

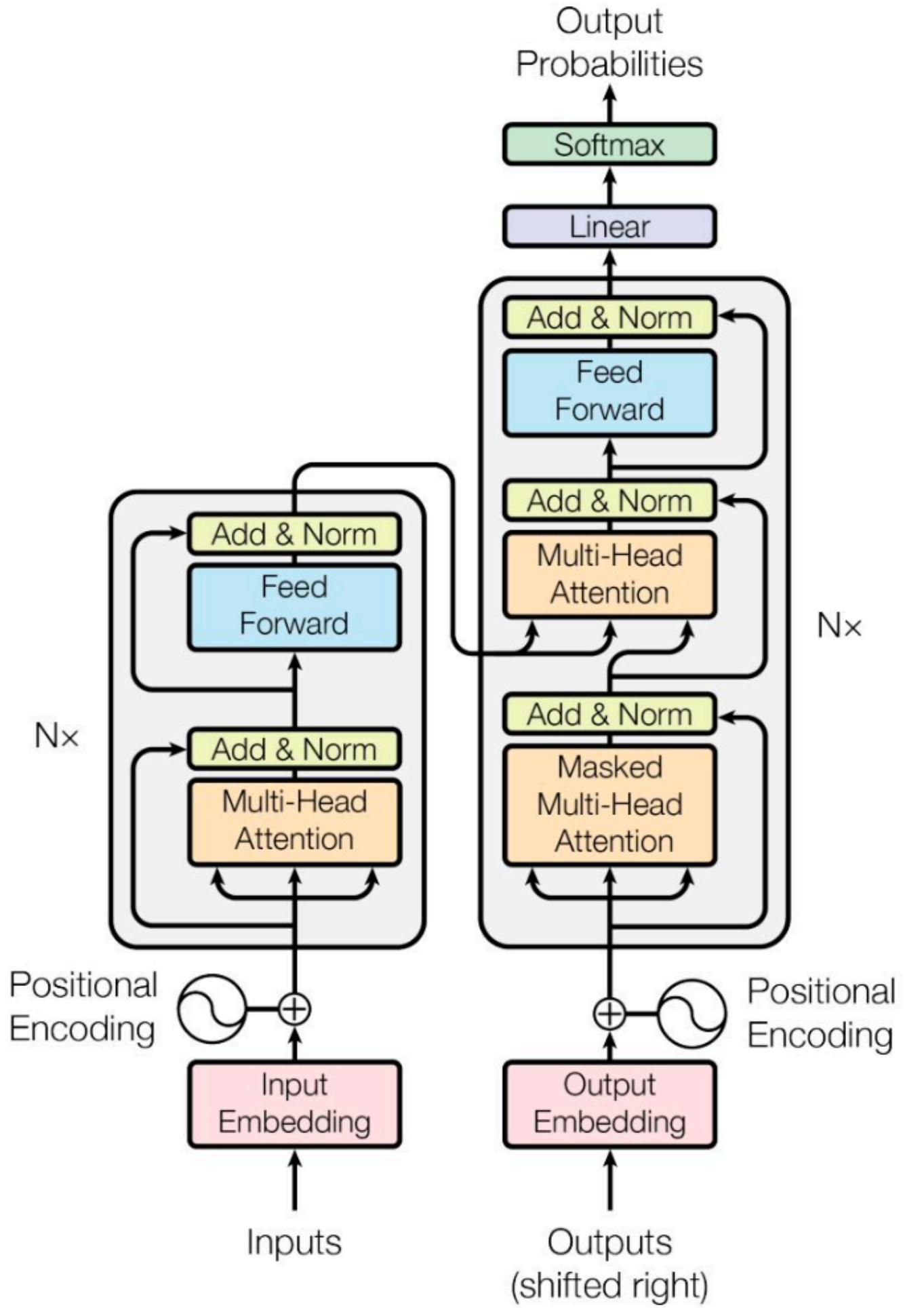


ESPER - HELEN - APARNA - ADAM

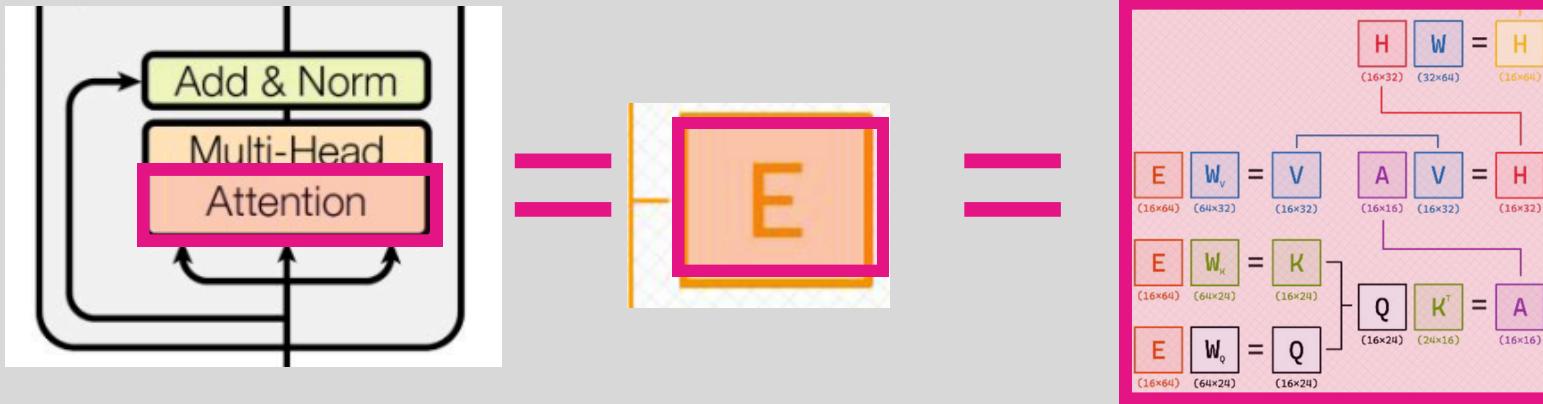
# RECURRENT REBELS

TRANSFORMERS AND MNIST

# ARCHITECTURE



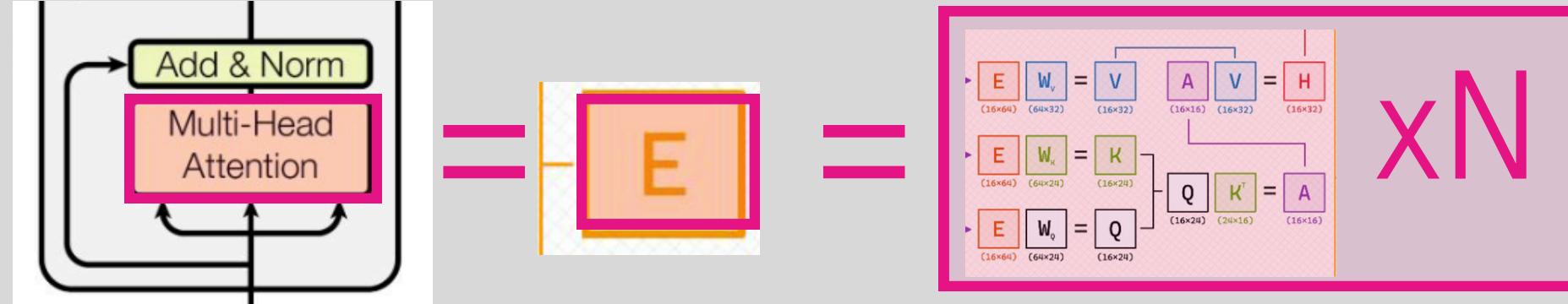
# GOING INTO THE DEPTHS OF THE CODE



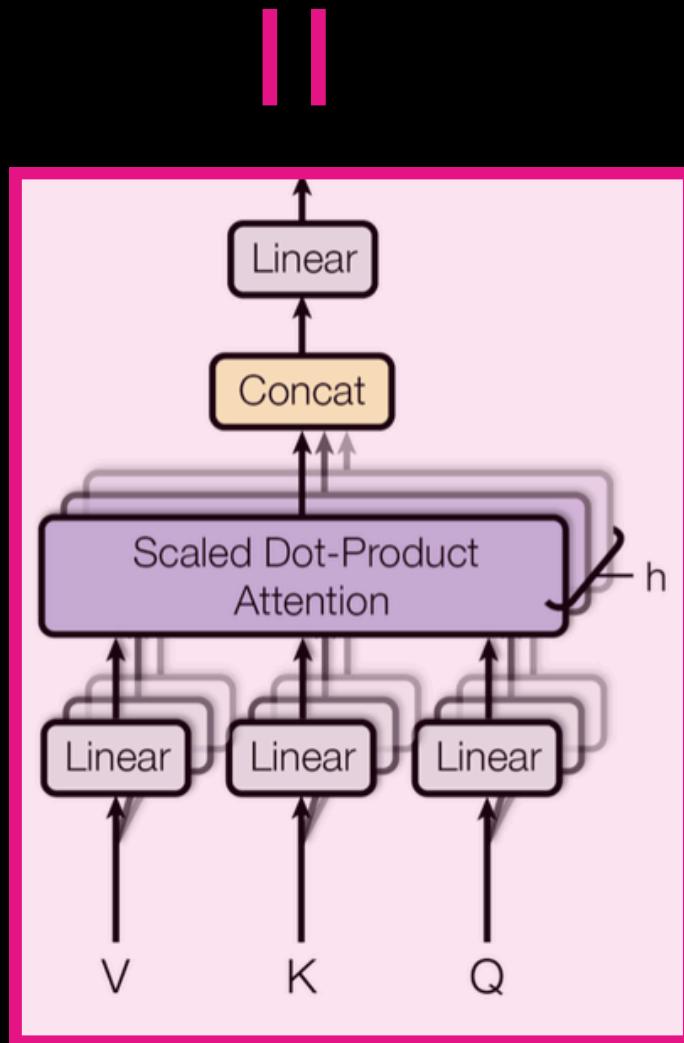
## SINGLE HEADED ATTENTION

```
>  $E \cdot W_q = Q$       =  
  (16x64) (64x24)    (16x24)  
>  $E \cdot W_k = K$       =  
  (16x64) (64x24)    (16x24)  
>  $E \cdot W_v = V$       =  
  (16x64) (64x32)    (16x32)  
  
>  $Q \cdot K^T = A$       =  
  (16x24) (24x16)    (16x16)  
>  $A \cdot V = H$         =  
  (16x16) (16x32)    (16x32)  
>  $H \cdot W = H$         =  
  (16x32) (32x64)    (16x64)  
  
= Q = self.query_proj(x)  # is nn.Linear(dim_input, dim_k)  
= K = self.key_proj(x)    # is nn.Linear(dim_input, dim_k)  
= V = self.val_proj(x)    # is nn.Linear(dim_input, dim_k)  
  
{ attention = Q @ K.transpose(-2, -1) * self.dim_k ** -0.5  
  attention = torch.softmax(attention, dim=-1)  
= out = torch.matmul(attention, V)  
= return self.out_proj(out)  # is nn.Linear(dim_k, dim_input)
```

