



# Why Transformers Work.



I'm Vincent  
Ask Me Anything [tm]

Think about all  
the actions  
in this →  
dialogue

😊 Hello

Hi 🤖

😊 I'd like a pizza!

What kind? 🤖

😊 By the way,  
are you a human?

No, I'm a bot. 🤖

But what kind of pizza? 🤖

intent



Hello

entity



Hi



intent




I'd like a pizza!


What kind? 

intent



By the way,  
are you a human?

No, I'm a bot. 

But what kind of pizza? 

intent



Hello

entity



Hi



action

intent



I'd like a pizza!

What kind?



action

intent



By the way,  
are you a human?

No, I'm a bot.



action

But what kind of pizza?



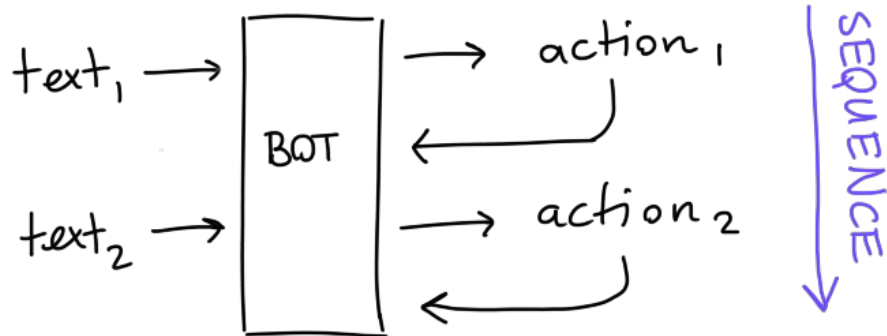
action

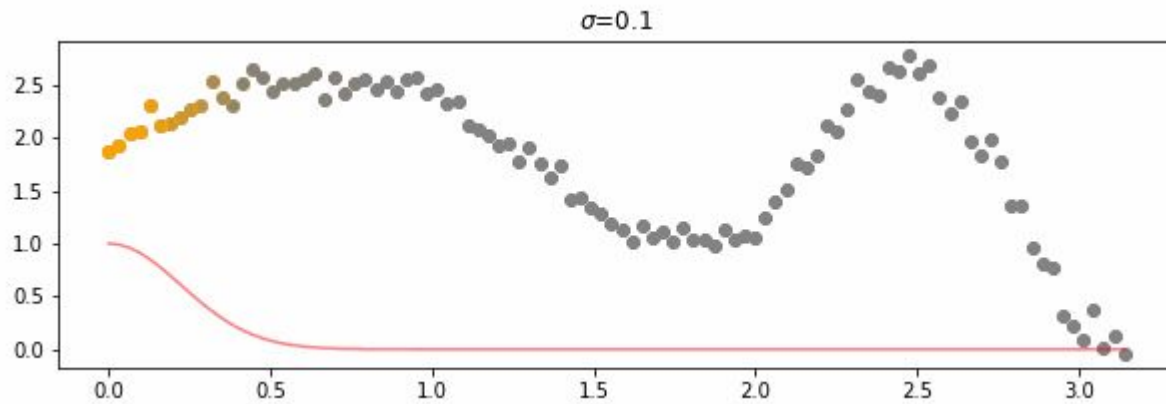
I wanna buy a pizza.

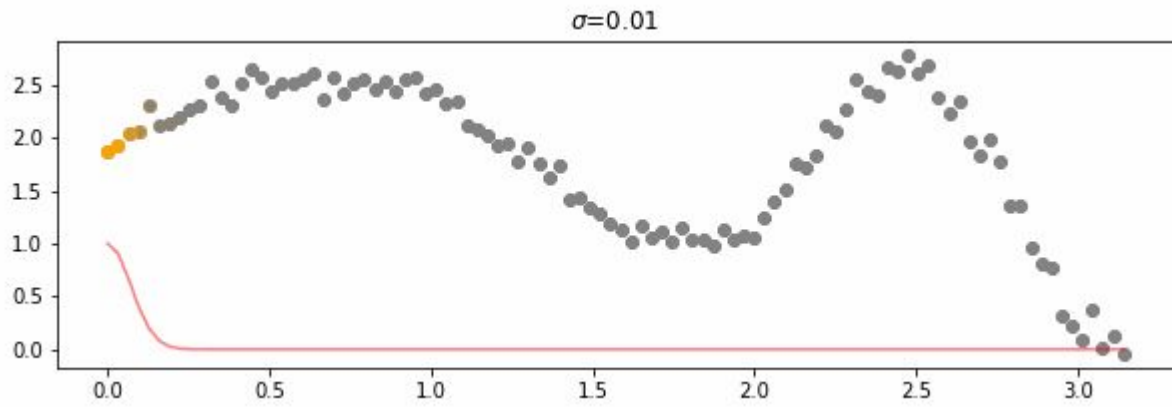
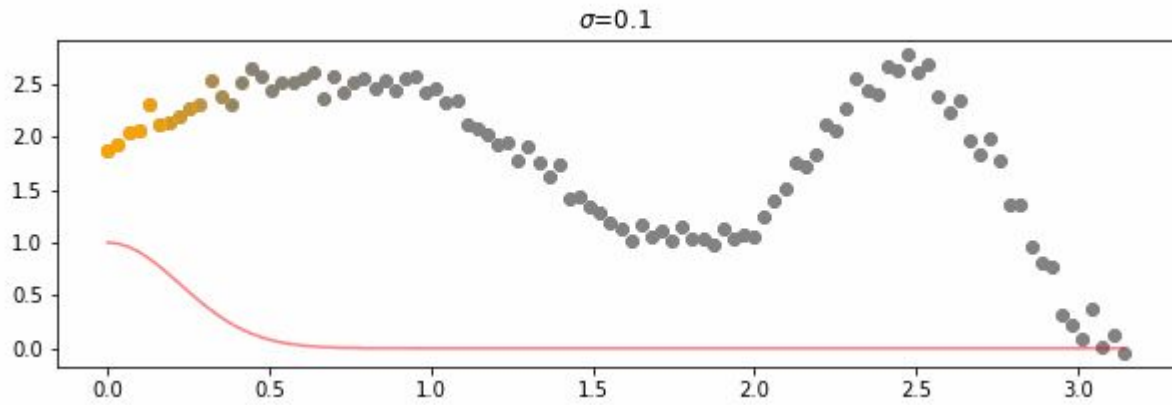
→  
SEQUENCE

I wanna buy a pizza.

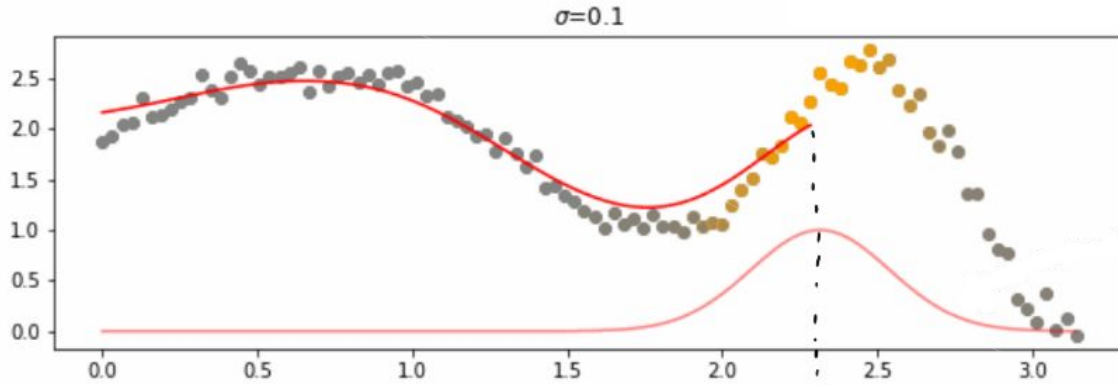
SEQUENCE



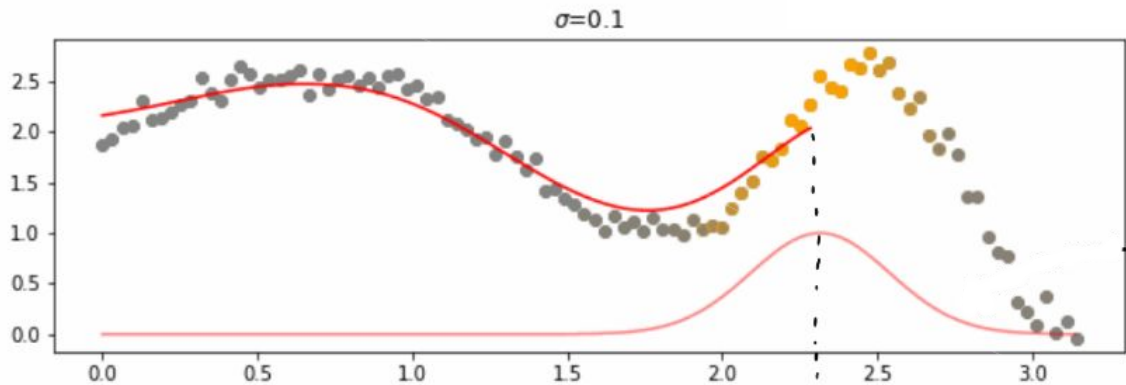








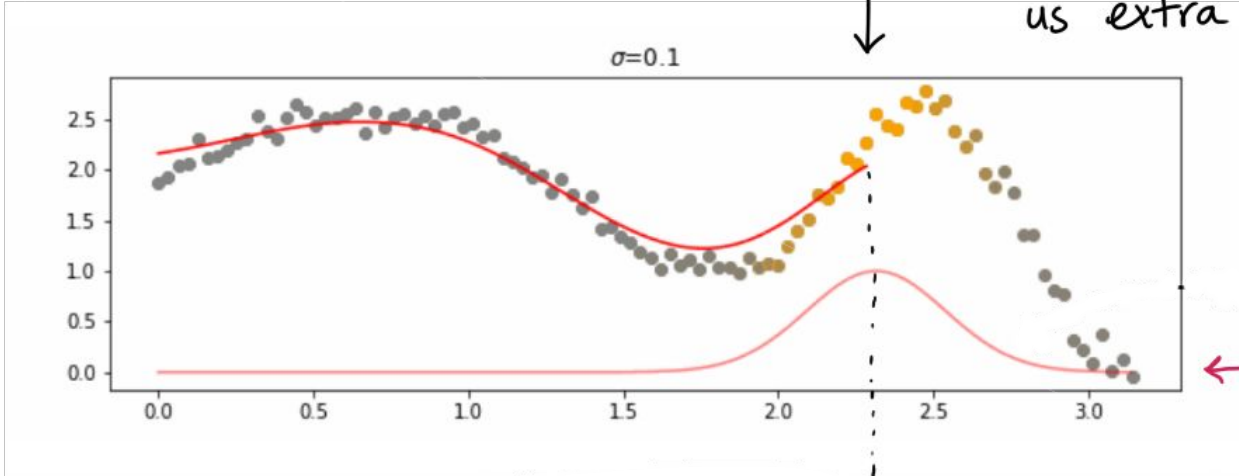
1. as far as  $x_i$  is concerned



1. as far as  $x_i$  is concerned

2. this is what we need to pay attention to

3. this allows us to re-weight and this operation hopefully gives us extra context



1. as far as  $x_i$  is concerned

2. this is what we need to pay attention to

Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$V_1$   $V_2$   $V_3$   $V_4$



Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

$V_1$   $V_2$   $V_3$   $V_4$



Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

If  $v_1$  and  $v_2$  share info

$v_1 \cdot v_2$  is big

If  $v_1$  and  $v_2$  don't

$v_1 \cdot v_2 \approx 0$

Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

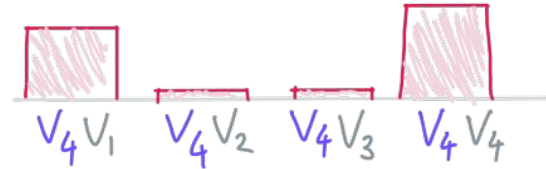
$v_1$   $v_2$   $v_3$   $v_4$



Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

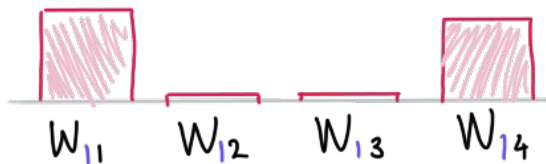
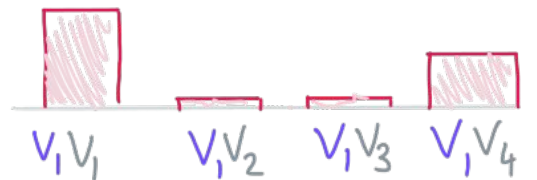
$v_1$   $v_2$   $v_3$   $v_4$



Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

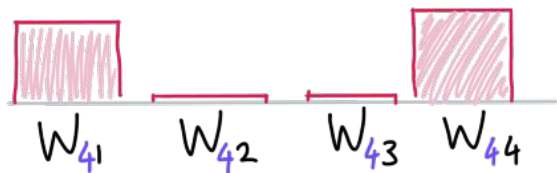
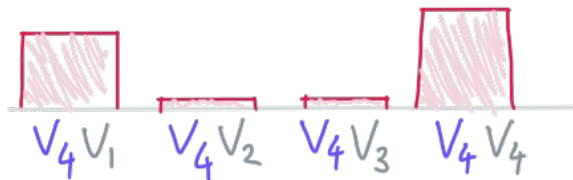
$v_1$   $v_2$   $v_3$   $v_4$



Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

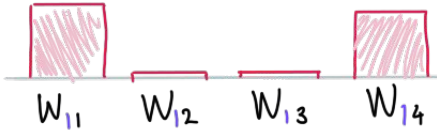
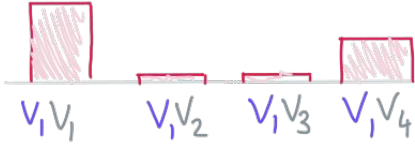


NORMALISE

Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

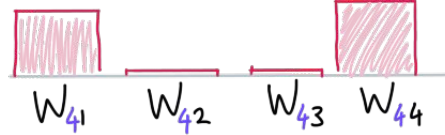
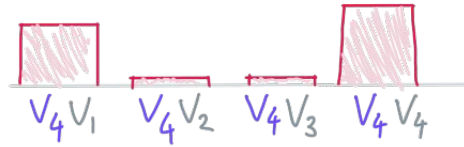


$$v_1^* = \sum_j w_{1j} v_j$$

Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$



$$v_4^* = \sum_j w_{4j} v_j$$

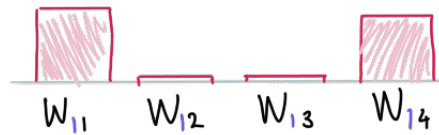
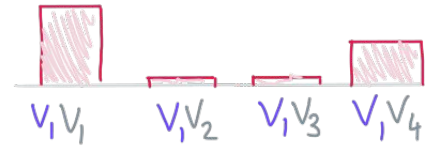
NORMALISE



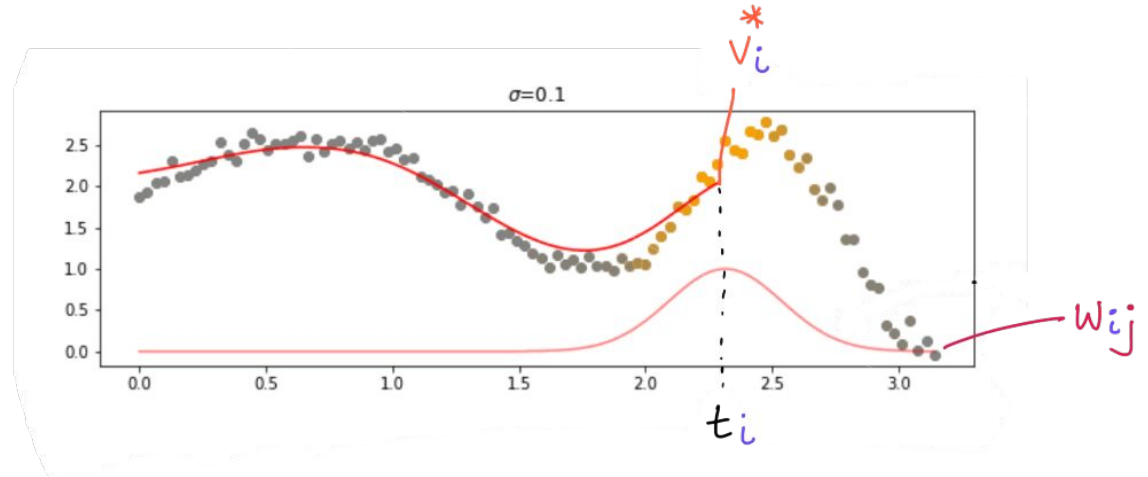
Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$



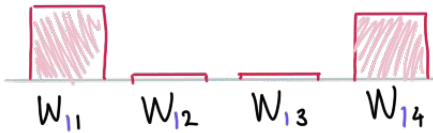
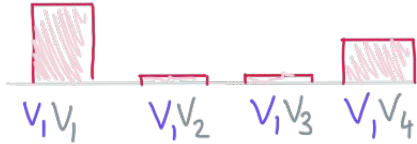
$$v_i^* = \sum_j w_{ij} v_j$$



Bank of the river.

