

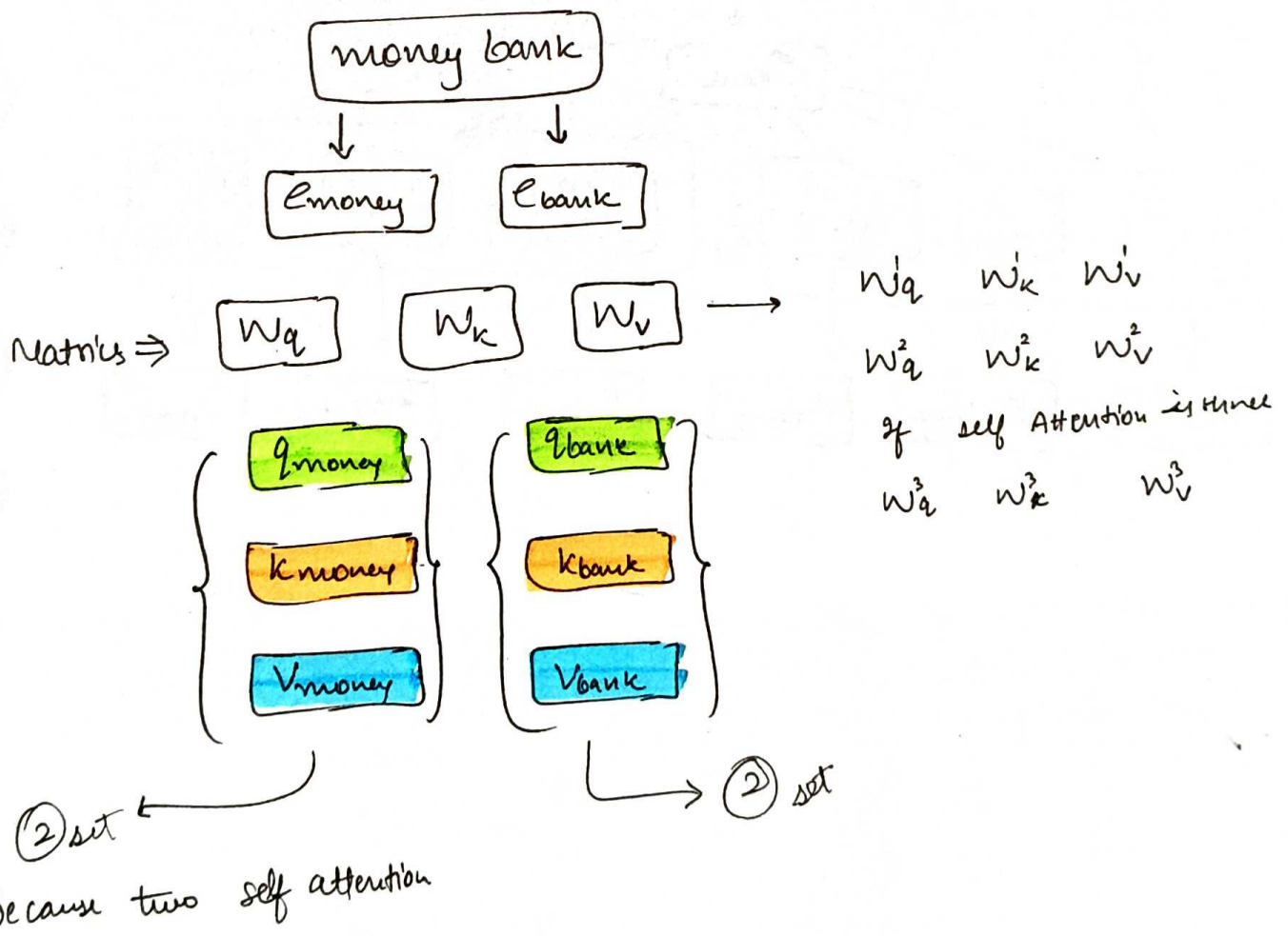
In NLP, there are multiple scenarios to capture multiple perspectives.

