

## Question 2 - Scale Space / SIFT Detection (20 points)

The matrices in the left column are the output of applying Gaussian filters with different bandwidths for a single octave in the SIFT detection algorithm. There is a single sift keypoint in the scale space. Your job is to find it. (This is before removing the edges and low contrast points, and sub-pixel tuning). Report the x, y and scale of the key point. As a reference, for the highlighted cell at scale=2, the scale-space location is (x=3, y=1, scale=2).

To find the keypoint, you first need to build the Difference of Gaussian Images in the scale-space. The key points are found at the locations of extrema in scale-space as explained in the class. Fill in the Difference of Gaussian values and locate the key point. Why is this a SIFT key point?

**Hint:** The keypoints do not exist in the in the first and last scale.

The keypoint is located at  $x=1, y=1, \text{scale}=3$ , since

Compared to neighbours at scale 2  
 $-3 > 1, 2, 0, 1, 0, 1, 1$

Compared to neighbours at scale 4  
 $-3 > 2, 5, 0, 4, -2, 8, 0, 8, 2$

Compared to neighbours in the same scale (scale 3)  
 $-3 > 2, 3, 0, 1, 2, 0, 1$

Gaussian Filtered Images			
25	22	20	17
25	28	19	17
20	19	19	17
15	15	15	15

Difference of Gaussian Images			
scale = 1			
0	2	0	1
0	-2	0	1
1	2	0	1
2	2	1	1

25	20	20	16
25	30	19	16
19	17	19	16
13	13	14	14

scale = 2			
1	2	0	2
0	-2	0	1
1	1	0	2

24	18	20	14
25	32	19	15
18	16	19	14
12	12	13	13

scale = 3			
2	3	0	2
1	3	1	1
2	0	1	1

22	15	20	12
24	35	18	14
16	16	18	13
11	11	11	12

scale = 4			
2	5	0	4
4	-2	8	6
0	8	-2	3

20	37	10	8
16	8	20	10
12	10	9	10

scale = 5			
1	1	2	2
1	1	2	2
1	1	2	2

### TAYLOR SERIES

Quadratic approximation simplifies to

$$E(u, v) \approx [u \ v] \left( \sum w(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \right) \begin{bmatrix} u \\ v \end{bmatrix}$$

**Example**

$$\begin{aligned} & \text{Image } \xrightarrow{\text{Gradients } I_x, I_y} \begin{bmatrix} 1 & 2 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \\ & I_x^2: \begin{bmatrix} 1^2 & 2^2 & 1^2 \\ 2^2 & 2^2 & 1^2 \\ 1^2 & 1^2 & 1^2 \end{bmatrix} = \begin{bmatrix} 4 & 1 \\ 4 & 1 \\ 1 & 1 \end{bmatrix} \Rightarrow \sum = 14 \\ & I_y^2: \begin{bmatrix} 2^2 & 1^2 & 2^2 \\ 2^2 & 2^2 & 0^2 \\ 1^2 & 1^2 & 0^2 \end{bmatrix} = \begin{bmatrix} 4 & 1 & 4 \\ 1 & 4 & 0 \\ 1 & 0 & 1 \end{bmatrix} \Rightarrow \sum = 16 \\ & I_x I_y: \begin{bmatrix} 1 & 2 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 4 & 0 \\ 0 & 1 & 1 \end{bmatrix} \Rightarrow \sum = 16 \\ & E(1, 1) = [1 \ 1] \begin{bmatrix} 14 & 12 \\ 12 & 16 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 26 \\ 28 \end{bmatrix} \\ & E(1, 1) = \frac{1}{54} \begin{bmatrix} 14 * 1 + 12 * 1 \\ 12 * 1 + 16 * 1 \\ 1 * 26 + 1 * 28 \end{bmatrix} = \begin{bmatrix} 26 \\ 28 \end{bmatrix} \end{aligned}$$

$$1x26 + 1x28 = 54$$

### Magnitude - Angle Range (degrees)

1	$0^\circ - 20^\circ$
2	$20^\circ - 40^\circ$
3	$40^\circ - 60^\circ$
4	$60^\circ - 80^\circ$
5	$80^\circ - 100^\circ$
6	$100^\circ - 120^\circ$
7	$120^\circ - 140^\circ$
8	$140^\circ - 160^\circ$
9	$160^\circ - 180^\circ$

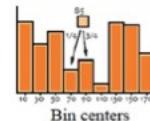
bin edges are  $0^\circ$  and  $60^\circ$ .

weight for lower bin

$$w_{\text{low}} = \frac{\theta_2 - \theta}{\theta_2 - \theta_1}$$

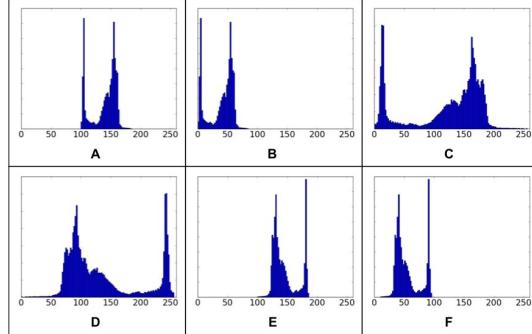
weight for higher bin

$$w_{\text{high}} = \frac{\theta - \theta_1}{\theta_2 - \theta_1}$$



### Question 4 - Histogram (18 points)

Look at the following histograms for images A, B, C, D, E and F, all images of the same scene. Zero and 255 represent black and white intensities respectively.



- A) Which of the above histograms correspond to the following image? Explain your answer. (6 points)

The histogram must have two peaks. A relatively narrow peak at darker intensities (camera man's coat and tripod's legs), and a wider peak mostly consisting of brighter intensities (the grass, the sky, the building in the background, etc.). This is because there are fewer dark pixels compared to the bright ones. Thus the answer is C.

- B) Which of the images D and E has a higher contrast? Why? (4 points)

D has a higher contrast since it has a wider histogram, covering a wider range of intensities.

- C) How does image B look like compared to image A? Explain. (4 points)

B is generally darker than A, as its histogram has been shifted to the left. But, both have about the same contrast.

- D) What is the relation between images C and D? Explain your answer. (4 points)

C is D with inverted intensities (negative of D). Mathematically C = 255 - D.

HOG

Instead of assigning gradient entirely to one bin, we use bilinear (fractional) interpolation

Means gradient contributes to two neighboring bins based on how close the angle is to them

For each pixel, we compute angle & magnitude

### Example A

Gradient magnitude M = 10

Gradient angle  $\theta = 50^\circ$

it lies closer to bin#3 and naturally the next bin would be bin#4

- [1] To compute weights:-

$$w_{\text{high}}: \frac{(60 - 50)}{(60 - 40)} = \frac{10}{20} = 0.5$$

$$w_{\text{low}}: \frac{(50 - 40)}{(60 - 40)} = \frac{10}{20} = 0.5$$

### Example B

Gradient magnitude M = 8

Gradient angle  $\theta = 47^\circ$

it lies closer to bin#3

and naturally the next bin would be bin#4

- [1]

$$w_{\text{high}}: \frac{60 - 47}{60 - 40} = \frac{13}{20} = 0.65$$

- [2]

To distribute magnitude:-

$$\text{Bin3} \rightarrow 0.65 \times 8 = 5.2$$

$$\text{Bin4} \rightarrow 0.35 \times 8 = 2.8$$

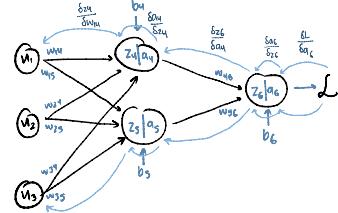
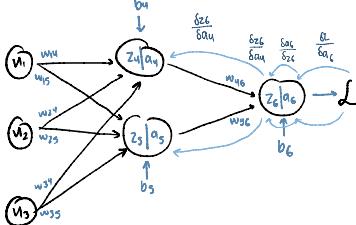
## Question - Backpropagation squared

$$L = \frac{1}{2} (y - a_6)^2$$

$$\frac{\delta L}{\delta a_6} = \frac{1}{2} [2(-1)(y - a_6)]$$

$$= -(y - a_6)$$

$$+ a_6 - y$$



$$Z_4 = W_{11}v_1 + W_{12}v_2 + W_{13}v_3 + b_4$$

$$a_4 = \text{Sigmoid}(z_4)$$

$$Z_5 = W_{21}a_4 + W_{22}a_5 + b_5$$

$$a_5 = \text{Sigmoid}(z_5)$$

$$Z_6 = W_{31}a_4 + W_{32}a_5 + b_6$$

$$a_6 = \text{Sigmoid}(z_6)$$

$$L = \frac{1}{2} (y - a_6)^2$$

$$\frac{\delta L}{\delta w_{46}} = \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta w_{46}}$$

$$= -(y - a_6) a_6 (1 - a_6) a_4$$

$$L = \frac{1}{2} (y - a_6)^2$$

$$\frac{\delta L}{\delta w_{56}} = \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta w_{56}}$$

$$= -(y - a_6) a_6 (1 - a_6) a_5$$

$$L = \frac{1}{2} (y - a_6)^2$$

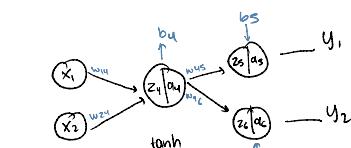
$$\frac{\delta L}{\delta b_6} = \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta b_6}$$

$$= -(y - a_6) a_6 (1 - a_6)$$

$$w_{46}^t = w_{46} - \eta \frac{\delta L}{\delta w_{46}}$$

$$w_{56}^t = w_{56} - \eta \frac{\delta L}{\delta w_{56}}$$

$$b_6^t = b_6 - \eta \frac{\delta L}{\delta b_6}$$



## Feed Forward

$$Z_4 = v_1 w_{11} + v_2 w_{12} + b_4 \rightarrow a_4 = \tanh(z_4)$$

$$Z_5 = a_4 w_{21} + b_5 \rightarrow a_5 = \sigma(z_5) = \frac{1}{1 + e^{-z_5}}$$

$$Z_6 = a_4 w_{31} + b_6 \rightarrow a_6 = \sigma(z_6) = \frac{1}{1 + e^{-z_6}}$$

## Backward

$$L = \text{Cross Entropy loss}$$

$$= -[y \log a_6 + (1-y) \log(1-a_6)]$$

$$L = -[y \log a_6 + (1-y) \log(1-a_6)] - [y \log a_6 + (1-y) \log(1-a_6)]$$

$$\frac{\delta L}{\delta s} = \frac{-y}{a_6} + \frac{1-y}{1-a_6} \quad \frac{\delta L}{\delta a_5} = \frac{-y}{a_6} + \frac{1-y}{1-a_6}$$

$$\frac{\delta L}{\delta w_{45}} = \frac{\delta L}{\delta a_5} \frac{\delta a_5}{\delta z_5} \frac{\delta z_5}{\delta w_{45}}$$

$$= \left[ -\frac{y}{a_6} + \frac{1-y}{1-a_6} \right] \left[ a_5 (1-a_5) \right] a_4$$

$$= -\frac{(y)(1-a_5) + a_5(1-y)}{a_6(1-a_6)} \left[ a_5 (1-a_5) \right] a_4$$

$$= -y \cdot a_5 a_6 + a_5 - y a_6$$

$$= a_5 (1 - y) a_6$$

$$\frac{\delta L}{\delta w_{45}} = (a_5 - y) a_4 \quad \frac{\delta L}{\delta w_{46}} = (a_6 - y) a_4$$

$$\frac{\delta L}{\delta b_5} = a_5 \cdot y_1 \quad \frac{\delta L}{\delta b_6} = a_6 \cdot y_2$$

$$H_1 = \frac{\delta L}{\delta y} = \frac{\delta L}{\delta a_5} \frac{\delta a_5}{\delta z_5} \frac{\delta z_5}{\delta b_5} + \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta b_6}$$

$$\frac{\delta L}{\delta y} = \frac{\delta L}{\delta a_5} \frac{\delta a_5}{\delta z_5} \frac{\delta z_5}{\delta w_{45}} + \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta w_{46}}$$

$$= (a_5 - y) w_{45} (1 - a_5 a_6) w_{46} \quad (1 - a_4 a_5)^2 \cdot y_1$$

$$\frac{\delta L}{\delta b_4} = \frac{\delta L}{\delta a_5} \frac{\delta a_5}{\delta z_5} \frac{\delta z_5}{\delta w_{45}} + \frac{\delta L}{\delta a_6} \frac{\delta a_6}{\delta z_6} \frac{\delta z_6}{\delta w_{46}}$$

$$= (a_5 - y) w_{45} (1 - a_5 a_6) w_{46} \quad (1 - a_4 a_5)^2 \cdot y_2$$

A) First, compute the derivatives using the differentiation kernels shown above. No normalization (division by 2) is needed. (5 points).

$$I_x = d/dx$$

x	x	x	x	x
x	<b>4</b>	<b>7</b>	<b>6</b>	x
x	<b>8</b>	<b>8</b>	<b>7</b>	x
x	<b>8</b>	<b>6</b>	<b>5</b>	x
x	x	x	x	x

$$I_y = d/dy$$

x	x	x	x	x
x	<b>4</b>	<b>8</b>	<b>8</b>	x
x	<b>8</b>	<b>6</b>	<b>7</b>	x
x	<b>6</b>	<b>6</b>	<b>4</b>	x
x	x	x	x	x

B) Now compute the Harris Matrix based on the derivative matrices. (10 points).

$$H = \sum_{(x,y) \in W} \begin{bmatrix} I_x(x,y)^2 & I_x(x,y) I_y(x,y) \\ I_x(x,y) I_y(x,y) & I_y(x,y)^2 \end{bmatrix}$$

$$\sum_{(x,y) \in W} I_x(x,y)^2 = 4^2 + 7^2 + 6^2 + 8^2 + 8^2 + 7^2 + 8^2 + 6^2 + 5^2 = 403$$

$$\sum_{(x,y) \in W} I_y(x,y)^2 = 4^2 + 8^2 + 8^2 + 6^2 + 7^2 + 6^2 + 4^2 = 381$$

$$\sum_{(x,y) \in W} I_x(x,y) I_y(x,y) = 4 * 4 + 7 * 8 + 6 * 8 + 8 * 8 + 8 * 6 + 7 * 7 + 8 * 6 + 6 * 6 + 5 * 4 = 385$$

$$H = [403 \ 385] \ [385 \ 381]$$

C) Compute the Harris cornerness score  $C = \det(H) - k \text{trace}(H)^2$  for  $k = 0.04$ . What do we have here? A corner? An edge? Or a flat area? Why? (5 points)

$$C = \det(H) - k \text{trace}(H)^2 = 5318 - 0.04 * (784)^2 = -19268.24$$

A negative Harris score indicates an edge.