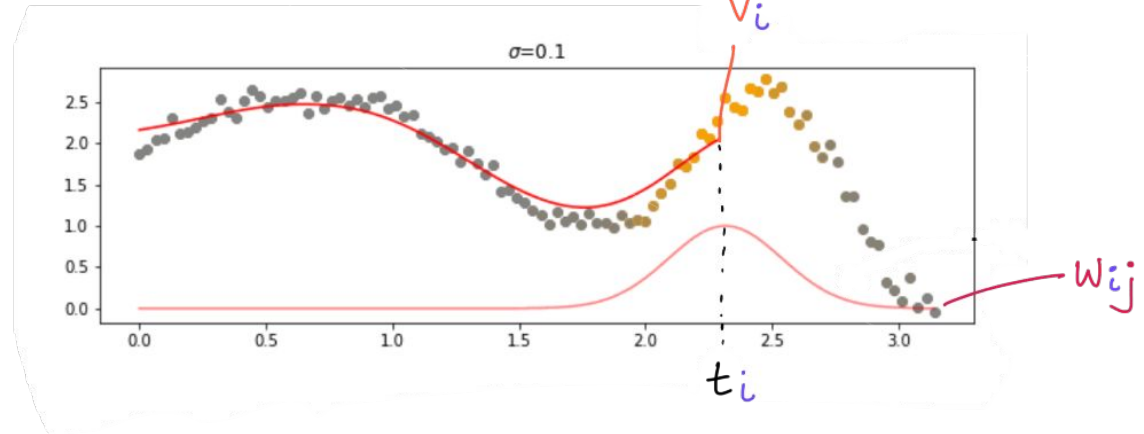
$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$



$$v_i^* = \sum_j w_{ij} v_j$$

reweighing based on  
time-distance



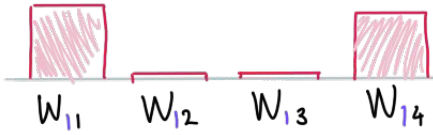
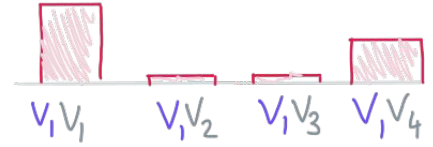
Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

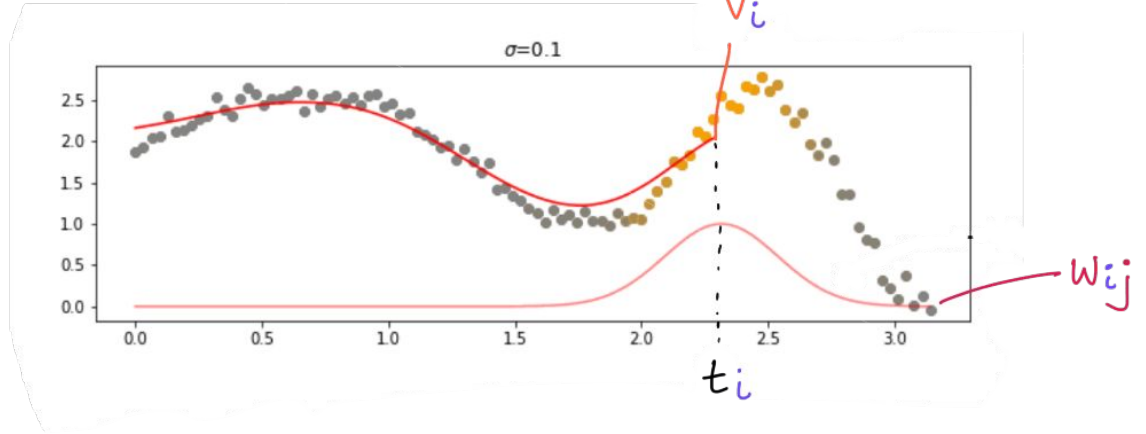
$v_1$   $v_2$   $v_3$   $v_4$

reweighing based on  
embedding similarity

reweighing based on  
time-distance



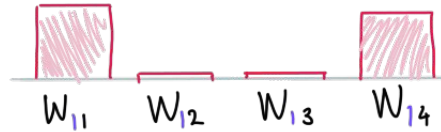
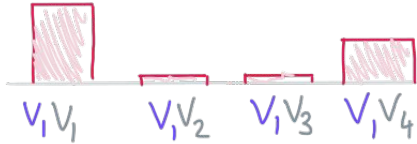
$$v_i^* = \sum_j w_{ij} v_j$$



Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

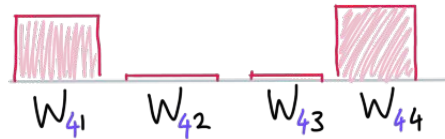
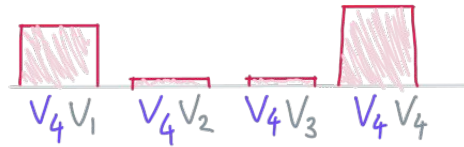


$$v_1^* = \sum_j w_{1j} v_j$$

Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

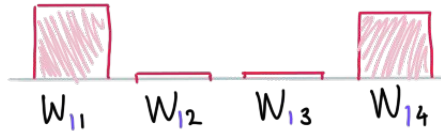
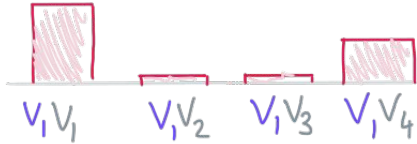


$$v_4^* = \sum_j w_{4j} v_j$$

Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$

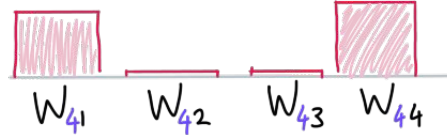
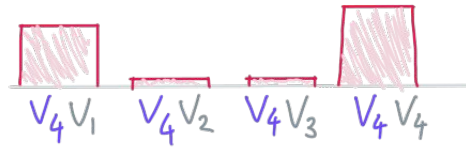


$$v_1^* = \sum_j w_{1j} v_j$$

Money on the bank.

$t_1$   $t_2$   $t_3$   $t_4$

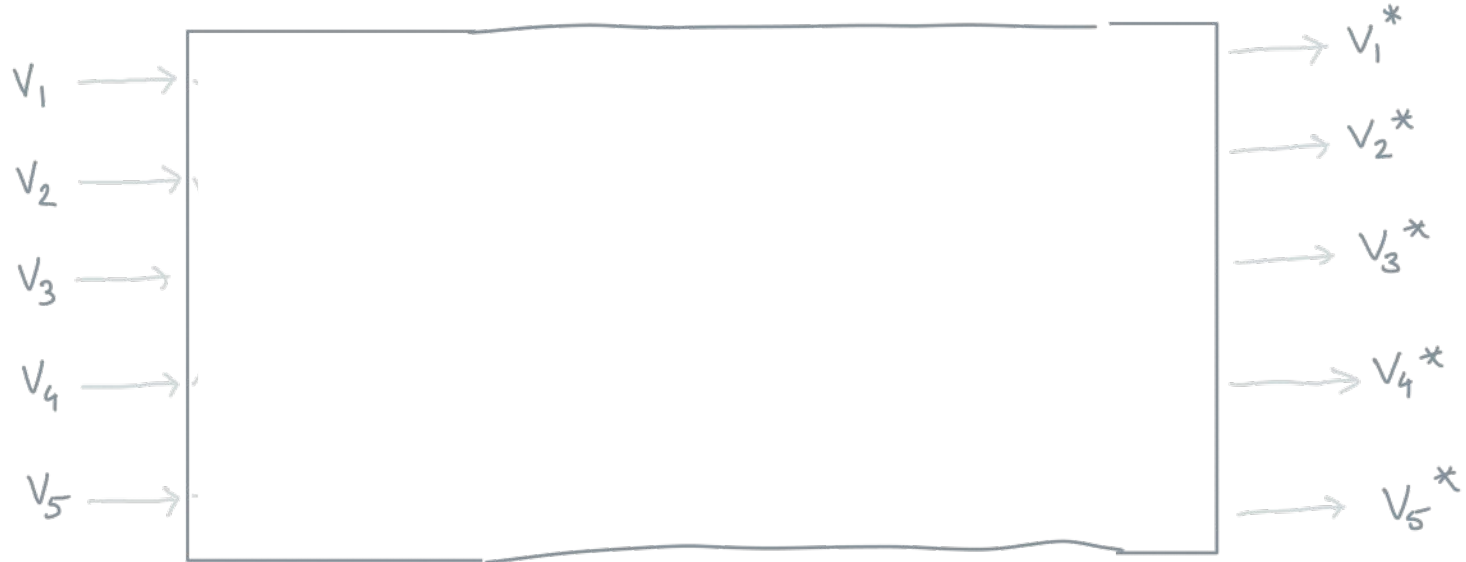
$v_1$   $v_2$   $v_3$   $v_4$



$$v_4^* = \sum_j w_{4j} v_j$$

Let's move this idea into a layer.

