

28 November 2025

TRANSLATE ENGLISH TO KOREAN WITH TRANSFORMER

Minh Le



THE UNIVERSITY
of ADELAIDE

CONTENTS

01. Motivation

02. Transformer Architecture

03. Implementation

04. Conclusion

05. References

1. MOTIVATION



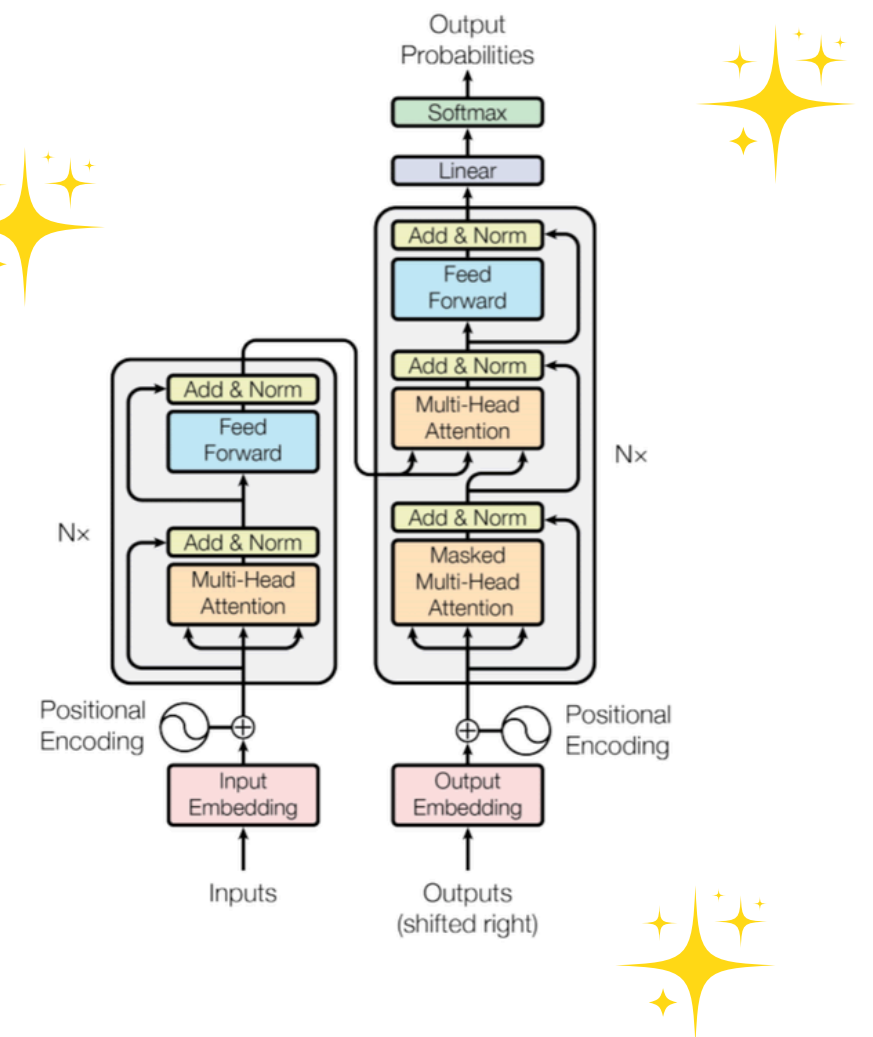
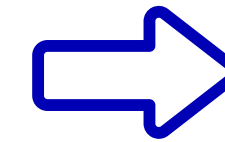
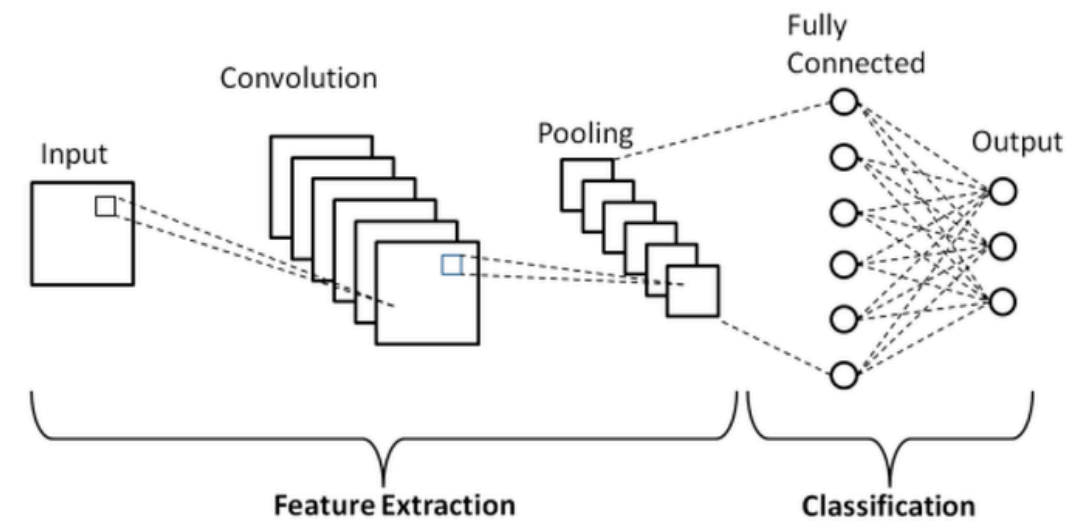
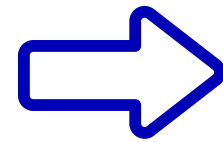
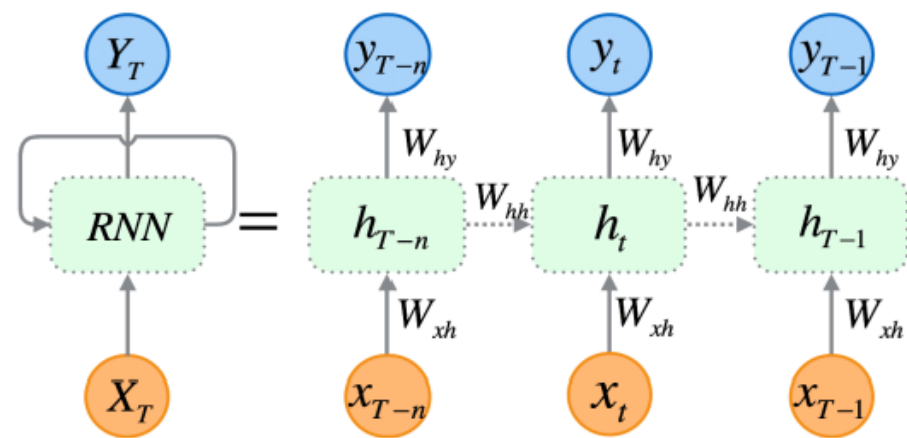
Hello, we are TWICE