# GOING INTO THE DEPTHS OF THE CODE



## MULTI HEADED ATTENTION

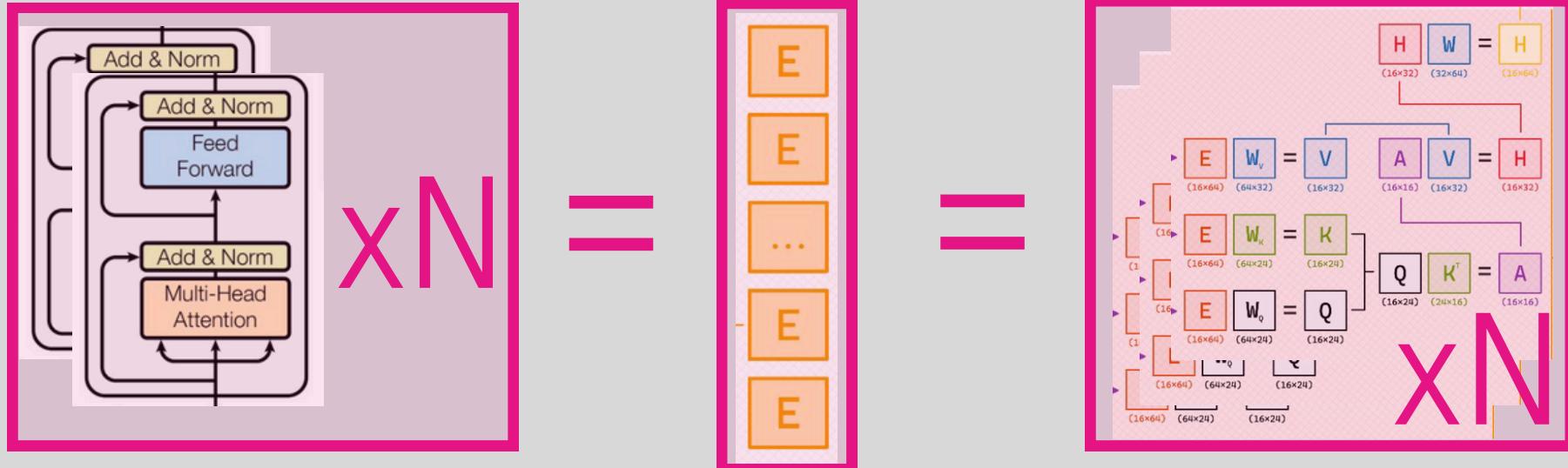


```
batch_size = x.size(0)
qkv = self.qkv_proj(x) # is nn.Linear(dim_input, dim_k * 3)
Q, K, V = qkv.chunk(3, dim=-1)
attention = Q @ K.transpose(-2, -1) / (self.head_dim ** 0.5)
attention = nn.functional.softmax(attention, dim=-1)
attention = attention @ V # shape: (B, num_heads, num_patches, head_dim)
# Below concat all outputs of all attention heads into a single embedding per token (patch)
attention = attention.transpose(1, 2).reshape(batch_size, self.num_patches, self.dim_k)
# Below is final projection after concatenating heads
attention = self.out_proj(attention) # (B, num_patches, dim_k)
```

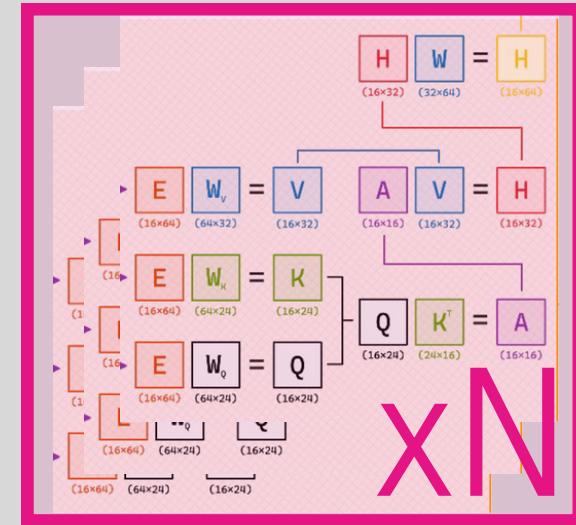
or

```
attention, _ = self.attn(x, x, x) # nn.MultiheadAttention(dim_k, num_heads...)
```

# GOING INTO THE DEPTHS OF THE CODE



## MULTIPLE ENCODER BLOCKS



```
# EncoderBlock has nn.MultiheadAttention & FNN (Linear, ReLU, Linear)
self.encoder_blocks = nn.Sequential([
    EncoderBlock(dim_k, num_heads) for _ in range(num_encoder_blocks)
])
...
x = self.encoder_blocks(x)
```

or

```
self.encoder_blocks = nn.ModuleList([
    EncoderBlock(dim_model, num_heads)
    for _ in range(num_layers)
])
...
for encoder_block in self.encoder_blocks:
    x = encoder_block(x. ...)
```