## Question 2 - Scale Space / SIFT Detection (20 points)

The matrices in the left column are the output of applying Gaussian filters with different bandwidths for a single octave in the SIFT detection algorithm. There is a **single** sift keypoint in the scale space. Your job is to find it. (This is **before** removing the edges and low contrast points, and sub-pixel tuning). Report the x, y and scale of the key point. As a reference, for the highlighted cell at scale=2, the scale-space location is (x=3, y=1, scale=2).

To find the keypoint, you first need to build the Difference of Gaussian Images in the scale-space. The key points are found at the locations of extrema in scale-space as explained in the class. Fill in the Difference of Gaussian values and locate the key point. Why is this a SIFT key point?

**Hint:** The keypoints do not exist in the in the first and last scale.

The keypoint is located at x=1, y=1, scale=3, since

Compared to neighbours at scale 2  
 $-3 > 1, 2, 0, 1, 1, 0, 1, 1, 1$

Compared to neighbours at scale 4  
 $-3 > 2, 5, 0, 4, -2, 8, 0, 8, -2$

Compared to neighbours in the same scale (scale 3)  
 $-3 > 2, 3, 0, 1, 1, 2, 0, 1$

Gaussian Filtered Images

25	22	20	17
25	28	19	17
20	19	19	17
15	15	15	15

Difference of Gaussian Images

scale = 1

0	2	0	1
0	-2	0	1
1	2	0	1
2	2	1	1

I picked 3 in DoG3 at (x=1, y=1) because it is significantly larger than surrounding values.

It's in the middle of the matrix (not on the edge).

It's in a valid scale (not 1st or last).

It has a strong chance to be a local maximum in 3D space.

→ 9

scale = 2

1	2	0	2
0	-2	0	1
1	1	0	2
1	1	1	1

→ 8

scale = 3

2	3	0	2
1	-3	1	1
2	0	1	1
1	1	2	1

→ 9

scale = 4

2	5	0	4
4	-2	8	6
0	8	-2	3
-1	1	2	2

24	18	20	14
25	32	19	15
18	16	19	14
12	12	13	13

22	15	20	12
24	35	18	14
16	16	18	13
11	11	11	12

20	10	20	8
20	37	10	8
16	8	20	10
12	10	9	10

# applying Gaussian

4. The matrices in the left column below, are the output of applying Gaussian filters with different bandwidths for a single octave in the SIFT detection algorithm. There is a single sift ~~keypoint~~ in the scale space. Your job is to find it. (This is before removing the edges and low contrast points, and sub-pixel tuning). Report the x, y and scale of the key point. As a reference, for the highlighted cell at scale=2, the scale-space location is (x=2, y=1, scale=2).

To find the ~~keypoint~~, you first need to build the Difference of Gaussian Images in the scale-space. The key points are found at the locations of extrema in scale-space. Fill in the Difference of Gaussian values and locate the key point. Why is this a SIFT key point?

*Hint: The ~~keypoints~~ do not exist in the first and last scale.*

Gaussian Filtered Images

1	25	22	20	17
	25	28	19	17
	20	19	19	17
	15	15	15	15

2	25	20	20	16
	25	30	19	16
	19	17	19	16
	13	13	14	14

3	24	18	20	14
	25	32	19	15
	18	16	19	14
	12	12	13	13

4	22	15	20	12
	24	35	18	14
	16	16	18	13
	11	11	11	12

5	20	10	20	8
	20	37	10	8
	16	8	20	10
	12	10	9	10

Difference of Gaussian Images

scale = 1

0	2	0	1
0	-2	0	1
1	2	0	1
2	2	1	1

$$\text{DOG} : G_1 - G_2$$

0	1	2	0	2
0	-2	0	1	
1	1	0	2	
2	1	1	1	

$$\text{DOG} : G_2 - G_3$$

2	3	0	2
1	-3	1	1
2	0	1	1
1	1	2	1

$$\text{DOG} : G_3 - G_4$$

2	5	0	4
4	-2	8	6
0	8	-2	3
-1	1	2	2

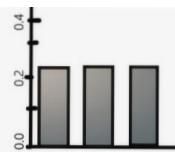
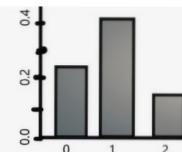
$$\text{DOG} : G_4 - G_5$$

ishma hafeez

notes

repst  
sheet

## KL Divergence



The above is an example consisting of the binomial distribution and uniform probability distribution, The binomial is considered as P and the uniform distribution is considered as Q. One can see for the binomial distribution it has got you may see 3 events 0,1 and 2 and having the probability of occurrence of the divine points 0.36,0.48 and 0.16 respectively, the other one is the uniform distribution so the probability is the same 0.333. So now let us see how can we compute the KL-Divergence for this, the computation is as follows:

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in X} P(x) \ln \left( \frac{P(x)}{Q(x)} \right)$$

$$= 0.36 \ln \left( \frac{0.36}{0.333} \right) + 0.48 \ln \left( \frac{0.48}{0.333} \right) + 0.16 \ln \left( \frac{0.16}{0.333} \right)$$

$$= 0.0852996$$

$$D_{\text{KL}}(Q \parallel P) = \sum_{x \in X} Q(x) \ln \left( \frac{Q(x)}{P(x)} \right)$$

$$= 0.333 \ln \left( \frac{0.333}{0.36} \right) + 0.333 \ln \left( \frac{0.333}{0.48} \right) + 0.333 \ln \left( \frac{0.333}{0.16} \right)$$

$$= 0.097455$$

A) Describe each of the three major stages of the SIFT algorithm. You just have to state their purpose, not the steps taken in each of them. (6 points)

a) SIFT detection

Detects the locations of the keypoints (points that can be identified in different images) and their associated scale.

b) SIFT description

For each keypoint finds a descriptor (a feature vector of size 128) based on how image looks like in a local neighbourhood around the keypoint. The descriptor is sought to be scale- and rotation-invariant.

c) SIFT matching

Tries to find matches between the keypoints in two (or more) images based on their descriptors.

B) Which of the above steps can be done independently for each image? Which ones require both images? (3 points)

a, b) can be done independently for each image.

c) requires both images

C) RANSAC can be used to improve which of the above stages? How? (3 points)

RANSAC can improve the matching stage. It can be applied when there is a global model relating the keypoints in one image to their true matches in the other image.

RANSAC improves matching by finding this model and removing the outliers at the same time.

D) Orientation assignment is done in which of the above stages? What is its purpose? (3 points)

It can be considered as the first step in the description stage, or alternatively, last part of the detection stage. Or it can be considered as a separate stage. All these answers are acceptable. The main purpose is to make the SIFT descriptor invariant to orientation. But it has other applications as well

## KL Divergence

### KL-Divergence.

- (a) Measures of how one probability distribution
- (b) diverges from a second distribution  $P$  (usually true)
- (c)  $\Theta$  is a subset of  $P$ .

Given two distributions, KL-Divergence is -

$$KL(\Theta(z) || P(z)) = \sum_z \Theta(z) \log \left( \frac{\Theta(z)}{P(z)} \right)$$

Approx: posterior  
Prior:  
Used in numerical.

### Properties:

$$KL(\Theta || P) \geq 0. \quad (\text{Non negativity})$$

$$KL(\Theta || P) \neq KL(P || \Theta) \quad (\text{Asymmetric})$$

Alternate form using joint Distribution.

$$KL(\Theta(z) || P(z)) = -\sum_z \Theta(z) \log \left( \frac{P(z)}{\Theta(z)} \right)$$

$$KL(\Theta(z) || P(z)) = -\sum_z \Theta(z) \log \left( \frac{P(x, z)}{P(x)\Theta(z)} \cdot \frac{P(x)}{P(z)} \right)$$

$$KL(\Theta(z) || P(z)) = \sum_z \Theta(z) \left[ \log \left( \frac{P(x, z)}{\Theta(z)} \right) + \log \frac{P(x)}{P(z)} \right]$$

$$= -\sum_z \Theta(z) \log \left( \frac{P(x, z)}{\Theta(z)} \right) + \log P(x)$$

$$\log P(x) = \mathbb{E}_{z \sim \Theta} \left[ \log \left( P(x, z) \right) \right]$$

$$\log P(x) = \mathbb{E}_{z \sim \Theta} \left[ \log \left( \frac{P(x, z)}{\Theta(z)} \right) \right] + KL(\Theta(z) || P(z))$$

C) Evidence Lower Bound (ELBO) used in variational inference.

[Measures the distance from  $\Theta$  to  $P$ .]

Example:  $\Theta$  smaller means closer.

$$P = [0.5, 0.4, 0.1] \quad (\text{True Distribution})$$

$$\Theta = [0.4, 0.4, 0.2]$$

$$KL(\Theta || P) = \sum_{x \in P, \Theta} \Theta(x) \cdot \log \left( \frac{\Theta(x)}{P(x)} \right)$$

$$= 0.4 \cdot \log_2 \left( \frac{0.4}{0.5} \right) + 0.4 \cdot \log_2 \left( \frac{0.4}{0.4} \right) + 0.2 \cdot \log_2 \left( \frac{0.2}{0.1} \right)$$

$$\approx 0.4 \left[ \frac{\log_2(0.4)}{\log_2(2)} - \frac{\log_2(0.5)}{\log_2(2)} \right] + 0 + 0.2$$

$$\approx 0.1288 + 0^2 \approx 0.0712$$

## SIFT invariance

Explain how rotation invariance is achieved in SIFT using the following scenario:

You have calculated the dominant edge angle of a particular cell, which is 70 degrees. The size of the cell is 6x6. In the next frame, the image is rotated by 25 degrees. Demonstrate that rotation invariance is preserved in the rotated image at pixel locations A and B, where the edge angle at A is 15 degrees, and at B, the edge angle is 35 degrees.



Before:

dominant angle = 70°

rotation = 25°

Before rotation

at A → 15°, B = 35°

Normalized angle at A' → 1170 - 1511 = 55°

" " " " B' → 1170 - 3511 = 35°

Rotational Invariance Preserved

After rotation

at A → 15° + 25° = 40°

at B → 35° + 25° = 60°

Normalized angle also becomes 70 + 25 - 95

at A' → 1170 - 4011 = 55°

" " " " B' → 1170 - 6011 = 35°

Normalized angle at A' → 1170 - 4011 = 55°

" " " " B' → 1170 - 6011 = 35°

# Generative Models

**Supervised learning:** Model is trained on example  $x$  paired with output  $y$ . used in Classification, regression, semantic segmentation, object detection.

**Unsupervised learning:** Model is trained on example  $x$  without  $y$  label, used to learn underlying structure of data. Used in Clustering, Dimensionality reduction.

**Semi Supervised learning:** Model is trained on some examples of  $x$  paired with output  $y + a$  large amount of unlabeled data.

**Self-Supervised learning:** Model creates labels from the data itself  $\Leftrightarrow$  **Weakly Supervised learning:** Training on noisy, incomplete/imprecise labels

**Generative modeling:** Process of learning the underlying distribution of data so the model can generate new, similar data samples. We try to learn by making  $P_{model}(x)$  super close to  $P_{data}(x)$ .

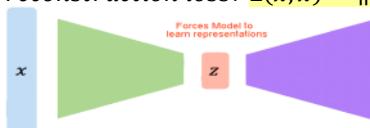
**Use cases:** [1] can be used for outlier detection [2] Allows sampling from data distribution

**Latent variable:** is a n-dimensional datapoint that is not directly observed, but rather inferred from the data through a model.

If you wanted to describe your appearance to someone and the other person doesn't know what you look like, you wouldn't start stating color of pixel 1, pixel 2, ..., pixel n. you would describe high level features that describe a group of pixels such as I have curly hair and green eyes. In other words, each point in the latent space is a representation of some high-dimensional observation

**Autoencoders:** were solely created to address unsupervised representation learning and data compression.

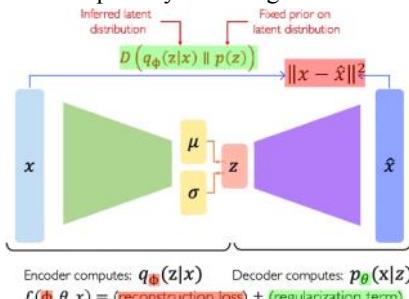
**reconstruction loss:**  $L(x, \hat{x}) = \|x - \hat{x}\|^2$



After the Model is trained, decoder is discarded and encoder is used to generate representations for most usecases

**Use cases:** [1] Image denoising by learning to reconstruct the image [2] Non-linear PCA, latent space capture most important features [3] Text retrieval by mapping same text to same representation [4] Similar image search by mapping similar images to same representation

**Variational Autoencoders:** learns probabilistic mapping from data to a latent space by modeling the latent variables as distributions

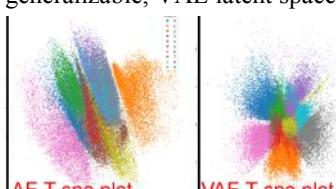


**KL Divergence:** Gaussian is a common Prior distribution choice for latent variables, Gaussian is imposed using a KL Divergence Regularization term;  $\oplus$  KL Term prevents overfitting  $\oplus$  Ensures latent space approximates the Gaussian properly  $\oplus$  Forces the latent space to be continuous and smooth  $\oplus$  Ensures at sampling time the latent space is well formed and capable of generating diverse outputs

Regularization term is actually KL Divergence

$$D(q_\phi(z|x) \parallel p(z)) = -\frac{1}{2} \sum_{j=1}^{k-1} (\sigma_j + \mu_j^2 - 1 - \log \sigma_j)$$

**Why use VAE over AE:** VAE can generate new data unlike AE, latent space is structured allowing better sampling, VAE is more generalizable, VAE latent space is more interpretable



**VAE overlapping t-sne benefit?**

- a. Model learns common features shared by all classes, can generalize better
- b. provides smooth transition between one-to-another class
- c. aesi hi tu now images aengi, beech mese sampling karli

**Reparameterization technique:** In VAEs we used to sample the latent vector before passing it to decoder, encoder outputs two values: Mean and variance of  $x$ . Then latent vector was sampled as:  $z = N(\mu(x), \sigma^2(x))$   $N$ : normal distribution

This direct sampling is non-differentiable. We reparameterized the

latent variable as:  $z_i = \mu(x_i) + \sigma(x_i) \cdot \epsilon_i$   $\epsilon_i$ : random noise  $N(0,1)$

Now randomness is moved outside the network making the sampling process differentiable wrt mean and variance.