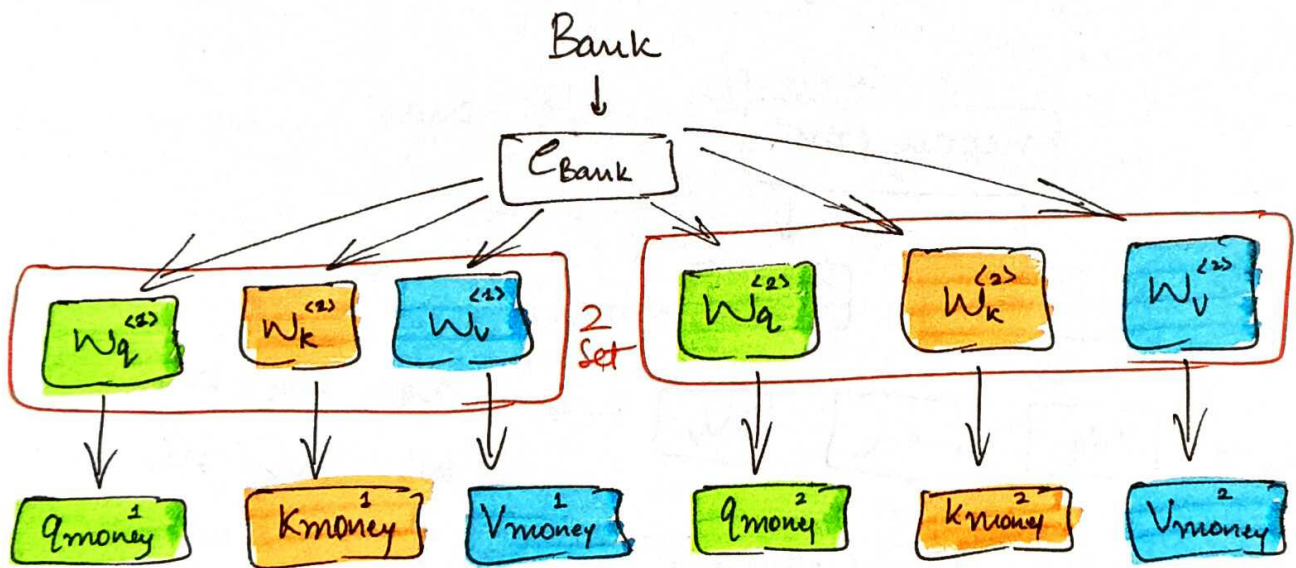
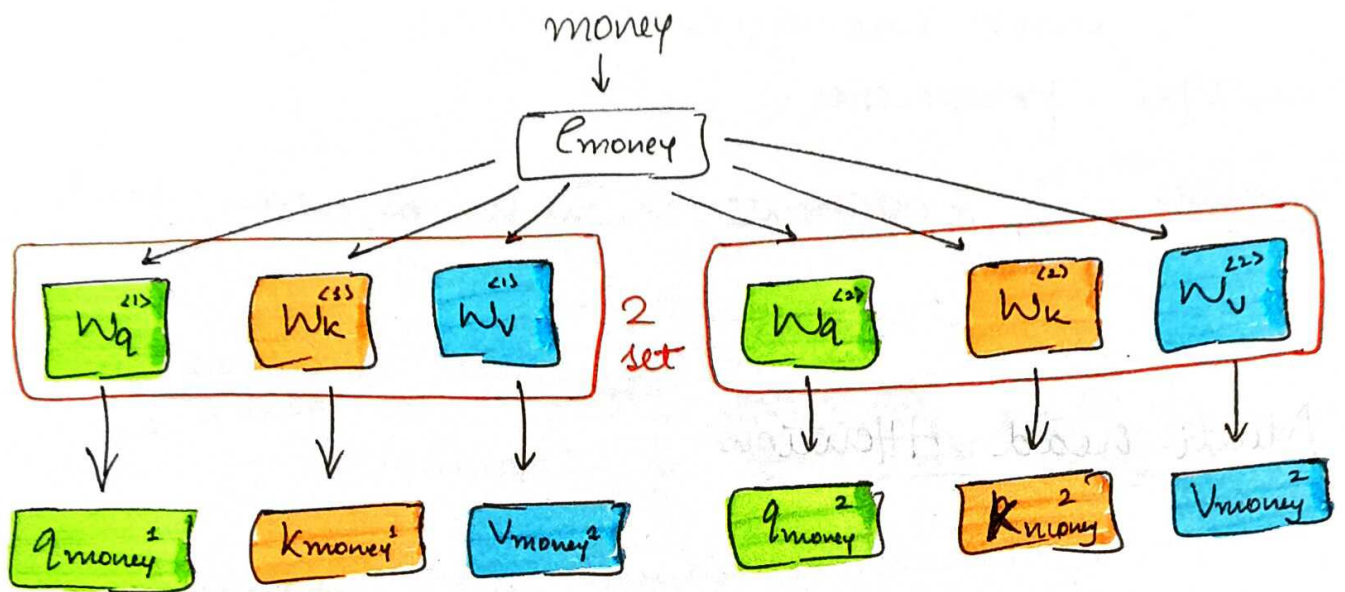
example → Document Summarization tool.

