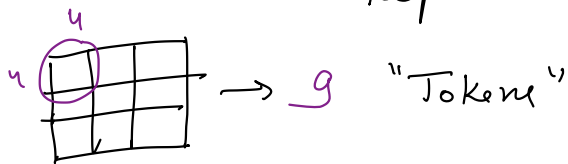


Agenda:

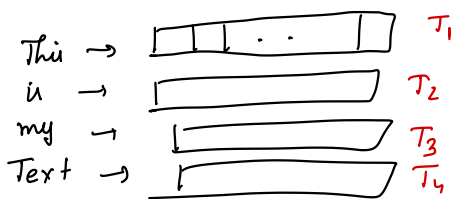
- ① Transformer (✓)
- ② Foundational / Multimodal (✓)
- ③ Foundational (Unimodal)
 - ↓
 - GPT
 - ↓
 - VIT
- ④ Recipes for Multimodality
- ⑤ Research Works

→ Transformer

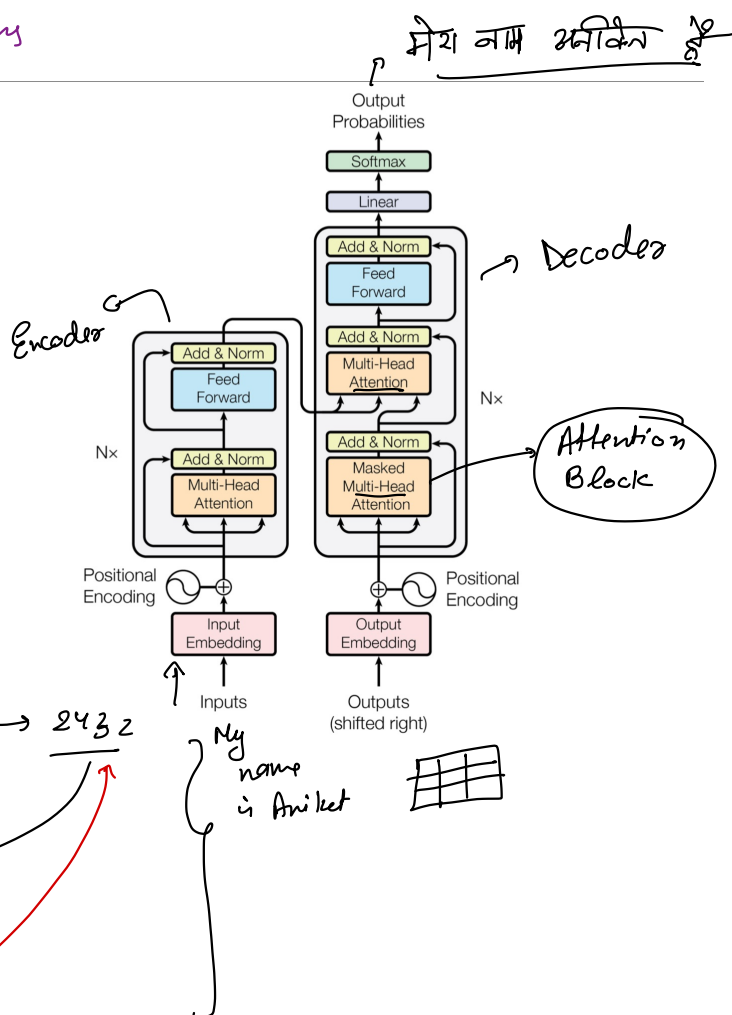
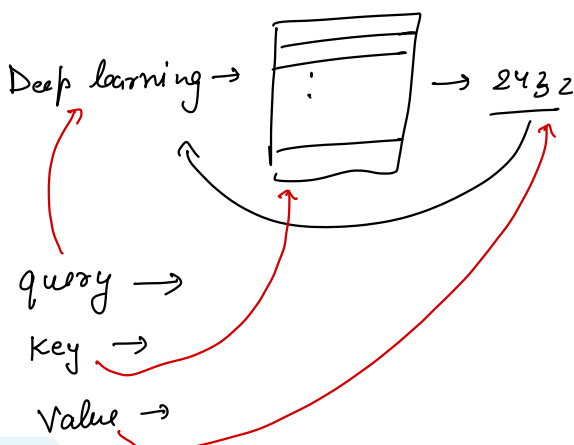
Representation

