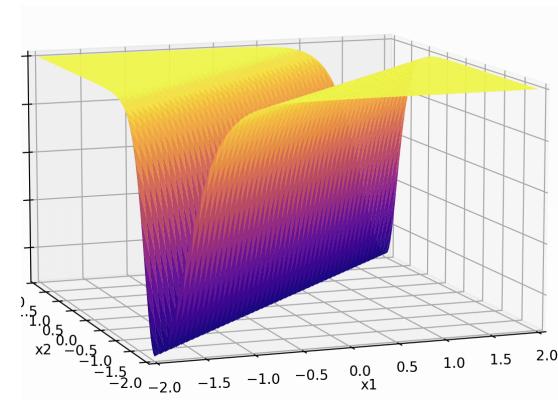
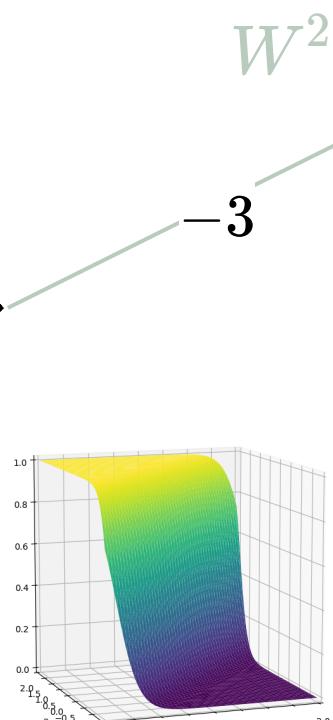
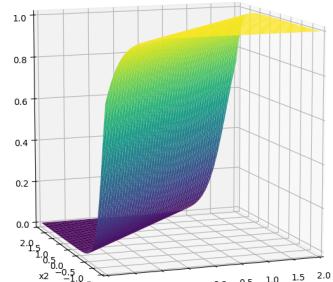
$$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$$

Recall



$$\sigma_2 = \sigma(-5x_1 + 5x_2 + 1)$$

$f(\cdot)$ identity function



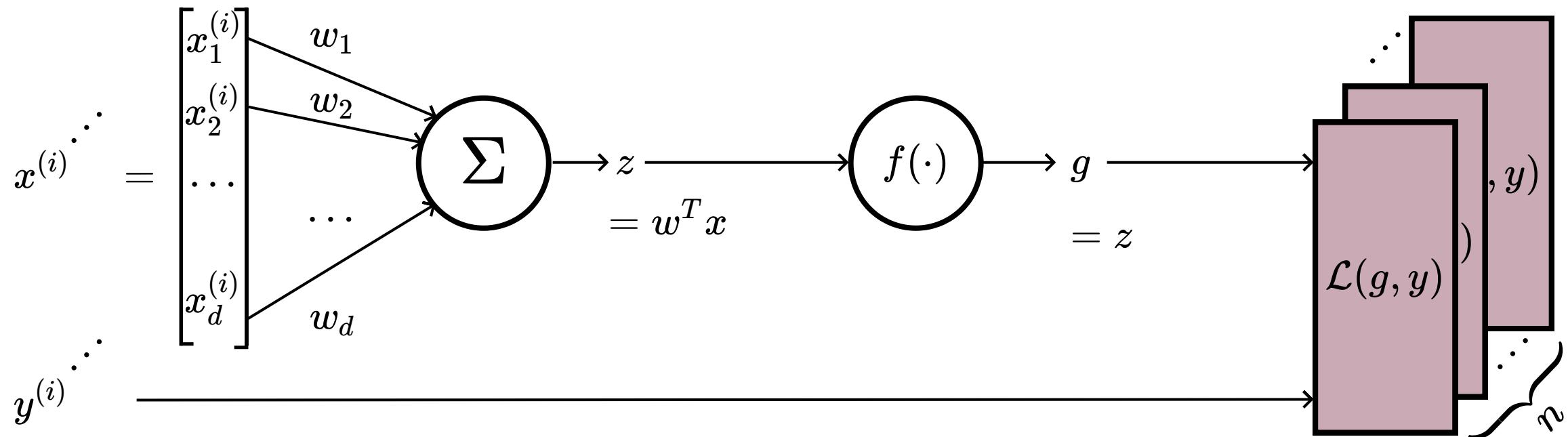
$$-3(\sigma_1 + \sigma_2)$$

<https://shenshen.mit.edu/demos/2layers.html>

Week 6 - Neural Networks

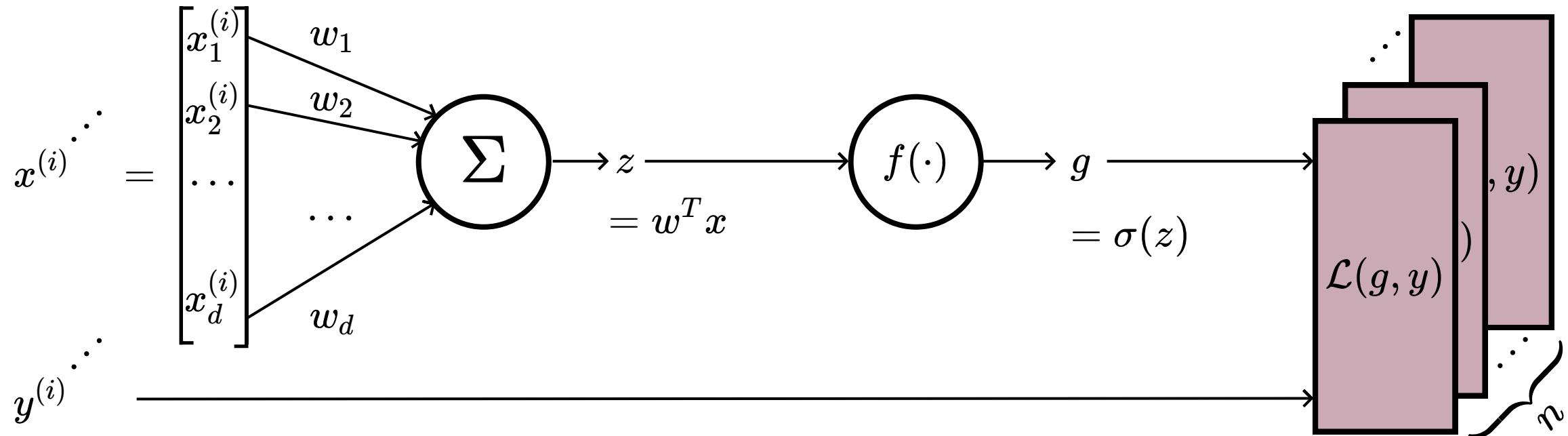
- Backward-pass (via backpropogation, for optimization)
- Source of expressiveness
- Output layer design
 - dimension, activation, loss
- Hand-designing weights
 - to match some given function form
 - achieve some goal (e.g. separate a given data set)

e.g. forward-pass of a linear regressor



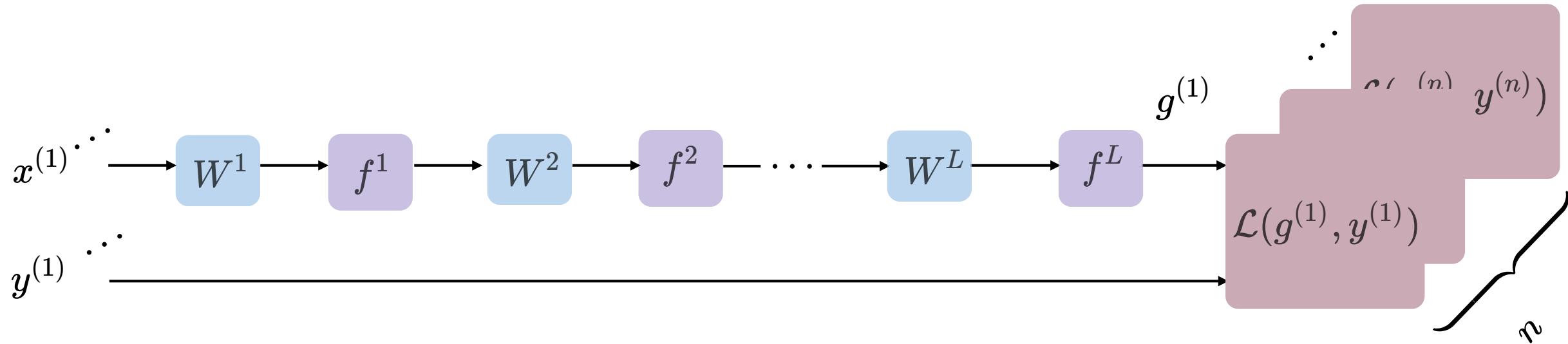
- Activation f is chosen as the identity function
- Evaluate the loss $\mathcal{L} = (g^{(i)} - y^{(i)})^2$
- Repeat for each data point, average the sum of n individual losses

e.g. forward-pass of a linear logistic classifier



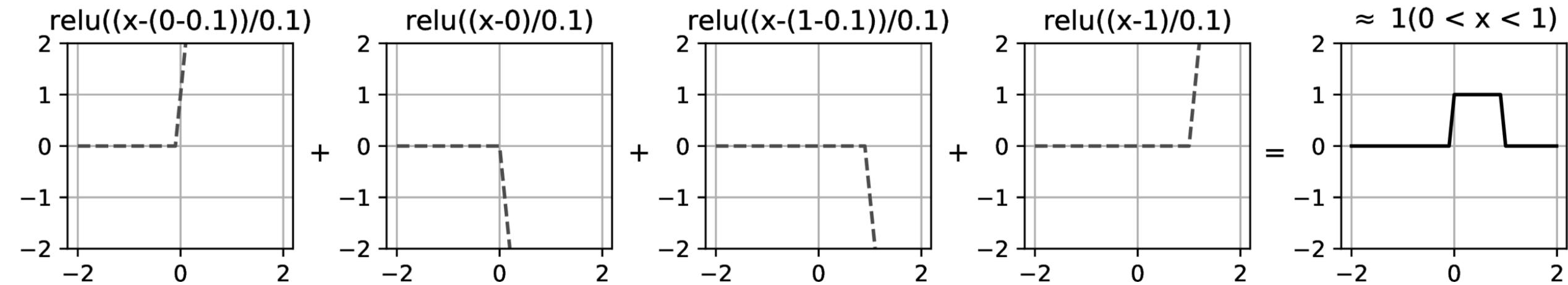
- Activation f is chosen as the sigmoid function
- Evaluate the loss $\mathcal{L} = -[y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log (1 - g^{(i)})]$
- Repeat for each data point, average the sum of n individual losses

Forward pass: evaluate, *given* the current parameters

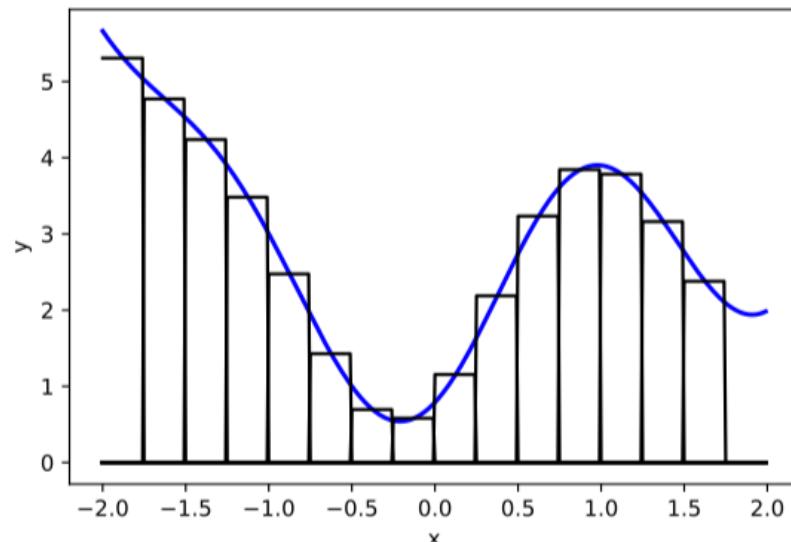


- the model outputs $g^{(i)} = f^L \left(\dots f^2 \left(f^1(\mathbf{x}^{(i)}; \mathbf{W}^1); \mathbf{W}^2 \right); \dots \mathbf{W}^L \right)$
- the loss incurred on the current data $\mathcal{L}(g^{(i)}, y^{(i)})$ linear combination
- the training error $J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(g^{(i)}, y^{(i)})$ (nonlinear) activation
loss function

Recall: compositions of ReLU(s) can be quite expressive



in fact, asymptotically, can approximate any function!

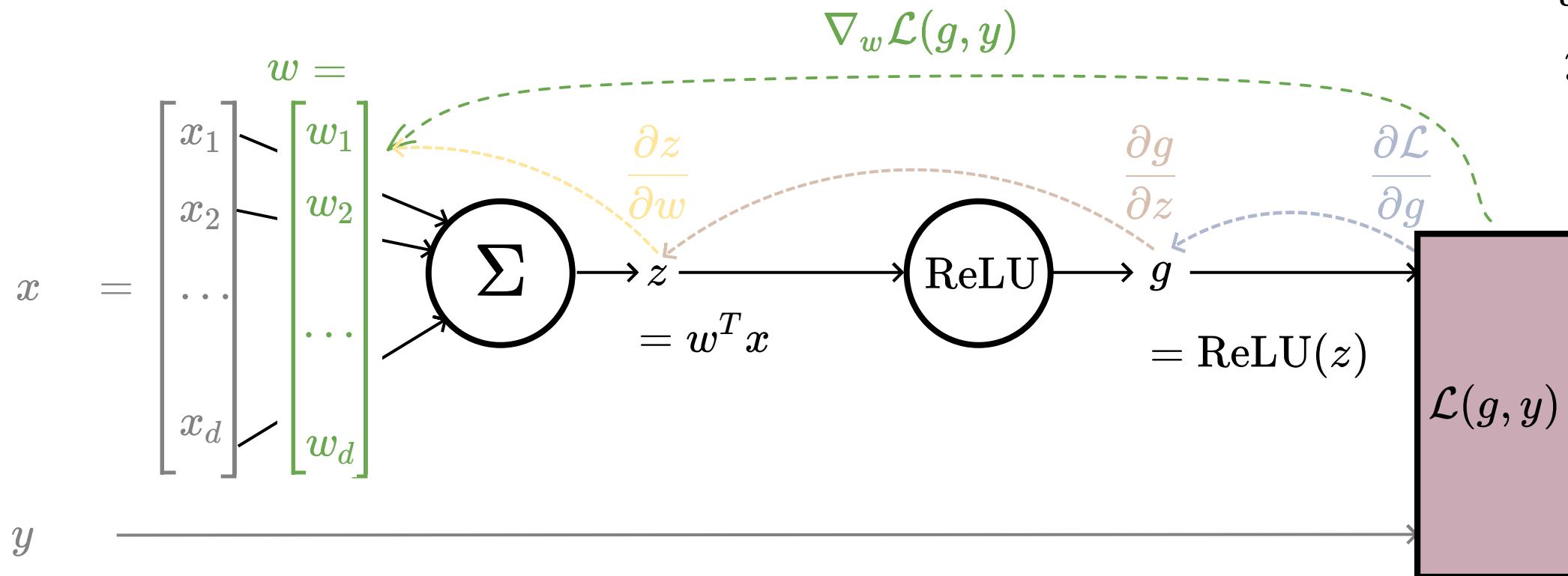


now, slightly more interesting activation:

$$x \in \mathbb{R}^d$$

$$w \in \mathbb{R}^d$$

$$y \in \mathbb{R}$$

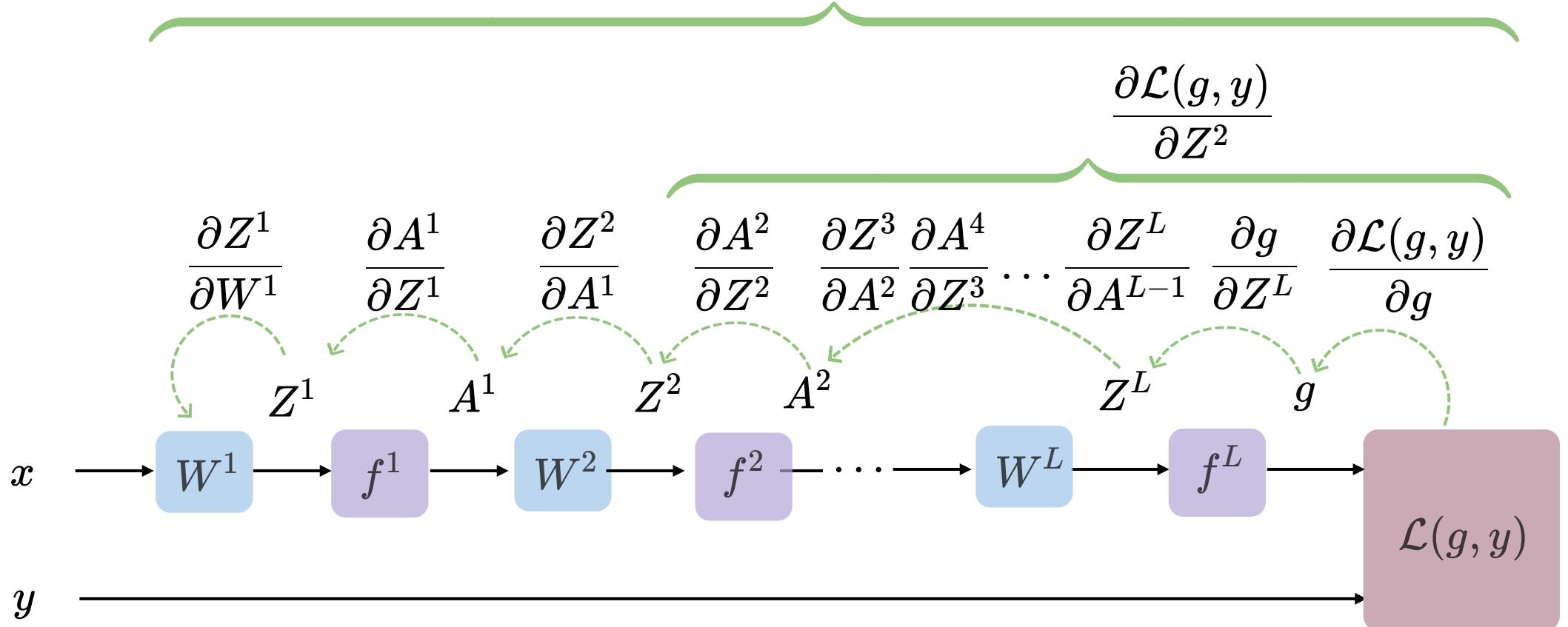


$$\nabla_w \mathcal{L}(g, y) = \frac{\partial \mathcal{L}(g, y)}{\partial w} = \frac{\partial[(g - y)^2]}{\partial w} = x \cdot \frac{\partial[\text{ReLU}(z)]}{\partial z} \cdot 2(g - y)$$