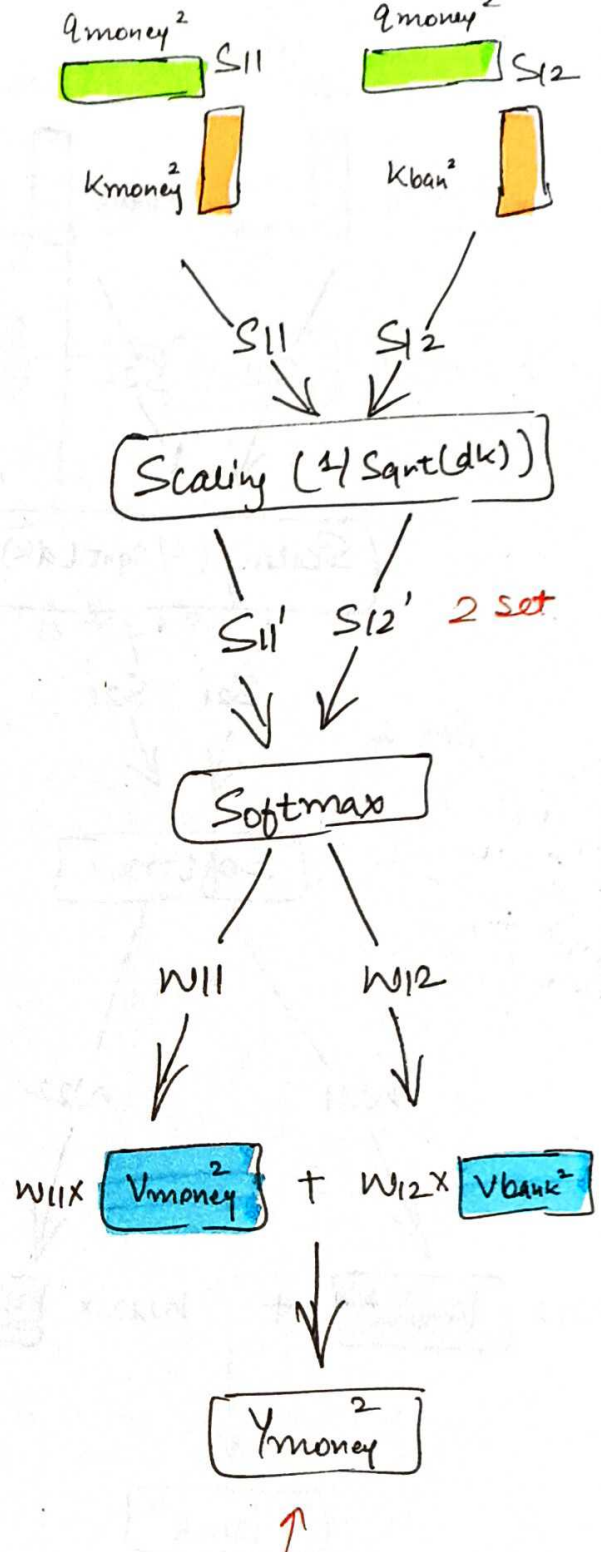
