

# How Transformers Generate and Translate Sequences: The Nuts and Bolts

Brandon Rohrer

files	find	my
1	0	0
0	1	0
0	0	1

	find	my	files
0	0	0	1
1	1	0	0
0	0	1	0









0	0	1	0
---	---	---	---

A

.2
.7
.8
.1

B

=

0
0
1
0

•

.2
.7
.8
.1

=

.8



1	0	0	0
0	0	1	0

A

.2
.7
.8
.1

B

=

Dot product of the first row of A with B
Dot product of the second row of A with B

=

.2
.8

0	0	1	0
---	---	---	---

A

.2	.9
.7	0
.8	.3
.1	.4

B

=

Dot product of A with the first column of B	Dot product of A with the second column of B
---	--

=

.8	.3
----	----

1	0	0	0
0	0	0	1
0	0	1	0

A

.2	.9
.7	0
.8	.3
.1	.4

B

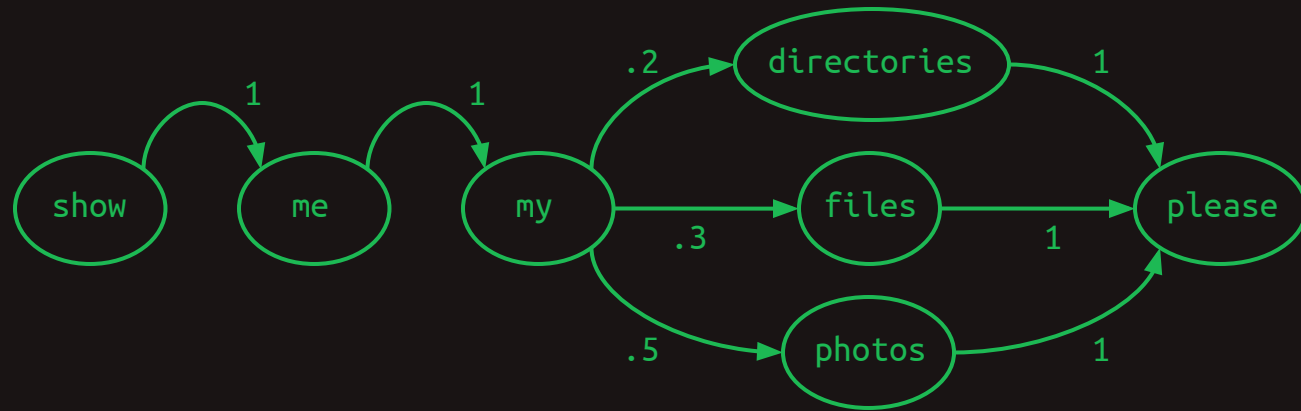
=

Dot products of

row 1 of A, col 1 of B	row 1 of A, col 2 of B
row 2 of A, col 1 of B	row 2 of A, col 2 of B
row 3 of A, col 1 of B	row 3 of A, col 2 of B

=

.2	.9
.1	.4
.8	.3



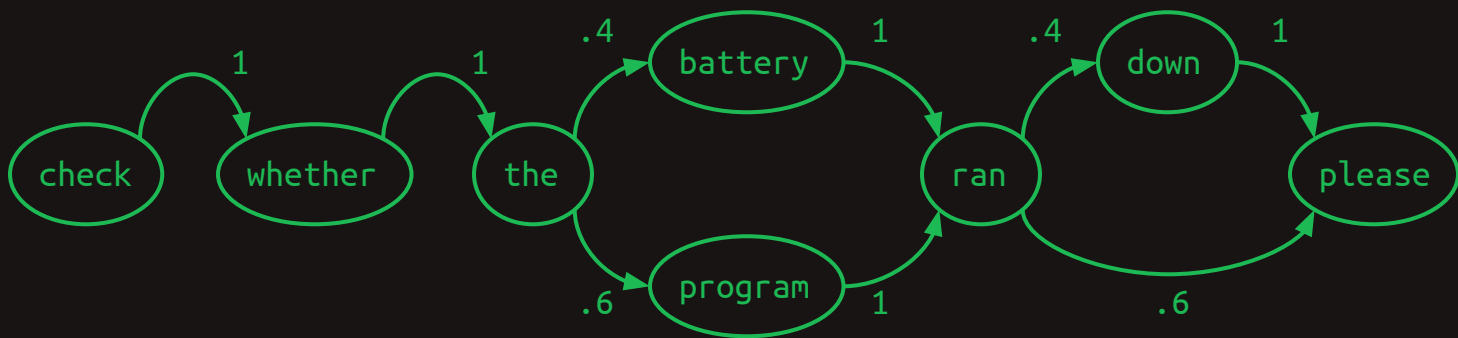
	directories	files	me	my	photos	please	show
directories	0	0	0	0	0	1	0
files	0	0	0	0	0	1	0
me	0	0	0	1	0	0	0
my	.2	.3	0	0	.5	0	0
photos	0	0	0	0	0	1	0
please	0	0	0	0	0	0	0
show	0	0	1	0	0	0	0

directories	files	me	my	photos	please	show
0	0	0	1	0	0	0

directories	files	me	my	photos	please	show
0	0	0	0	0	1	0
0	0	0	0	0	1	0
0	0	0	1	0	0	0
.2	.3	0	0	.5	0	0
0	0	0	0	0	1	0
0	0	0	0	0	0	0
0	0	1	0	0	0	0