# SIDE NOTE: WHY DOES MULTIPLYING = MORE UNDERSTANDING

**IMAGINE:** LEGO Encoder Teams trying to understand a LEGO instruction book

**INSTRUCTION:** Attach the red flag to the highest tower

## MULTI HEADED ATTENTION

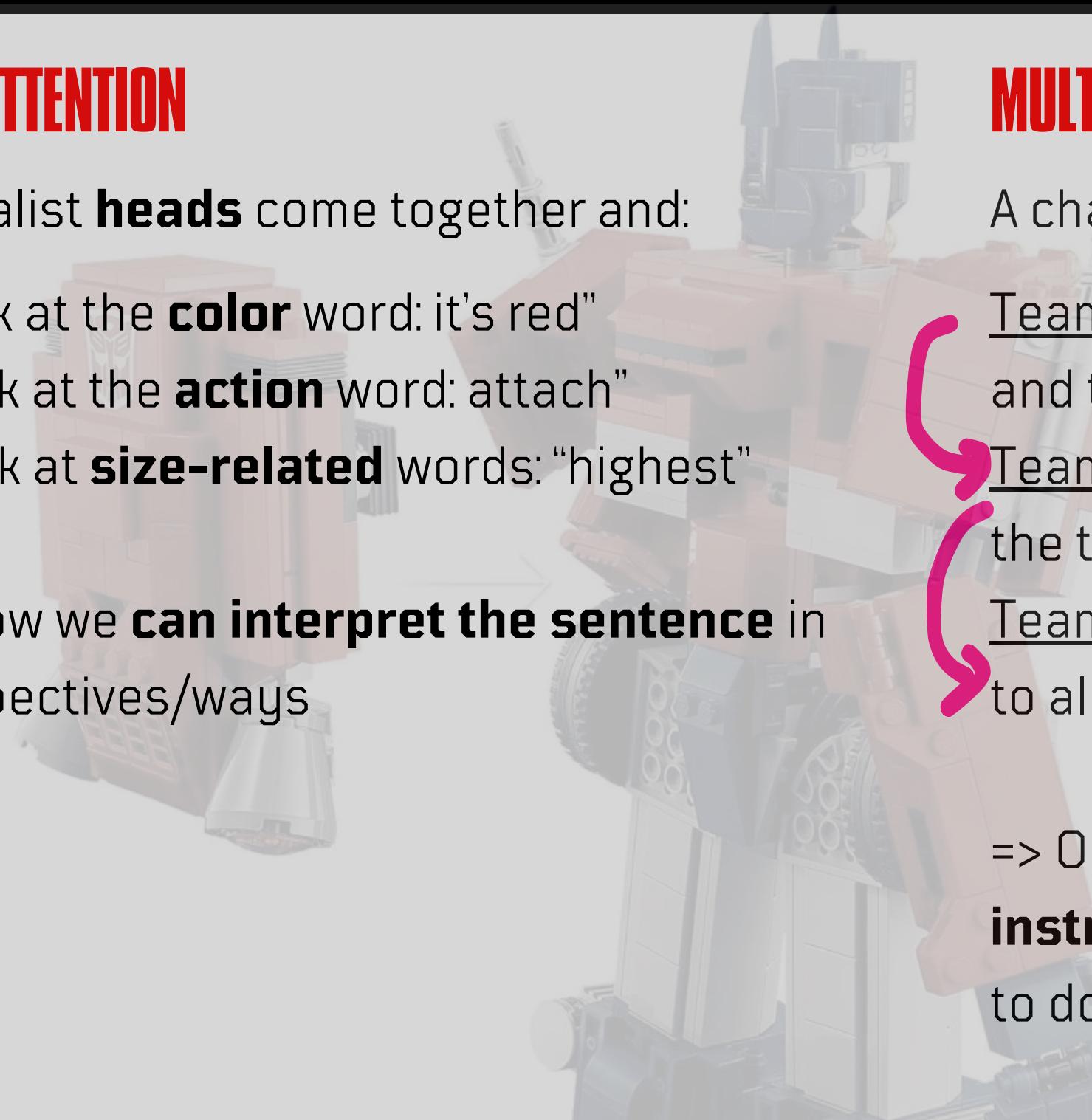
Multiple specialist **heads** come together and:

Head 1: "I'll look at the **color** word: it's red"

Head 2: "I'll look at the **action** word: attach"

Head 3: "I'll look at **size-related** words: "highest"

=> OK great, now we **can interpret the sentence** in different perspectives/ways



## MULTIPLE ENCODER BLOCKS

A chain of **Encoder teams** come together:

Team 1. "Looks like this task involves a tower, flag and the flag going on top"