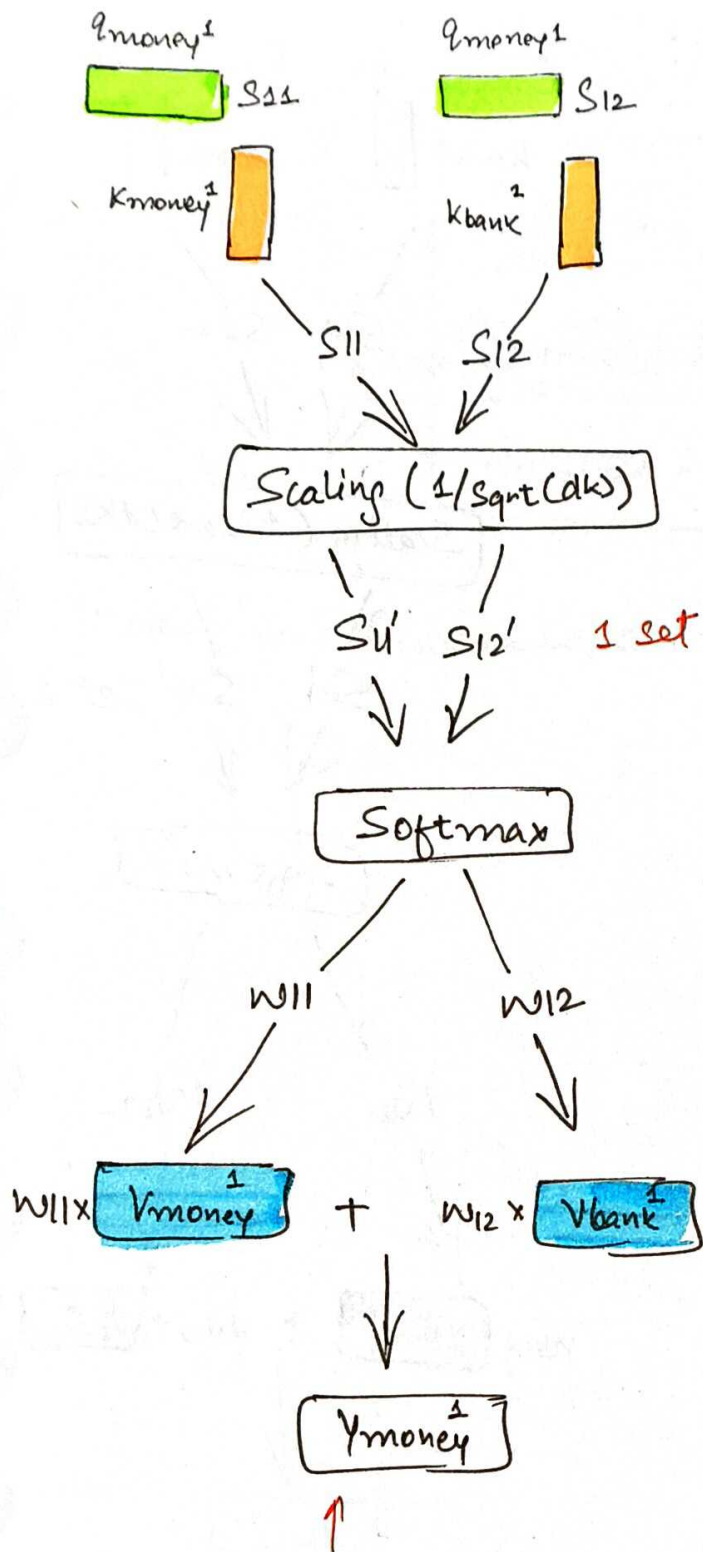
Multi-head Attention