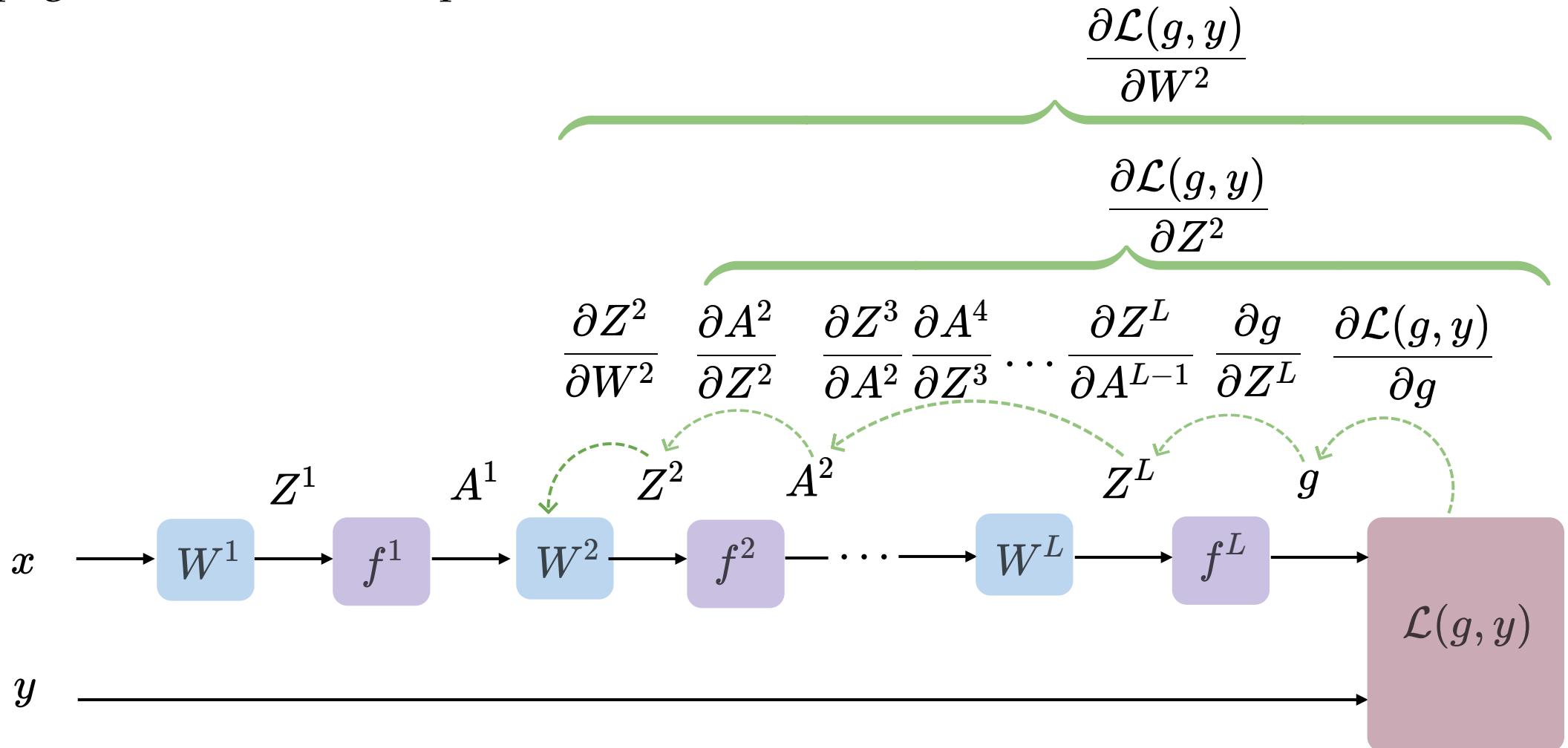
example on black-board 45

how to find $\frac{\partial \mathcal{L}(g, y)}{\partial W^1}$

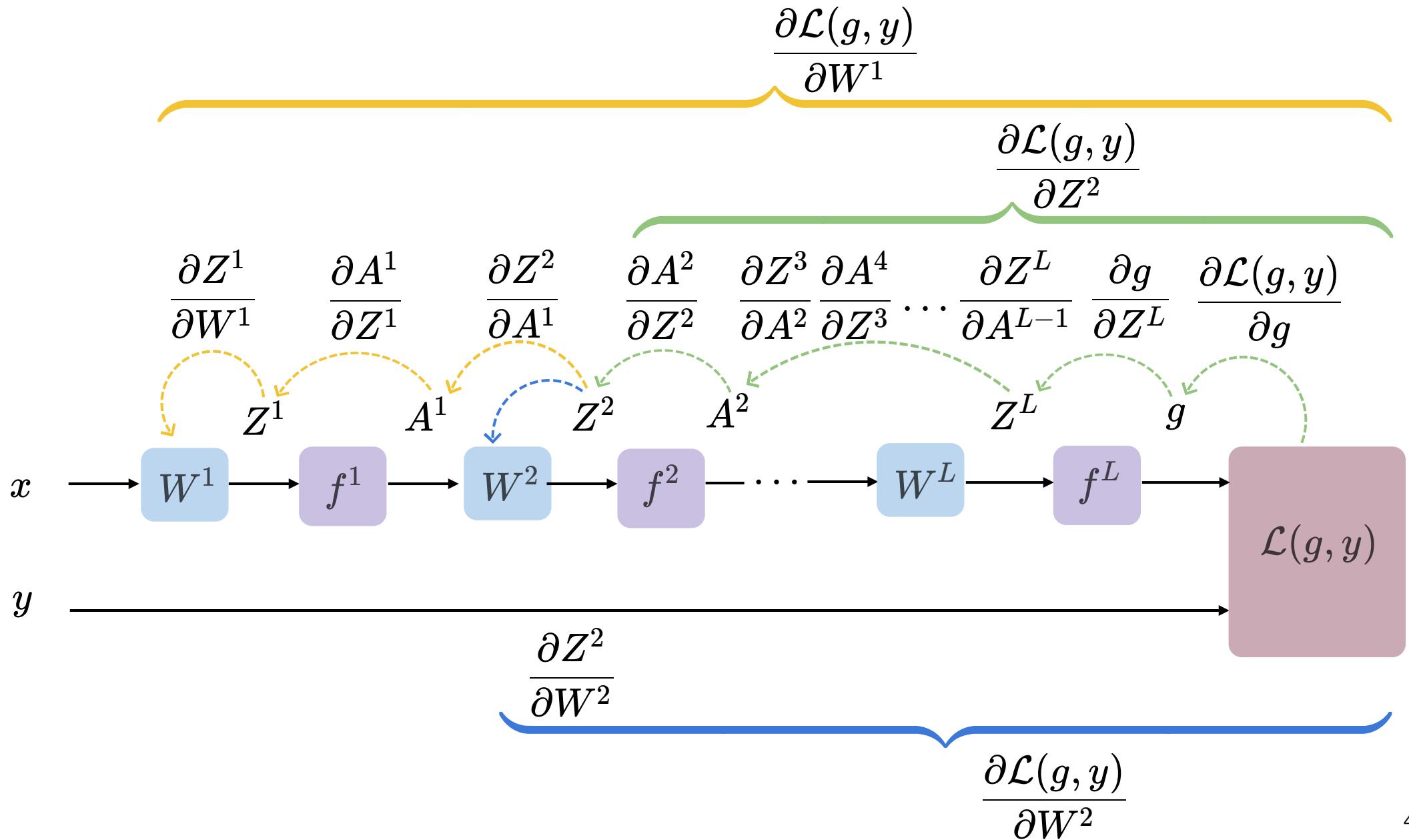
Now $\frac{\partial \mathcal{L}(g, y)}{\partial W^1}$



back propagation: reuse of computation



back propagation: reuse of computation

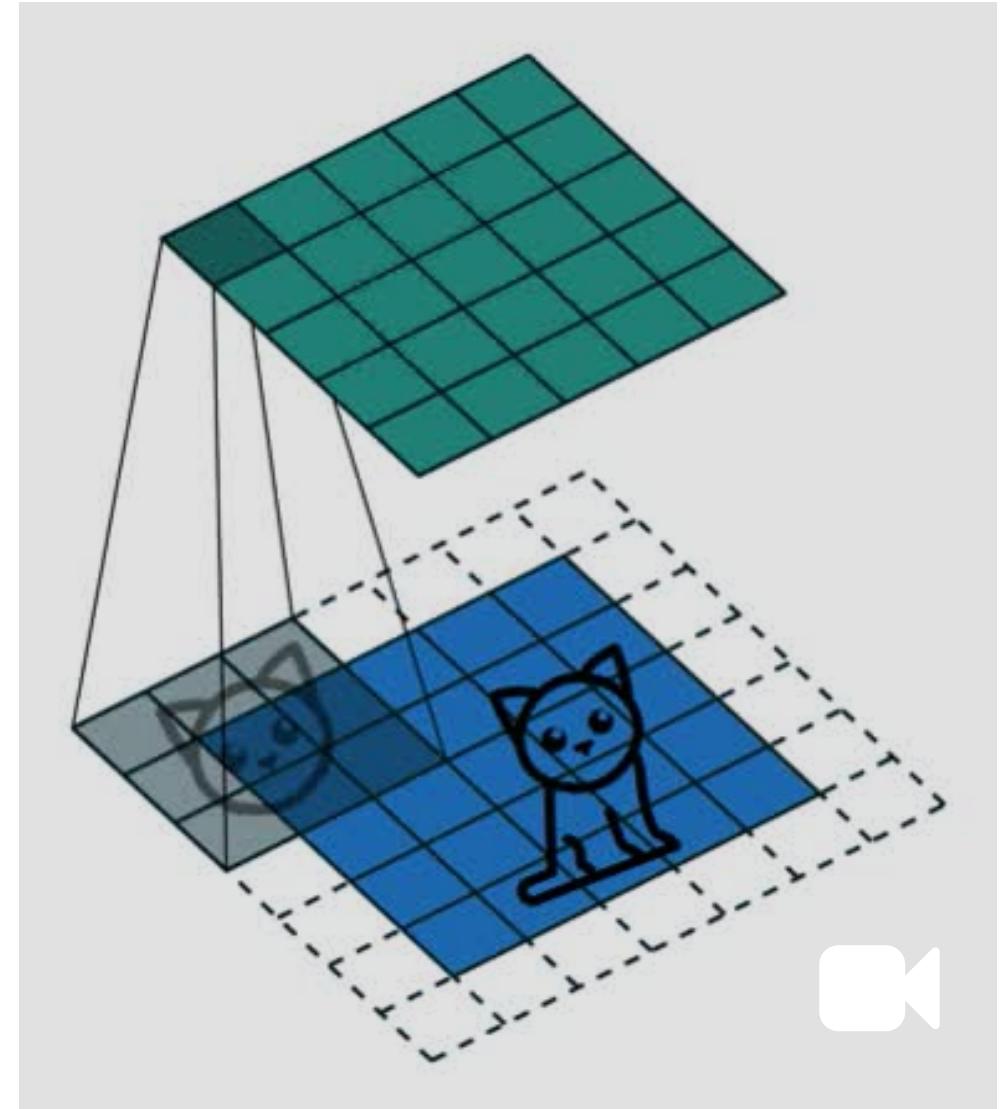


Week 7 - CNN

- Forward pass: convolution operation; max-pooling and the typical "pyramid" stack.
- Backward pass: back-propagation to learn filter weights/bias.
- The convolution/max-pooling operation
 - various hyper-parameters (filter size, padding size, stride) in spatial dimension;
 - the 3rd channel/depth dimension
 - reason about in/out shapes
- Conceptually: weight sharing, "pattern matching" template, independent and parallel processing.

convolution interpretation:

- Looking locally
- Parameter sharing
- Template matching
- Translational equivariance

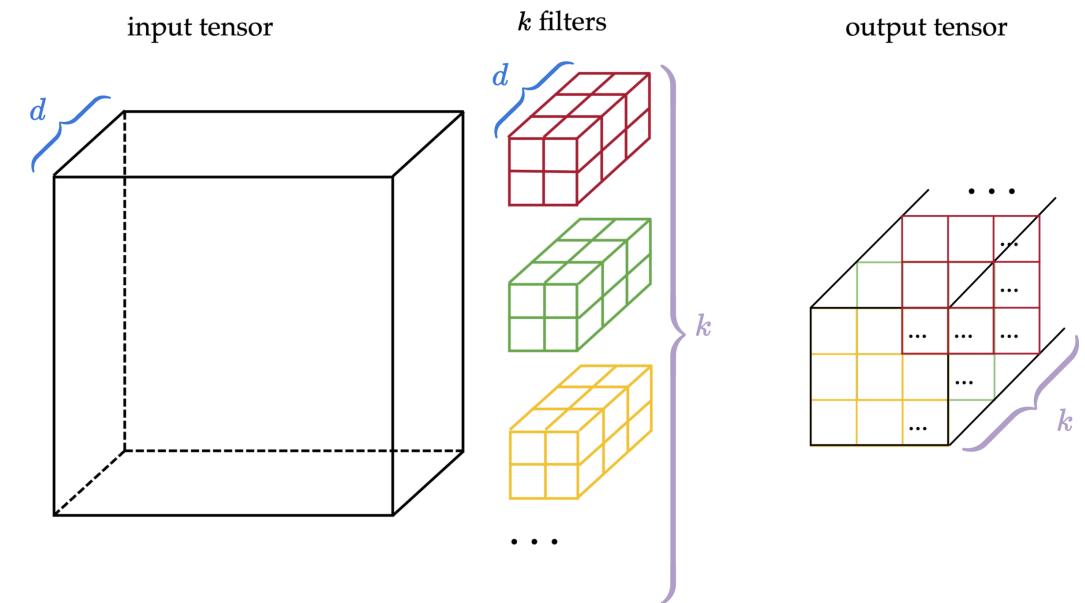
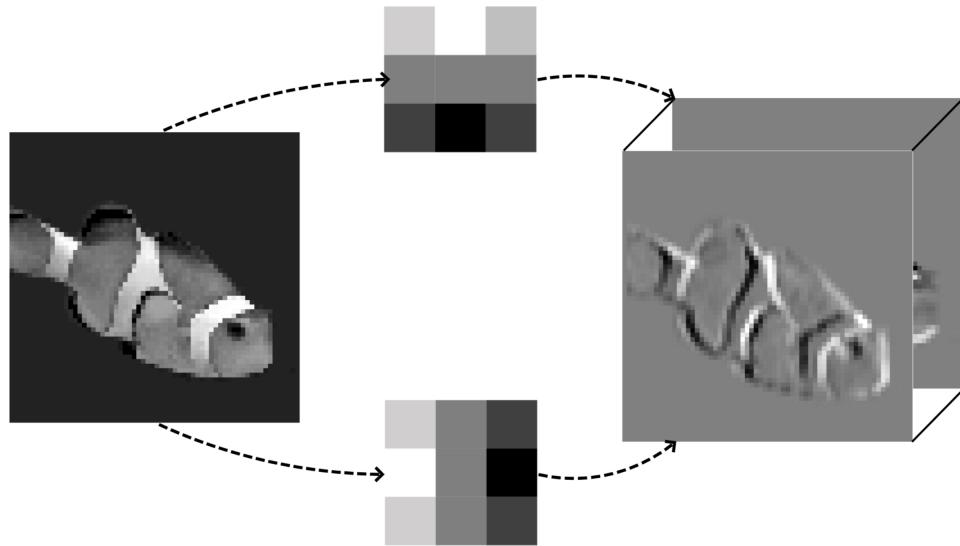


We'd encounter 3d tensor due to:

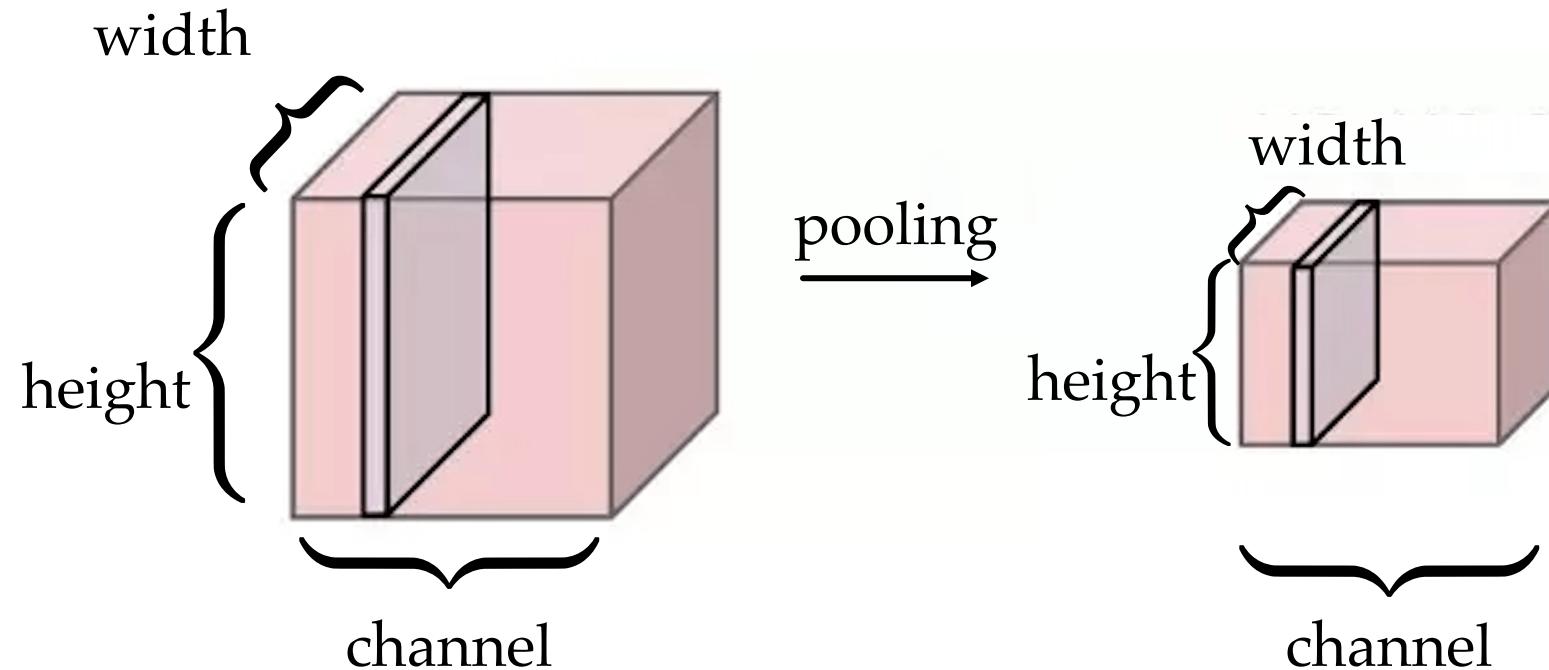
1. color input



2. the use of multiple filters -- in doing 2-dimensional convolution



Pooling across *spatial* locations achieves invariance w.r.t. small translations:



so the *channel* dimension remains *unchanged* after pooling.

Week 8 - Representation Learning

- Unsupervised learning setup
- Auto-encoder: compression and reconstruction
- The concept of embedding space:
 - Similarity comparison
 - vector arithmetic

$$g = \text{softmax}(z_2)$$

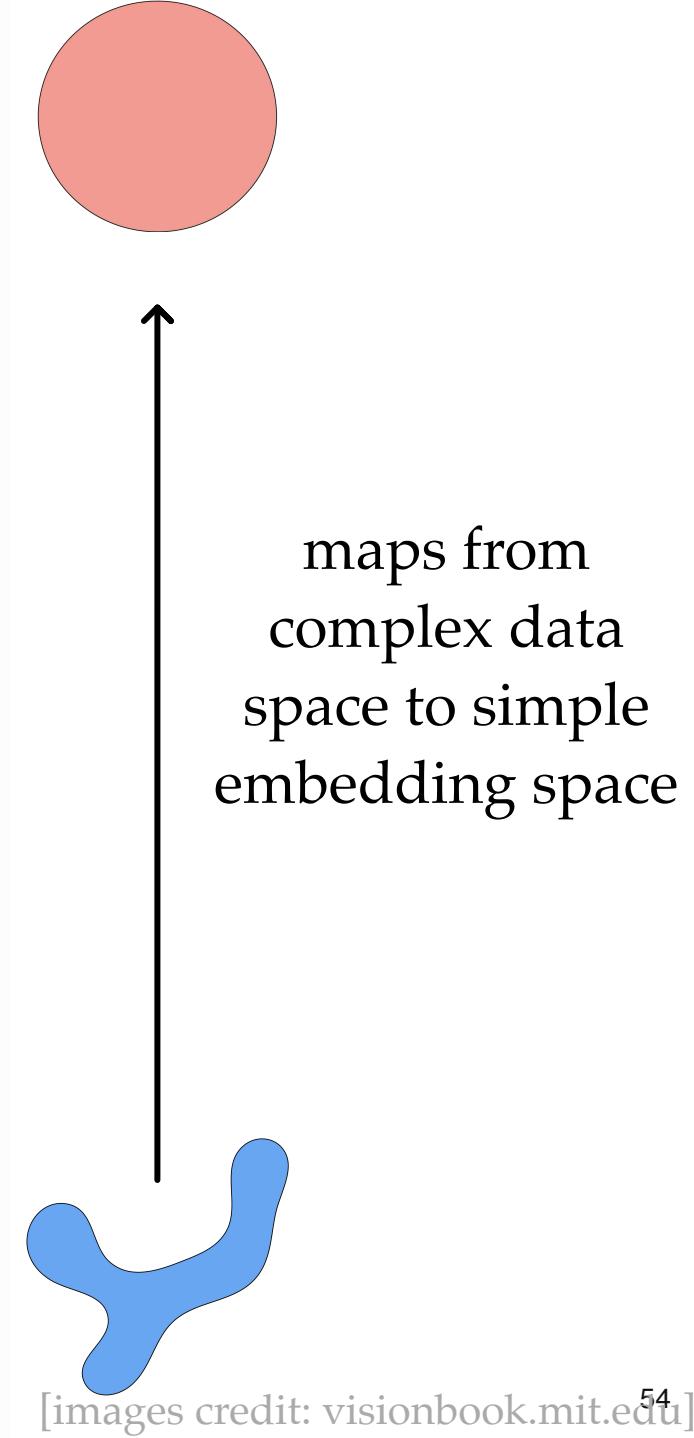
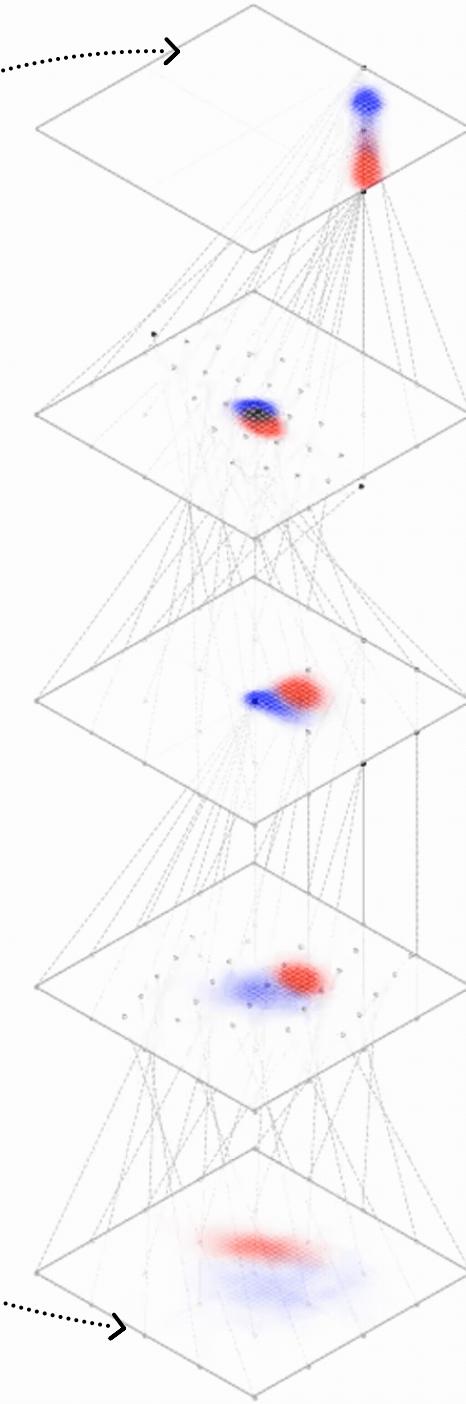
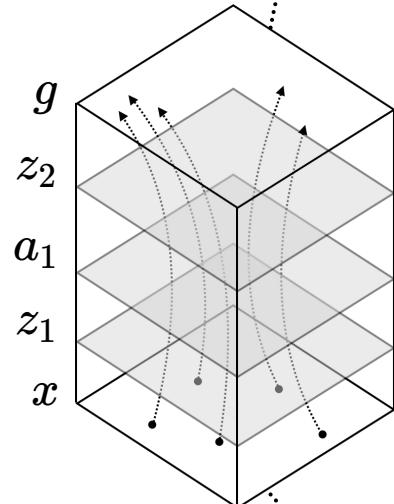
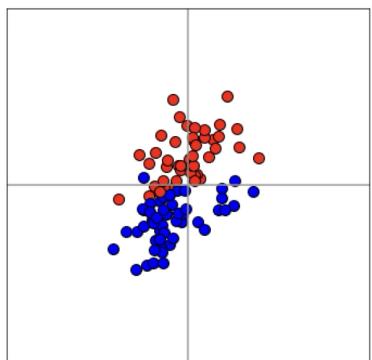
$$z_2 = \text{linear } (a_1)$$