**Stochastic Gradient Descent:** When the dataset is too large, performing gradient descent with respect to all data points can be computationally expensive and slow. SGD samples a single data point to compute the gradient to update the model.

**Process:** [1] Shuffle the dataset to avoid bias in training. [2] pick data point  $i$  (or mini batch) [3] Calculate error and update weights of whole network based on point  $i$  (or mini batch)

**Epoch:** one complete pass through entire training dataset

**Iteration:** number of iterations per epoch is determined by the batch size and total number of training samples  $t - \gamma \frac{\partial L}{\partial \theta}$

$$\text{With momentum: } \theta_{t+1} = \theta_t - \gamma \left( \alpha \frac{\partial L}{\partial \theta} + \beta \theta_t \right)$$

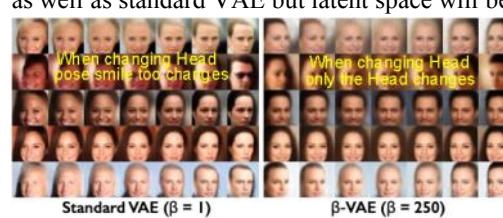
**Latent Perturbation:** slowly changing only one dimension of latent variable to explore effect on the generated output

**Disentanglement:** process of learning Disentangled latent representation where each latent variable encodes a distinct, independent, meaningful variation in the data. (e.g. one latent variable can correspond to head rotation, another to smile another to lightning)

**B-VAE (Beta-VAE):** encourages the model to learn disentangled latent representations.

$$L_{\text{total}} = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \beta \cdot D_{\text{KL}}[q(z|x)||p(z)]$$

$\beta$  is a hyperparameter which can modify the intensity of KL divergence term,  $\beta \geq 1$  usually and enforces independence and disentanglement KL divergence measures how much posterior distribution diverges from prior distribution, increasing  $\beta$  encourages the model to sway away the focus from strictly reconstructing and model now focuses more on the posterior to match prior more strictly forcing the latent space to follow a distribution that is more uniform.  $\beta$ -VAEs might not reconstruct input data as well as standard VAE but latent space will be more disentangled



We want  $q(z|x)$  to be as close as possible to the true posterior  $p(z|x)$ , measured using Kullback Leibler (KL) divergence:

$$\text{KL}(q(z|x) \parallel p(z|x)) = \mathbb{E}_{q(z|x)} \left[ \log \frac{q(z|x)}{p(z|x)} \right]$$

we want the posterior  $p(z|x)$

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

$x$ : observed data

$z$ : latent (hidden) variables

$p(x|z)$ : likelihood

$p(z)$ : prior

$$p(x) = \int p(x|z)p(z) dz$$

We can't directly minimize the KL divergence because  $p(z|x)$  contains  $p(x)$ , which is intractable. So we use the Evidence Lower Bound (ELBO) instead:

$$\log p(x) = \text{ELBO}(x) + \text{KL}(q(z|x) \parallel p(z|x))$$

$$\text{ELBO}(x) = \mathbb{E}_{q(z|x)} [\log p(x, z) - \log q(z|x)]$$

if we maximize ELBO, we indirectly minimize the KL divergence and bring  $q(z|x)$  closer to the true posterior.

Neurons on a surface level have local connectivity (connected to a small region of input image) but are depth wise global (connected to each channel value)

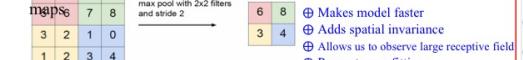
**Local connectivity:** neuron is connected to small, local region of the input, not the entire input

**Spatial connectivity:** how neighboring pixels in image are connected helps identify shapes, edges

**Receptive field:** region of the input image that affects output neuron

**Depth map:** represents distance from camera, used in 3D vision

**Pooling:** down sampling technique to reduce spatial size of feature maps



**Saddle points:** gradient = 0 but not a minima, optimization gets stuck

**Anisotropic loss surfaces:** loss changes fast in one direction, slow in another

$$\text{SGD equation: } x_{t+1} = x_t - \alpha \cdot \nabla f(x_t)$$

$\oplus$  Faster updates  $\ominus$  Unstable training  $\ominus$  Can stuck in saddle points or local minima

**D**  $\oplus$  Large gradient change  $\ominus$  Small gradient change

Jittery updates due to large change in y and small change in x

Our gradients come from minibatches so they can be noisy

Hessian matrix has a high condition number

If gradients keep pointing in the same direction, inertia increases

$$v_{t+1} = p v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - a \cdot v_{t+1}$$

**SGD+Momentum:** If gradients keep pointing in the same direction, inertia increases

$$v_{t+1} = p v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - a \cdot v_{t+1}$$

**Adan**

**Drop**

**Types of backprop time**

**a. Online/stochastic Gradient descent:** updates weight after every training sample, good for streaming data, noisy updates

**b. Mini-batch Gradient descent:** updates weight after every small batch, allows for training large data on small GPUs

**c. Batch Gradient descent:** computes gradients over entire training set, slow and memory intensive

**NN Train**

1. Define

2. Define

Why do VAEs typically produce blurry images?

VAEs learn an explicit distribution of the data by fitting it into a multivariate Gaussian, and the output is blurry because of the conditional independence assumption of the samples given latent variables. This is not true for most realistic distributions, like natural images. Indeed, since the real posterior is often far from a multivariate Gaussian, we observe a lot of variance/noise added to the model, and this causes blurriness when we decode since Maximum Likelihood training will distribute the probability mass diffusely over the data space.

**Image with 2 channels, with 2x2 filter, and output of 3 channels**

- $i = 2$  (2 channels)
- $f = 2$
- $o = 3$
- $\text{num\_params} = [i \times (f \times f) \times 3] + 3 = [2 \times (2 \times 2) \times 3] + 3 = 27$

# GANs

## Applications

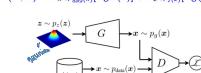
- Labels to street scenes
- How to make
- Labels to code
- How to crop
- How to move
- Labels to photo

## Generative Adversarial Network GANs

- ↳ An unsupervised architecture
- ↳ Uses 2 neural networks, competing against each other in order to generate new synthetic instances of data
- ↳ They can learn to mimic the distribution of data
- ↳ Can be used for:
  - ↳ Image generation
  - ↳ Voice generation
  - ↳ Video generation
- ↳ Steps
  - ↳ The generator takes in random numbers
  - ↳ The discriminator is fed:
    - ↳ Real images
    - ↳ Fakes generated by the generator
  - ↳ A stream of images will be generated, some fake images and others previously representing  $\xrightarrow{G}$  a Predictor of differentials
  - ↳ To fool

$$\min_{\theta_G} \max_{\theta_D} V(D, G)$$

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(G(x)))]$$



QUESTION

1. GANs are hard to train → Div. convergence hard to detect → how realistic do I want it to be?
2. Hard to meet convergence criteria
3. Vanilla GANs suffer from vanishing gradient problem
4. Overwhelming adversary
5. Can't count objects
6. Sensitive to choice of hyperparameters

as it makes learning O stand because O is hard to learn. Can't be done.

## Mode Collapse Problem in GAN

- The mode collapse problem is a common problem in GAN models.
- The generator learns to map several different  $z$  values to the same generated data point  $x$ .
- Mode collapse usually happens in GAN when the distribution of training data,  $p_{\text{data}}(x)$ , has multiple modes.

### SOLUTIONS

- Minibatch discrimination**
  - resolve the mode collapse problem
  - When mode collapses, all images created look similar.
  - Feed real images and generated images into the discriminator separately in different batches and compute the similarity of the image  $x$  with images in the same batch.
  - append the similarity  $\omega(x)$  in one of the dense layers in the discriminator to classify whether this image is real or generated.
  - If the mode starts to collapse, the similarity of generated images increases.
  - The discriminator can use this score to detect generated images and penalize the generator if mode is collapsing.

## Alternative Strategies to Prevent Mode Collapse

- Uncoupled GANs:** They allow the generator to update itself using a copy of the discriminator from several steps ahead, preventing the generator from overfitting to the current discriminator.
- Experience Replay:** This involves keeping a memory bank of past generated images and occasionally showing them to the discriminator, which helps maintain diversity in the generator's output.
- Modified Training Objectives:** Alternative loss functions, like Wasserstein loss or square losses, provide more stable gradients and can help prevent mode collapse.
- Regularization:** Adding noise to inputs or labels, or using dropout in the discriminator, can prevent it from making overly confident decisions, which in turn pressures the generator to be more diverse.
- Architecture Tweaks:** Adjusting the network architecture, such as adding more layers or changing activation functions, can help the generator explore a wider range of outputs.
- Feature Matching:** The generator is trained to match the statistical features of the real data in some intermediate layer of the discriminator, promoting diversity in the generated data.
- Penalizing the Discriminator:** Introducing penalties for the discriminator when it gets too confident can prevent it from overpowering the generator, leading to a more varied generation.
- Two-Time-Scale Update Rule (TTUR):** Using different learning rates for the generator and the discriminator can help balance their training and prevent the generator from collapsing to two modes.

## Application: Image to Image Translation by GAN

- Interactive GAN (IGAN)
- PatchGAN
- CycleGAN (Cycle-Consistent Generative Adversarial Networks)
- Simulated Gan (SimGAN)
- DeepFaceDrawing



Transforming zebra to horse and vice versa



Generating a facial image from a facial sketch.

## Training Generative Adversarial Networks

$$\max_{\theta_G} V(D, G) = \max_{\theta_D} \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x))] + p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$x = G(z) \Rightarrow z = G^{-1}(x) \Rightarrow d\pi(z) = d\pi(G^{-1}(x)) dx$$

$$\Rightarrow p_G(x) = p_G(G^{-1}(x)) \cdot |\det(G^{-1}(x))|$$

$$\int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(x))]$$

### Conditional GAN

- Assume that the dataset with which GAN is trained has  $c$  number of classes.
- The original GAN generates points from any class, and we do not have control to generate a point from a specific class.
- Conditional GAN, also called the conditional adversarial network, gives the user the opportunity to choose the class of generation of points.

For the dataset  $\{(x_i \in \mathbb{R}^{d_x})_{i=1}^n, \{y_i \in \mathbb{R}^{c_y}\}_{i=1}^n\}$ , let the one-hot encoded class labels be  $\{y_i \in \mathbb{R}^{c_y}\}_{i=1}^n$ . In conditional GAN, we can use the following loss function instead:

$$\min_{\theta_G} \max_{\theta_D} V_c(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x|y))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(x|y))]$$

$$\min_{\theta_G} \max_{\theta_D} V_c(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x|y))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(G(x|y)))]$$

the discriminator and generator are both conditioned on the labels

### Inference from Conditional GAN

- In the inference phase, the user chooses the desired class label and the generator generates a new point from that class.

### Application: Text to Image Generation

- Image is generated from some descriptive text.



### StackGAN

text-to-image generation.

## MixMatch, FineGAN

Designed to combine attributes (like background, texture, shape) from different images or labels.

### Other Applications

- In some structures of GAN, learned latent space is meaningful and we can do vector arithmetic in the latent space.
- For example in Deep Convolutional GAN (DCGAN)
- DCGAN uses an all-convolutional network for both generator and discriminator.
- **Inpainting**
- GAN learns to inpaint the lost part based on the available pixels in the image.
- **Medical application**
- Generating histopathology images which can give insight into cancer diagnosis from pathology whole slide images
- **NLP**
- **Speech processing**
- **Network embedding**
- **Logic**
- **Sketch retrieval**

## Understanding the objective function

$$\max_{\theta_G} V(D, G) = \max_{\theta_D} \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x))] + p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$\frac{\partial}{\partial D(x)} (\rho_{\text{data}}(x) \log(D(x)) + p_G(x) \log(1 - D(G(x)))) = 0$$

$$\Rightarrow \frac{\rho_{\text{data}}(x) - p_G(x)}{1 - D(x)} = 0$$

$$\Rightarrow D(x) = \frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)}$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(\frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)})] dx + p_G(x) [\log(\frac{p_G(x)}{\rho_{\text{data}}(x) + p_G(x)})] dx - \log(4)$$

$$C(G) = \int_{\mathcal{X}} p_{\text{data}}(x) [\frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)}] dx + \int_{\mathcal{X}} p_G(x) [\frac{p_G(x)}{\rho_{\text{data}}(x) + p_G(x)}] dx - \log(4) \geq 0$$

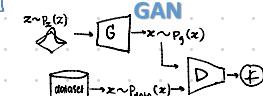
$$\min_{\theta_G} C(G) = 0 + 0 - \log(4) = -\log(4)$$

$$KL(p_{\text{data}}(x) || \frac{\rho_{\text{data}}(x) + p_G(x)}{2}) = 0$$

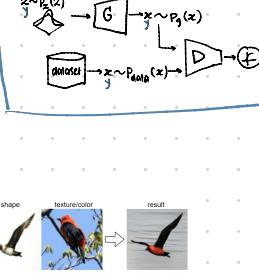
$$\text{when } \rho_{\text{data}}(x) = \frac{\rho_{\text{data}}(x) + p_G(x)}{2} \Rightarrow \rho_{\text{data}}(x) = p_G(x)$$

## Implementation of Conditional GAN

- Concatenate the one-hot encoded label  $y$  to the point  $x$  for the input to the discriminator.
- Concatenate the one-hot encoded label  $y$  to the noise  $z$  for the input to the generator.
- For these, the input layers of discriminator and generator are enlarged to accept the concatenated inputs



## Conditional GAN



### Application: Mixing Image Characteristics

- FineGAN
- An unsupervised GAN model which disentangles the features of the generated image to background, shape, and color/texture.
- FineGAN generates an image hierarchically. It starts with generating the background.