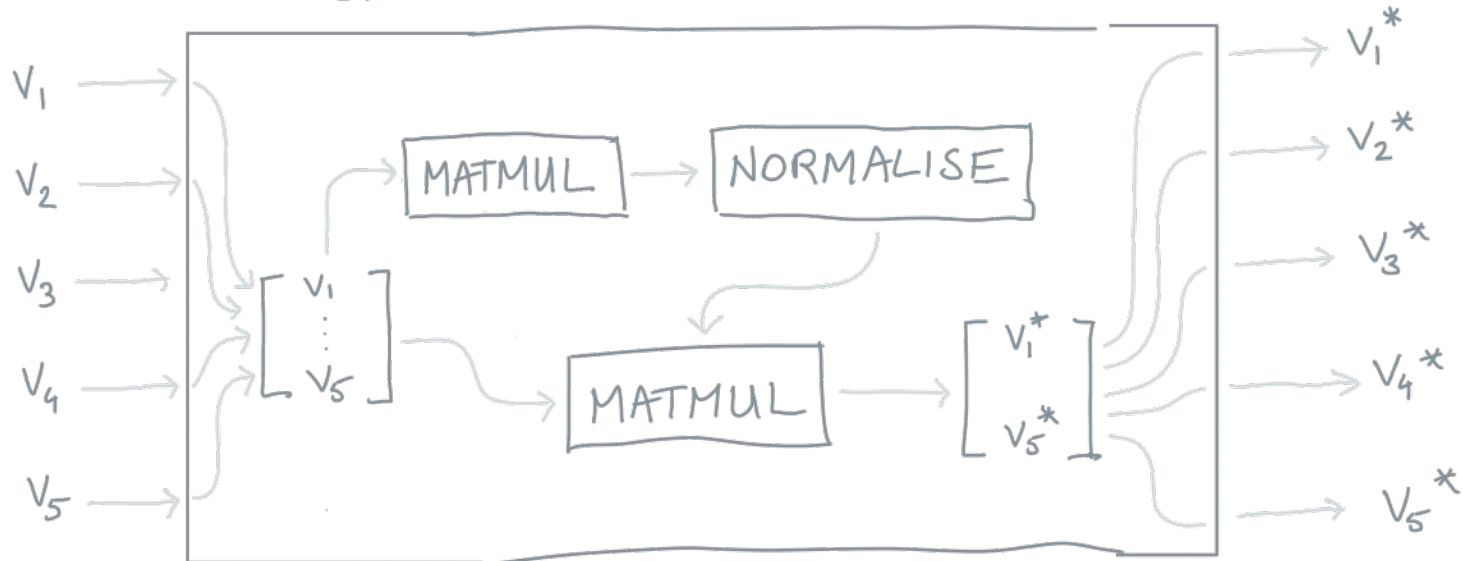
# SELF ATTENTION BLOCK



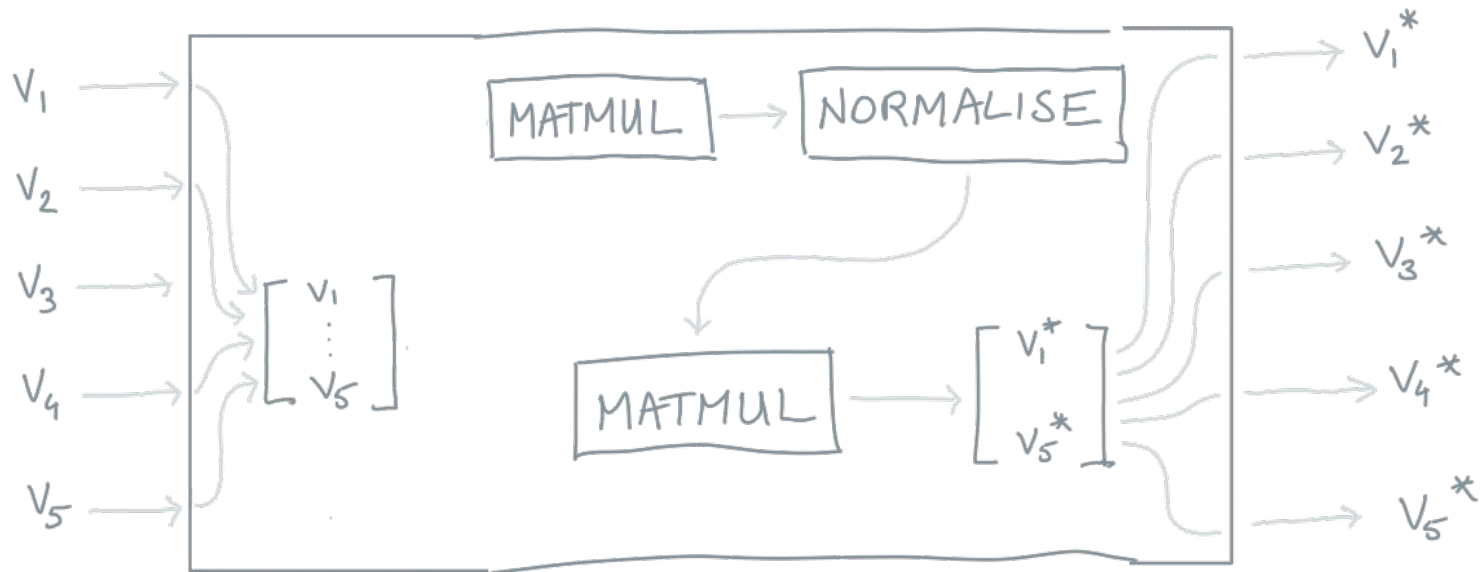
# SELF ATTENTION BLOCK

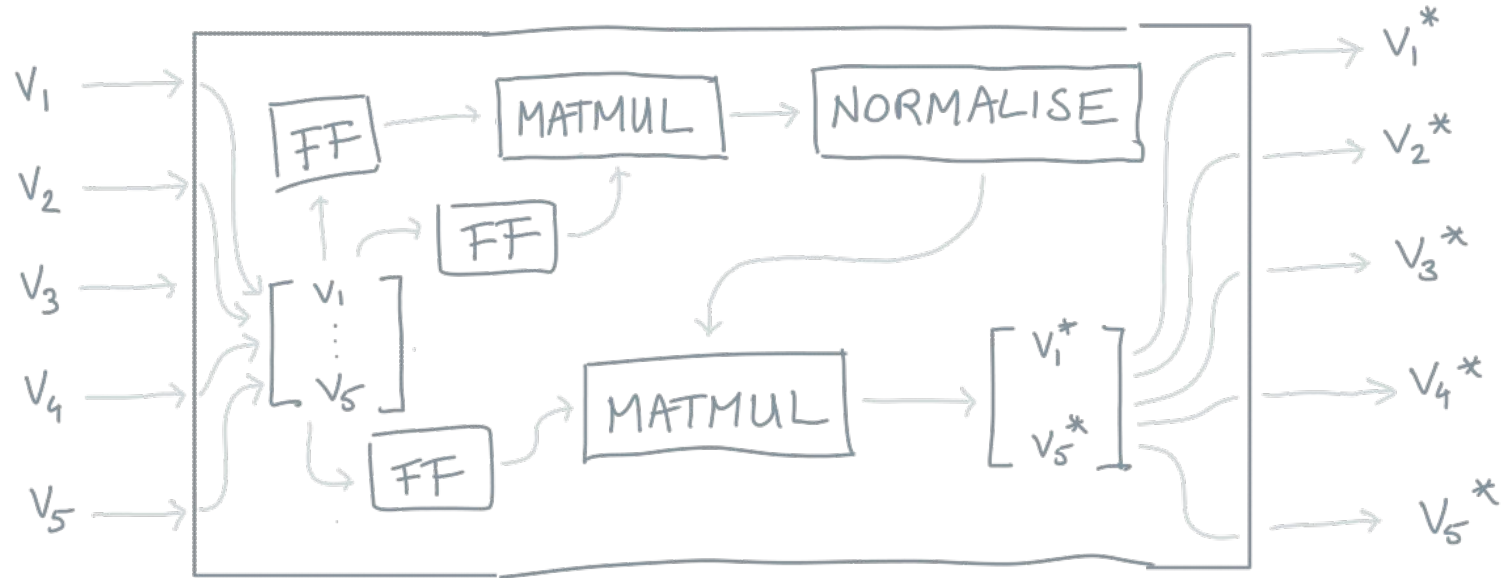


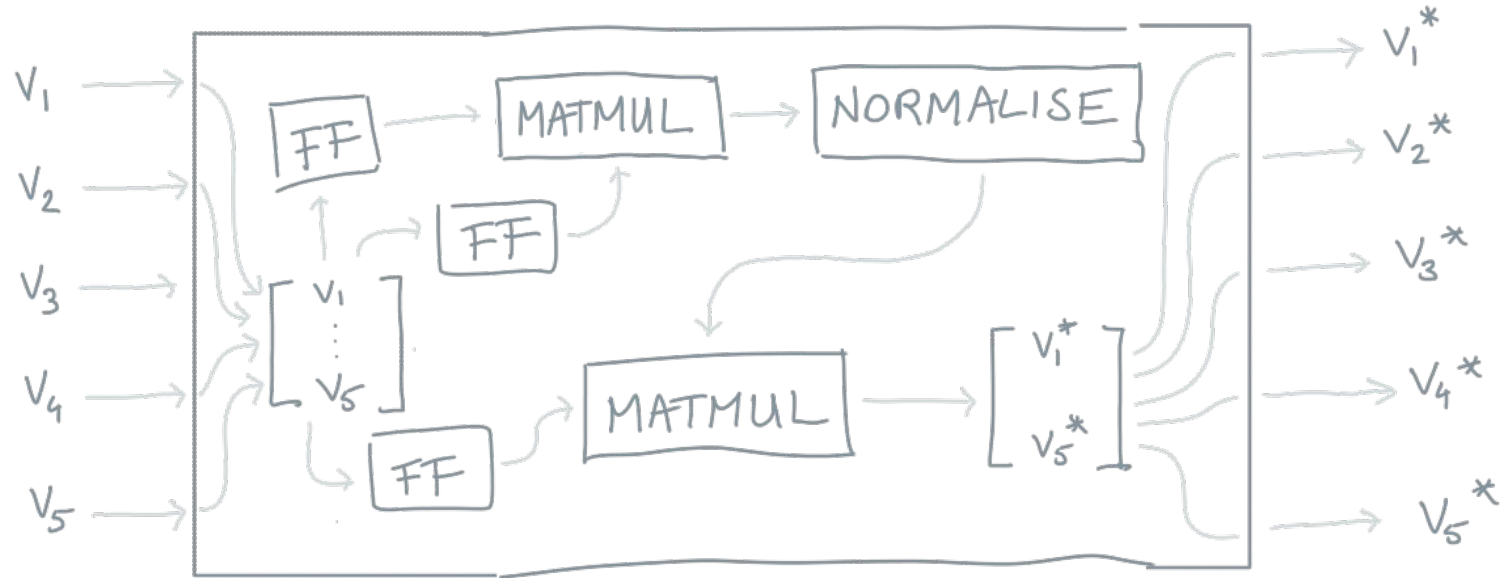
# SELF ATTENTION BLOCK



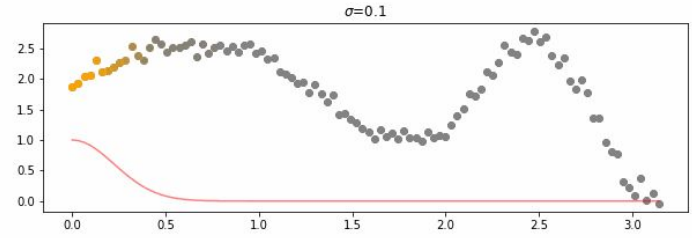








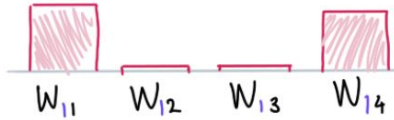
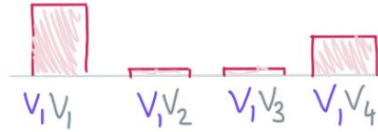
It's a more elaborate way to do stuff like this →



Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

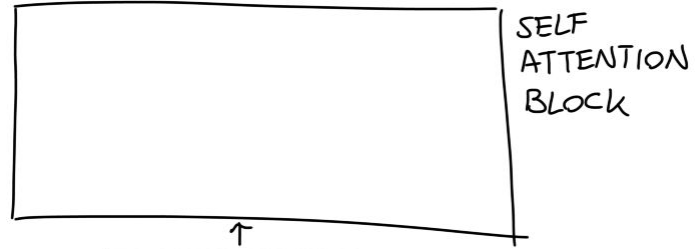
$v_1$   $v_2$   $v_3$   $v_4$



$$v_i^* = \sum_j w_{ij} v_j$$

$$\begin{bmatrix} v_1^* \\ \vdots \\ v_n^* \end{bmatrix}$$

↑

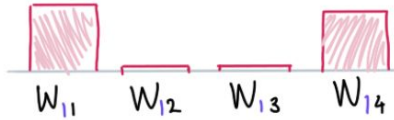


$$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

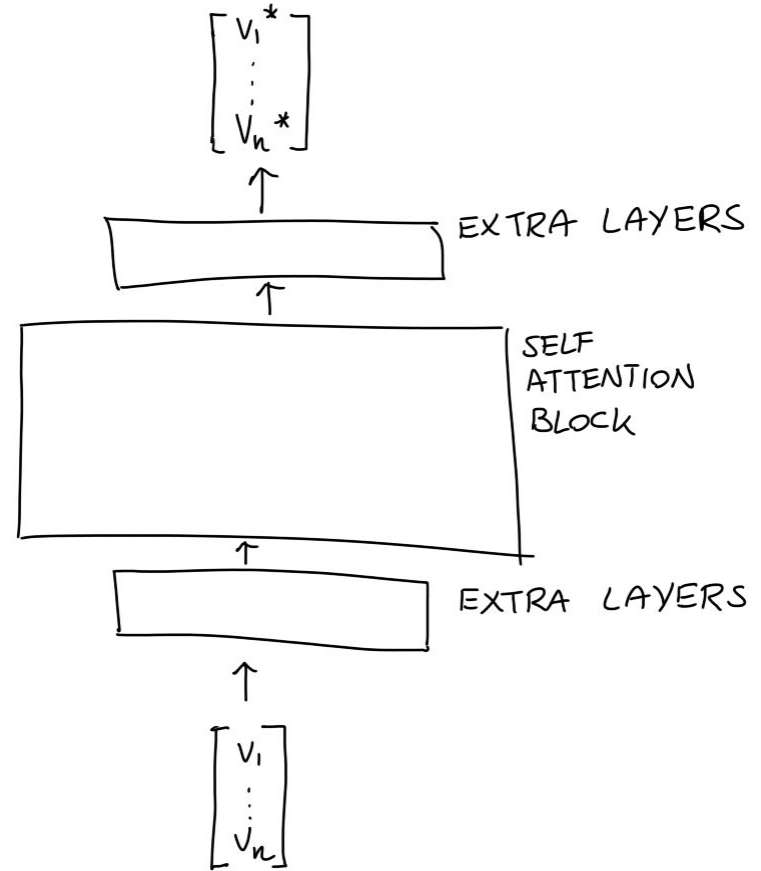
Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

$v_1$   $v_2$   $v_3$   $v_4$



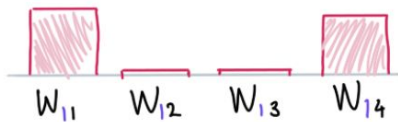
$$v_1^* = \sum_j w_{1j} v_j$$



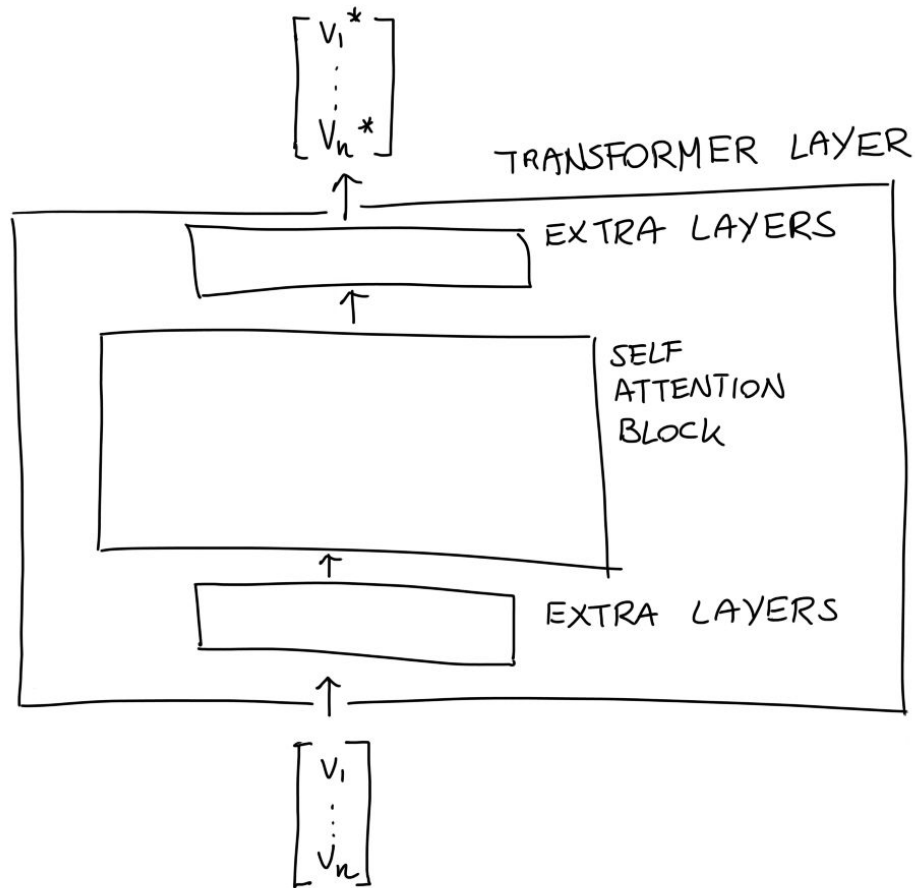
Bank of the river.

$t_1$   $t_2$   $t_3$   $t_4$

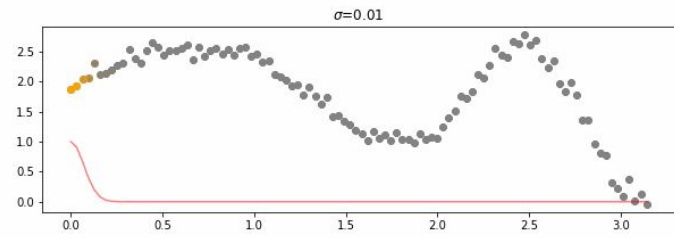
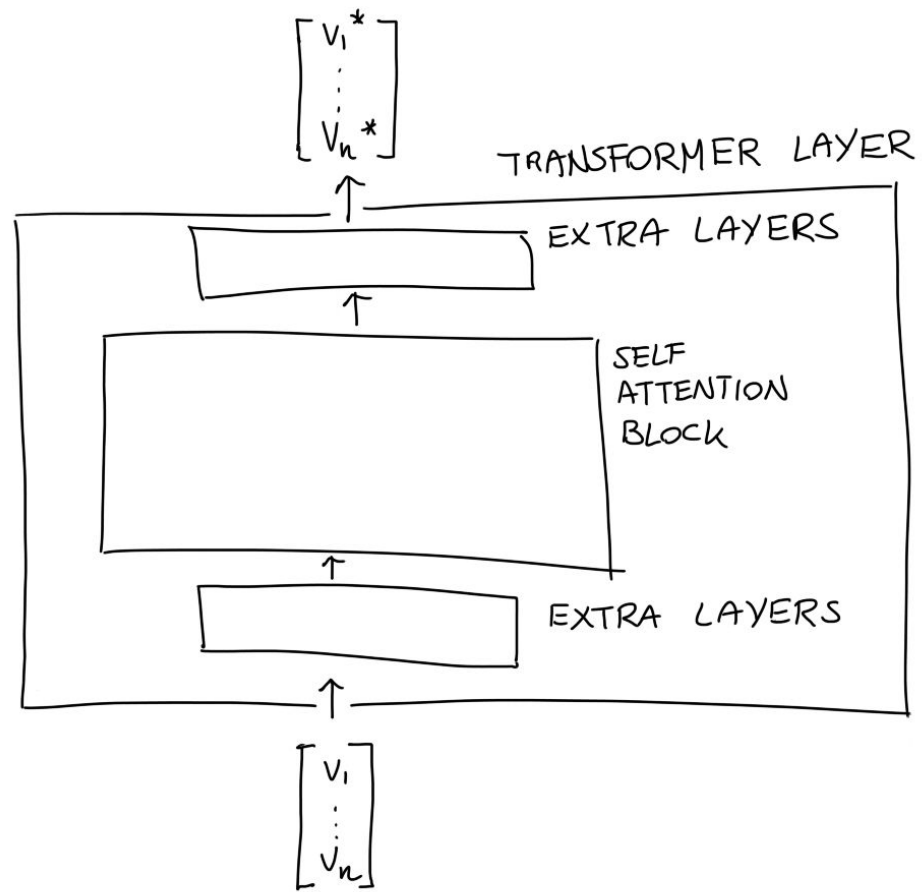
$v_1$   $v_2$   $v_3$   $v_4$

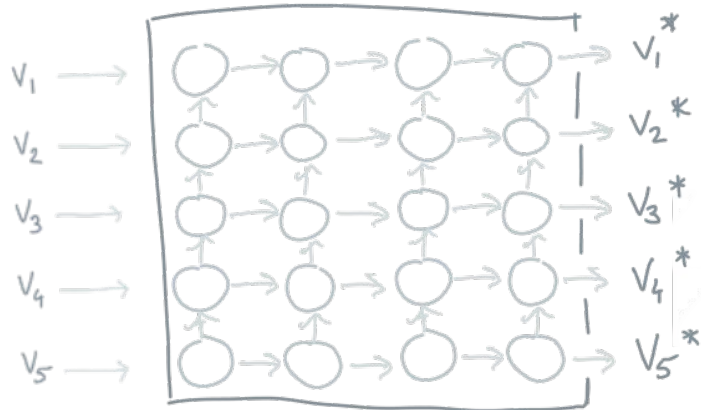


$$v_i^* = \sum_j w_{ij} v_j$$

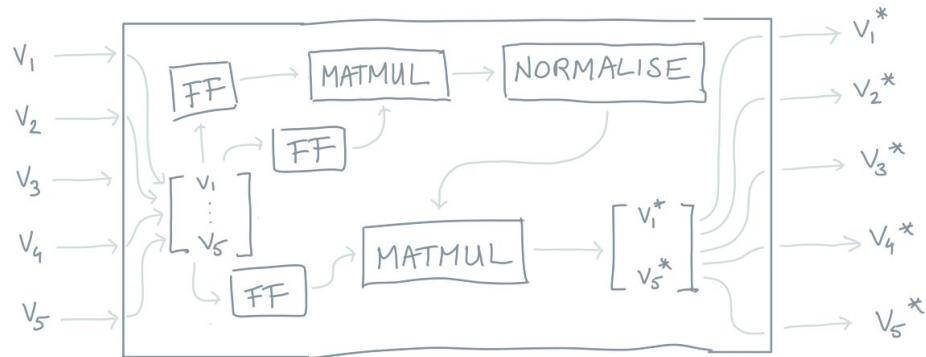


\*glancing over some details  
the goal here is intuition  
in the interest of time