$$a_1 = \text{ReLU}(z_1)$$

$$z_1 = \text{linear } (x)$$

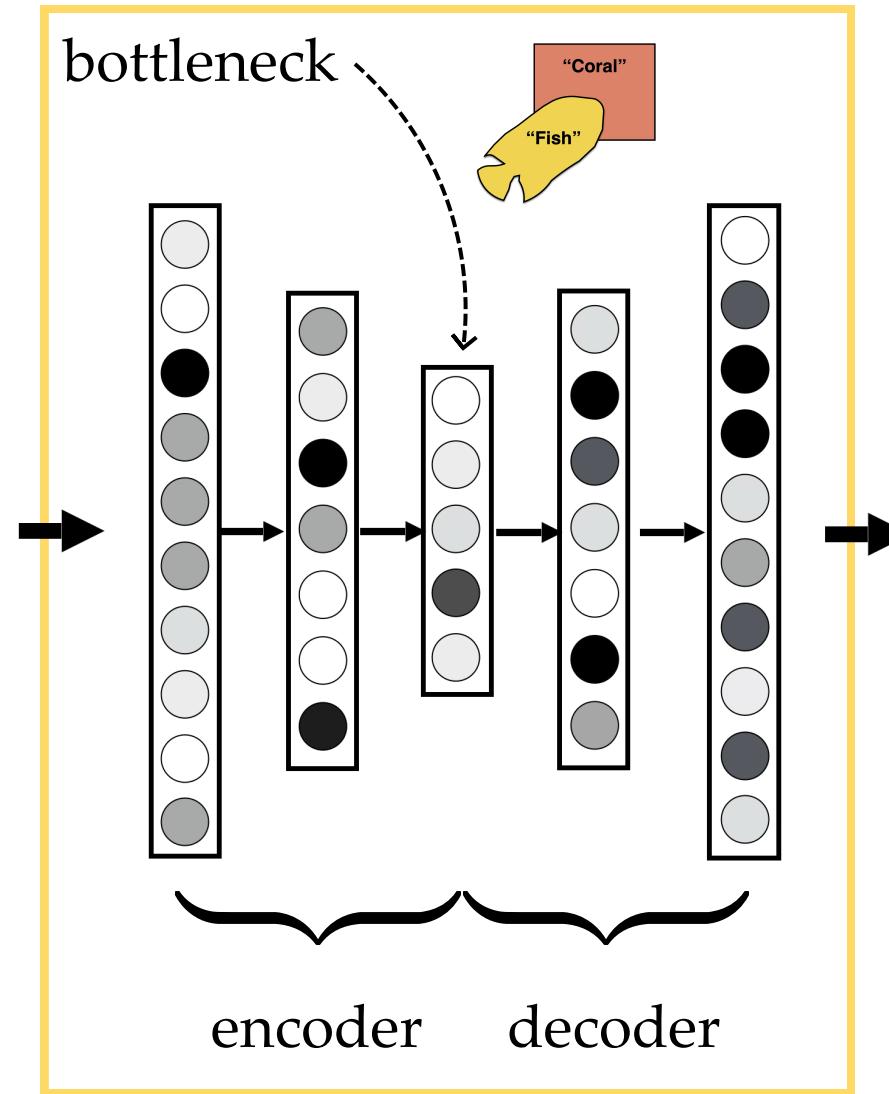
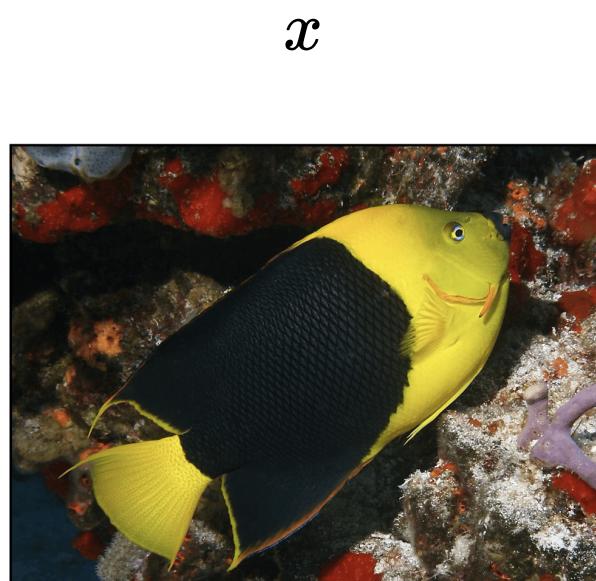
$$x \in \mathbb{R}^2$$

Training data



Auto-encoder

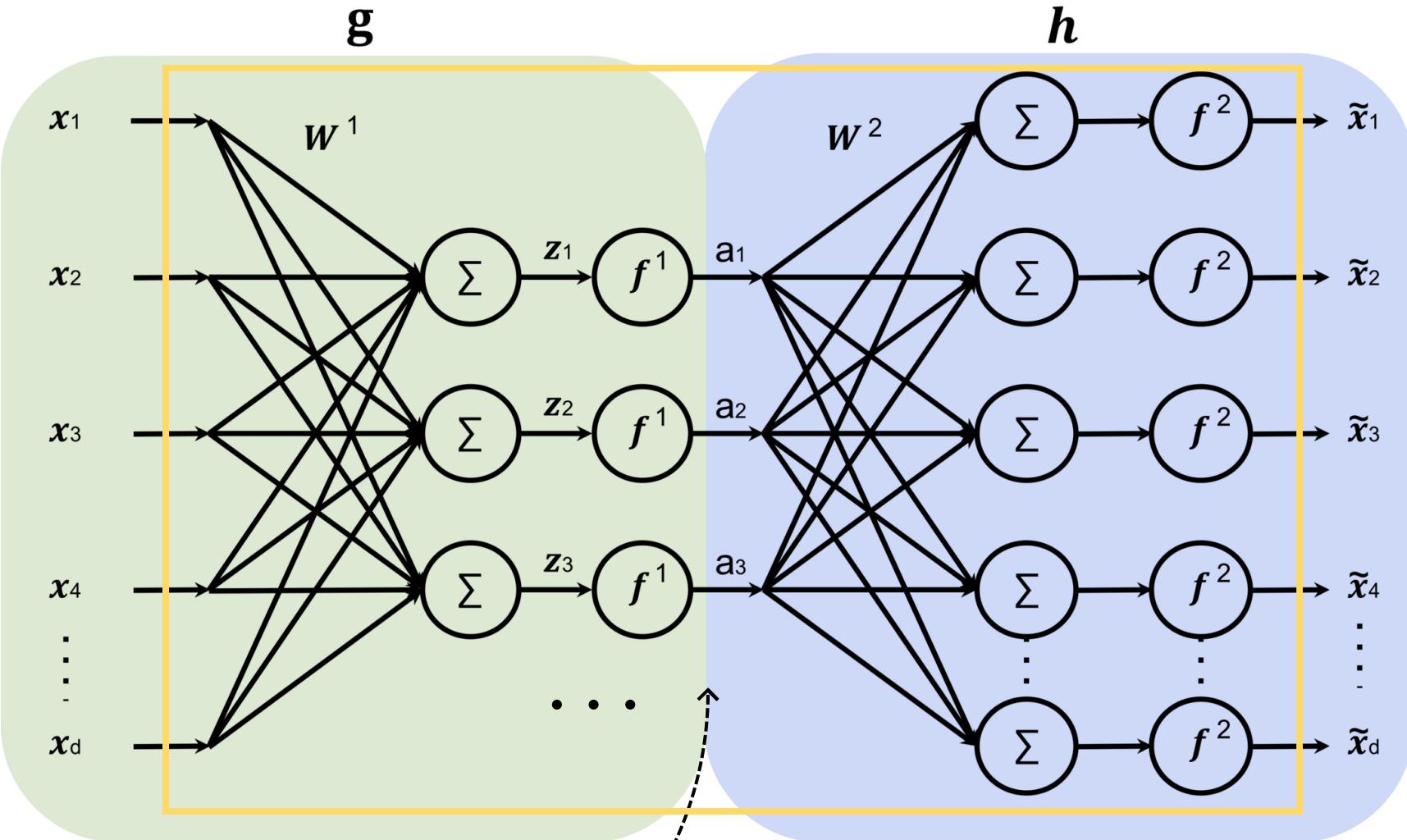
$$\min_W ||x - \tilde{x}||^2$$



$$\tilde{x} = \text{NN}(x; W)$$



input
 $x \in \mathbb{R}^d$

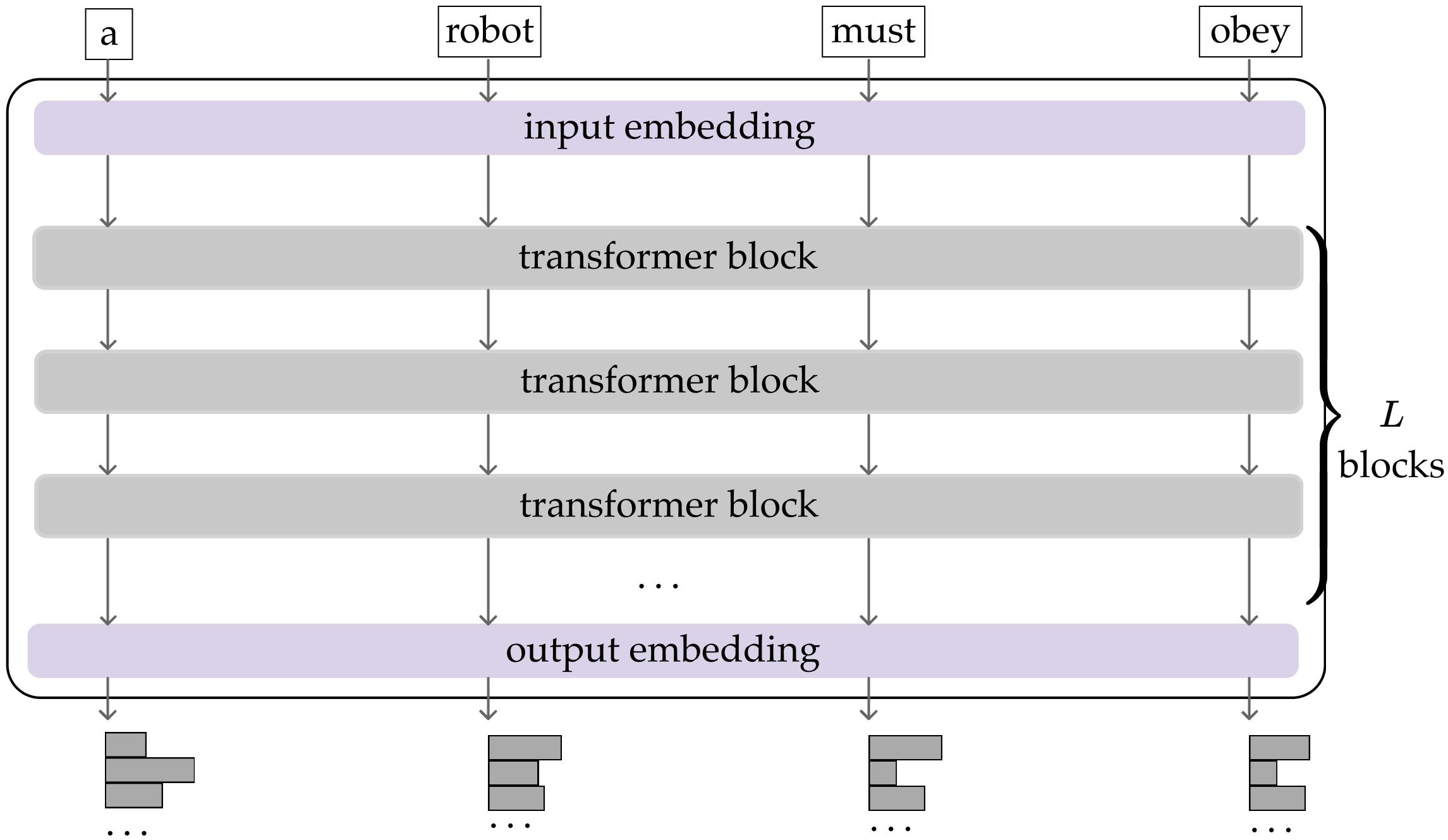


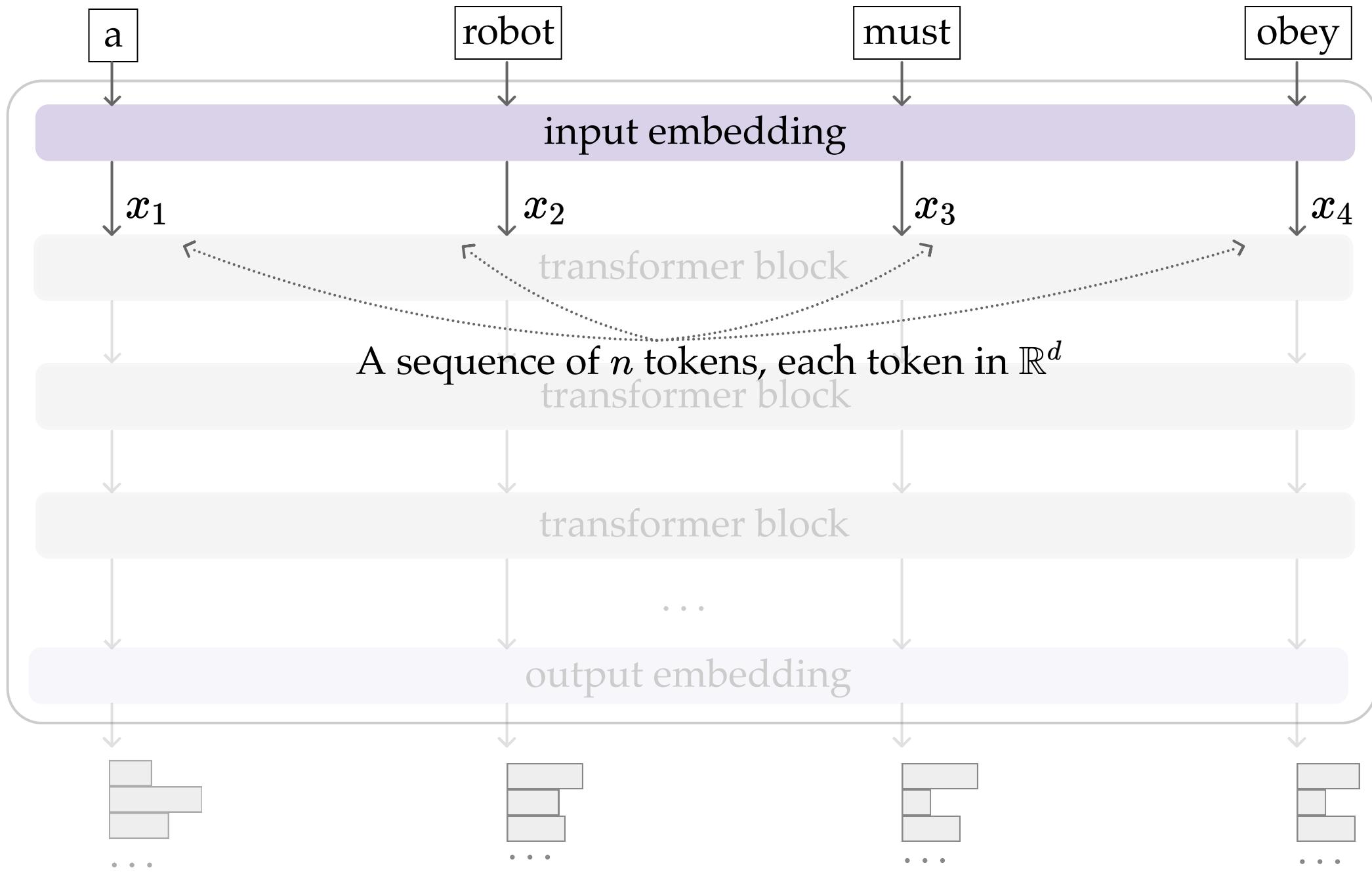
bottleneck
typically, has lower dimension than d