=

directories	files	me	my	photos	please	show
.2	.3	0	0	.5	0	0

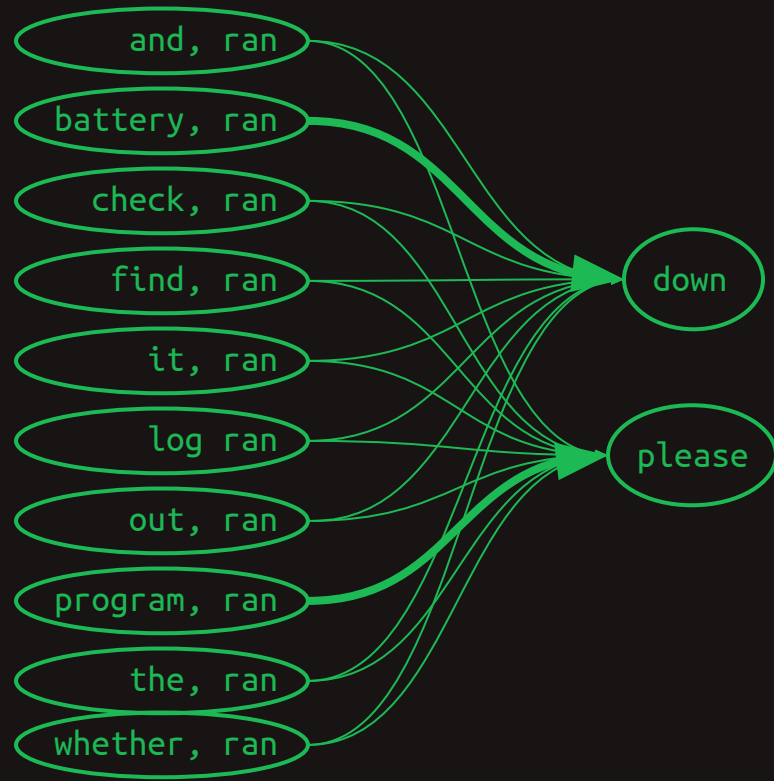


	battery	check	down	please	program	ran	the	whether
battery	0	0	0	0	0	1	0	0
check	0	0	0	0	0	0	0	1
down	0	0	0	1	0	0	0	0
please	0	0	0	0	0	0	0	0
program	0	0	0	0	0	1	0	0
ran	0	0	.4	.6	0	0	0	0
the	.4	0	0	0	.6	0	0	0
whether	0	0	0	0	0	0	1	0

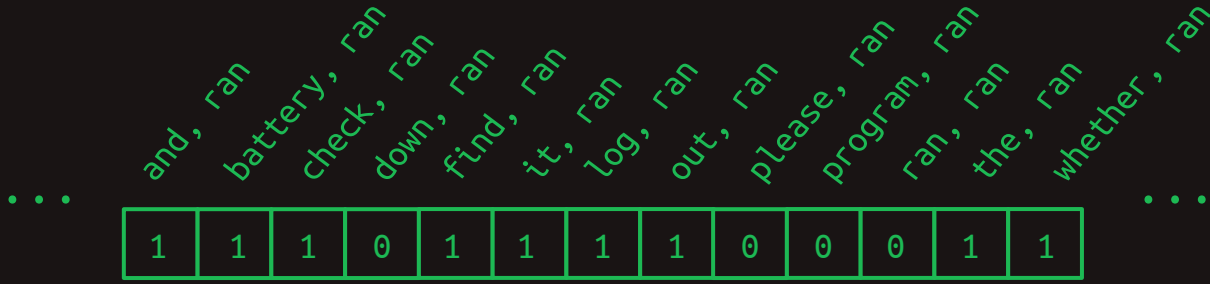




	battery	check	down	please	program	ran	the	whether
battery ran	0	0	1	0	0	0	0	0
check whether	0	0	0	0	0	0	1	0
program ran	0	0	0	1	0	0	0	0
the battery	0	0	0	0	0	1	0	0
the program	0	0	0	0	0	1	0	0
ran down	0	0	0	1	0	0	0	0
whether the	.4	0	0	0	.6	0	0	0
:	0	0	0	0	0	0	0	0
:								
:								



	and	battery	check	down	find	it	log	out	please	program	ran	the	whether
and, ran			.5					.5					
battery, ran			1					0					
check, ran			.5					.5					
down, ran													
find, ran			.5					.5					
it, ran			.5					.5					
log, ran			.5					.5					
out, ran			.5					.5					
please, ran													
program, ran			0					1					
ran, ran													
the, ran			.5					.5					
whether, ran			.5					.5					



feature  
activities

1	1	1	0	1	1	1	1	0	0	0	1	1
---	---	---	---	---	---	---	---	---	---	---	---	---

X

mask

[illegible]

masked  
feature  
activities

[illegible]

[illegible]

all the masks (keys:  $K^T$ )

0	1	0	0	0	0	0	0	0	0	0	1	0	0

one-hot feature vector (query:  $Q$ )