The man saw the astronomer with a telescope
↳ 2 meaning

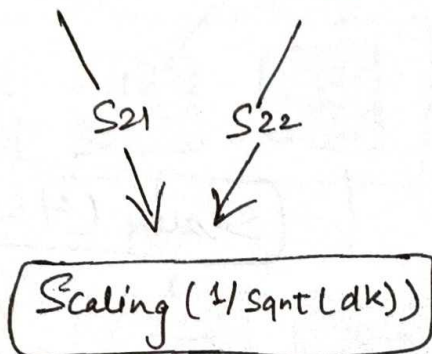
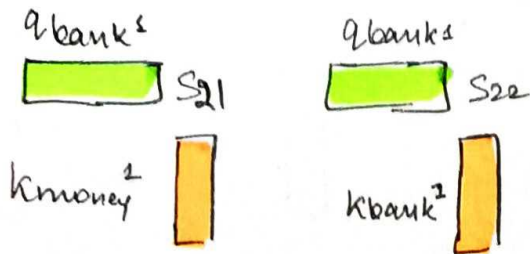
So, we use two self Attention.



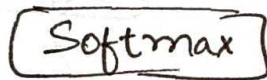
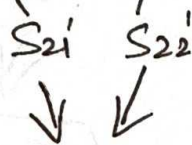




All these steps for money word

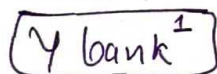
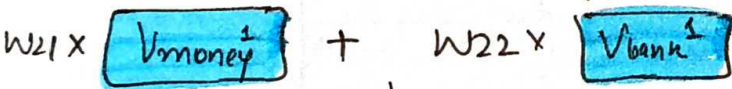


Set 1

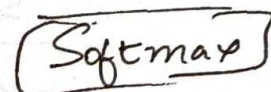
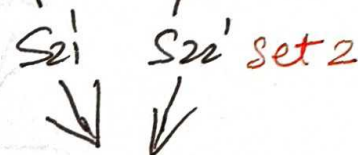
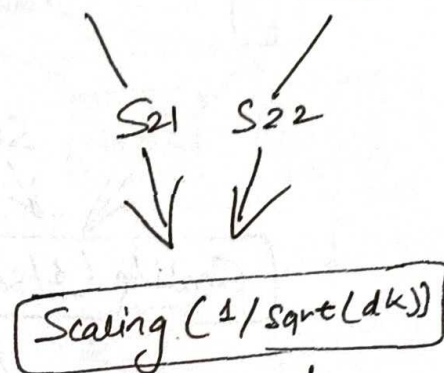
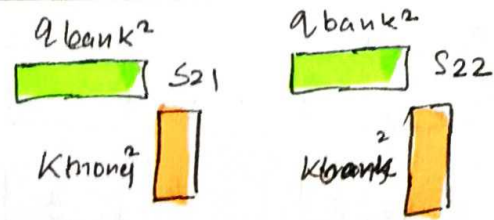


w_{21}

w_{22}

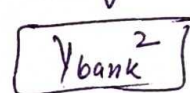
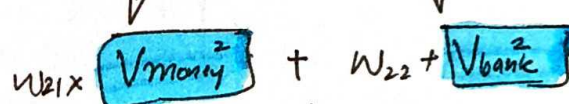