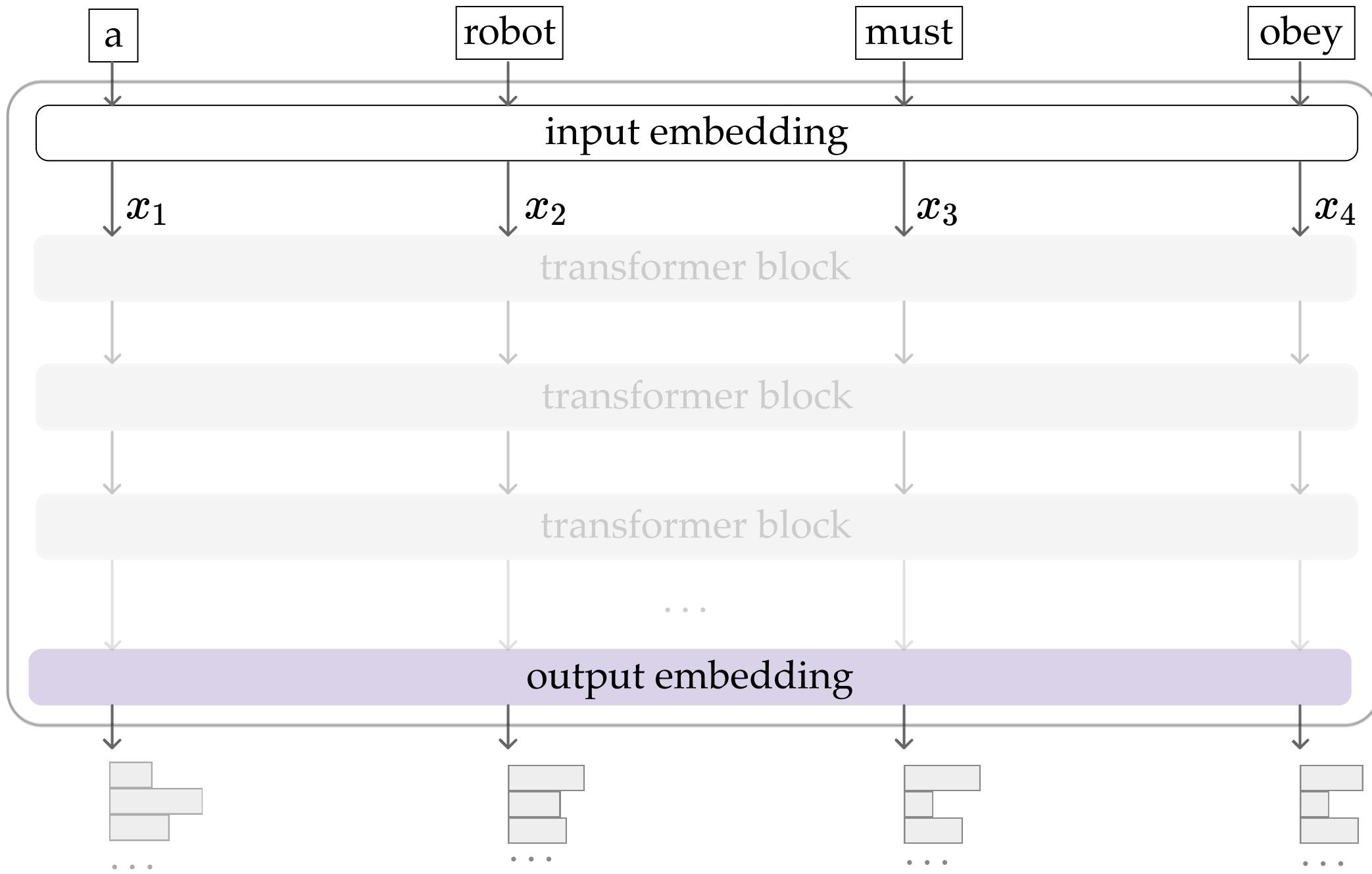
output
 $\tilde{x} \in \mathbb{R}^d$

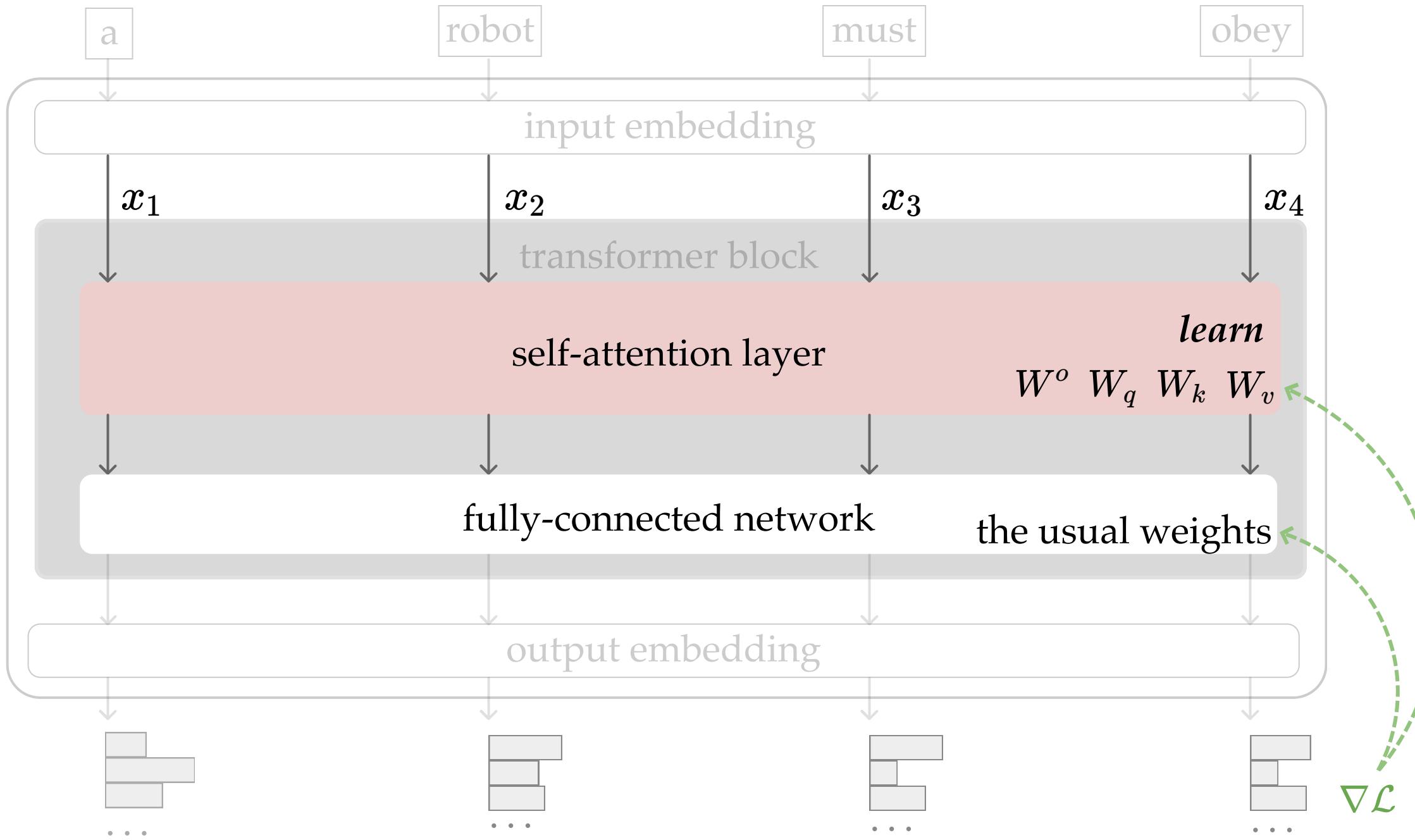
Week 9 - Transformers

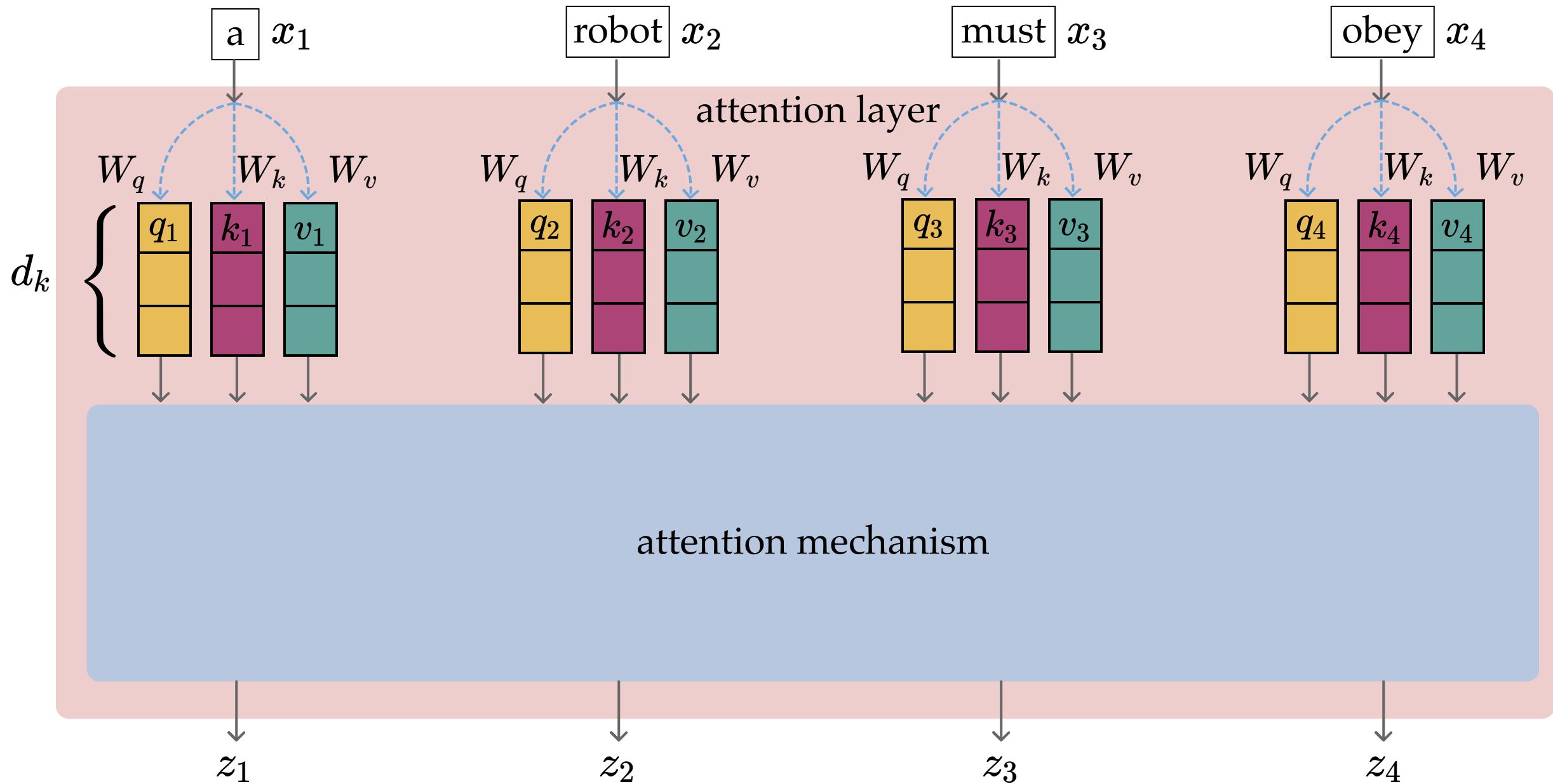
- A sequence of tokenized inputs: n tokens, each token x is d dimensional
- the attention mechanism (one head)
 - learn weights W_q, W_k, W_v to turn raw x inputs into (query, key, value)
 - the mechanics, softmax(raw attention score), shapes
 - masking: why and how
- parallel-processing machines
 - each head is processed in parallel
 - inside a head, each token is processed in parallel

