

Attention based Encoder Decoder Calculations

CMR@BRACU

Encoder

I go to school

Start with embeddings of I, go, to ,school

Add positional encoding to determine the order of the words

Form Query, Key and Value matrix for encoder

Transform each embedding into K, Q and V embedding by multiplying with Query, Key and value matrix

Calculate self attention of I, between I and go, I and to, I and school and update the V embedding of I. Now calculate the self attention of go, between go and I, go and to, go and school and update the value embedding of go. Do the same for to and school

Add the V embedding of each token to their original embedding

Decoder

<EOS>

Form a new set of Query, Key and Value matrix fpr decoder and transform <EOS> embedding into Q embedding, K embedding and V embedding

Calculate self attention of EOS , update the V embedding of EOS and add the V embedding to the original embedding of EOS

Form a new set of Query, Key and Value matrix for cross attention between encoder and Decoder and transform the encoding to I, go, to, school and EOS

Calculate the attention between EOS and I, EOS and go, EOS and to and EOS and school. Update the V embedding of EOS. Add this V embedding of EOS to its latest embedding

Feed latest embedding of EOS to a dense network to produce the embedding of আমি

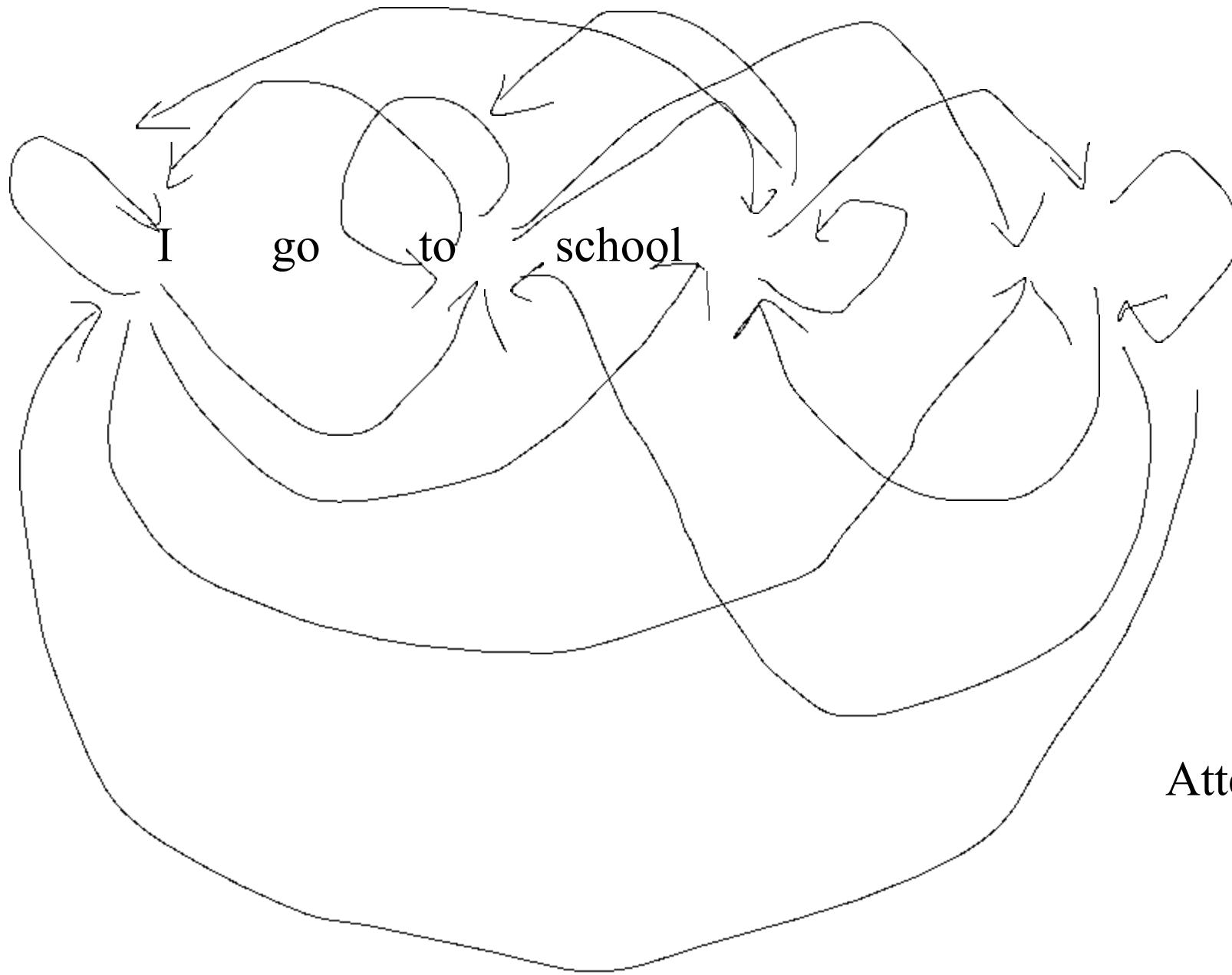
Embedding of আমি is converted to Q, K and V embedding using matrices for decoder

Calculate self attention of আমি and between আমি and <EOS>. Update the V embedding of আমি and add this V embedding with its original embedding

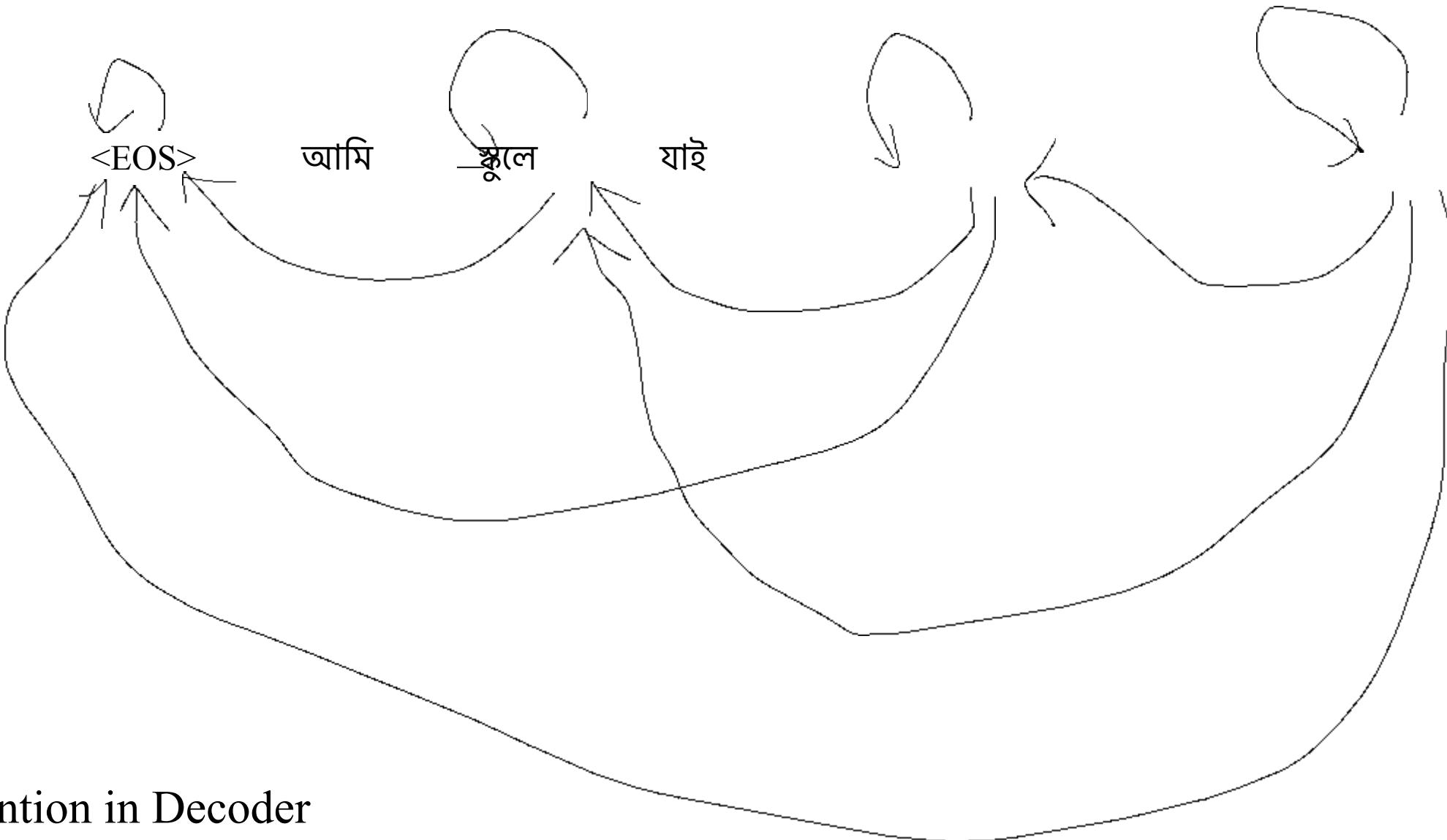
Now convert the latest embedding of আমি to Q, K and V embedding using the matrices of cross attention

Calculate the attention between আমি and I, আমি and go, আমি and to and আমি and School. Update the V encoding of আমি and add this with the latest encoding of আমি

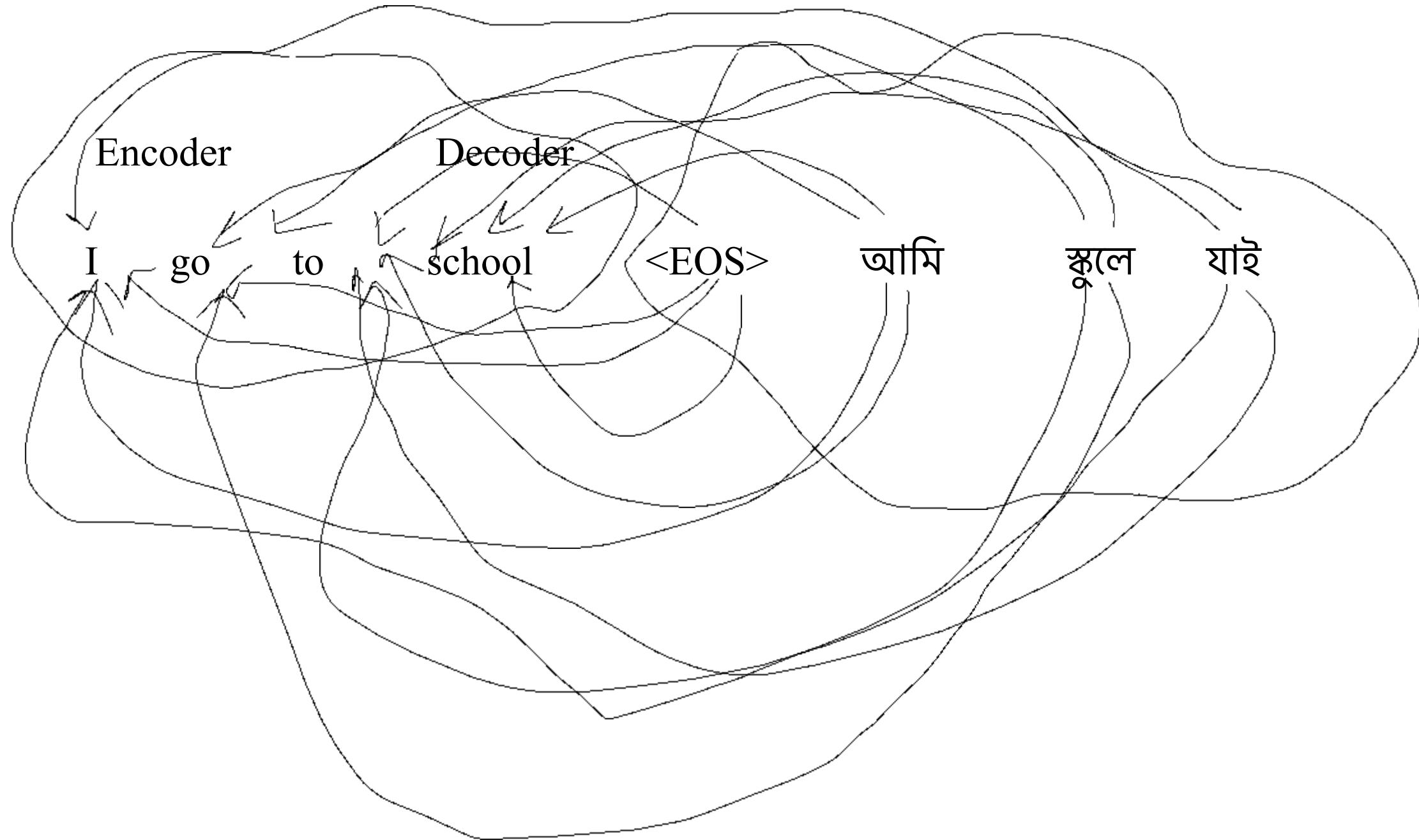
Feed the latest encoding of আমি to a dense neural network and produce the second word



Attention in Encoder



Attention in Decoder



Decoder encoder cross attention

Embedding of କୁଳେ is converted to Q, K and V embedding using matrices for decoder

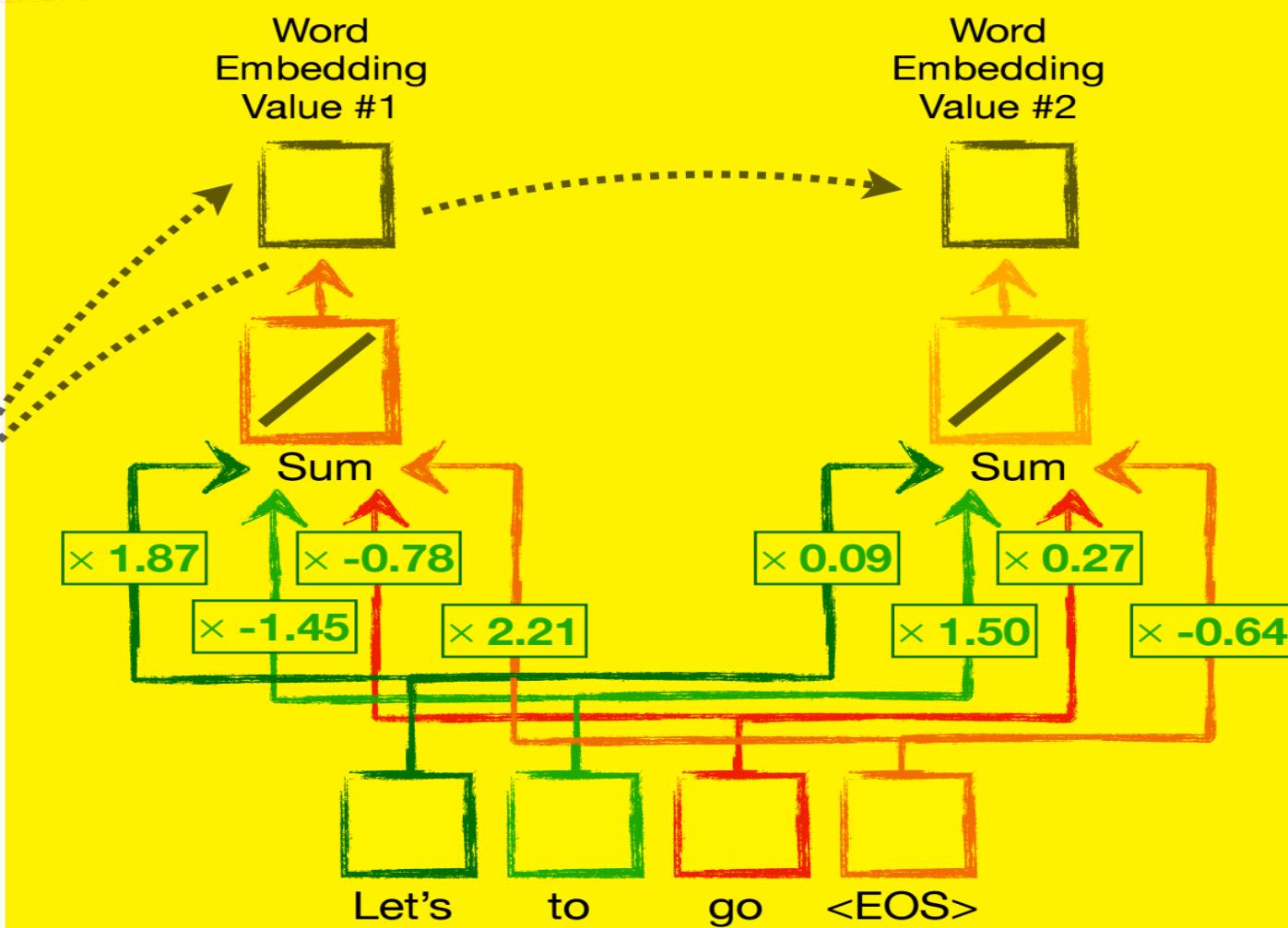
Calculate self attention of କୁଳେ and between କୁଳେ and ଆମି and କୁଳେ and <EOS>.

