# https://github.com/TAMIDSpiyalong/Gen-AI

Lecture 1: 9:00-9:40

# Yalong Pi (Texas A&M Institute of Data Science)
## Address: Office 221C, John R. Blocker Building
## Email: piyalong@tamu.edu

- ❑ B.S., Mechanical Engineering, 2007-2011
- ❑ M.S., Civil Engineering, 2011-2013
- ❑ Ph.D., Architecture, 2017-2020

- ❑ Research Scientist, 2020-present
- ❑ Architect, 2016-2017
- ❑ Project manager, 2013-2016

# Agenda

- **Part 1: Background and Theory on Generative AI**
- Fundamentals of machine learning    9:00 — 10:00
- Tokenization and word embedding    10:00 — 10:30
- Transformers for Language Models    10:30 — 12:00
- **Part 2: Applications and Hands-On Exercises**
- Prompt Engineering
- Generative AI Applications
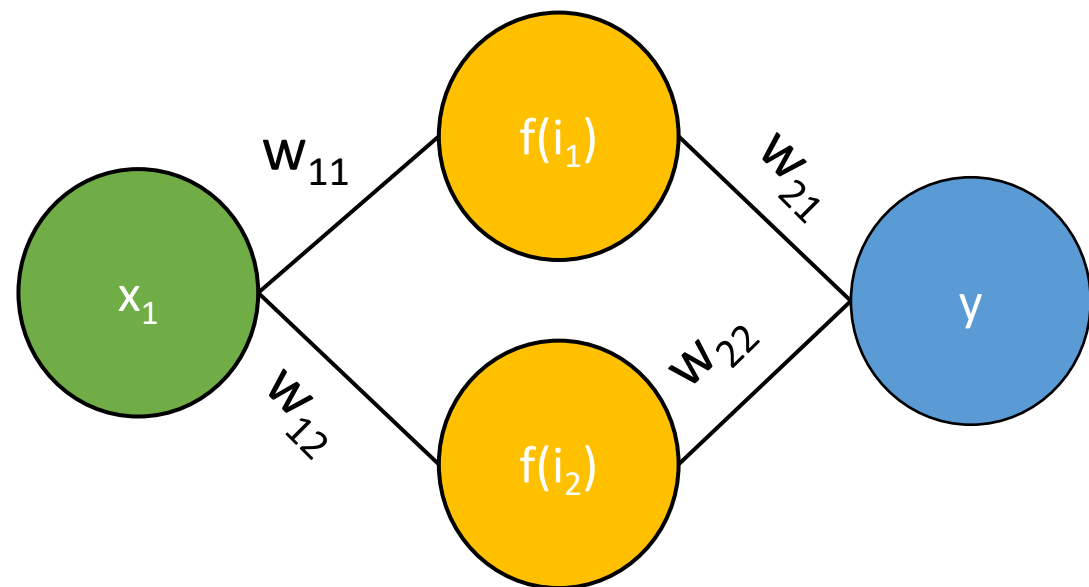- Evolution of Generative AI and Future Directions