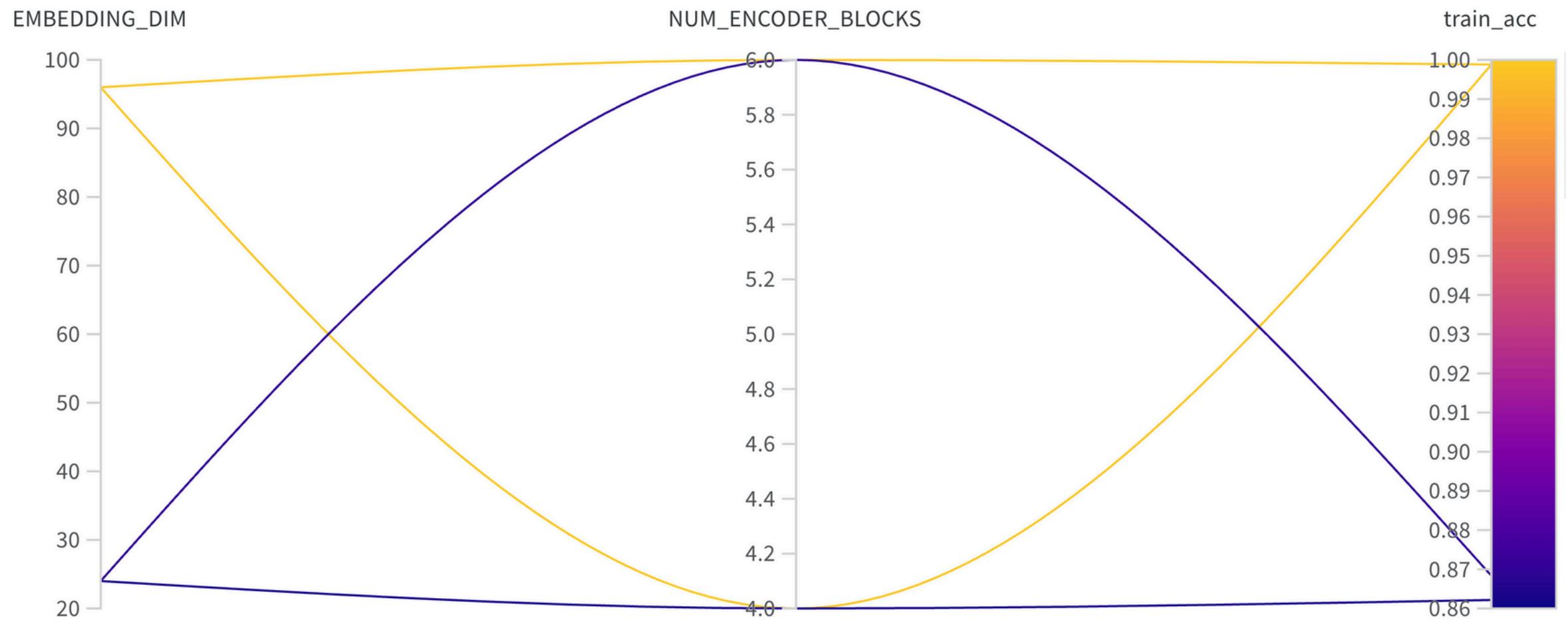
Team 2. "Wait, 'highest' must mean flag must go on the tallest tower, not just any tower"

Team 3. "Ok, so we find the tallest tower compared to all the towers we've built. and put the flag on top."

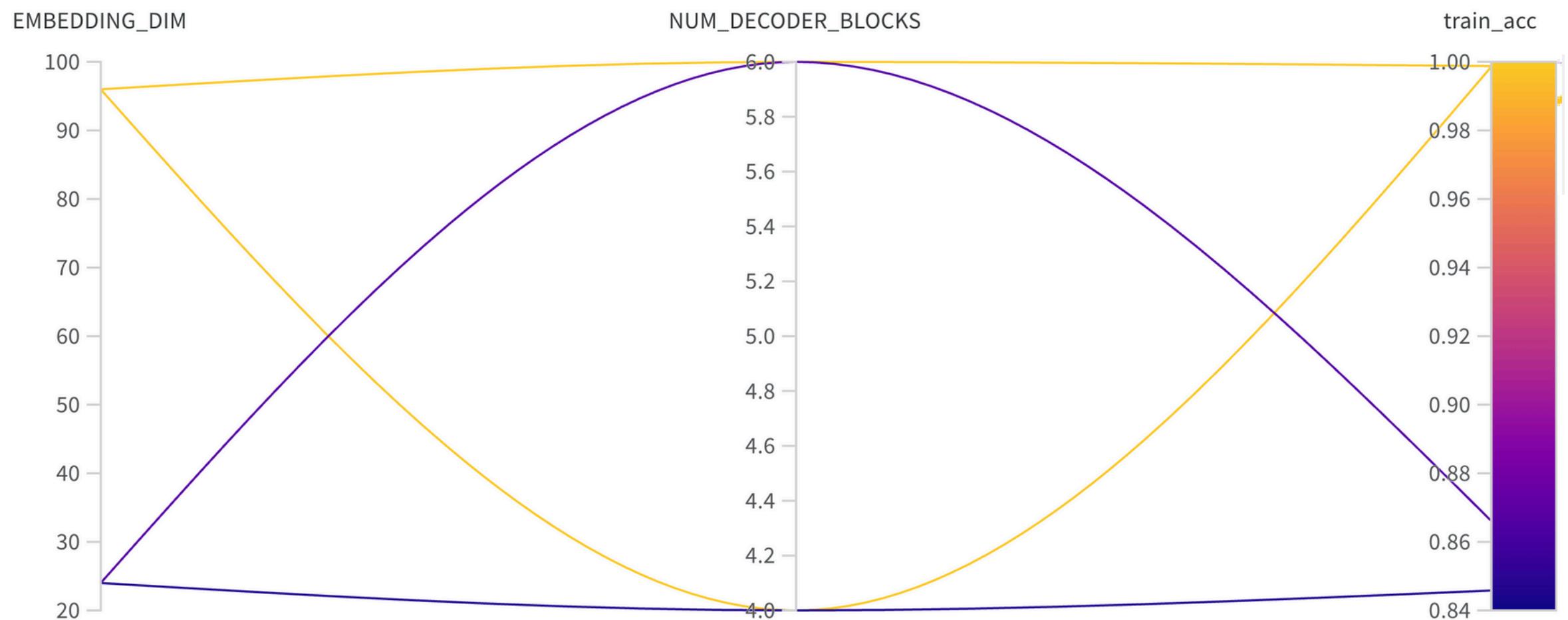
=> Ok now we have truly **understand the instruction**, and the final team knows exactly what to do to execute the instruction 🔧