### MixMatch

- built upon FineGAN
- It gives the user the opportunity to choose the background, shape, and colour/texture from three pictures and it generates an image with the chosen characteristics.

generating an image by borrowing its characteristics from three images using MixMatch

generating an image by borrowing its characteristics from different domains using improved MixMatch.

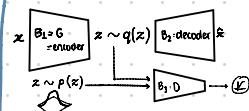


Vector arithmetic in the latent space on images by DCGAN.

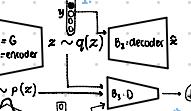
### Adversarial Autoencoder (AAE)

- In contrast to variational autoencoder which uses KL divergence and evidence lower bound, AAE uses adversarial learning for imposing a specific distribution on the latent variable in its coding layer.

### Unsupervised AAE

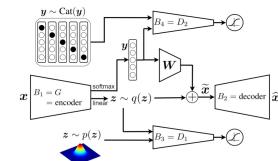


### Supervised AAE



### Dimensionality Reduction with AAE

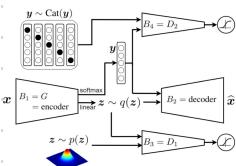
The low-dimensional representation of.



### Semi-Supervised AAE

- Consider a partially labelled dataset.
- The labelled part of data has  $n$  number of classes.
- AAE can be used for semi-supervised learning with the partially labelled dataset.
- We can use the same structure for Semi-Supervised.
- But rather than the clusters, we assume we have  $n$  number of classes.
- We have a partially labelled part of the dataset.

### Semi-Supervised AAE

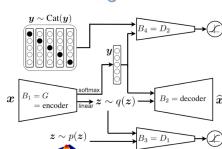


- If the point  $x$  has a label use it. Labels are one-hot vectors.
- If the point  $x$  does not have any label, randomly sample a label  $y \in \mathbb{R}^n$  from a categorical distribution, i.e.,  $y \sim \text{Cat}(y)$
- This categorical distribution gives a one-hot encoded vector where the prior probability of each class is estimated by the proportion of the class's population to the total number of labelled points.

### Clustering with AAE

- Assume we have  $c$  number of clusters.
- All points are unlabeled.
- The cluster indices are sampled randomly by the categorical distribution.
- The cluster labels and the latent code are both trained.

### Clustering with AAE



### Comparisons

#### GAN:

Think of a forger ( $G$ ) and a detective ( $D$ ).  $G$  tries to fool  $D$  with fake images.

#### VAE:

No explicit understanding of how  $G$  makes data — it just learns to trick  $D$ . Learns a smooth, interpretable "latent space" (compressed version of data). Generates new samples by decoding latent vectors.

Loss is based on how well it reconstructs and how close latent vectors are to a normal distribution. Tends to produce blurrier images but better control.

#### AAE:

Combines VAE's structure (encoder-decoder) with GAN's loss in the latent space. Instead of using KL divergence (like VAE), uses a GAN to make sure latent codes look like the prior (e.g., normal distribution). Better at clustering and generating structured latent representations.

#### Conditional GAN:

Like GAN but with labels added as input. For example, "generate a 7" instead of any digit. Helps solve the lack of control in vanilla GAN.

Model	Pros	Cons
GAN	Realistic generation, no labels needed	Training instability, mode collapse, no inference
AAE	Generative + encoding, flexible latent priors	Less sharp images, adversarial tuning needed
Supervised AAE	Combines generation & classification	Needs labeled data, may overfit
Unsupervised AAE	No labels, useful for clustering	No semantic control, weaker generation quality
Conditional GAN	Controlled generation by labels	Needs labels, can overfit, GAN issues remain
Semi-Supervised AAE	Uses limited labels effectively	Complex to train, balancing losses is tricky
Clustering AAE	Finds groups without labels, useful latent representations	May need to predefine cluster count, cluster quality depends on prior

### When to Use

#### 1. GANs (Generative Adversarial Networks)

Use when: You want to generate realistic data (e.g., images, speech) from noise.

Input: Random noise  $z$

Output: Synthetic sample

Labels needed?  No

Use cases: Image generation, super-resolution, style transfer, image inpainting

#### 2. AAE (Adversarial Autoencoder)

Use when: You want to impose a prior distribution on latent space and reconstruct input data.

Combines: Autoencoder + GAN

Input: Data sample  $x$

Output: Reconstructed  $x$  + latent vector  $z$

Labels needed?  No

Use cases: Dimensionality reduction, clustering, anomaly detection

#### 3. Supervised AAE

Use when: You have labeled data and want the latent space to reflect class structure.

Labels needed?  Yes

Learn: Both reconstruction and class-discriminative latent codes

Use cases: Semi-supervised learning, structured latent spaces

#### 4. Unsupervised AAE

Use when: You want to learn structure in data without labels

Labels needed?  No

Learn: Prior matching of latent space + reconstruction

Use cases: Clustering, representation learning, generative modeling

#### 5. Conditional GAN (cGAN)

Use when: You want to generate data conditioned on labels or attributes

Input: Noise  $z$  + condition  $y$  (e.g., class label)

Output: Sample matching the label

Labels needed?  Yes (for training)

Use cases: Text-to-image, image-to-image translation, label-controlled generation

#### 6. Semi-Supervised AAE

Use when: You have both labeled and unlabeled data

Use: A small amount of labeled data to guide training; unlabeled data to improve generalization

Labels needed?  Partially

Use cases: When labeled data is scarce but unlabeled data is abundant (e.g., medical imaging, fraud detection)

#### 7. Clustering with AAE

Use when: You want to automatically cluster data in latent space

Technique: Sample latent codes from a mixture of distributions (e.g., Gaussian + Categorical)

Use cases: Discovering groups in data, unsupervised classification, data exploration

B) Assume that we want to classify images in 3 categories of Apple, Orange and Banana for which the prior probabilities are  $P(A) = 0.3$ ,  $P(O) = 0.4$ ,  $P(B) = 0.2$ . We see a new image  $x$  from which we extract a feature vector  $x$ . The likelihoods are as:

$$P(x | A) = 0.24$$

$$P(x | O) = 0.15$$

$$P(x | B) = 0.38$$

What should image  $x$  be classified as? An apple, an orange, or a banana? Why?

Write down the complete derivations. (7 points)

$$P(A | x) = p(x | A) p(A) / p(x) = 0.24 * 0.3 / p(x) = 0.72 / p(x)$$

$$P(O | x) = p(x | O) p(O) / p(x) = 0.15 * 0.4 / p(x) = 0.75 / p(x)$$

$$P(B | x) = p(x | B) p(B) / p(x) = 0.38 * 0.2 / p(x) = 0.76 / p(x)$$

$$\Rightarrow P(B | x) > P(O | x) > P(A | x)$$

$\Rightarrow$  It has to be classified as **Banana**



### 1. ELBO (Evidence Lower Bound)

ELBO is the function VAEs maximize during training.  
It is a lower bound on the log-likelihood of the data ( $\log p(x)$ ).  
VAEs can't directly maximize  $\log p(x)$  (too complex), so they optimize ELBO instead.

### Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_t | x_0)] - D_{KL}(q(x_t | x_0) || p_{\theta}(x_t | x_0)) \\ &= \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_0 | x_t)] - \int q(x_t | x_0) \log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} dx_0 \\ &= \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_0 | x_t)] - \mathbb{E}_{q(x_t | x_0)} \left[ \log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[ \log p_{\theta}(x_0 | x_t) - \log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[ \log p_{\theta}(x_0 | x_t) + \log \frac{p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[ \log \frac{p_{\theta}(x_0 | x_t) p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] = \mathbb{E}_{q(x_t | x_0)} \left[ \log \frac{p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{p_{\theta}(x_0 | x_t) p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{p(x_T)}{\prod_{t=1}^T p(x_t | x_{t-1})} \right] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p_{\theta}(x_{t-1} | x_t)] - \mathbb{E}_q [\log q(x_t | x_{t-1})] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p_{\theta}(x_{t-1} | x_t)] - \mathbb{E}_q \left[ \log \prod_{t=1}^T q(x_t | x_{t-1}) \right] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p_{\theta}(x_{t-1} | x_t)] - \sum_{t=1}^T \mathbb{E}_q [\log q(x_t | x_{t-1})] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q \left[ \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right] \end{aligned}$$

### 2. Decomposing the ELBO

ELBO has two main parts:

ELBO = Reconstruction Term - KL Divergence

Reconstruction Term: Measures how well the decoder reconstructs input  $x$  from latent code  $z$ . It encourages accurate generation.

KL Divergence: Regularizes the latent distribution to be close to a standard normal ( $N(0, 1)$ ), preventing overfitting.

### Decomposing the ELBO

$\mathbb{E}_q[\log p(x_T)]$ : Validates the model's calibration to the true noise at the last diffusion step, ensuring the reverse process starts accurately.

$\sum_{t=1}^T \mathbb{E}_q \left[ \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$ : Acts as a regularizer by comparing the model's backward predictions to the known forward diffusion, guiding the model to correct any discrepancies.

$$\mathcal{L}(x_0) = \mathbb{E}_{q(x_{1:T} | x_0)} \left[ \log p(x_T) + \sum_{t=1}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$$

By the Markov property:

$$q(x_t | x_{t-1}) = q(x_t | x_{t-1}, x_0),$$

and by Bayes' rule:

$$q(x_t | x_{t-1}, x_0) = \frac{q(x_t | x_1, x_0) q(x_1 | x_0)}{q(x_{t-1} | x_0)}.$$

Plugging this equation into the ELBO, we get

$$\begin{aligned} \mathcal{L}(x_0) &= \mathbb{E}_{q(x_{1:T} | x_0)} \left[ \log p(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \right. \\ &\quad \left. + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t)}{q(x_t | x_0)} + \log q(x_1 | x_0) \right] \end{aligned}$$

\* =  $-\log q(x_T | x_0) + \log q(x_1 | x_0)$

Hence the negative ELBO (variational upper bound) becomes

$$\begin{aligned} & -\mathbb{E}_q \left[ \log \frac{p(x_T)}{q(x_T | x_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} + \log p_{\theta}(x_0 | x_1) \right] \\ &= \underbrace{D_{KL}(q(x_T | x_0) || p(x_T))}_{L_T(x_0)} \\ &+ \sum_{t=2}^T \underbrace{\mathbb{E}_{q(x_{1:T} | x_0)} D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t))}_{L_{t-1}(x_0)} \\ &- \underbrace{\mathbb{E}_{q(x_{1:T} | x_0)} \log p_{\theta}(x_0 | x_1)}_{L_0(x_0)} \end{aligned}$$

### 3. Reverse Process (in Diffusion Models)

Starts with pure noise and gradually denoises it using a trained model.

At each step, the model predicts the noise and removes it to recover the original data.

It's the generation phase of diffusion models.

### Revers Process

We need to learn a neural network to approximate the conditioned probability distributions in the reverse diffusion process.

$$p_{\theta}(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)).$$

$$q(x_t | x_0)$$

Recall that the forward process is defined by the equation:

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t | \sqrt{1-\alpha_t} x_{t-1}, \beta_t I)$$

We can sample  $x_t$  at any time step  $t$  in a closed form.

Let  $\alpha_t = 1 - \beta_t$  and  $\tilde{\alpha}_t = \prod_{i=t}^T \alpha_i$ :

$$q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\tilde{\alpha}_t} x_0, (1 - \tilde{\alpha}_t) I)$$

### 4. Mean ( $\mu$ ) and Variance ( $\sigma^2$ ) in VAEs

The encoder outputs two vectors: mean ( $\mu$ ) and log-variance ( $\log \sigma^2$ ). They define a Gaussian distribution in latent space:

$$z \sim N(\mu, \sigma^2)$$

Sampling from this gives diverse and continuous latent representations.

### Mean of $q(x_{t-1} | x_t, x_0)$

The reverse conditional probability is tractable when conditioned on  $x_0$ :

$$q(x_{t-1} | x_t, x_0) = q(x_t | x_{t-1}, x_0) = \mathcal{N}(x_{t-1} | x_0)$$

Using Bayes' rule, we have:

$$\begin{aligned} q(x_{t-1} | x_t, x_0) &= q(x_t | x_{t-1}, x_0) = \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \\ &\propto \exp \left( -\frac{1}{2} \left( \frac{(x_t - \sqrt{\alpha_{t-1}} x_{t-1})^2}{\beta_t} + \frac{(x_{t-1} - \sqrt{1-\alpha_{t-1}} x_0)^2}{1-\alpha_{t-1}} - \frac{(x_t - \sqrt{\alpha_{t-1}} x_0)^2}{1-\alpha_{t-1}} \right) \right) \\ &= \exp \left( -\frac{1}{2} \left( \frac{x_t^2 - 2\sqrt{\alpha_{t-1}} x_t x_{t-1} + \alpha_{t-1} x_{t-1}^2}{\beta_t} + \frac{x_{t-1}^2 - 2\sqrt{1-\alpha_{t-1}} x_{t-1} x_0 + (1-\alpha_{t-1}) x_0^2}{1-\alpha_{t-1}} - \frac{x_t^2 - 2\sqrt{\alpha_{t-1}} x_t x_0 + (1-\alpha_{t-1}) x_0^2}{1-\alpha_{t-1}} \right) \right) \\ &= \exp \left( -\frac{1}{2} \left( \left( \frac{\alpha_t}{\beta_t} - \frac{1}{1-\alpha_{t-1}} \right) x_{t-1}^2 + \left( \frac{2\sqrt{\alpha_{t-1}}}{\beta_t} x_t + \frac{2\sqrt{1-\alpha_{t-1}}}{1-\alpha_{t-1}} x_0 \right) x_{t-1} + C(x_t, x_0) \right) \right) \end{aligned}$$

where  $C(x_t, x_0)$  is some function not involving  $x_{t-1}$ .