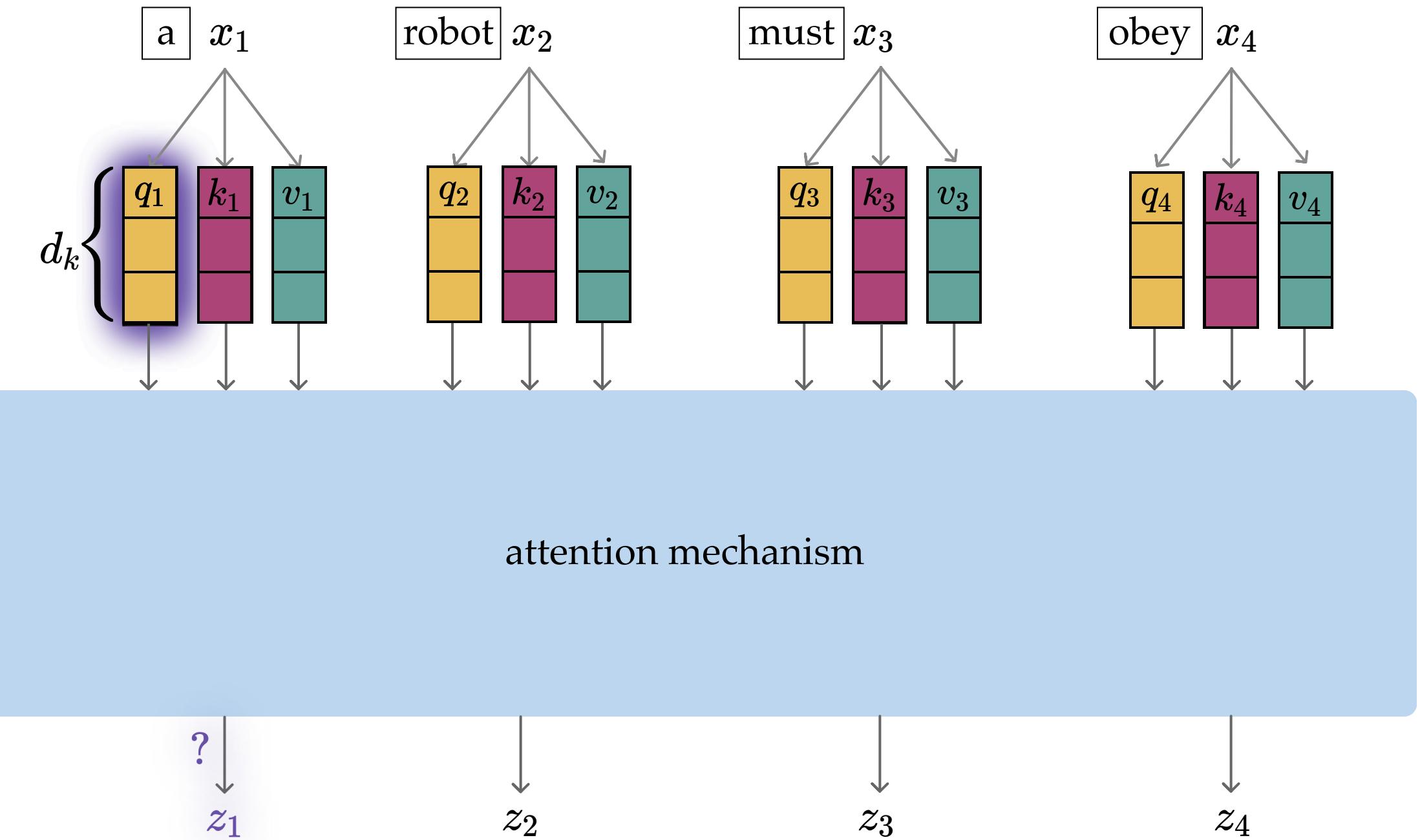


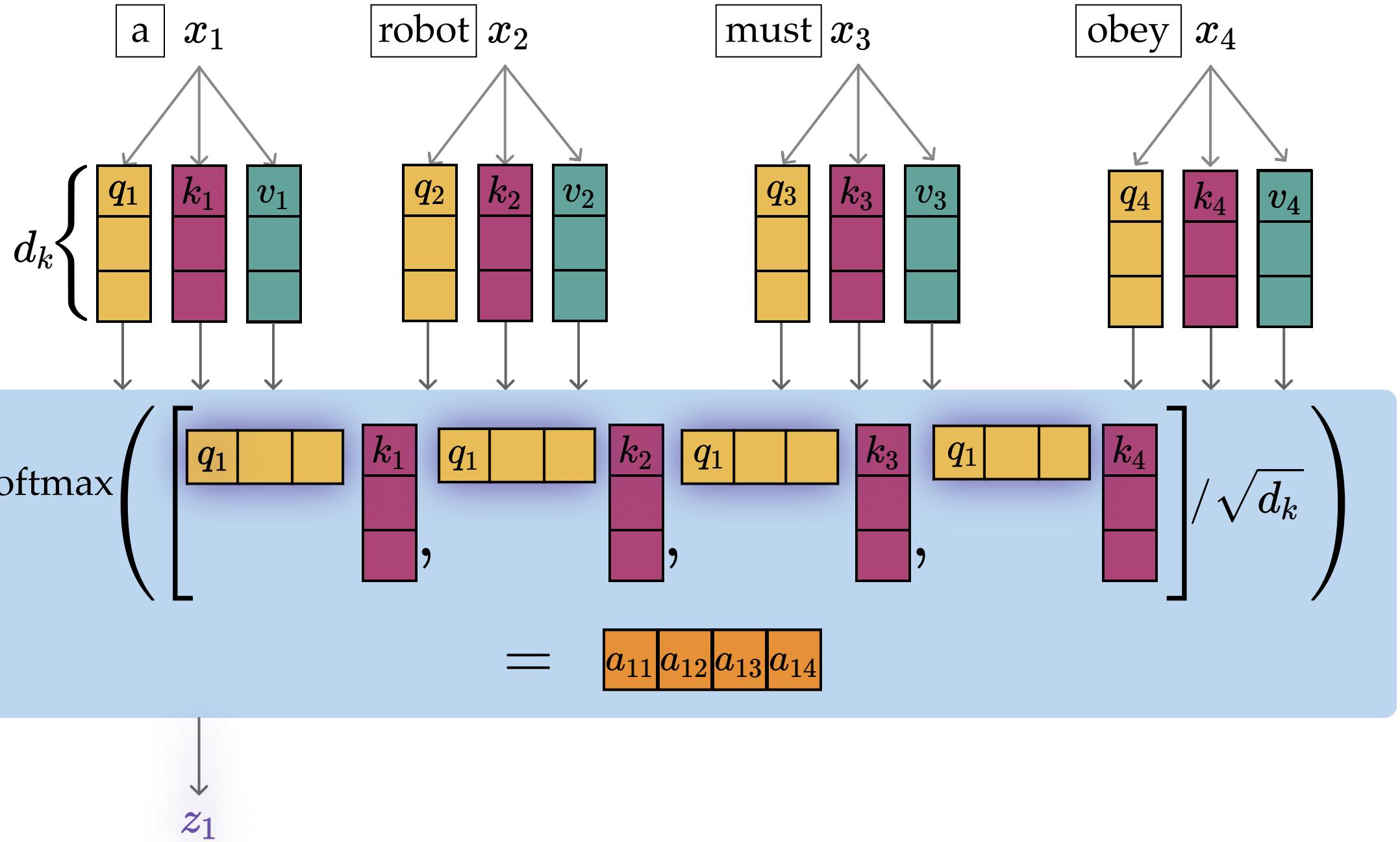


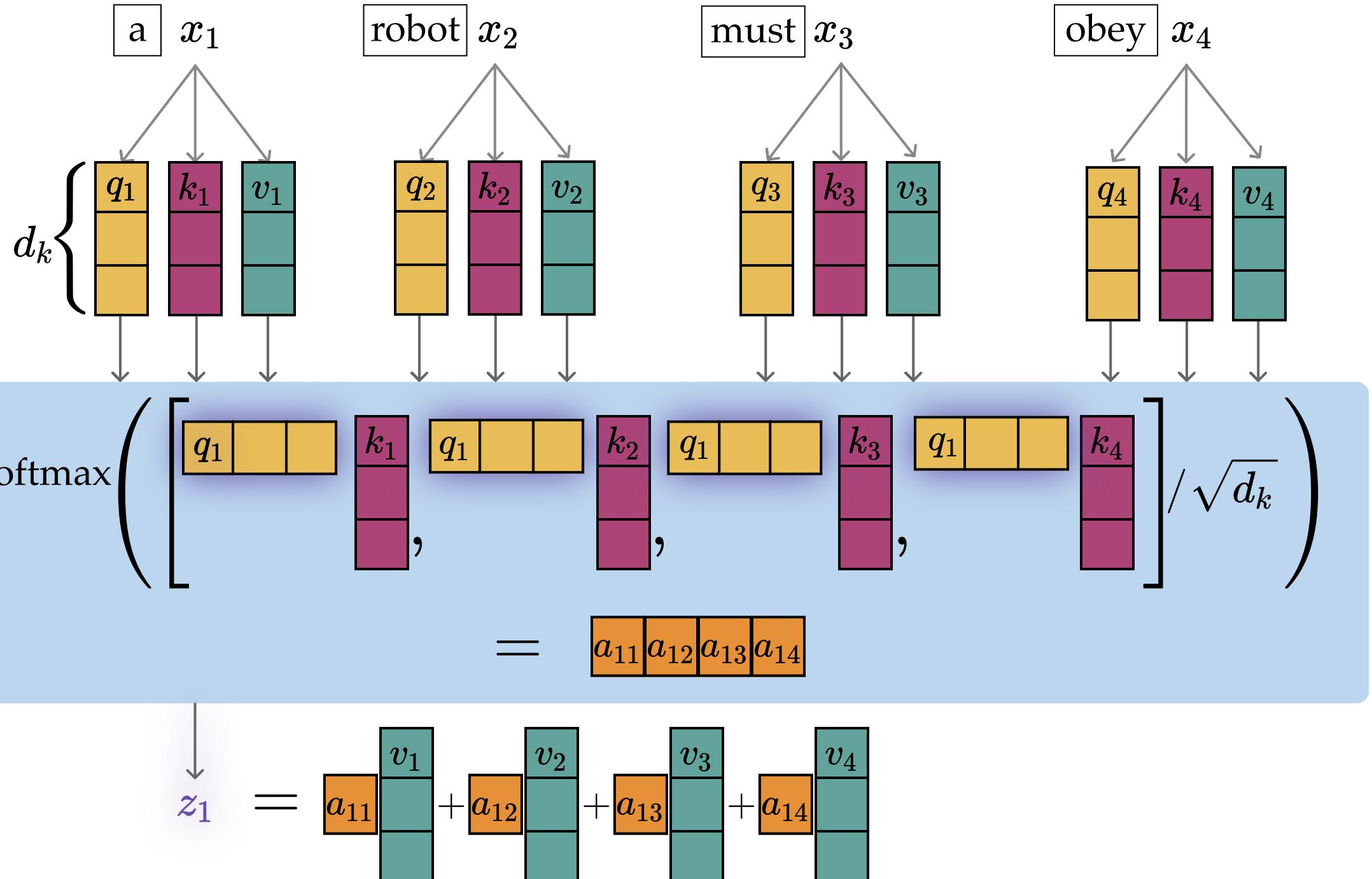




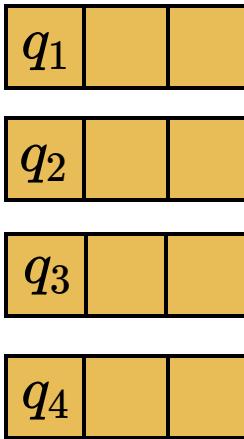


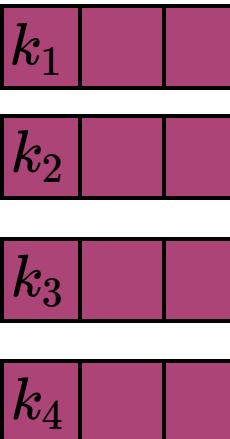




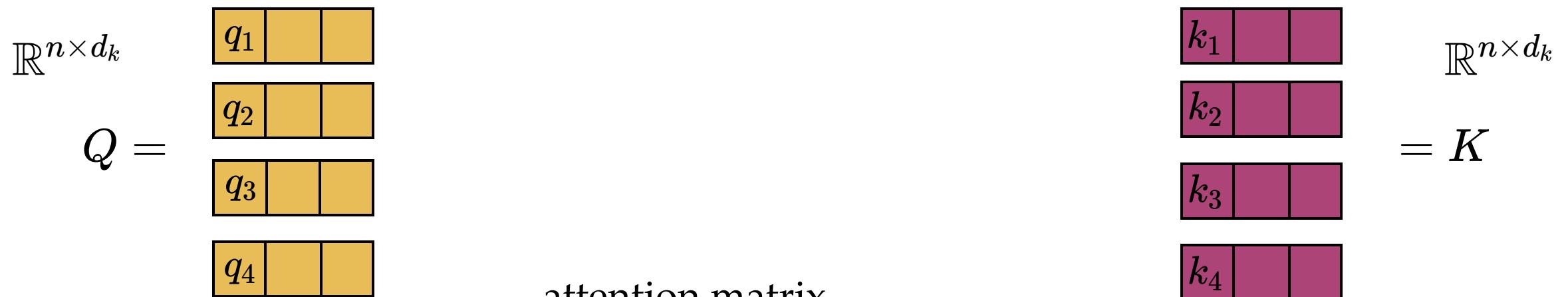


$\mathbb{R}^{n \times d_k}$

$$Q = \begin{matrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{matrix}$$


$$\begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix}$$


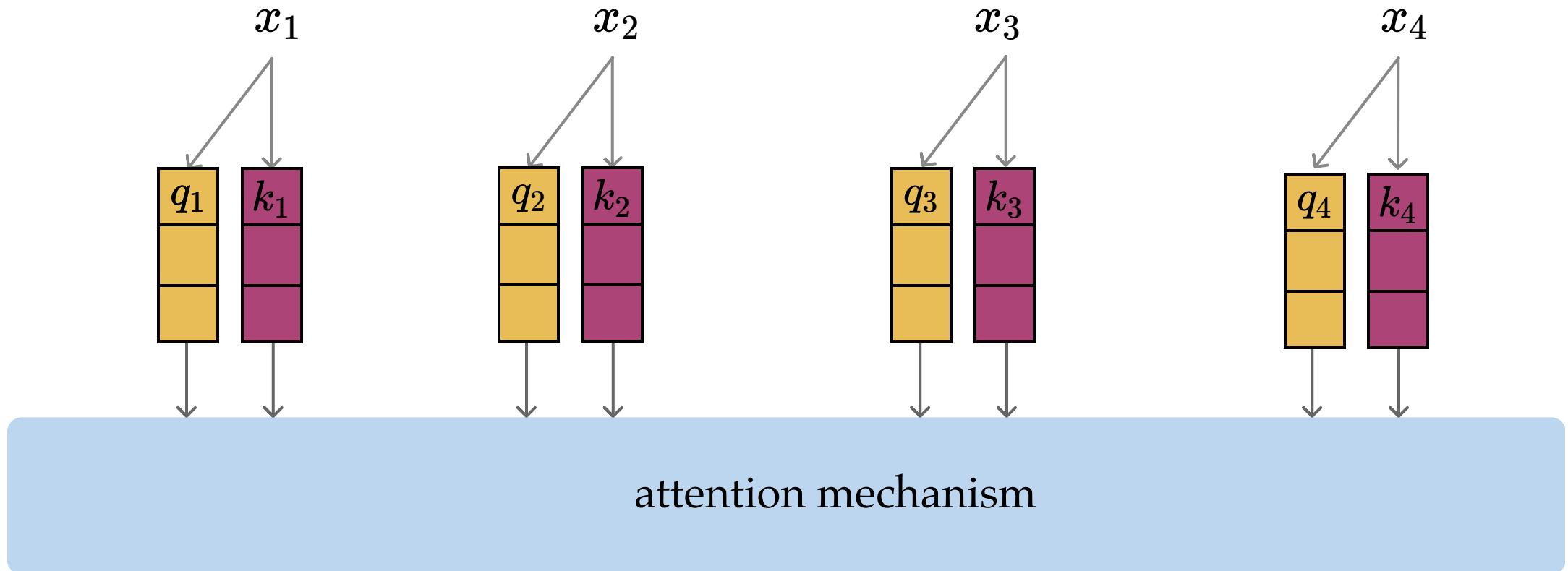
 $\mathbb{R}^{n \times d_k}$ $= K$



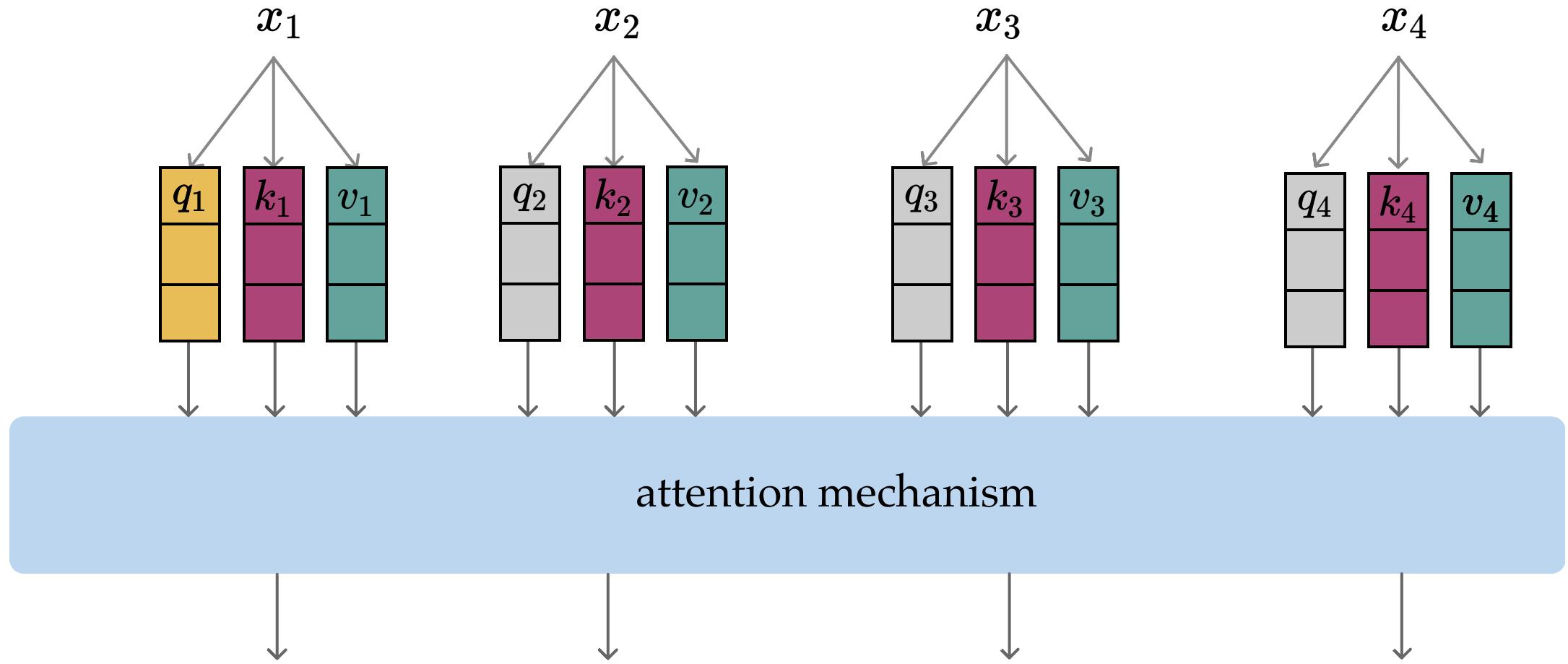
attention matrix

each row sums up to 1

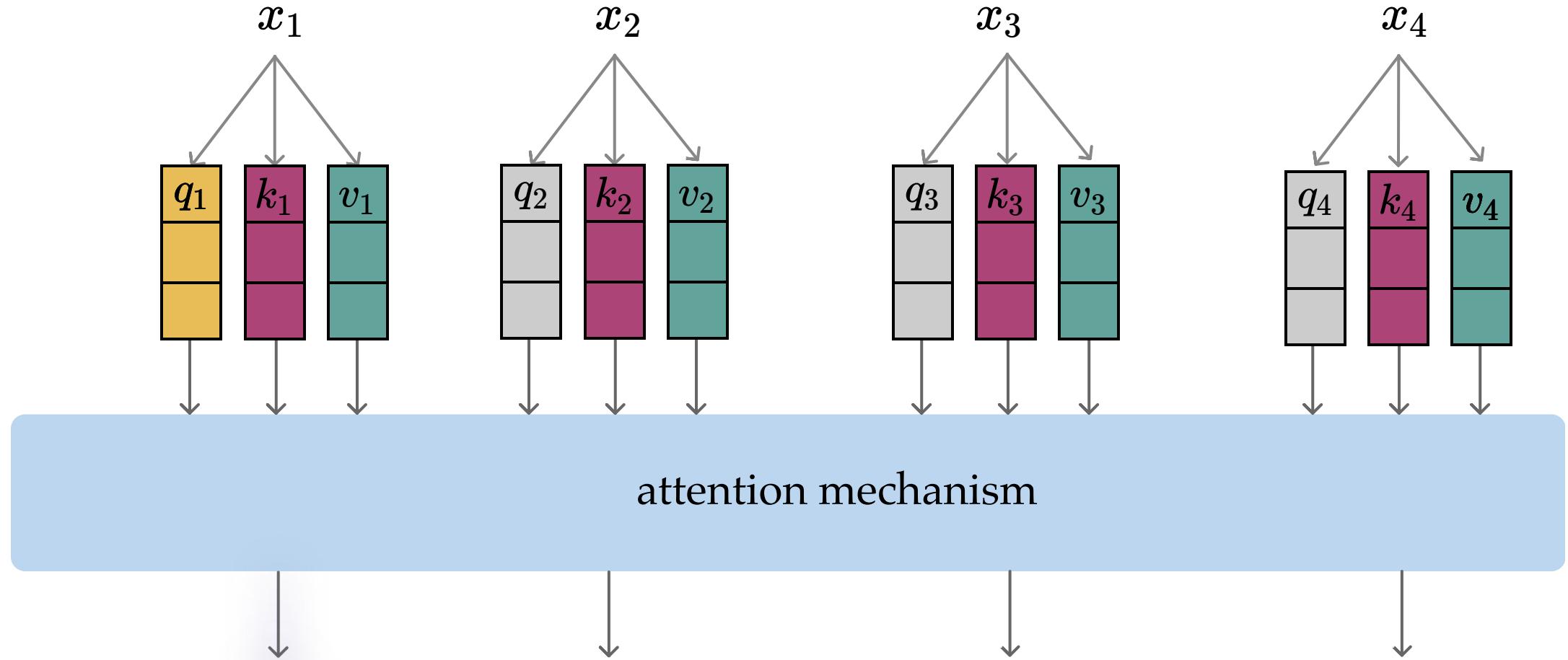
$$A = \begin{matrix} \mathbb{R}^{n \times n} \\ \downarrow \\ A = \begin{bmatrix} \text{softmax} \left(\begin{matrix} \text{orange} & \text{orange} & \text{orange} & \text{orange} \end{matrix} / \sqrt{d_k} \right) \\ \text{softmax} \left(\begin{matrix} \text{orange} & \text{orange} & \text{orange} & \text{orange} \end{matrix} / \sqrt{d_k} \right) \\ \text{softmax} \left(\begin{matrix} \text{orange} & \text{orange} & \text{orange} & \text{orange} \end{matrix} / \sqrt{d_k} \right) \\ \text{softmax} \left(\begin{matrix} \text{orange} & \text{orange} & \text{orange} & \text{orange} \end{matrix} / \sqrt{d_k} \right) \end{bmatrix} = \begin{matrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{matrix} \end{matrix}$$



a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}



a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}



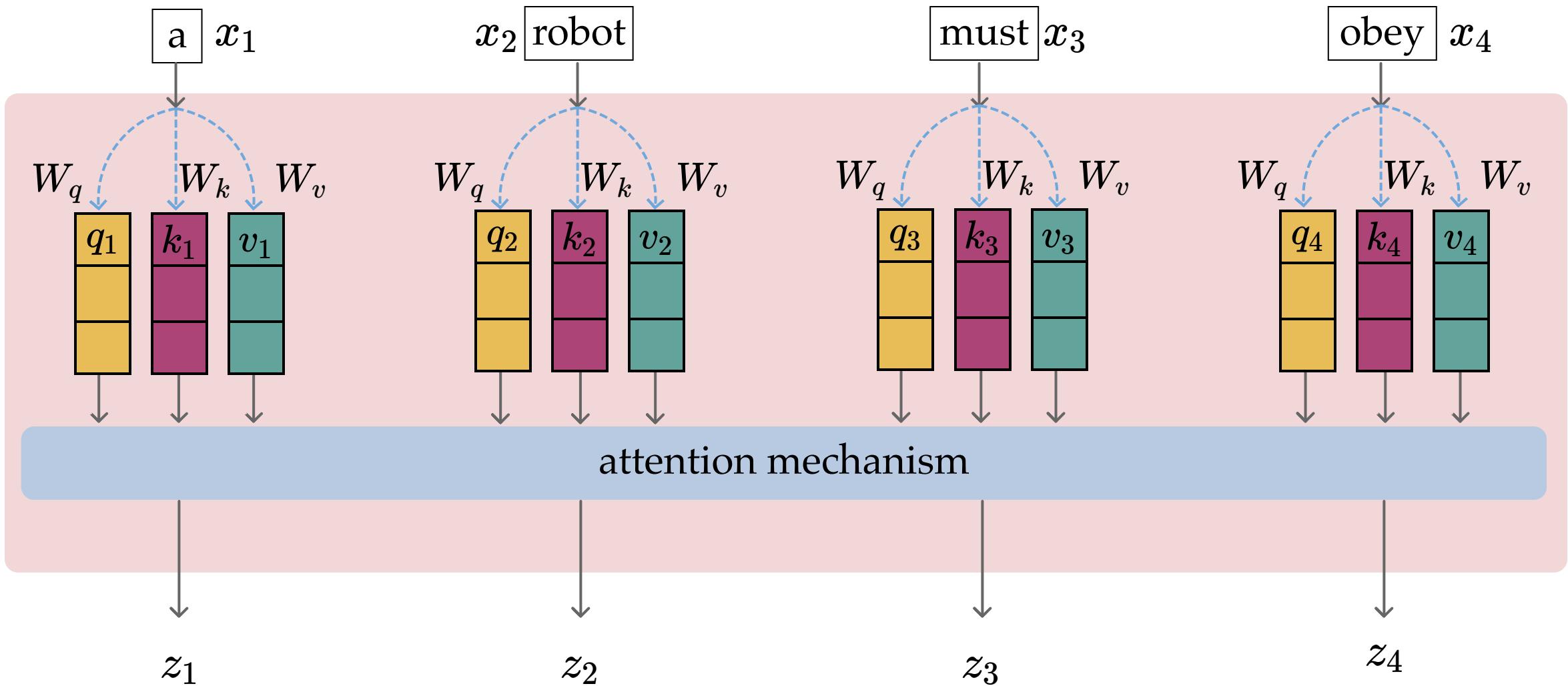
a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}

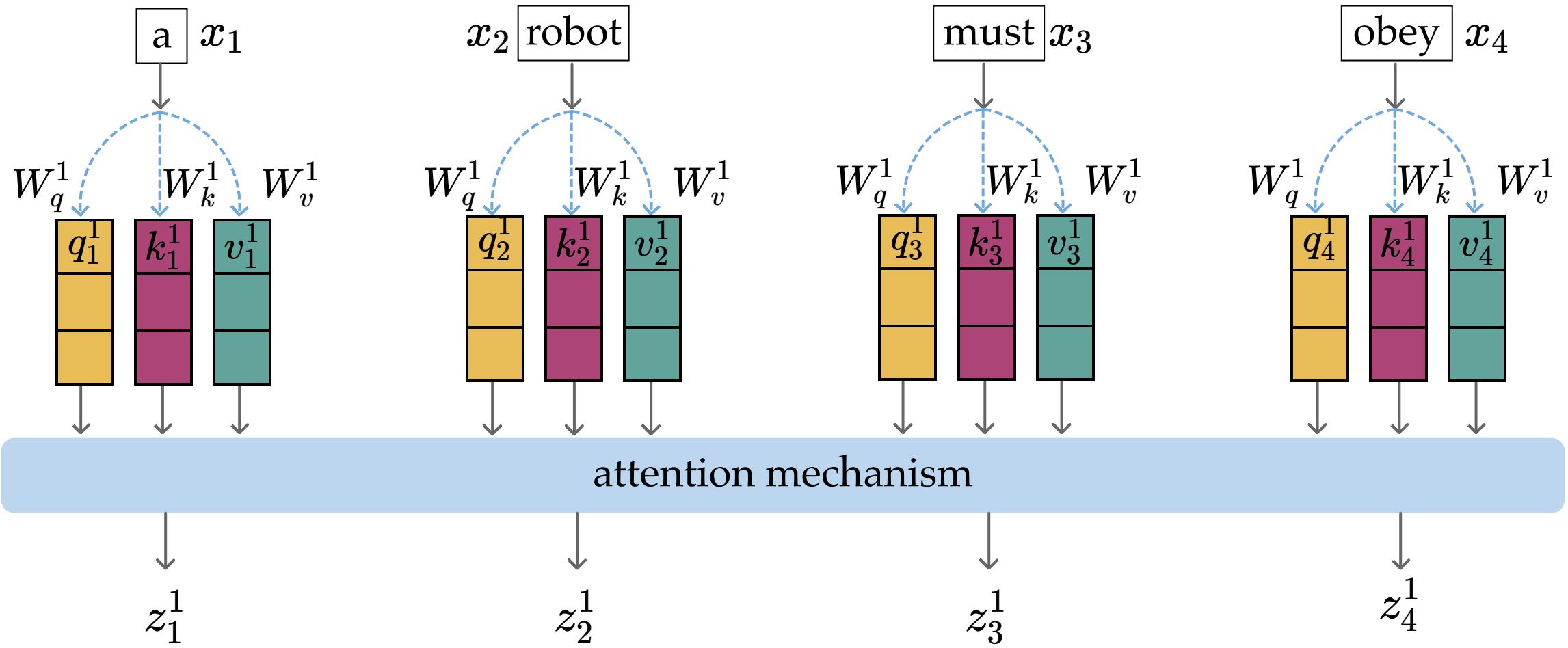
z_1

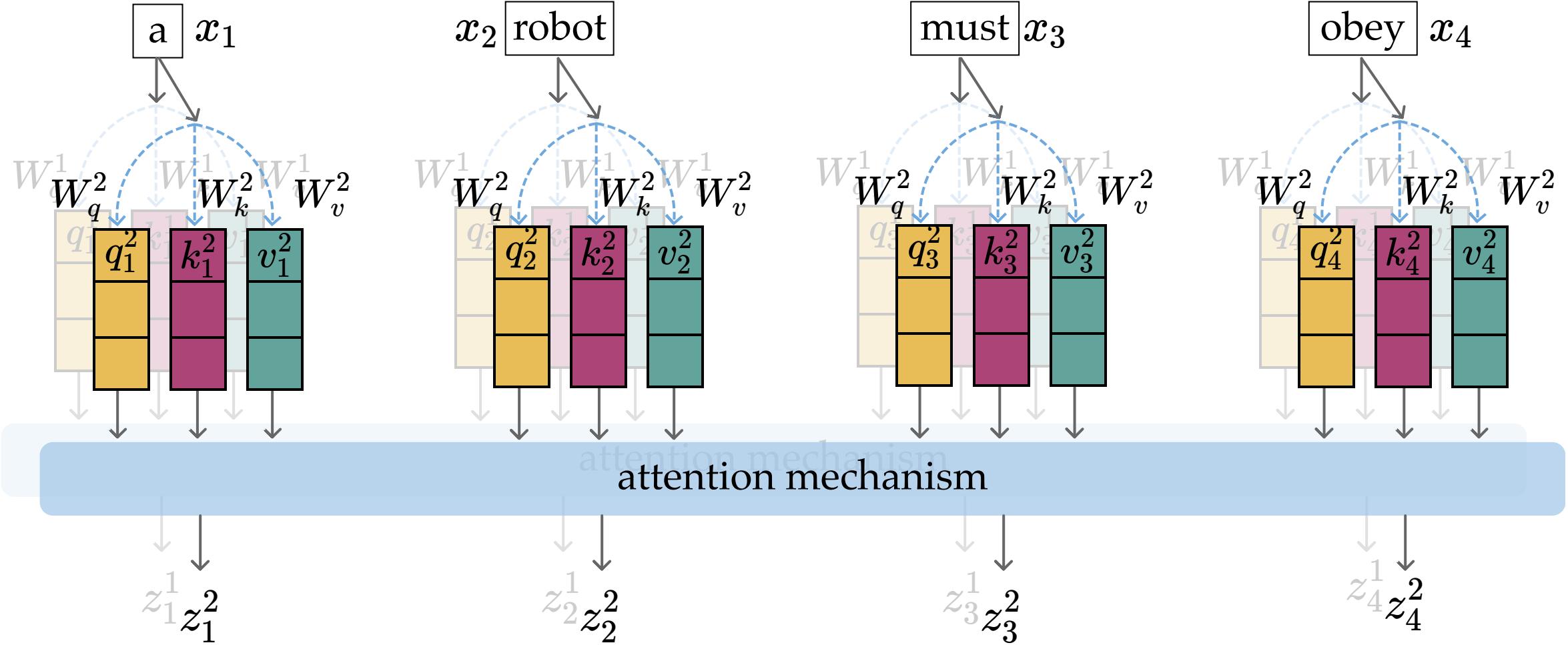
attention mechanism

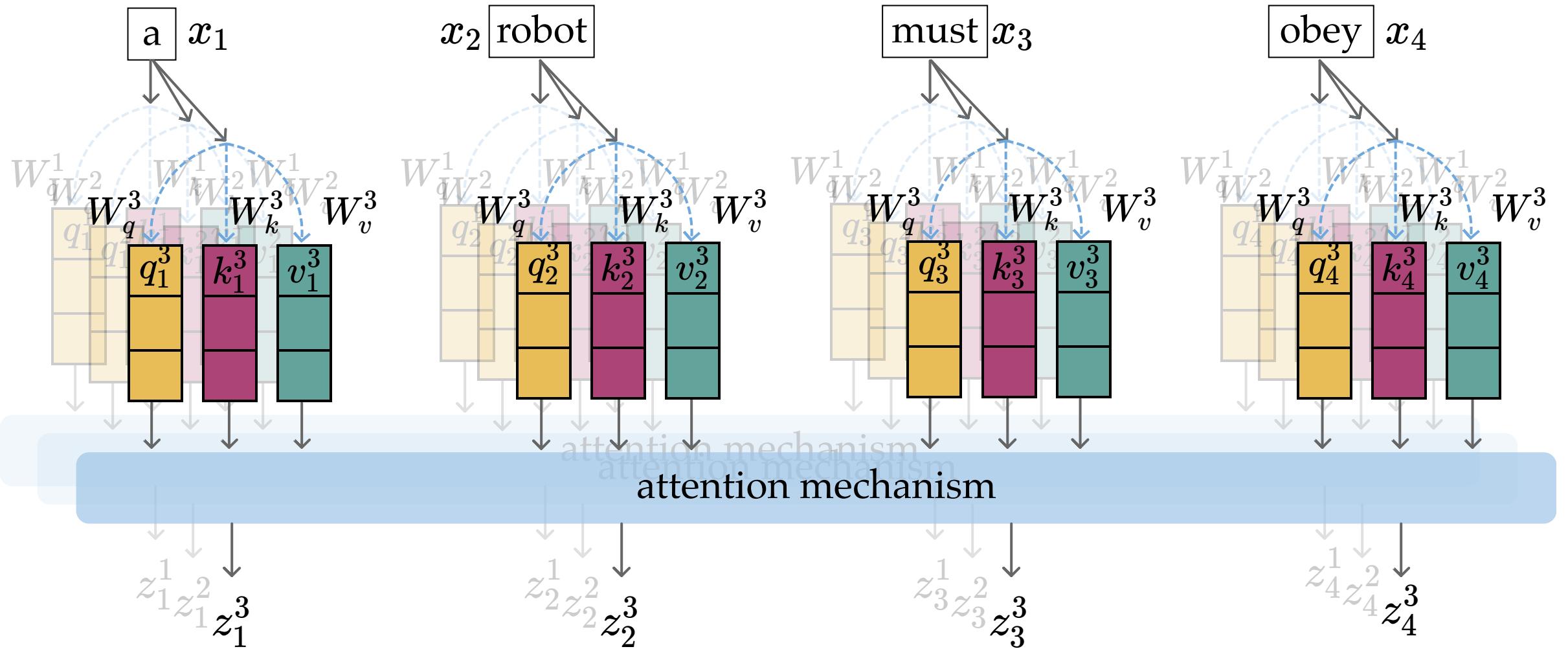
$$= \begin{matrix} a_{11} \\ a_{12} \\ a_{13} \\ a_{14} \end{matrix} + \begin{matrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{matrix} \in \mathbb{R}^{d_k}$$