Update the V embedding of କୁଳେ and add this V embedding with its original embedding

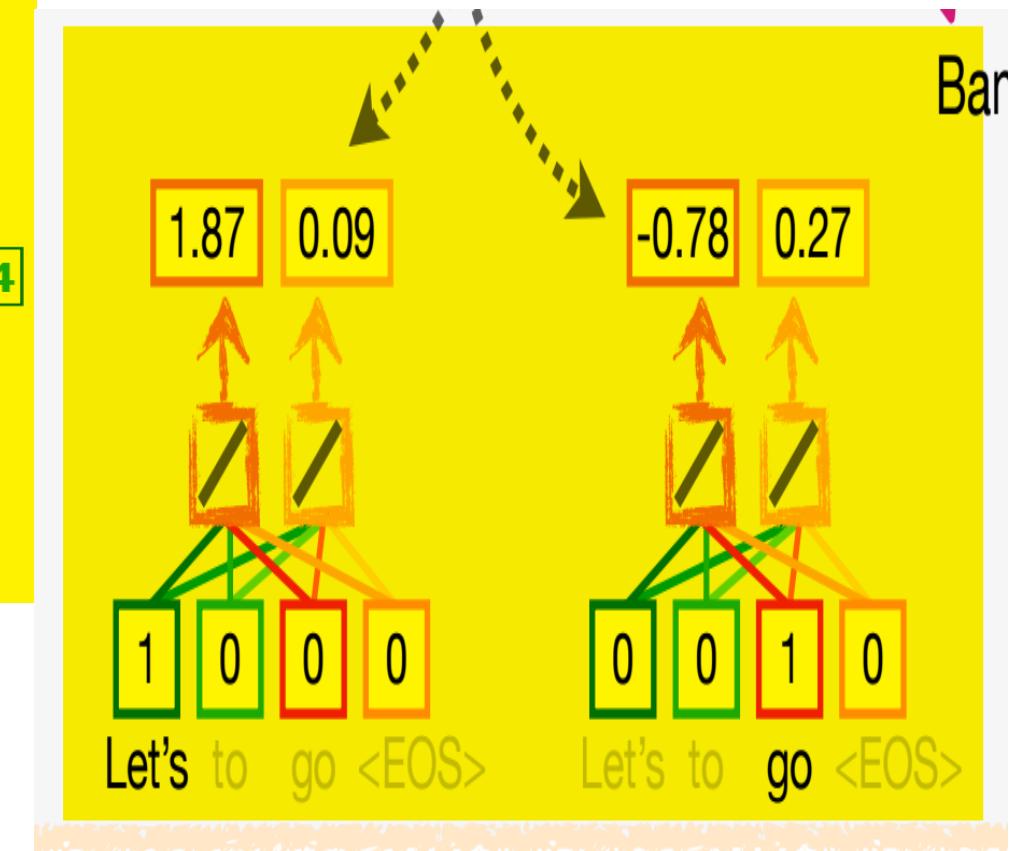
Now convert the latest embedding of କୁଳେ to Q, K and V embedding using the matrices of cross attention

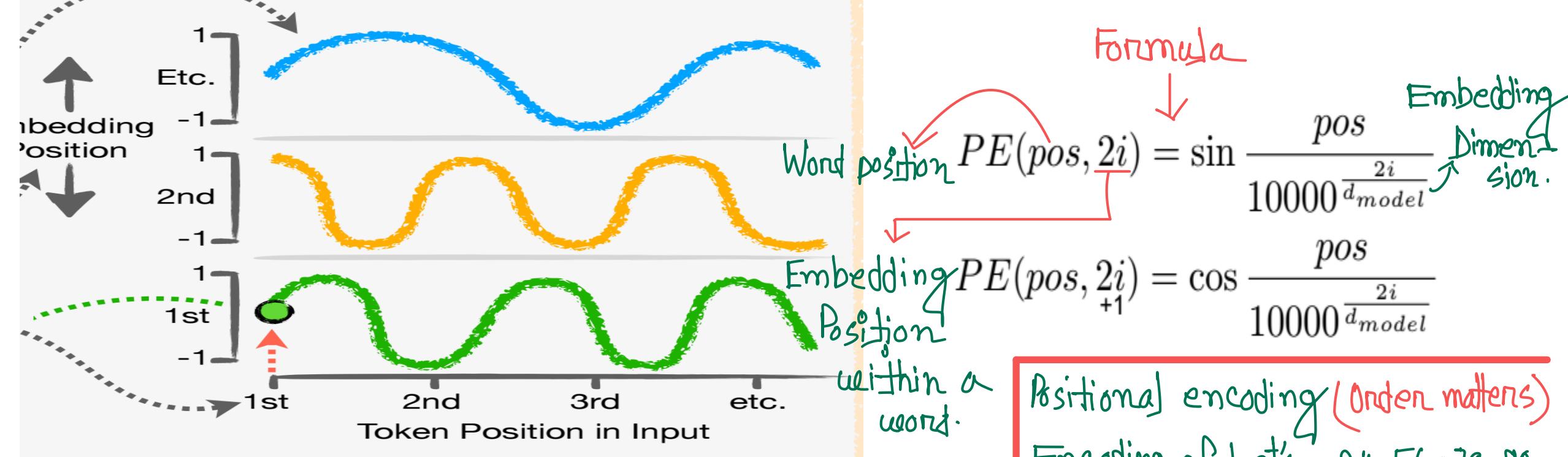
Calculate the attention between କୁଳେ and I, କୁଳେ and go, କୁଳେ and to and କୁଳେ and School. Update the V encoding of କୁଳେ and add this with the latest encoding of କୁଳେ

Feed the latest encoding of କୁଳେ to a dense neural network and produce the third word Embedding for ସାଇ



⇒ Word embedding
⇒ 2 dimensional word embedding





$$PE(0,0) = \sin(0) = 0, PE(0,1) = \cos(0), PE(0,2) = \sin(0), PE(0,3) = 1$$

$$PE(1,0) = \sin(1), PE(1,1) = \cos(1), PE(1,2) = \sin(1/10000^{(2/4)}), PE(1,3) = \cos(1/10000^{(2/4)})$$

$\cancel{pos \rightarrow 0}$

word1	word2	word3
.34 .56 .78 .90 $i=0$.56 .64 .88 .76 $i=1$.34 .12 .39 .98
.34+0, .56+1, .78+0, .90+1	.56+0.02, .64+1, .88+0, .76+1	

Formula

$$PE(pos, 2i) = \sin \frac{pos}{10000^{\frac{2i}{d_{model}}}}$$

$$PE(pos, 2i+1) = \cos \frac{pos}{10000^{\frac{2i+1}{d_{model}}}}$$

Positional encoding (Order matters)

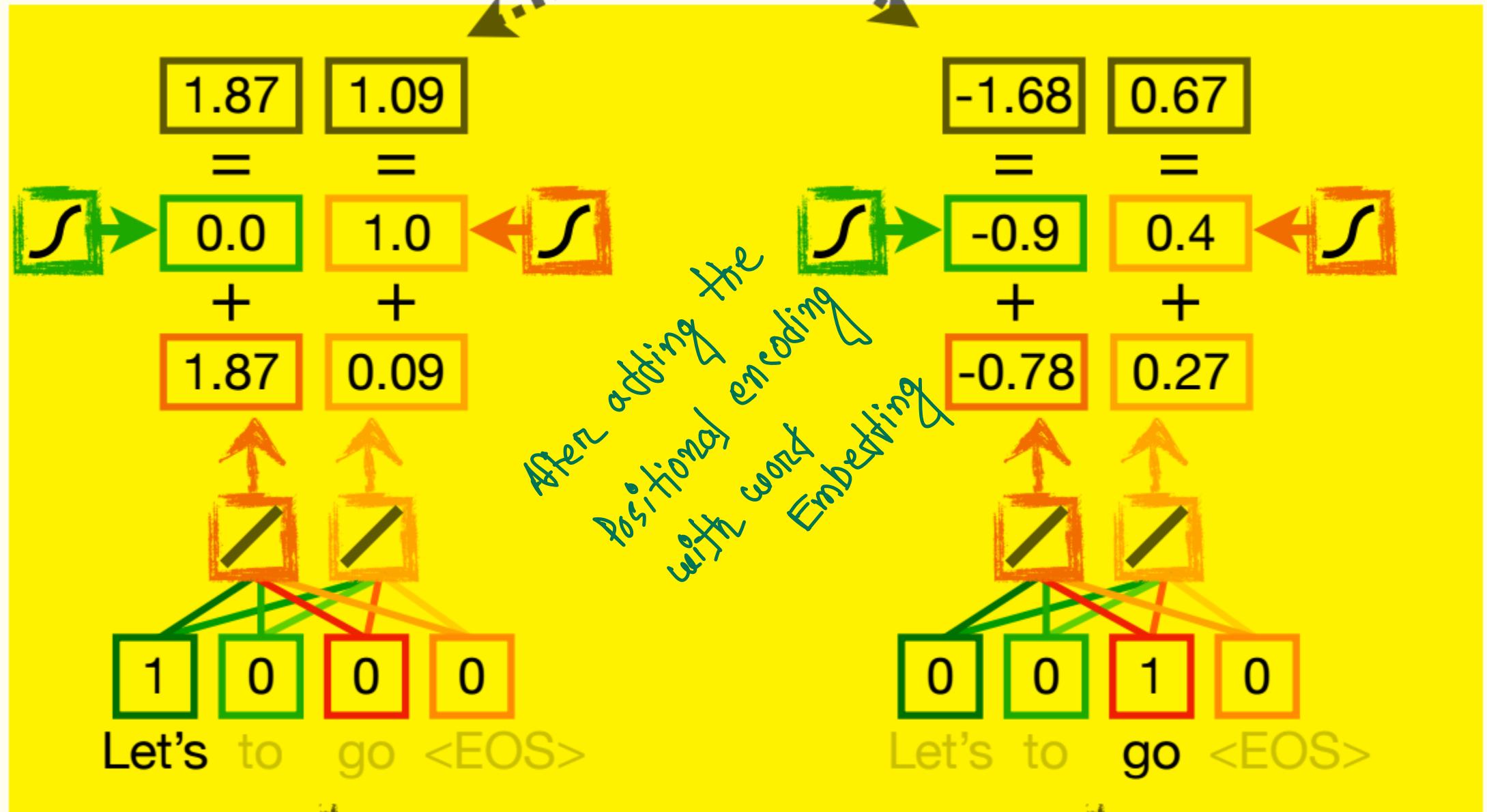
Encoding of Let's = .34, .56, .78, .90

$$\begin{aligned} PE &\rightarrow 2^{x0} \\ PE &\rightarrow (0, 0) \end{aligned}$$

$$= \sin \frac{0}{10000^{\frac{2x0}{4}}} = 0$$

$$= .34 + 0$$

PE(3, 3)?

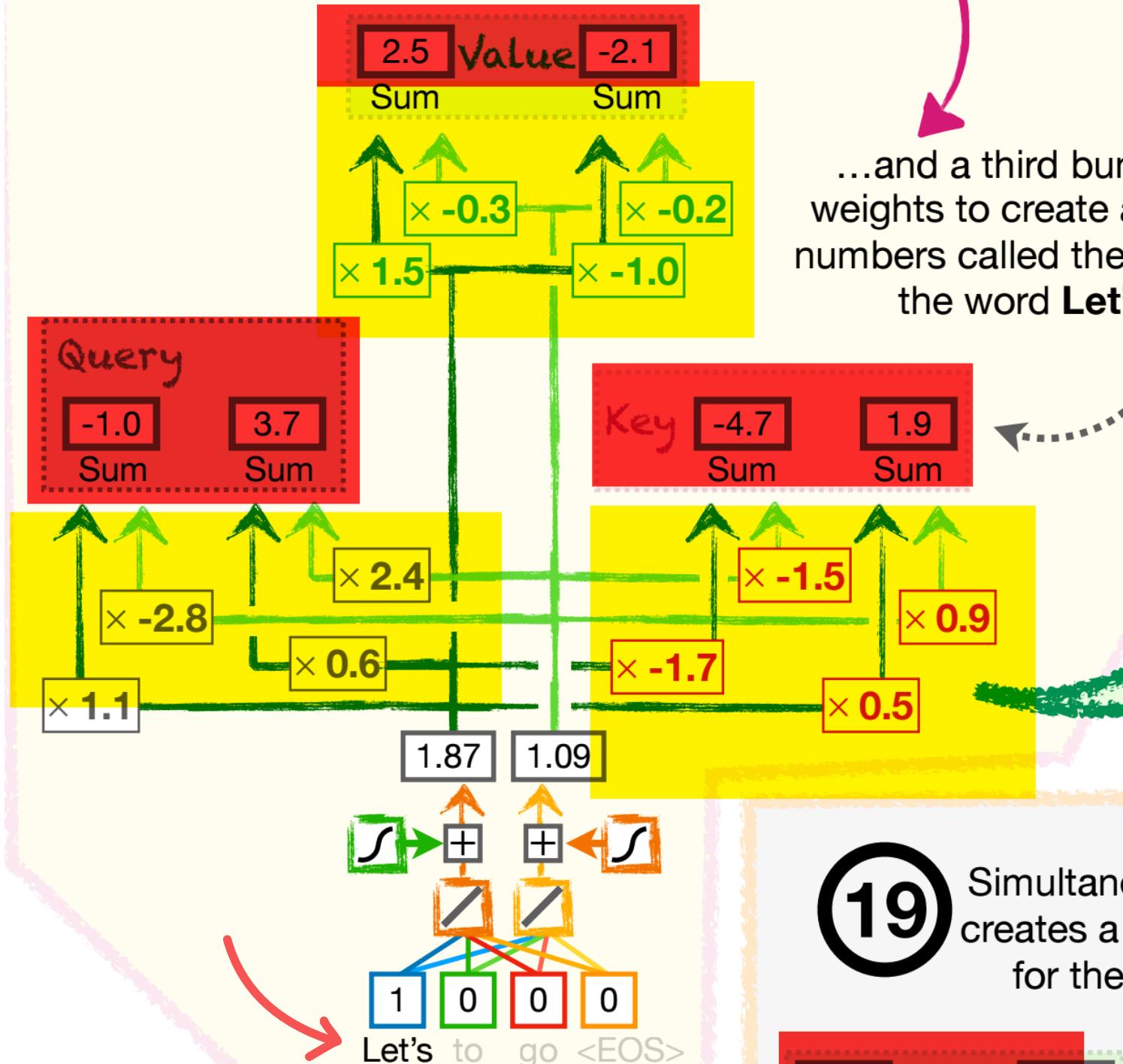




...and a third burst of weights to create a third set of numbers called the word **'Let'**