1st contextual representation
for bank word



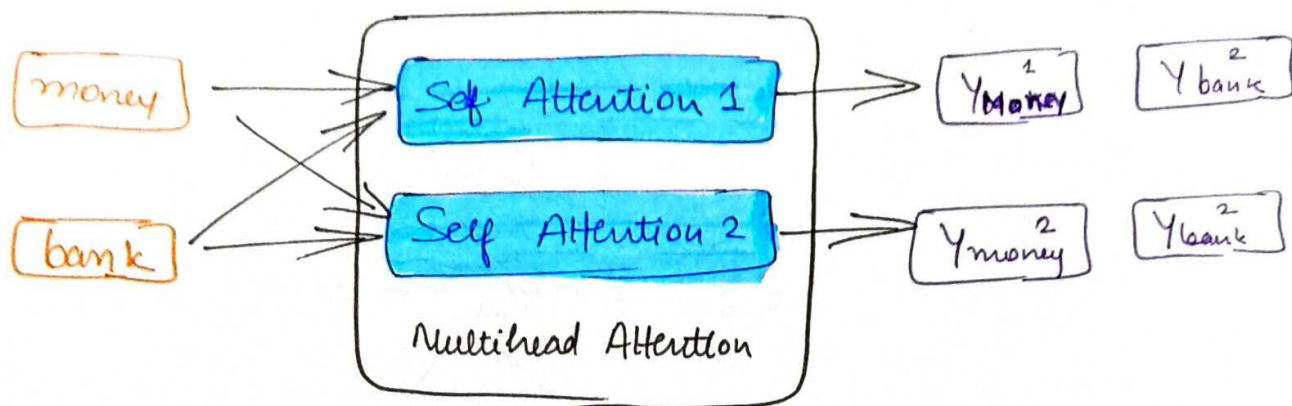
w_{21}

w_{22}



2nd contextual
representation for
bank word.

bank word



Transformer Research paper \rightarrow 8 heads

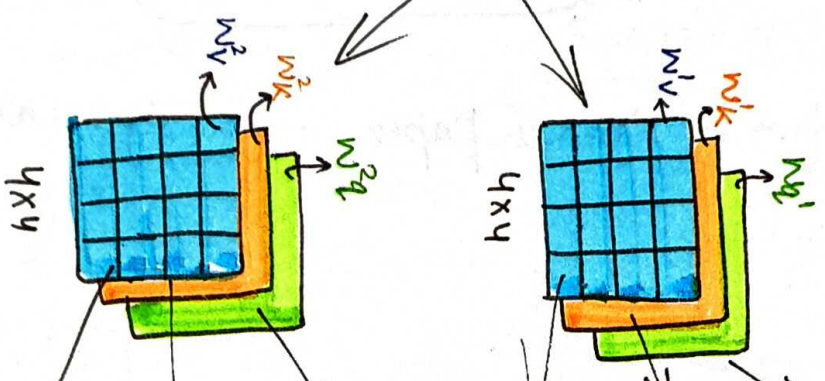
money
bank

(Embedding)

2x4 dim

4 dim

4 dim



4 dim

Q¹
money₁
bank₁
(2x4)

K¹
money₁
bank₁
(2x4)

V¹
money₁
bank₁
(2x4)

Self
Attention

Z¹
money₁
bank₁
(2x4)

concatenate

(2x8)

Z²
money₂
bank₂
(2x4)

(Z¹)

X

linear transform

W⁰
money₀
bank₀
(8x4)

weights
(learn during
training
time)

Z²
money₂
bank₂
(2x4)

V²
money₂
bank₂
(2x4)

K²
money₂
bank₂
(2x4)