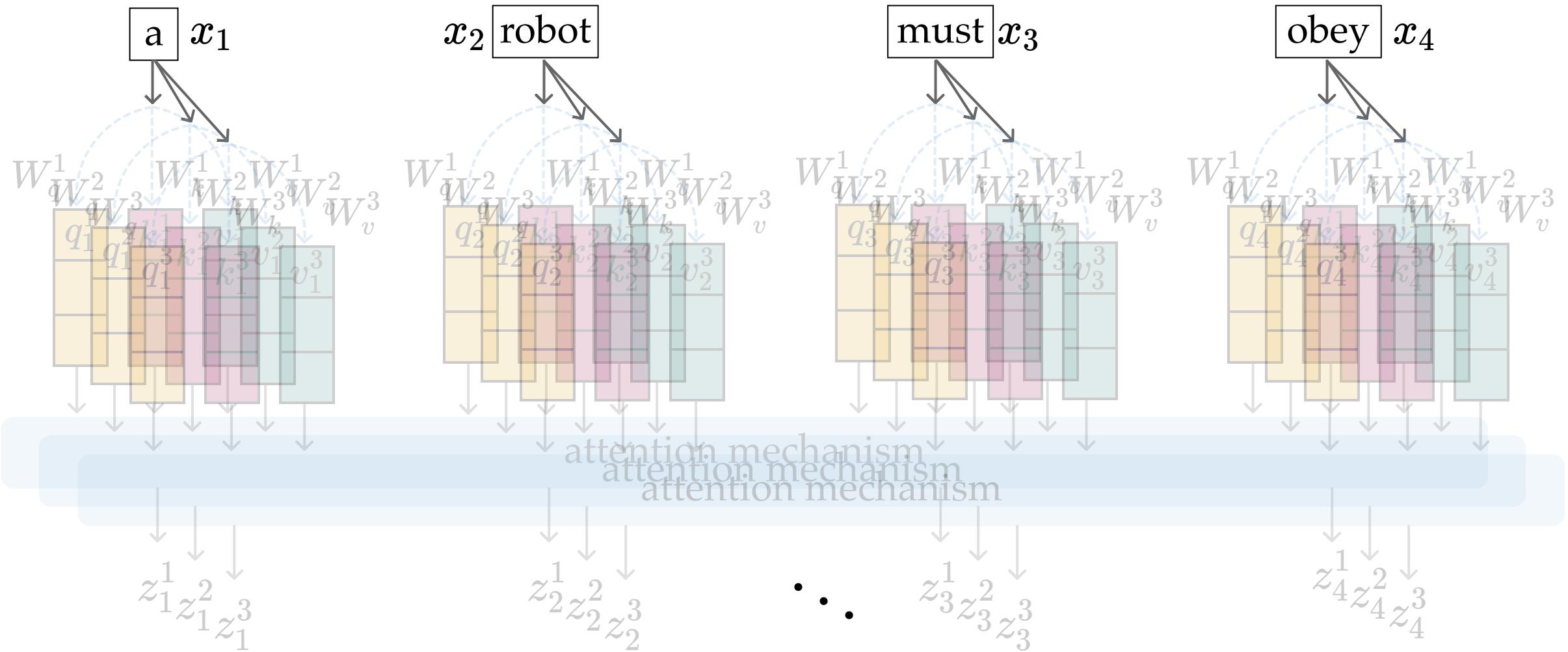
one attention head

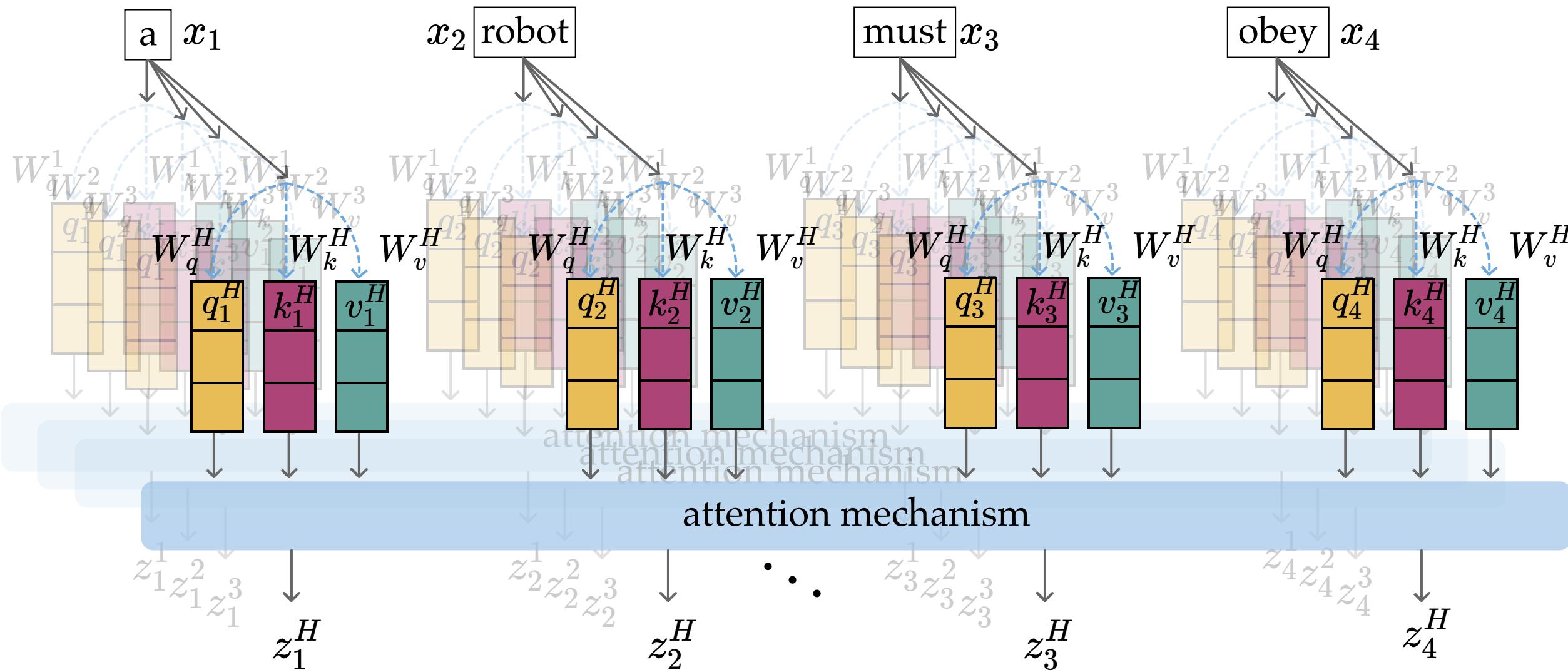




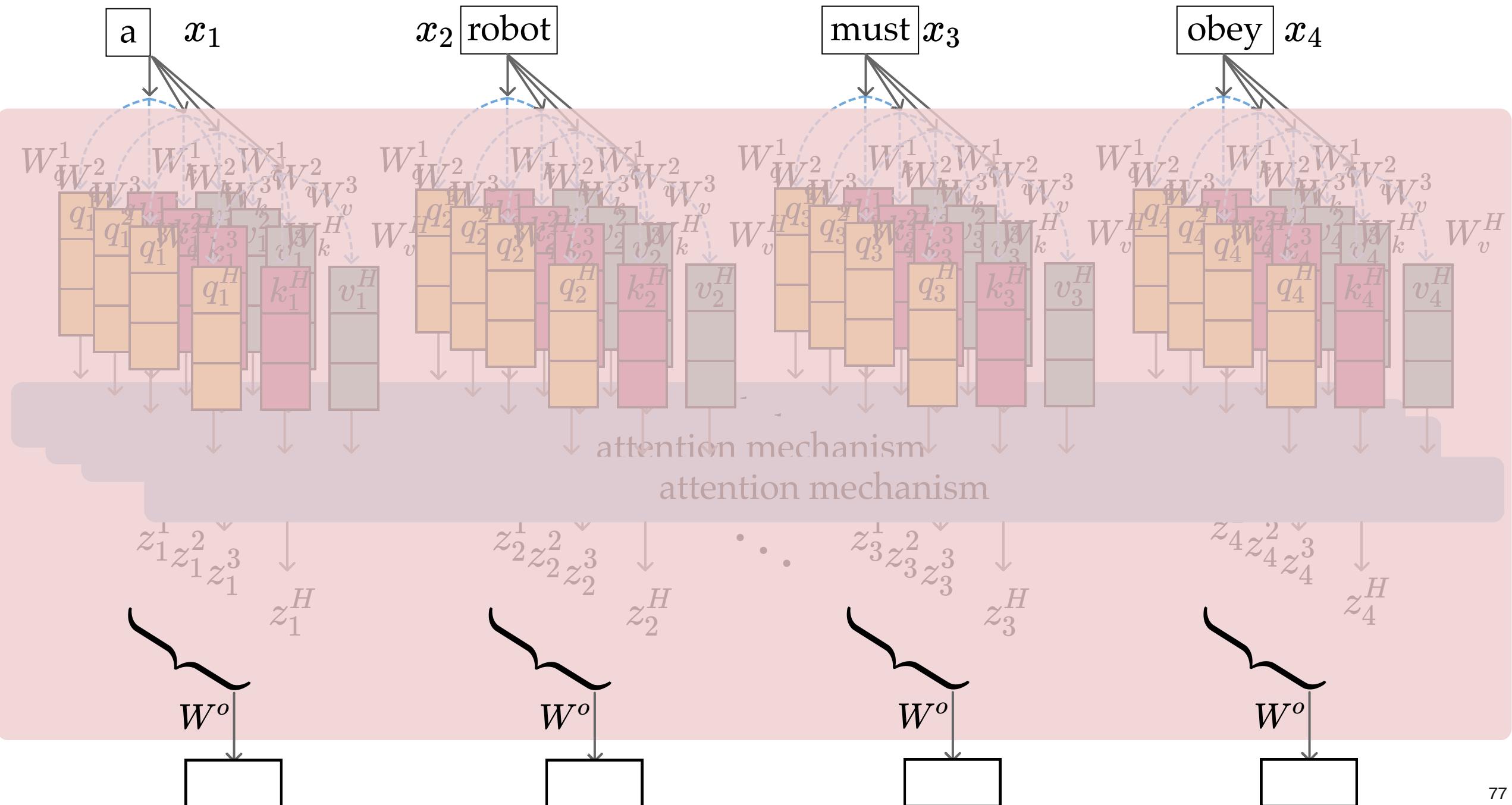








multi-head attention



Shape Example

n	num tokens	2
d	token dim	4
d_k	(qkv) dim	3
H	num heads	5

W_q^h	query proj	$d \times d_k$	4×3
W_k^h	key proj	$d \times d_k$	4×3
W_v^h	value proj	$d \times d_k$	4×3
W^o	output proj	$d \times Hd_k$	4×15
input	-	$n \times d$	2×4
Q^h	query	$n \times d_k$	2×3
K^h	key	$n \times d_k$	2×3
V^h	value	$n \times d_k$	2×3
A^h	attn matrix	$n \times n$	2×2
Z^h	head out.	$n \times d_k$	2×3
output	-	$n \times d$	2×4

learned

Week 10 - MDPs

- Definition (the five tuple)
 - π , V and Q : definition and interpretation
- Policy evaluation: given $\pi(s)$, calculate $V^\pi(s)$
 - via summation, or via Bellman recursion or equation
- Policy optimization: finding optimal policy $\pi^*(s)$
 - toy setup: solve via heuristics; more generally: Q value-iteration
- Interpretation of optimal policy
 - how R, γ, h can change optimal policy

Markov Decision Processes - Definition and terminologies

- \mathcal{S} : state space, contains all possible states s .
- \mathcal{A} : action space, contains all possible actions a .
- $T(s, a, s')$: the probability of transition from state s to s' when action a is taken.
- $R(s, a)$: reward, takes in a (state, action) pair and returns a reward.
- $\gamma \in [0, 1]$: discount factor, a scalar.

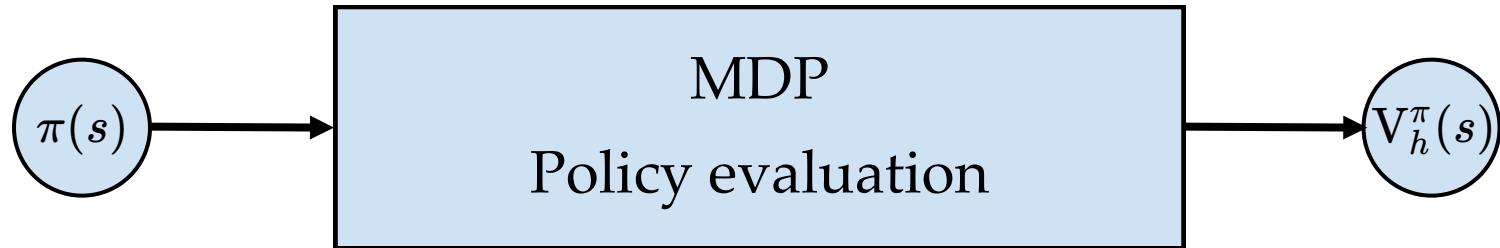
- $\pi(s)$: policy, takes in a state and returns an action.

The goal of an MDP is to find a "good" policy.

In 6.390,

- \mathcal{S} and \mathcal{A} are small discrete sets, unless otherwise specified.
- s' and a' are short-hand for the next-timestep
- $R(s, a)$ is deterministic and bounded.
- $\pi(s)$ is deterministic.

Quick summary



1. By summing h terms:

Recall: For a *given* policy $\pi(s)$, the (state) **value functions**

$$V_h^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{h-1} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi \right], \forall s, h$$

2. By leveraging structure:

finite-horizon Bellman recursions

$$V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s'), \forall s$$

infinite-horizon Bellman equations

$$V_\infty^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\infty^\pi(s'), \forall s$$

Optimal state-action value functions $Q_h^*(s, a)$

$Q_h^*(s, a)$: the expected sum of discounted rewards for

- starting in state s ,
- take action a , for one step
- act **optimally** there afterwards for the remaining $(h - 1)$ steps

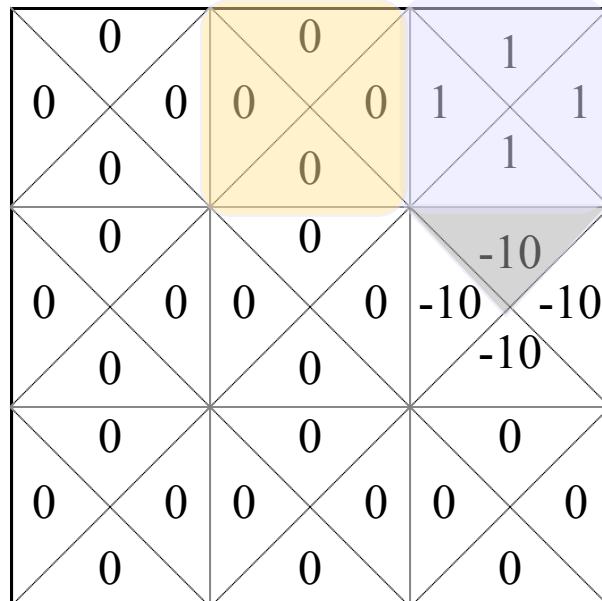
$$Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a'), \forall s, a, h$$



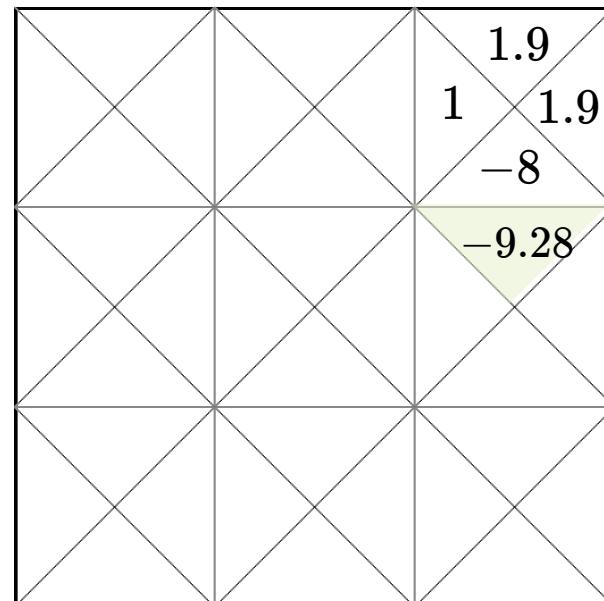
$Q_h^*(s, a)$: the value for

- starting in state s ,
- take action a , for one step
- act **optimally** there afterwards for the remaining $(h - 1)$ steps

$$Q_1^*(s, a) = R(s, a)$$



$$Q_2^*(s, a)$$



Recall:

$$\gamma = 0.9$$

States and
one special
transition:

1	2	3
4	5	6
7	8	9

States and one special transition:

States: 1, 2, 3, 4, 5, 6, 7, 8, 9

Transitions:

- From 1 to 2: 20%
- From 1 to 3: 80%
- From 2 to 5: 20%
- From 2 to 6: 80%
- From 3 to 6: 100%

- receive $R(6, \uparrow)$
- act **optimally** for one more timestep, at the next state s'
 - 20% chance, $s' = 2$, act optimally, receive $\max_{a'} Q_1^*(2, a')$
 - 80% chance, $s' = 3$, act optimally, receive $\max_{a'} Q_1^*(3, a')$

Let's consider $Q_2^*(6, \uparrow)$

$$\begin{aligned}
 &= R(6, \uparrow) + \gamma[.2 \max_{a'} Q_1^*(2, a') + .8 \max_{a'} Q_1^*(3, a')] \\
 &= -10 + .9[.2 \times 0 + .8 \times 1] = -9.28
 \end{aligned}$$

Given the recursion $Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a')$

we can have an infinite horizon equation

$$Q_\infty^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_\infty^*(s', a')$$

Value Iteration

1. **for** $s \in \mathcal{S}, a \in \mathcal{A}$:

2. $Q_{\text{old}}(s, a) = 0$

3. **while** True:

4. **for** $s \in \mathcal{S}, a \in \mathcal{A}$:

5. $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

6. **if** $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$.