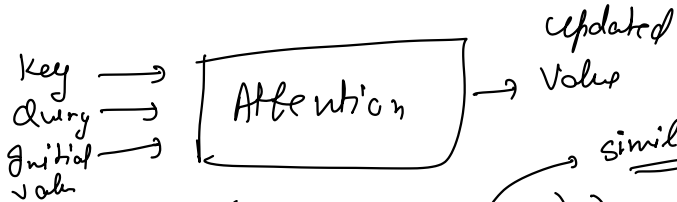
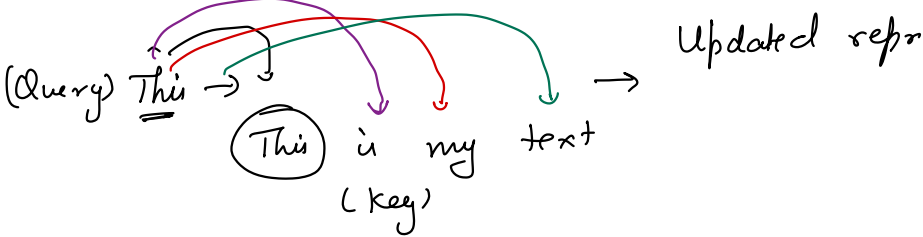
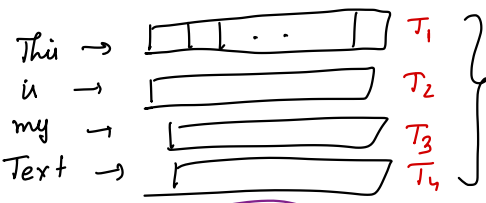


"This is a class" → 4 tokens



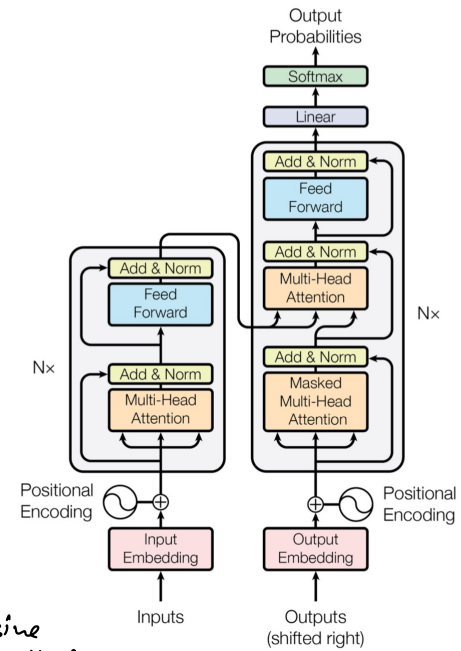
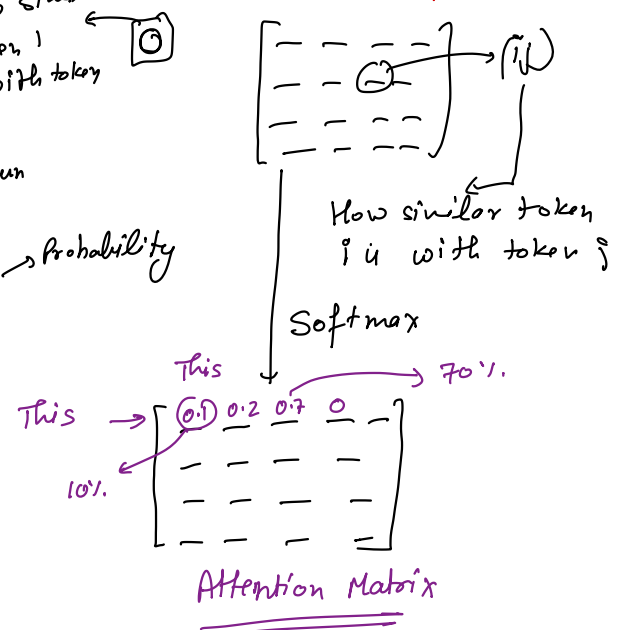
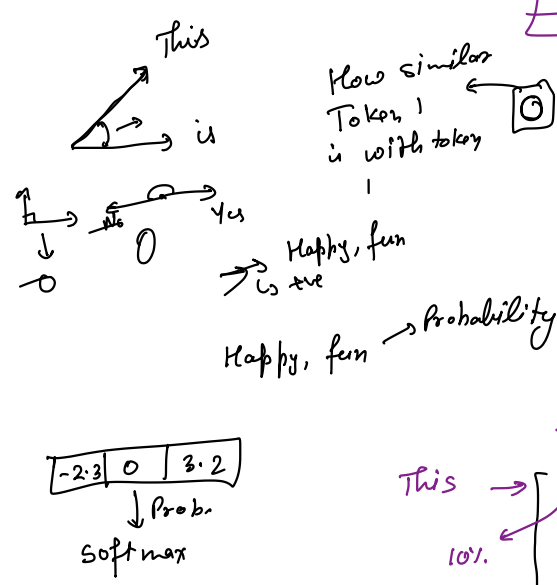
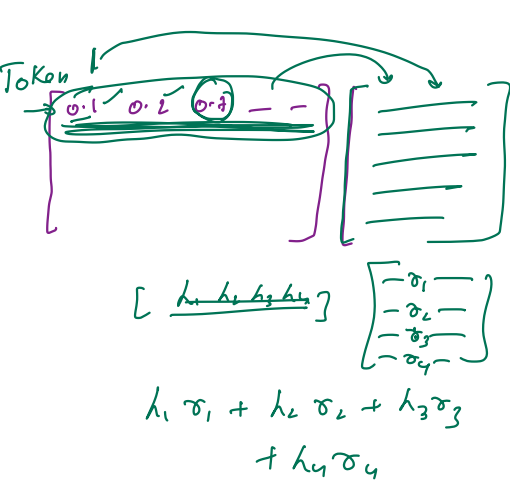
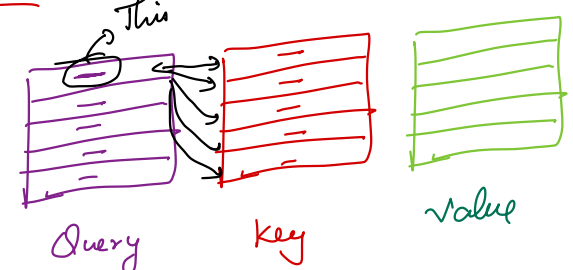
o Attention (Self-Attention)