We can express the mean and variance in the following manner: (recall that  $\alpha_t = 1 - \beta_t$  and  $\tilde{\alpha}_t = \prod_{i=t}^T \alpha_i$ ):

$$\begin{aligned} \hat{\beta}_t &= 1 / \left( \frac{\alpha_t}{\beta_t} + \frac{1}{1-\alpha_{t-1}} \right) = 1 / \left( \frac{\alpha_t - \tilde{\alpha}_{t-1} \beta_t}{\beta_t (1 - \alpha_{t-1})} \right) = \frac{1 - \tilde{\alpha}_{t-1}}{1 - \alpha_t} \cdot \beta_t \\ \hat{\mu}_t(x_t, x_0) &= \left( \frac{\sqrt{\alpha_t}}{\beta_t} x_t + \frac{\sqrt{1-\alpha_{t-1}}}{1-\alpha_{t-1}} x_0 \right) / \left( \frac{\alpha_t}{\beta_t} + \frac{1}{1-\alpha_{t-1}} \right) \\ &= \left( \frac{\sqrt{\alpha_t}}{\beta_t} x_t + \frac{\sqrt{1-\alpha_{t-1}}}{1-\alpha_{t-1}} x_0 \right) \frac{1 - \tilde{\alpha}_{t-1}}{1 - \alpha_t} \beta_t \\ &= \frac{\sqrt{\alpha_t} (1 - \tilde{\alpha}_{t-1}) x_t + \sqrt{1-\alpha_{t-1}} x_0}{1 - \alpha_t} \end{aligned}$$

Recall  $x_0 = \frac{1}{\sqrt{\tilde{\alpha}_t}} (x_t - \sqrt{1-\tilde{\alpha}_t} x_0)$

$$\begin{aligned} \hat{\mu}_t &= \frac{\sqrt{\alpha_t} (1 - \tilde{\alpha}_{t-1}) x_t + \sqrt{1-\alpha_{t-1}} \beta_t}{1 - \alpha_t} \frac{1}{1 - \tilde{\alpha}_t} (x_t - \sqrt{1-\tilde{\alpha}_t} x_0) \\ &= \frac{1}{\sqrt{\tilde{\alpha}_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} x_0 \right) \end{aligned}$$

### 5. KL Divergence

Measures how much one distribution (e.g.,  $q(z|x)$ ) differs from another (e.g.,  $N(0, 1)$ ).

In VAEs, it keeps the latent space well-behaved and prevents overfitting.

Lower KL = closer to normal distribution.

### KL for two Gaussian

$$D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t))$$

$$q(x) = \mathcal{N}(x; \mu_q, \Sigma_q) \quad \text{and} \quad p(x) = \mathcal{N}(x; \mu_p, \Sigma_p),$$

$$D_{KL}(q || p) = \frac{1}{2} \left( \text{tr}(\Sigma_p^{-1} \Sigma_q) + (\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q) - k + \ln \left( \frac{\det \Sigma_p}{\det \Sigma_q} \right) \right)$$

for  $q$

$$\hat{\mu}_t(x_t, x_0) = \frac{1}{\sqrt{\tilde{\alpha}_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \tilde{\alpha}_t}} \epsilon_t \right)$$

for  $p$

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(x_t, t) \right)$$

### 6. VAE Loss Function

Total loss = Reconstruction Loss + KL Divergence

Reconstruction Loss ensures the output matches the input.

KL Divergence regularizes the latent space.

### Loss function

$$\begin{aligned} L_t &= E_{x_t, \epsilon} \left[ \frac{1}{2 \| \Sigma_{\theta}(x_t, t) \|_2^2} \| \hat{\mu}_t(x_t, x_0) - \mu_{\theta}(x_t, t) \|^2 \right] \\ &= E_{x_t, \epsilon} \left[ \frac{1}{2 \| \Sigma_{\theta}(x_t, t) \|_2^2} \left\| \frac{1}{\sqrt{\tilde{\alpha}_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \epsilon_t) - \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(x_t, t)) \right\|^2 \right] \\ &= E_{x_t, \epsilon} \left[ \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \alpha_t) \| \Sigma_{\theta}(x_t, t) \|_2^2} \| \epsilon_t - \epsilon_{\theta}(x_t, t) \|^2 \right] \\ &= E_{x_t, \epsilon} \left[ \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \alpha_t) \| \Sigma_{\theta}(x_t, t) \|_2^2} \| \epsilon_t - \epsilon_{\theta}(\sqrt{\tilde{\alpha}_t} x_0 + \sqrt{1 - \tilde{\alpha}_t} \epsilon_t, t) \|^2 \right] \end{aligned}$$

• Ho et al. (2020) discovered empirically that the diffusion model produces higher-quality images when using a simplified objective, omitting the weighting term

$$\begin{aligned} L_{\text{simple}} &= E_{t \sim [1, T], x_0, \epsilon_t} [\| \epsilon_t - \epsilon_{\theta}(x_t, t) \|^2] \\ &= E_{t \sim [1, T], x_0, \epsilon_t} [\| \epsilon_t - \epsilon_{\theta}(\sqrt{\tilde{\alpha}_t} x_0 + \sqrt{1 - \tilde{\alpha}_t} \epsilon_t, t) \|^2] \end{aligned}$$

## 9. Attention

A technique that allows input word information to be passed all the way to the decoder.

### What is Attention?

A weighted average of input states. Each output word is generated by giving different importance (weights) to input words. It helps the model focus on important parts of the input when producing each output.

$$Y = \text{Attention}(X)$$

$y_1$	$x_1$	$x_2$	$x_3$
0.3	0.2	0.1	0.1
0.2	0.4	0.1	0.1
0.1	0.3	0.3	0.1

Each output word is influenced by all input words — but not equally. The attention weights decide how much each input word contributes.

### Exact Match

Key	Value
Alice Johnson	565-1234
Bob Smith	555-5678
Carol Davis	555-3765
Dave Martin	555-4321
Eve Clark	555-3456

Example Query:

- ▶ Query: Find the phone number for Bob Smith.
- ▶ Result: 555-5678

### Approximate Match

#### Database of Box Sizes and Weights

Key (Box Size in cubic inches)	Value (Weight in pounds)
300	30.0
400	40.0

Query: Box Size "360" cubic inches

#### Calculate Similarity Scores:

- ▶ Similarity of size "300":  $\frac{1}{300} = \frac{1}{300}$
- ▶ Similarity of size "400":  $\frac{1}{400} = \frac{1}{400}$

#### Calculate Coefficients:

- ▶  $a_1 = \frac{30}{300} = \frac{3}{10}$
- ▶  $a_2 = \frac{40}{400} = \frac{1}{10}$

#### Calculate Weighted Average Weight:

$$\text{Weighted Average Weight} = (a_1 \cdot 30.0) + (a_2 \cdot 40.0) = 12.0 + 24.0 = 36.0$$

Estimation: The estimated weight of a box with a size of "360" cubic inches is "36.0" pounds.

### Calculating Attention for the Word "bird"

**Objective:** Calculate the attention weights for the word "bird" within the sentence.

**Procedure:**

- Focus on the word "bird" as our point of interest.
- Attention Calculation:

Calculate a similarity between the query vector of "bird" and the key vectors of other words in the sentence.

word: "The", "early", "bird", "catches", "the", "worm".

embeds: [0, 0], [1, 1], [0, 1], [0, 1], [1, 1], [1, 1]

### 2. Define weight matrices

### 3. Calculate Query Vector for Bird

$$W_Q = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, W_K = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, W_V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

### 4. Compute Key Vectors

Multiple each word by  $W_K$

$$K = W_K \cdot E$$

$$\begin{aligned} \text{The} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ \text{Early} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \\ \text{Bird} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \\ \text{Catches} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \\ \text{The} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ \text{Worm} &\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \end{aligned}$$

### 6. Softmax Applying to Scores

$$\text{softmax}(s_i) = \frac{e^{s_i}}{\sum e^{s_i}}$$

$$\begin{aligned} \text{The} &: e^1 = 2.72, \quad 2.72 / (2.72+1) = 0.707 \\ \text{Early} &: e^0 = 20.09, \quad 0.007 \\ \text{Bird} &: e^1 = 5.96, \quad 0.294 \\ \text{Catches} &: e^1 = 10.93, \quad 0.295 \\ \text{The} &: e^0 = 10.96, \quad 0.699 \\ \text{Worm} &: e^1 = 15.77, \quad 0.497 \end{aligned}$$

### Traditional Sequence-to-Sequence

#### What it did:

Used an encoder to process the entire input into one vector. The decoder then used this vector to generate output.

#### Problems:

Long sentences were hard to remember → Long-range dependencies. No focus mechanism → Equal importance to every word. Couldn't parallelize well.

#### Solutions:

Handled long sentences better by accessing all encoder outputs.

Improved translation quality and long-range dependency modeling.

What is Attention?

Used an encoder to process the input sequence into a series of hidden states (not just one vector).

The decoder generated each output step by using attention to focus on relevant parts of the input.

#### Improvements:

Attention mechanism: Let the decoder dynamically focus on specific input words at each time step.

Handled long sentences better by accessing all encoder outputs.

Improved translation quality and long-range dependency modeling.

What is Attention?

Used an encoder to process the input sequence into a series of hidden states (not just one vector).

The decoder generated each output step by using attention to focus on relevant parts of the input.

#### Generalized definition

To calculate the Attention of a target word with respect to the input word.

Use the Query of the target and the Key of the input.

Calculate a matching score.

These matching scores act as the weights of the Value vectors.

Database

Query  $\rightarrow$  Value 1

Key 1 Value 1

Key 2 Value 2

Key 3 Value 3

⋮ ⋮

Key n Value n

Query  $\rightarrow$  Value i

Value 1

⋮ ⋮

Value i

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

⋮ ⋮

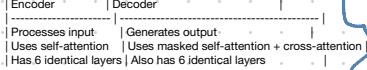
⋮ ⋮

# Transformer

## What is a Transformer?

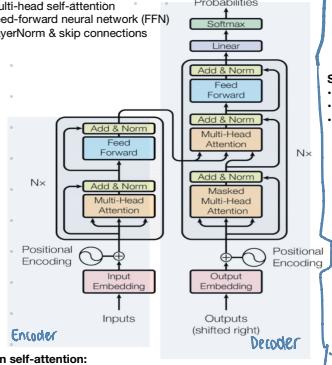
A Transformer is a deep learning architecture  
Unlike RNNs, it:  
Processes all words in parallel

Uses self-attention to understand relationships between words  
is the foundation for models like BERT, GPT, T5, etc.



## Each layer consists of:

- Multi-head self-attention
- Feed-forward neural network (FFN)
- LayerNorm & skip connections



## In self-attention:

Every word compares itself to every other word.  
It answers: "How important is word X to word Y?"

## Final Output Layer

### Linear Projection:

- Primary Role: Adjusting dimensionality.
- The linear layer serves to change the dimensionality of the feedforward network's output to match the size of the vocabulary.
- This ensures that the output has a dimension corresponding to every word in the dictionary.
- Softmax Activation:**

  - This function transforms the linear layer's output into probabilities.
  - It's applied to predict the next word.

## 1. Why do Transformers need Positional Encoding?

Transformers don't use RNNs or CNNs, so they have no built-in sense of order.

Positional Encoding adds information about the position of each word/token in the sequence to the input embeddings so the model can still understand sequence structure.

## 2. What is the difference between Self-Attention and Cross-Attention?

Self-Attention	Cross-Attention
Input = Output = Same seq,	Input ≠ Output sequences
Used in Encoder/Decoder	Used only in Decoder

Each token attends to all others in the same sequence | Each token in the decoder attends to the encoder output |

## 3. What is Multi-Head Attention and why is it useful?

- Multi-head attention:  
Splits queries, keys, and values into multiple smaller projections ("heads").
  - Each head learns different aspects of relationships (e.g., syntax, coreference).
  - The results of all heads are concatenated and passed through a final linear layer.
- It increases the model's ability to attend to multiple things at once.

## 4. What's the purpose of the Feed-Forward Network (FFN) after Attention?

- Attention captures global dependencies across tokens.
  - FFN refines each token individually using learned transformations.
- Together, they model both interactions between tokens and position-wise updates.

## 5. What does the softmax in scaled dot-product attention do?

- Converts raw attention scores into a probability distribution.
  - Ensures weights sum to 1 and highlight more relevant positions.
- The "scaling" (division by  $\sqrt{d}$ ) avoids extreme softmax outputs when dot products are large.

# ENCODER

A Transformer encoder has two main parts:

## Self-Attention

### Feedforward Neural Network (FFN)

The image focuses on the Self-Attention part

## Self-Attention:

- Input: Matrix  $X$
- Linear Transformations to generate Query ( $Q$ ), Key ( $K$ ), and Value ( $V$ ) matrices:

$$Q = W_Q^T X, \quad K = W_K^T X, \quad V = W_V^T X$$

Compute attention output  $Z$  using the formula:

$$Z = \text{softmax}\left(\frac{Q^T K}{\sqrt{d}}\right)$$

Residual Connection:

$$X + Z$$

Normalization:

$$(X + Z)$$

## Structure of the Feed Forward Network

- Linear Layer 1
- ReLU Activation
- Linear Layer 2

## FEED-FORWARD NETWORK

After attention, each token goes through a feedforward neural net, individually.