# DEMO TIME





Optimisation  
Sweeps

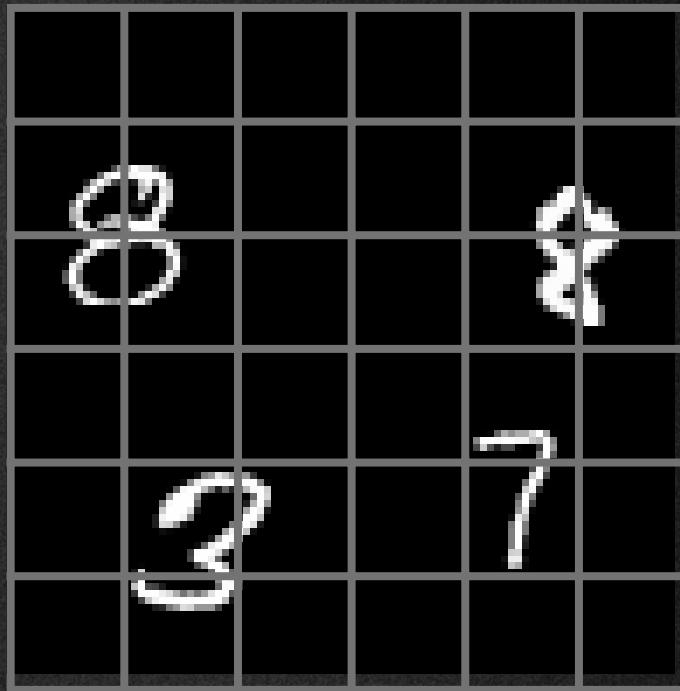




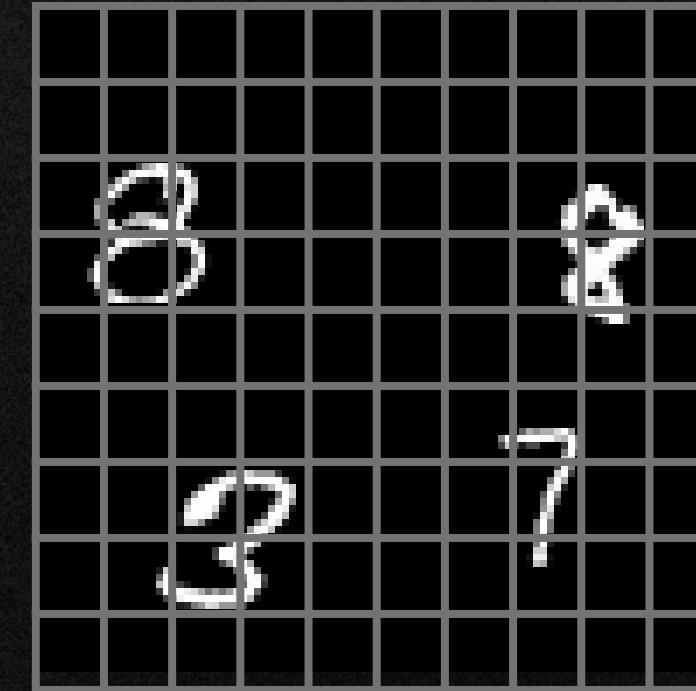
# GENERAL LEARNINGS & REFLECTIONS

- Patches
- Loss function adventures
- More Loss function adventures
- Masking
- Zeros and Sparsity