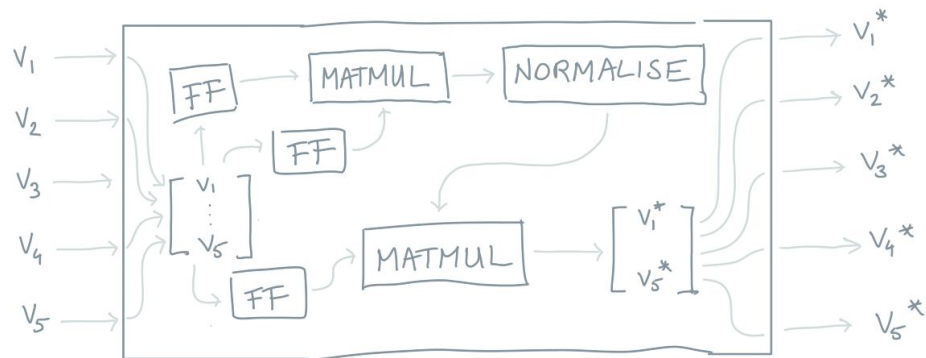
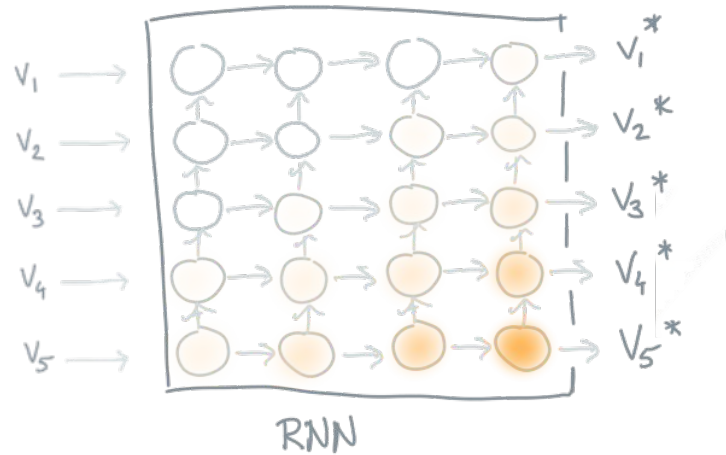


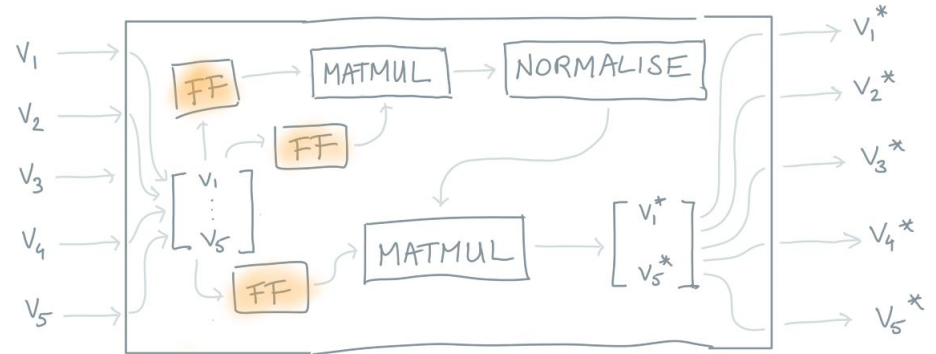
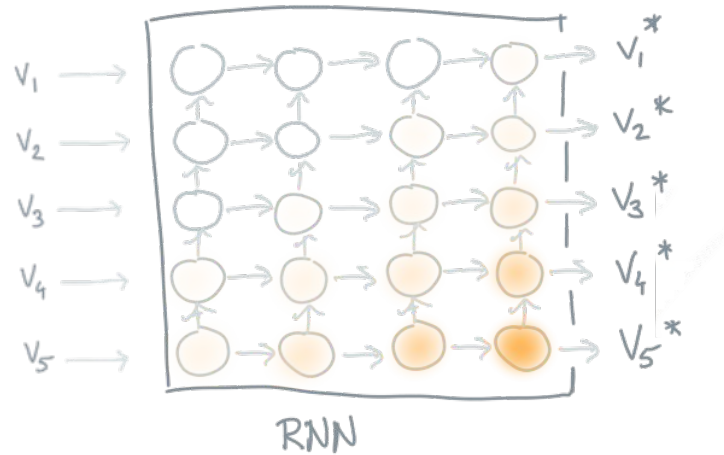


RNN

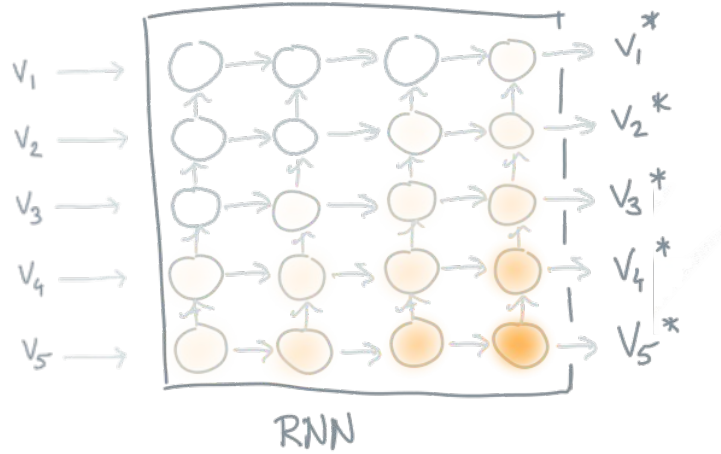




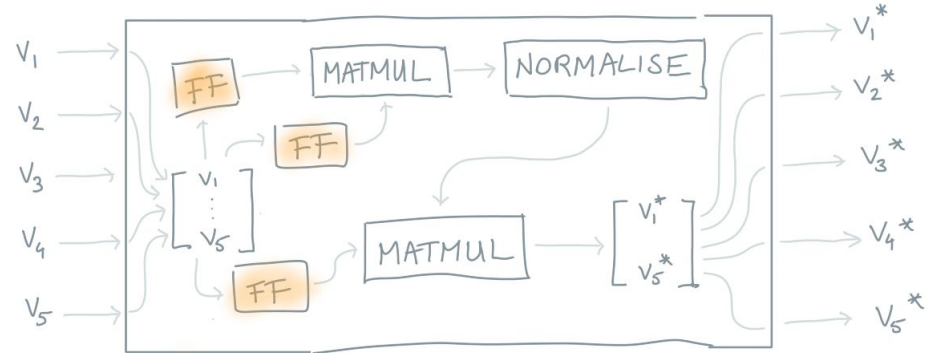




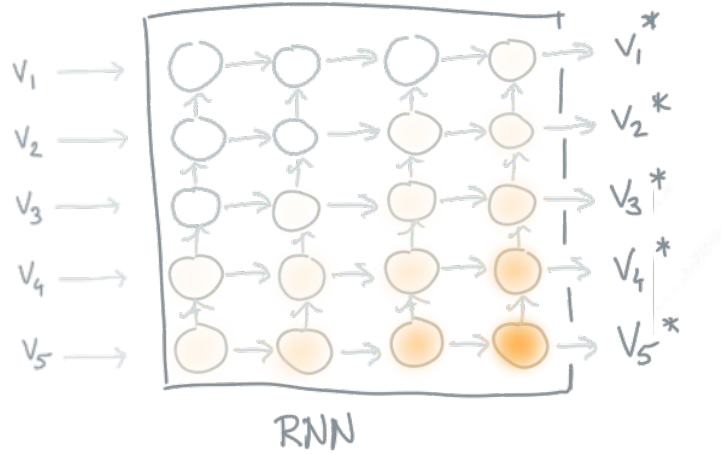
Needs a lot of data  
to stop looking at  
its previous neighbors.



Starts out by looking  
at similarity of  
pretrained embeddings.

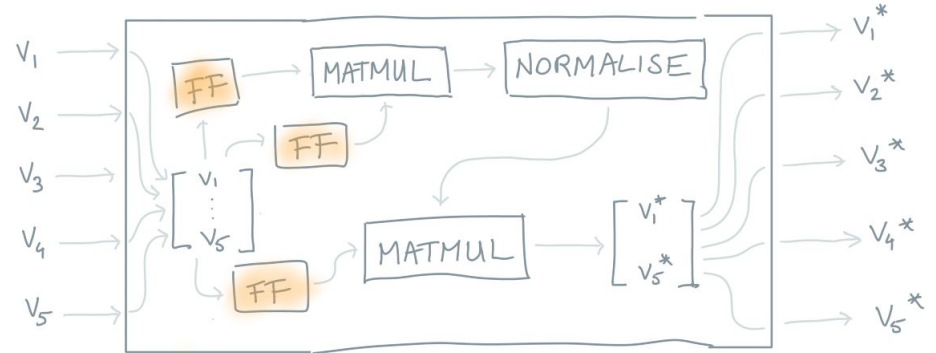


Needs a lot of data  
to stop looking at  
its previous neighbors.



Hard to parallelise.

Starts out by looking  
at similarity of  
pretrained embeddings.



More parallel options.

intent



Hello

entity



Hi



action

intent



I'd like a pizza!

What kind? 

action

intent



By the way,  
are you a human?

No, I'm a bot.



action

But what kind of pizza?



action

Q: So how does Rasa find these intents and entities?

A: DIET

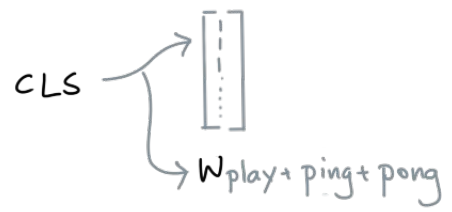
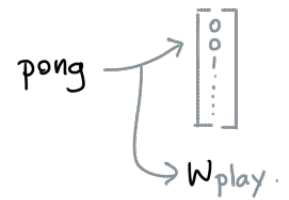
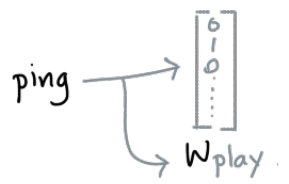
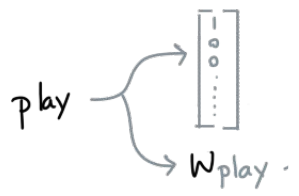
Dual Intent & Entity Transformer

play -

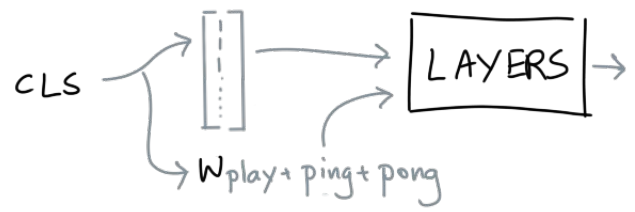
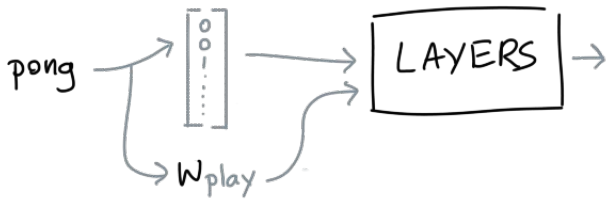
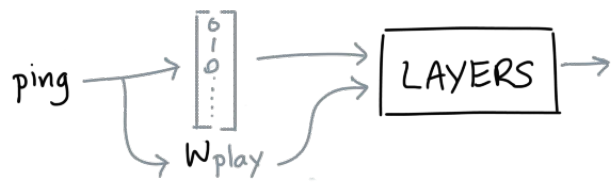
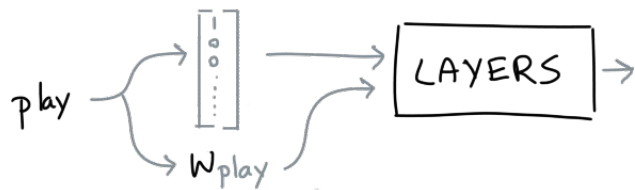
ping -