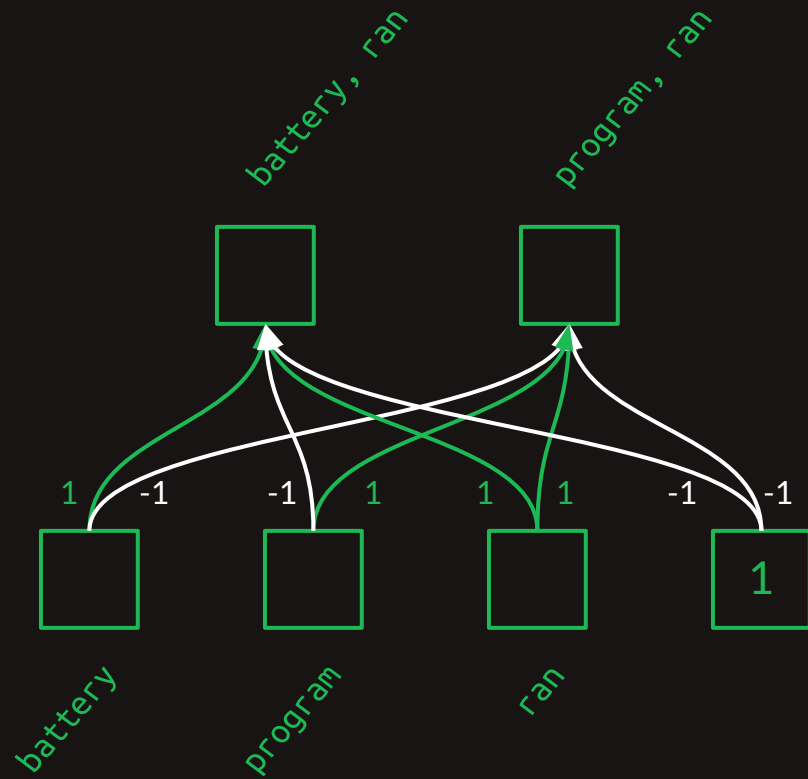
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

feature-specific mask

=	0	1	0	0	0	0	0	0	0	0	1	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



	<i>battery, ran</i>	<i>program, run</i>
<i>battery</i>	1	-1
<i>program</i>	-1	1
<i>run</i>	1	1
<i>bias</i>	-1	-1





The diagram shows the equation  $\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$  on a white background. Four green ovals highlight specific parts of the equation:  $\max(0, \cdot)$ ,  $x$ ,  $W_1 + b_1$ , and  $W_2 + b_2$ . Green lines connect these ovals to text labels: 'ReLU' points to the max function, 'masked word activities' points to  $x$ , 'multi-word feature creation matrix' points to  $W_1 + b_1$ , and 'selective second order transition matrix' points to  $W_2 + b_2$ .

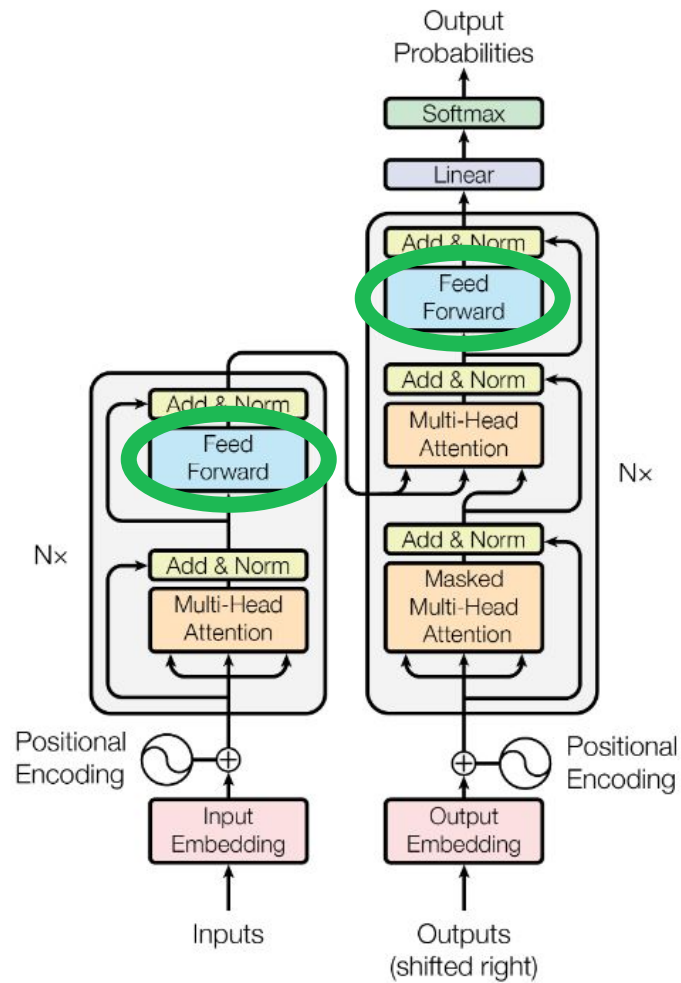
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