# PATCHES

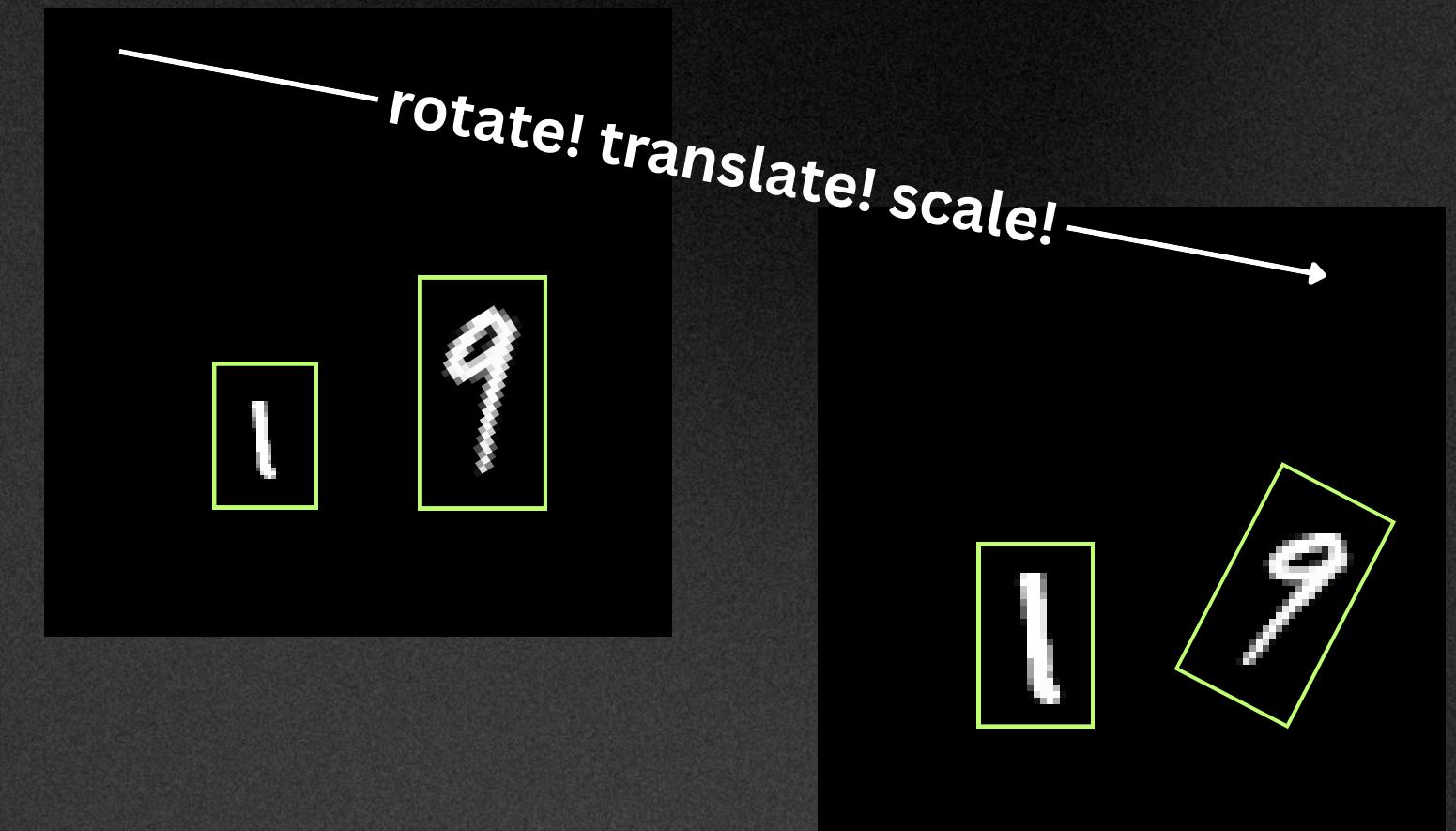


smaller patches → efficiency



smaller patches → accuracy

# AUGMENTATION



**Token:**

[<sos>, x, x, <eos>, <pad>, <pad>]

# LOSS ADVENTURES PT. 1

**Target:** [10, 0, 0, 11, 12, 12]

**Predicted:** [10, 1, 0, 11, 11, 12]

**Not ignoring <pad>,  
no custom class weighting**

Penalized for predicting <pad> tokens  
→ **distraction from digits**  
→ **misleading loss**



**Ignoring <pad>,  
no custom class weighting**

9 digits and <eos> treated equally,  
except <pad>



**Ignoring <pad>,  
yes custom class weighting**

- Class 9 gets a 2x loss penalty
- <eos> and <pad> get 0.5x penalty

→ **did not learn when to stop**



# LOSS ADVENTURES

## PT. 2



We had some unfair loss where: **731 != 137**

I thought any loss function you could write would work, however it's anything that pytorch detects as differentiable, so some attempts at custom loss failed

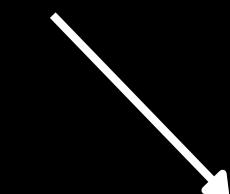
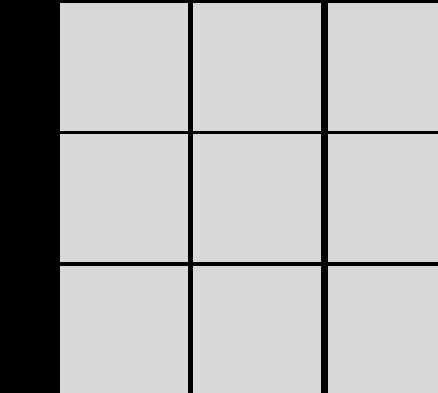
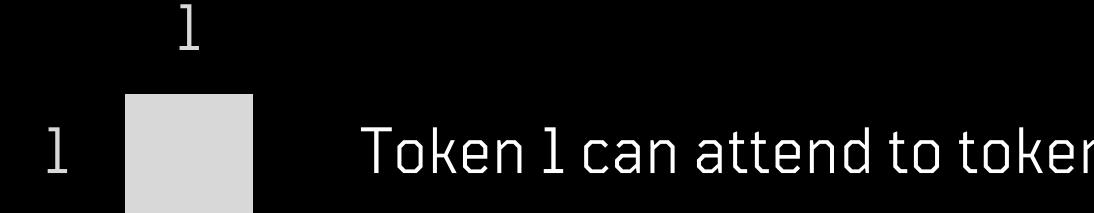
I converted the logits back into a probability distribution with softmax to lose the concept of sequence

It worked but wasn't any sort of obvious improvement (and may cause issues with Start and End tokens)- But could be useful elsewhere.