# Fundamentals of machine learning

NLP

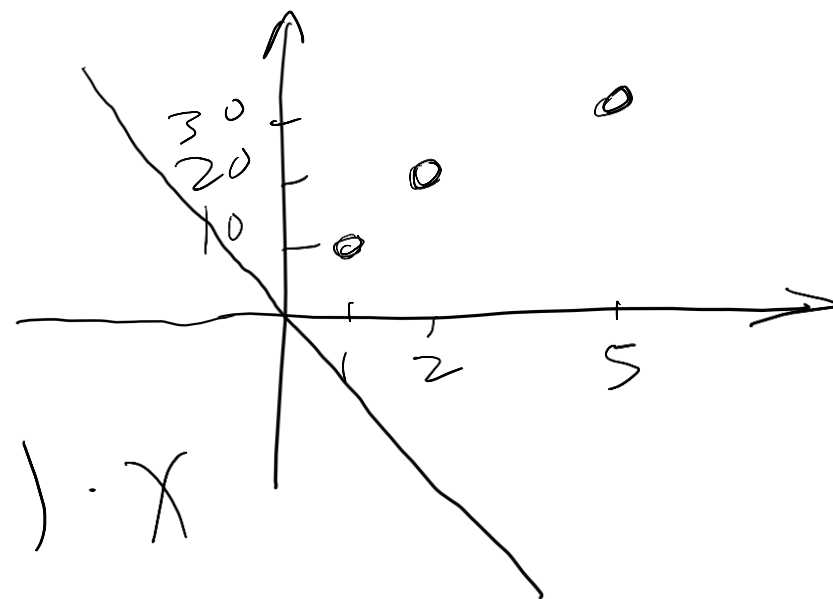CV

alfa fold

recomendation

# Back Propagation

| x | y |
|---|---|
| 1 | 10 |
| 2 | 20 |
| 5 | 30 |

$$y_{pred} = w_{21} * x_1 * w_{11} + w_{22} * x_1 * w_{12}$$

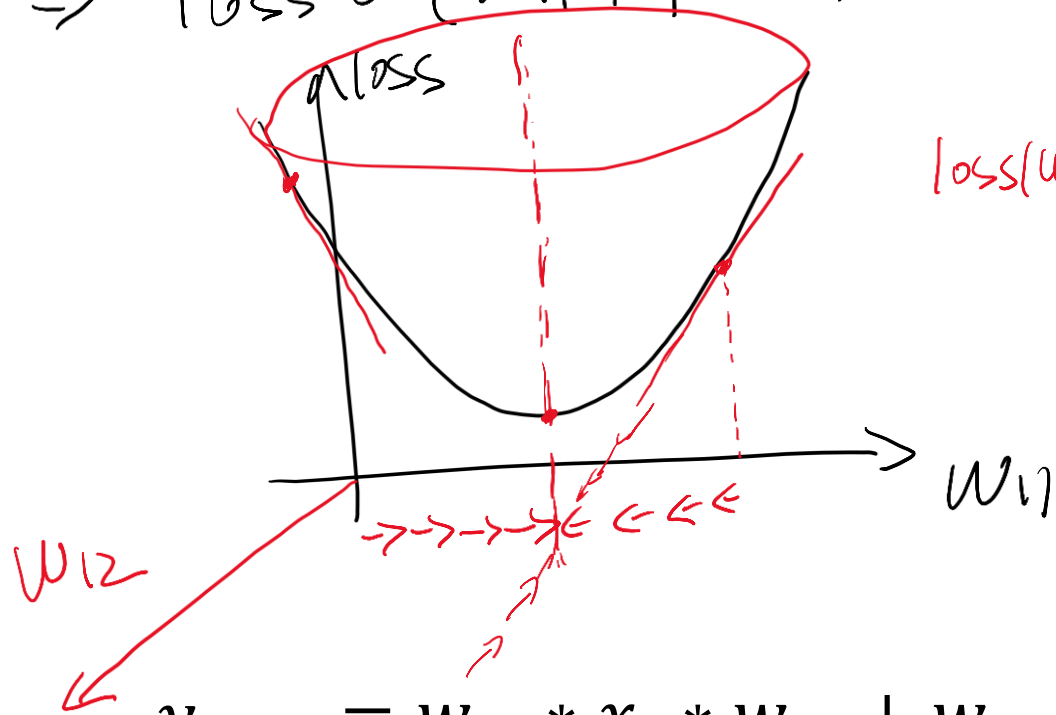$$Loss = (y_{pred} - y_{true})^2$$

$$W_n' = Wn - LR \, (\partial Loss / \partial W_n)$$

$$y = (W_{11} \cdot W_{21} + W_{22} \cdot W_{12}) \cdot x$$

$$x = 1 = W_{21} = W_{22} = W_{12}$$

$$\Rightarrow \ loss = (W_{11} + 1 - 10)^2 = (W_{11} - 9)^2 = W_{11}^2 - 18W_{11} + 81$$



$$loss(w_{11}) \ \frac{\Delta loss}{\Delta W_{11}} = \frac{\partial loss}{\partial W_{11}} = 2W_{11} - 18 = \frac{loss(w_{11} + \Delta w_{11})}{\Delta W_{11}}$$

$$W_{11}' = W_{11} - \frac{\partial loss}{\partial W_{11}} \cdot LR$$

$$W_{12}' = W_{12} - \frac{\partial loss}{\partial W_{11}}$$

$$y_{pred} = w_{21} * x_1 * w_{11} + w_{22} * x_1 * w_{12}$$

$$Loss = (y_{pred} - y_{true})^2$$

$$W_n' = Wn - LR \ (\partial Loss / \partial W_n)$$

| x | y | |
|---|---|---|
| 1 | 10 | ✓ |
| 2 | 20 | ✓ |
| 5 | 30 | ✓ |

$$\begin{bmatrix} w_{11}' \\ w_{12}' \\ w_{21}' \\ w_{22}' \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{12} \\ w_{21} \\ w_{22} \end{bmatrix} - \begin{bmatrix} \dfrac{\partial loss}{\partial w_{11}} \\ \dfrac{\partial loss}{\partial w_{12}} \\ \dfrac{\partial loss}{\partial w_{21}} \\ \dfrac{\partial loss}{\partial w_{22}} \end{bmatrix} \cdot LR$$

SGD    adam

$$\vec{W}' = \vec{W} - \nabla \cdot LR$$

0.001

batch

$$y_{pred} = w_{21} * x_1 * w_{11} + w_{22} * x_1 * w_{12}$$

$$Loss = (y_{pred} - y_{true})^2$$

$$W_n' = Wn - LR \, (\partial Loss / \partial W_n)$$

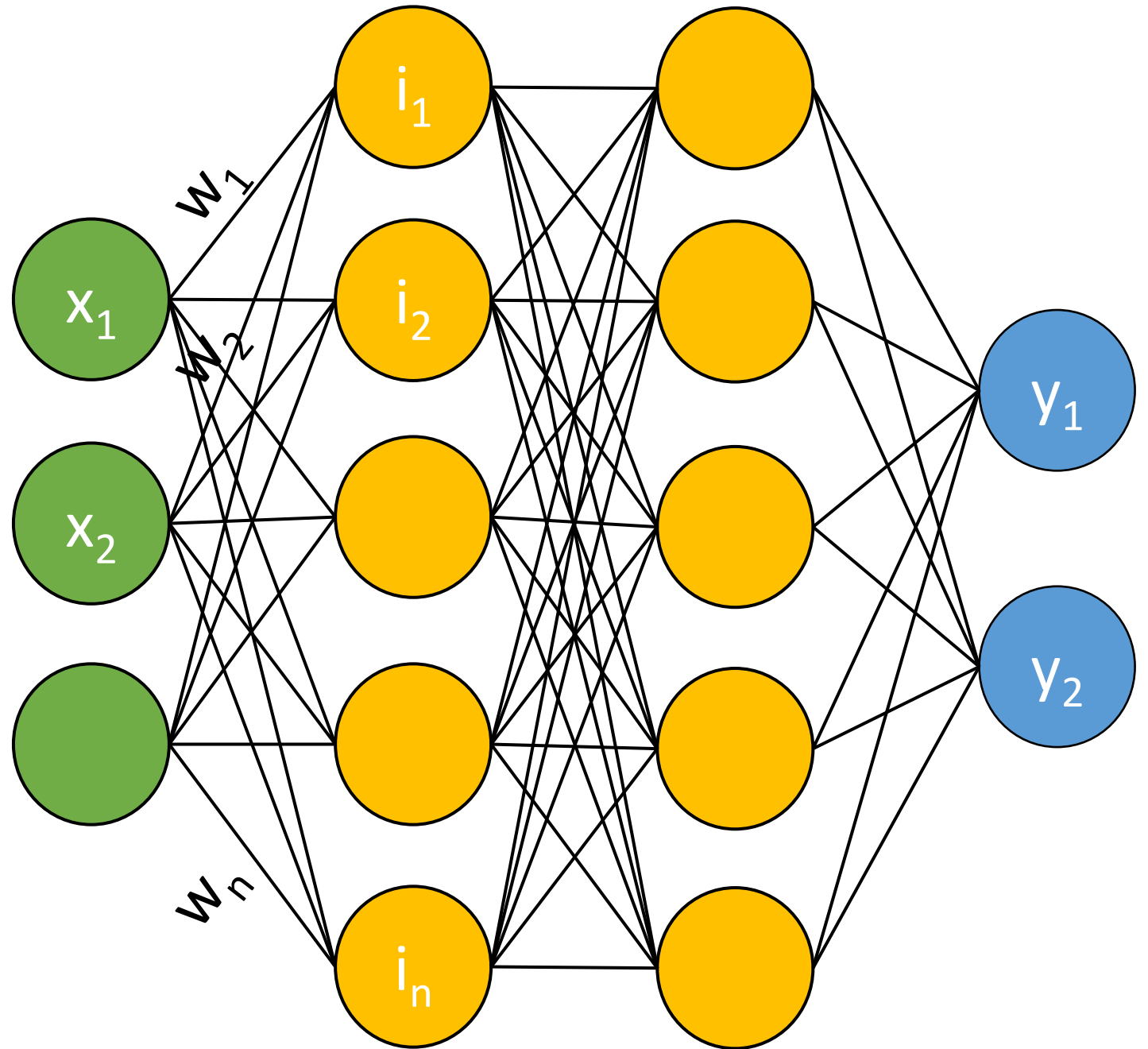| x | y | |
|---|---|---|
| 1 | 10 | $\nabla_1$ |
| 2 | 20 | $\nabla_2$ |
| 5 | 30 | $\nabla_3$ |