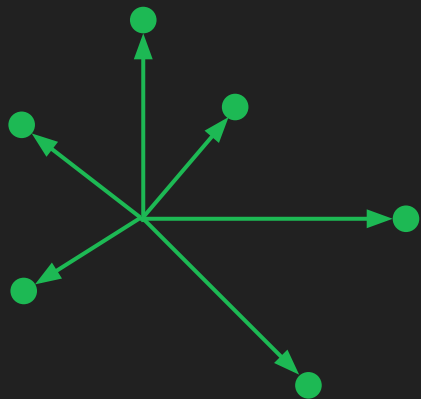
ReLU

masked word activities

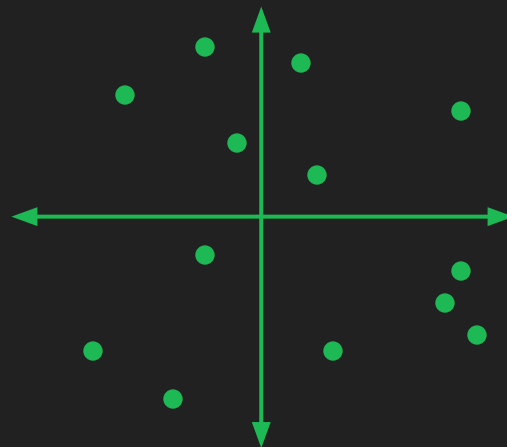
multi-word feature creation matrix

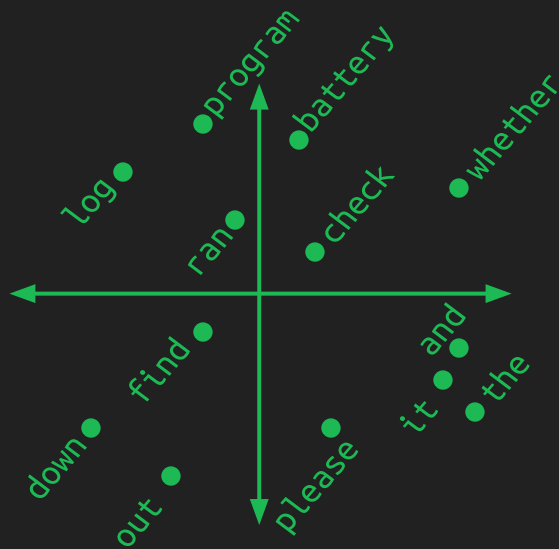
selective second order transition matrix





