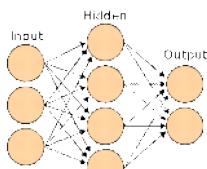
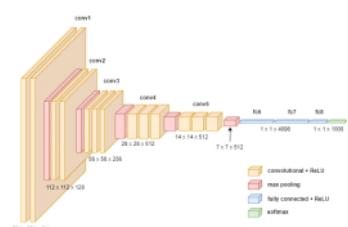


What is Transformer?

27 January 2024 18:41



ANN
tabular



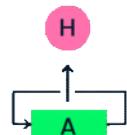
CNN
image

→ [Transformer] → seq2seq task

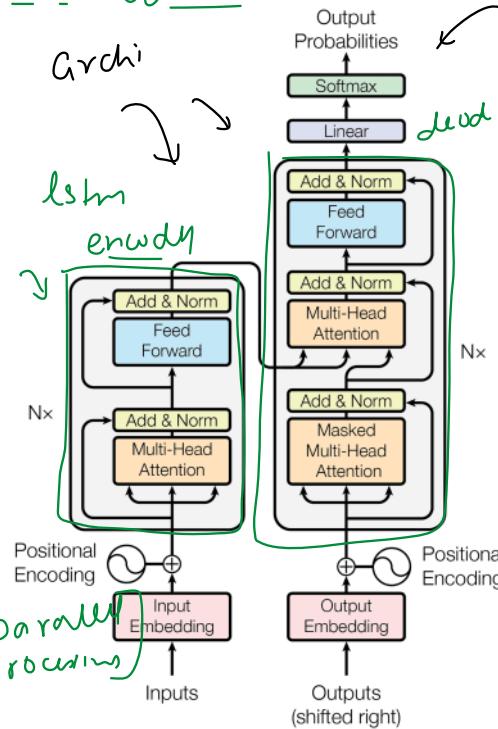
NN arch

+ machine translation
+ question ans
+ text summariz

→ Self attention stable



RNN
Sequence
→ Text



Attention Is All You Need

2017

→ Google Brain

Deep learnis

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

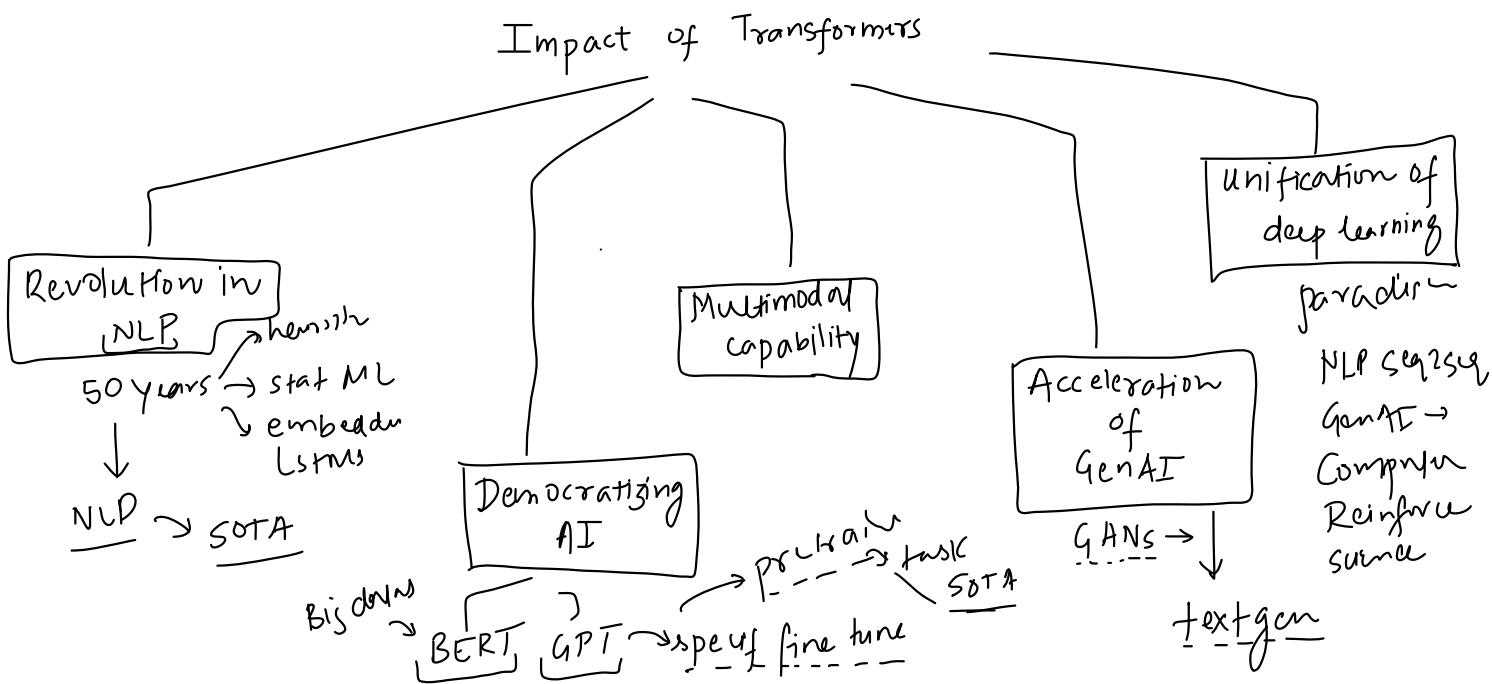
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Impact of Transformers

27 January 2024 20:03



The Origin Story!

27 January 2024 22:38

2014-15 → seq2seq machine

Sequence to Sequence Learning with Neural Networks

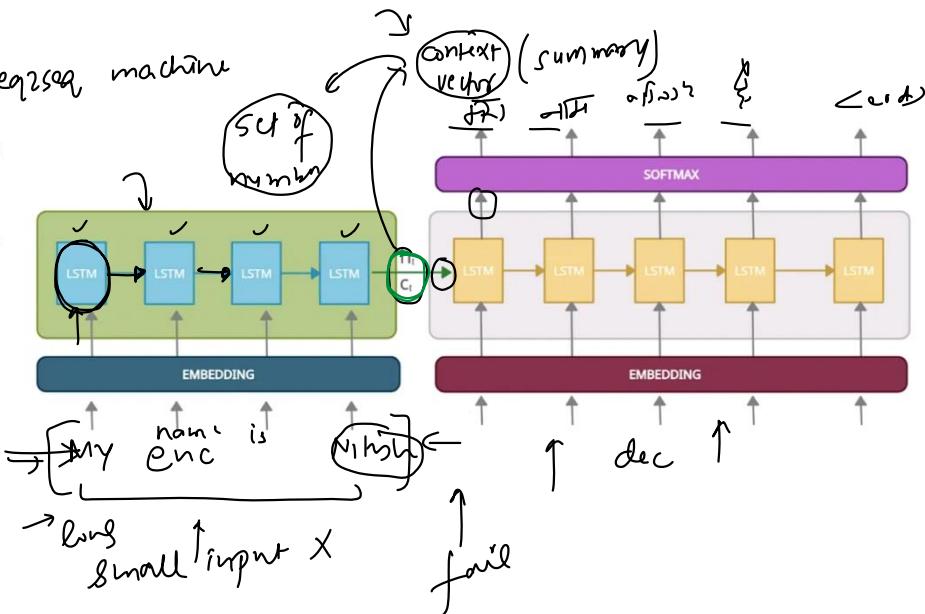
Ilya Sutskever
Google
ilyasu@google.com

Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
qvl@google.com

Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that an English to French translation task from the WMT'14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous best result on this task. The LSTM also learned sensible phrase and sentence representations that are sensitive to word order and are relatively invariant to the active and the passive voice. Finally, we found that reversing the order of the words in all source sentences (but not target sentences) improved the LSTM's performance markedly, because doing so introduced many short term dependencies between the source and the target sentence which made the optimization problem easier.



{ NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE }

Dzmitry Bahdanau
Jacobs University Bremen, Germany

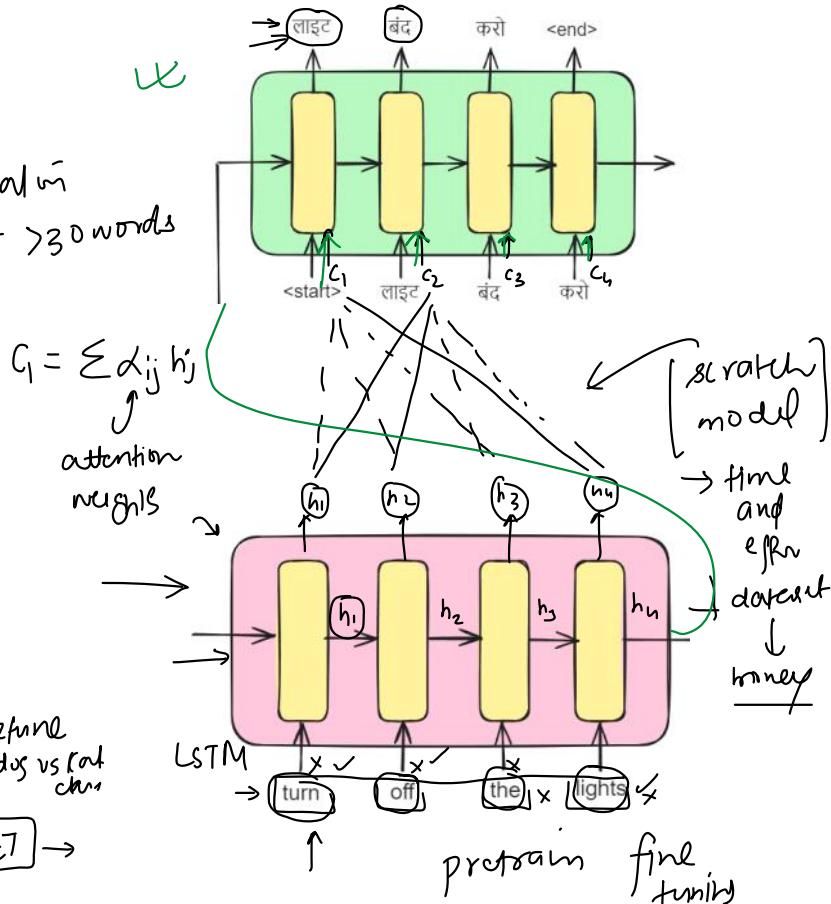
KyungHyun Cho Yoshua Bengio*
Université de Montréal

ABSTRACT

Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-alignments) found by the model agree well with our intuition.

(sequential) → slow → huge dataset

Transfer learning
→ CNN → dog vs cat
→ ImageNet →



research → incremental time travel

[Attention Is All You Need] → 2017

→ LSTM / self → parallelly
→ stable
→ hyperfine
→ NLP → BERT → fine tune
Output Probabilities

→ research → train

[Attention Is All You Need] → 2017



Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

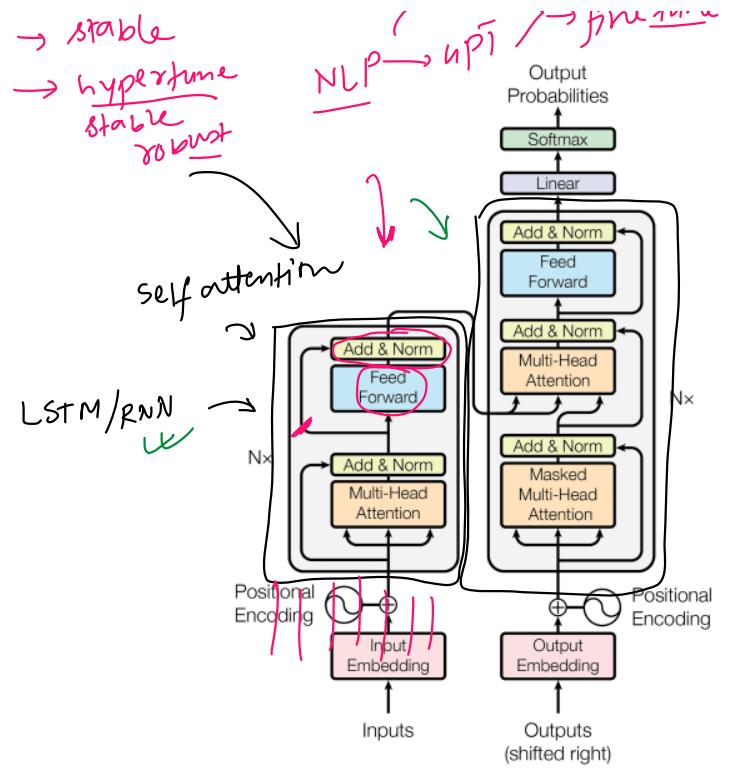
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.