- The Feed Forward Network (FFN) is applied independently to each position in the input sequence.
- Despite individual processing, all positions share the same set of weights and biases in the FFN.
- Key Points:**
  - Shared parameters ensure consistency in processing across all positions.
  - Enables the model to generalize learnings from one position to all positions.
  - Facilitates parallel processing of the sequence, enhancing computational efficiency.

# DECODER

$$Z := \text{maskedAttention}(Q, K, V) = \text{softmax}\left(\frac{1}{\sqrt{d}}(Q^T K + M)\right)$$

where the mask matrix  $M \in \mathbb{R}^{n \times n}$  is:

$$M(i,j) := \begin{cases} 0 & \text{if } j \leq i, \\ -\infty & \text{if } j > i. \end{cases}$$

## Masked self-attention

- Cross Attention: allows each position in one sequence to attend over all positions in another sequence.
- Query ( $Q$ ): Originates from a position in the first sequence, i.e. the output of a previous layer in the decoder.

- Memory Keys ( $K$ ) and Values ( $V$ ): Both come from all positions in the second sequence, i.e. the output of the encoder.

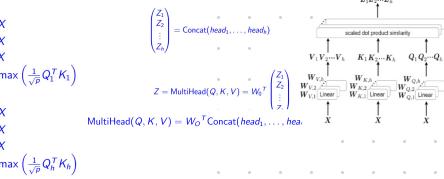
## Multi-Head Attention

Instead of doing self-attention once, we do it multiple times in parallel using different learned projections:

### Why?

Each head focuses on different parts of the sentence (e.g., grammar, position, syntax).

## Multi-Headed Attention



## Global vs Local

### Attention Mechanism:

- Global Understanding: Captures relationships among different positions in the sequence.
- Context Aggregation: Spreads relevant information across the sequence, enabling each position to see a broader context.

### Feed-Forward Networks (FFN):

- Local Processing: While attention looks across the entire sequence, FFN zooms back in to process each position independently.
- Individual Refinement: Enhances the representation of each position based on its own value, refining the information gathered so far.

### Synergy of Attention and FFN:

- Holistic Understanding:** The combination of global interaction observation (attention) and individual assessment (FFN) provides a holistic understanding of both group dynamics and individual performances.
- Balanced Processing:** A balanced approach to processing global relationships and local, position-specific information, leading to richer representations and enhanced learning.

In the decoder:

### MASKED SELF-ATTENTION

While generating output word-by-word, we mask future tokens.

The model can't "cheat" by looking at words it hasn't generated yet.

### CROSS-ATTENTION

After masked self-attention:

Each decoder word uses cross-attention to attend to encoder outputs. So the decoder knows:

"What parts of the input should I focus on to generate this word?"

### Positional Encoding

- Problem:** Transformers don't have recurrence, so they don't know order. To fix this, we add positional information to token embeddings.
- Solution:** It creates unique values for each position, letting the model differentiate:

"This word came first, this one came third..."

$$\text{PE}(pos_0, 2, +1) = \cos\left(\frac{\pi \cdot pos_0 \cdot 2}{10000}\right) \quad \text{PE}(pos_0, 2) = \sin\left(\frac{\pi \cdot pos_0 \cdot 2}{10000}\right)$$

# MULTI HEAD attention Scaled Dot Product

## 8) $x = [1, 2, 3, 4]$ , use 2 attention heads

Each projection matrix is size  $(4 \times 2) \rightarrow$  reducing dimension from 4 to 2.

### Head 1:

$$W_Q^{(1)} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_K^{(1)} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_V^{(1)} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

### Head 2:

$$W_Q^{(2)} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_K^{(2)} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad W_V^{(2)} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}$$

## For Head 1

### Step 1: Compute Q, K, V

$$Q^{(1)} = x \cdot W_Q^{(1)} = [1, 2, 3, 4] \cdot W_Q^{(1)} = [1 \times 1 + 2 \times 0 + 3 \times 1 + 4 \times 0, 1 \times 0 + 2 \times 1 + 3 \times 0 + 4 \times 1] = [4, 6]$$

$$K^{(1)} = x \cdot W_K^{(1)} = [1, 2, 3, 4] \cdot W_K^{(1)} = [1 \times 1 + 2 \times 1 + 3 \times 0 + 4 \times 1, 1 \times 0 + 2 \times 1 + 3 \times 1 + 4 \times 0] = [6, 6]$$

$$V^{(1)} = x \cdot W_V^{(1)} = [1, 2, 3, 4] \cdot W_V^{(1)} = [1 \times 0 + 2 \times 1 + 3 \times 1 + 4 \times 0, 1 \times 1 + 2 \times 0 + 3 \times 1 + 4 \times 1] = [5, 4]$$

### Step 2: Scaled Dot-Product Attention

$$\text{Dot}(Q, K) = Q \cdot K = 4 \times 4 = 16 \times 4 = 28 + 36 = 64$$

$$\text{Scale: } \sqrt{2} \approx 1.414 \approx 45.24$$

$$\text{Softmax}(64/4) = 1.0 \quad (\text{Only one value, so weight = 1})$$

So output = 1.0 · V = [5, 4]

## For Head 2

### Step 1: Compute Q, K, V

$$Q^{(2)} = x \cdot W_Q^{(2)} = [1, 2, 3, 4] \cdot W_Q^{(2)} = [1 \times 0 + 2 \times 1 + 3 \times 1 + 4 \times 0, 1 \times 1 + 2 \times 0 + 3 \times 1 + 4 \times 1] = [5, 8]$$

$$K^{(2)} = x \cdot W_K^{(2)} = [1, 2, 3, 4] \cdot W_K^{(2)} = [1 \times 1 + 2 \times 1 + 3 \times 0 + 4 \times 1, 1 \times 0 + 2 \times 1 + 3 \times 1 + 4 \times 1] = [7, 8]$$

$$V^{(2)} = x \cdot W_V^{(2)} = [1, 2, 3, 4] \cdot W_V^{(2)} = [1 \times 1 + 2 \times 0 + 3 \times 1 + 4 \times 0, 1 \times 1 + 2 \times 1 + 3 \times 0 + 4 \times 0] = [4, 3]$$

### Step 2: Scaled Dot-Product Attention

$$\text{Dot}(Q, K) = Q \cdot K = 5 \times 7 + 8 \times 8 = 35 + 72 = 107$$

$$\text{Scale: } \sqrt{2} \approx 1.414 \approx 75.64$$

Again, softmax(107/4) = 1  
→ Output = 1.0 · V = [4, 3]

## Final Answer

### Final Multi-Head Attention Output:

We concatenate the outputs of both heads:

$$\text{Head 1 output} = [5, 4] \text{ Head 2 output} = [4, 3]$$

$$\text{Final Output} = [5, 4, 4, 3]$$

## Scaled Dot Product Attention

$$\theta) q = [2, 1], k = [1, 0], v = [0, 1], z = [1, 1]$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$

$$z_1 = \sqrt{2}, z_2 = \sqrt{6}, z_3 = \sqrt{3}, z_4 = \sqrt{1}$$

$$w_1 = \sqrt{2}, w_2 = \sqrt{6}, w_3 = \sqrt{3}, w_4 = \sqrt{1}$$

$$v_1 = \sqrt{4}, v_2 = \sqrt{2}, v_3 = \sqrt{3}, v_4 = \sqrt{1}$$



# 12. CLIP

## 1. Introduction

Motivated by NLP's success with task-agnostic web-scale pretraining (e.g., GPT-3).

Goal: Can we do the same for vision, by training on natural language supervision instead of human-labeled datasets?

CLIP achieves zero-shot performance (without training on the target task).

Caveat: Align images and texts in a joint embedding space using contrastive learning.

## 2. Approach

### 2.1 Natural Language Supervision

Uses image-text pairs from the internet.

Treats language as a rich source of supervision, more scalable than labels.

Leverages contrastive learning to associate correct (image, text) pairs and separate mismatched ones.

### 2.2 Dataset Construction

400 million image-text pairs scraped from the internet.

Different from prior datasets like ImageNet – unfiltered and noisy but massive in size.

No manual labeling.

### 2.3 Model Architecture

Dual-encoder architecture:

Image Encoder: ResNet or Vision Transformer (ViT)

Text Encoder: Transformer (e.g., 12-layer, 8-head, 63M params)

Each encoder maps its input to a 512-dimensional embedding.

Uses dot product similarity in the shared embedding space.

## 3. Experiments

### 3.1 Zero-Shot Transfer

Instead of evaluating representation quality (like in self-supervised learning), they measure task generalization.

Zero-shot: CLIP predicts by choosing the class text whose embedding is closest to the image.

### 3.2 Robustness

CLIP models are more robust to distribution shifts than supervised baselines.

Shows better generalization to datasets outside training distribution.

Figure 2: Shows CLIP's 3x to 4x greater efficiency compared to Bag-of-Words and language-modeling baselines.

### 3.4 Results

Evaluated on 30+ datasets: OCR, action recognition, geo-localization, emotion recognition.

CLIP outperforms supervised models on many.

Example: Matches ResNet-50 on ImageNet without any ImageNet training.

Few-shot learning is possible via fitting linear classifiers on frozen CLIP features.

Figure 15: Shows few-shot performance drops compared to zero-shot due to lack of sample efficiency.

Table 2: Human vs CLIP accuracy on Oxford Pets. CLIP zero-shot beats human zero-shot.

### 3.5 Data Overlap Analysis

Checks if evaluation datasets leaked into the training set.

Built a duplicate detector using ResNet + custom augmentation.

### 3.6 Bias and Fairness (FairFace Analysis)

Evaluated CLIP on bias across race, gender, and age.

CLIP shows significant biases (especially in zero-shot mode).

Example: Underperforms on non-white FairFace categories compared to white.

Table 3 & 4: CLIP's accuracy on race/gender/age for white and non-white categories.

### 3.7 Broader Impacts

CLIP can classify arbitrary categories (e.g., shoplifters) which raises ethical concerns.

Promotes flexibility but also risk of misuse.

Large models like CLIP require societal accountability.

### 3.8 Related Work

Builds on contrastive learning, weakly supervised learning, image-text retrieval (e.g., ViTEx, ConVIRT, etc.)

Distinct from works that rely on complex attention-based architectures like ViLBERT or VisualBERT.

### 3.9 Conclusion

CLIP demonstrates the feasibility of language-supervised pretraining for vision.

Offers strong zero-shot generalization and task-agnostic transfer.

Promising direction, but has biases, poor few-shot performance, and data limitations.

## 4 PRACTICE QUESTIONS & ANSWERS

Q1: What is CLIP and what problem does it solve?

A: CLIP is a vision-language model that learns to associate images and text using contrastive learning on a web-scale dataset. It enables zero-shot learning for image classification tasks without needing task-specific training.

Q2: How does CLIP perform zero-shot classification?

A: By embedding a set of text prompts (like "A photo of a cat") and selecting the one whose embedding is closest (dot product similarity) to the image embedding.

Q3: What are the two types of encoders used in CLIP?

A: Image encoder (ResNet or ViT) and Text encoder (Transformer).

Q4: What loss function does CLIP use?

A: Symmetric InfoNCE loss (cross-entropy over pairwise cosine similarities).

## 5 Technical / Architecture Questions

Q5: What is the dimensionality of the embedding space?

A: 512-dimensional shared embedding space.

Q6: How is the temperature parameter  $\tau$  used?

A: Scales the logits before softmax to control sharpness in the similarity distribution.

Q7: Why does few-shot performance drop compared to zero-shot in CLIP?

A: Because CLIP isn't optimized for few-shot settings, unlike humans who benefit significantly from even a single example.

Q8: What are the image encoder variants used in CLIP?

A: ResNet-50 family (e.g., RN50x4, RN50x16, RN50x64).

Vision Transformer (ViT-B/32, ViT-B/16, ViT-L/14).

ViTs are found to be 3x more compute efficient for CLIP's objectives

Q9: How is attention used in CLIP's ResNet variant?

A: Replaces global average pooling with a single-layer multi-head attention pooling. This improves performance by letting the network learn which regions to emphasize for the image representation.

Q10: Why is the temperature  $\tau$  learned during training?

A: It controls the scale of logits before softmax in the contrastive loss. Learning it avoids manual tuning and helps balance similarity scores across batches.

## 6 Dataset and Training

Q11: How many image-text pairs are used in CLIP's training?

A: 400 million unfiltered internet (image, text) pairs.

Q12: What was the largest model trained?

A: ViT-L/14@336px and RN50x64, taking up to 592 GPUs and 18 days.

Q13: What is the loss function used in CLIP, and how does it work?

A: CLIP uses the InfoNCE loss, a symmetric contrastive loss. For a batch of image-text pairs:

It maximizes cosine similarity for matching (image, text) pairs. Minimizes it for all non-matching pairs in the batch (N-N negative pairs).

Implements it as a cross-entropy loss over similarity logits, scaled by a learnable temperature parameter  $\tau$ .

Q14: What augmentations are applied to images during training?

A: Only a random square crop from resized images. The training is kept minimal and straightforward to avoid data-specific tuning.

Q15: How does CLIP differ from prior contrastive-image-text models like Visual N-Grams?

A: CLIP trains on a dataset 10x larger, uses 100x more compute, leverages transformers, and shows a 95% top-5 ImageNet accuracy, outperforming Visual N-Grams on multiple benchmarks.

## 7 BIAS, FAIRNESS & ETHICAL IMPLICATIONS

Q16: What benchmark was used to analyze bias in CLIP?

A: FairFace – diverse face dataset categorized by race, gender, and age

Q17: What are CLIP's bias-related issues?

Zero-shot CLIP underperforms on non-white races compared to white.

Class design bias: If the user defines biased class labels (e.g., "criminal"), CLIP may reinforce harmful stereotypes.

Representation harms: Model classifies people of color more often into "animal" or "crime" categories in label sets containing such terms

Q18: How was intersectional bias studied?