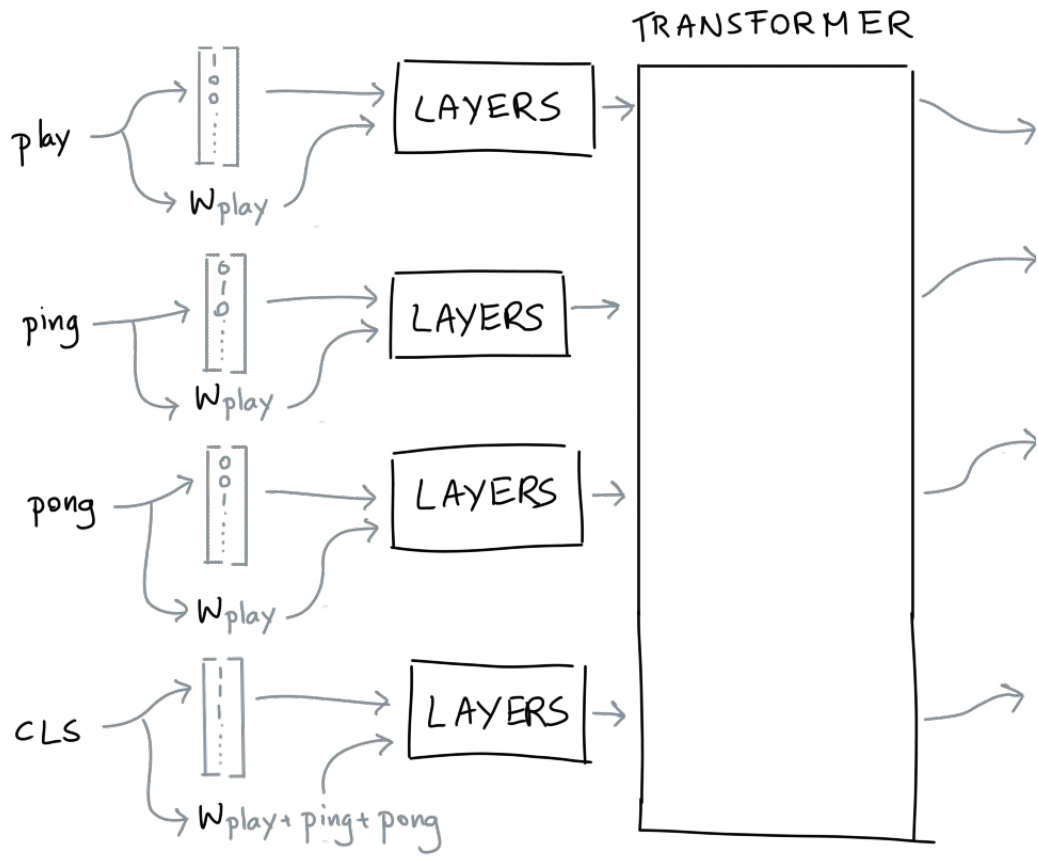
pong -

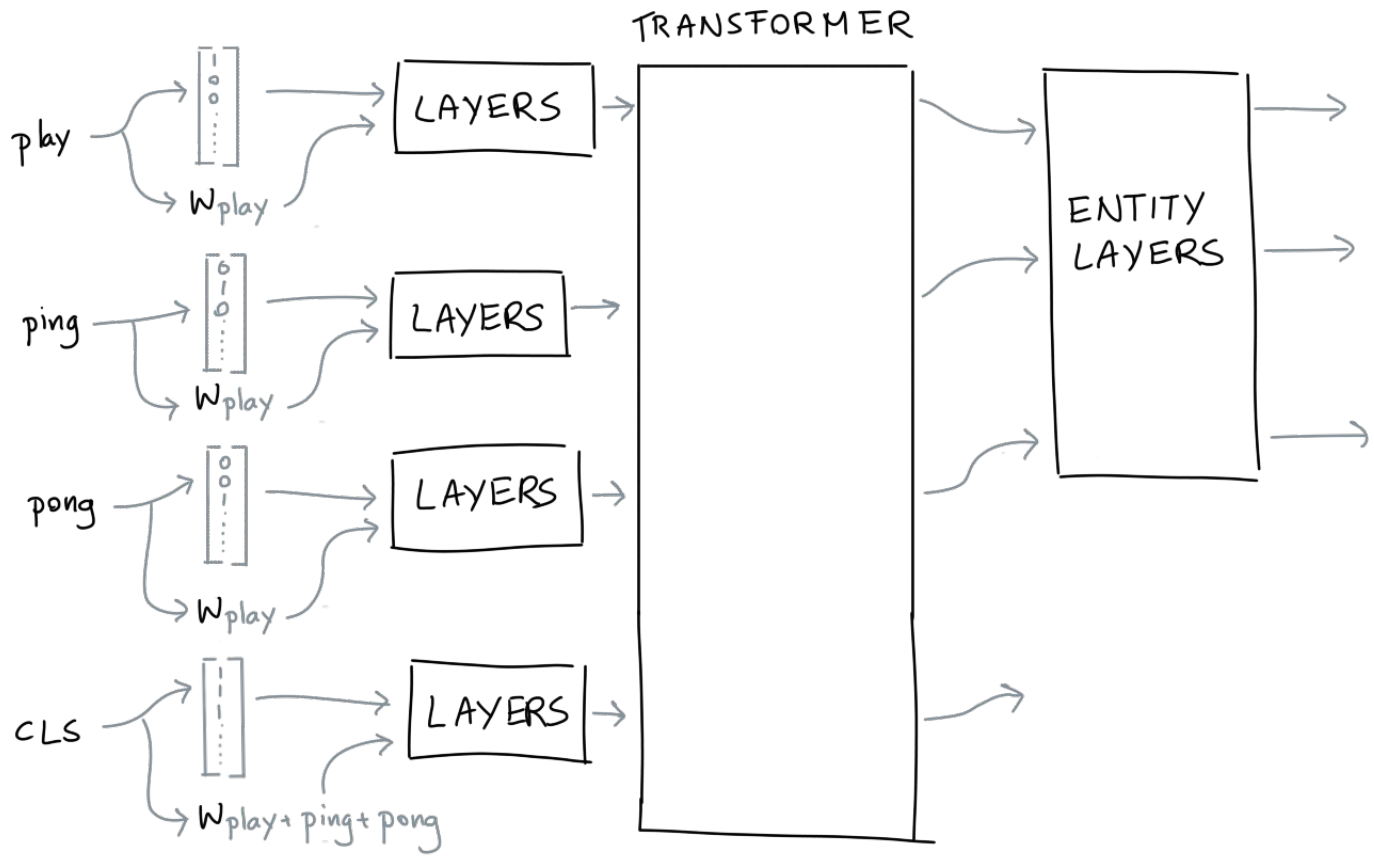
CLS -

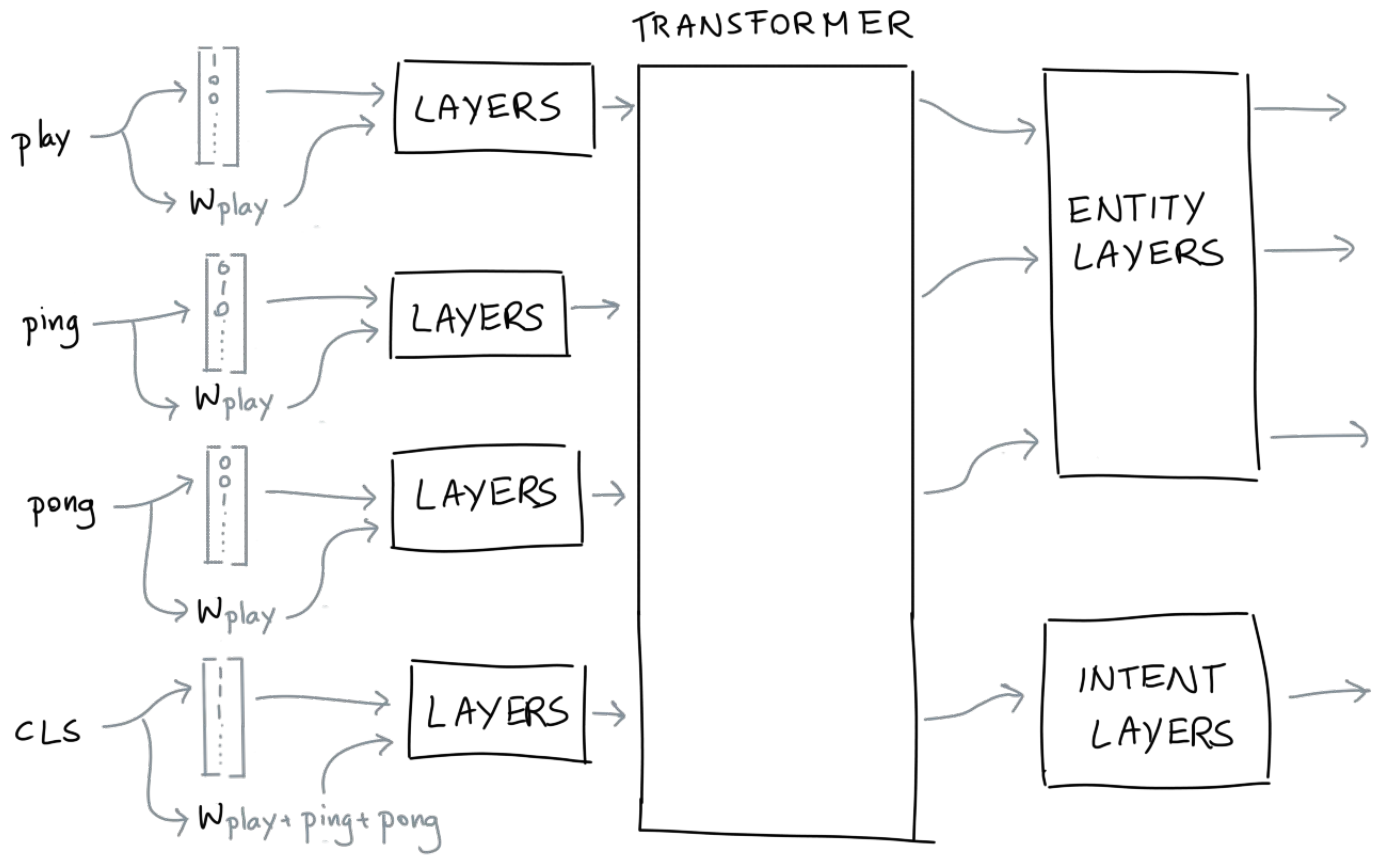


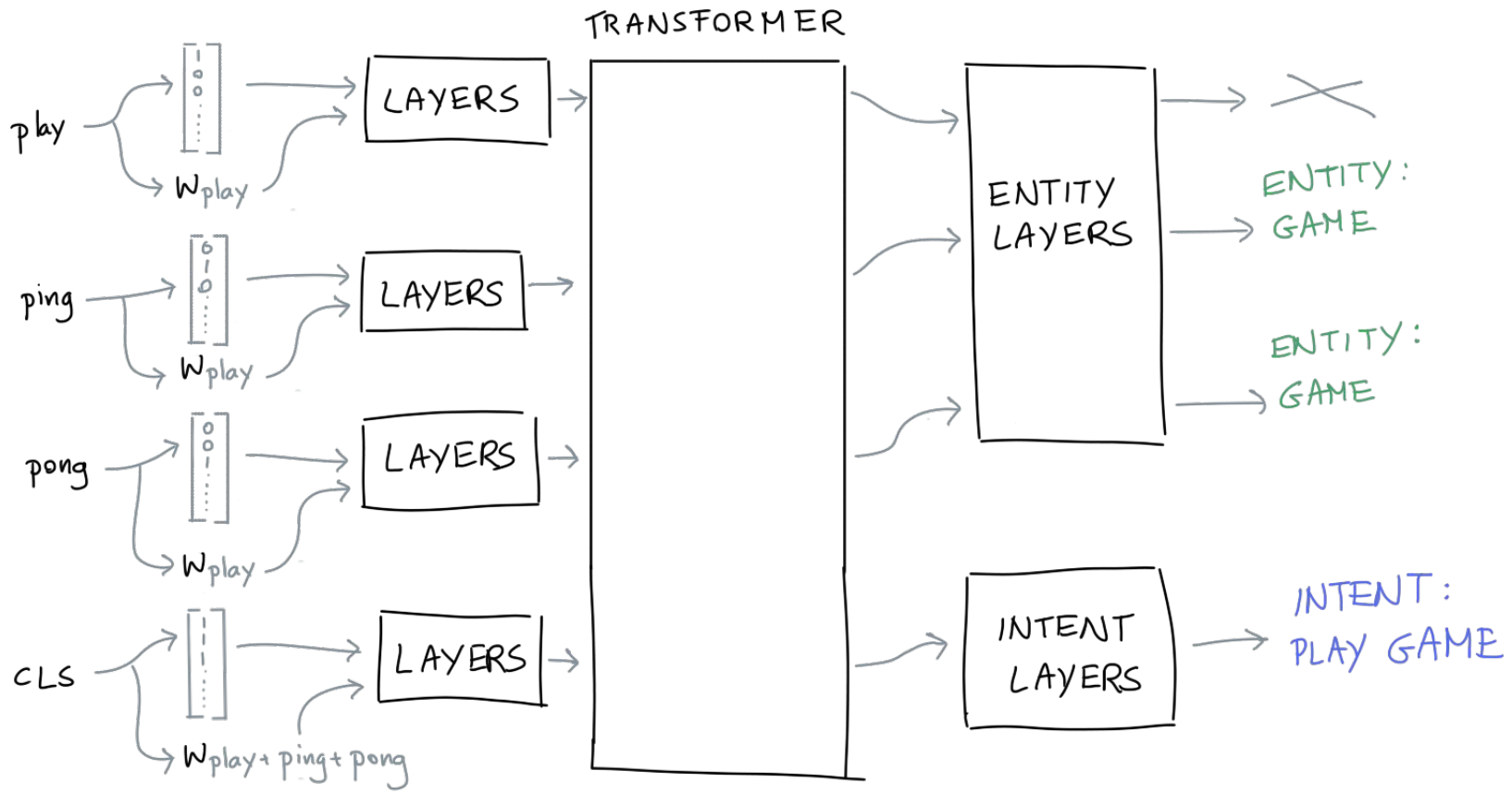




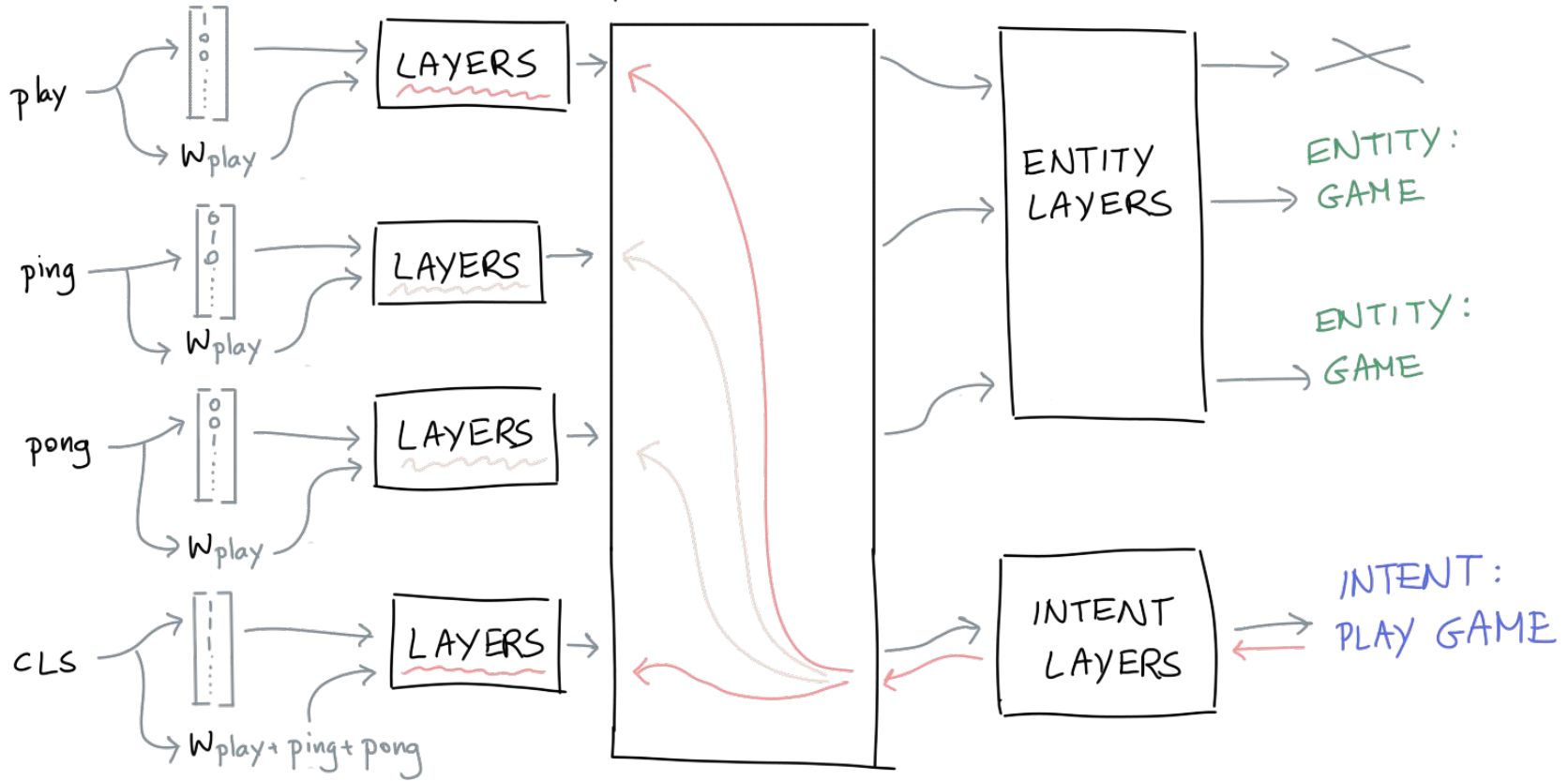








# TRANSFORMER



Q: So how does Rasa determine what actions to take?

A: TED  
Transformer Embedding Dialogue

# TED policy

"I'd like a pizza."



# TED policy

"I'd like a pizza." →

intent

•  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

order

# TED policy

"I'd like a pizza." →

intent

•  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

order

entity

•  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

pizza

# TED policy

"I'd like a pizza."



intent

•  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

order

entity

•  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

pizza

slot

•  $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

adress

# TED policy

"I'd like a pizza."