The Timeline

28 January 2024 00:55

2000 - 2014 → RNNs / LSTMs



2014 → Attention



2017 → Transformer

2018 → BERT / GPT (Transfer learning)

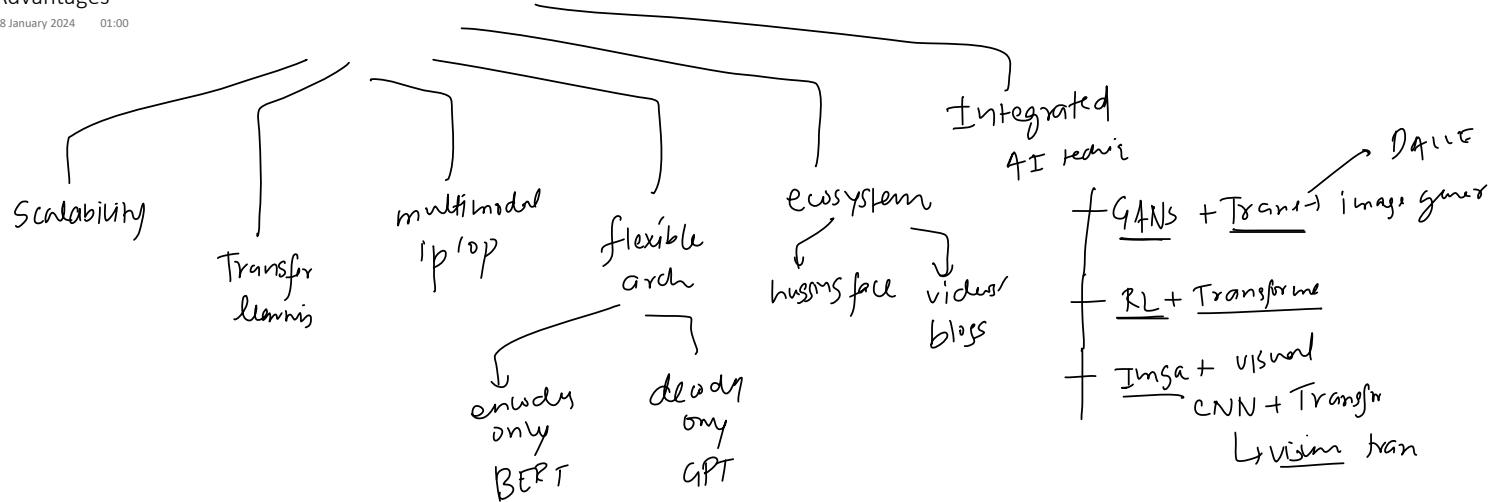
2018 - 2020 → Vision Transfer / AlphaFold-2

2021 → Gen AI

2022 → ChatGPT / Stable Diffusion

Advantages

28 January 2024 01:00

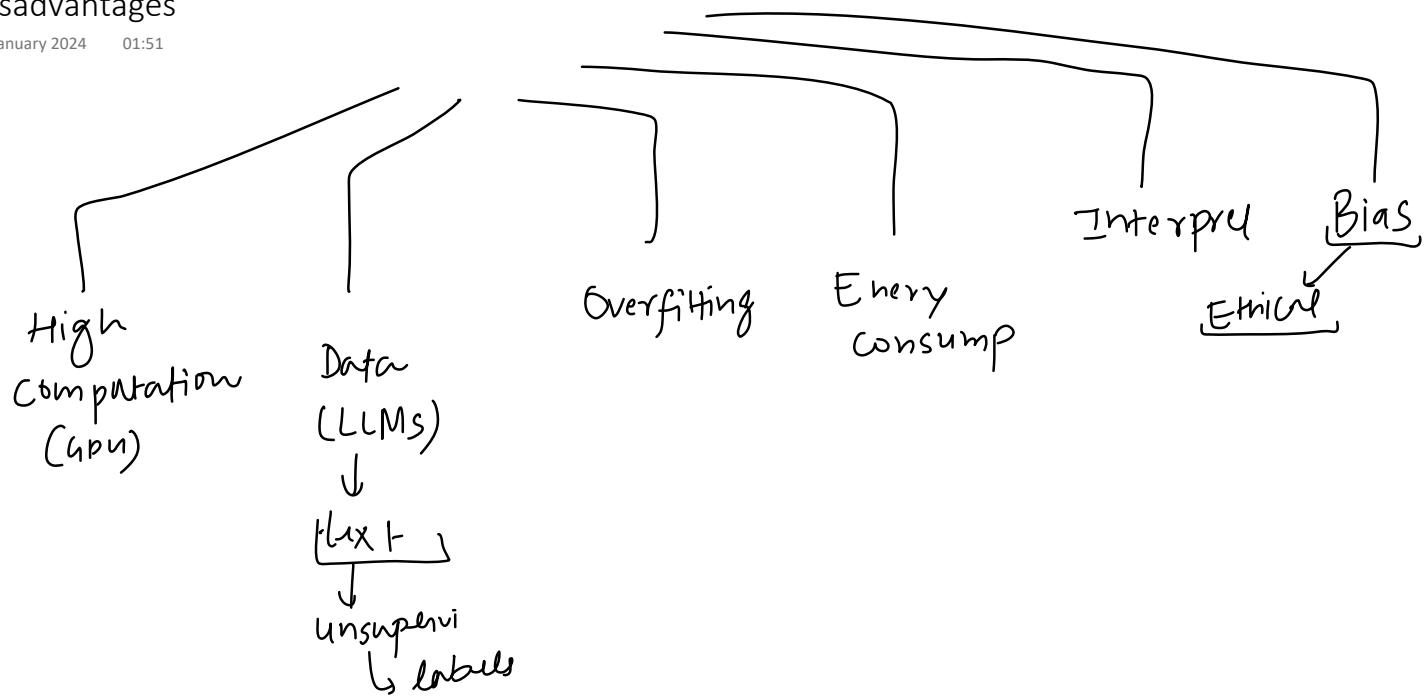


Famous Applications

28 January 2024 01:17

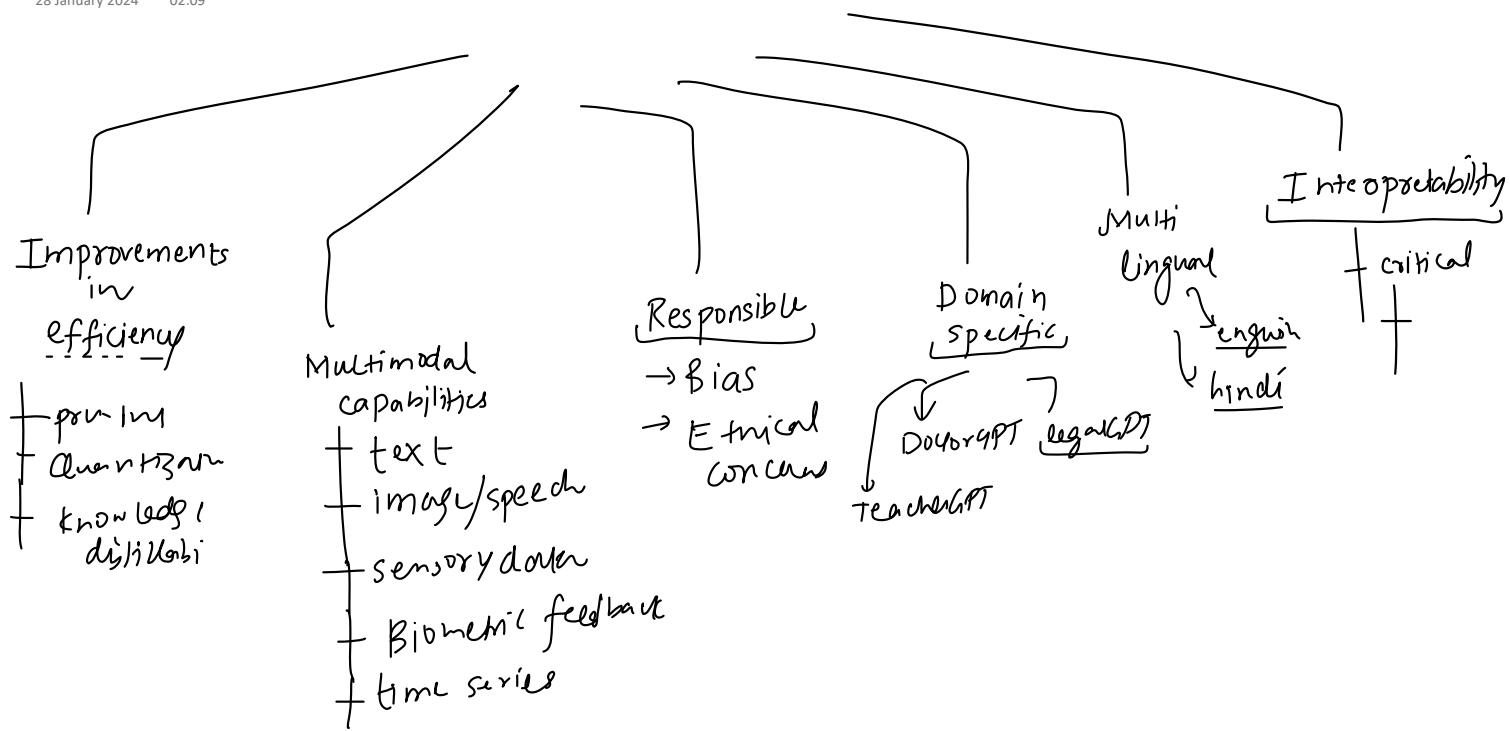
Disadvantages

28 January 2024 01:51



Future

28 January 2024 02:09



The What

Self attention

computer X

NLP → words → numbers → vectorization

OHE	mat	cat	mat
1)	mat	cat	mat
2)	cat	rat	rat

num → [1 0 0] [0 1 0] [1 0 0]

Bow
SI
S2

mat rat cat
[2 0 1]
→ [0 2 1]

	mat	cat	rat
mat	1	0	0
cat	0	1	0
rat	0	0	1

Tf Idf

word embedding

number

semantic meaning

training with → [0 0 0 0 0] → [0 0 0 0 0] 5 dim

256, 512
64

[0.2 0.9 0.9 0.9 0.9] n-dim vector

cricket

king → [0.9 0.1 1 0 0.9]
queen → [0.9 0.2 0.4 1 0] → similar

royalty

meaning → vector

Apple → vector

The problem of "Average Meaning"