embedding  
in 2 dimensions





.2	1.0
----	-----

 battery

-.3	1.1
-----	-----

 program

-.9	-.8
-----	-----

 down

-.5	-1.2
-----	------

 out

⋮





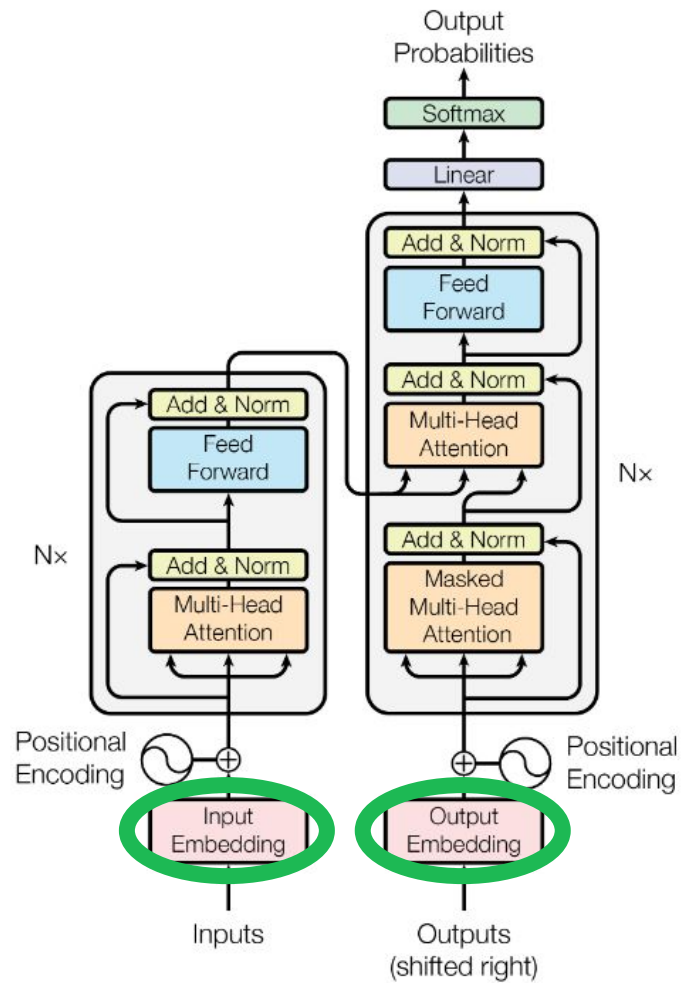
embedding  
projection  
matrix

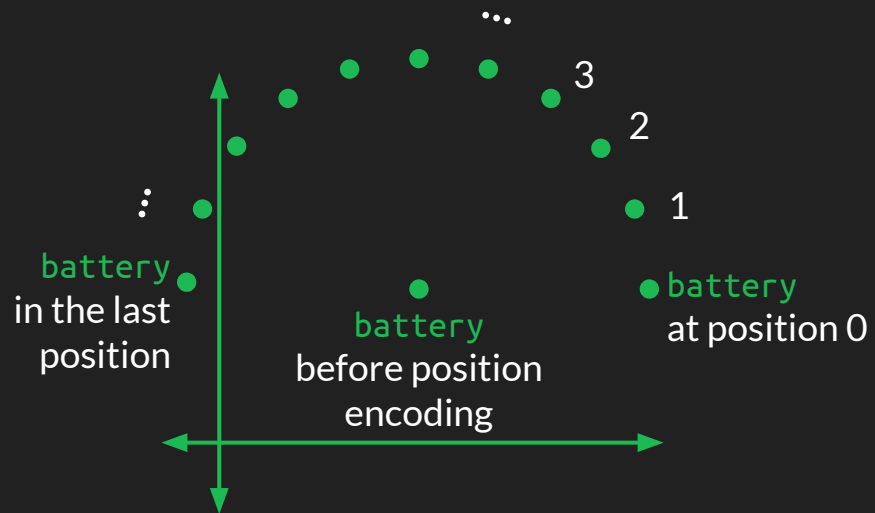