For each connection:

$$I_n = f(\sum_n x_n w_n + b)$$
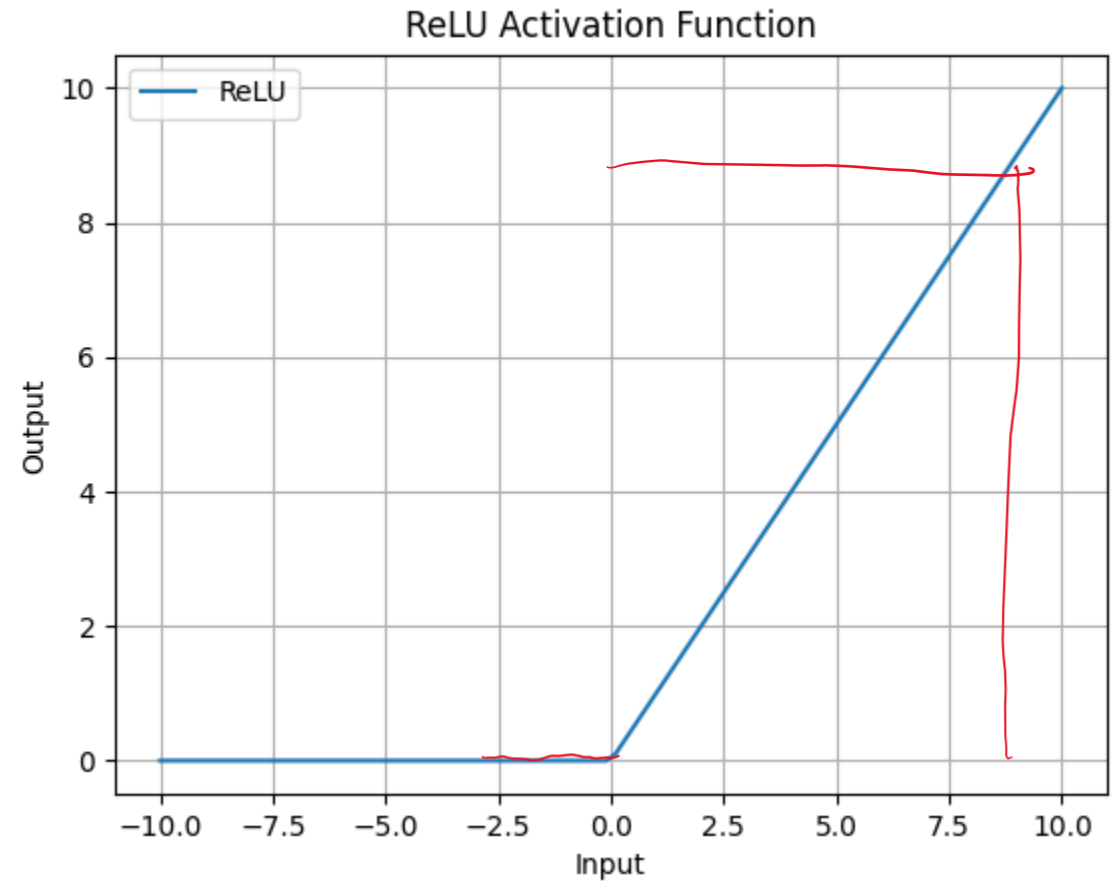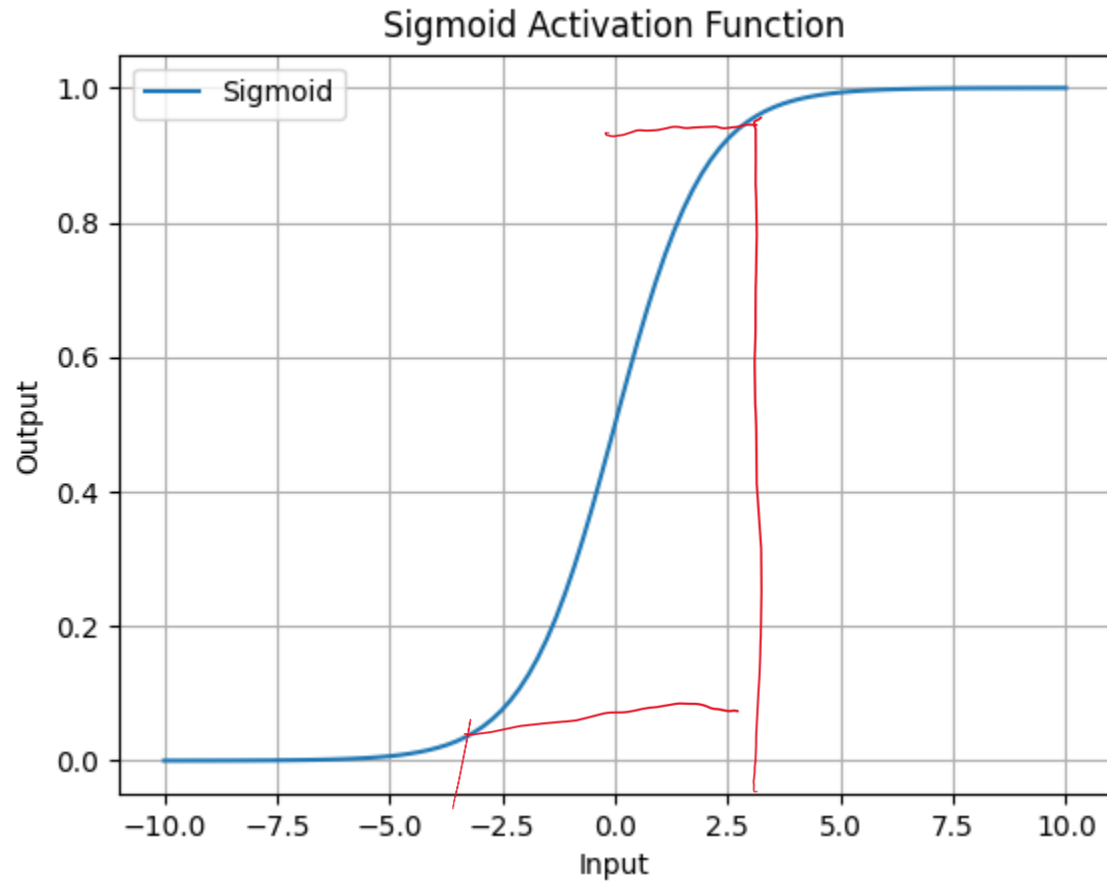
- ☐ $f$ is the activation function
- ☐ $w_n$ is the weight
- ☐ $b$ is the bias.
- ☐ A DNN has millions of weights and biases
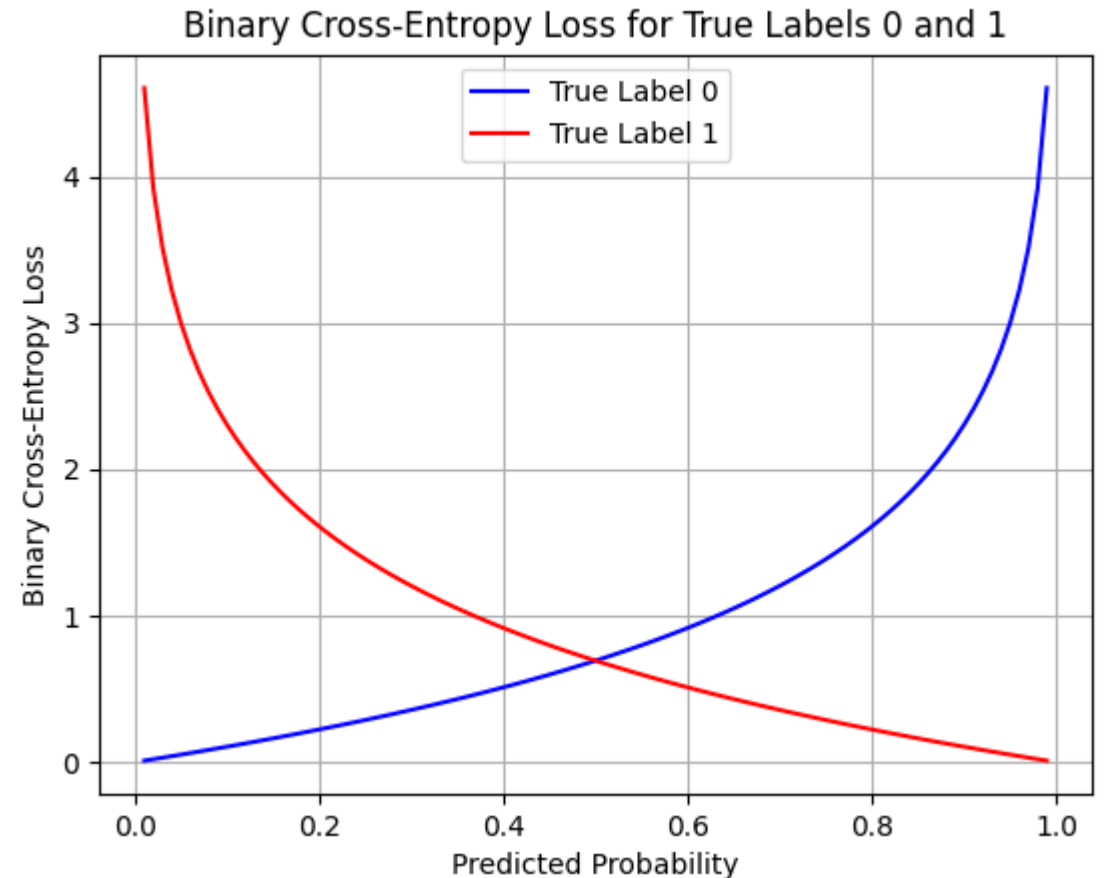
# Activation Functions