# TRIANGLE HOW?

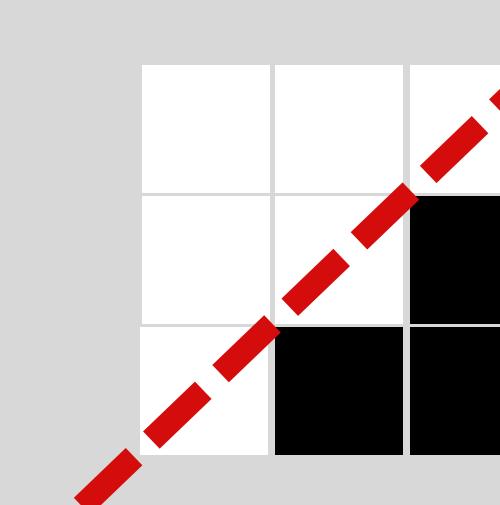


Past X Past looks like a square



Wait this aint a triangle?

Or if you draw a line here...

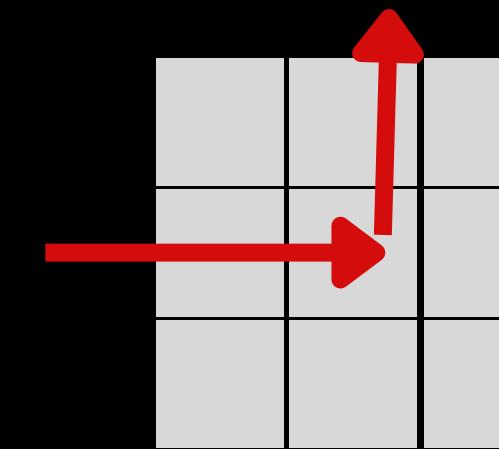


...You get a triangle, that sweeps right, y no sweep?

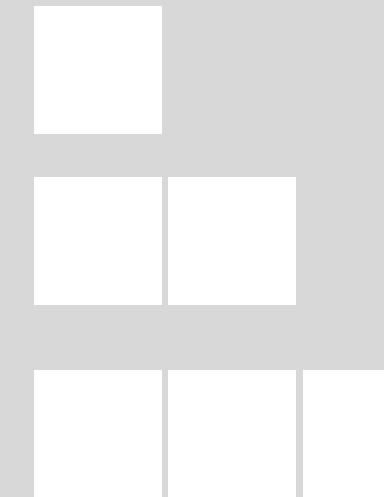
# TRIANGLE HOW?



**It's only valid to read the triangle matrix left then up.**



**It's actually this**



**Console has small vertical spacing**

# ZEROS AND SPARSITY

- “zeroing out” matrixes is a natural term but beware if your tensors went from -1 to +1 and actual 0 is 50%
- +inf and -inf if building yourself or be very careful normalizing before/after
- These elements in Pytorch can accept variable lengths:
  - nn.Transformer
  - nn.TransformerEncoder/ nn.TransformerDecoder
  - nn.MultiheadAttention
  - nn.Embedding
  - torch.nn.utils.rnn.pack\_padded\_sequence / pad\_packed\_sequence



# ZEROS AND SPARSITY

- Under the hood, variable length seems to be just fixed length with padding
- You don't get to save compute by making out stuff, unless you have a "Architecturally sparse" models
- A tensor with 1000000 inputs, all but one masked, would still compute things for the remaining 999999
- In models like MLP 0 or -inf is as valid an input as any other and is processed with everything else, potentially influencing things with connected layers.

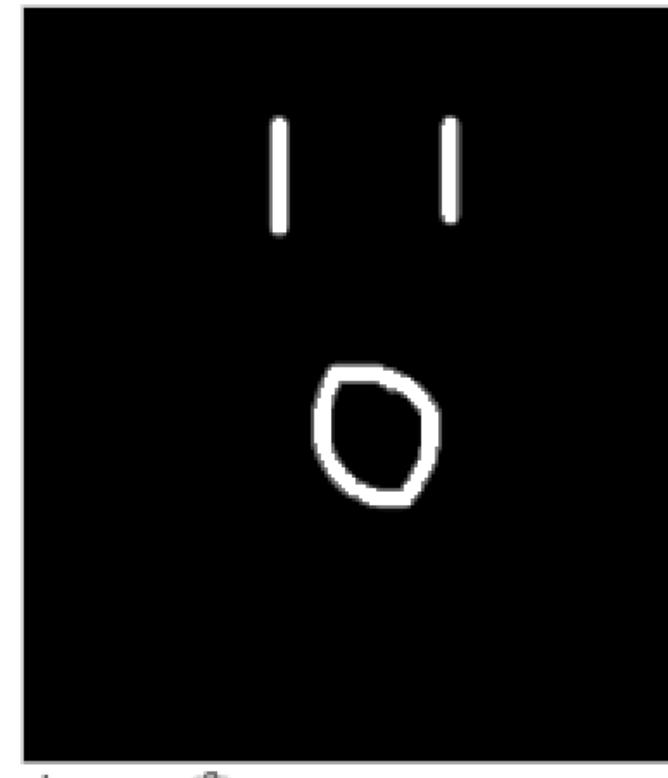


# Digit Transformer ✨

A mini transformer model for recognizing sequences of handwritten digits.

Draw a multi-digit  
image

Classify



Prediction

Raw tokens: [1, 0, 1, 11, 12, 12]

Digits: 1 0 1



# Digit Transformer ✨

A mini transformer model for recognizing sequences of handwritten digits.



ion

[1, 0, 1, 11, 12, 12]

: 1 0 1