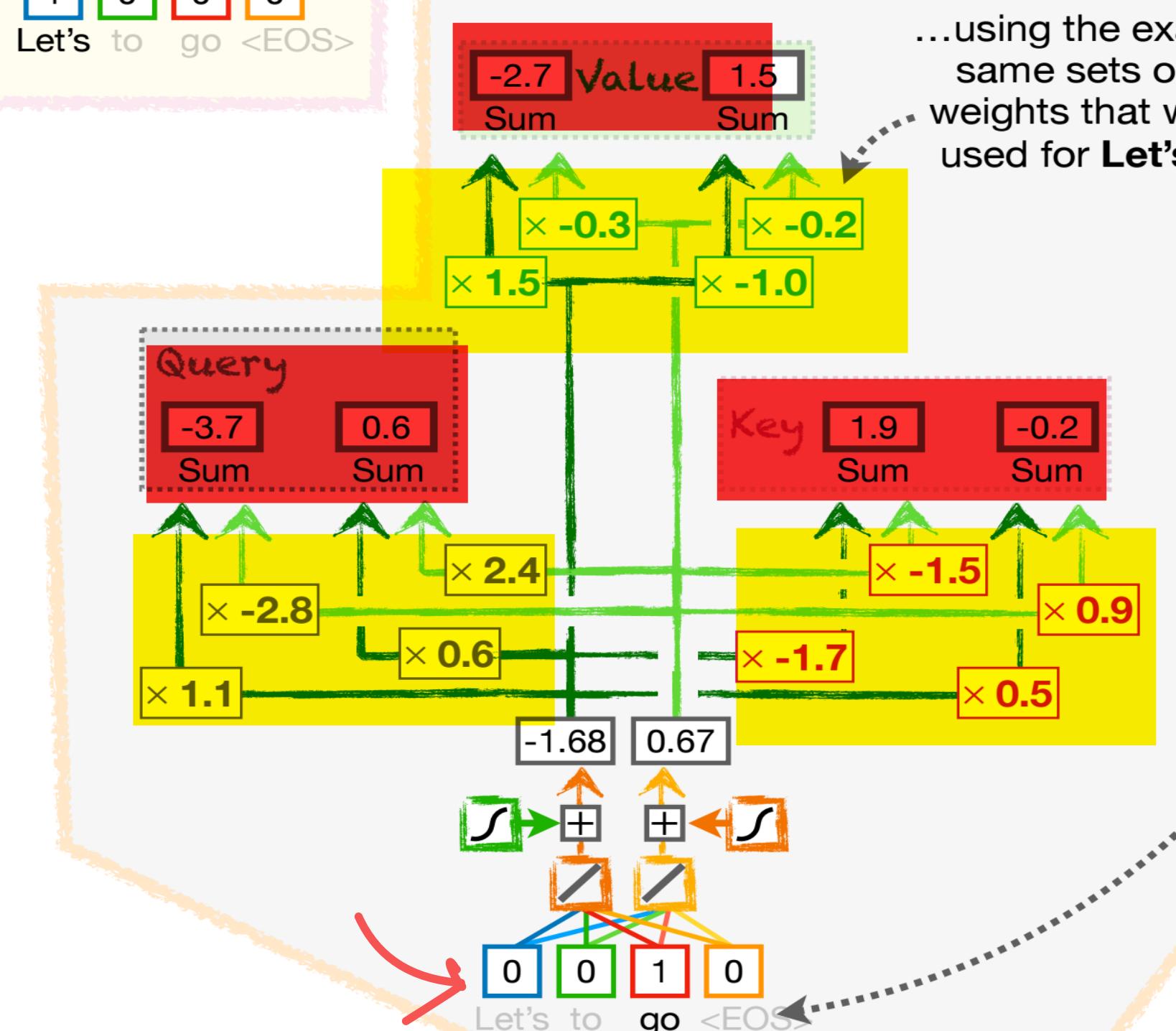
19

Simultaneously creates a vector for the word

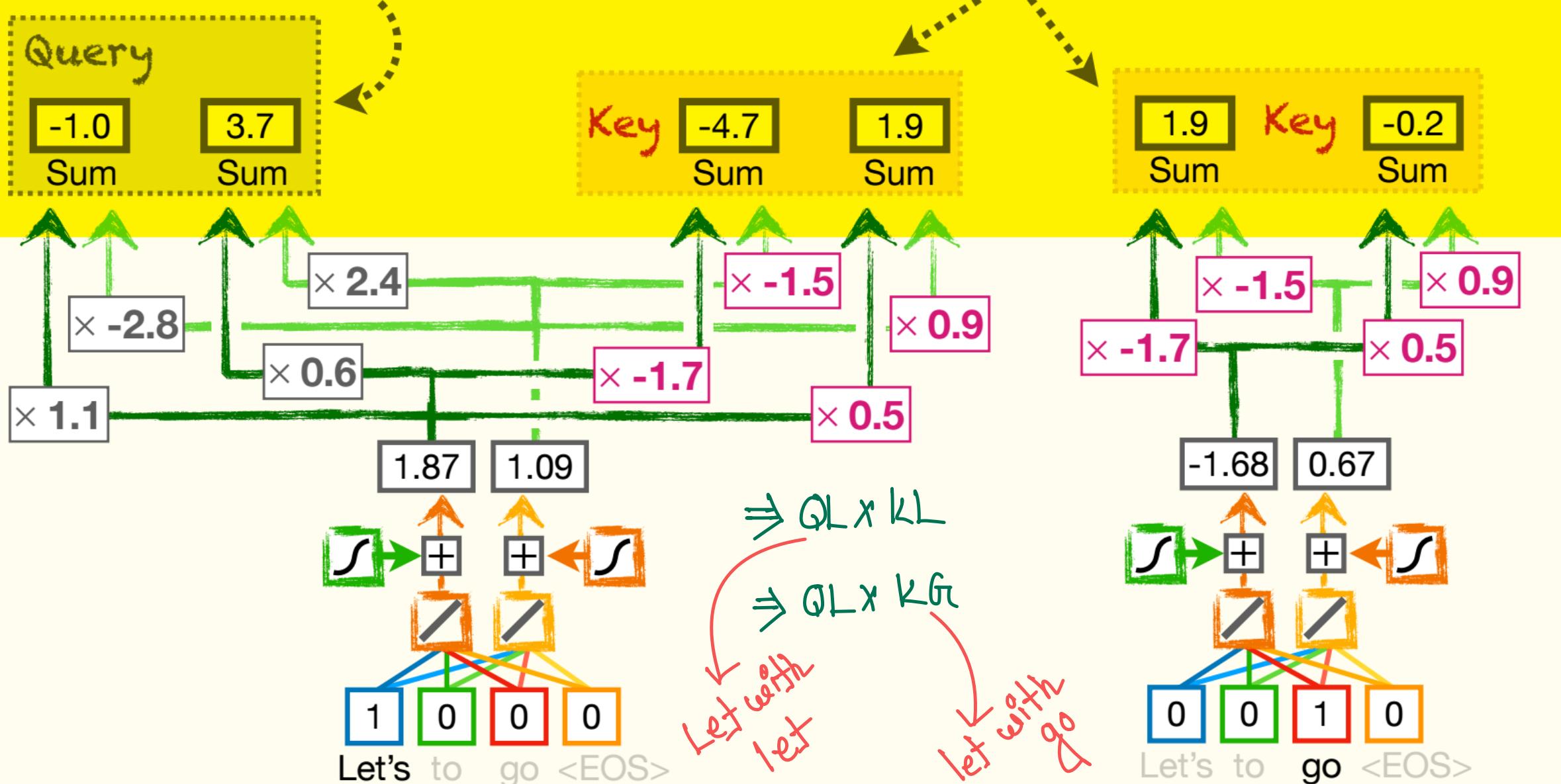


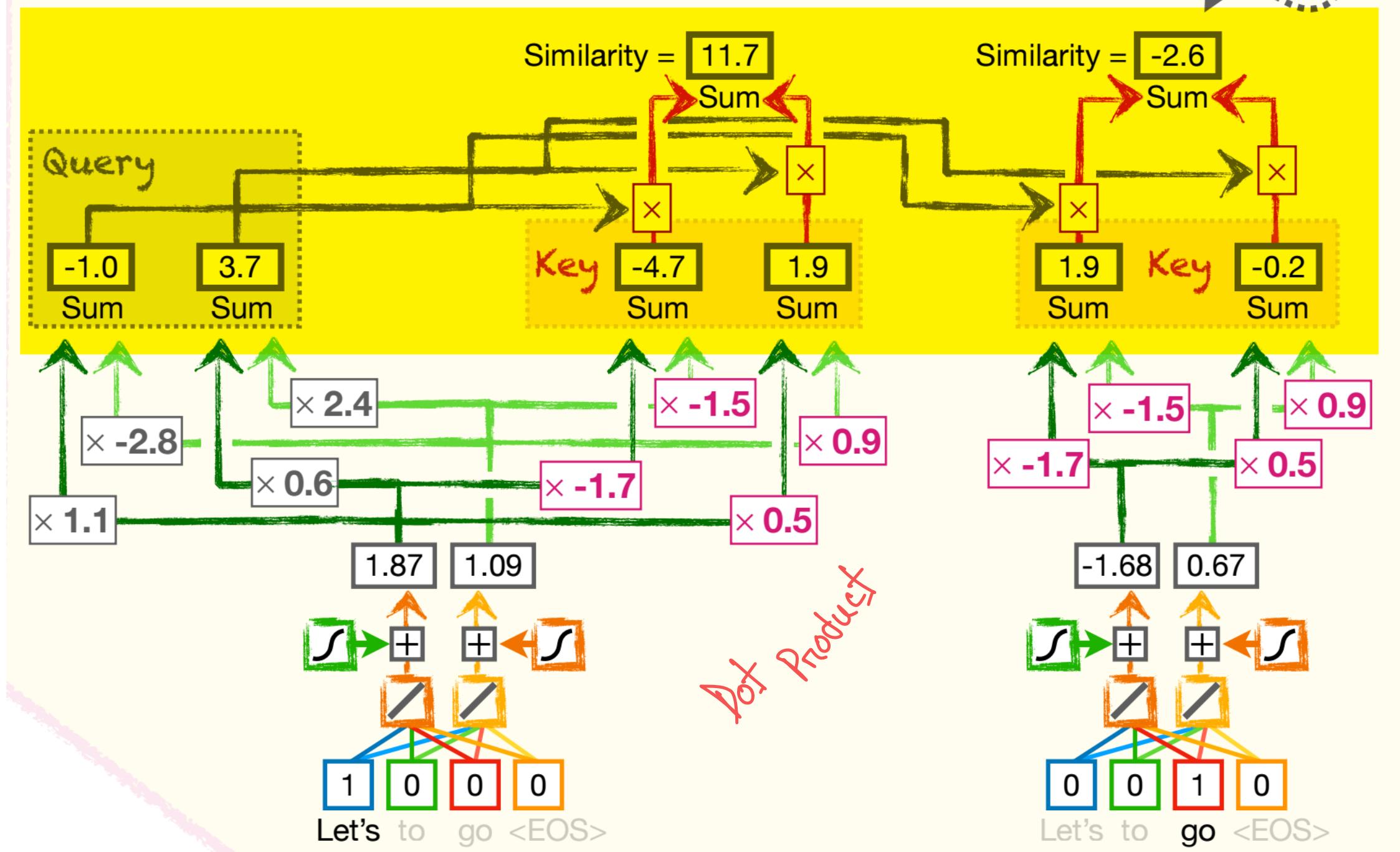
Let's to go <EOS>

Let's to go <EOS>

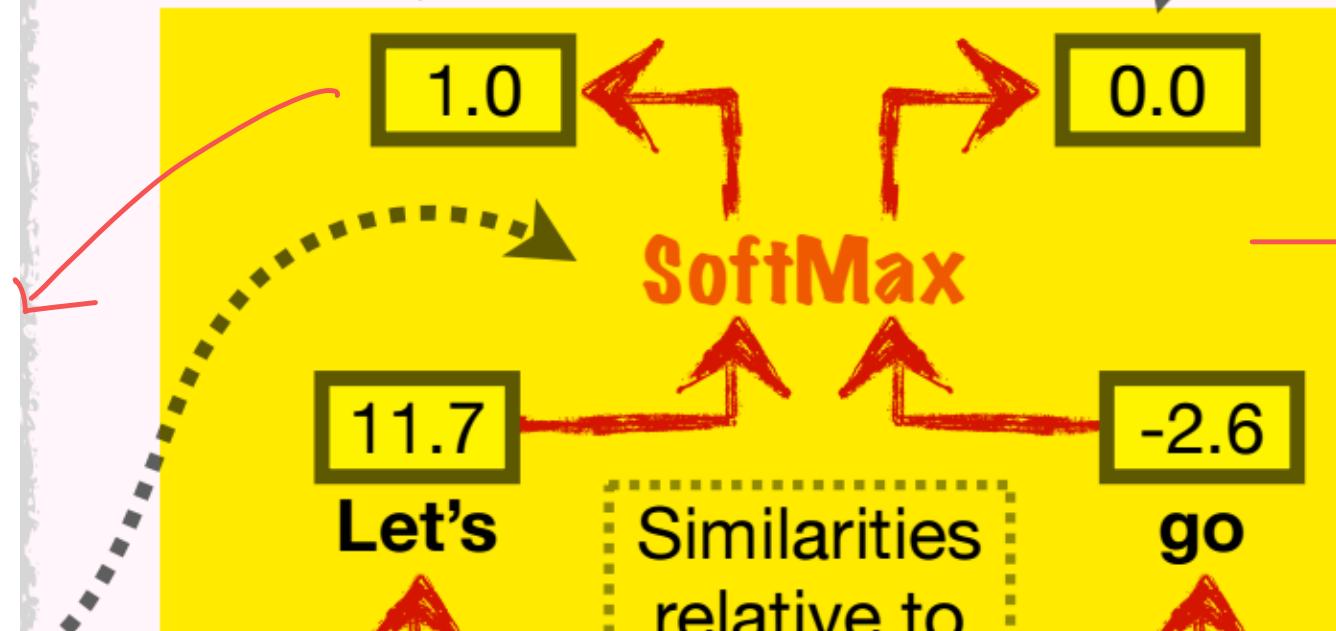


the query for Let's...



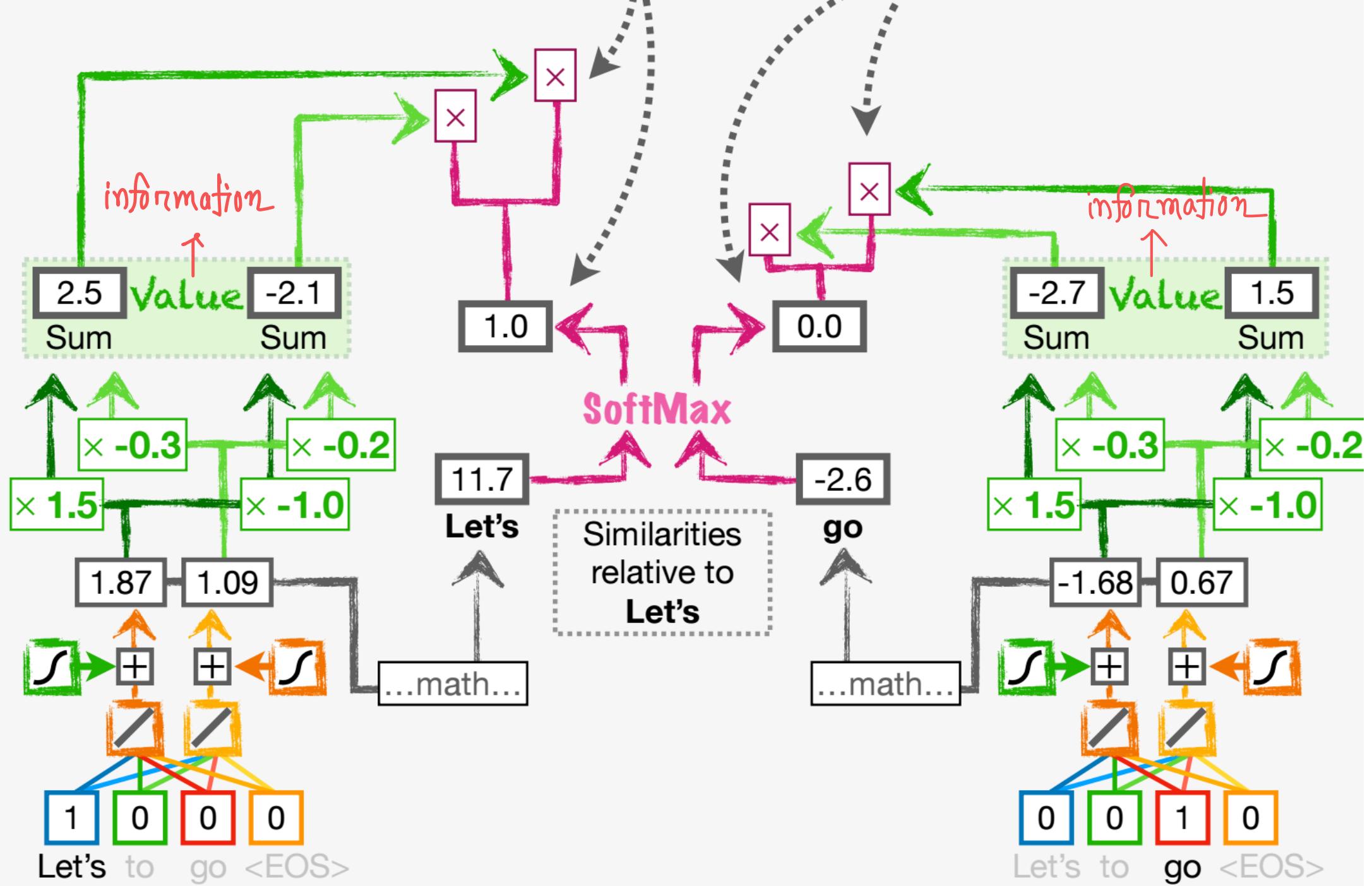


$$\frac{e^{11.7}}{e^{11.7} + e^{-2.6}} \approx 0.99$$



Pass the similarity score into SoftMax function.

→ What 1. of each input word should we use to encode the word Let.



Transformers: Encoding Details

inner + inner [Self
outer + outer [attention]

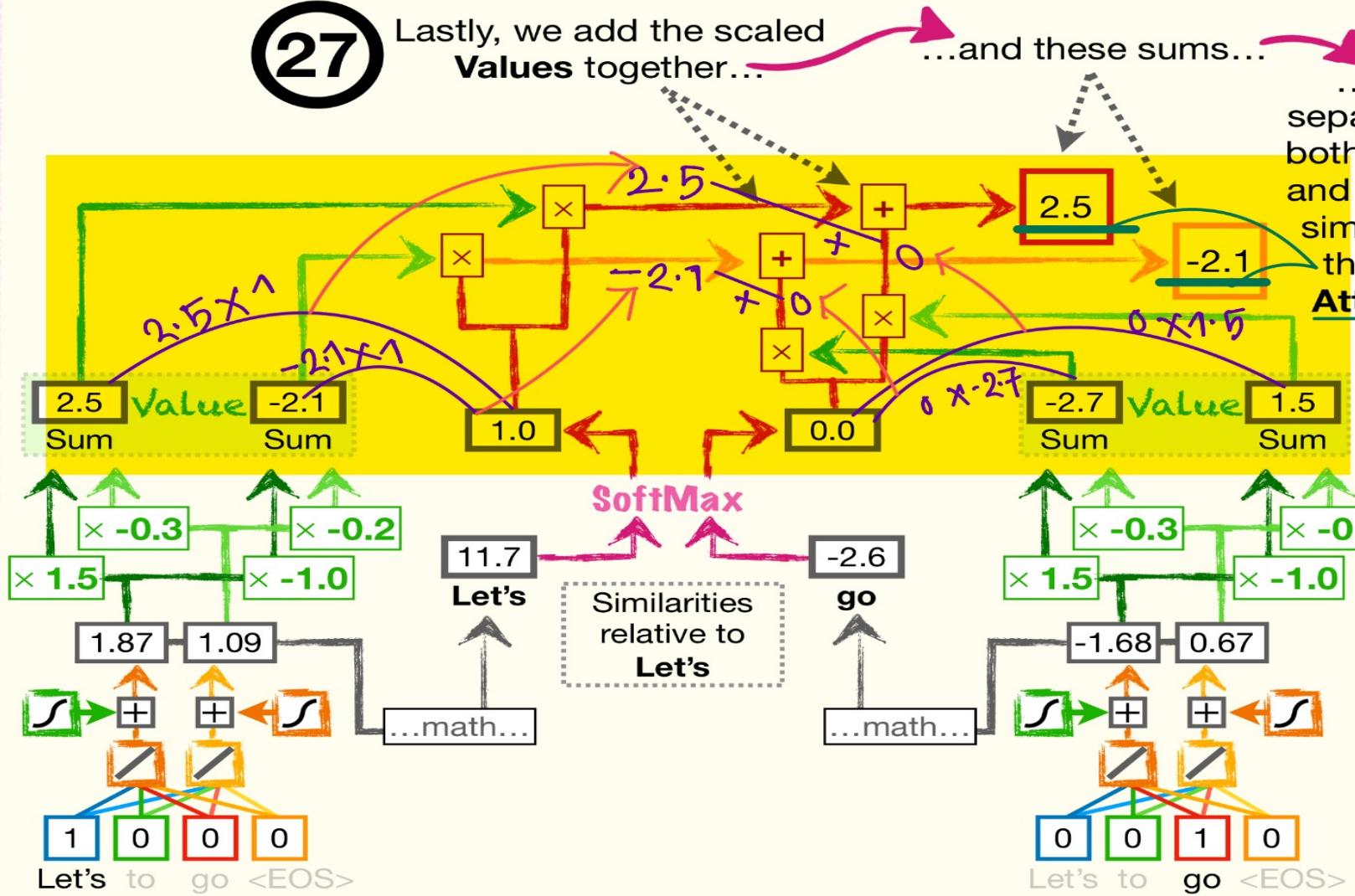
27

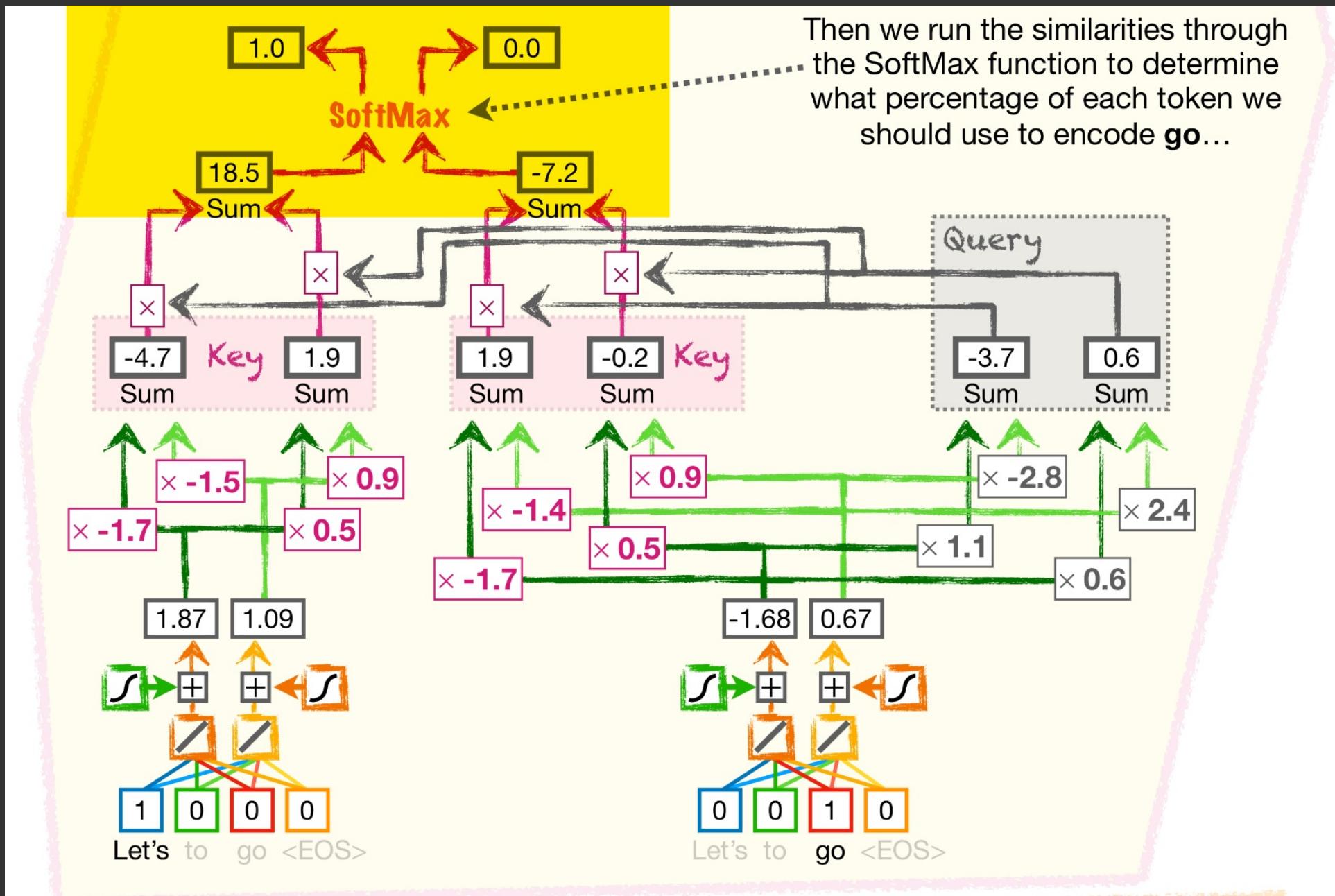
Lastly, we add the scaled
Values together...