# Binary Cross Entropy

$$L(y, \hat{y}) = -[\, y * \log(\hat{y}) + (1 - y) * \log(1 - \hat{y})\,]$$

Where:
- y is the true label (0 or 1)
- $\hat{y}$ is the predicted probability (between 0 and 1)
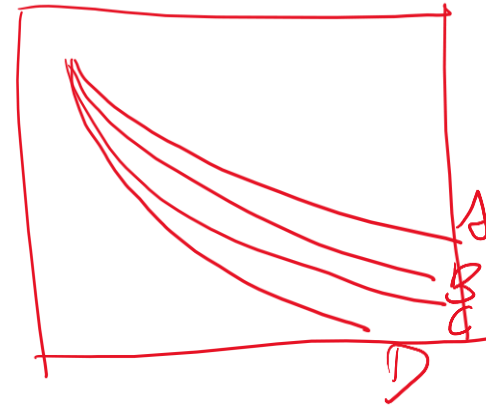- log can be ln



Binary Cross-Entropy Loss for True Labels 0 and 1

# Softmax

$$\sigma(Z)_i = \exp(z_i) / \Sigma(\exp(z_j))$$

Where:

- $z_i$ is the $i$ input score (logits)

- $\Sigma$ is the sum over all input scores (logits)

- $\sigma(z)_i$ is the probability assigned to class $i$

*CE*

*not after softmax*

Example:

[-0.37, -1.06, -0.07, -1.47, -0.90] -> [0.265, 0.133, 0.358, 0.088, 0.155]

*A    B    C    D    E*

| Prediction\Ground Truth | Positive | Negative |
|---|---|---|
| Positive | TP | FP |
| Negative | FN | TN |