$$\text{Probab} = \text{Softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$$

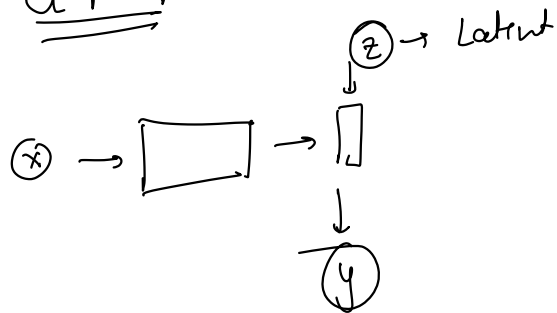
similarity (via cosine similarity)
 Normalization
 Updated values



o Foundation Model

- ① Trained on a large-scale data
- ② It is able to perform multiple tasks

⇒ GPT

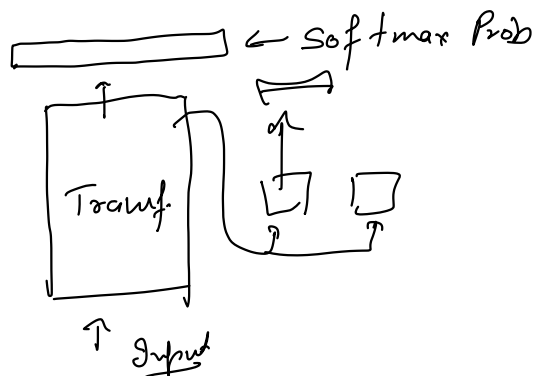


• Self-supervised

↳ Labels are generated from data itself.

→ Next word prediction

"Labels are generated itself, then it is called self-supervision → output."



Prediction: Generative

Stackoverflow GMDP
↓ ↓
Q: — ?
↳ Answer

Labels are gen — then it is called

Yes!

Mixture of Experts.