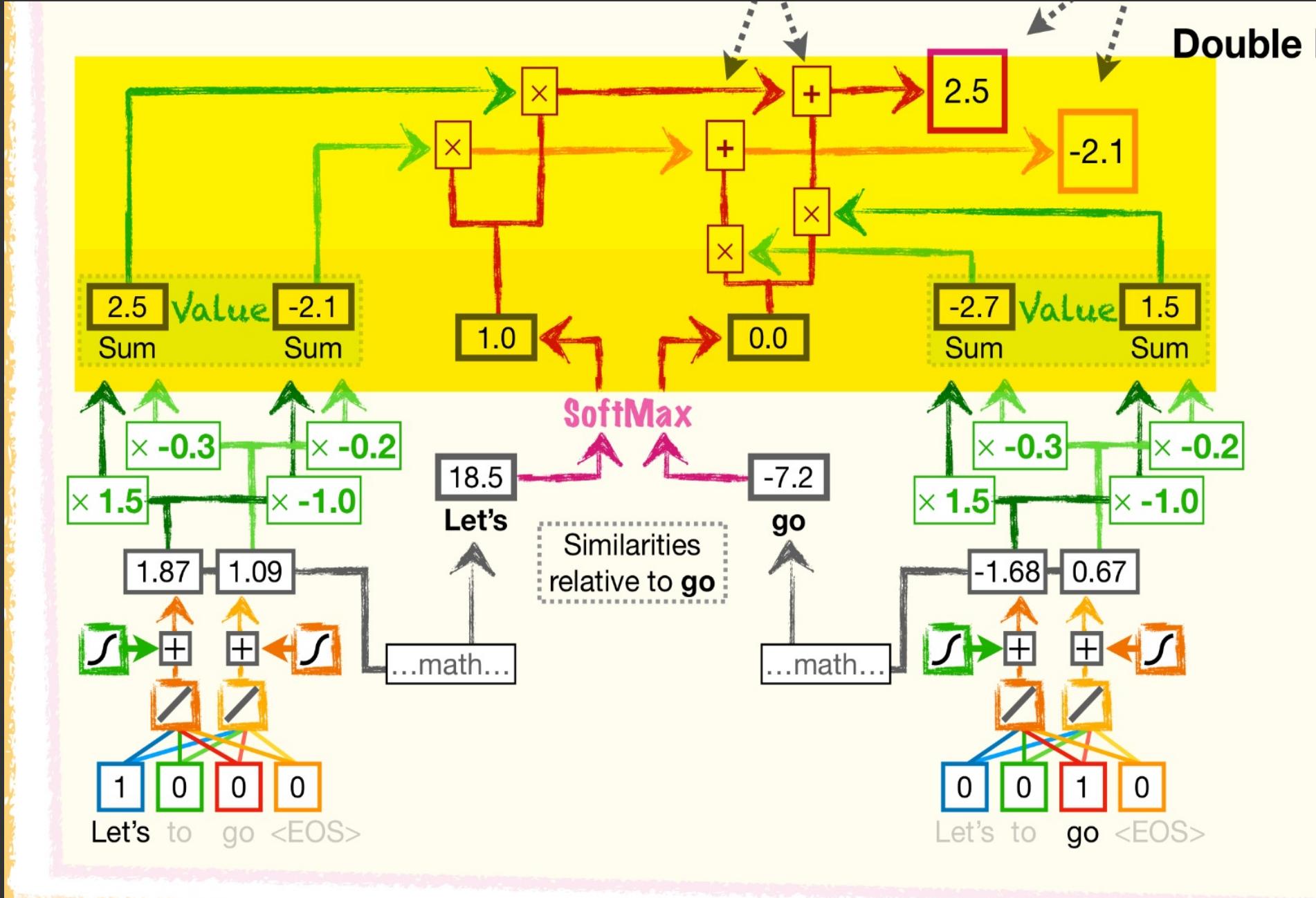
...and these sums...

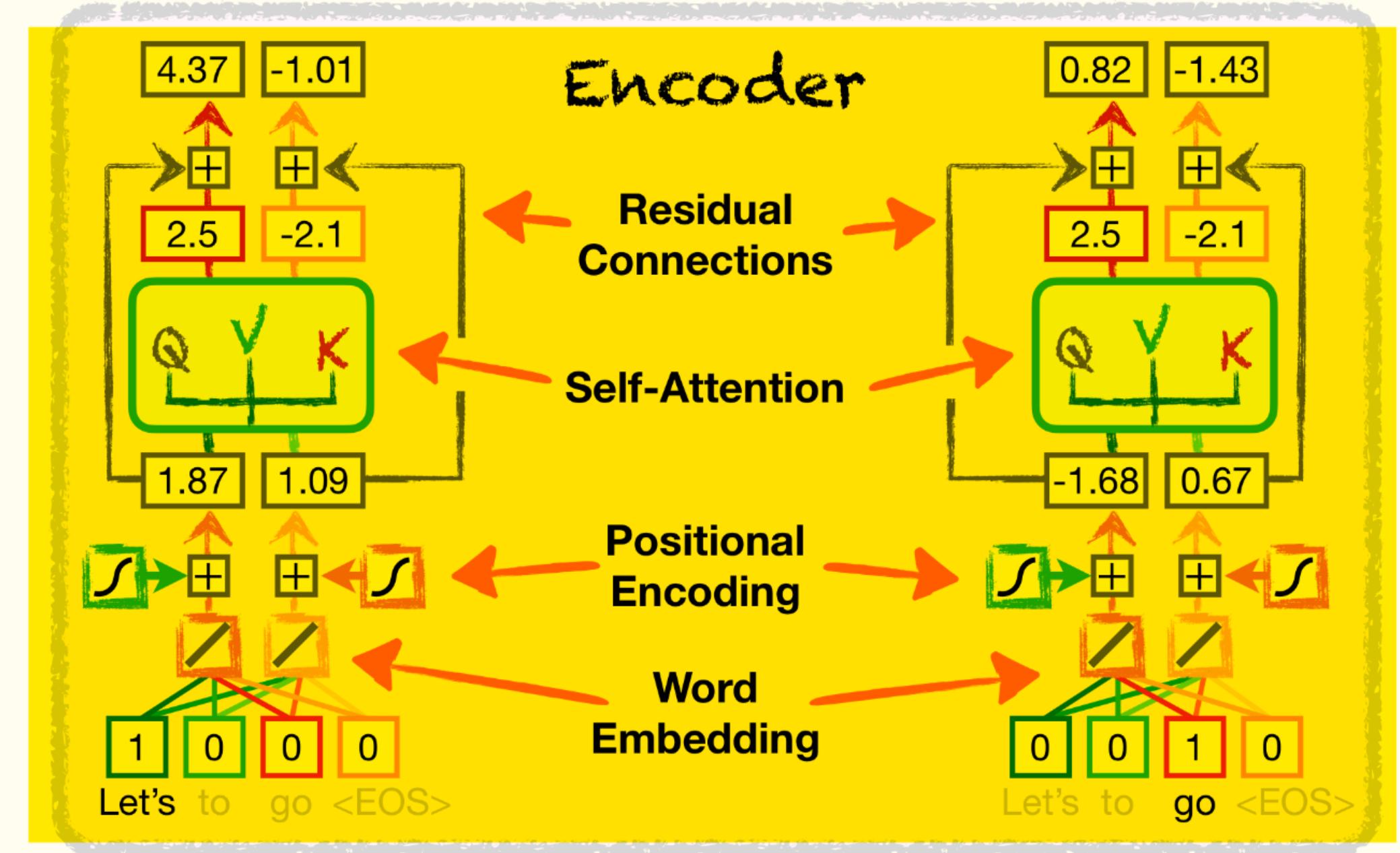
...which combine
separate encodings for
both input words, **Let's**
and **go**, relative to their
similarity to **Let's**, are
the individual Self-
Attention values for
Let's.

Bam!