Q²

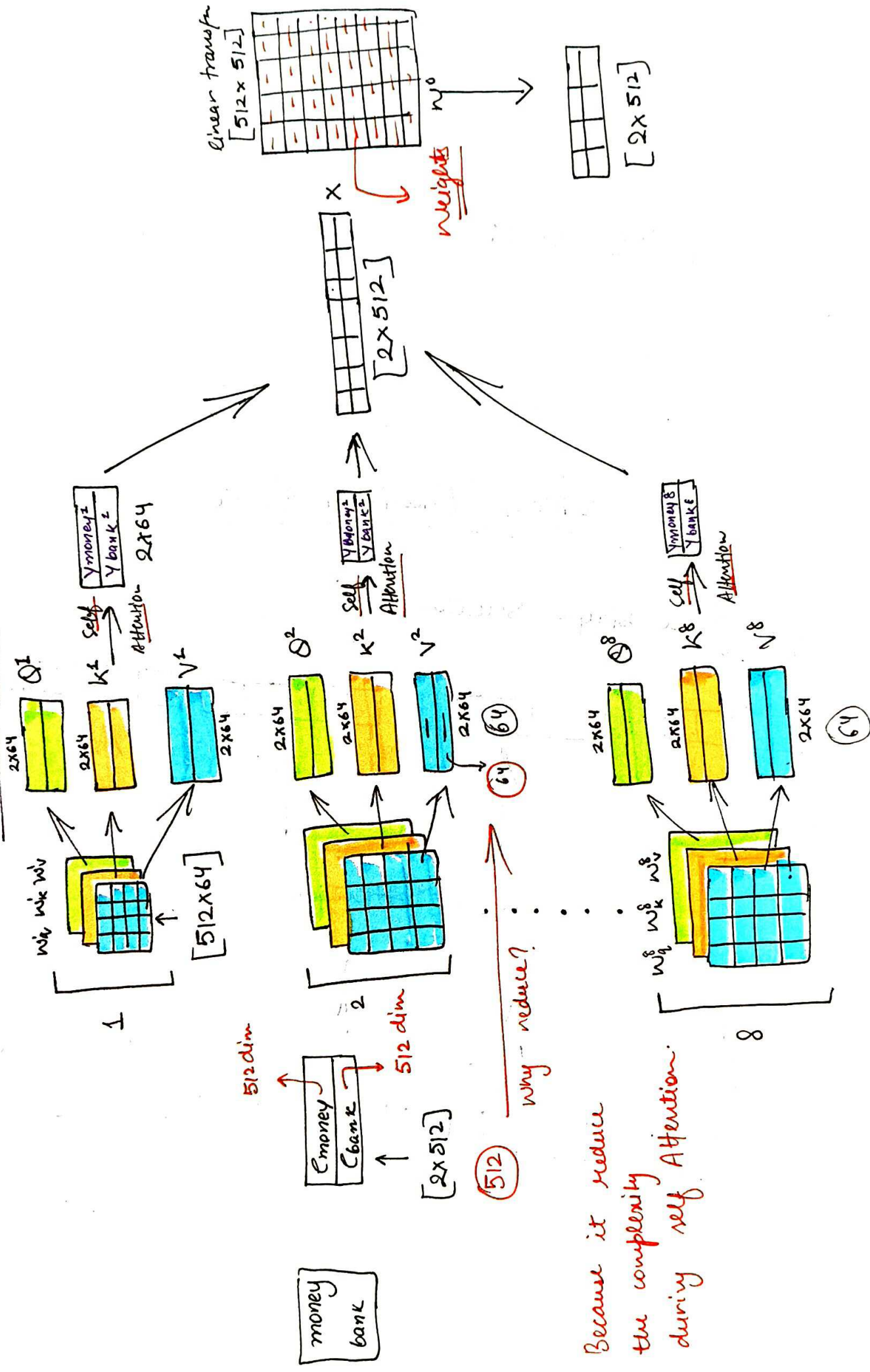
(to balance
Z¹ and Z²)

Main
goal

Z²
money₂
bank₂
(2x4)

(Z¹)
money₁
bank₁
(2x4)

Attention in all you need



Because it reduce the complexity during self Attention.

(512) why reduce?

(64)

(64)

Benefit of Self Attention

Parallel Training

Drawback of Self Attention

word order

eg:- "Nitish killed lion" is same "lion killed Nitish"

Positional Encoding in Transformer

Proposing a Simple Solution