- 1) An apple a day keeps the doctor away
 - 2) Apple is healthy
 - 3) Apple is better than orange
 - 4) Apple makes great phones
 - ⋮
 - ⋮
- 10000
9000
fruits
1000 phones
9000 term
10000 fruits

dataset



0 0 0 → [0 0 0]

task

2 dim

x y
task technology

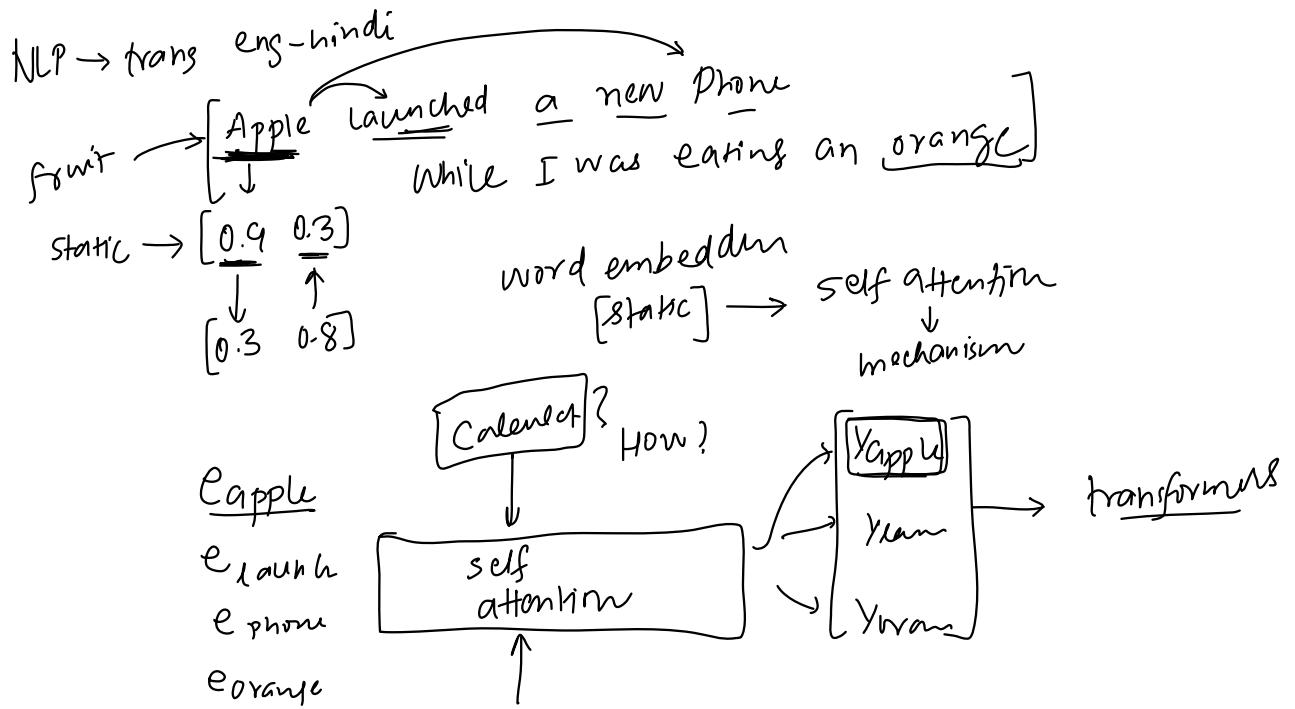
[0.8 0.2] → [0.9 0.3]

word embeddings

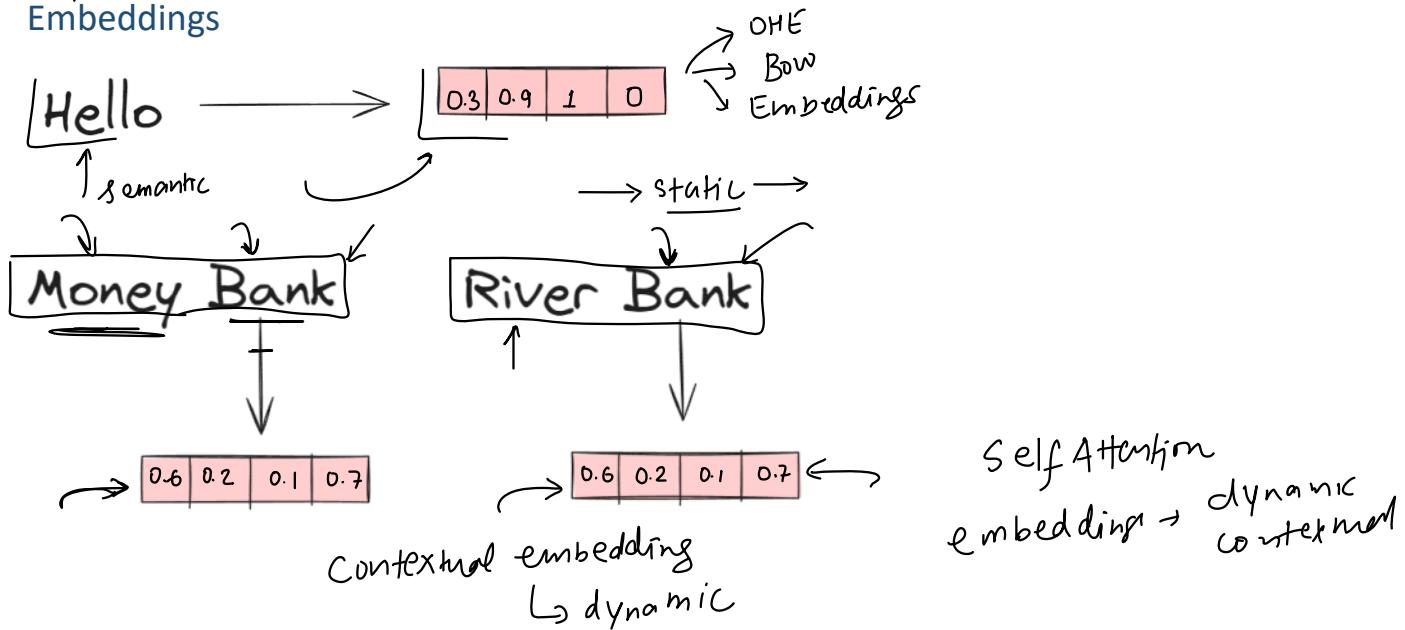
create → use → static

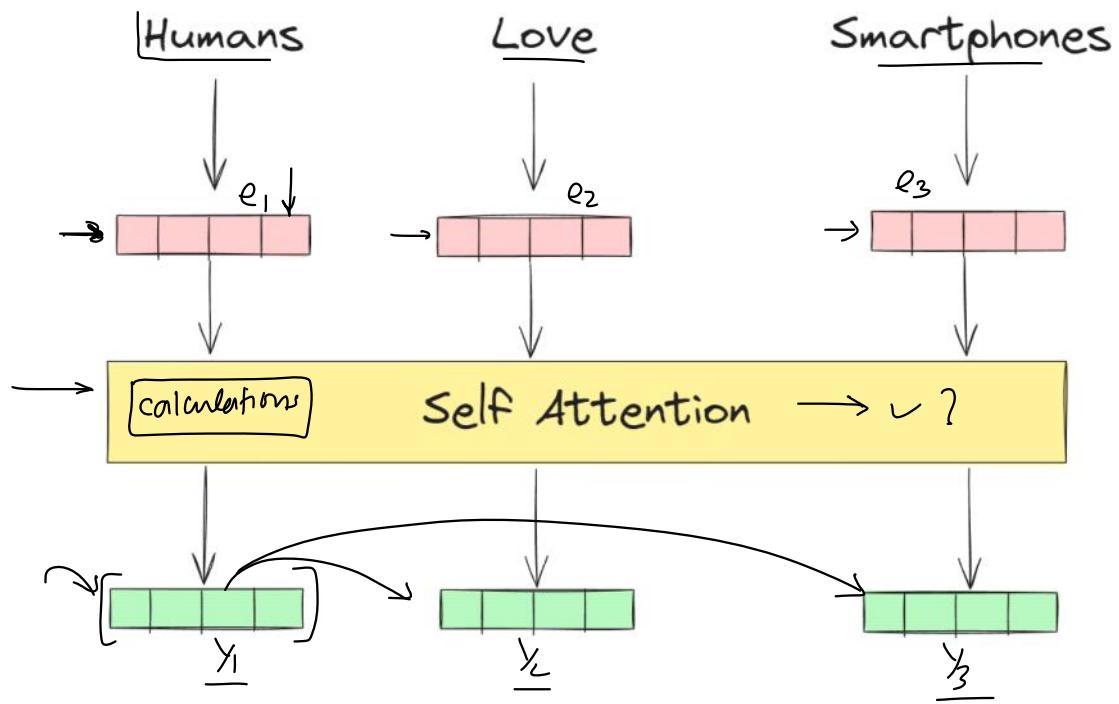
static embedding → [8 mark-
contextual embedding]