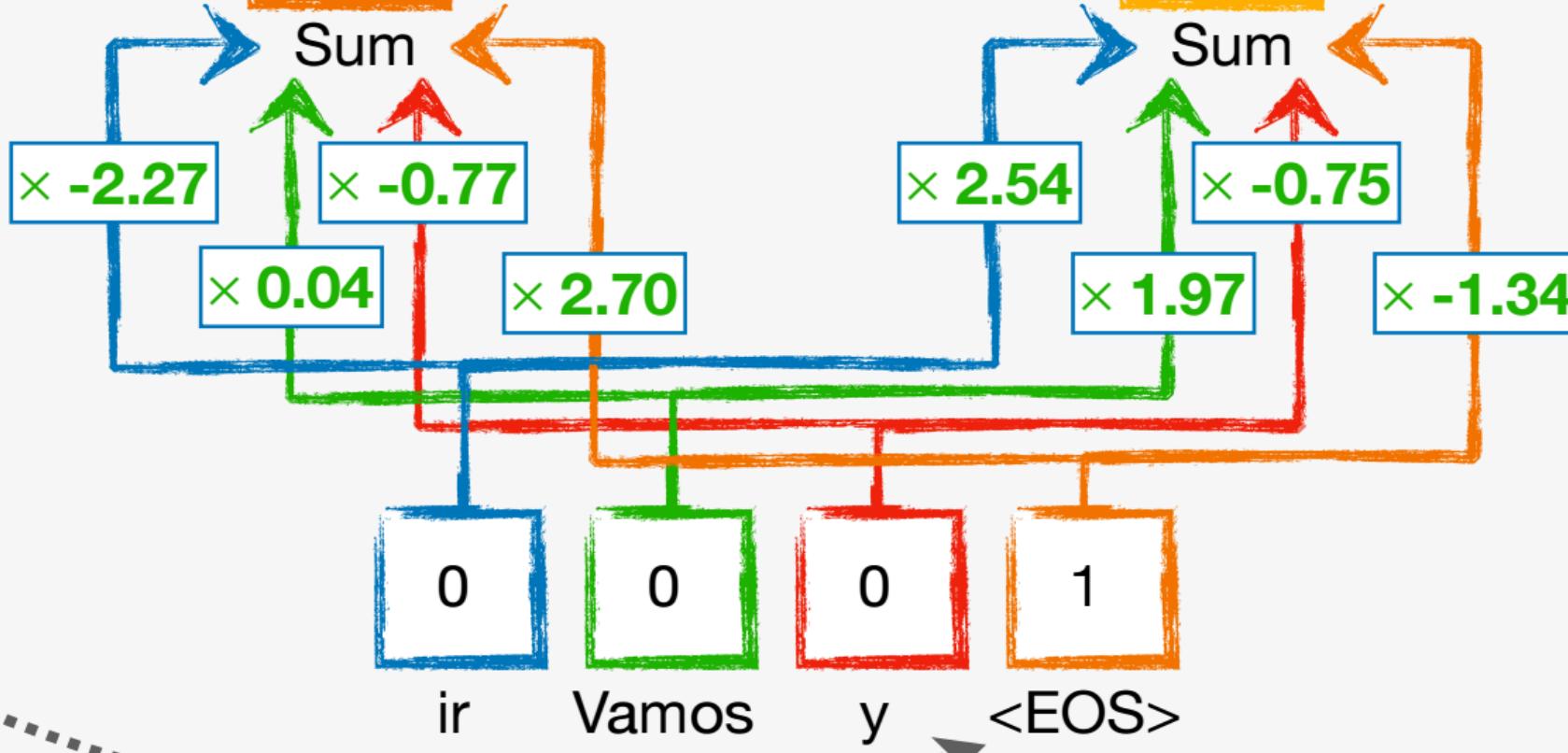


ter 10
odels
ntion,
Word

2.70

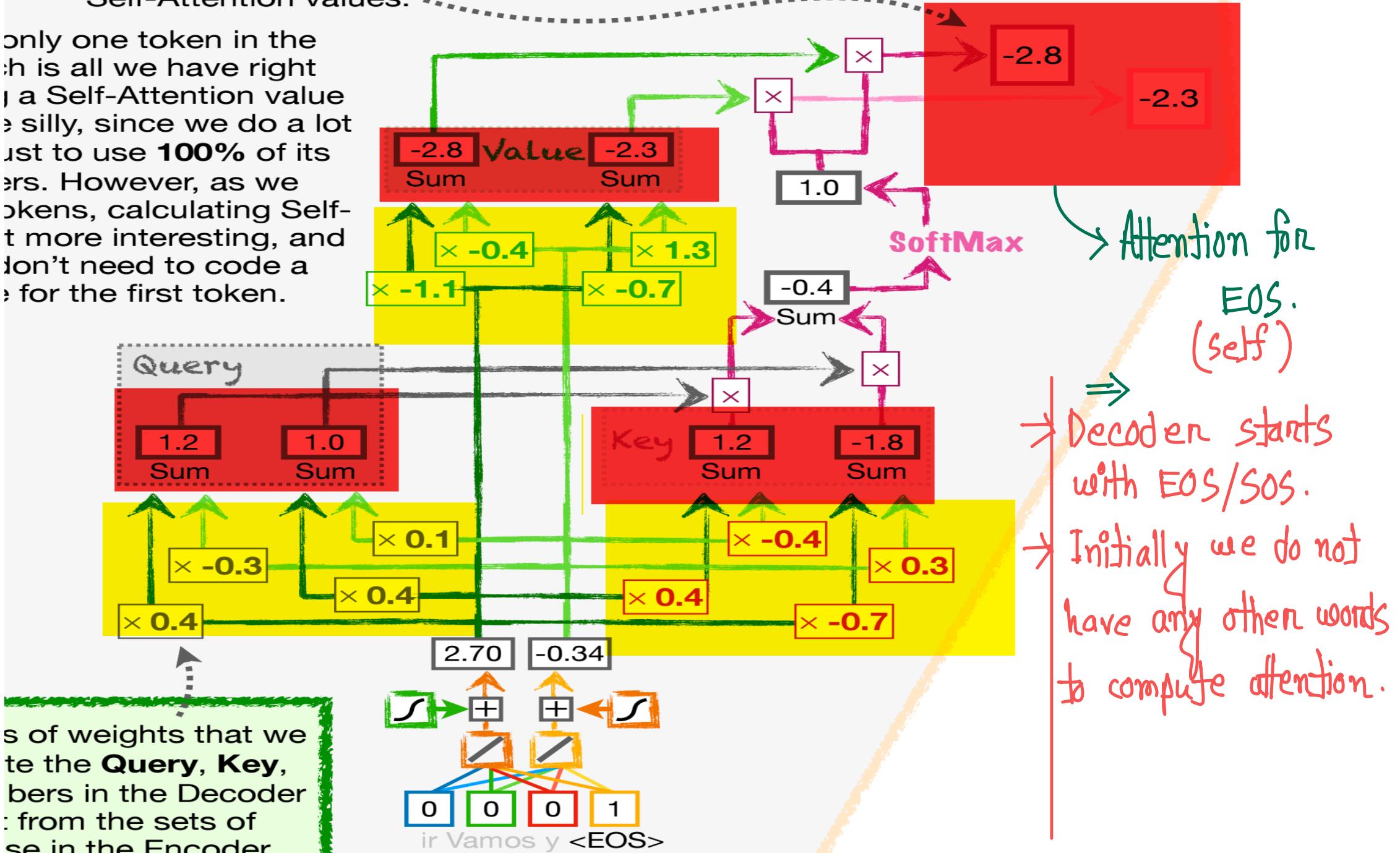
-1.34

Now we work
with
Decoder

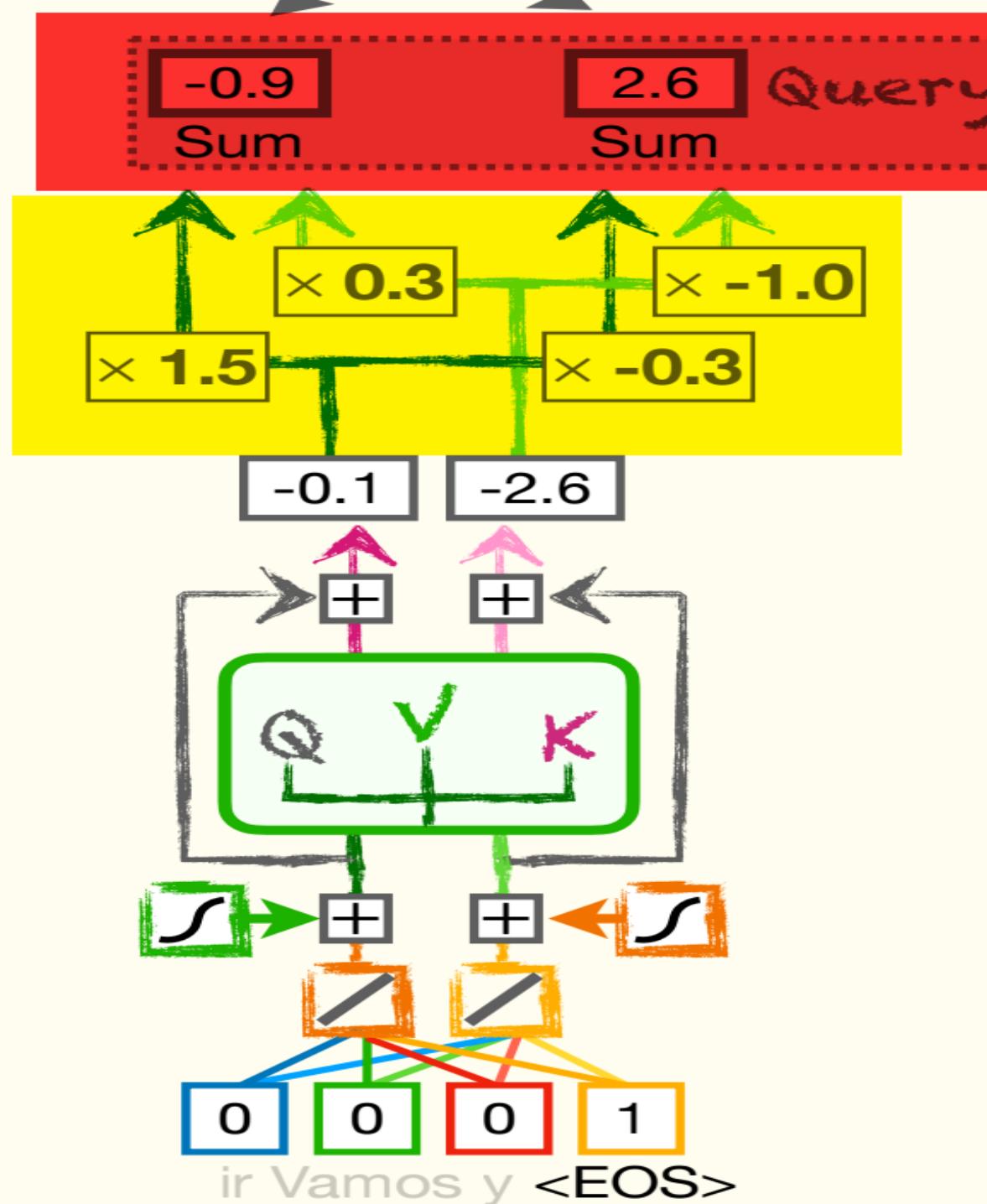


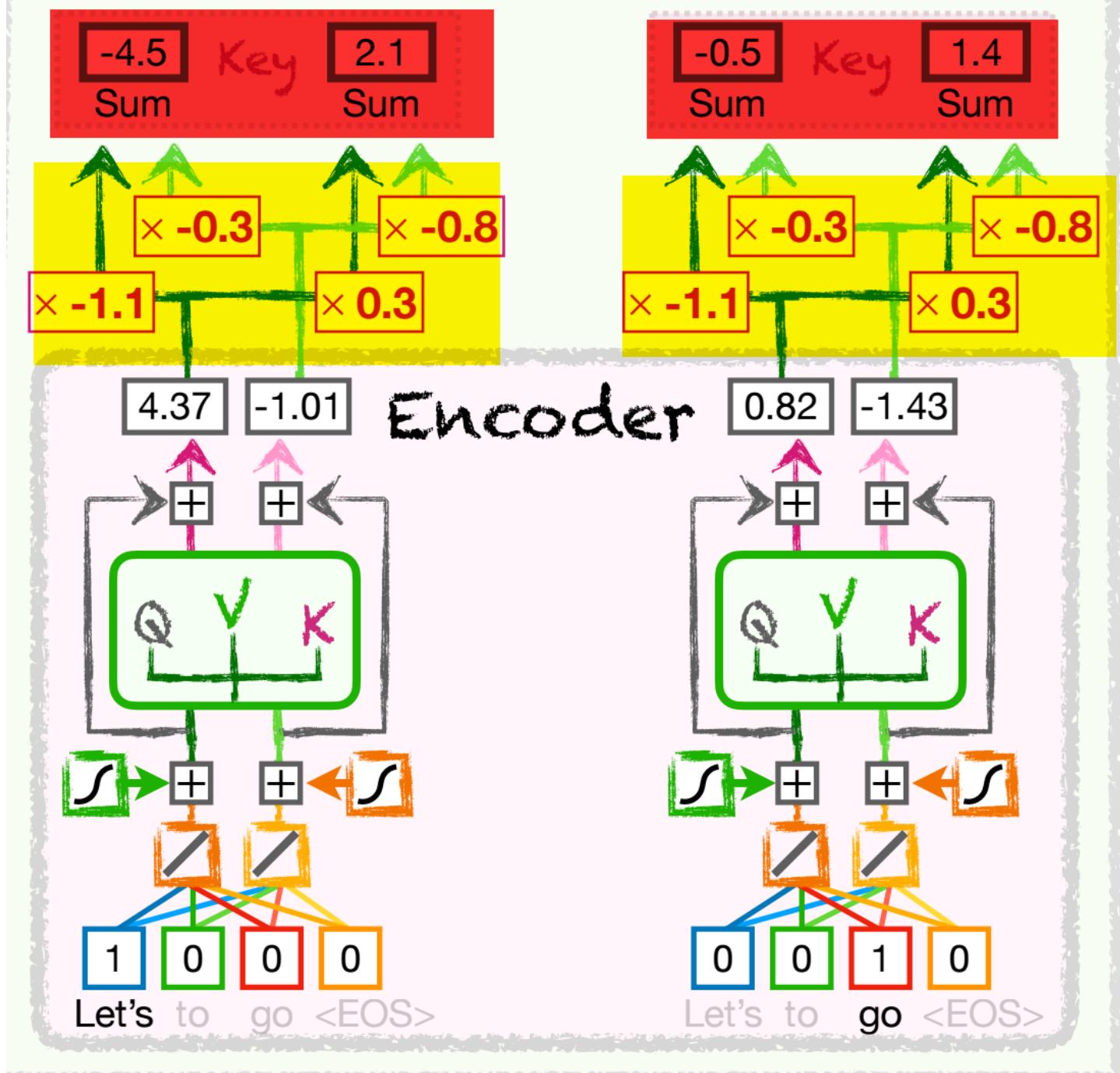
ansformer, we calculate the individual Self-Attention values.

only one token in the batch is all we have right now. A Self-Attention value is silly, since we do a lot just to use 100% of its powers. However, as we add tokens, calculating Self-Attention becomes more interesting, and we don't need to code a case for the first token.

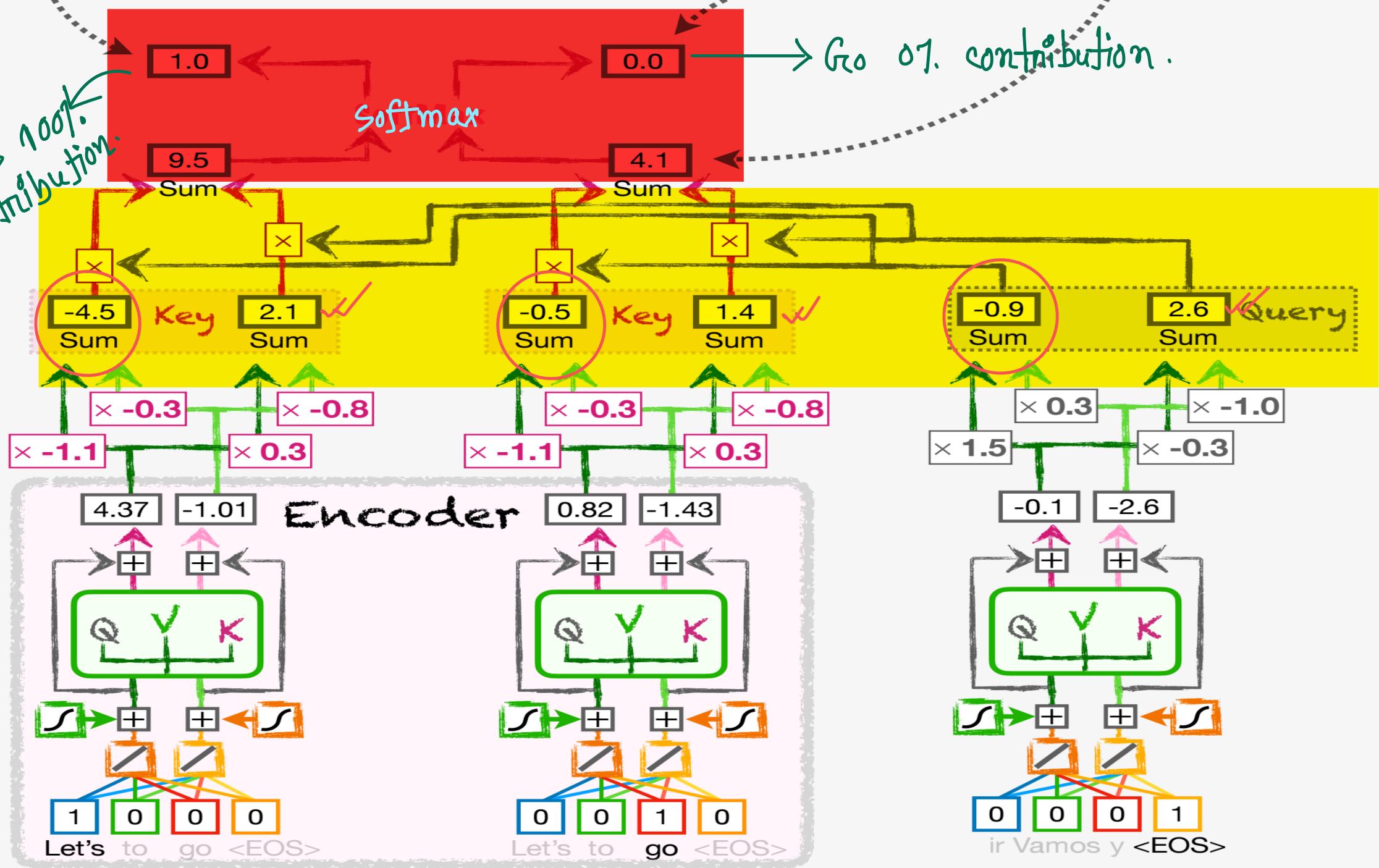


Now we will compute
attention: $\text{EOS} \rightarrow \text{Go}$

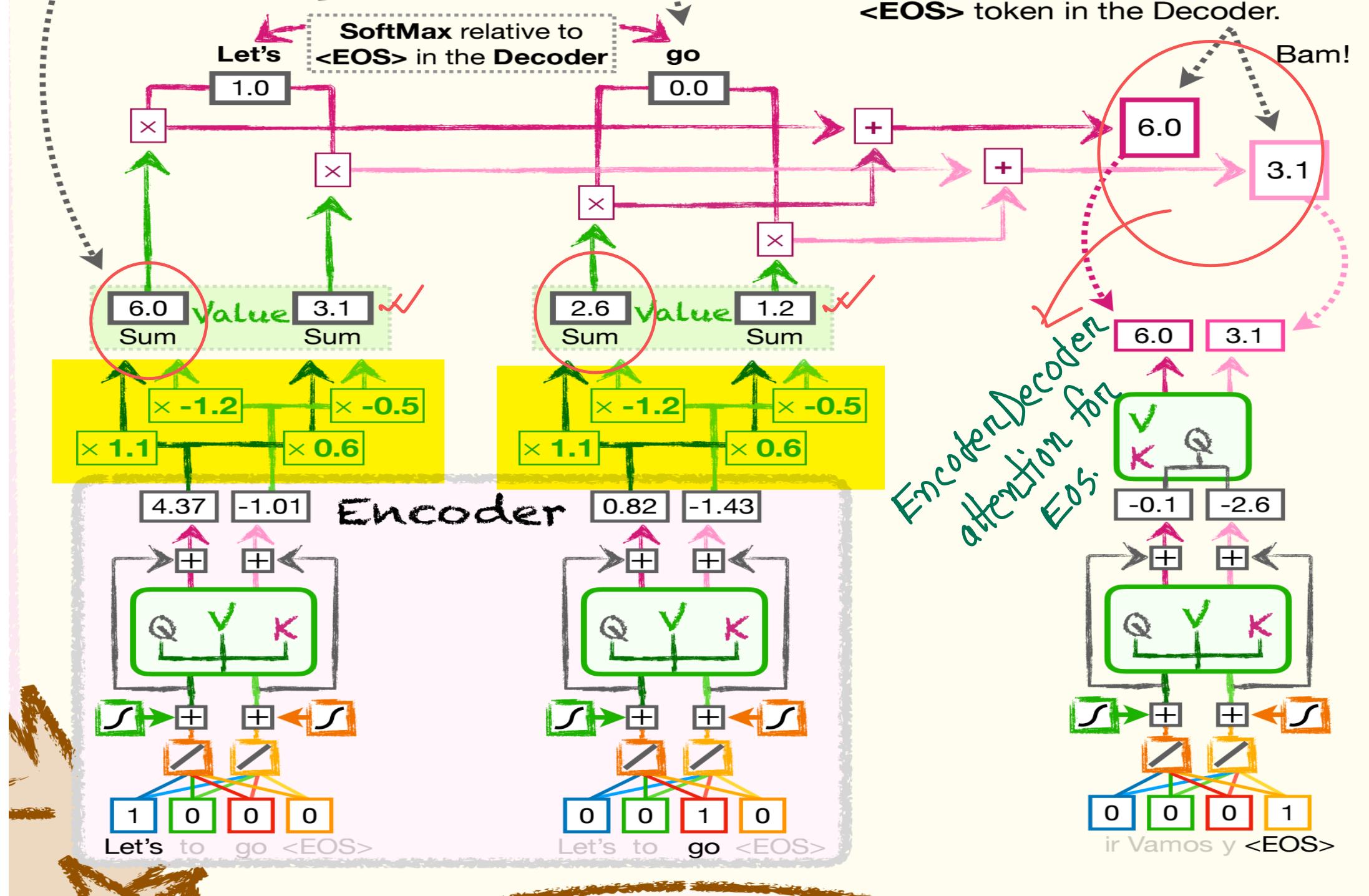




Let's 100% contribution.

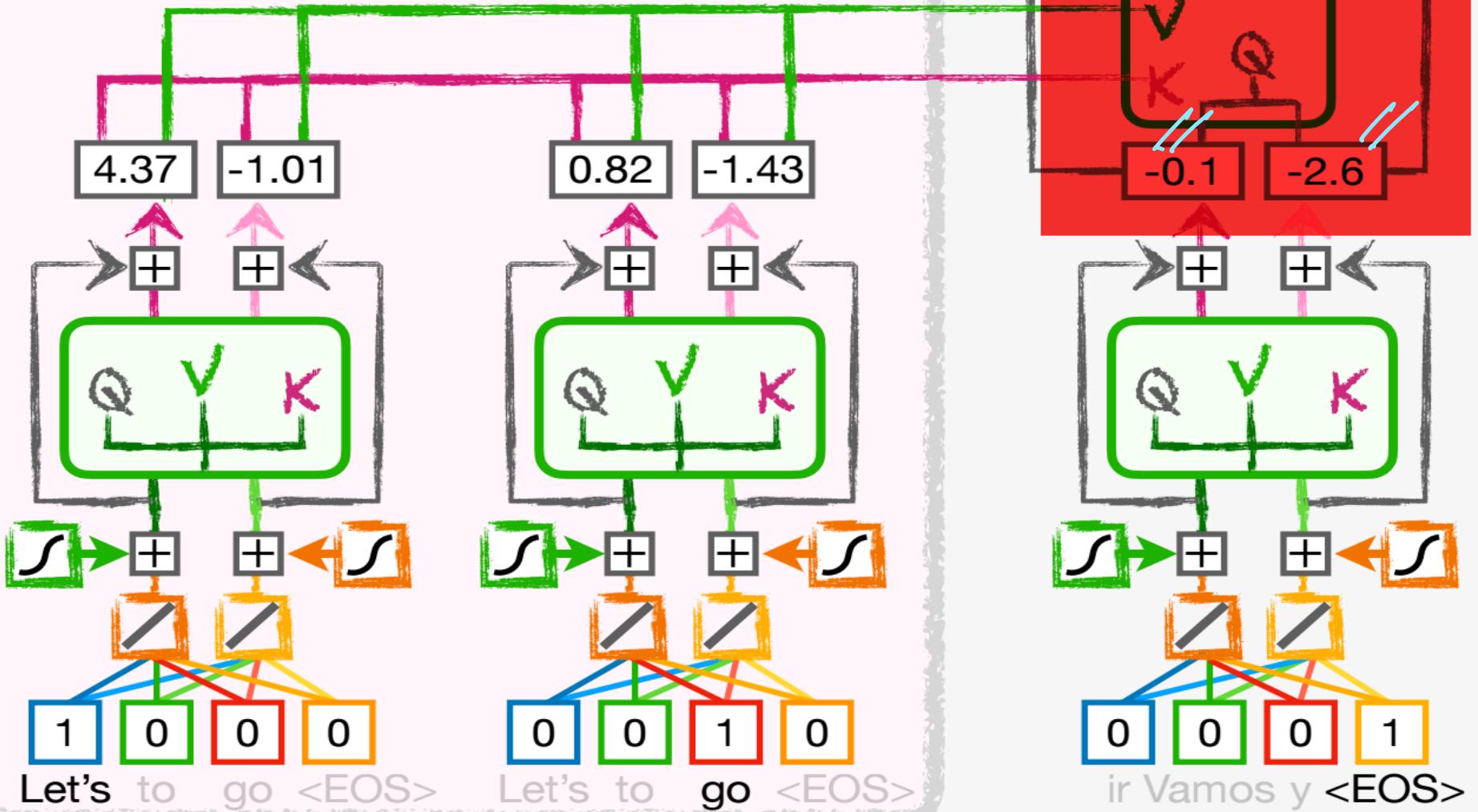


Go 07. contribution.