.2	1.0

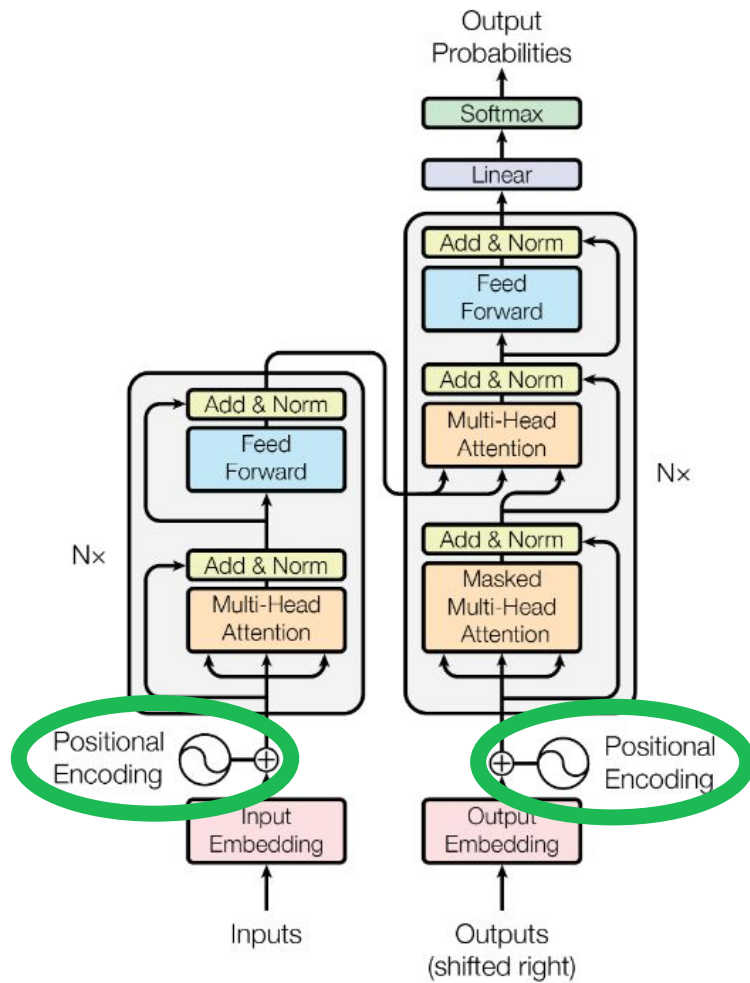
=

battery

.2	1.0
----	-----







program

-.3	1.1
-----	-----

									-.1			
									.9			

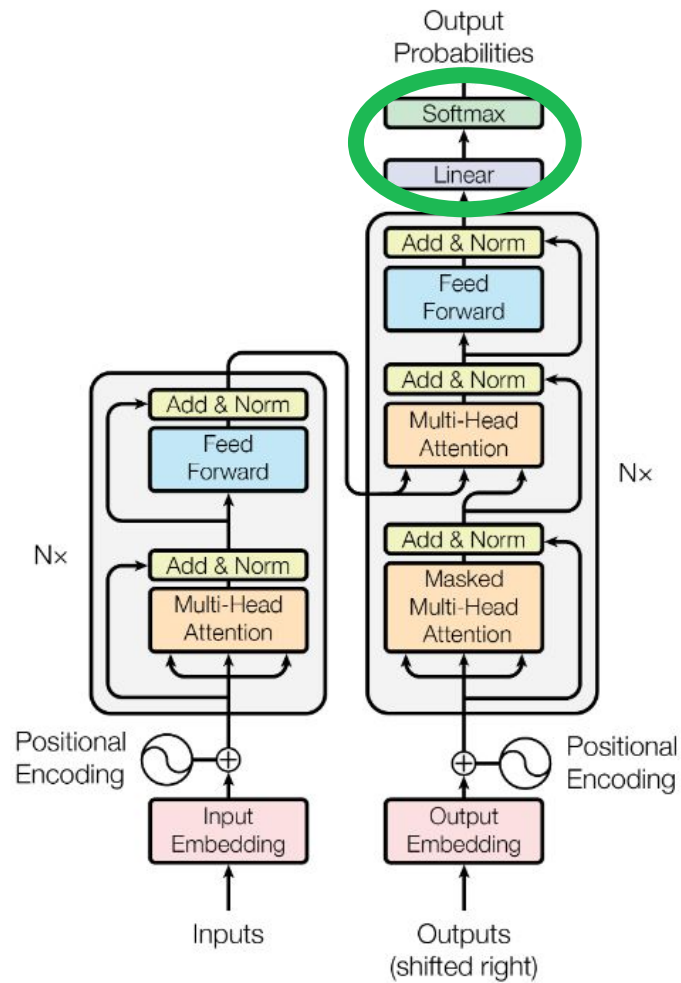
de-embedding  
projection matrix

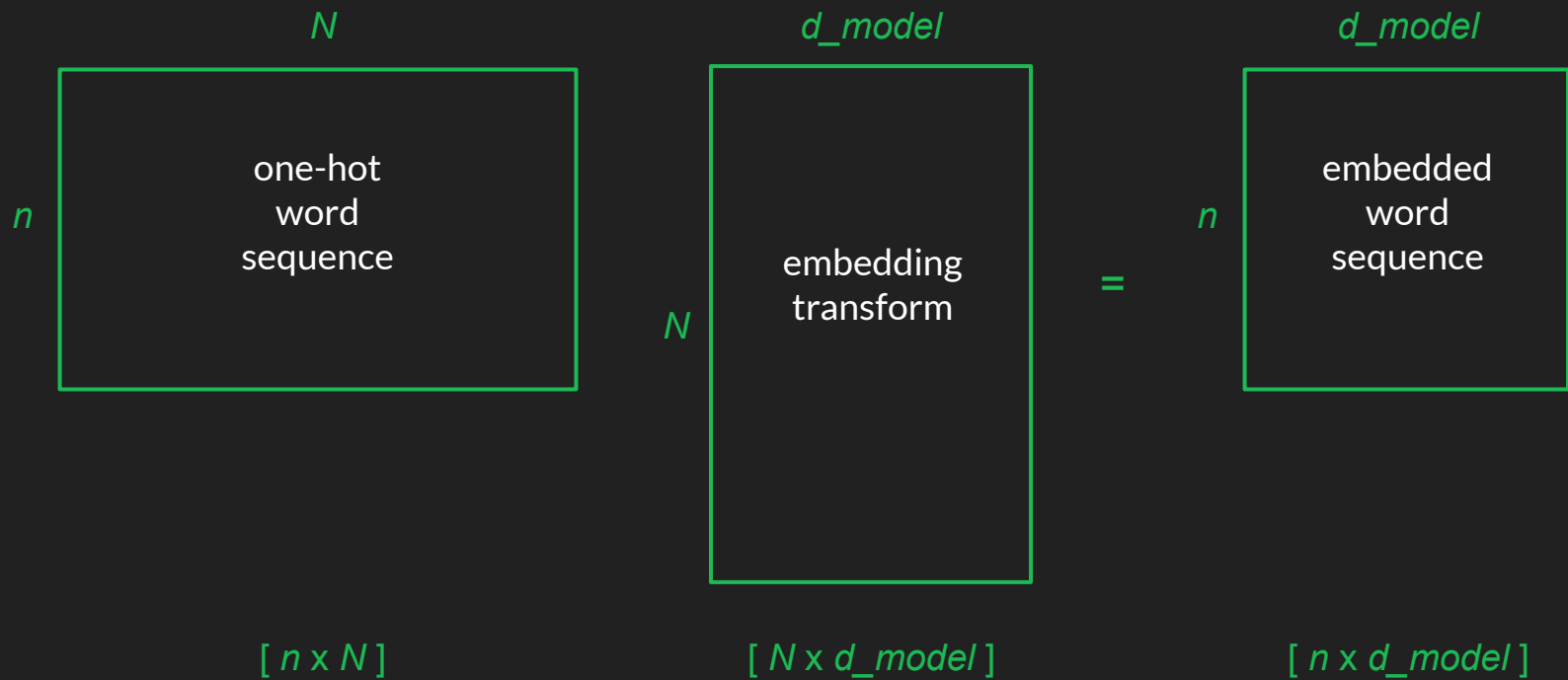
=