7. **return** Q_{new}

8. $Q_{\text{old}} \leftarrow Q_{\text{new}}$

$Q_\infty^*(s, a)$

if run this block h times
and break, then the
returns are exactly Q_h^*

{

Week 11 - Reinforcement Learning

- How RL setup differs from MDP
- Q-learning algorithm
 - Forward thinking: given experiences, work out Q-values.
 - Backward thinking: given realized Q-values, work out experiences.
 - Two new hyper-parameters (compared with MDP value iteration):
 - ϵ -greedy action selection
 - α the learning rate
- The idea of fitting parameterized Q-functions via regression, can handle larger or continuous state / action space

~~Markov Decision Processes~~ - Definition and terminologies

Reinforcement Learning

- \mathcal{S} : state space, contains all possible states s .
- \mathcal{A} : action space, contains all possible actions a .
- $T(s, a, s')$: the probability of transition from state s to s' when action a is taken.
- $R(s, a)$: reward, takes in a (state, action) pair and returns a reward.
- $\gamma \in [0, 1]$: discount factor, a scalar.
- $\pi(s)$: policy, takes in a state and returns an action.

The goal of an ~~MDP~~ problem is to find a "good" policy.

RL

Value Iteration($\mathcal{S}, \mathcal{A}, T, R, \gamma, \epsilon$)

```

1. for  $s \in \mathcal{S}, a \in \mathcal{A}$  :
2.    $Q_{\text{old}}(s, a) = 0$ 
3. while True:
4.   for  $s \in \mathcal{S}, a \in \mathcal{A}$  :
5.      $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$ 
6.   if  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$  :
7.     return  $Q_{\text{new}}$ 
8.    $Q_{\text{old}} \leftarrow Q_{\text{new}}$ 

```

"calculating"

Q-Learning ($\mathcal{S}, \mathcal{A}, \gamma, \alpha, s_0$ max-iter)

```

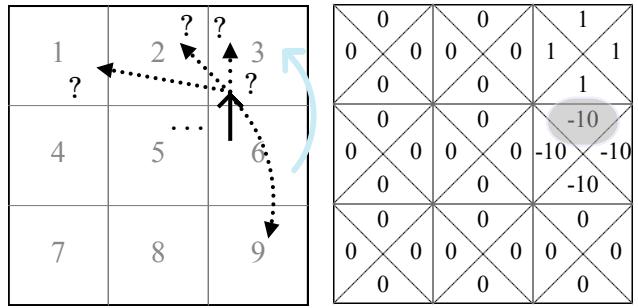
1.  $i = 0$ 
2. for  $s \in \mathcal{S}, a \in \mathcal{A}$  :
3.    $Q_{\text{old}}(s, a) = 0$ 
4.    $s \leftarrow s_0$ 
5. while  $i < \text{max-iter}$  :
6.    $a \leftarrow \text{select\_action}(s, Q_{\text{old}}(s, a))$ 
7.    $r, s' \leftarrow \text{execute}(a)$ 
8.    $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$ 
9.    $s \leftarrow s'$ 
10.   $i \leftarrow (i + 1)$ 
11.   $Q_{\text{old}} \leftarrow Q_{\text{new}}$ 
12. return  $Q_{\text{new}}$ 

```

"learning" (estimating)

$\gamma = 0.9$

rewards now known



Q-learning update

e.g. pick $\alpha = 0.5$

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha) Q_{\text{old}}(s, a) + \alpha \left(r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

 $Q_{\text{old}}(s, a)$

0	0	0	0	1	1
0	0	0	0	1	1
0	0	0	0	-10	-10
0	0	0	0	-10	-10
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

 $Q_{\text{new}}(s, a)$

0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

To update the estimate of $Q(6, \uparrow)$:

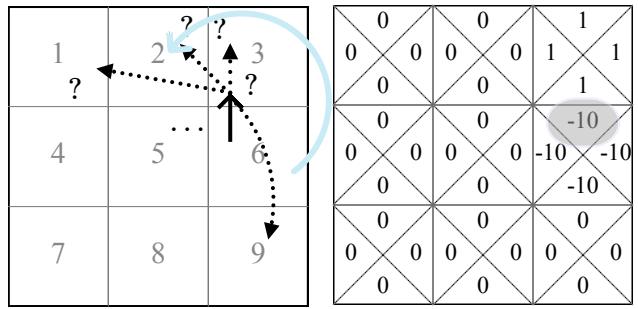
- execute $(6, \uparrow)$
- observe reward $r = -10$, next state $s' = 3$
- update $Q(6, \uparrow)$ as:

$$(1 - 0.5) * -10 + 0.5(-10 + 0.9 \max_{a'} Q_{\text{old}}(3, a'))$$

$$= -5 + 0.5(-10 + 0.9) = -9.55$$

$\gamma = 0.9$

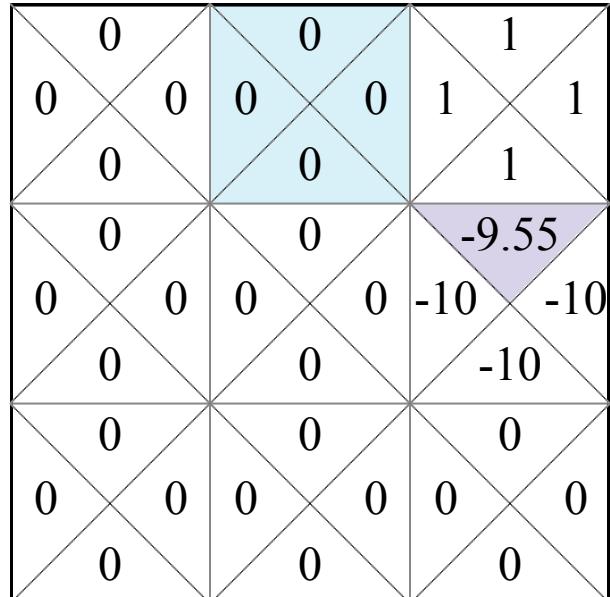
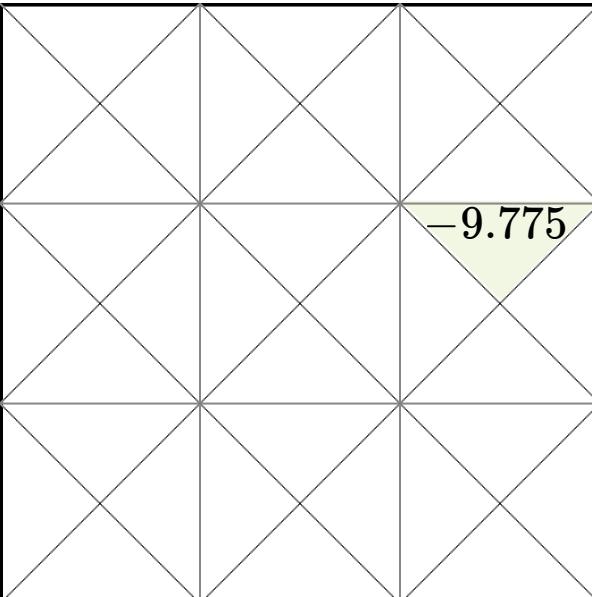
rewards now known



Q-learning update

e.g. pick $\alpha = 0.5$

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha \left(r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

 $Q_{\text{old}}(s, a)$  $Q_{\text{new}}(s, a)$ To update the estimate of $Q(6, \uparrow)$:

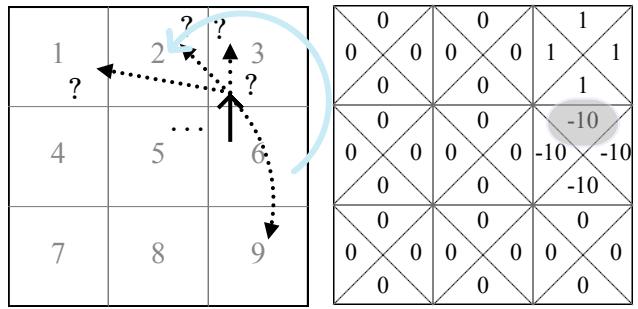
- execute $(6, \uparrow)$
- observe reward $r = -10$, next state $s' = 2$
- update $Q(6, \uparrow)$ as:

$$(1 - 0.5) * -9.55 + 0.5(-10 + 0.9 \max_{a'} Q_{\text{old}}(2, a'))$$

$$= -4.775 + 0.5(-10 + 0) = -9.775$$

$\gamma = 0.9$

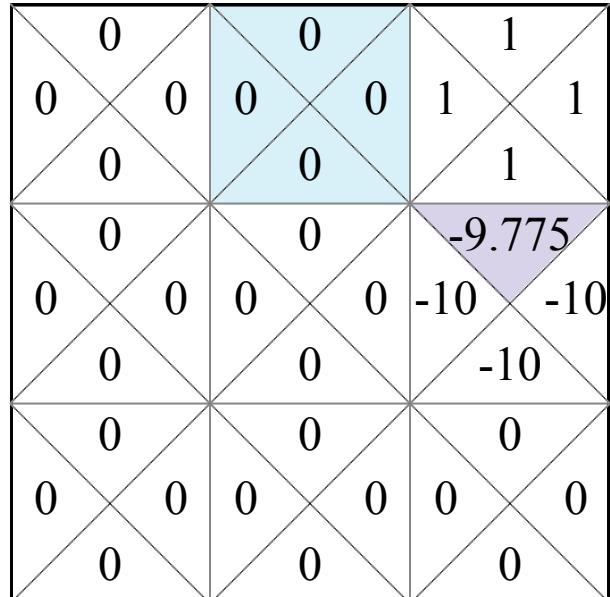
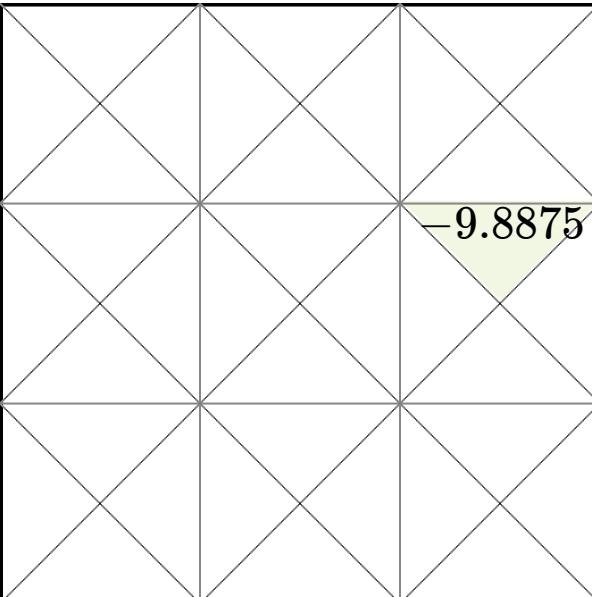
rewards now known



Q-learning update

e.g. pick $\alpha = 0.5$

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha \left(r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

 $Q_{\text{old}}(s, a)$  $Q_{\text{new}}(s, a)$ To update the estimate of $Q(6, \uparrow)$:

- execute $(6, \uparrow)$
- observe reward $r = -10$, next state $s' = 2$
- update $Q(6, \uparrow)$ as:

$$(1 - 0.5) * -9.775 + 0.5(-10 + 0.9 \max_{a'} Q_{\text{old}}(2, a'))$$

$$= -4.8875 + 0.5(-10 + 0) = -9.8875$$

Q-Learning ($\mathcal{S}, \mathcal{A}, \gamma, \alpha, s_0$ max-iter)

1. $i = 0$
2. **for** $s \in \mathcal{S}, a \in \mathcal{A}$:
3. $Q_{\text{old}}(s, a) = 0$
4. $s \leftarrow s_0$
5. **while** $i < \text{max-iter}$:
6. $a \leftarrow \text{select_action}(s, Q_{\text{old}}(s, a))$
7. $r, s' \leftarrow \text{execute}(a)$
8. $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$
9. $s \leftarrow s'$
10. $i \leftarrow (i + 1)$
11. $Q_{\text{old}} \leftarrow Q_{\text{new}}$
12. **return** Q_{new}

"learning"

- ϵ -greedy action selection strategy:
 - with probability ϵ , choose an action $a \in \mathcal{A}$ uniformly at random
 - with probability $1 - \epsilon$, choose $\arg \max_a Q_{\text{old}}(s, a)$

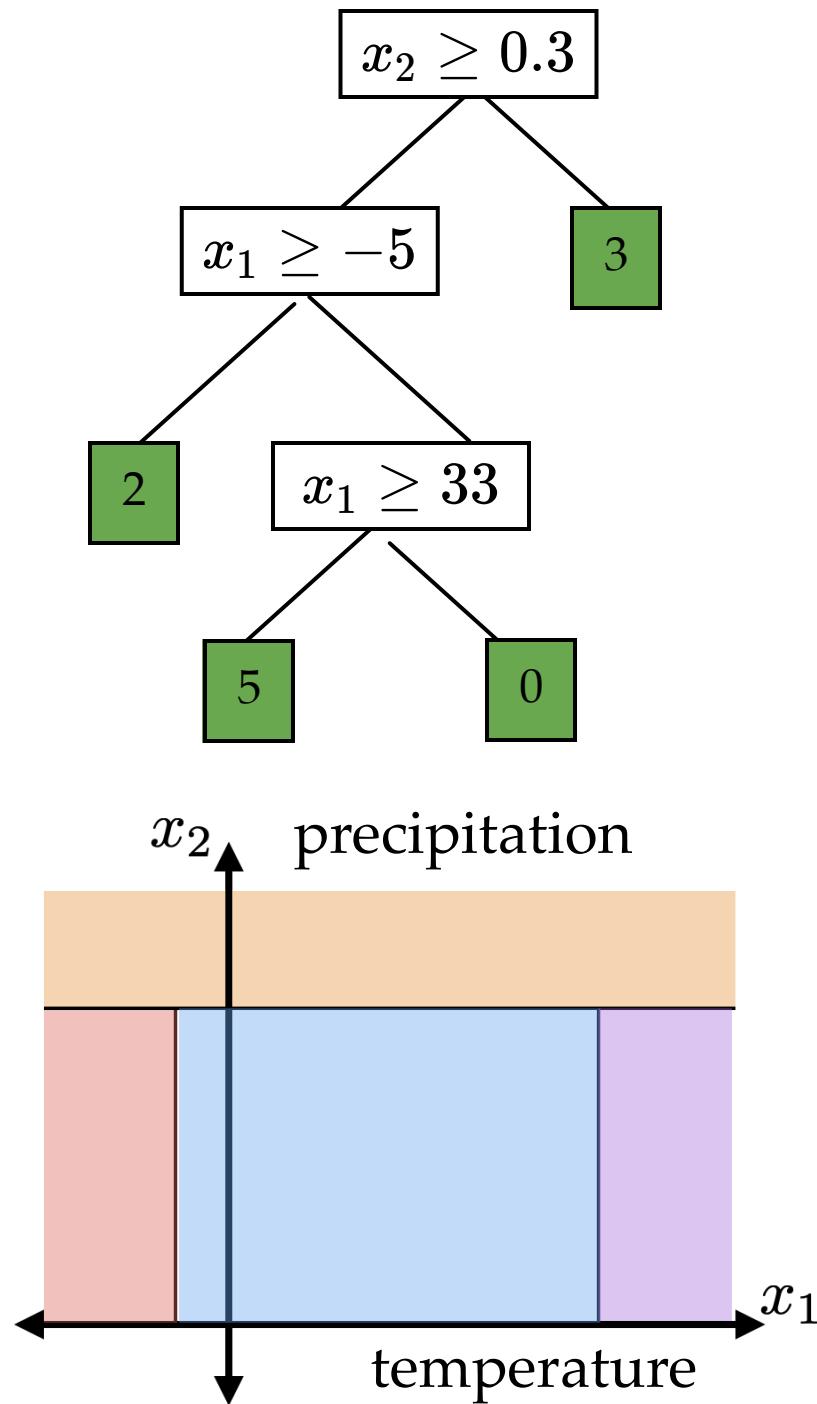


the current estimate of Q values

- ϵ controls the trade-off between *exploration vs. exploitation*.

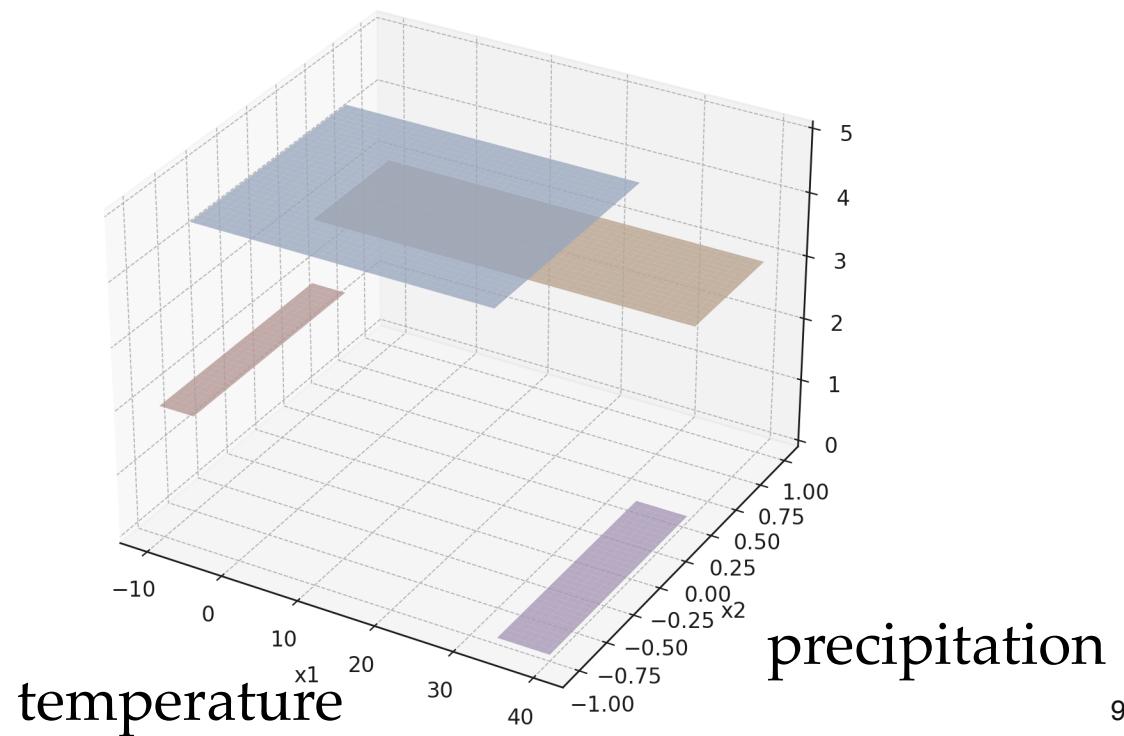
Week 12 - Non-parametric methods

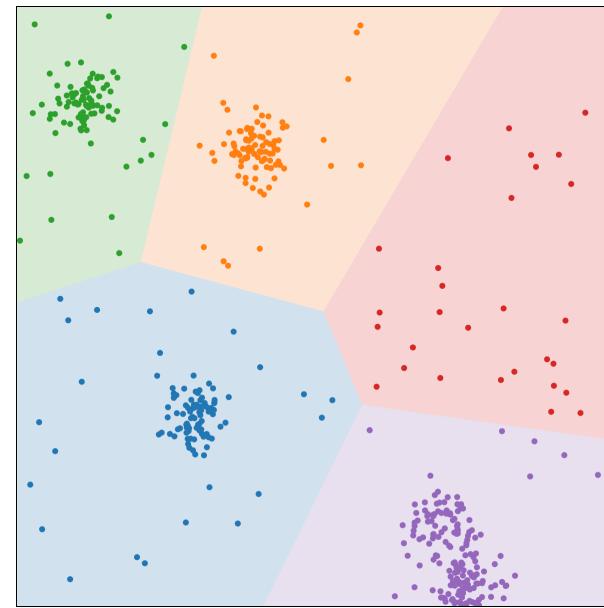
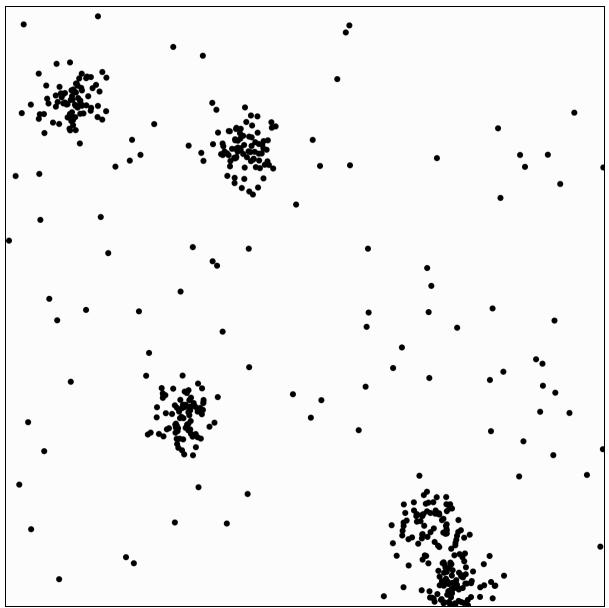
- Decision trees:
 - Split dimension, split value, tree structure (root/decision node and leaf), recursive algorithm, bagging, ensemble, largest leaf size k
- k -nearest neighbors:
 - memorizes data, scaling matters, k matters, inefficient in test/prediction time
- k -means clustering: (cluster assignment + cluster center updates), local convergence, initialization matters
- how choice of hyper-parameter k matters in each case