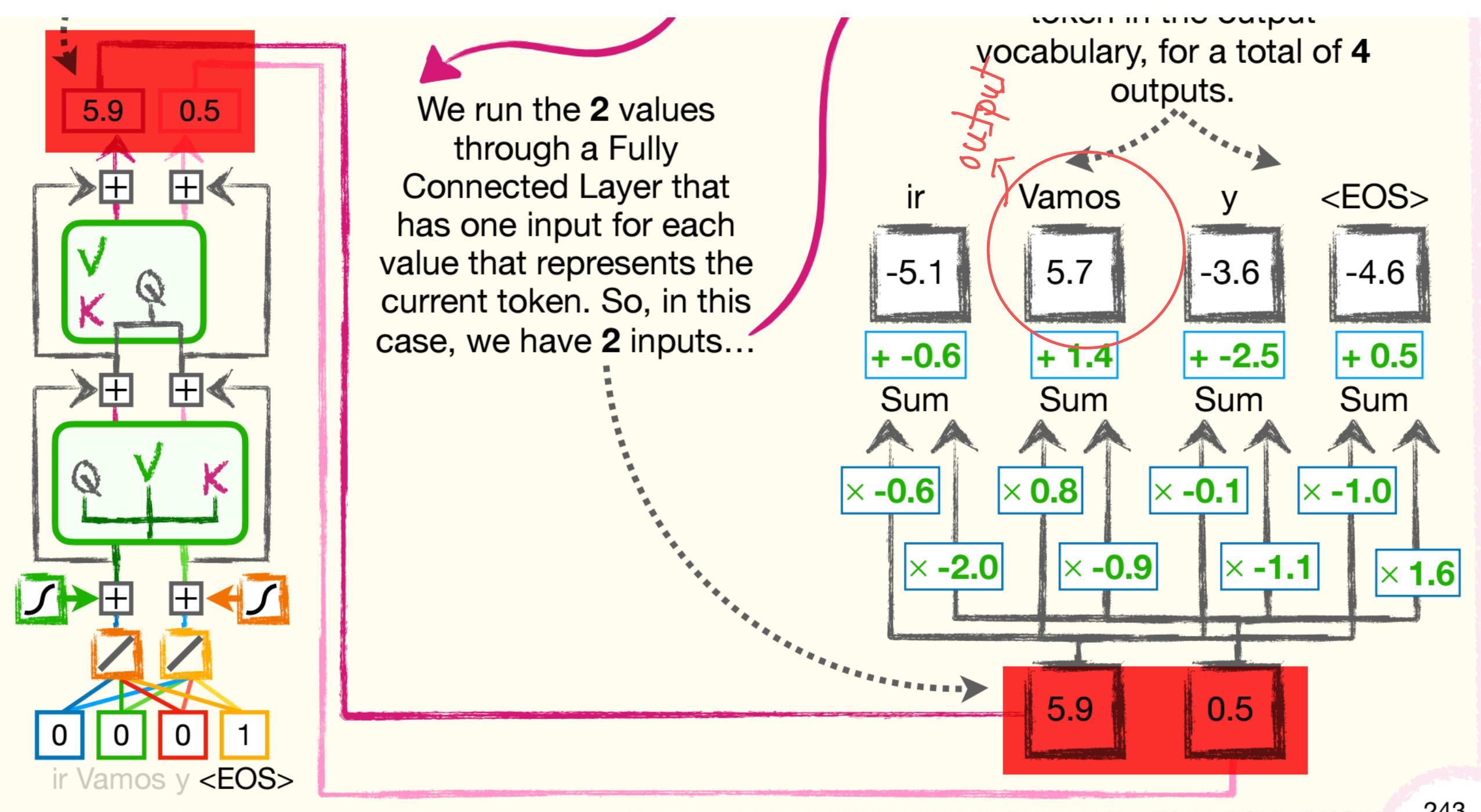
for the <EOS> token...