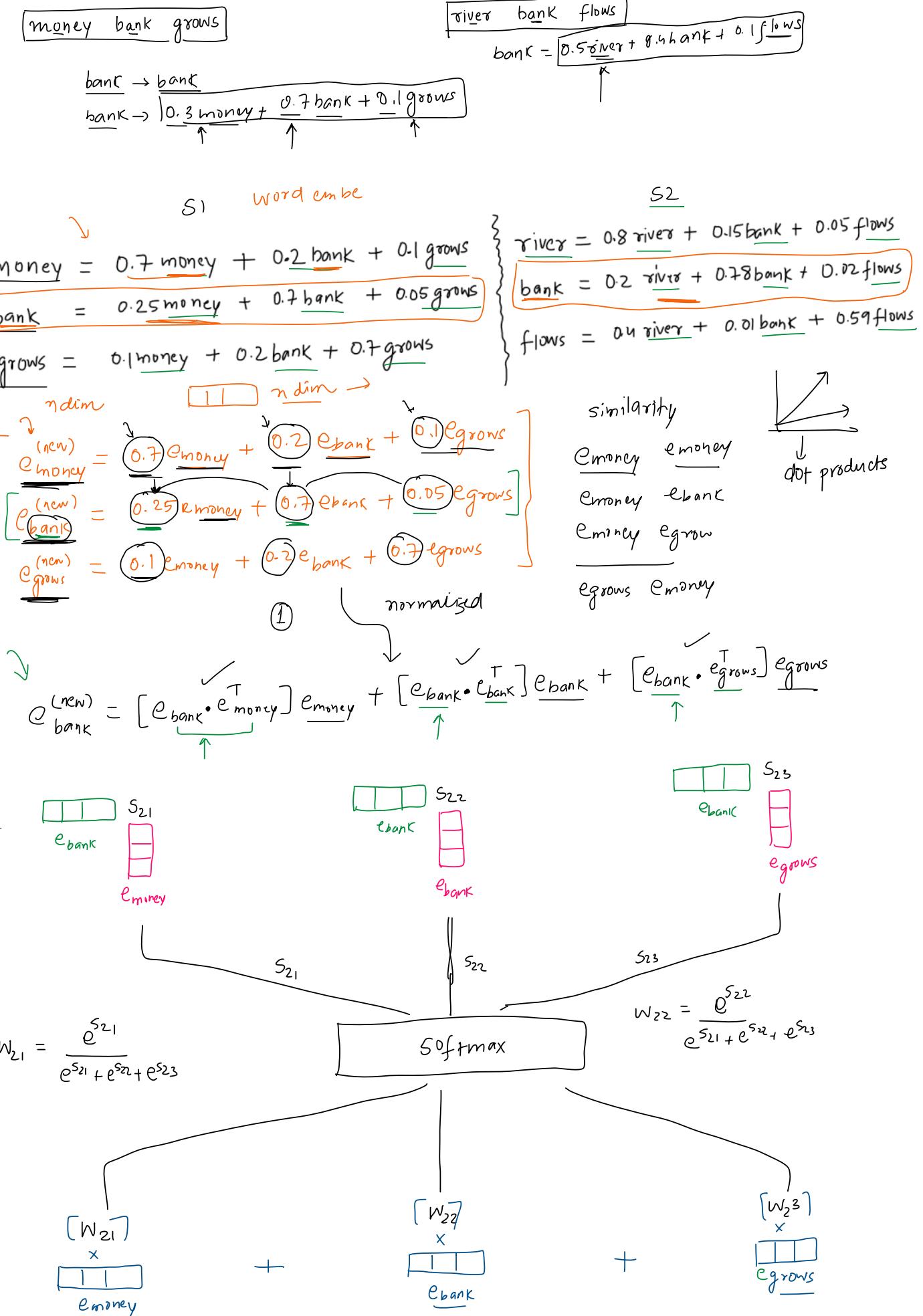
110 - drane eng-hindi

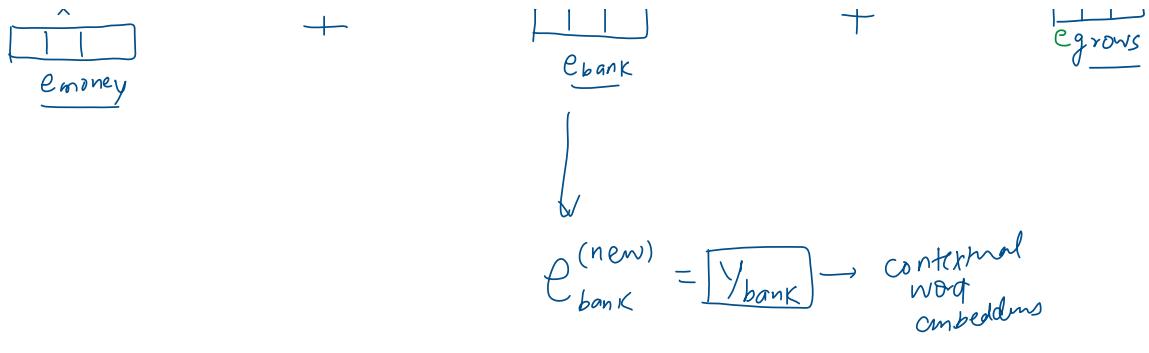


Embeddings



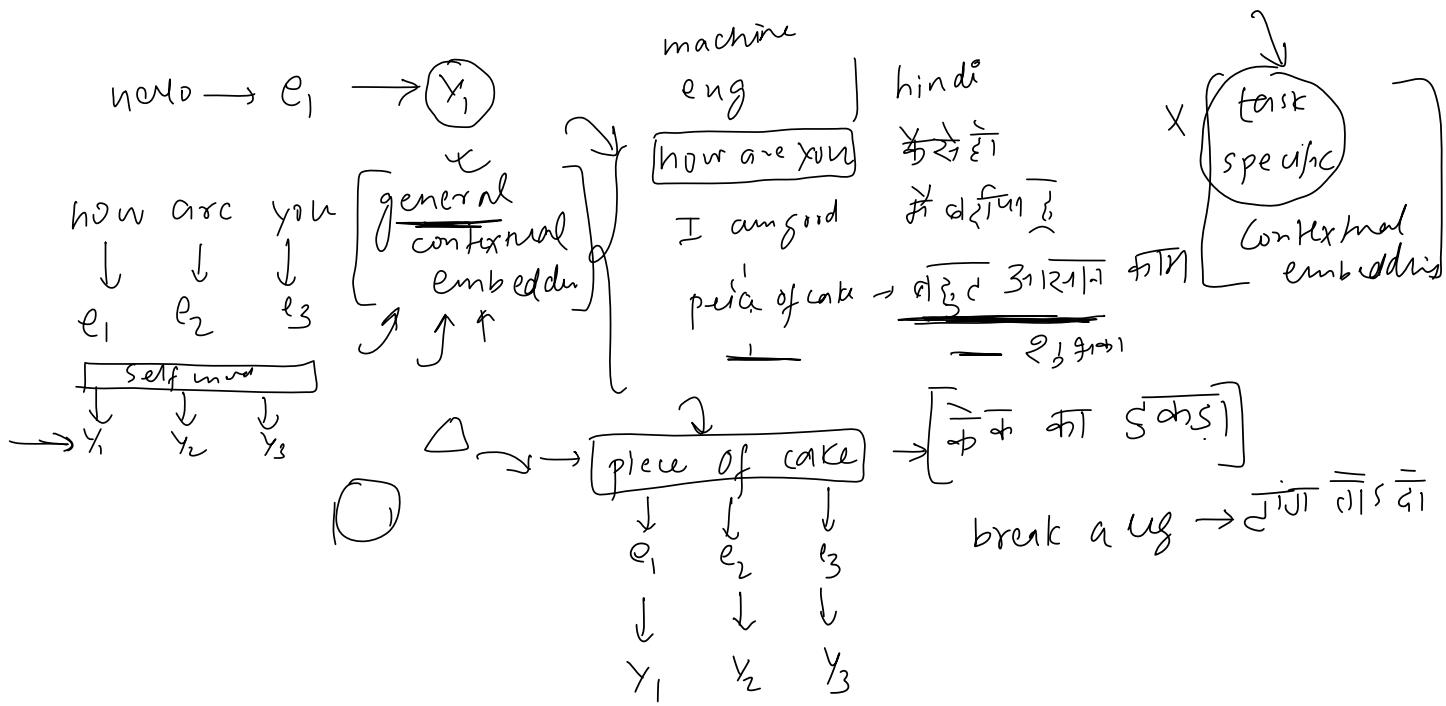
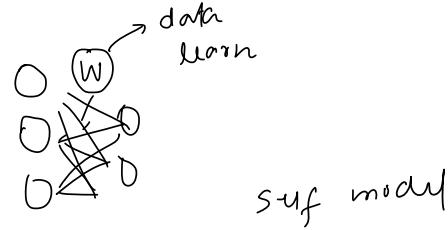