A: Accuracy measured across race x gender categories. CLIP performs above 95% across most gender categories but shows disparity in age and race

Q19: What other methods were used to detect duplicates or data contamination?

A: A custom near-duplicate detector using synthetic augmentations (zoom, rotation, crop, JPEG compression, etc.) trained using a modified ResNet-50 and InfoNCE loss

## 8 BROADER IMPACTS

Q20: What are potential misuses of CLIP?

A: Tasks like:

Surveillance (e.g., classifying "shoplifters")

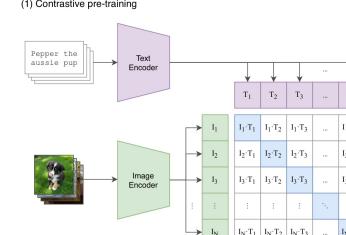
Biased classification based on race/gender

Unethical applications in facial recognition or social profiling

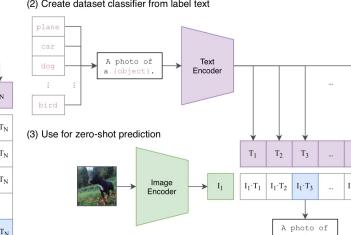
Q21: Why is CLIP's interface considered "flexible but dangerous"?

A: Users can easily define arbitrary categories using natural language. This makes CLIP powerful but vulnerable to misuse if it prompts encode bias or offensive semantics

(1) Contrastive pre-training



(2) Create dataset classifier from label text



ishma hafeez  
notes  
repsht  
hect

# 13. TVLT

## Abstract & Introduction

TVLT is a transformer-based model that learns visual-linguistic representations from raw visual and audio inputs—without relying on text, tokenization, or ASR (Automatic Speech Recognition). Instead of converting speech to text, TVLT directly aligns video frames with audio spectrograms using two core training strategies: masked autoencoding and vision-audio matching. The model is efficient (1/3 the parameters, 28x faster) and achieves comparable or better results than text-based counterparts on tasks like VQA, sentiment analysis, and retrieval.

## Related Work

The paper contrasts TVLT with:

- Text-based models like BERT, which use written language.
- Audio-based models that use spectrograms and modality-specific architectures.

VL (vision-language) models that rely on ASR for aligning text and vision. TVLT fills the gap as the first modality-agnostic, homogeneous transformer handling raw video + audio without text.

## Architecture: TVLT Overview

### Input Embeddings:

Video: 224x224 images  $\rightarrow$  16x16 patches  $\rightarrow$  Linear projection  $\rightarrow$  768-dim embeddings.

Audio: Log-mel spectrograms (128 freq bins)  $\rightarrow$  patches  $\rightarrow$  768-dim embeddings.

Added: modality embeddings + temporal/spatial (for video) or temporal/frequency (for audio) embeddings

Transformer Encoder: 12-layer, modality-agnostic.

Transformer Decoder: 8-layer, used only for masked autoencoding reconstruction.

## Training Objectives

Two losses:

Vision-Audio Matching (VAM): Predict if paired audio + video belong together. Uses [CLS] token with a sigmoid classifier. Binary cross-entropy loss

Masked Autoencoding (MAE): Randomly mask 75% of patches in both video and audio inputs. Decoder reconstructs original data. Uses MSE loss on only masked patches

## ARCHITECTURE & ENCODING

Q1: What are the input types accepted by TVLT?

A1: TVLT accepts raw video frames and audio spectrograms, not text. These are processed into patch embeddings with added modality, spatial/temporal, and frequency encodings.

Q2: What is the size of each visual patch, and how are visual embeddings created?

A: Each 224x224 frame is divided into 16x16 patches. Each patch is linearly projected into a 768-dimensional embedding.

Q3: How are audio signals transformed for use in TVLT?

A: 1D waveforms are converted into 128-dimensional log mel-spectrograms (Tx128), divided into patches (16x16 or 2x128), and projected to 768-dimensional embeddings.

Q4: What is the purpose of modality embeddings in TVLT?

A: They are trainable vectors that inform the model whether the input patch is from vision or audio.

Q5: What is the architecture of the encoder and decoder in TVLT?

A: The encoder is a 12-layer transformer (hidden size 768). The decoder is an 8-layer transformer (hidden size 512), used only during pretraining for masked autoencoding.

## PRETRAINING OBJECTIVES

Q6: Why does TVLT use a dual-objective setup for pretraining?

A: To learn strong unimodal and cross-modal representations: Vision-Audio Matching (VAM) captures cross-modal alignment, while Masked Autoencoding (MAE) strengthens unimodal reconstructions.

Q7: How are positive and negative examples generated for the VAM task?

A: Positive: matched video and audio from the same clip. Negative: video from one clip paired with audio from a different one.

Q8: What is the exact loss used for VAM?

A: Binary cross-entropy loss applied to the [CLS] token output of the encoder after passing through a sigmoid classifier.

Q9: How much of the input is masked during MAE?

A: 75% of both visual and audio patches are masked independently.

Q10: What is the final pretraining loss formula?

A: loss =  $\lambda VAM \times loss_{VAM} + \lambda MAE \times loss_{MAE}$  (where  $\lambda VAM = 1.0$  and  $\lambda MAE = 0.3$ )

## Efficiency and Design

ASR modules in traditional pipelines dominate inference time (e.g., 2890ms just for ASR).

TVLT runs in 103ms for same input due to bypassing ASR. 1/3 the parameters (88M vs. 283M)

## Results on Benchmark Tasks

TVLT outperforms text-based counterparts in audio-to-video retrieval and sentiment analysis (e.g., on CMU-MOSEI).

Slightly underperforms on audio-to-image retrieval and VQA but still competitive

Handles emotion categories better than text-based models except for "disgust"

## Ablation Studies

Joint encoder > separate encoders (better fusion).

Separate decoders > joint decoders (better reconstruction).

VAM + MAE > each alone

## Diagram Explanations

### Figure 1 – Efficiency Diagram

Shows traditional VL models (with ASR) vs. TVLT:

ASR alone takes ~2.6s.

TVLT Pipeline takes ~100ms.

TVLT is 27-28x faster for inference.

### Figure 2 – TVLT Architecture

(a) Vision-Audio Matching: Both inputs go through the encoder. A [CLS] token predicts if they are matched (label 0 or 1).

(b) Masked Autoencoding: Masked input  $\rightarrow$  encoder  $\rightarrow$  reconstructed output. Applies to both spectrogram and video.

### Figure 3 & 4 – Reconstruction Visualization

Left: masked input. Middle: reconstruction. Right: original.

Shows TVLT can successfully reconstruct masked video and audio using its learned embeddings

## Exam Practice Questions & Answers

Q1: What is the main novelty of the TVLT model compared to traditional VL models?

A: TVLT removes reliance on text and ASR by using a modality-agnostic transformer to learn directly from video frames and audio spectrograms.

Q2: What are the two main training objectives used in TVLT?

A: Vision-Audio Matching (VAM) and Masked Autoencoding (MAE).

Q3: How does Vision-Audio Matching work?

A: It creates positive (matched) and negative (random) video-audio pairs and uses binary cross-entropy loss to train a [CLS] token classifier.

Q4: How is Masked Autoencoding applied in TVLT?

A: Random patches from both audio and video are masked; the decoder reconstructs the original input using MSE loss only on masked patches.

Q5: What does the efficiency comparison show about TVLT?

A: TVLT is over 27x faster in inference than text-based VL models due to removing ASR, and it uses 1/3 the parameters.

Q6: Why is TVLT considered more "green" than other models?

A: It reduces computation by eliminating ASR modules and uses a simpler, unified architecture, although large-scale pretraining is still required.

Q7: What happens when you use joint vs. separate encoders/decoders?

A: Joint encoders perform better than separate ones, while separate decoders perform better than joint ones in masked autoencoding.

Q8: How does TVLT perform on CMU-MOSEI emotion analysis?

A: TVLT outperforms text-based models in recognizing emotions like happy, sad, angry, and fear, thanks to its ability to learn tone and loudness from raw audio.

Q9: What kind of input preprocessing is used for audio in TVLT?

A: Raw audio is converted to 128-bin log-mel spectrograms; patches are created and embedded with temporal/frequency encodings.

Q10: Why does masking speech spans help in audio MAE?

A: Because speech spans contain semantic information; masking them improves the model's ability to learn useful acoustic features.

## ABLATION STUDIES

Q11: What did the decoder architecture ablation show?

A: Using separate decoders (one for audio, one for vision) outperformed using a joint decoder.

Q12: What did the encoder ablation study reveal?

A: A joint encoder (shared for vision and audio) outperformed separate encoders followed by a fusion layer.

Q13: How did speech-span masking affect performance?

A: Masking only the speech-active regions in the spectrogram improved performance on both MSR-VTT and VQAv2.

Q14: Between 16x16 and 2x128 patch sizes for audio, which is better?

A: Mixed results. 2x128 was better on MSR-VTT; 16x16 was better on VQAv2. Default chosen is 16x16 for modality alignment.

## PERFORMANCE & RESULTS

Q15: Which benchmark tasks were used for evaluation?

A: Audio-to-video retrieval (MSR-VTT, YouCook2, CrossTask), multimodal sentiment (CMU-MOSEI), audio-to-image (Places-400k), and visual QA (VQA1/VQA2).

Q16: How does TVLT compare to its text-based version on video tasks?

A: TVLT consistently outperforms the text-based version on video tasks like retrieval and sentiment/emotion analysis.

Q17: Why is TVLT more robust in emotion classification?

A: It directly captures acoustic features like tone and loudness from raw audio, unlike text-based models limited to transcript content.

Q18: In which areas does the text-based model slightly outperform TVLT?

A: On tasks with clean, written text like VQA and audio-to-image retrieval, text-based models have a slight edge.

Q19: What efficiency gains does TVLT offer over traditional ASR pipelines?

A: TVLT is up to 28x faster in inference and uses only 88M parameters vs. 283M+ in ASR-based pipelines.

Q20: What task head is used during finetuning on retrieval tasks?

A: A two-layer MLP maps the encoder's [CLS] token to a score  $\in [0, 1]$  used for match classification.

## DATASETS & TRAINING

Q21: Which datasets were used for pretraining?

A: HowTo100M and a 20% subset of YTtemporal180M (called YTT-S).

Q22: How long does pretraining TVLT take?

A: 200k steps over 2 weeks using 4 NVIDIA RTX A6000 GPUs (49GB each).

Q23: How is ASR used in the text-based version?

A: ASR (like SpeechBrain) transcribes speech audio to text, which is then tokenized and embedded.

Q24: How is TTS used in the paper?

A: Text questions from VQA1/VQA2 are converted to synthetic speech via WaveNet and used as audio queries.

Q25: What are the hyperparameters used during pretraining?

A: LR = 1e-5, batch size = 4096, weight decay = 0.001, cosine learning rate schedule.

## BROADER INSIGHTS

Q26: What key insight about vision-language learning does TVLT challenge?

A: That high-quality vision-language representations require written language. TVLT shows this can be achieved with raw audio instead.

Q27: Why is TVLT considered more compact?

A: Because it avoids heavy modules like ASR and text tokenizers, reducing overall architecture complexity and compute.

Q28: What advantage does TVLT have in real-time systems like smart assistants?

A: Faster inference and direct audio input makes it ideal for deployment in systems with limited latency budgets.

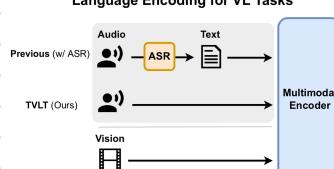
Q29: How does TVLT handle image-only tasks like VQA?

A: By converting the text questions into synthetic speech via TTS, enabling consistent input modality (audio + vision).

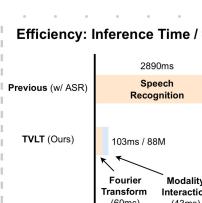
Q30: What is a practical use case where TVLT would outperform text-based models?

A: Emotion recognition from speech and video in real-time, such as in therapeutic AI or expressive virtual agents.

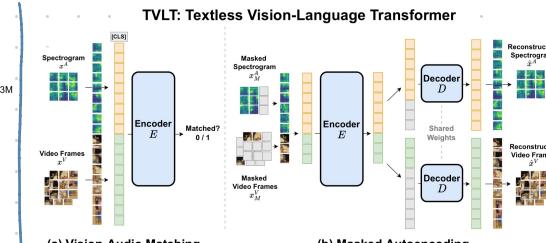
## Language Encoding for VL Tasks



## Efficiency: Inference Time / #Parameters



## TVLT: Textless Vision-Language Transformer



## Q2: Attention Masking

You are training GPT on this sequence:

"The dog barked loudly."

Which tokens can the word "barked" attend to during training?

We are predicting the word "barked".

It can only attend to:

"The"

"dog"

It cannot attend to:

"loudly" (it comes after)

## Q3: Token Flow

Draw or describe the token flow in a BERT-based classification task using this input:

[CLS] The movie was great [SEP]

What is the role of the [CLS] token?

What happens to the token embeddings after self-attention?

Tokens: [CLS], The, movie, was, great, [SEP]

All tokens go into self-attention.

Each token attends to all other tokens (because it's bidirectional).

The final representation of the [CLS] token becomes the sentence embedding.

For classification tasks (like sentiment analysis), the [CLS] token is passed to a classifier head.

Feature	BERT	GPT	T5
Architecture	Encoder-only	Decoder-only	Encoder-Decoder
Direction	Bidirectional	Left-to-right	Bidirectional (enc), causal (dec)
Pretraining Task	MLM + NSP	Next token prediction	Span corruption
Input Format	Tokenized text + [CLS]	Plain text	Text-to-text
Token Masking?	Yes	Yes (causal)	Yes (with sentinel tokens)
Can generate text?	No (not directly)	Yes	Yes

what's projection?