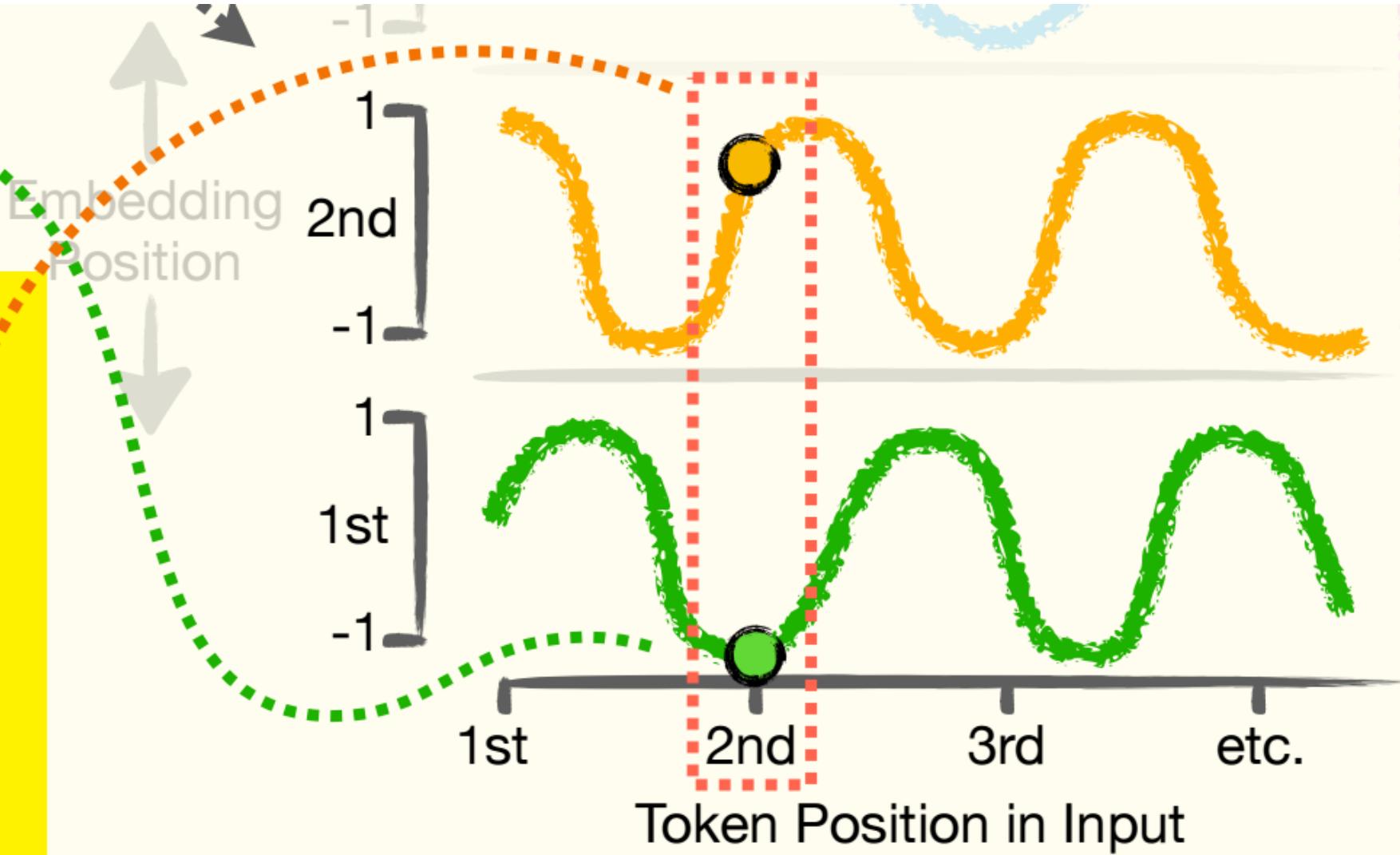
Encoder

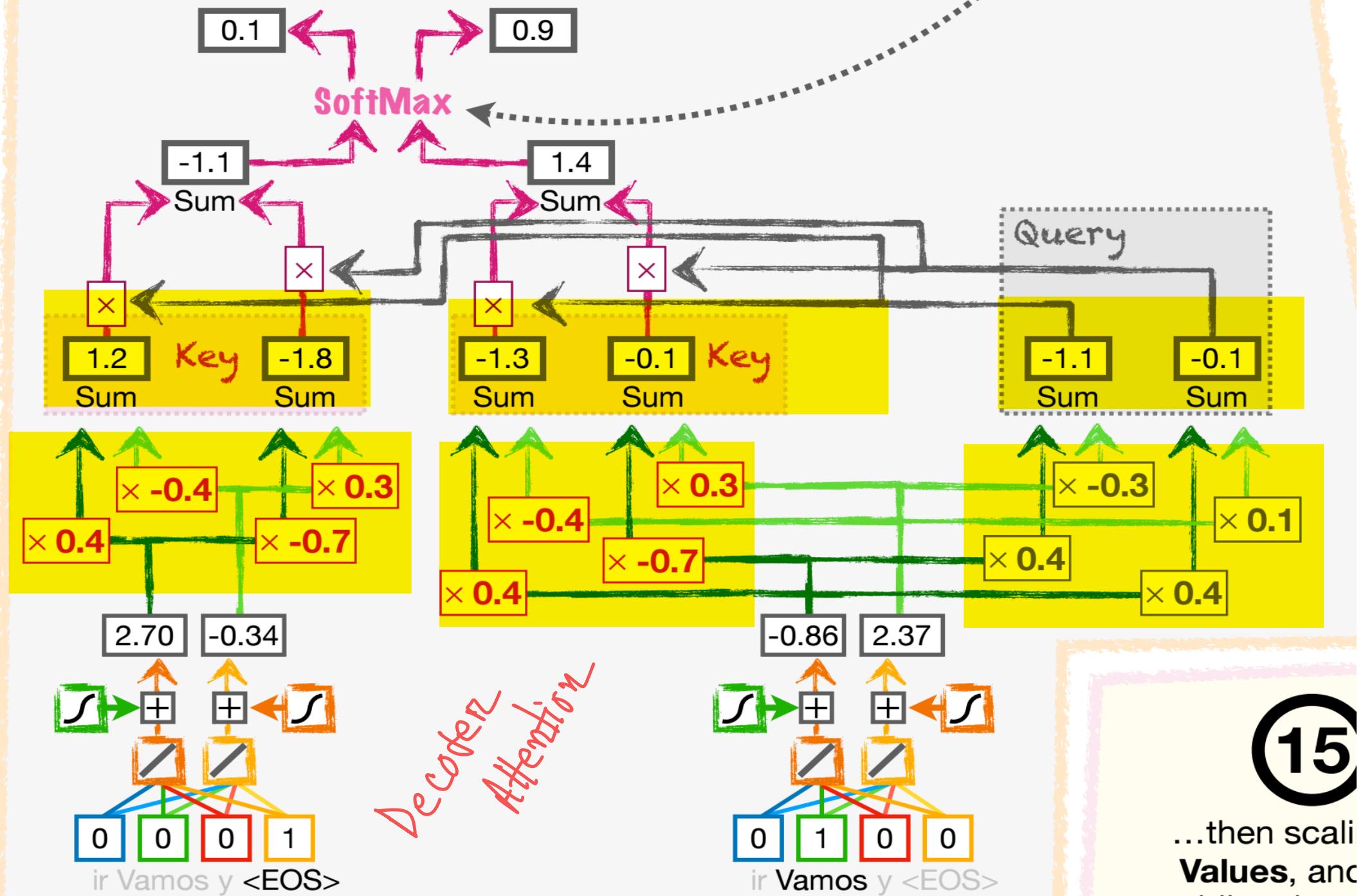


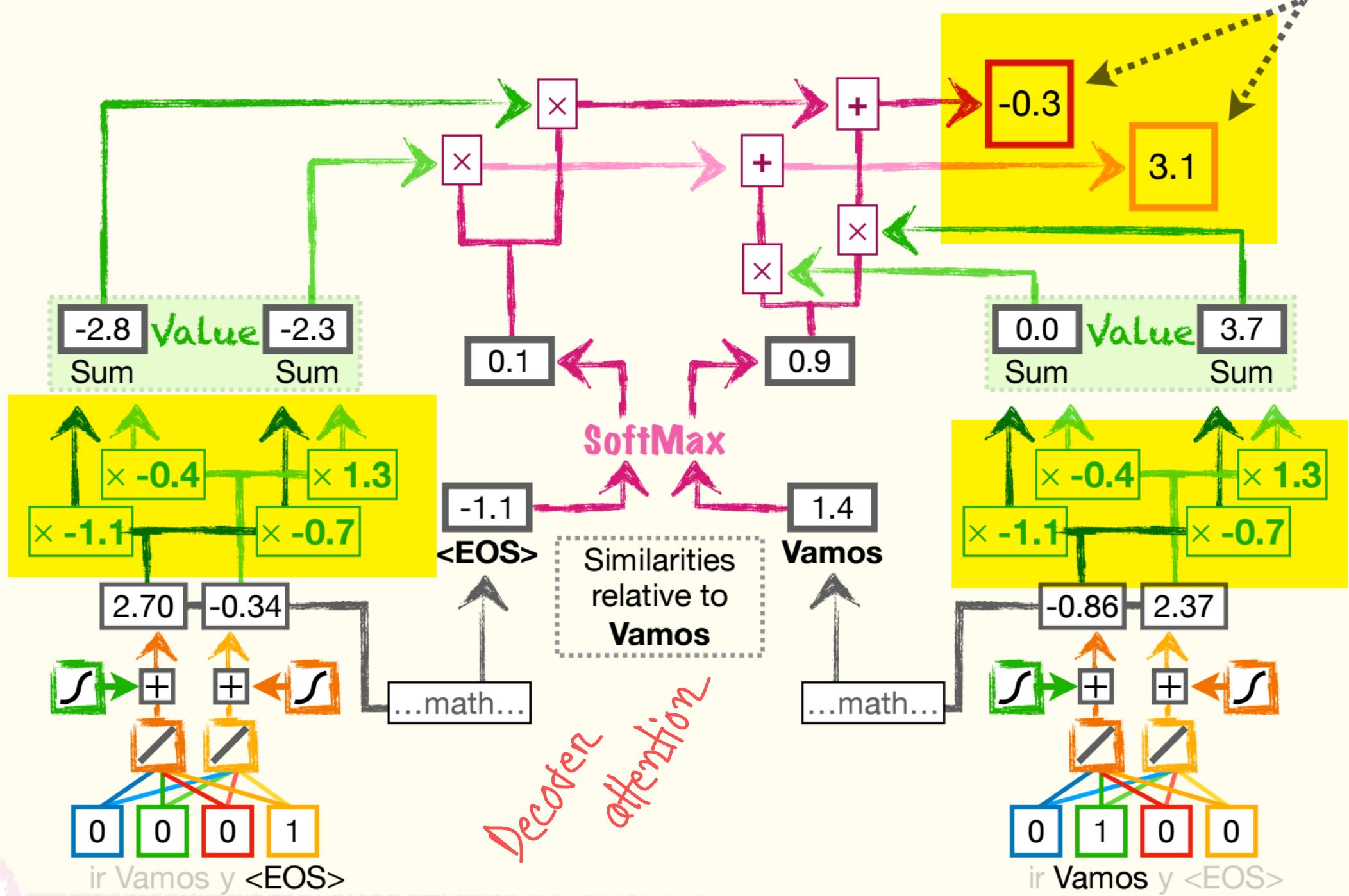
Decoder
of the
pre
op

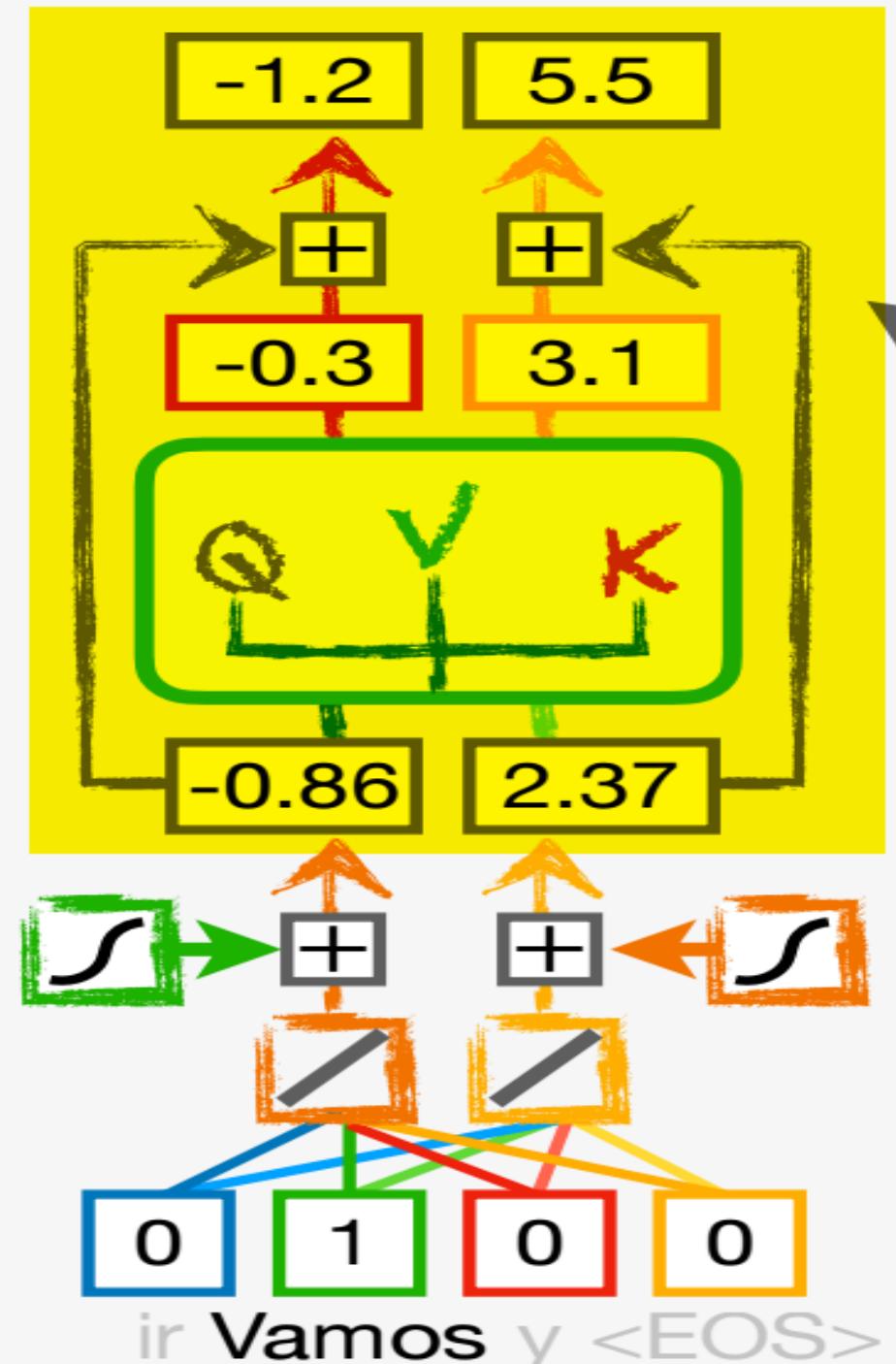
Attention +
PC

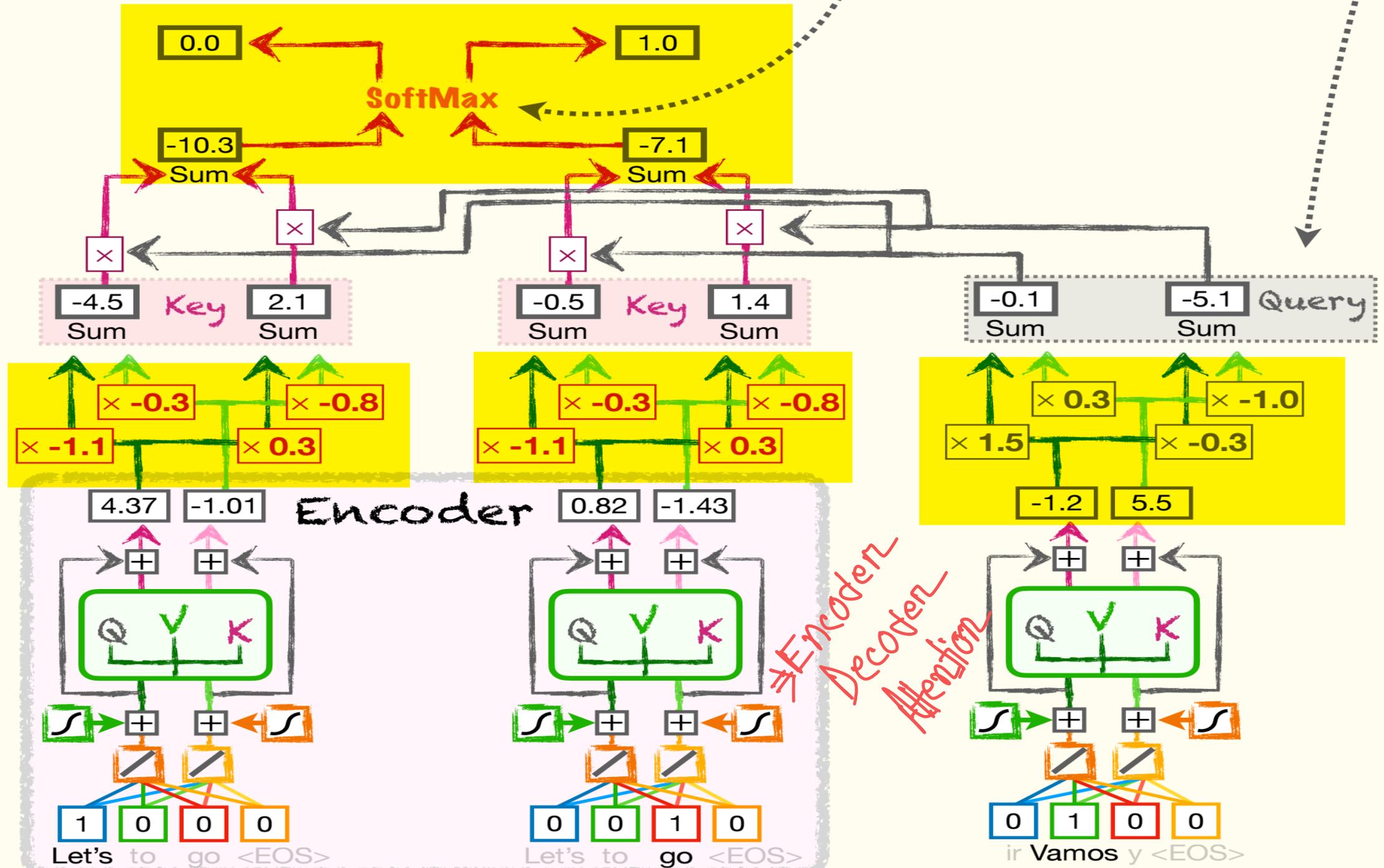


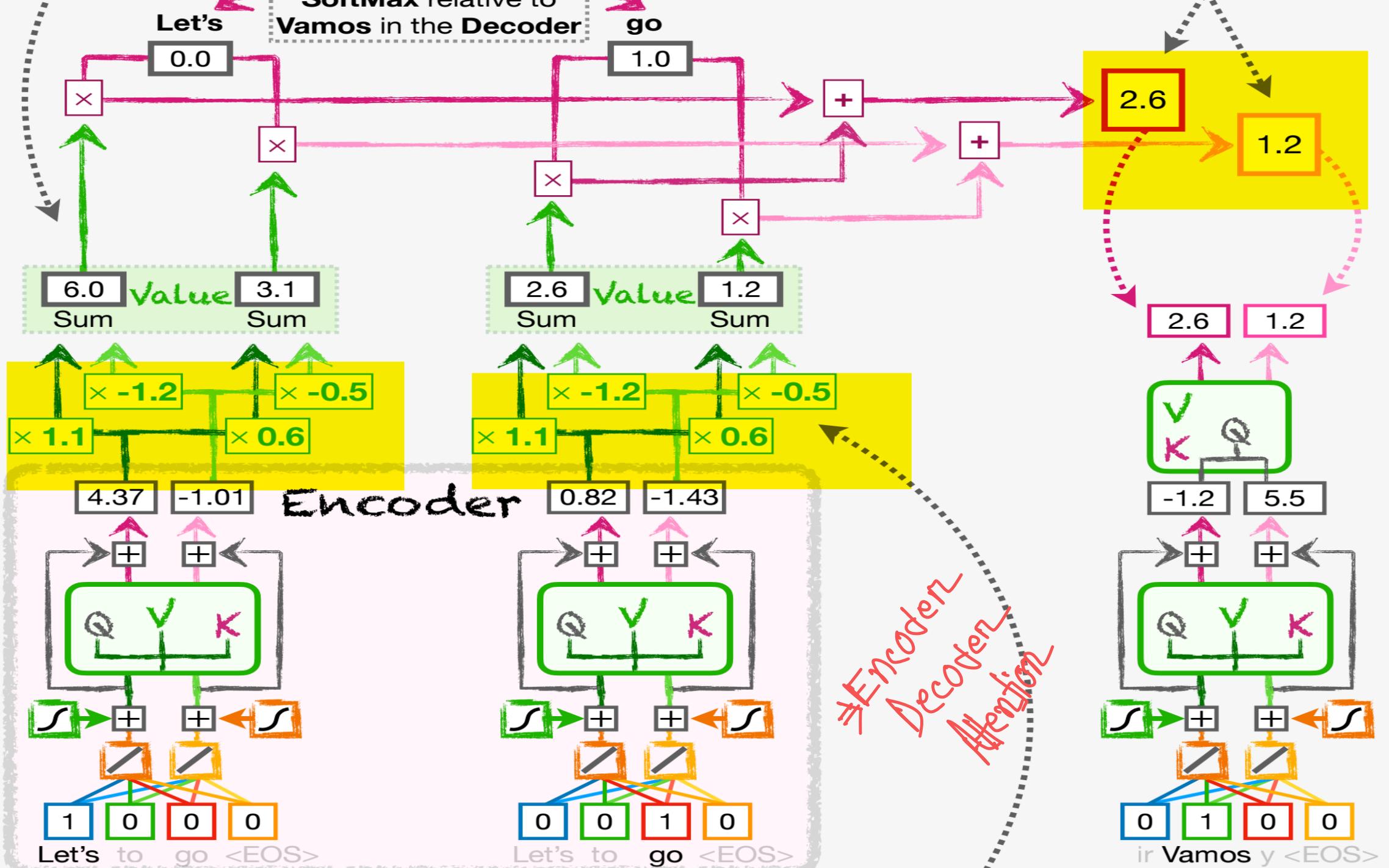




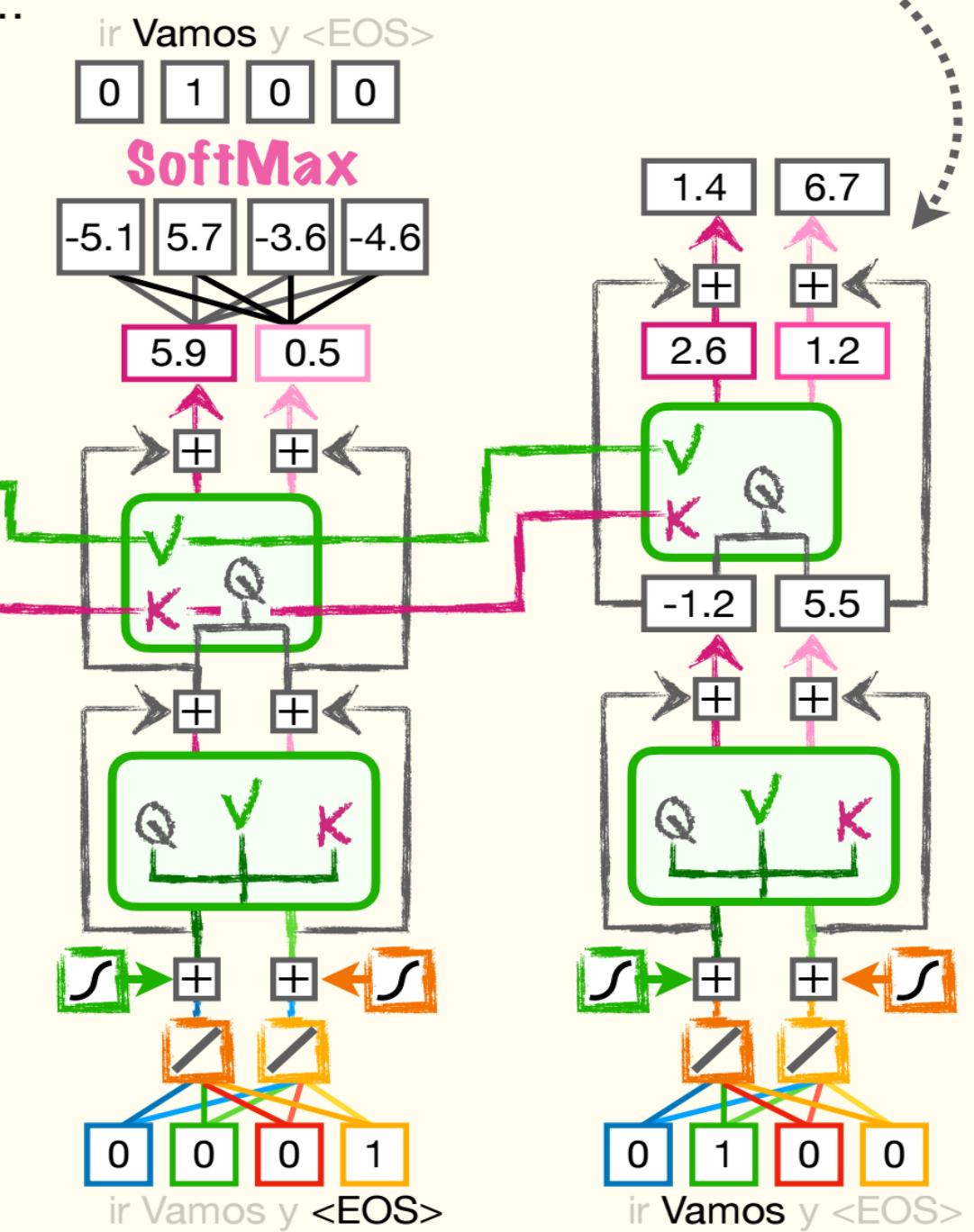
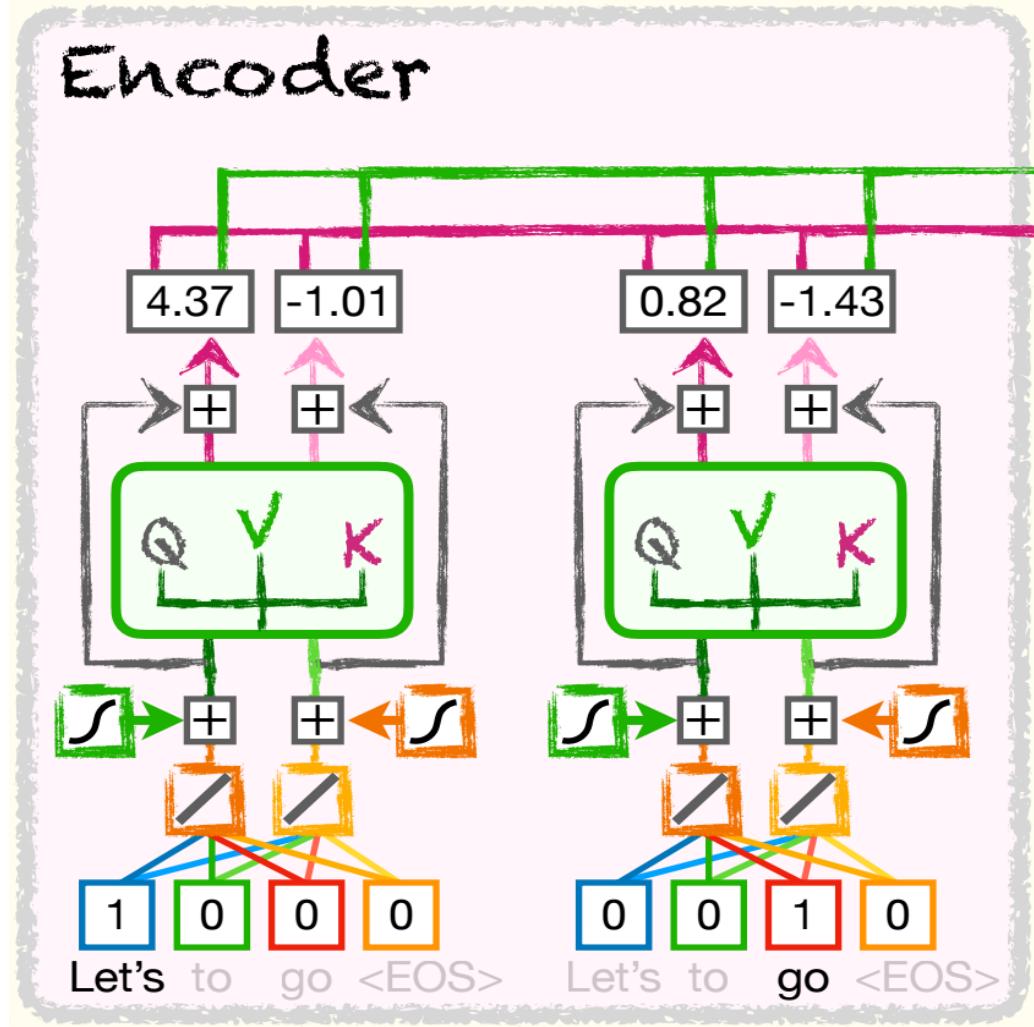