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

$$x_1 W_{11} + x_2 W_{12} + \ldots + b_1$$
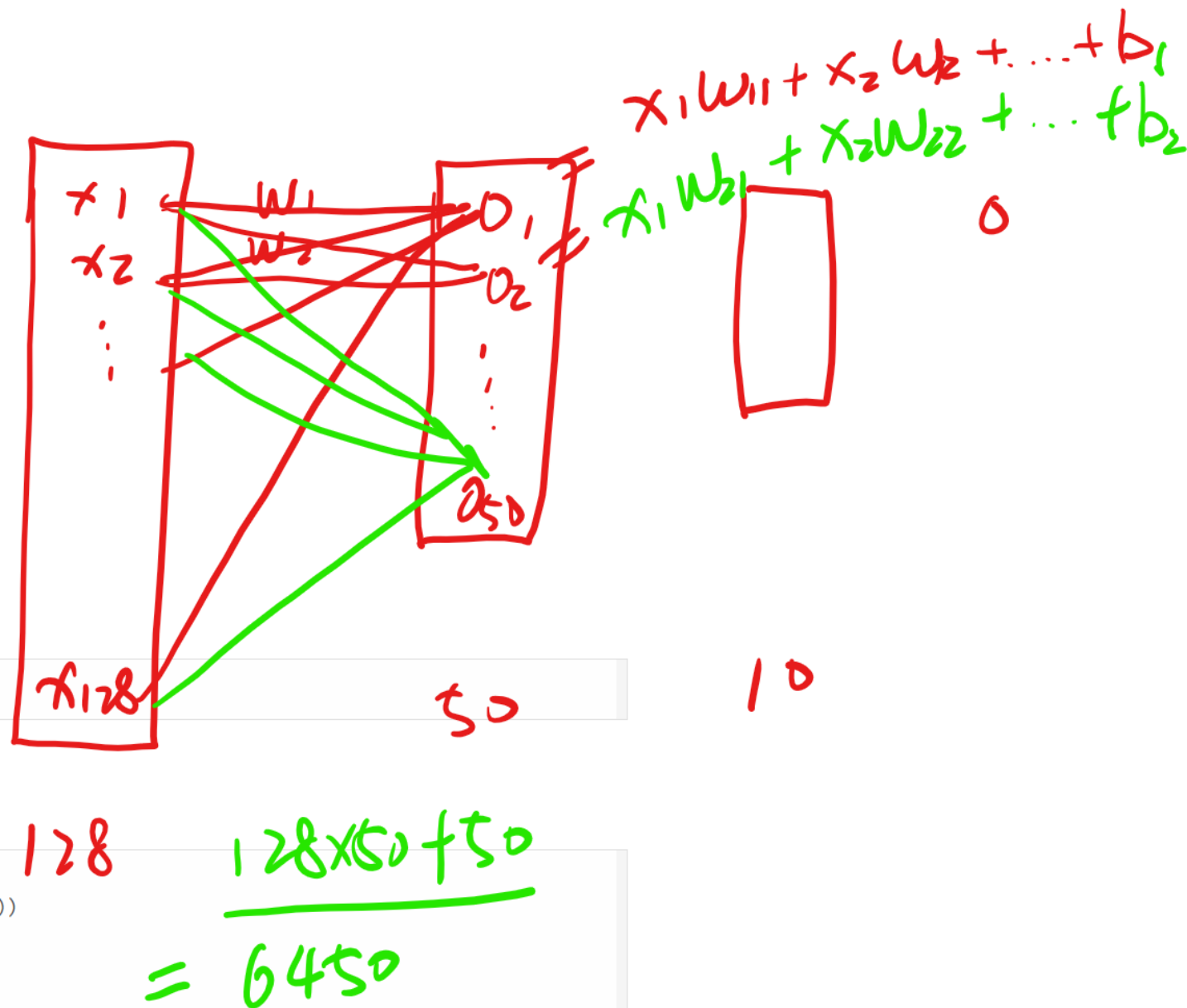$$x_1 W_{21} + x_2 W_{22} + \ldots + b_2$$
$$0$$

```
-4.37806658e-02, -9.67954174e-02, -1.28794938e-01,
 1.45295113e-01,  3.34722281e-01, -2.49092355e-02,
-7.21996874e-02,  7.60333985e-02,  1.09788142e-01,
-5.91083243e-02,  1.77998230e-01,  1.05147131e-01,
 2.73706466e-01,  1.63680017e-01,  2.92986393e-01,
 1.62288636e-01,  1.92936987e-01, -7.25108087e-02,
 1.48647577e-01,  1.20697133e-01,  1.75806686e-01,
-8.22802186e-02,  3.19161601e-02,  9.75683853e-02,
-2.27390900e-01, -1.89130962e-01, -7.75573701e-02,
 7.51652941e-02, -9.91581455e-02,  9.62962583e-03,
 6.42622411e-02, -1.50064066e-01,  1.14945382e-01,
 5.34672337e-03,  1.96428418e-01, -2.10412573e-02,
 5.44419959e-02, -3.22782189e-01,  5.69203123e-03,
-1.00528084e-01, -7.29111880e-02, -1.84138656e-01,
-1.51076904e-02, -6.30587935e-02]], dtype=float32)
```

$x_1$
$x_2$
$W_1$
$W_2$
$O_1$
$O_2$
$250$
$x_{128}$
$50$
$10$

In [ ]: `hub_layer(train_examples[:3]).shape`

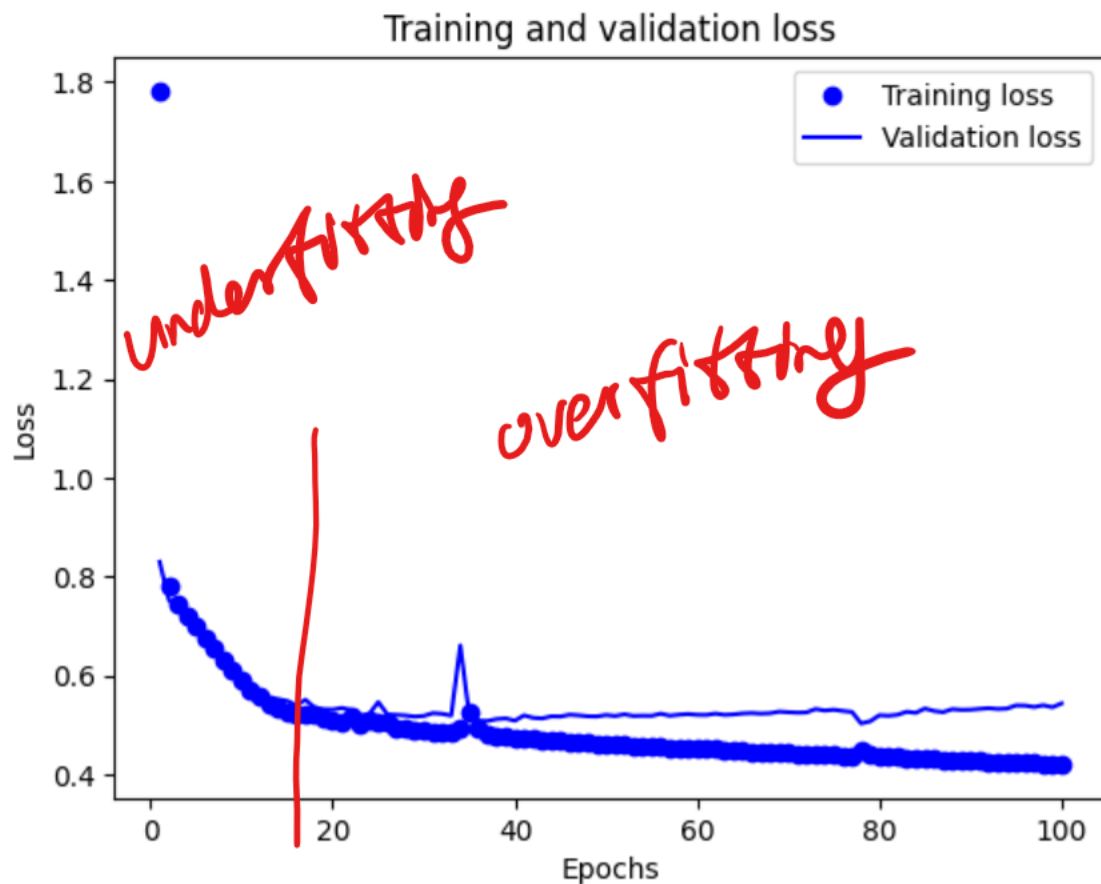Out[ ]: `TensorShape([3, 128])`

Let's now build the full model:

In [ ]:
```python
model = tf.keras.Sequential()
model.add(tf.keras.layers.Dense(50, activation='relu',input_dim=128))
model.add(tf.keras.layers.Dense(10, activation='relu'))
model.add(tf.keras.layers.Dense(1))
```

$128$

$$\frac{128 \times 50 + 50}{= 6450}$$

80%

training valid test

Training and validation loss

underfitting

overfitting

• Training loss
— Validation loss

early stop

save the best

data augmentation
regularization
drop out

# Tokenization and word embedding

# Tokenization

"I want pizza"

"我想要披萨"

"ピザが欲しいです"

"Eu quero pizza"

"أريد بيتزا"

"मुझे पिज़्ज़ा चाहिए"

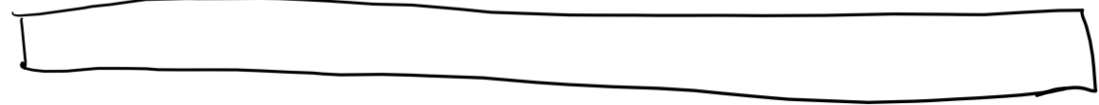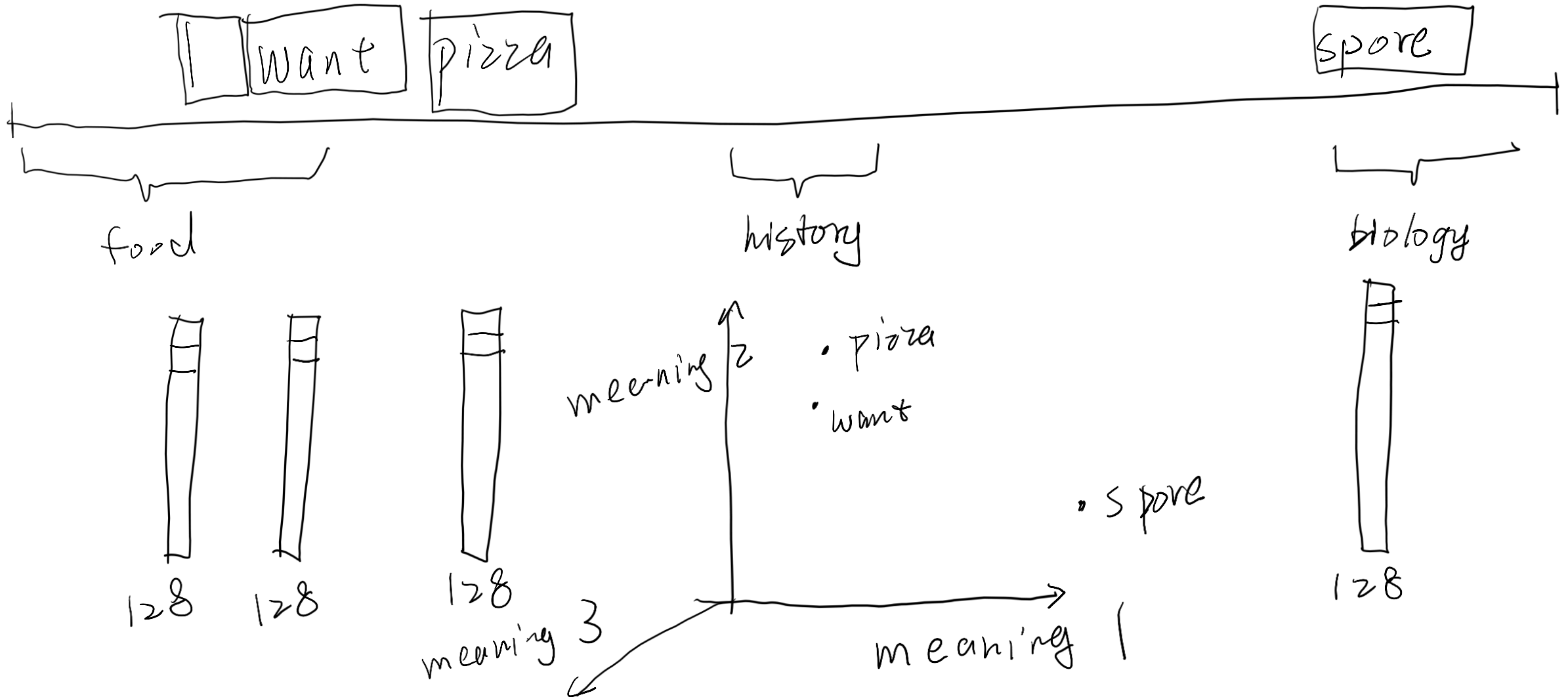"Quiero pizza"

"피자가 먹고 싶어요"

pizzas

Wanted

token vocabulary

ASCII

40k-45k

50k   GPT

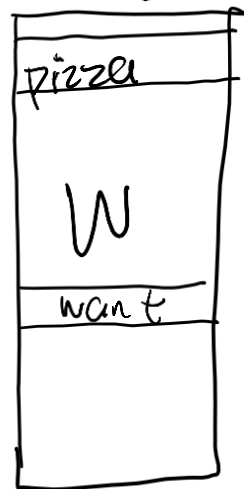# Word2Vec: Skig-gram and Negative Sampling

maybe I [want] pizza from
        target       context
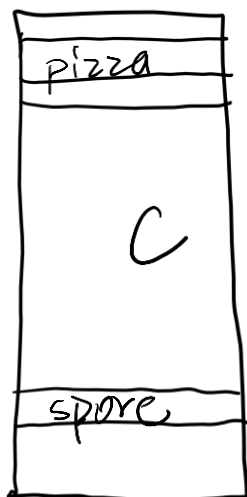
spore