									1			
--	--	--	--	--	--	--	--	--	---	--	--	--

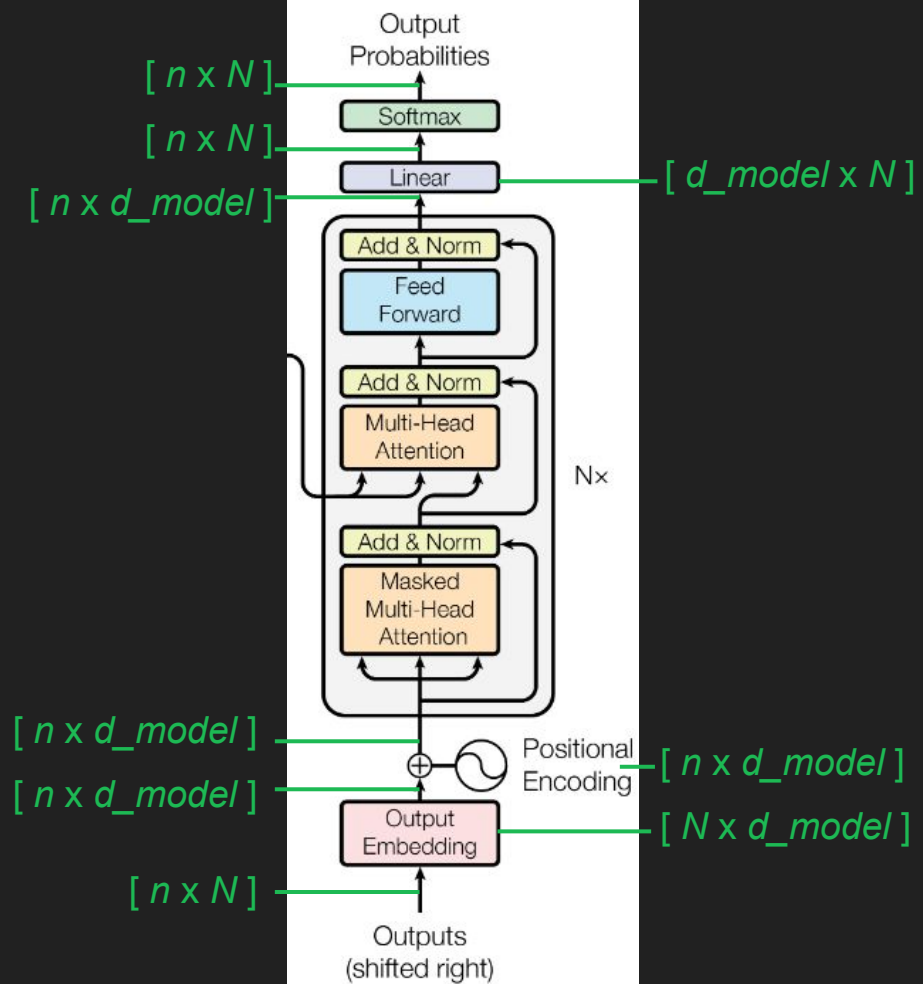
and battery check down find it log out please program ran the whether

and	battery	check	down	find	it	log	out	please	program	ran	the	whether
-.7	.9	.1	-.2	0	-.7	.8	-.3	-.6	1	.4	-.8	0