intent

•  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

order

slot

•  $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

address

entity

•  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

pizza

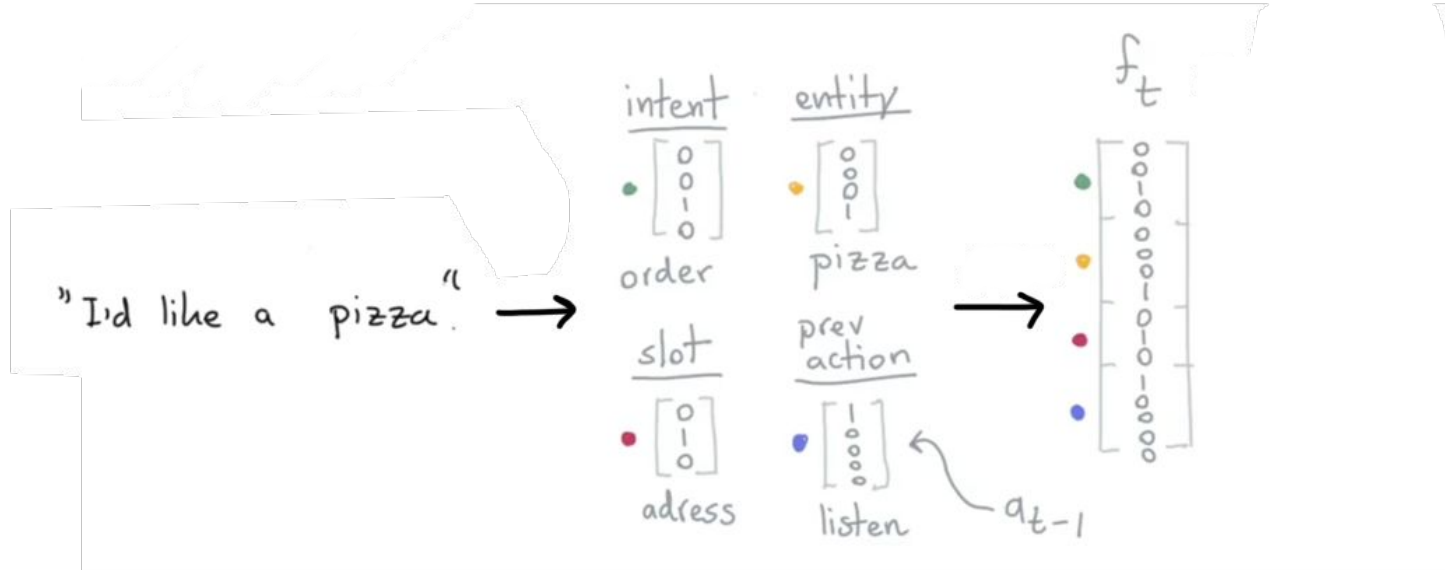
prev  
action

•  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

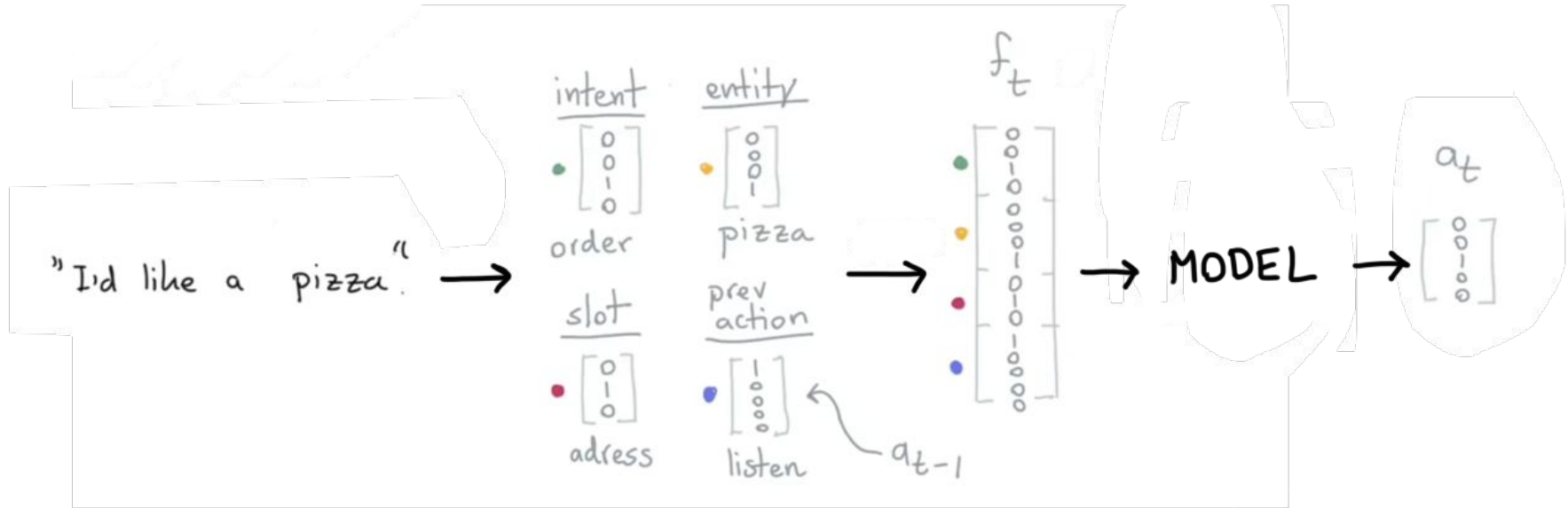
listen

$a_{t-1}$

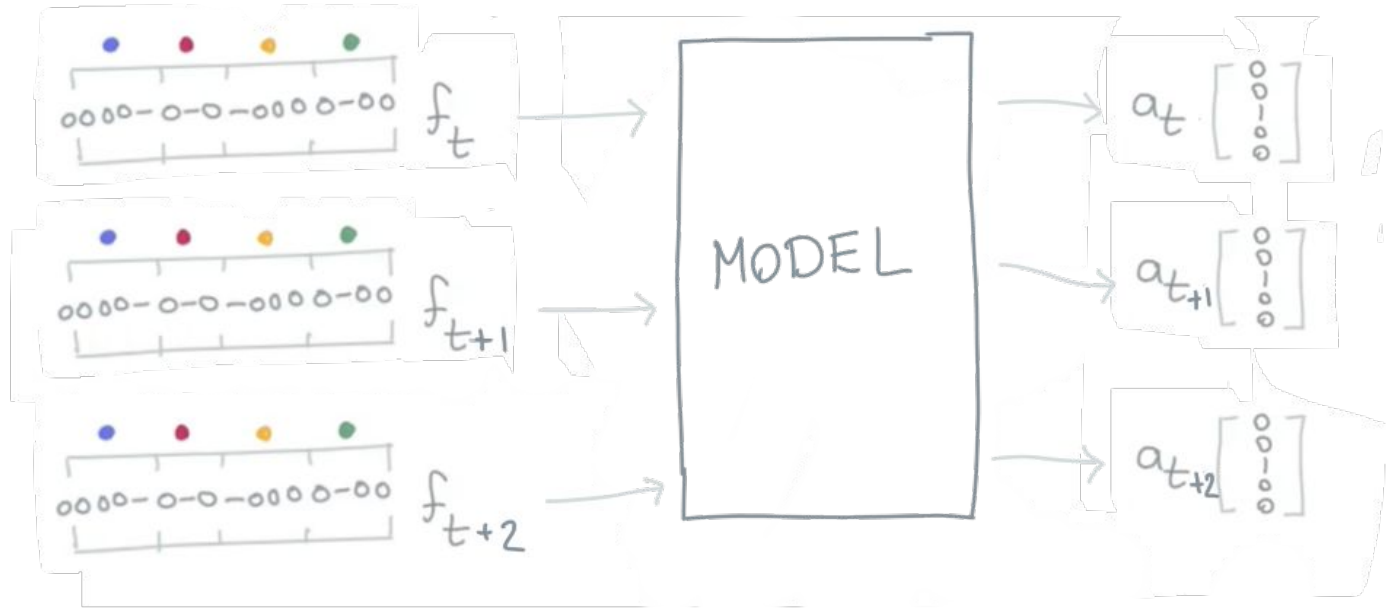
# TED policy

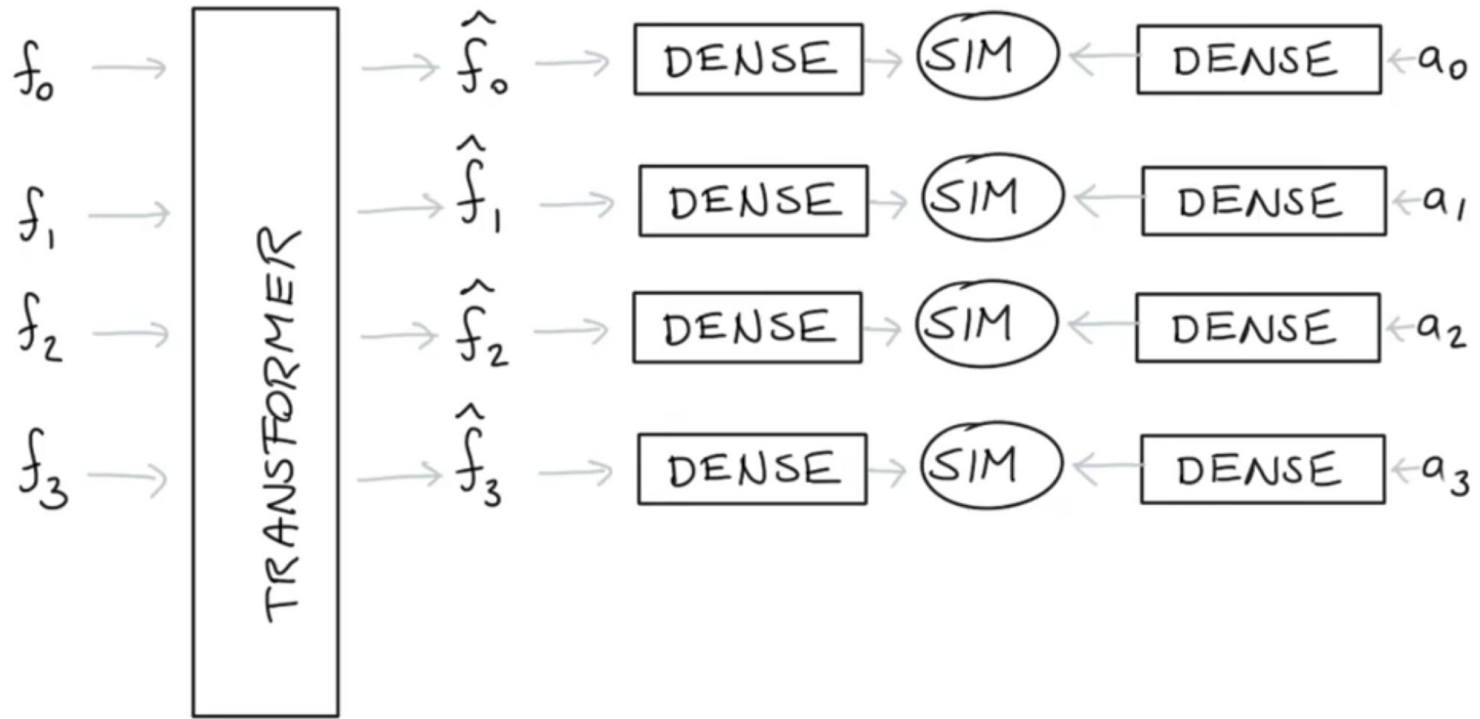


# TED policy

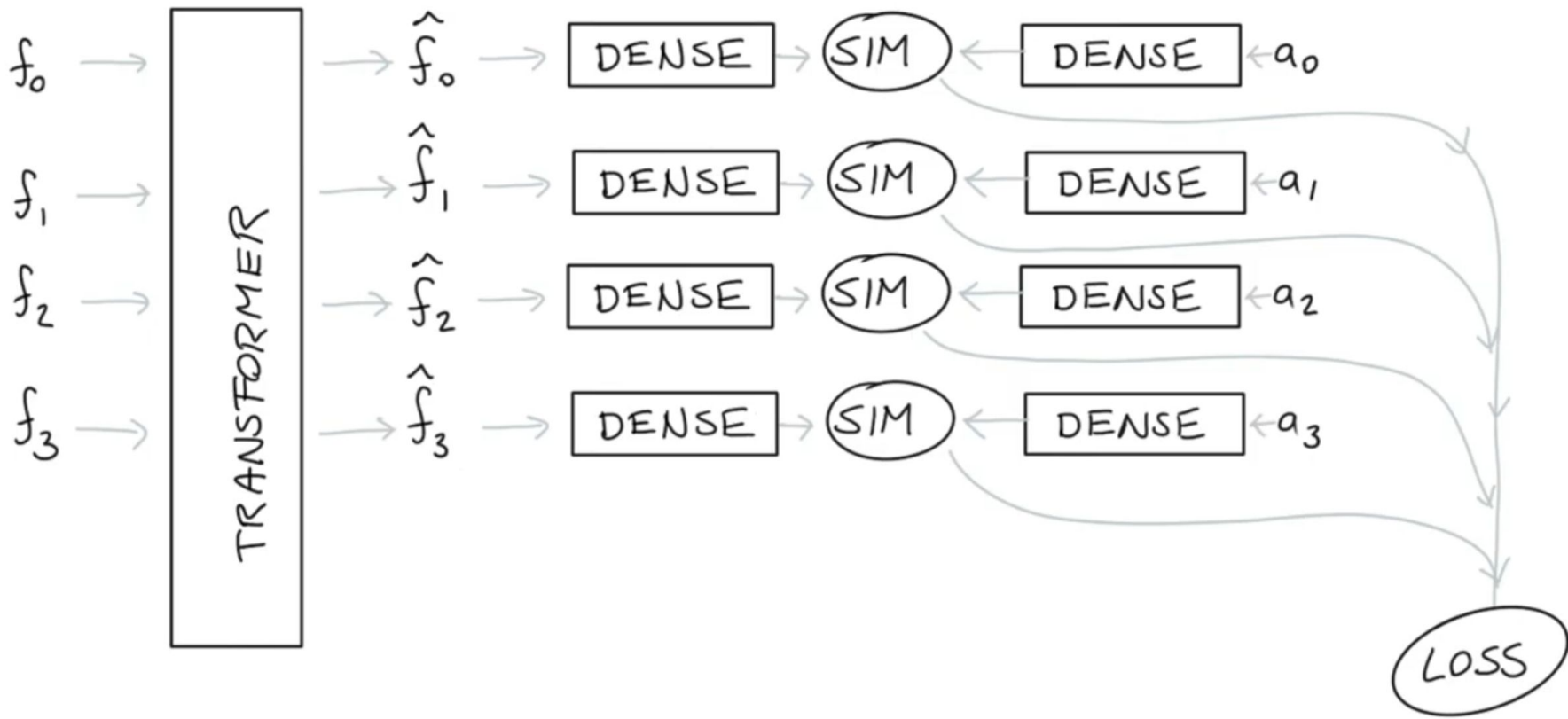


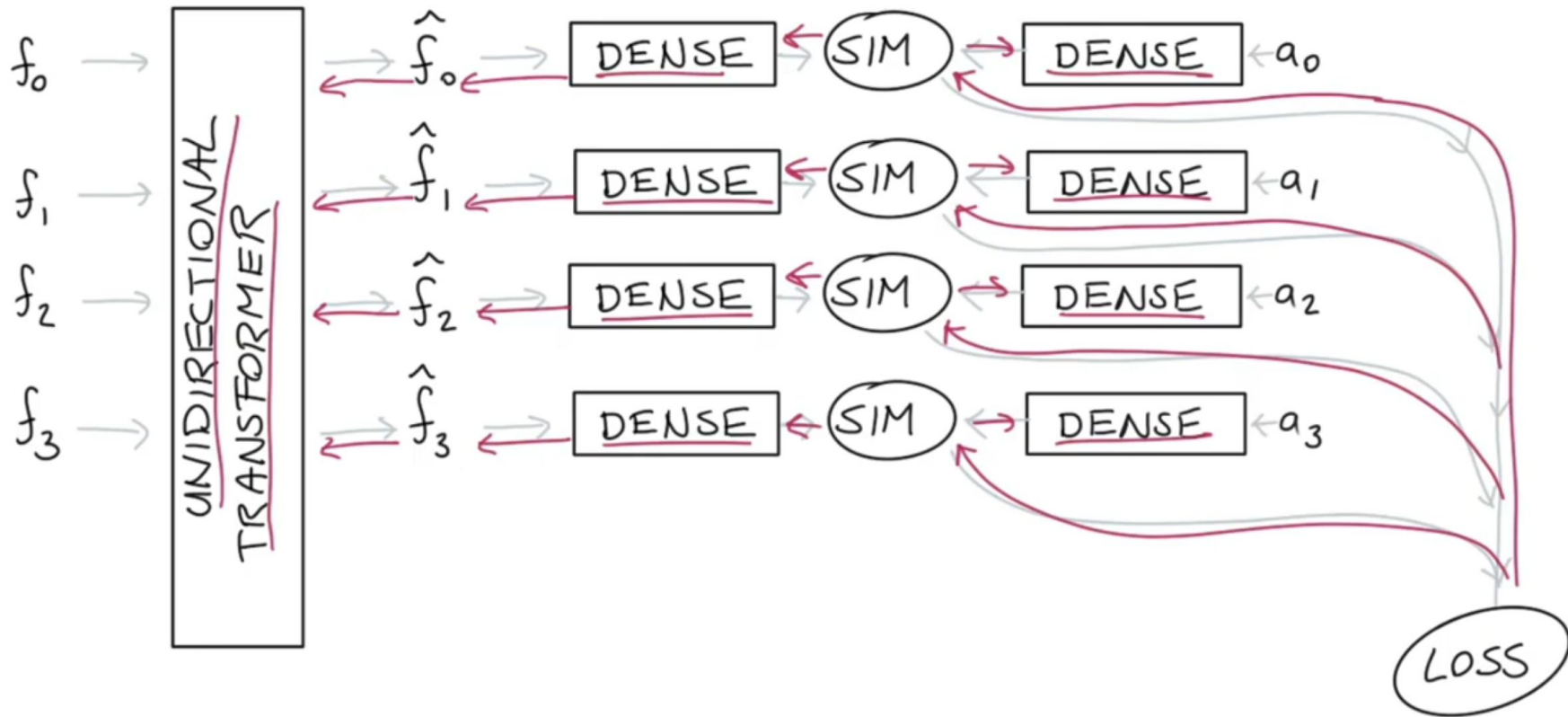
# TED policy

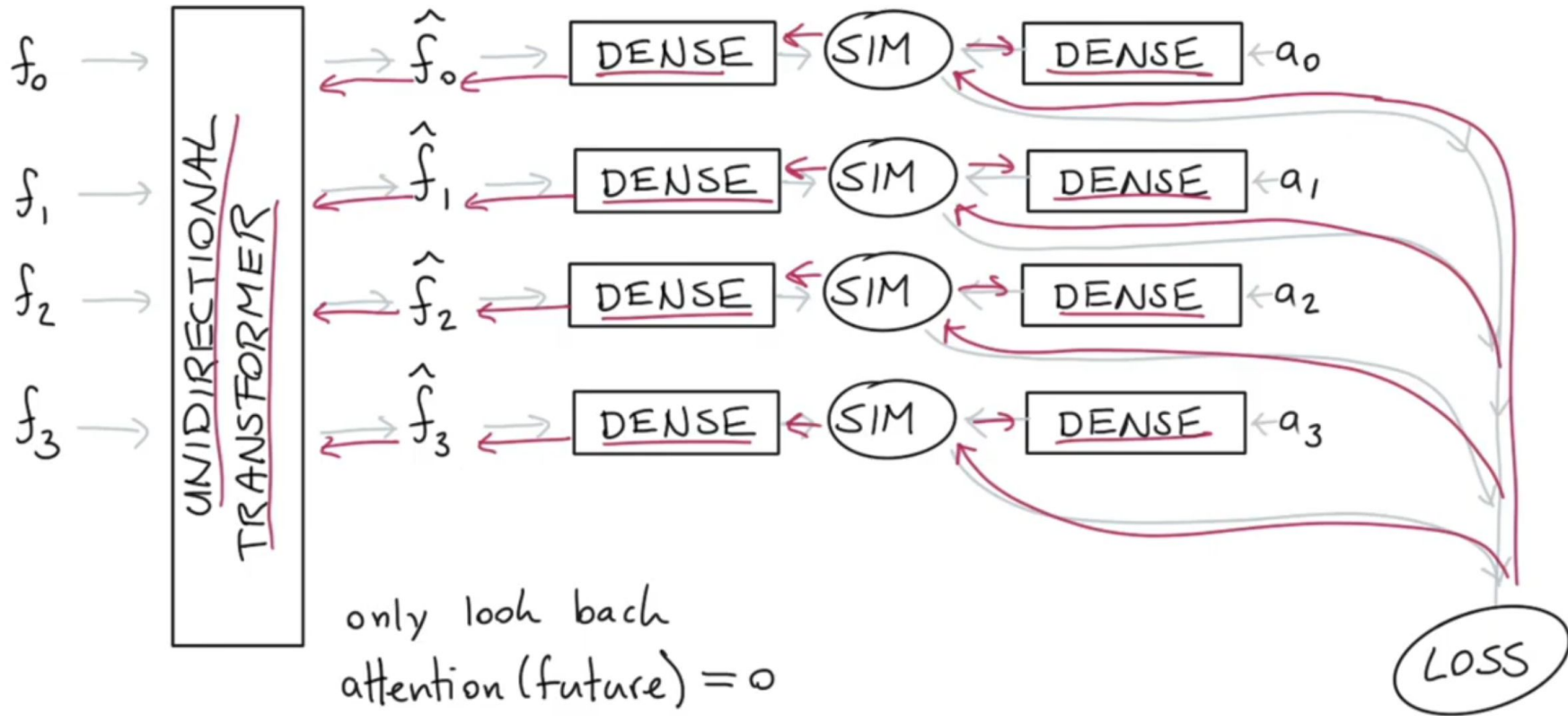












Demotime

## config.yml

policies:

- name: MemoizationPolicy
- name: TEDPolicy  
epochs: 200  
max\_history: 1
- name: MappingPolicy

Your input -> start counting

Countdown! ETA 10 🎵

Your input -> ok

Countdown! ETA 5 🎵

Your input -> ok

Countdown! ETA 5 🎵

Your input -> ok

Countdown! ETA 5 🎵

Your input -> ok

Countdown! ETA 5 🎵

## config.yml

### policies:

- name: MemoizationPolicy
- name: TEDPolicy  
epochs: 200  
max\_history: 3
- name: MappingPolicy

Your input -> count  
Countdown! ETA 10 🎵

Your input -> ok  
Countdown! ETA 9 🎵

Your input -> ok  
Countdown! ETA 8 🎵

Your input -> ok  
Countdown! ETA 7 🎵

Your input -> are you a bot?  
I am a bot, not a human, powered by Rasa.  
Countdown! ETA 6 🎵

Your input -> ok  
Countdown! ETA 5 🎵

Your input -> are you a bot?  
I am a bot, not a human, powered by Rasa.

Your input ->  
Countdown! ETA 10 🎵

## config.yml

### policies:

- name: MemoizationPolicy
- name: TEDPolicy  
epochs: 200  
max\_history: 10
- name: MappingPolicy

```
Your input -> count
    Countdown! ETA 10 🎵
Your input -> ok
    Countdown! ETA 9 🎵
Your input -> ok
    Countdown! ETA 8 🎵
Your input -> are you a bot?
    I am a bot, not a human, powered by Rasa.
    Countdown! ETA 7 🎵
Your input -> ok
    Countdown! ETA 6 🎵
Your input -> are you a bot?
    I am a bot, not a human, powered by Rasa.
    Countdown! ETA 5 🎵
Your input -> ok
    Countdown! ETA 4 🎵
Your input -> are you a bot?
    I am a bot, not a human, powered by Rasa.
    Countdown! ETA 3 🎵
Your input -> ok
    Countdown! ETA 2 🎵
Your input -> are you a bot?
    I am a bot, not a human, powered by Rasa.
    Countdown! ETA 1 🎵
Your input -> ok
    End of      countdown!
```

## config.yml

### Using LSTM

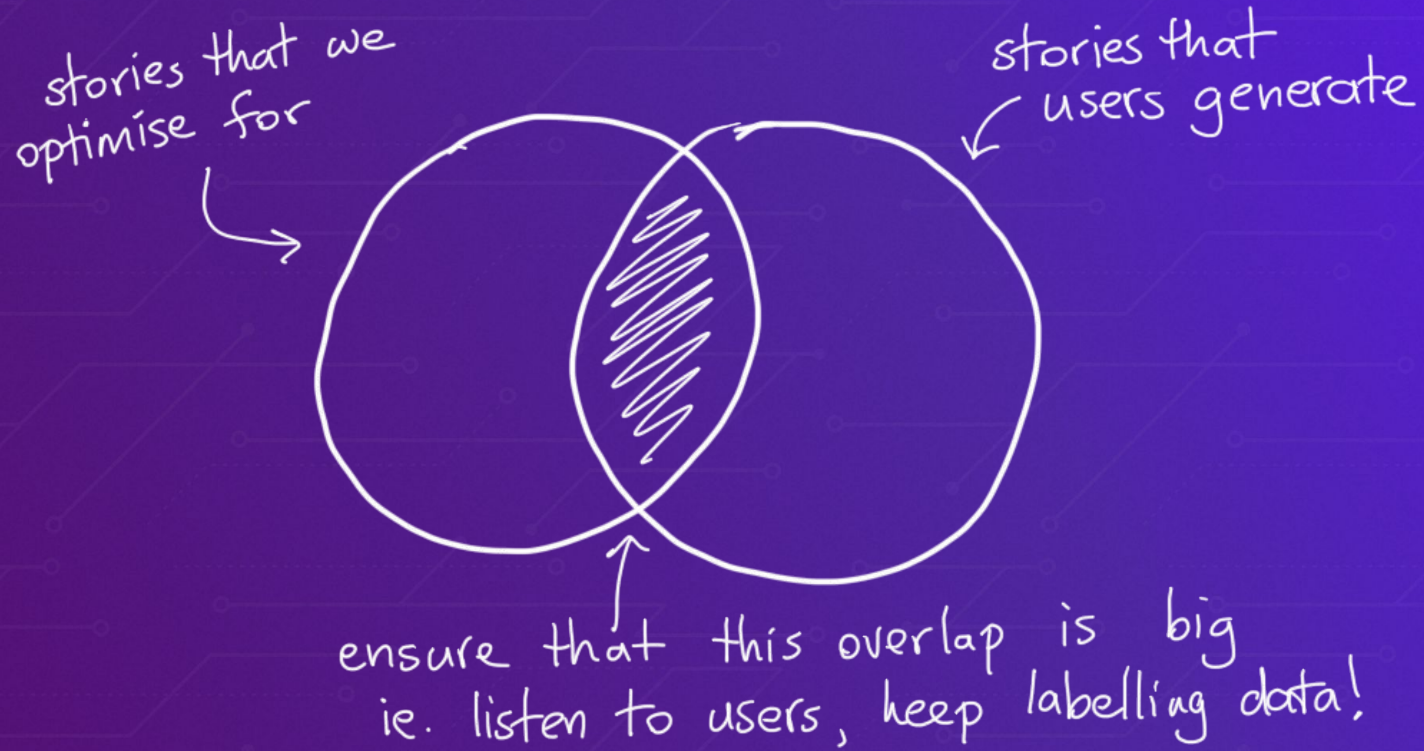
policies:

- name: MemoizationPolicy
  - name: KerasPolicy
- epochs: 200
- name: MappingPolicy

```
Your input -> count
Countdown! ETA 10 🎵
Your input -> ok
Countdown! ETA 9 🎵
Your input -> ok
Countdown! ETA 8 🎵
Your input -> are you a bot?
I am a bot, not a human, powered by Rasa.
Countdown! ETA 7 🎵
Your input -> are you a bot
I am a bot, not a human, powered by Rasa.
Countdown! ETA 6 🎵
Your input -> are you a bot?
Countdown! ETA 5 🎵
Your input -> ok
Countdown! ETA 4 🎵
Your input -> are you a bot?
I am a bot, not a human, powered by Rasa.
Countdown! ETA 3 🎵
Your input -> are you a bot?
I am a bot, not a human, powered by Rasa.
```



# The real problem though.



While I have you here: check the [Rasa Algorithm Whiteboard!](#)

Our algorithms are explained in more detail! I also take requests for content.

## Self Attention

RASA ALGORITHM WHITEBOARD

"Naa" can be annoying, but she is a great cat.

think of a way to do this, maybe automatically

cannot blindly use proximity



## GloVe

RASA ALGORITHM WHITEBOARD

dot prod

$$v_i \cdot v_j = a_1 b_1 + a_2 b_2 + a_3 b_3 = \dots$$



## CBOW and Skip Gram

RASA ALGORITHM WHITEBOARD

t<sub>1</sub> → t<sub>2</sub> → t<sub>3</sub> → t<sub>4</sub> → t<sub>5</sub> → t<sub>6</sub> → t<sub>7</sub> letters

t<sub>1</sub> → t<sub>2</sub> → t<sub>3</sub> → t<sub>4</sub> → t<sub>5</sub> → t<sub>6</sub> → t<sub>7</sub> words



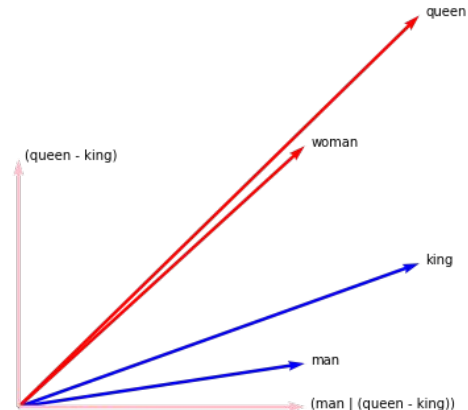
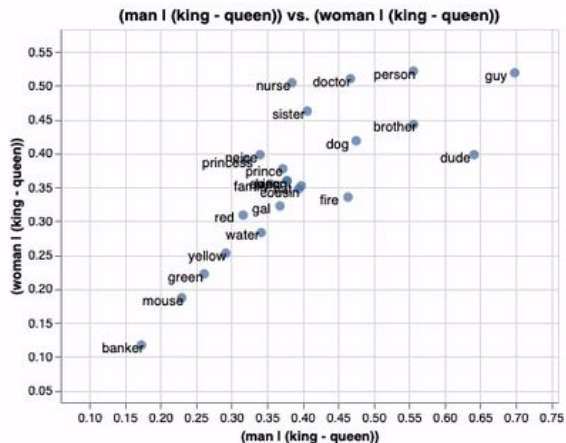
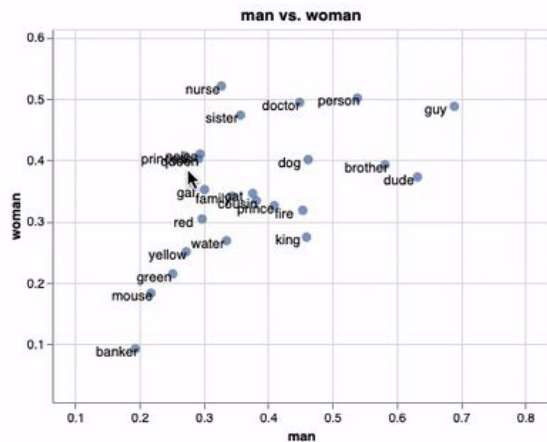
## Multi Head Attention

RASA ALGORITHM WHITEBOARD



## While I have you here: check out [WHATLIES](#)

```
orig_chart = emb.plot_interactive('man', 'woman')
new_ts      = emb | (emb['king'] - emb['queen'])
new_chart  = new_ts.plot_interactive('man', 'woman')
```



It's an open source package for visualising word embeddings.  
Soon: features for detecting bias. Feedback is appreciated!

# Get in touch!



**Vincent D. Warmerdam**

*Research Advocate*

[v.warmerdam@rasa.com](mailto:v.warmerdam@rasa.com)

# Appendix

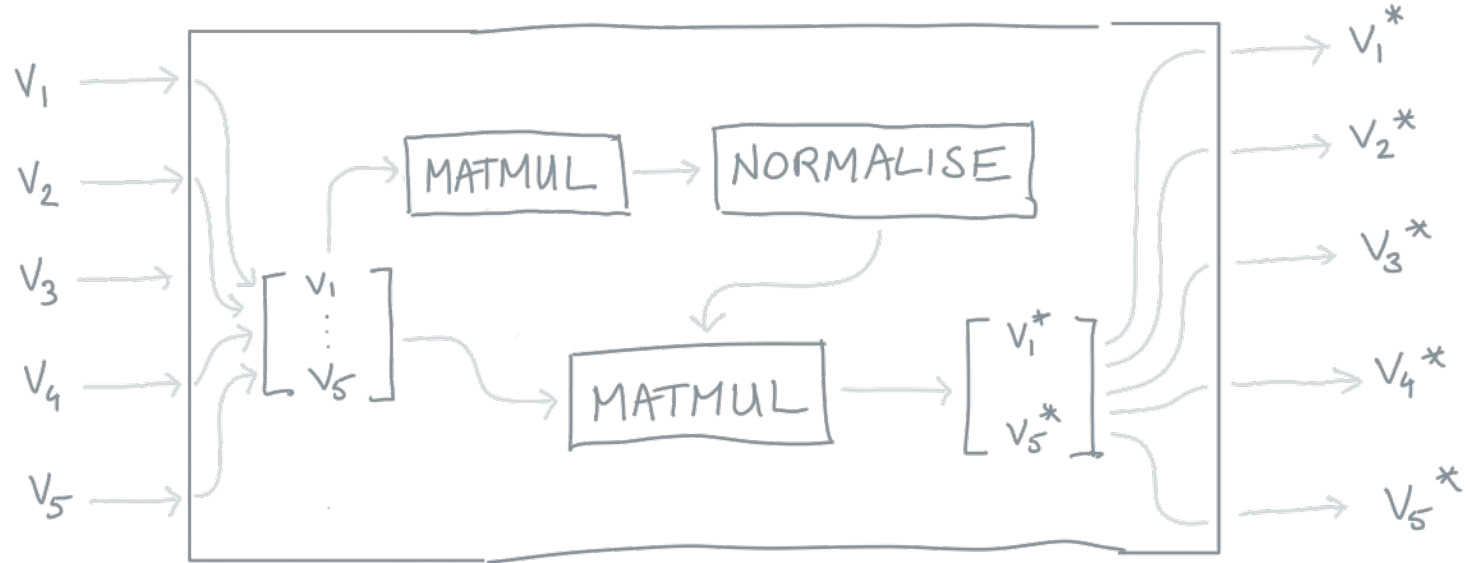


**Vincent D. Warmerdam**

*Research Advocate*

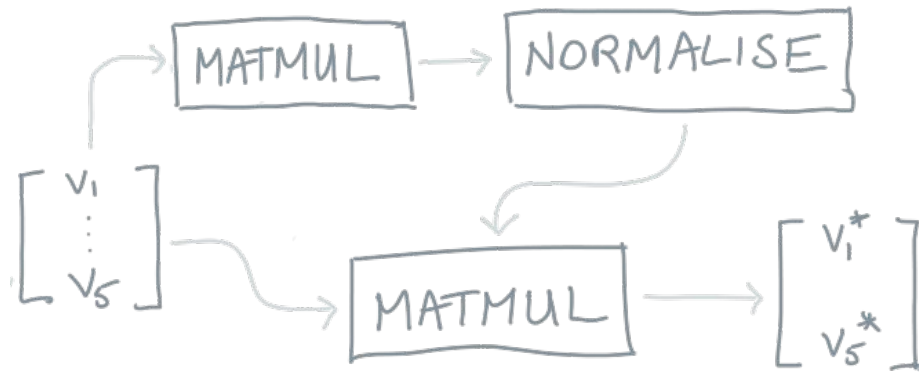
[v.warmerdam@rasa.com](mailto:v.warmerdam@rasa.com)