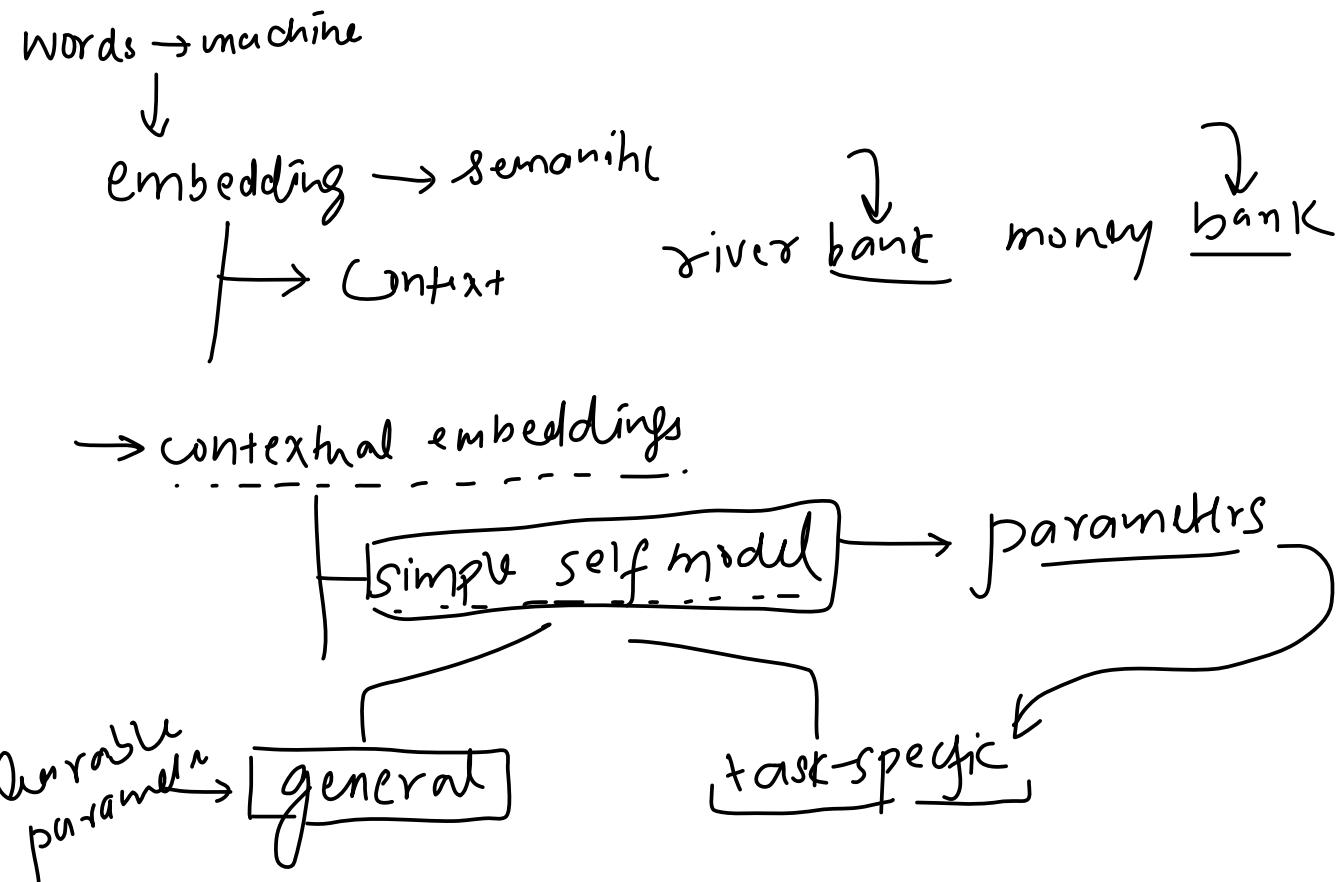
Points to consider

- This operation is a parallel operation
- There are no parameters involved



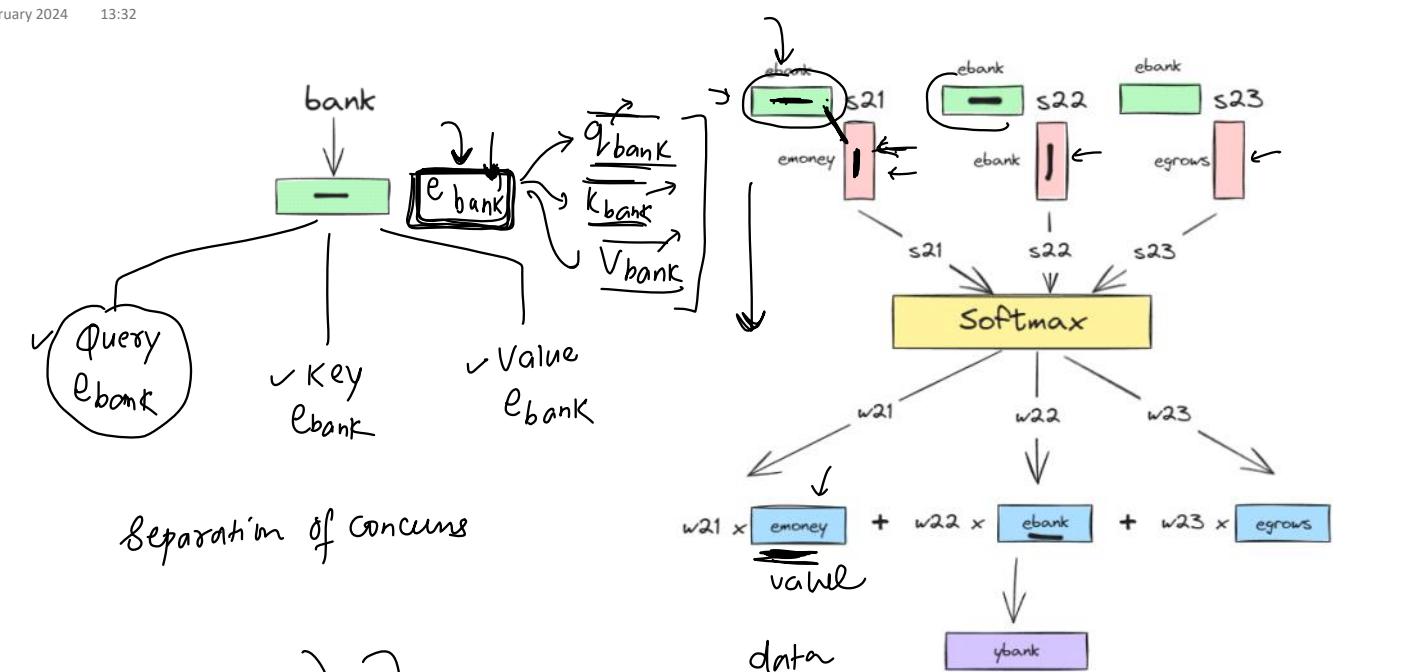
Progress

06 February 2024 00:42

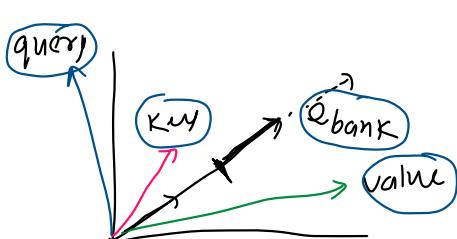
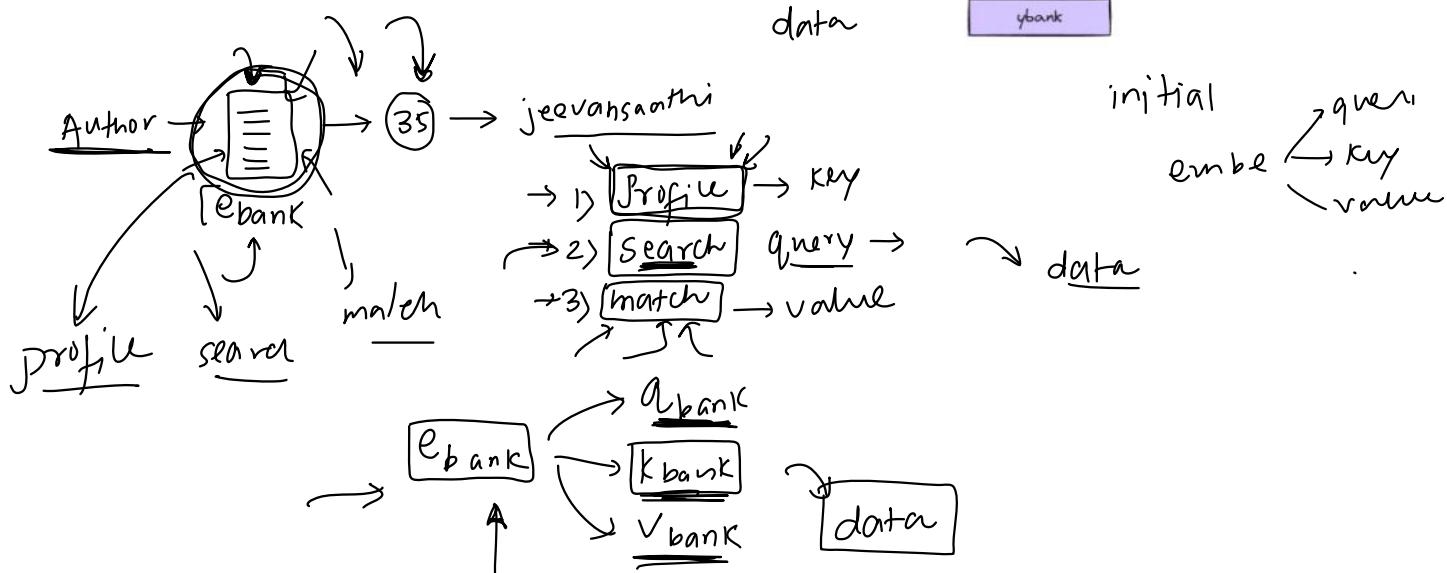


Query, Key & Value Vectors

06 February 2024 13:32



Separation of Concurrence

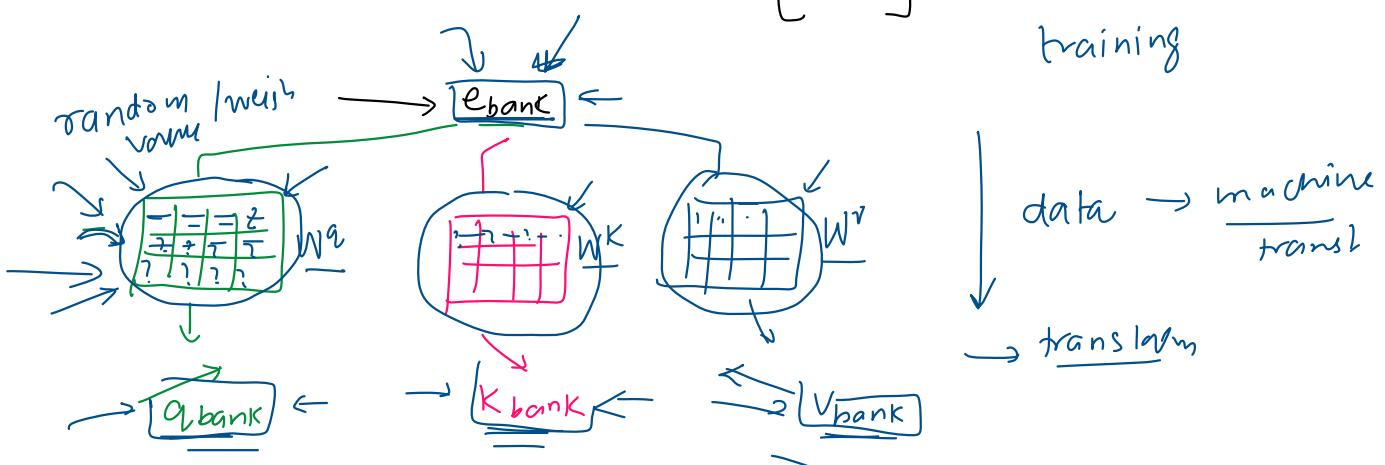


d magnitude (Scaling)

linear transform

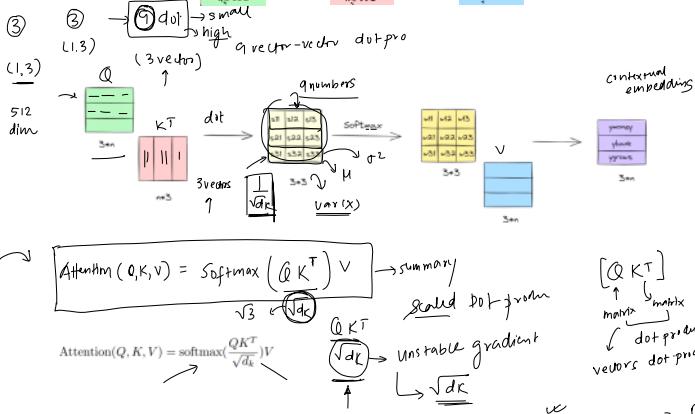
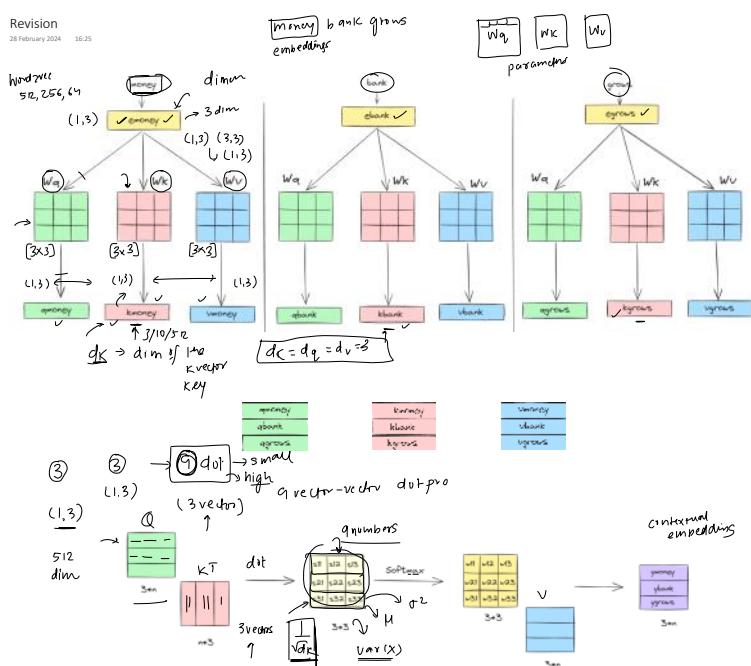
$$\rightarrow \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

training

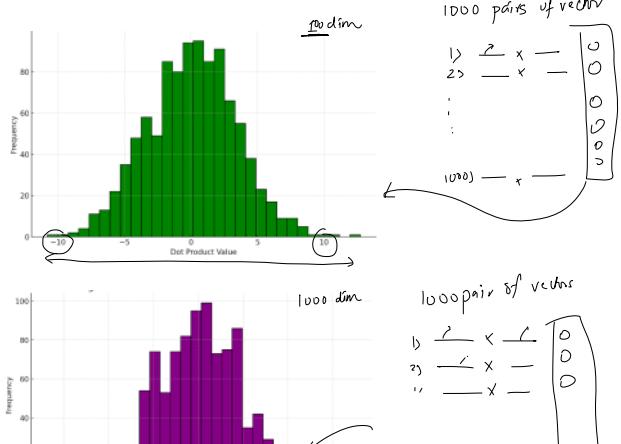
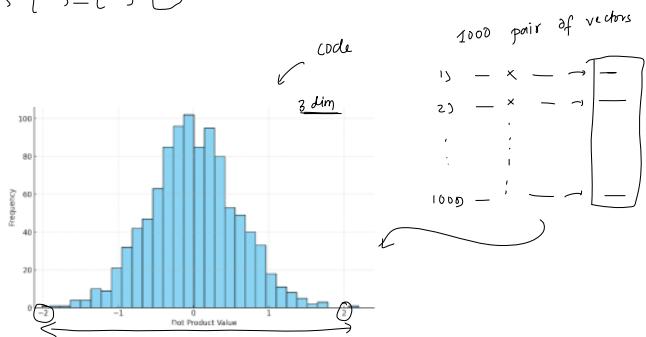


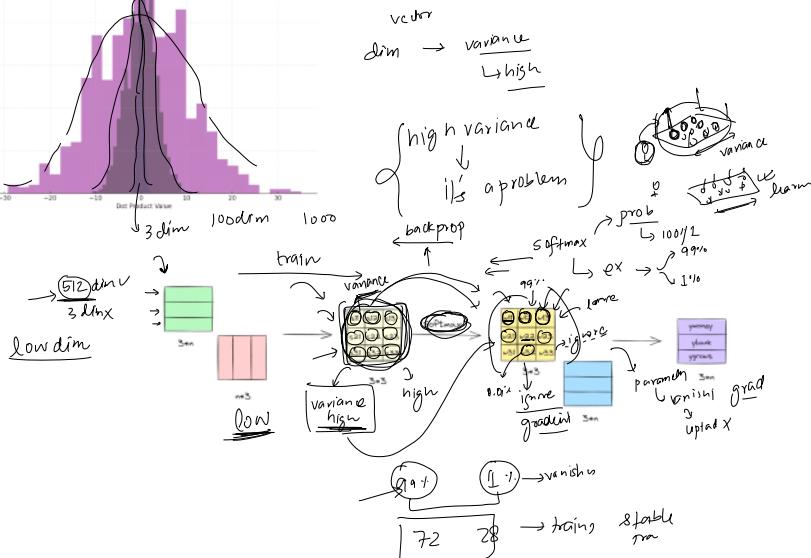
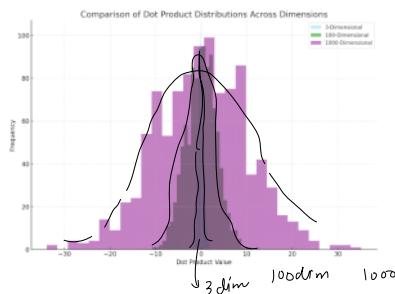
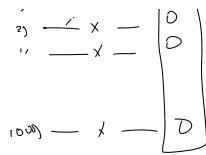
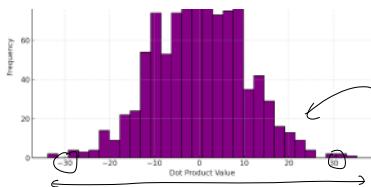
$3 \times 3 \rightarrow \text{dot} \rightarrow \textcircled{9}$ values