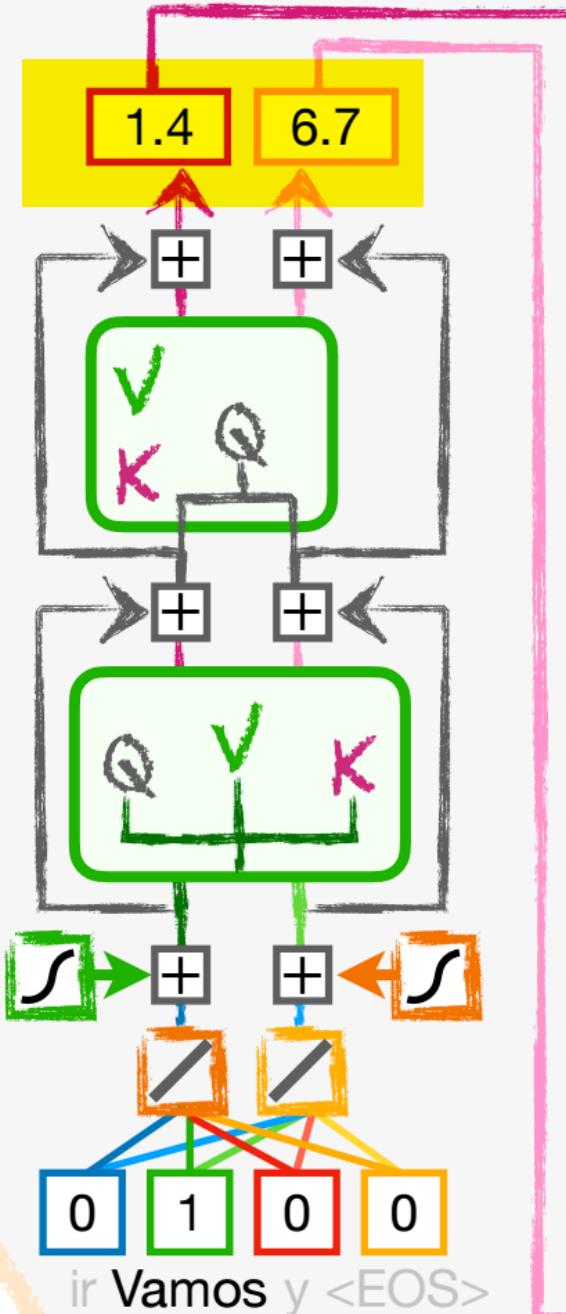




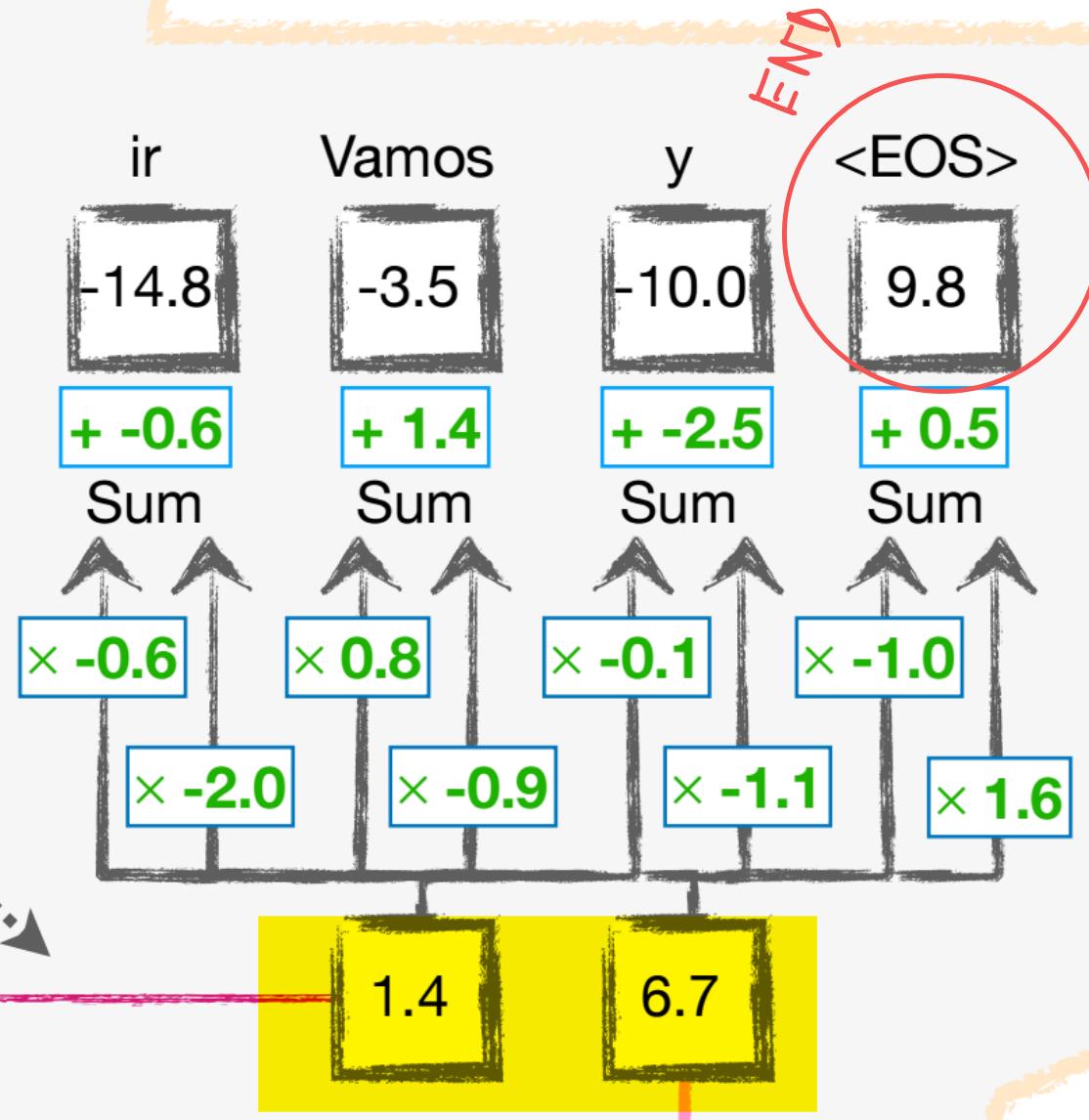


16 Residual Connections...

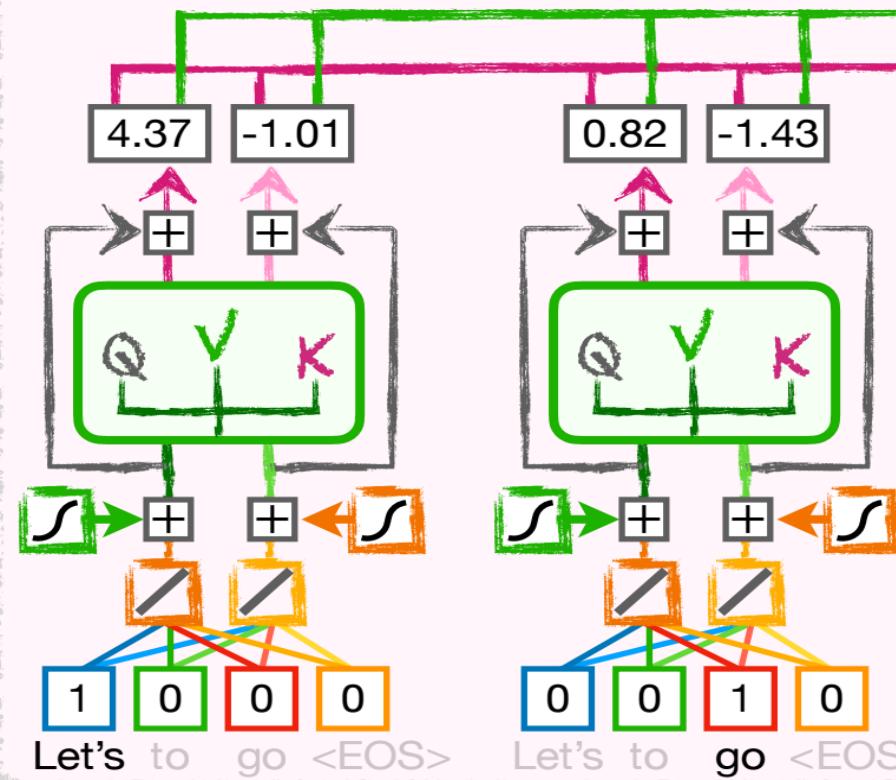




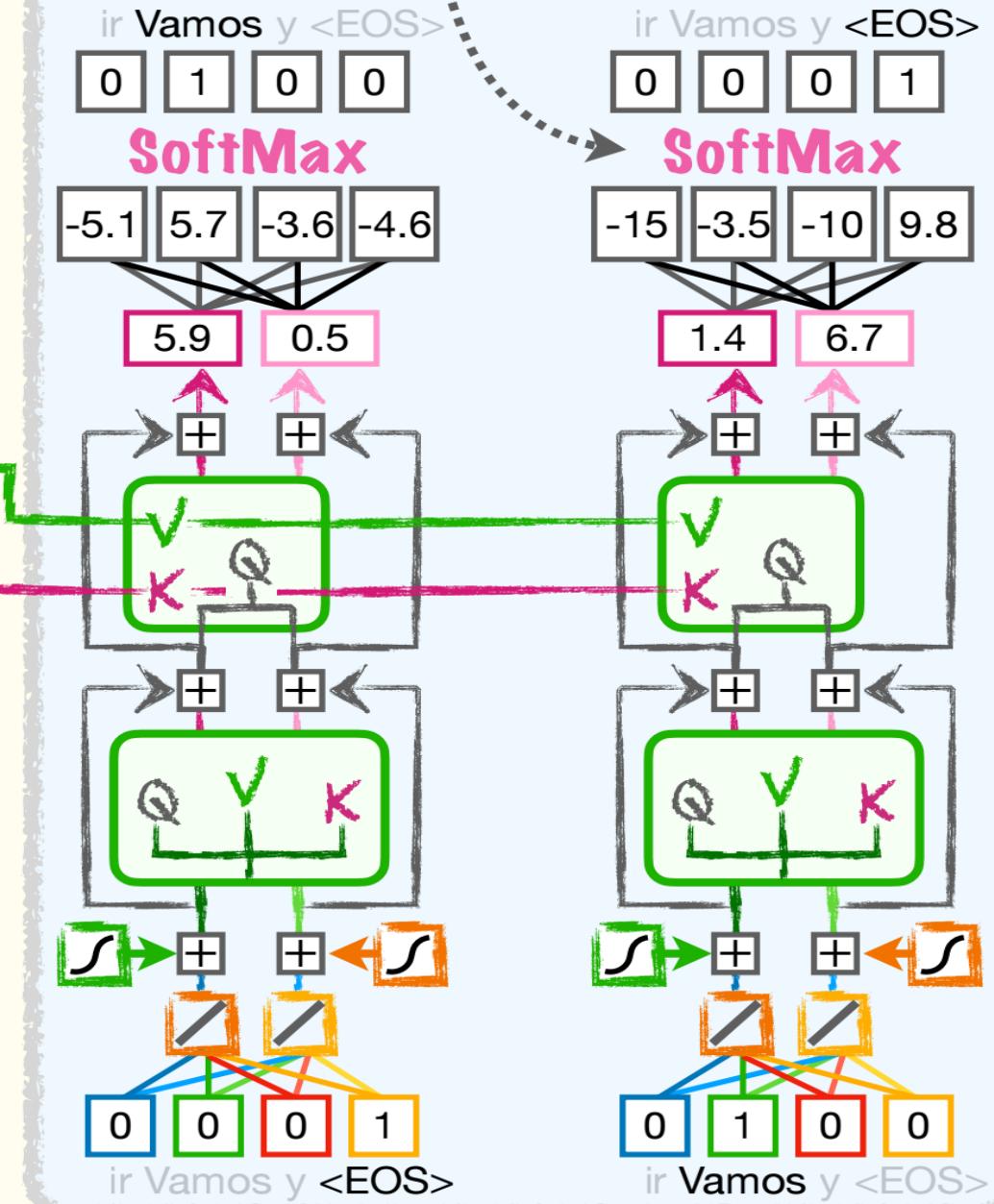
...and then run the sums
through the Fully
Connected Layer...



Encoder



Decoder



Encoder

Lets	go
[1.87, 0.09]	[-0.78, 0.27]
Positional encoding	
+[0.0, 1, 0]	+[-0.9, 0.4]
[1.87, 1.09]	[-1.68, 0.67](1)
Attention calculation in encoder (Lets go)	
Q=1.1 0.6 -2.8 2.4	K = -1.7 0.5 -1.5 0.9
V = 1.5 -1.0 -0.3 -0.2	
QL = Lets X Q	KL = Lets X K
VL = Lets X V	
QG = go X Q	KG = go X K
VG = go X V	
QL = [-1.0, 3.7]	KL = [-4.7, 1.9]
VL = [2.5, -2.1]	
QG = [-3.7, 0.6]	KG = [1.9, -0.2]
VG = [-2.7, 1.5]	
QLxKL = 11.7	QLxKG = -2.6
QGxKL = 18.5	
QGxKG = -7.2	
Softmax(11.7, -2.6) 1.0, 0.0	Softmax(18.5, -7.2) 1.0, 0.0
VLx1.0 + VGx0 [2.5, -2.1]	VLx1.0 + VGx0.0 [2.5, -2.1]

From (1) [1.87, 1.09] [-1.68, 0.67]

+ [2.5, -2.1] + [2.5, -2.1]

[4.37, -1.01] [0.82, -1.43]

Decoder

Attention calculation in decoder (EOS)

EOS [2.70, -1.34] + positional encoding [0.0, 1.0] = [2.70, -0.34]

Q=0.4 0.4
-0.3 0.1
K = 0.4 -0.7
-0.4 0.3
V = -1.1 -0.7
-0.4 1.3

QE = EOSXQ KE = EOSXK VE = EOSXV

QE = [1, 2, 1, 0] KE = [1.2, -1.8] VE = [-2.8, -2.3]

QEXKE = -0.4 Softmax(-0.4) = 1.0

EOS[2.70, -0.34] + VE X 1.0 = [-0.1, -2.6](2)

$$PE(pos, 2i) = \sin \frac{pos}{10000^{\frac{2i}{d_{model}}}}$$

$$PE(pos, 2i+1) = \cos \frac{pos}{10000^{\frac{2i+1}{d_{model}}}}$$

Decoder Encoder cross attention calculation (EOS, Lets, go)

$$\begin{array}{lll} Q = \begin{pmatrix} 1.5 & -0.3 \\ 0.3 & -1.0 \end{pmatrix} & K = \begin{pmatrix} -1.1 & 0.3 \\ -0.3 & -0.8 \end{pmatrix} & V = \begin{pmatrix} 1.1 & 0.6 \\ -1.2 & -0.5 \end{pmatrix} \end{array}$$

$$\begin{aligned} QE &= EOS \times Q & KL &= Lets \times K & KG &= go \times K \\ VL &= Lets \times V & VG &= go \times V \end{aligned}$$

$$\begin{aligned} QE &= [-0.9, 2.6] & KL &= [-4.5, 2.1] & VL &= [6.0, 3.1] \\ KG &= [-0.5, 1.4] & VG &= [2.6, 1.2] \end{aligned}$$

$$QExKL = 9.5 \quad QExKG = 4.1$$

Softmax(9.5, 4.1)

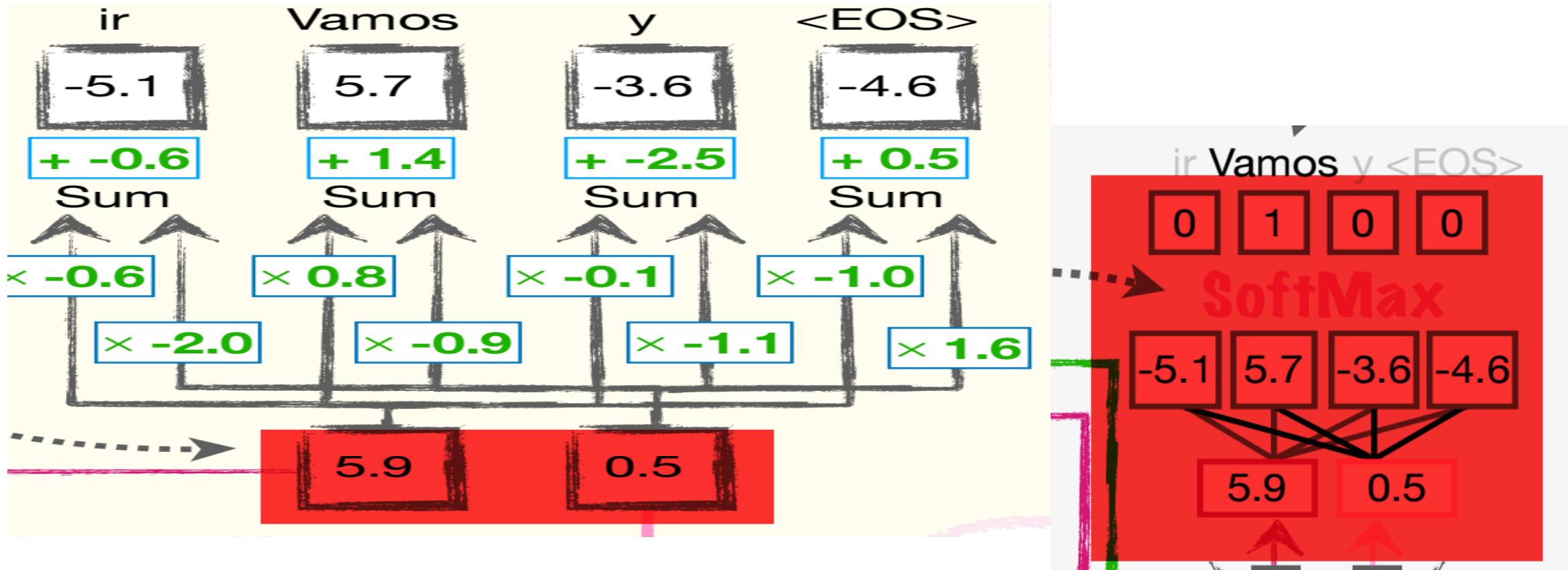
1.0, 0.0

VLx1.0 + VGx0

[6.0, 3.1]

$$\text{From (2) EOS } [-0.1, -2.6] + [6.0, 3.1] = [5.9, 0.5] \dots\dots\dots (3)$$

Generating first token in decoder after processing EOS



Decoder

Attention calculation in decoder (Vamos, EOS)

Vamos [-0.86,2.37] + positional encoding [-0.9,0.4] = [-0.86,2.37]

$$Q = \begin{pmatrix} 0.4 & 0.4 \\ -0.3 & 0.1 \end{pmatrix}, \quad K = \begin{pmatrix} 0.4 & -0.7 \\ -0.4 & 0.3 \end{pmatrix}, \quad V = \begin{pmatrix} -1.1 & -0.7 \\ -0.4 & 1.3 \end{pmatrix}$$

KE = EOSXK VE = EOSXV QV = Vamos X Q VV = vamos
X V

$$KE = [1.2, -1.8] \quad VE = [-2.8, -2.3] \quad QV = [-1.3, -0.1] \quad VV = [0.0, 3.7]$$

QVXKE = -1.1 QVXKV = 1.4

$$\text{Softmax}(-1.1, -1.4) = 0.1, 0.9$$

varmos[-0.86, 2.37] + VE X 0.1 + VVX0.9 = [-1.2, 5.5]

.....(4)

Decoder Encoder cross attention (Vamos, Lets, go)

$$Q = \begin{pmatrix} 1.5 & -0.3 \\ 0.3 & -1.0 \end{pmatrix}, \quad K = \begin{pmatrix} -1.1 & 0.3 \\ -0.3 & -0.8 \end{pmatrix}, \quad V = \begin{pmatrix} 1.1 & 0.6 \\ -1.2 & -0.5 \end{pmatrix}$$

QV = vamos X Q KL = Lets X K KG = go X K
VL = Lets X V VG = go X V

$$QV = [-0.1, -5.1] \quad KL = [-4.5, 2.1] \quad VL = [6.0, 3.1]$$

$$KG = [-0.5, 1.4] \quad VG = [2.6, 1.2]$$

$$\begin{aligned} \text{QVxKL} &= -10.3 & \text{QVxKG} &= -7.1 \\ \text{Softmax}(-10.3, -7.1) & \\ 0.0, 1.0 & \\ \text{VLx}0.0 + \text{VGx}1.0 & \\ [2.6, 1.2] & \end{aligned}$$

From (4) Vamos $[-1.2, 5.5] + [2.6, 1.2] = [1.4, 6.7]$
.....(5)