The same prediction applies to an axis-aligned ‘box’ or ‘volume’ in the feature space

label: km run





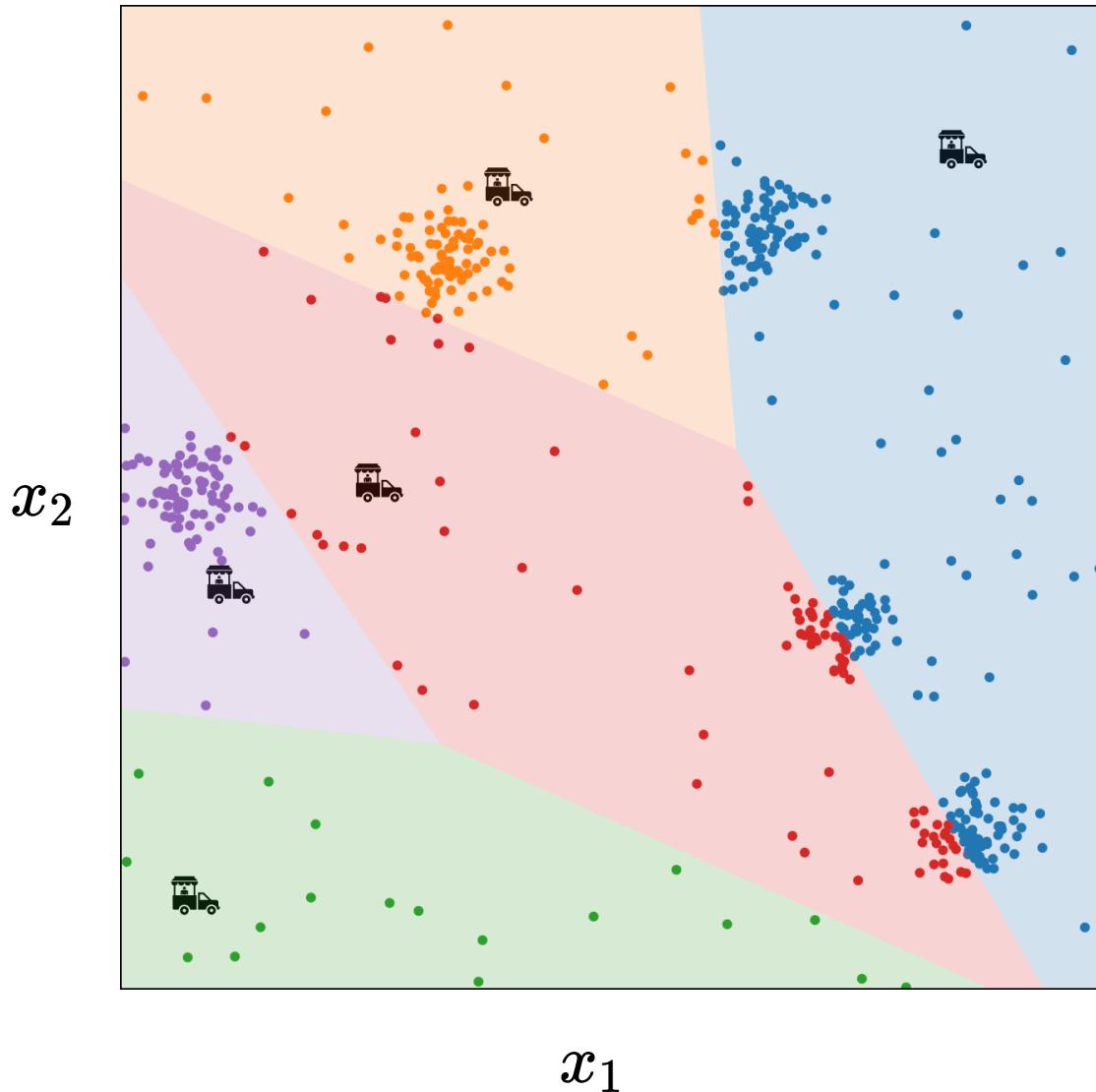
$$\{x^{(1)}\}$$

$$\{x^{(2)}\}$$

$$\{x^{(3)}\}$$

...

K-MEANS($k, \tau, \{x^{(i)}\}_{i=1}^n$)

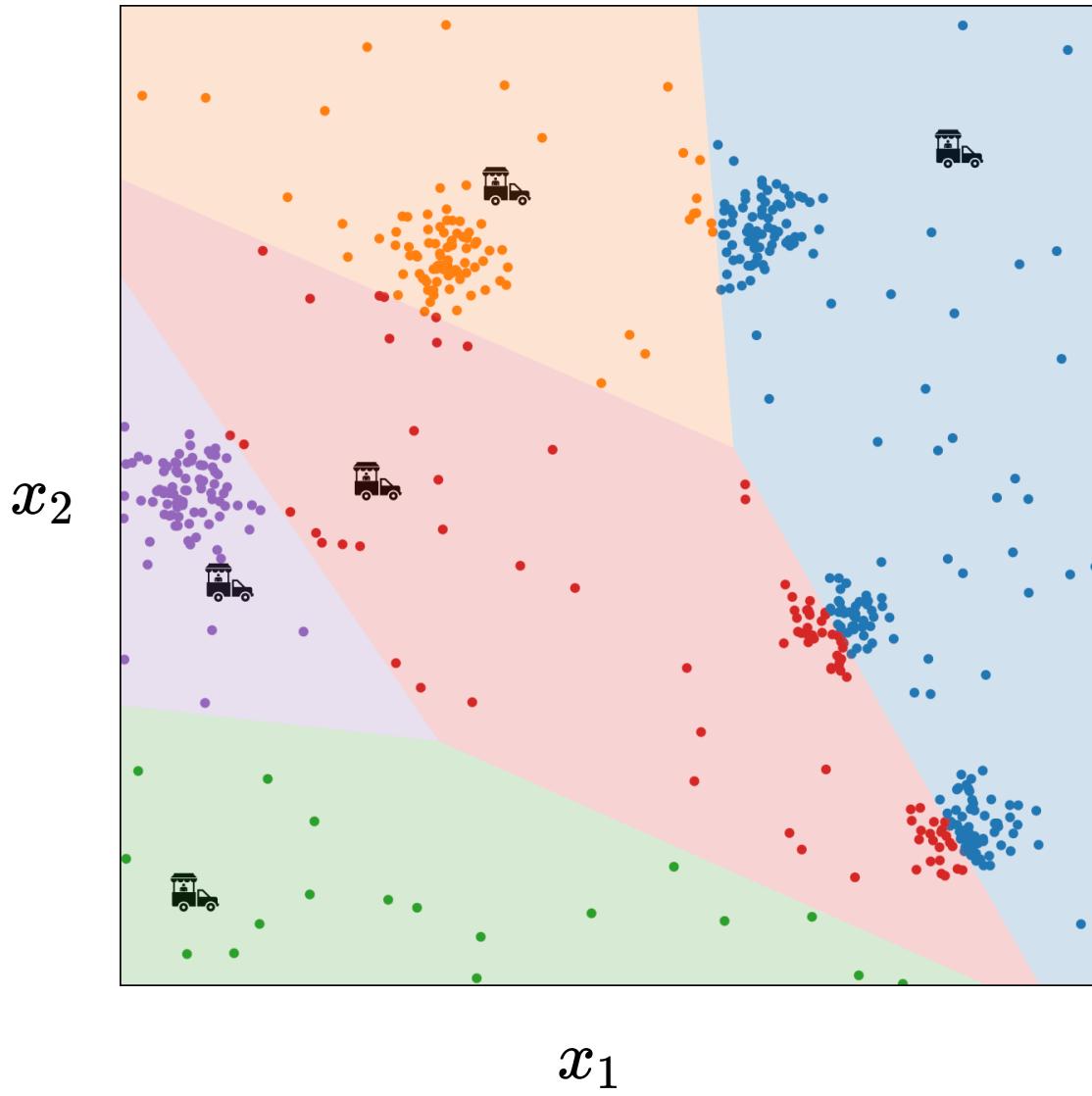


```
1  $\mu, y = \text{random initialization}$ 
2 for  $t = 1$  to  $\tau$ 
3    $y_{\text{old}} = y$ 
4   for  $i = 1$  to  $n$ 
5      $y^{(i)} = \arg \min_j \|x^{(i)} - \mu^{(j)}\|^2$ 
    ...

```

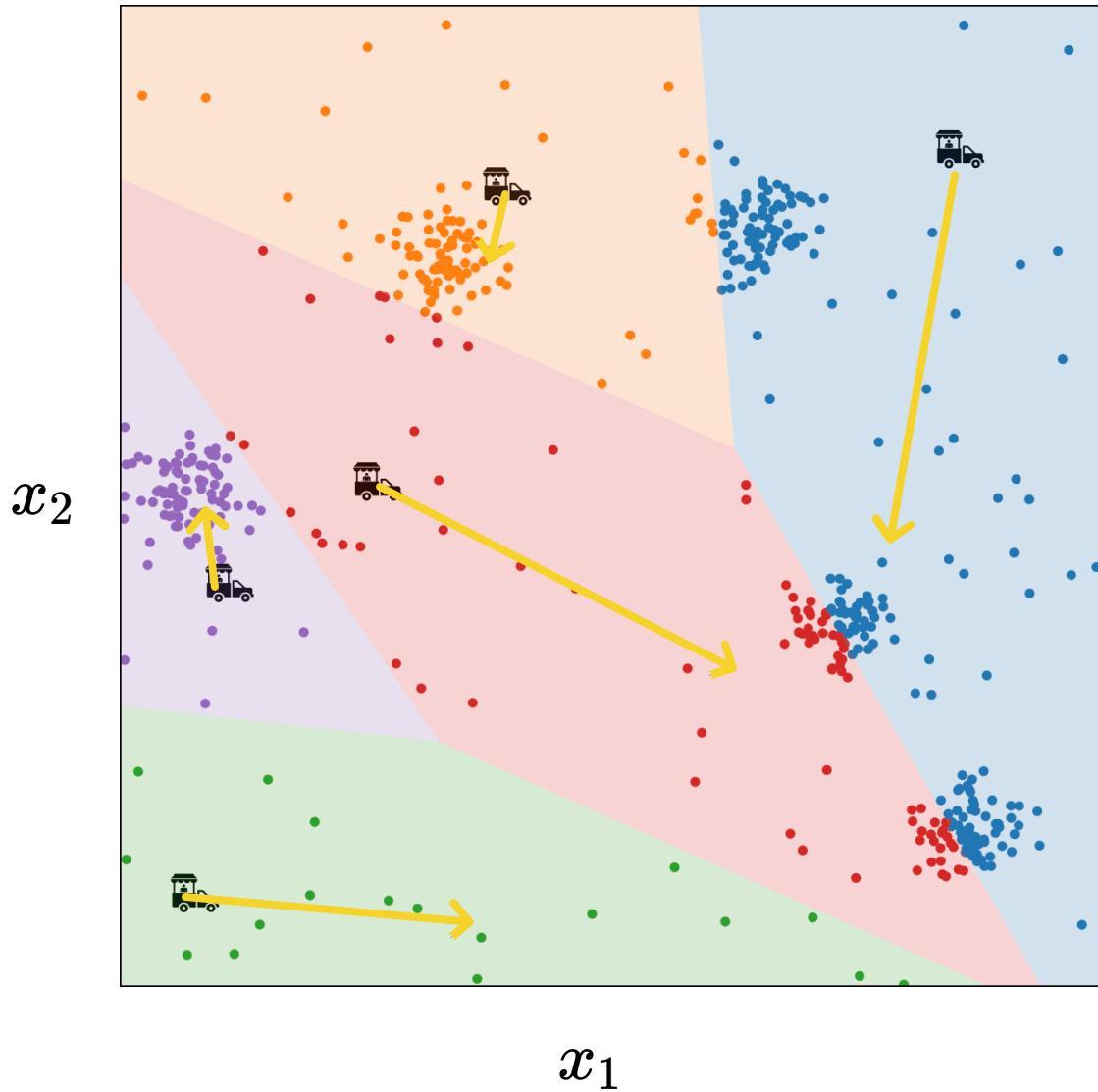
each person i gets assigned to food truck j , color-coded.

K-MEANS($k, \tau, \{x^{(i)}\}_{i=1}^n$)



```
1  $\mu, y = \text{random initialization}$ 
2 for  $t = 1$  to  $\tau$ 
3    $y_{\text{old}} = y$ 
4   for  $i = 1$  to  $n$ 
5      $y^{(i)} = \arg \min_j \|x^{(i)} - \mu^{(j)}\|^2$ 
6   for  $j = 1$  to  $k$ 
7      $\mu^{(j)} = \frac{1}{N_j} \sum_{i=1}^n \mathbf{1}(y^{(i)} = j) x^{(i)}$ 
8   if  $y == y_{\text{old}}$ 
9     break
10  return  $\mu, y$ 
```

K-MEANS($k, \tau, \{x^{(i)}\}_{i=1}^n$)

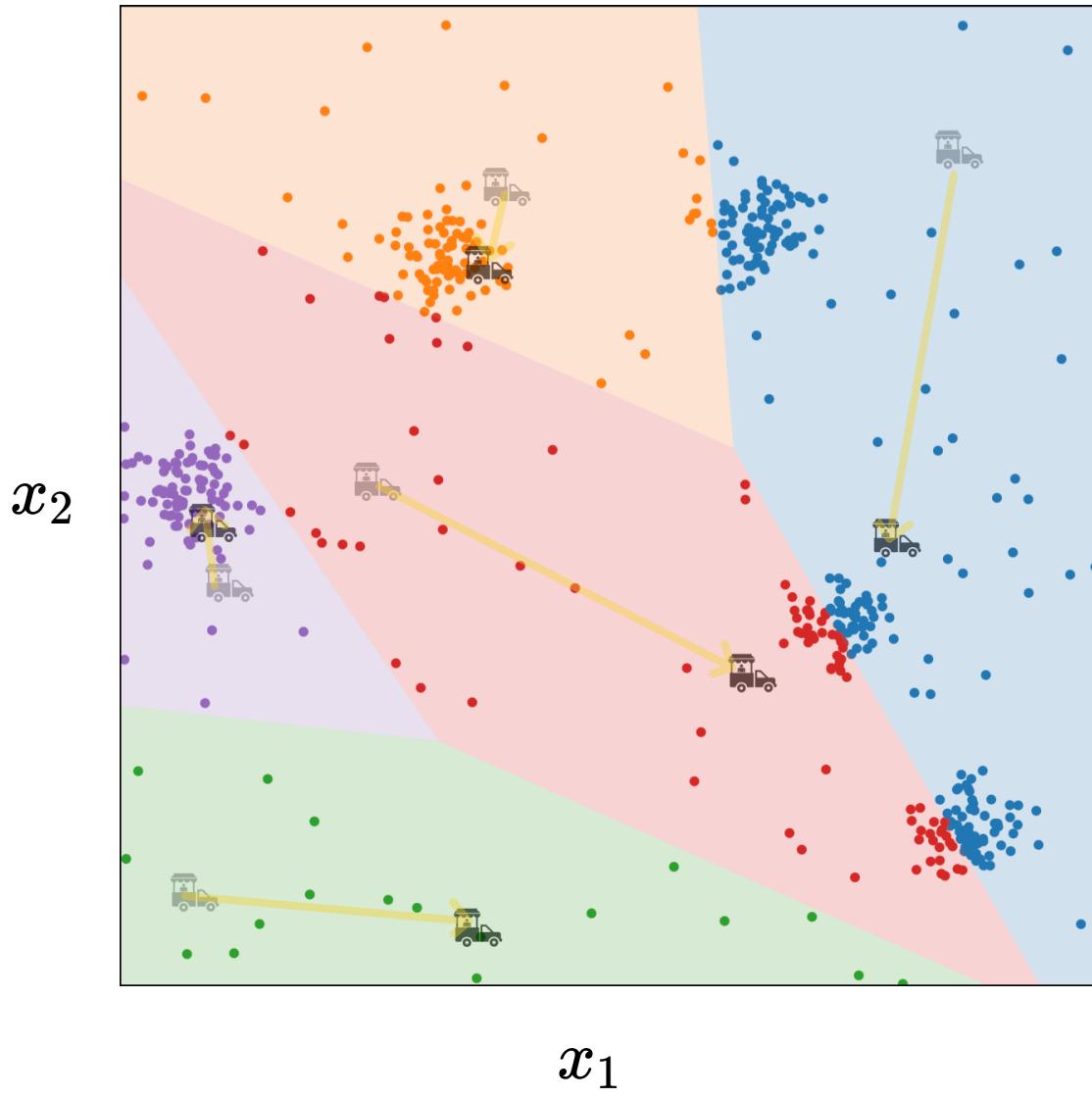


```

1  $\mu, y = \text{random initialization}$ 
2 for  $t = 1$  to  $\tau$ 
3    $y_{\text{old}} = y$ 
4   for  $i = 1$  to  $n$ 
5      $y^{(i)} = \arg \min_j \|x^{(i)} - \mu^{(j)}\|^2$ 
6   for  $j = 1$  to  $k$ 
7      $\mu^{(j)} = \frac{1}{N_j} \sum_{i=1}^n \mathbf{1}(y^{(i)} = j) x^{(i)}$ 
      ...
 $N_j = \sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}$ 
food truck  $j$  gets moved to the "central"
location of all ppl assigned to it

```

K-MEANS($k, \tau, \{x^{(i)}\}_{i=1}^n$)



```
1  $\mu, y = \text{random initialization}$ 
2 for  $t = 1$  to  $\tau$ 
3    $y_{\text{old}} = y$ 
4   for  $i = 1$  to  $n$ 
5      $y^{(i)} = \arg \min_j \|x^{(i)} - \mu^{(j)}\|^2$ 
6   for  $j = 1$  to  $k$ 
7      $\mu^{(j)} = \frac{1}{N_j} \sum_{i=1}^n \mathbf{1}(y^{(i)} = j) x^{(i)}$ 
8   if  $y == y_{\text{old}}$ 
9     break
10  return  $\mu, y$ 
```

Resources

- All course materials Week 1 - Week 12
- Past-exam Sampler

```
● ● ●

1 import random
2 terms= ["spring2024", "fall2023", "spring2023", "fall2022", "spring2022",
3         "fall2021", "fall2019", "fall2018", "fall2018"]
4
5 qunums = range(1,10)
6 base_URL = "https://introml.mit.edu/_static/fall24/final/review/final-"
7
8 term = random.choice(terms)
9 num = random.choice(qunums)
10 print("term:", term)
11 print("question number:", num)
12 print(f"Link: {base_URL+term}.pdf")
```

- Video walk-through of recent final exams
- Piazza and OHs

Course Evaluations

<https://registrar.mit.edu/subjectevaluation>

<https://eduapps.mit.edu/subjeval/studenthome.htm>

We'd love to hear your thoughts on the course: this provides valuable feedback for us and other students, for future semesters! Thank you! 

<https://www.youtube.com/embed/kRPAGejCAeY?enablejsapi=1>

<https://shenshen.mit.edu/tree>

(The demo won't embed in PDF. But the direct link below works.)

<https://shenshen.mit.edu/tree>

General problem-solving tips

Polya's Problem Solving Techniques

In 1945 George Polya published the book *How To Solve It* which quickly became his most prized publication. It sold over one million copies and has been translated into 17 languages. In this book he identifies four basic principles of problem solving.

Polya's First Principle: Understand the problem

This seems so obvious that it is often not even mentioned, yet students are often stymied in their efforts to solve problems simply because they don't understand it fully, or even in part. Polya taught teachers to ask students questions such as:

- Do you understand all the words used in stating the problem?
- What are you asked to find or show?
- Can you restate the problem in your own words?
- Can you think of a picture or diagram that might help you understand the problem?
- Is there enough information to enable you to find a solution?

Polya's Second Principle: Devise a plan

Polya mentions that there are many reasonable ways to solve problems. The skill at choosing an appropriate strategy is best learned by solving many problems. You will find choosing a strategy increasingly easy. A partial list of strategies is included:

- Guess and check
- Make an orderly list
- Eliminate possibilities
- Use symmetry
- Consider special cases
- Use direct reasoning
- Solve an equation
- Look for a pattern
- Draw a picture
- Solve a simpler problem
- Use a model
- Work backwards
- Use a formula
- Be ingenious

Polya's Third Principle: Carry out the plan

This step is usually easier than devising the plan. In general, all you need is care and patience, given that you have the necessary skills. Persist with the plan that you have chosen. If it continues not to work discard it and choose another. Don't be misled, this is how mathematics is done, even by professionals.

Polya's Fourth Principle: Look back

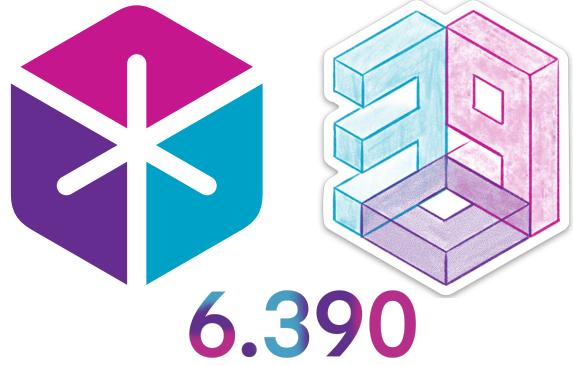
Polya mentions that much can be gained by taking the time to reflect and look back at what you have done, what worked, and what didn't. Doing this will enable you to predict what strategy to use to solve future problems.

Exam-taking tips

- Arrive 5min early to get settled in.
- Bring a pencil (and **eraser**), a watch, and some water.
- Look over whole exam and strategize for the order you do problems.



Good luck!



<https://introml.mit.edu/>

Thanks for the semester!