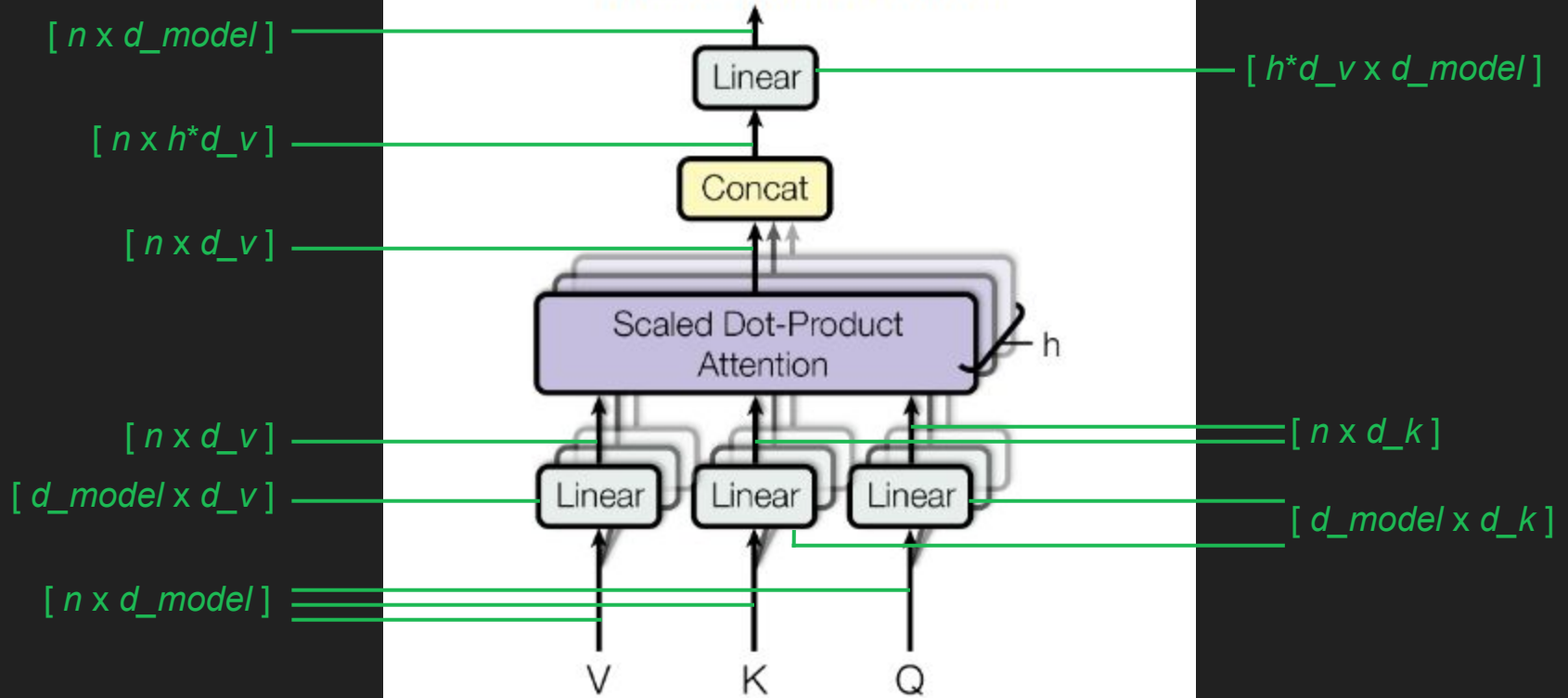




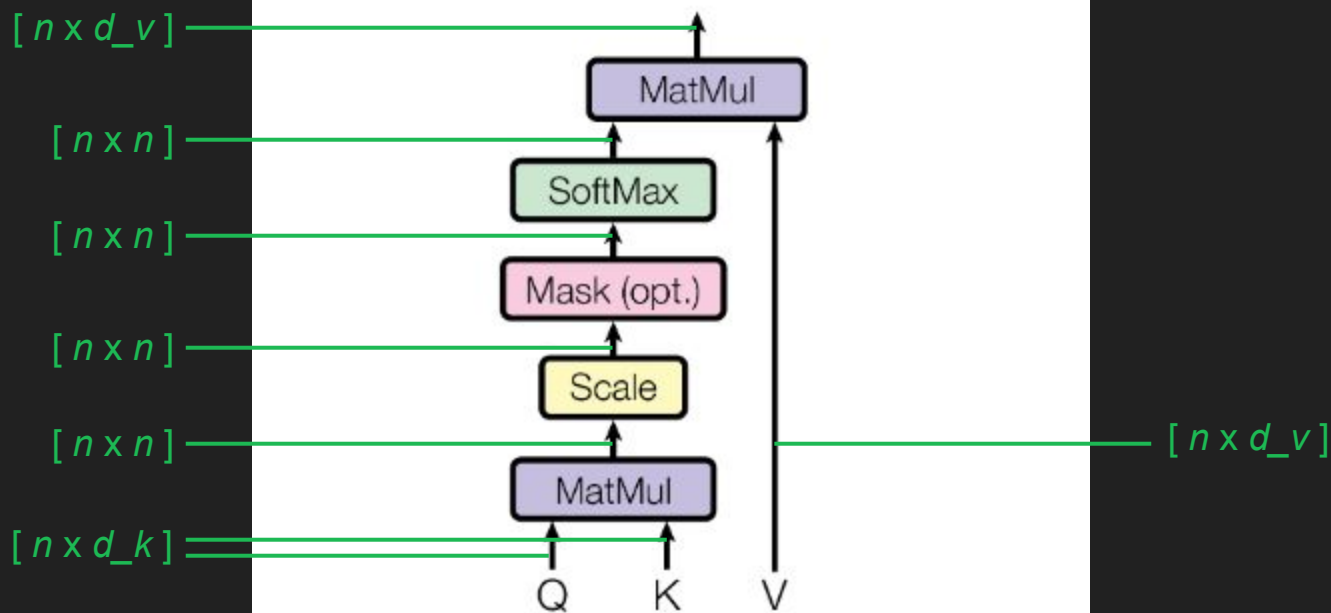


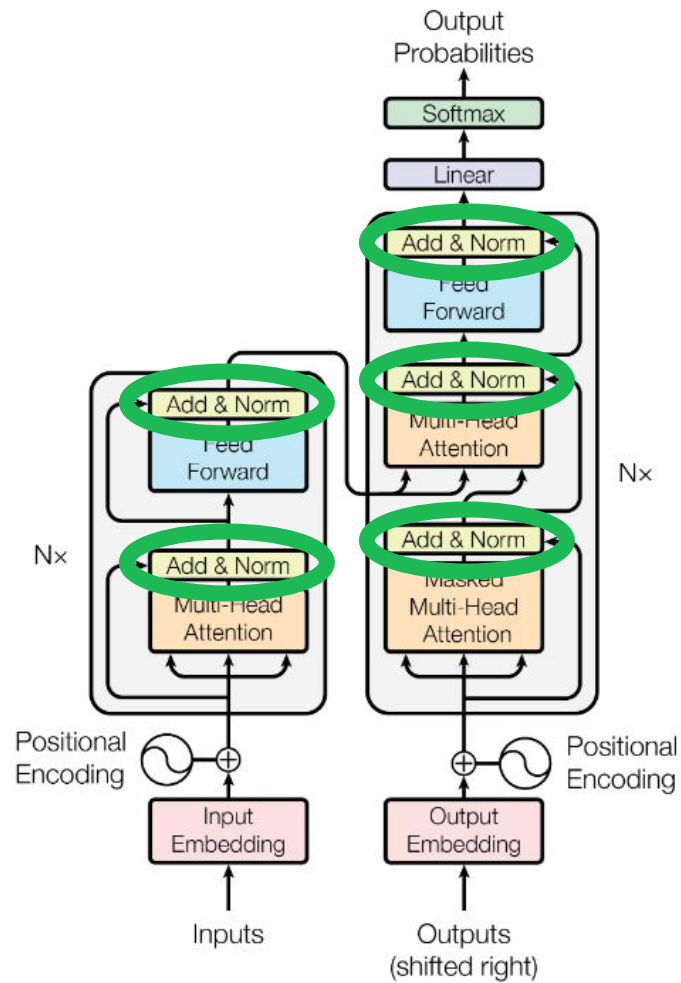


## Multi-Head Attention



## Scaled Dot-Product Attention





Layer 1

Layer 2

...