Suppose you have an input word embedding:

$x \in \mathbb{R}^{512}$  (i.e., 512-dimensional vector)

You apply a projection matrix  $W_Q \in \mathbb{R}^{64 \times 512}$

$$q = W_Q x$$

Then  $q \in \mathbb{R}^{64}$  is your projected query vector projection from 512-D into a 64-D

## All Model Compositions

Model	Architecture	Training Objective	Strengths	Limitations	Typical Use Cases
Diffusion Models	Markov chain noise addition + iterative denoising	Denoising score matching	Stable training, diverse & high-fidelity samples	Slow sampling, computationally heavy	High-quality image generation, inpainting
GANs	Generator + Discriminator (adversarial)	Adversarial loss	Fast generation, sharp images	Training instability, mode collapse	Real-time image/video synthesis, deepfakes
VAEs	Encoder-decoder with probabilistic latent space	Variational lower bound	Stable, diverse outputs, latent space	Blurry outputs, posterior collapse	Representation learning, anomaly detection
AEs	Encoder-decoder (deterministic)	Reconstruction loss	Simple, fast, unsupervised learning	No explicit generative modeling	Dimensionality reduction, denoising
BERT	Transformer encoder (bidirectional)	Masked Language Modeling + NSP	Deep understanding, strong NLP performance	Not generative, heavy compute	NLP understanding, classification
GPT	Transformer decoder (unidirectional)	Next-token prediction	Strong text generation	Limited context, large compute	Text generation, chatbots, summarization

## Formula Sheet

### Convolutional Layers (CNN)

#### Conv2D Output Size:

$$\text{Output size} = \left\lfloor \frac{\text{Input size} - 2 \cdot \text{padding} - (\text{kernel size} - 1)}{\text{stride}} \right\rfloor + 1$$

#### Pooling Output:

$$\text{Output size} = \left\lfloor \frac{\text{Input size} - \text{kernel size}}{\text{stride}} \right\rfloor + 1$$

#### Number of Weights (Conv Layer):

$$\text{Weights} = \text{Out channels} \times \text{In channels} \times \text{Kernel height} \times \text{Kernel width}$$

#### Number of Biases (Conv Layer):

$$\text{Biases} = \text{Out channels}$$

### Positional Encoding (Transformers)

$$\text{PE}(\text{pos}, i) = \begin{cases} \sin\left(\frac{\text{pos} \cdot \text{pos}_i}{10000^{(i-1)/d}}\right) & \text{if } i \text{ is even} \\ \cos\left(\frac{\text{pos} \cdot \text{pos}_i}{10000^{(i-1)/d}}\right) & \text{if } i \text{ is odd} \end{cases}$$

### Scaled Dot Product Attention

#### Attention Score:

$$\text{score} = \frac{\vec{q} \vec{v}^T}{\sqrt{d_k}}$$

#### Softmax:

$$\text{Softmax}(\vec{x}_i) = \frac{e^{\vec{x}_i}}{\sum_j e^{\vec{x}_j}}$$

#### Weighted Sum (Attention output):

$$\text{Output} = \sum_i \text{softmax}(\text{score}_i) \cdot v_i$$

### Diffusion/VAEs

#### KL Divergence (VAE Loss)

$$D_{KL}(q(z|x)||p(z)) = \frac{1}{2} \sum_i (\mu^2 + \sigma^2 - \log(\sigma^2) - 1)$$

$$D_{KL}(P||Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right)$$

Note: Use ln (natural log)

### Harris Corner Detection

Used for detecting corners in an image.

#### Step 1: Compute image gradients

$$I_x = \frac{\partial I}{\partial x}, \quad I_y = \frac{\partial I}{\partial y}$$

#### Step 2: Compute structure matrix MMM

$$M = \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \quad (\text{smoothed over a window})$$

#### Step 3: Harris Response

$$R = \det(M) - k \cdot (\text{trace}(M))^2$$

$$\bullet \det(M) = I_x^2 I_y^2 - (I_x I_y)^2$$

$$\bullet \text{trace}(M) = I_x^2 + I_y^2$$

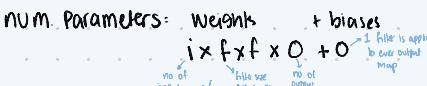
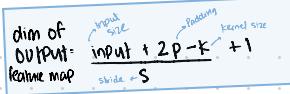
$$\bullet k \in [0.04, 0.06]$$

#### If R

#### Large positive $\rightarrow$ corner

#### Large negative $\rightarrow$ edge

#### Close to zero $\rightarrow$ flat region



grey scale image - 1 channel

RGB image - 3 channels

what's projection?

Suppose you have an input word embedding:

$x \in \mathbb{R}^{512}$  (i.e., 512-dimensional vector)

You apply a projection matrix  $W_Q \in \mathbb{R}^{64 \times 512}$

$$q = W_Q x$$

Then  $q \in \mathbb{R}^{64}$  is your projected query vector projection from 512-D into a 64-D

### Positional Encoding

position  $z_i = x_i + \text{Positional encoding}_i$

#### Step 1: Get token embeddings

Imagine the embedding layer gives us fixed vectors for tokens:

Token 1:  $[0, 1, 0, 2, 3, 0, 4, 0, 5, 0, 6]$

love:  $[0, 5, 0, 4, 0, 3, 2, 0, 1, 0, 0]$

AI:  $[0, 9, 0, 8, 0, 7, 0, 6, 0, 5, 0, 4]$

#### Step 2: Calculate positional encoding vectors

$$\text{formula: } \text{PE}(\text{pos}, i) = \sin\left(\frac{\text{pos} \cdot \text{dim}}{10000^{(i-1)/6}}\right) \text{PE}(\text{pos}, 2i+1) = \cos\left(\frac{\text{pos} \cdot \text{dim}}{10000^{(i-1)/6}}\right)$$

For pos = 6, dimensions  $i = 0, 1, 2$ :

Calculate denominator  $10000^0 = 1$  for each  $i$ :

$$1: 2i/d \quad \text{Denominator} = 10000^{0/6} = 1$$

$$0: 10000^1 = 1$$

$$1: 2/6 \cdot 333 = 10000^{0.333} \approx 21.54$$

$$2: 4/6 \cdot 666 = 10000^{0.666} \approx 46.16$$

Calculate each position:

$$\text{For pos = 0: } \sin(0/1) = 0, \quad \cos(0/1) = 1$$

$$\sin(1/21.54) \approx 0.0464, \cos(1/21.54) \approx 0.9989$$

$$\sin(1/46.16) \approx 0.0225, \cos(1/46.16) \approx 0.999977$$

$$\text{PE}_1 = [0.8415, 0.5403, 0.0464, 0.9989, 0.0215, 0.999997]$$

For pos = 1:

$$\sin(2/21.54) \approx 0.0903, \cos(2/21.54) \approx -0.4161$$

$$\sin(2/46.16) \approx 0.00231, \cos(2/46.16) \approx 0.999997$$

$$\text{PE}_2 = [0.9983, -0.4161, 0.0927, 0.9957, 0.00231, 0.999997]$$

#### Step 3: Add positional encoding to token embeddings

Token: Read input to position:

$$1: [0, 1, 0, 2, 3, 0, 4, 0, 5, 0, 6]$$

$$\text{love: } [1, 3415, 0, 9403, 0, 3464, 1, 1989, 0, 10215, 0, 999997]$$

$$\text{AI: } [1, 8093, 0, 3898, 0, 7927, 1, 5957, 0, 50431, 1, 3999997]$$

### Additional Attention/Rishanau Attention:

You are doing machine translation using seq2seq + attention. We are at decoder timestep 1.

We have 3 encoder hidden states:  $h_1, h_2, h_3$

We have 1 decoder hidden state from previous timestep:  $s_0$

$$h_1 = [1, 1], \quad h_2 = [0, 1], \quad h_3 = [1, 1], \quad s_0 = [1, 2]$$

We'll use a score function of the form:

$$s_{ij} = v^T \tanh(W_s s_i + W_h h_j)$$

$$W_s = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_h = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad v = [1, 1]^T$$

For  $h_1 = [1, 0]$

$$W_s h_1 = [1 \times 1 + 1 \times 0, 1 \times 1 + (-1) \times 0] = [1, 1]$$

$$W_s s_0 + W_h h_1 = [1 + 1, 2 + (-1)] = [2, 1]$$

$$\tanh([2, 1]) \approx [0.964, 0.995]$$

$$e_{11} = v^T \cdot [0.964, 0.995] = 0.964 + 0.995 = ** 1.959 **$$

For  $h_2 = [0, 1]$

$$W_s h_2 = [0 \times 1 + 1 \times 1, 0 \times 1 + (-1) \times 1] = [-1, 1]$$

$$W_s s_0 + W_h h_2 = [1 + 1, 2 + (-1)] = [2, 1]$$

$$\tanh([2, 1]) \approx [0.964, 0.994]$$

$$e_{12} = 0.964 + 0.994 = ** 0.964 **$$

For  $h_3 = [1, 1]$

$$W_s h_3 = [1 \times 1 + 1 \times 1, 1 \times 1 + (-1) \times 1] = [2, 0]$$

$$W_s s_0 + W_h h_3 = [1 + 2, 2 + (-1)] = [3, 2]$$

$$\tanh([3, 2]) \approx [0.995, 0.964]$$

$$e_{13} = 0.995 + 0.964 = ** 1.959 **$$

Apply softmax to get attention weights  $a_{ij}$

$$\text{softmax}(e_{11}, e_{12}, e_{13}) = \frac{e^{e_{11}}}{\sum e^{e_{ij}}}$$

$$e^{1.959} \approx 7.09, \quad e^{1.725} \approx 5.61$$

$$\text{Denominator} = 7.09 + 5.61 + 7.09 = 19.79$$

$$a_{11} = \frac{7.09}{19.79} \approx 0.358 \quad a_{12} = \frac{5.61}{19.79} \approx 0.285 \quad a_{13} = \frac{7.09}{19.79} \approx 0.358$$

Compute context vector  $c_i$

$$h_1 = [1, 0], \quad h_2 = [0, 1], \quad h_3 = [1, 1]$$

$$c_i = [0.358, 1, 0] + [0.283, 0, 1] + [0.358, 1, 1] = [0.358, 0] + [0, 283] + [0.358, 0.358] = [0.716, 0.641]$$

Attention weights  $[0.358, 0.283, 0.358]$

$$\text{Context vector } c_i = [0.716, 0.641]$$

$$v_1 \quad v_2 \quad v_3 \quad v_4 \quad [2] \quad \text{Ans}$$

# Transformers

$$8) \quad x = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_i = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, w_v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

1. Compute  $\theta, k, v$

$$\theta = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^T = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

2. Compute Attention Scores

$$\theta = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad k^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$k^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

scaled by  $\sqrt{2}$  (1/1)

$$\frac{\theta k^T}{\sqrt{2}} = \begin{bmatrix} 0.71 & 0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

3. Apply softmax

$$\text{softmax}(0.71, 0.71) = (0.5, 0.5)$$

$$\text{softmax}\left(\frac{\theta k^T}{\sqrt{2}}\right) = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

4. Multiply V

2. Softmax.V

$$\begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} = \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix}$$

5. Residual Connection

=  $x + z$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix}$$

$$= \begin{bmatrix} 1.25 & 0.25 \\ 0.25 & 1.25 \end{bmatrix}$$

Hello [1,2] world [3,M]

$$8) \quad x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, w_i = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, w_v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

1. Compute  $\theta, k, v$

$$\theta = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^T = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

2. Compute Attention Scores

$$\theta = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad k^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$k^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

scaled by  $\sqrt{2}$  (1/2)

$$\frac{\theta k^T}{\sqrt{2}} = \begin{bmatrix} 6.36 & 14.89 \\ 14.89 & 34.64 \end{bmatrix}$$

3. Apply softmax

$$\text{softmax}\left(\frac{\theta k^T}{\sqrt{2}}\right) = \begin{bmatrix} 0.00021 & 0.99979 \\ 0 & 1 \end{bmatrix}$$

4. Multiply V

$$2. \text{Softmax.V} = \begin{bmatrix} 0.5 & 1 \\ 1.5 & 2 \end{bmatrix}$$

$$2. \begin{bmatrix} 0.00021 & 0.99979 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.5 & 1 \\ 1.5 & 2 \end{bmatrix} = \begin{bmatrix} 1.49979 & 1.99979 \\ 1.5 & 2 \end{bmatrix}$$

5. Residual Connection

=  $x + z$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 1.49979 & 1.99979 \\ 1.5 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 2.5 & 4 \\ 4.5 & 6 \end{bmatrix}$$

18/07/2023  
Audio is aligned with the video

Do frame by frame but it has less accuracy than video based.

3D-CNN for videos but has more parameters.  
2D web-based has limited parameters.

Department of Education

- It's very slow

(b) The text is not bringing very big impact into the picture, remove it

title

By applying Fourier transformation you can map convert 1D into 2D signal

As Image encoder is 2D & audio encoder is 1D and we end up with 2 encoders slowing us down, using Fourier transform

we end up with only 1 encoder speed is really good.

It is called spectrogram, when you end up converting 1D signal into 2D signal

ASR (case study)  
Automatic Speech recognition 02/05/2023

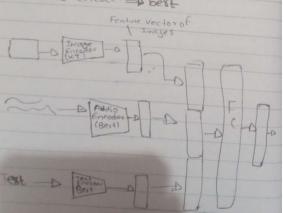
ASO

Kids with read autism

Web-based Emotion recognition  
It takes video as input

End goal off

Text encoder -> best  
Image Encoder -> VIT  
Audio encoder -> best



Paper is called CLIP  
Contrastive Language Pre-training

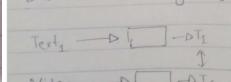
Learning Transferable Visual Models from Natural Language Supervision

Used to extract features from Text.

From features we would use Vision encoders.

CLIP (Research paper reading assistant)

Contrastive Learning Image Pre-training



If those 2 belong to the same sample we put the class in the feature space.

ishma halfeez notes

represent

CLIP  
Contrastive learning, very important used in computer vision  
Masked auto encoding