want spore

want pizza

want I

want from

want maybo

positive

K   negative

128

128

50k



50k
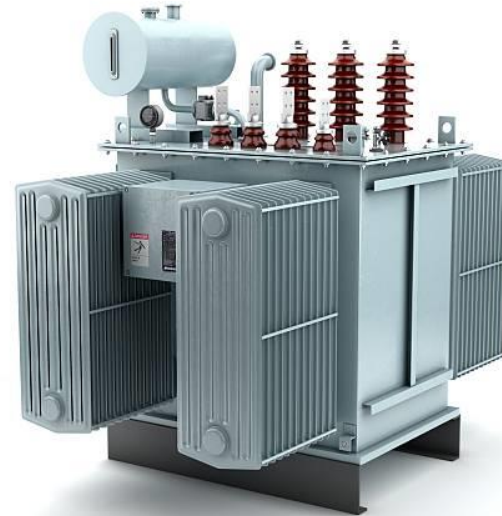
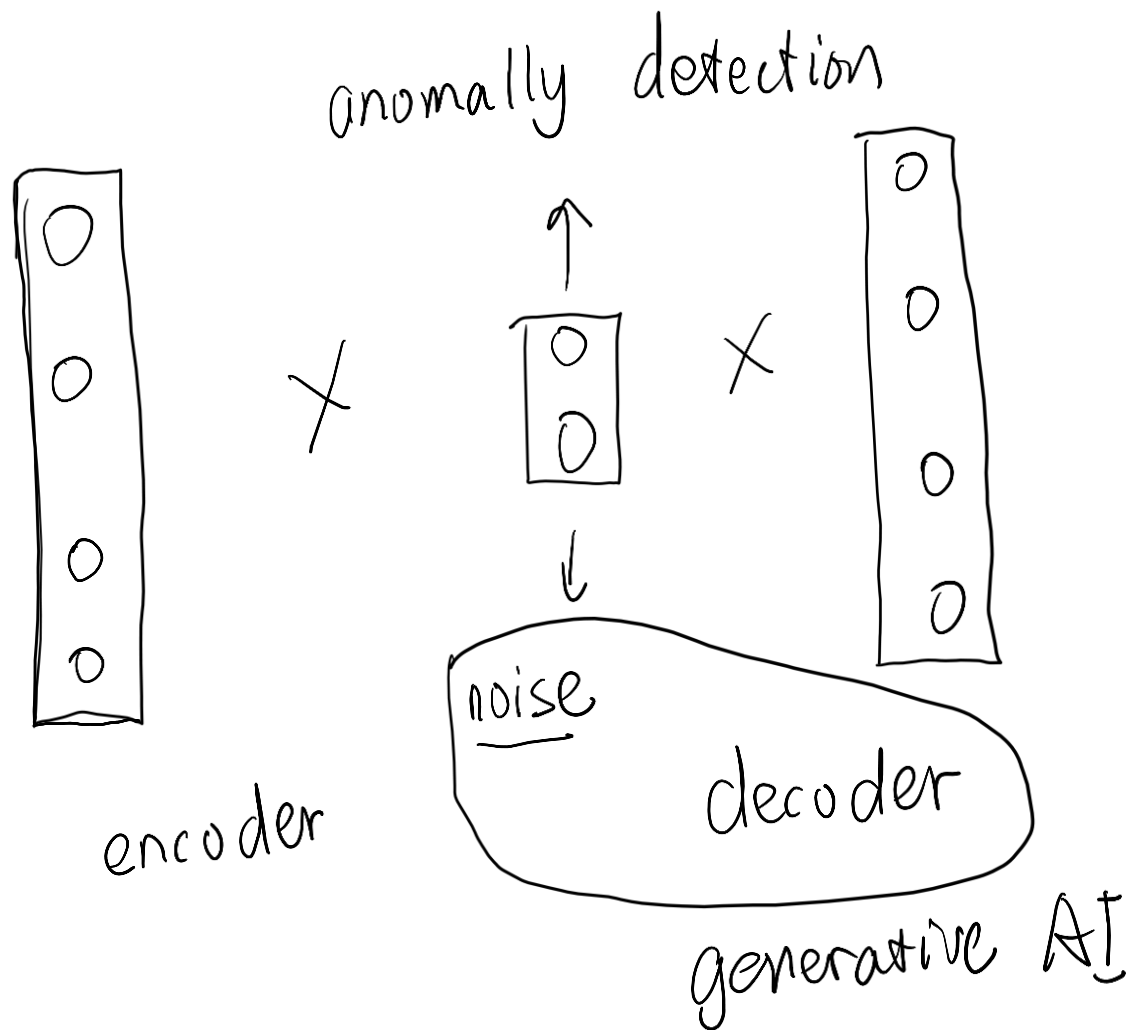$$\vec{w_1} \atop \vec{w_2}\} \vec{w}' = \vec{w} - \nabla \cdot LR \quad Q$$