Revision
28 February 2024 16:25



$$\begin{array}{l} 1) [1, 2, 3] - [3, 2, 1] \\ 2) [] - [] \\ 3) [] - [] \\ 4) [] - [] \\ 5) [] - [] \end{array}$$



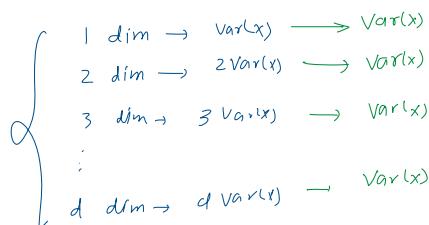
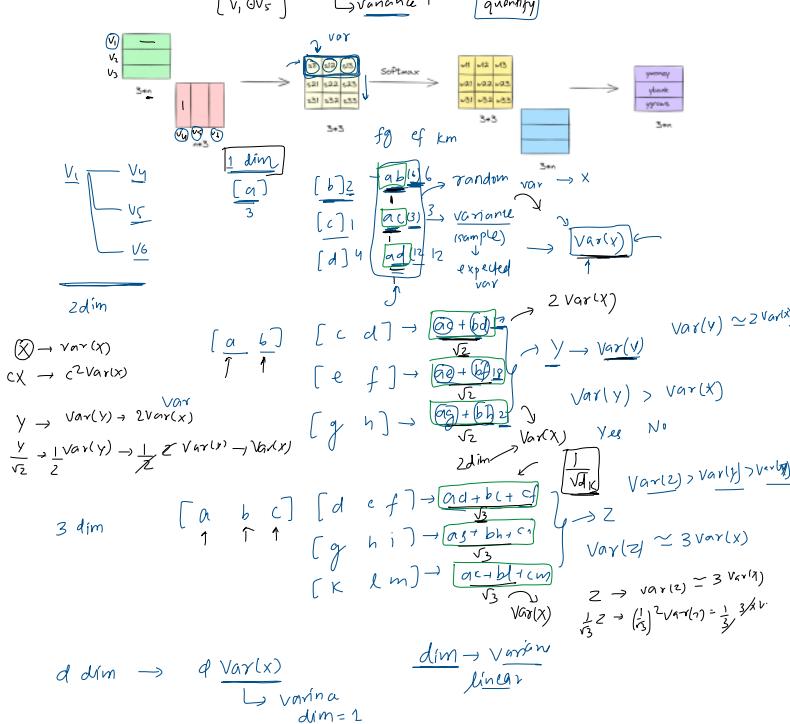


$c \propto \text{Var}(X)$

$$\begin{cases} X \rightarrow \text{Var}(X) \\ Y \rightarrow cX \rightarrow c^2\text{Var}(X) \end{cases}$$

If you have a random variable X with a variance of $\text{Var}(X)$, and you create a new variable Y by scaling X with a constant c , so that $Y = cX$, the variance of Y ($\text{Var}(Y)$) is related to the variance of X by the square of the scaling factor c . Mathematically, this relationship is expressed as:

$$\text{Var}(Y) = c^2\text{Var}(X)$$



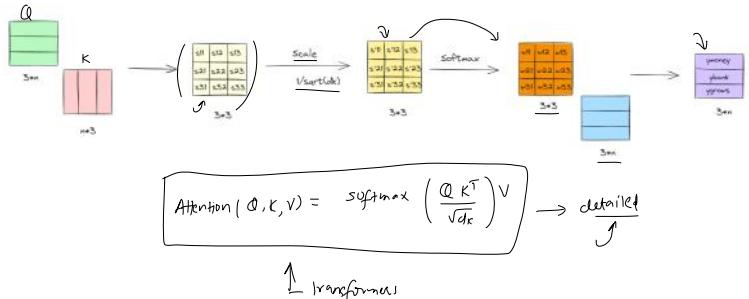
$$[\underbrace{\dots}_{d \text{ dim}}] \rightarrow \text{Var}(x)$$



$\overbrace{\dots}^d$

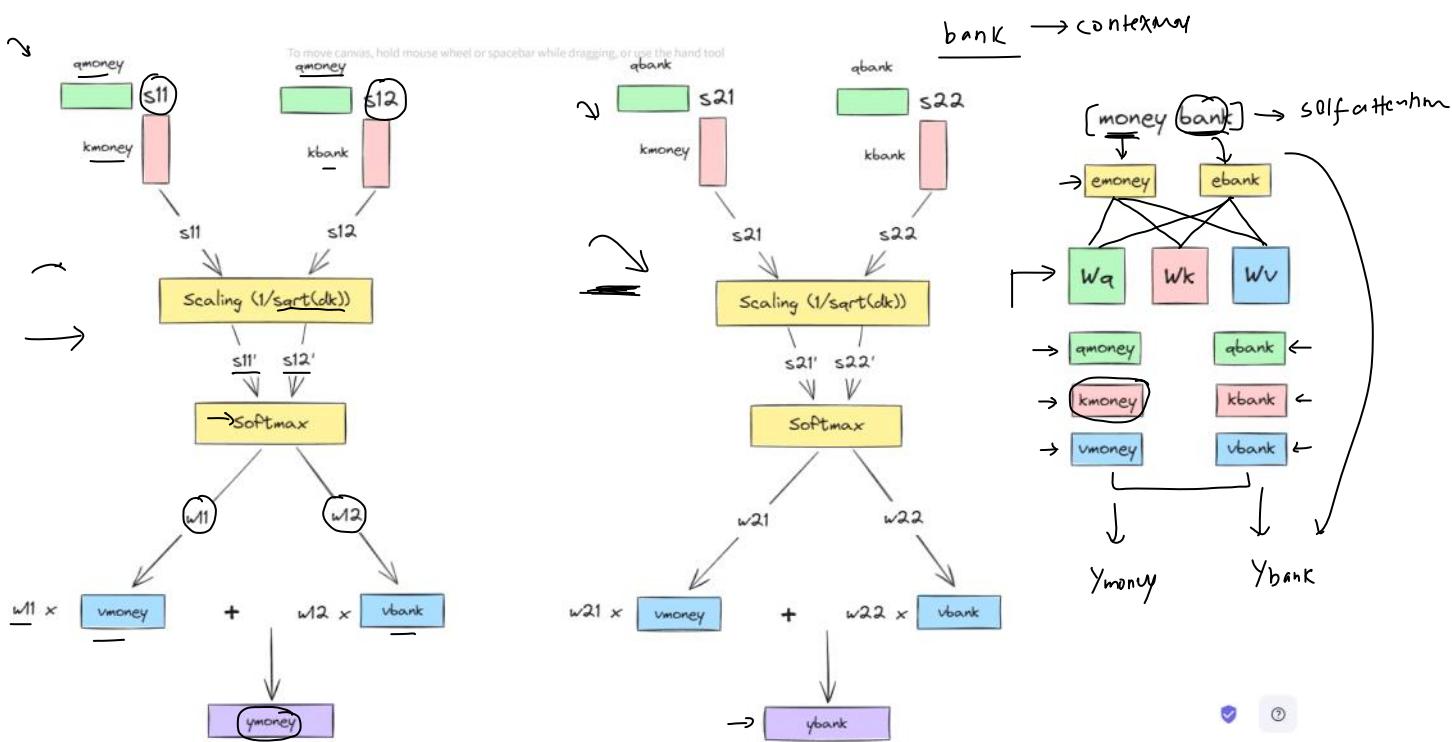
γ

$\sim \dots$



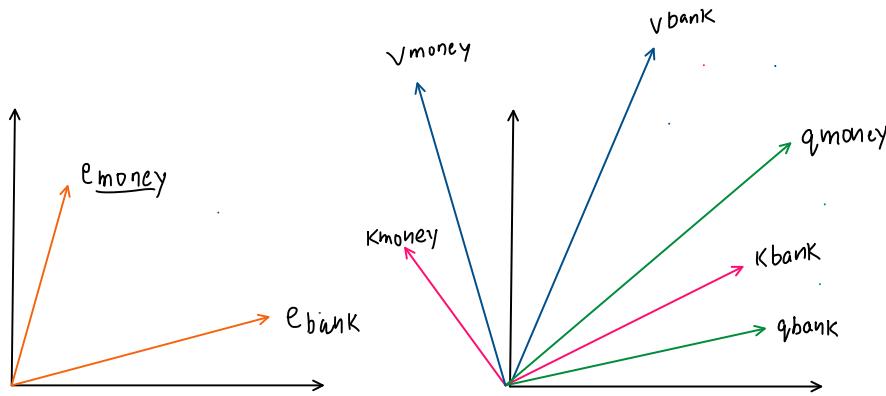
What is d_k

28 February 2024 16:59



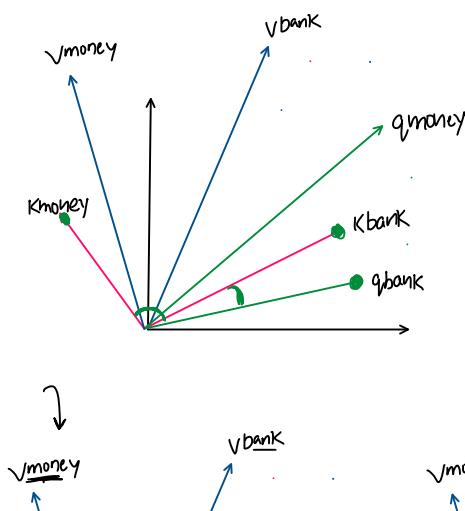
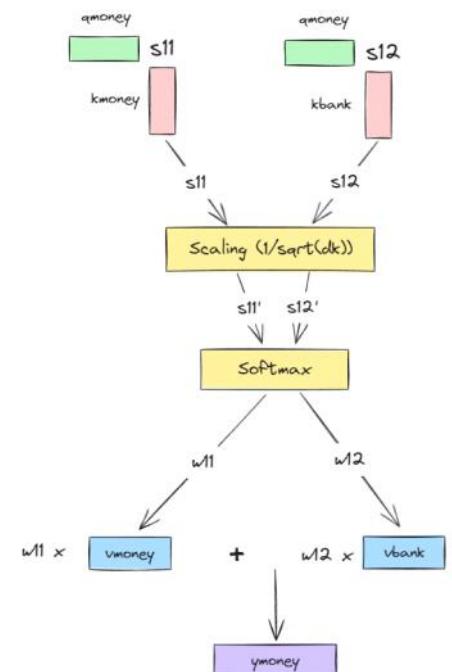
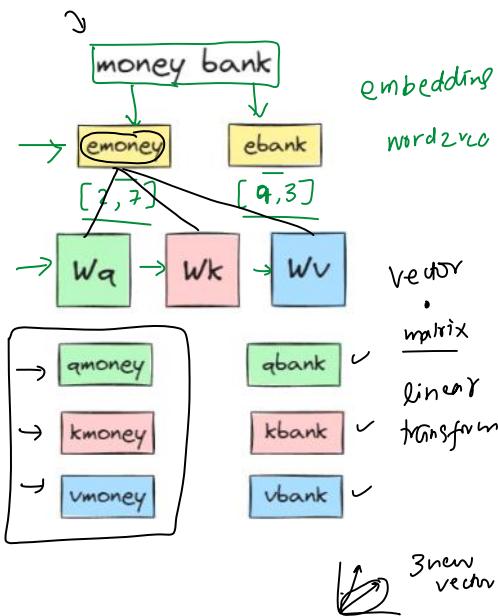
Geometric Intuition