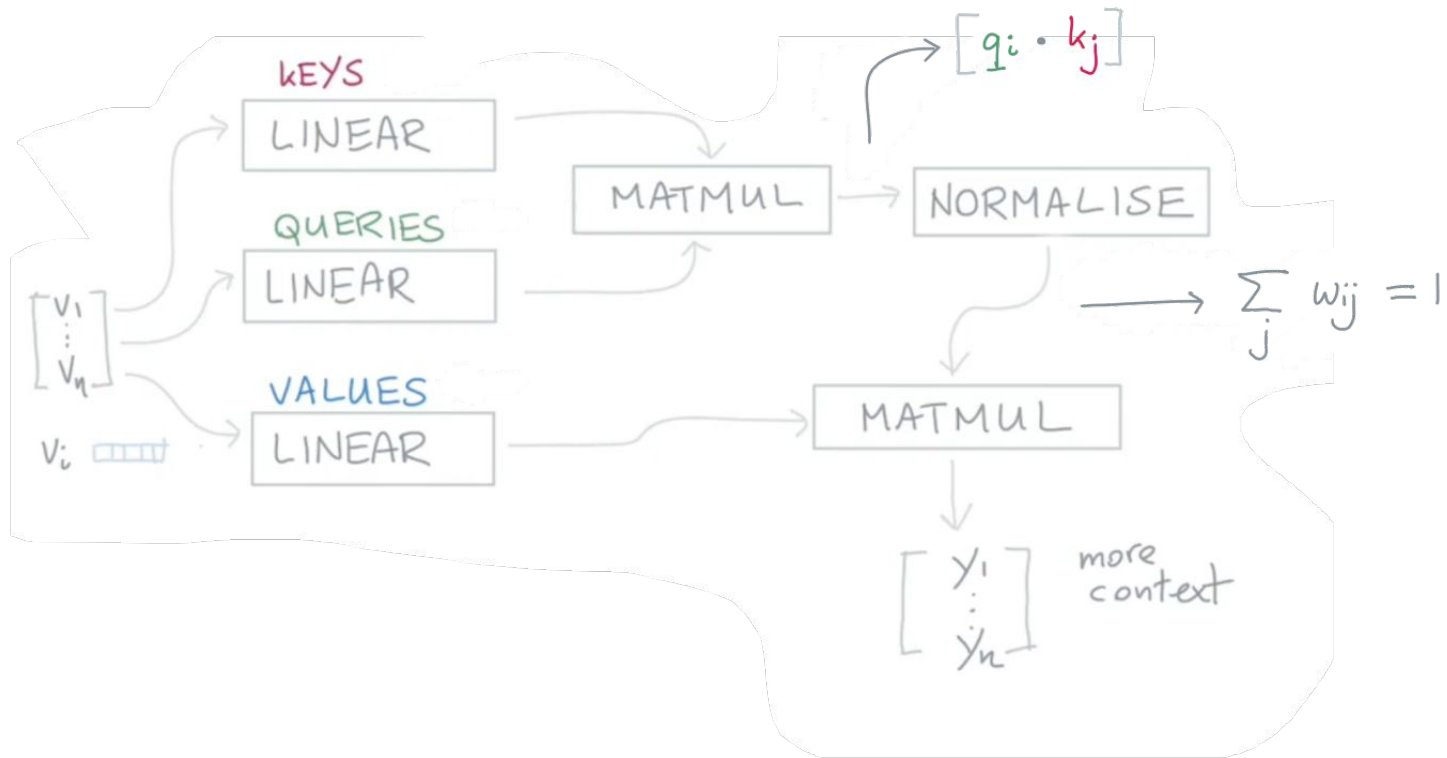
# SELF ATTENTION BLOCK



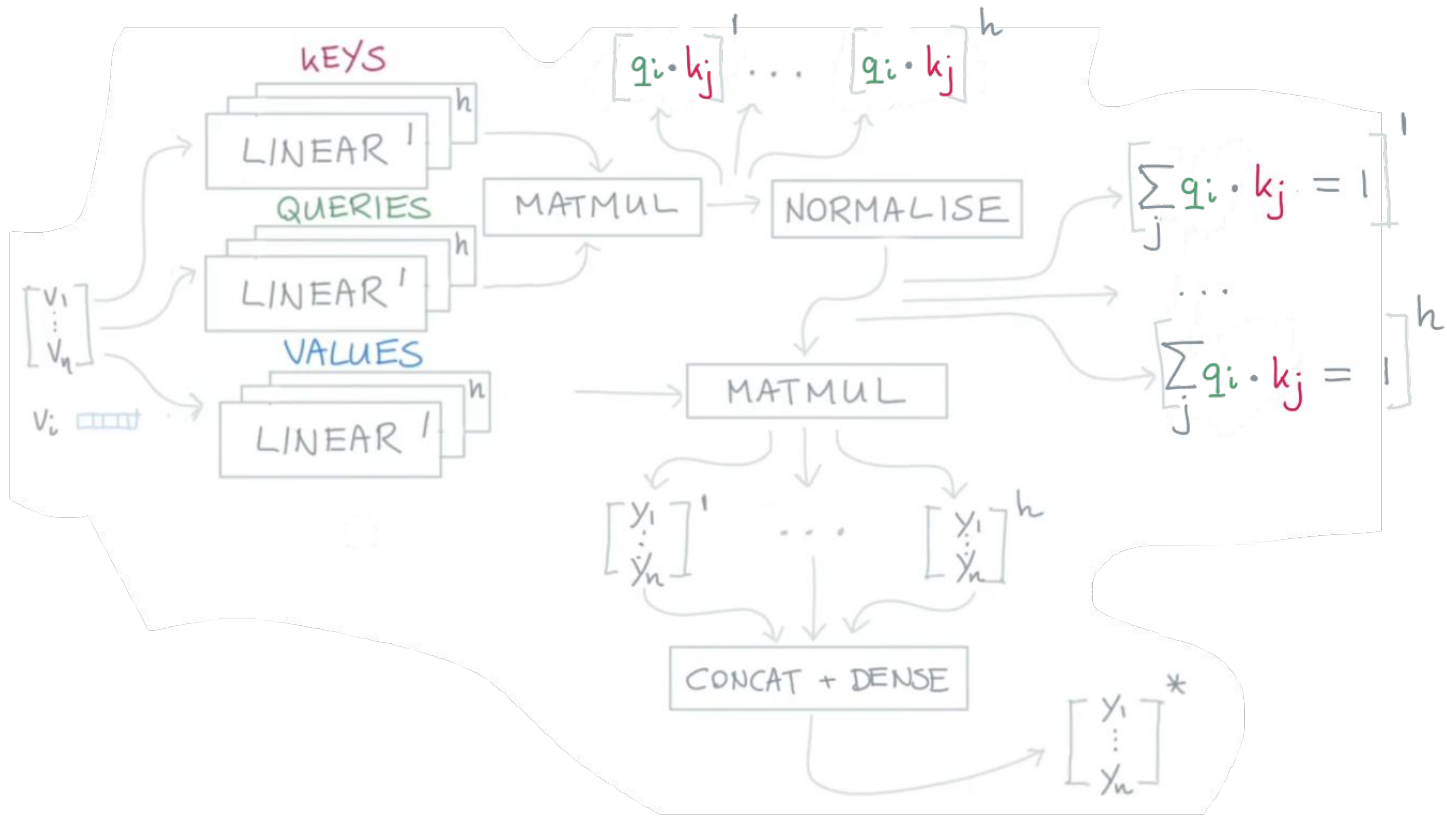
- ORDER DOES NOT MATTER
- NO TRAINABLE WEIGHTS (YET)
- DEPENDS ON PRETRAINED QUALITY
- TASK INDEPENDANT

How's about we add us  
some trainable weights?

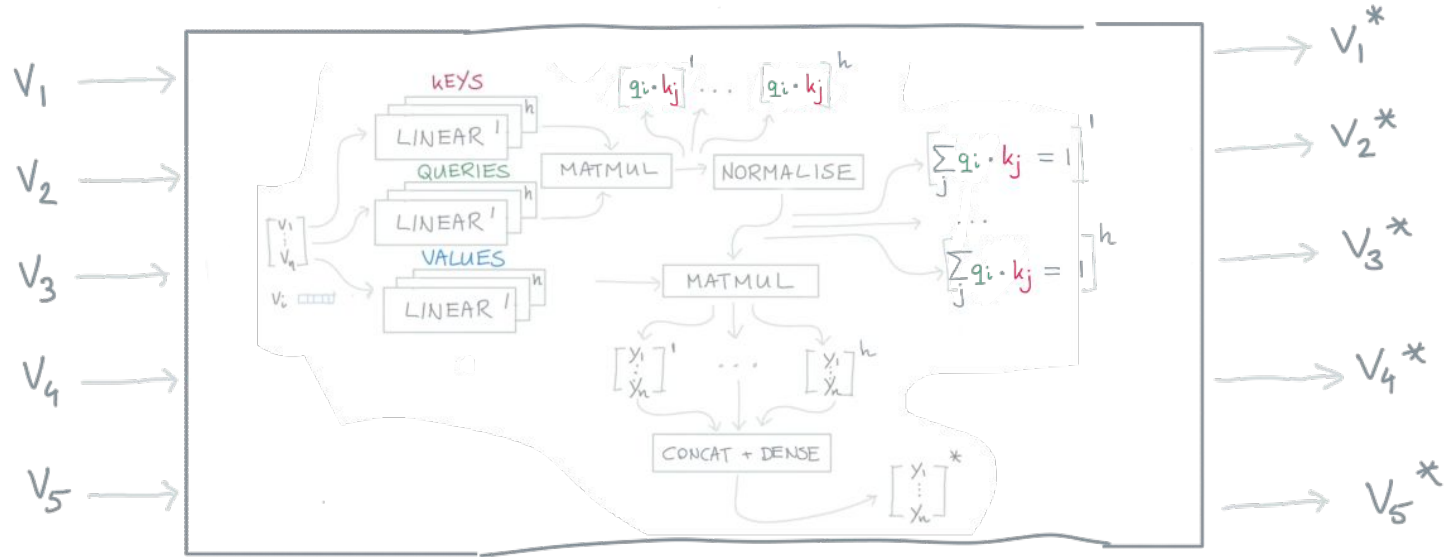




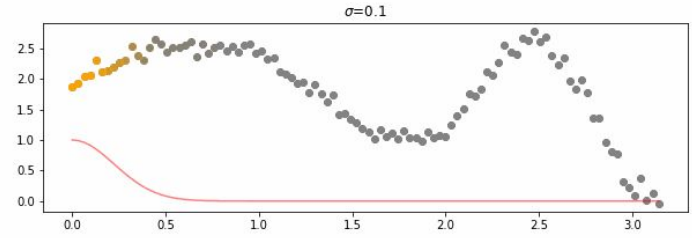


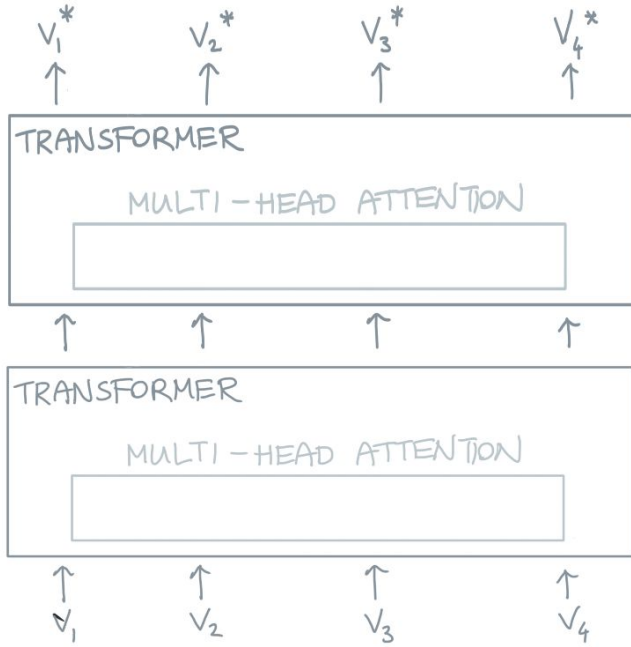


# MULTI HEAD ATTENTION



It's a more elaborate way to do stuff like this →





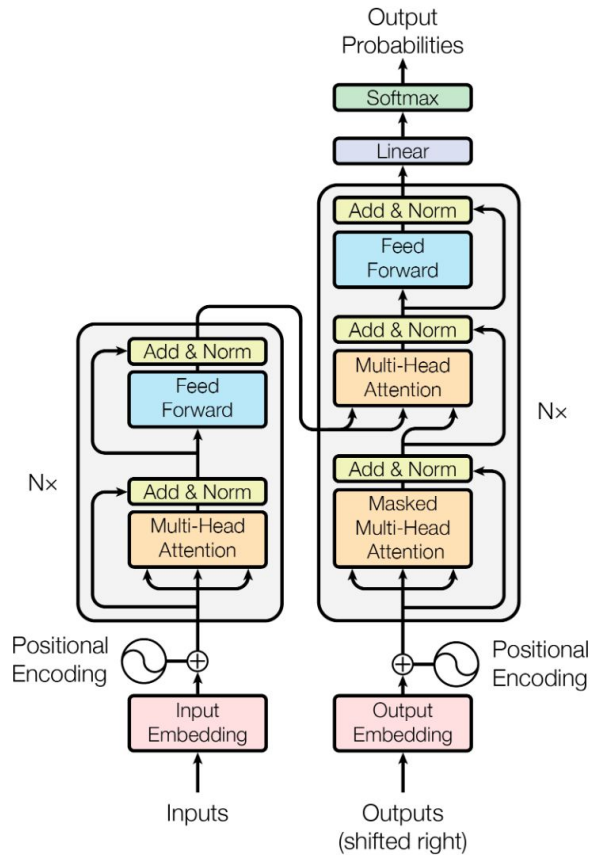
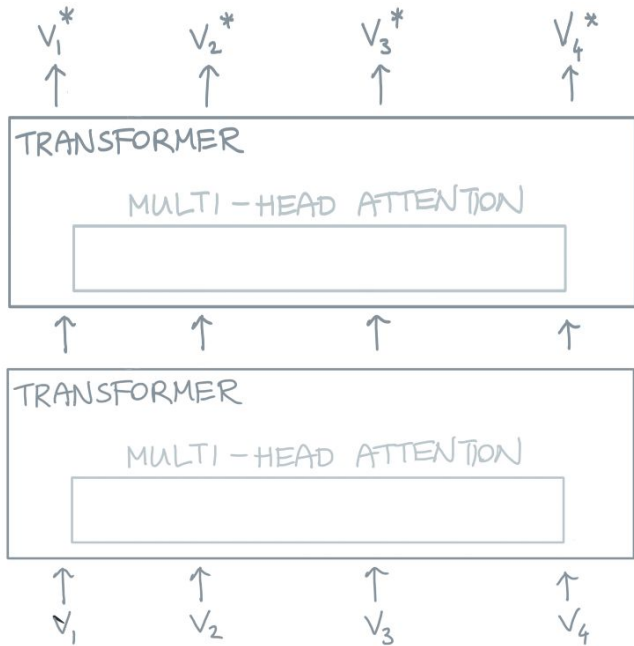
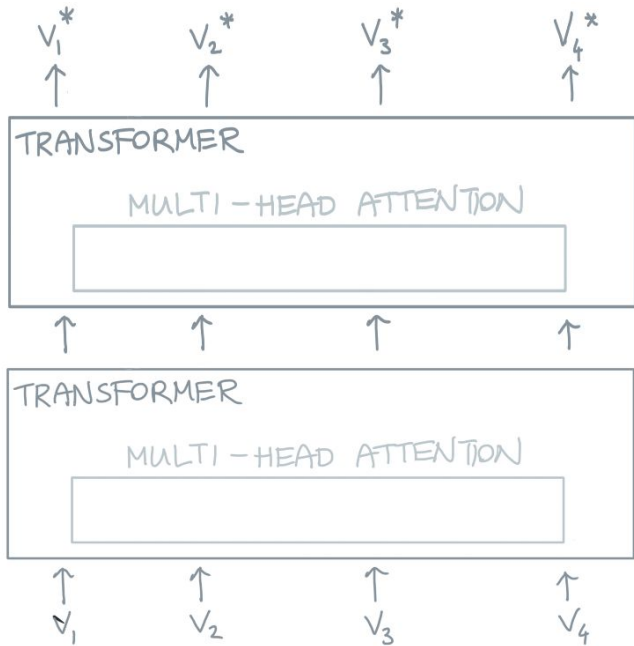


Figure 1: The Transformer - model architecture.



We use this **encoder bit**.

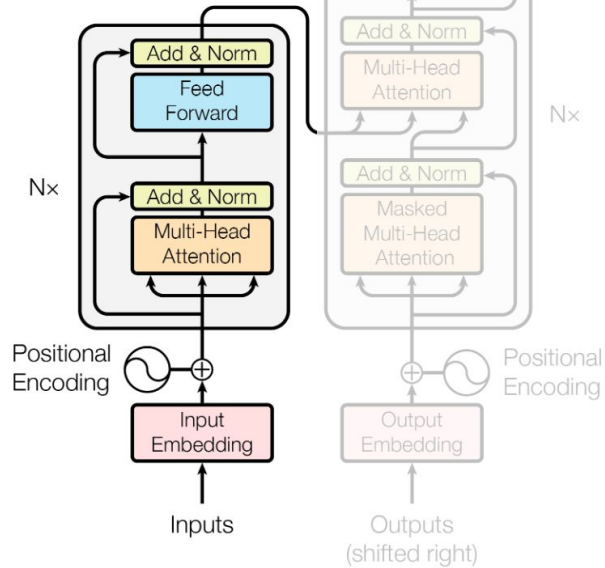


Figure 1: The Transformer - model architecture.