# Transformers for Language Models of machine learning

# Transformer

- GPT: Generative Pre-trained Transformer
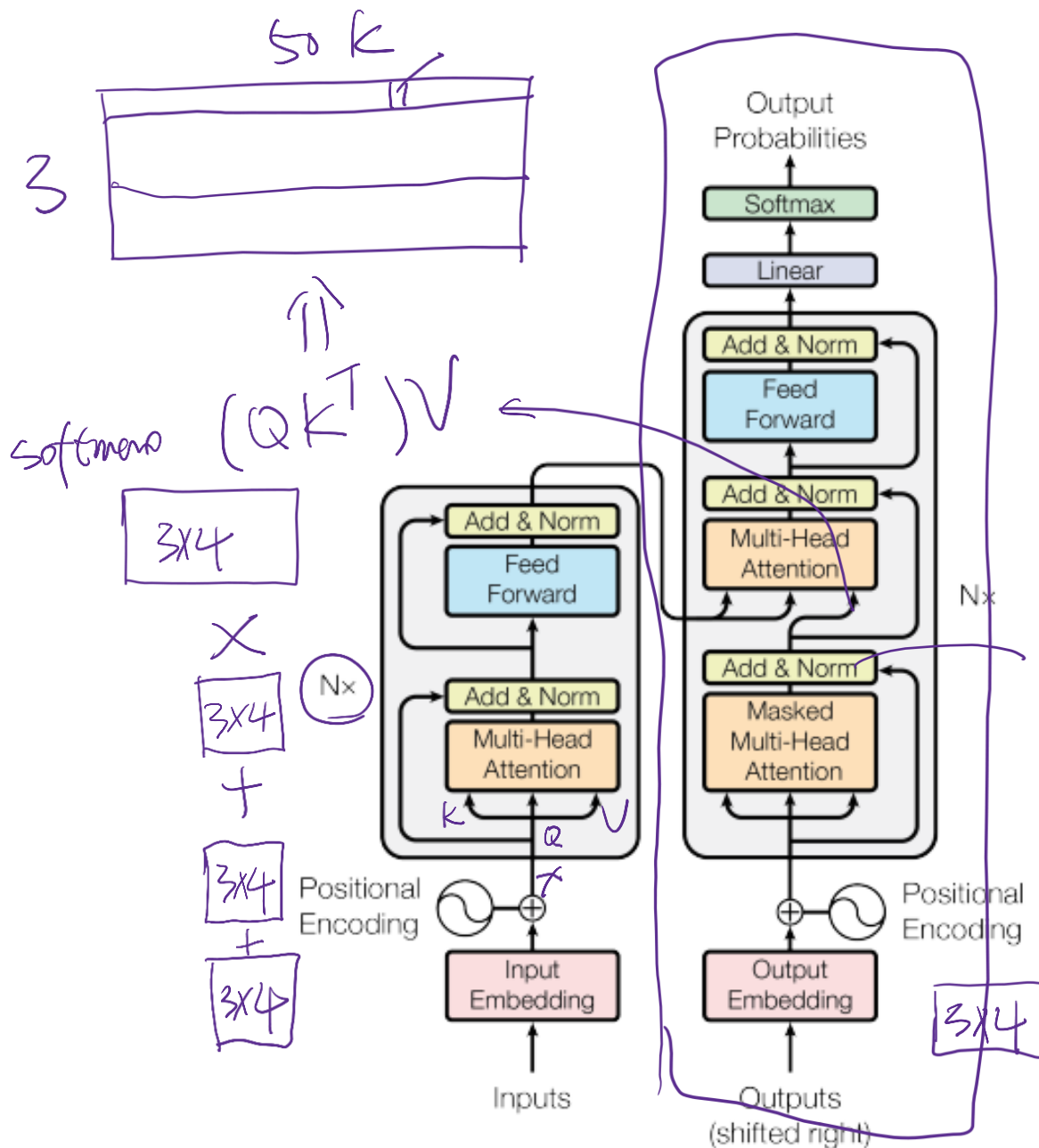- Vision Transformer: DETR
- Attention mechanism

anomally detection

MSE

CE

combination

encoder X noise X

decoder

generative AI

unsupervised

no ground truth

50 K

3

softmax $(QK^T)V$ ①

3x4

X

3x4 Nx

+

3x4

+

3x4

① translation

② there is ground truth labels to compare

3x4

GPT 1 & 2



Output Probabilities

Softmax

Linear

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

Nx

Add & Norm
Masked Multi-Head Attention

Positional Encoding

Output Embedding

Outputs (shifted right)

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

K    Q    V
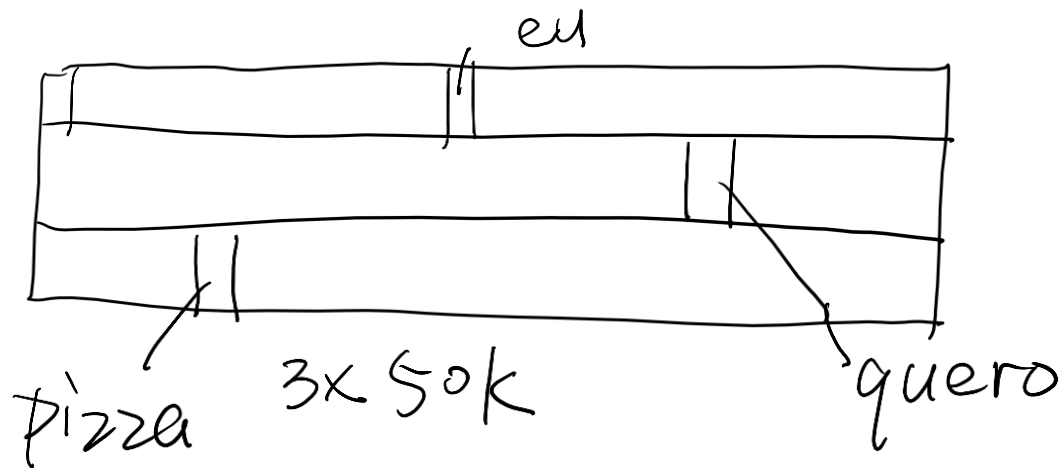
Nx

Positional Encoding

Input Embedding

Inputs

3x4
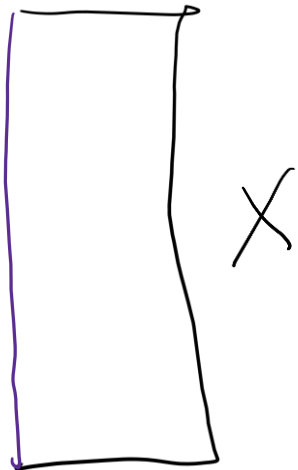
Vaswani, A. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*.
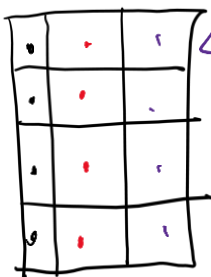
I
want
pizza

$3 \times 4$

K    Q    V

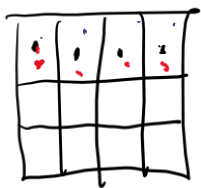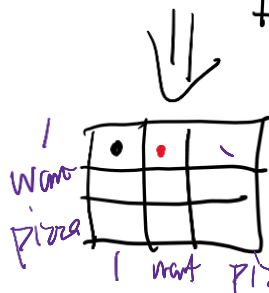$\rightarrow$    X

eu

pizza    $3 \times 50k$    quero

multihead head attention

softmax $(K \; Q^T) \; V$

$3 \times 4 \quad 4 \times 3 \qquad 3 \times 4$

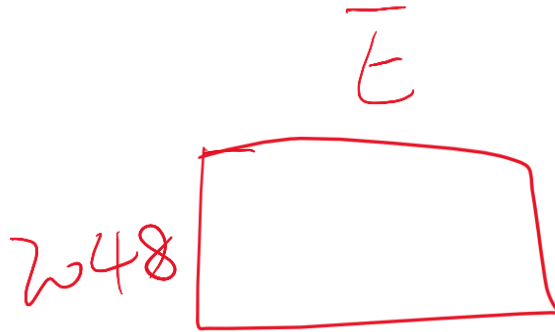$\overline{3 \times 3}$

$\boxed{3 \times 4} \times n$

I
want
pizza

I want pizza

attention

# GPT 3

- Parameters 175 B
- Dataset 45T
- 96 attention heads
- 2048 token size
- Learn from their chief scientist:
  https://www.youtube.com/watch?v=kCc8FmEb1nY