08 March 2024 15:16



$$Wq = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad Wk = \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} \quad Wv = \begin{bmatrix} 4 & 1 \\ 2 & 1 \end{bmatrix}$$

* All values are hypothetical



[Dot Product]

$$S_{21} = 10$$

$$S_{22} = \frac{32}{\sqrt{2}}$$

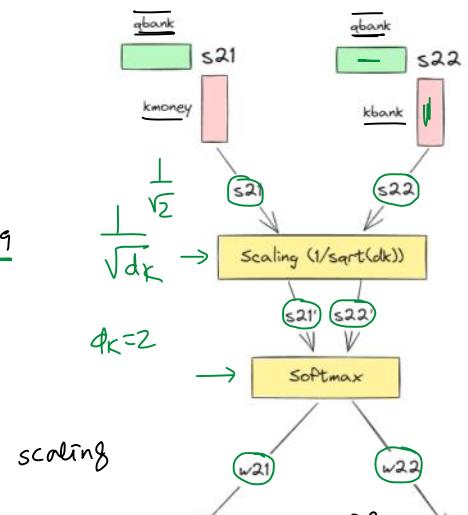
[Scaling]

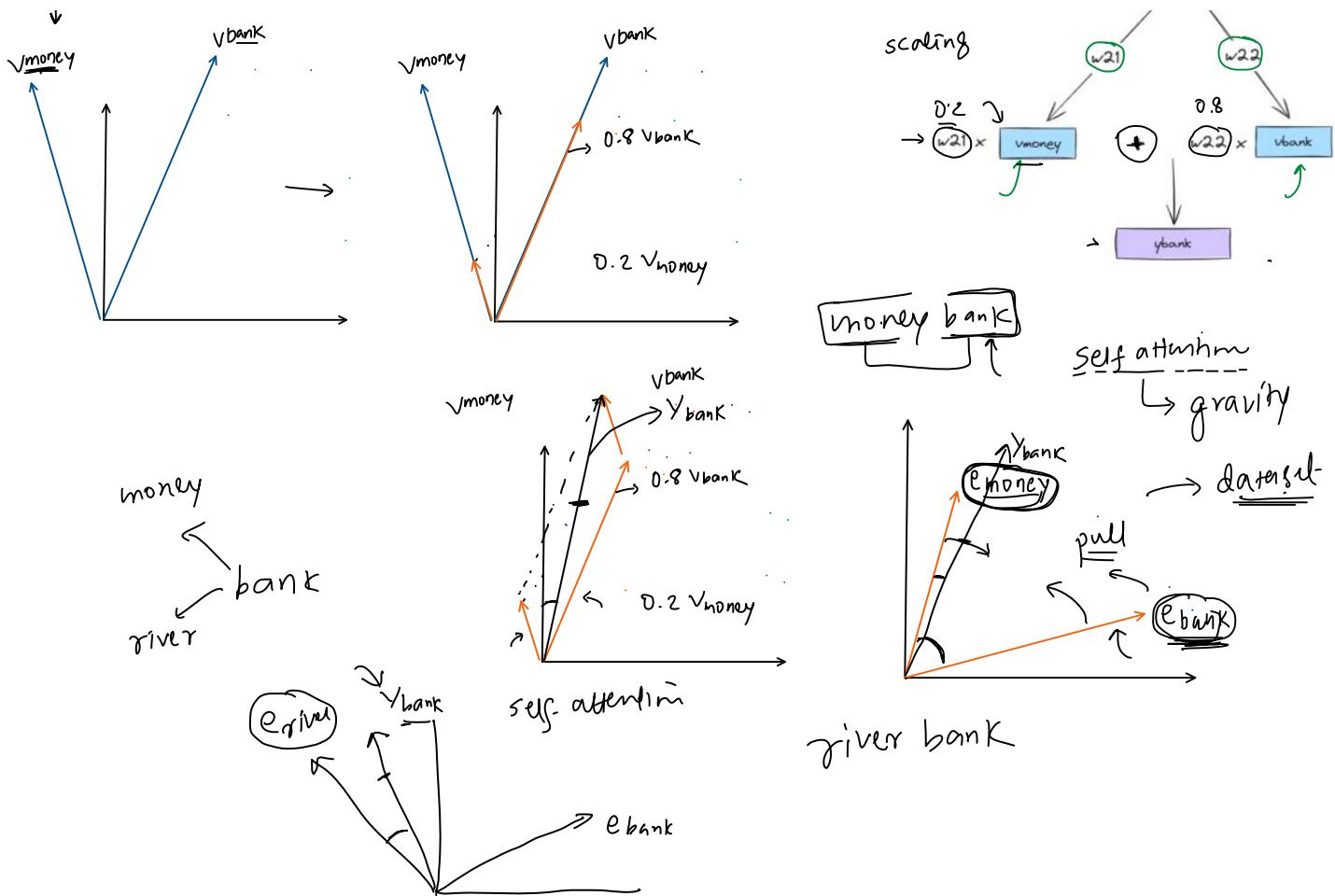
$$S_{21}' = \frac{10}{\sqrt{2}} = 7.09 \quad S_{22}' = 22.69$$

[Softmax]

$$w_{21} = 0.2$$

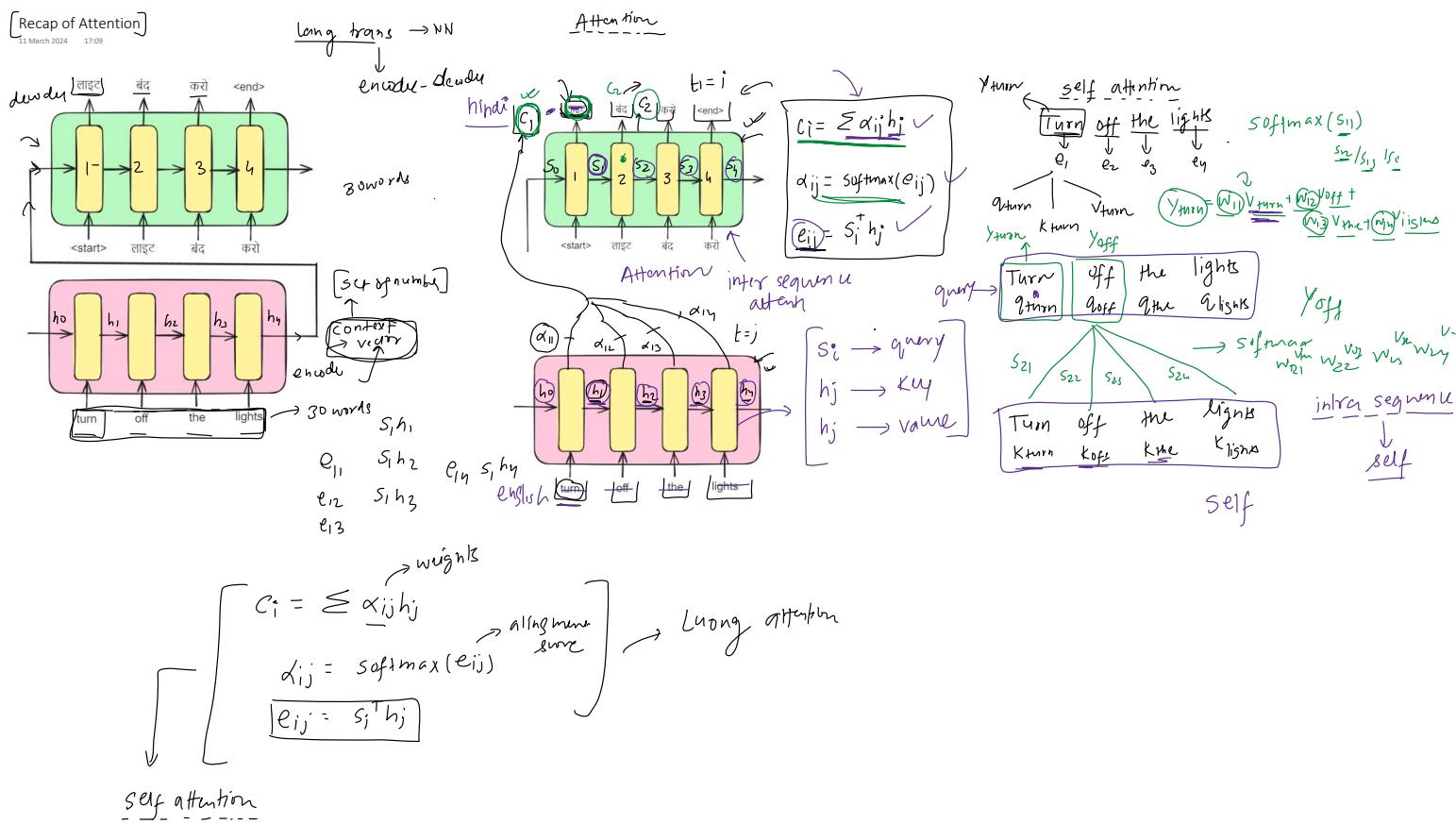
$$w_{22} = 0.8$$





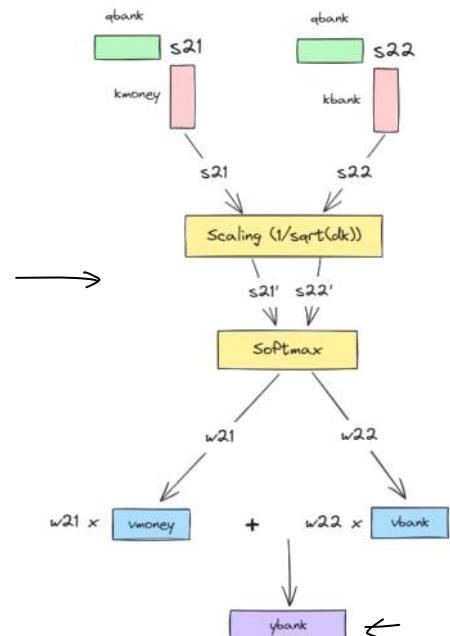
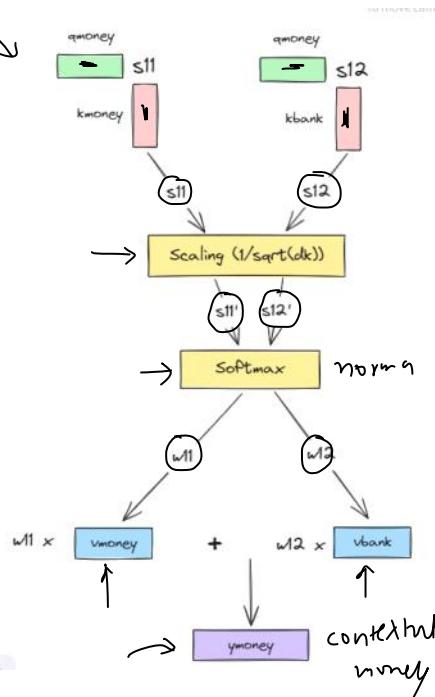
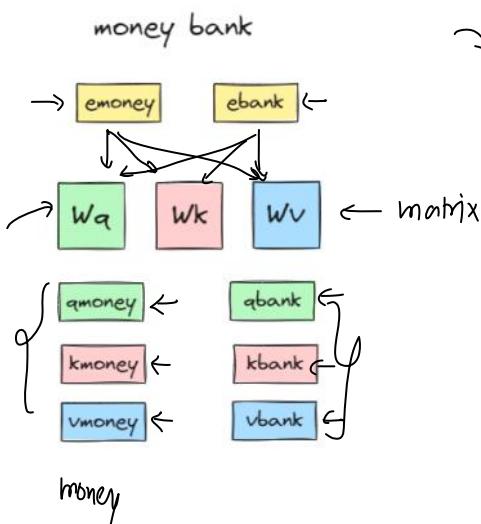
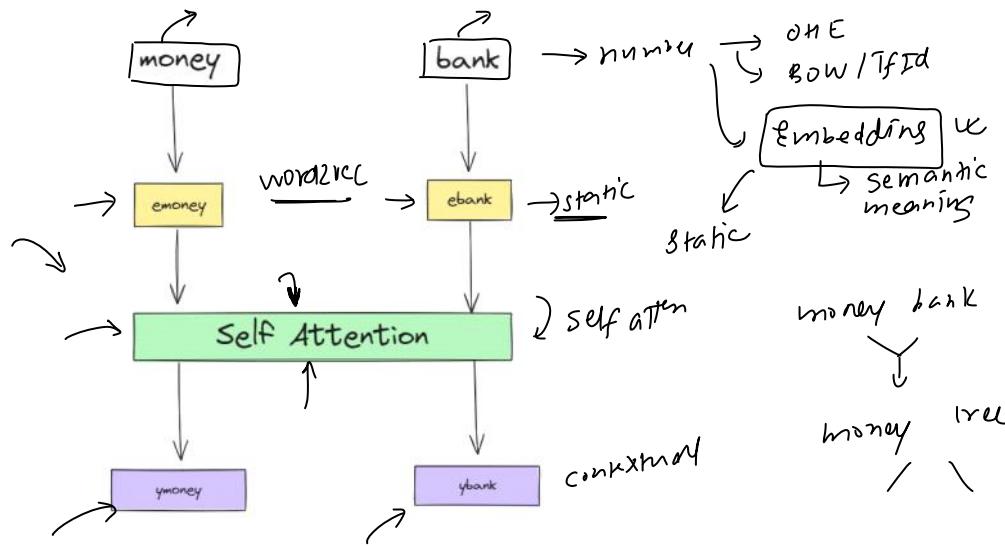
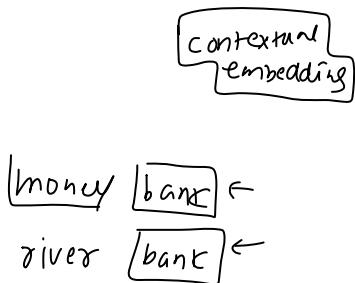
Recap of Attention

11 March 2024 17:09



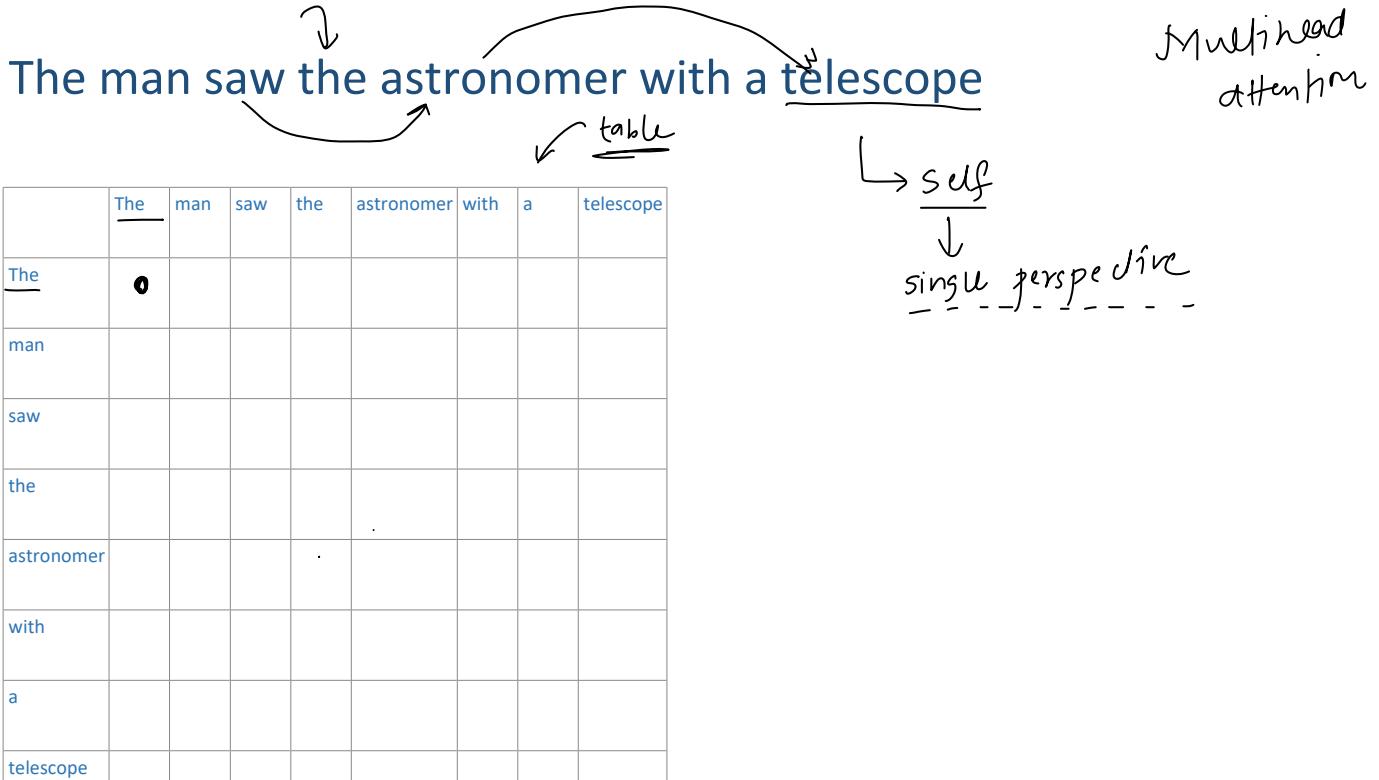
Recap of Self Attention

08 April 2024 14:46



Problem with Self Attention

08 April 2024 14:47



doc summarization

Self attention

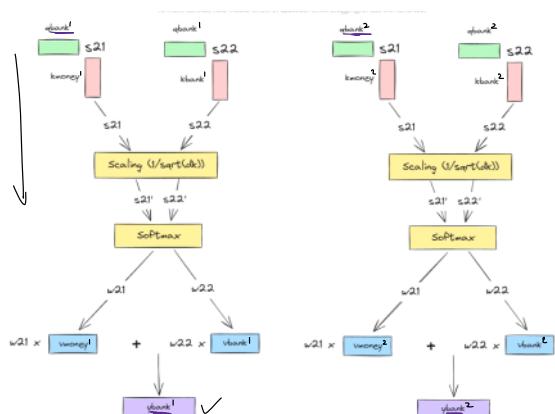
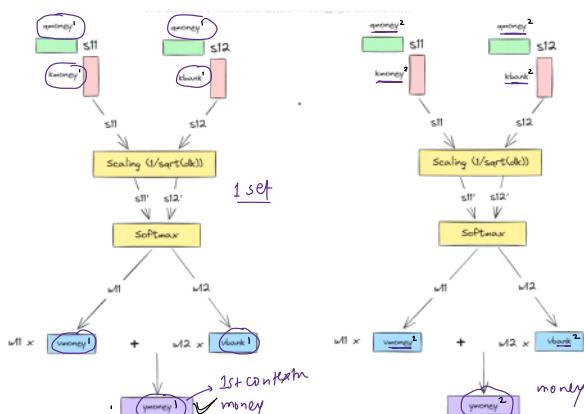
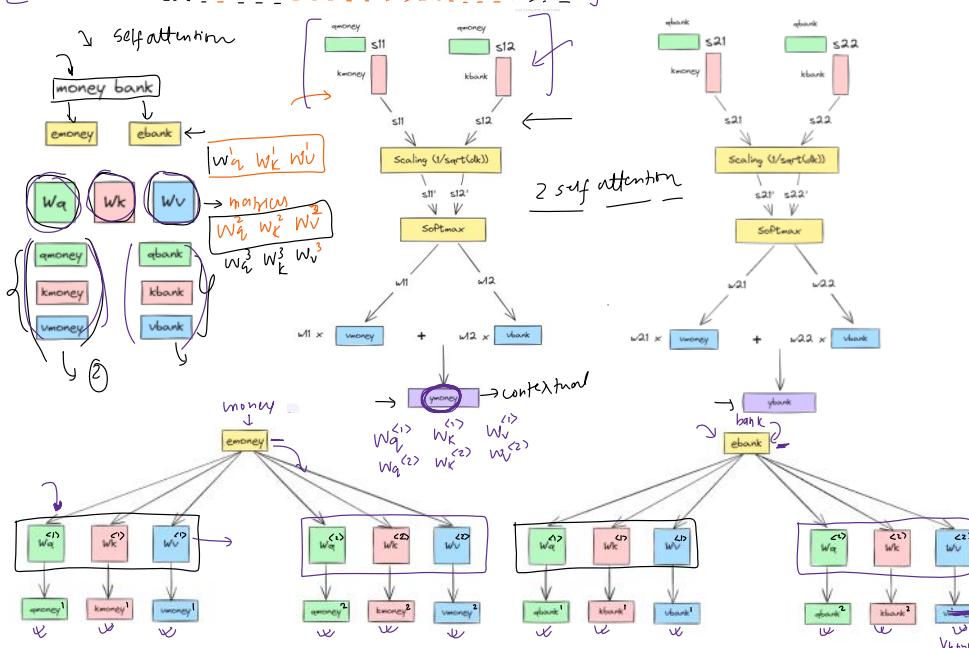
The future of AI in India presents a dynamic and promising landscape, marked by rapid advancements and a burgeoning ecosystem of innovation. With a robust talent pool of engineers and IT professionals, India is poised to become a significant player in the global AI arena. The government's proactive stance on AI, exemplified by initiatives like the National Strategy for Artificial Intelligence, aims to harness AI's potential across various sectors, including healthcare, education, agriculture, and urban infrastructure. Indian startups and tech giants are increasingly incorporating AI to solve complex societal challenges, improve efficiency, and enhance service delivery. Moreover, India's focus on ethical AI and data security aims to create a sustainable and responsible growth trajectory. As AI becomes more integrated into daily life and industry, India's unique blend of technological prowess, entrepreneurial spirit, and societal needs will likely shape a distinctive path in the AI domain, fostering innovation that is not only technologically advanced but also socially inclusive and impactful.

India is poised to become a key player in the AI domain, leveraging its skilled workforce and government initiatives to apply AI across various sectors like healthcare and education. With a focus on innovation, ethical AI, and data security, India aims to integrate AI to address societal challenges, enhance efficiency, and promote inclusive growth. This approach positions India to uniquely contribute to global AI advancements while ensuring sustainable and responsible development within its own borders.

India's AI future holds promise as it harnesses a burgeoning talent pool and government initiatives to pioneer AI-driven innovation. With a focus on sectors like healthcare and education, India aims to leverage AI for societal development. The emphasis on ethical AI and data security underscores India's commitment to responsible technological advancement. This approach positions India not only as a global AI hub but also as a trailblazer in addressing societal challenges through cutting-edge technology.

The man saw the astronomer with a telescope

self attention → ②



mult