안녕하세요, 저희는 트와이스입니다

1. MOTIVATION: WHY TRANSFORMERS?

Attention Is All You Need (2017).



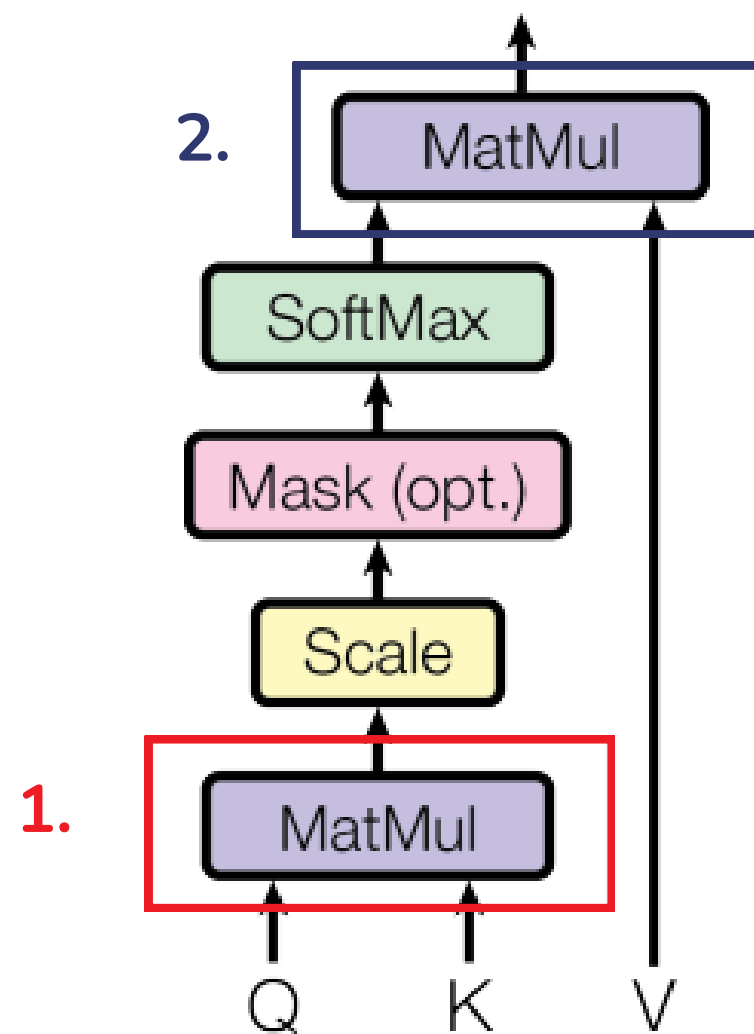
- ✗ Prone to **vanishing gradient**.
- ✗ Sequential computation **inhibits parallelization**.

- ✗ Requires many layers to capture long-term dependencies
→ **Large network**

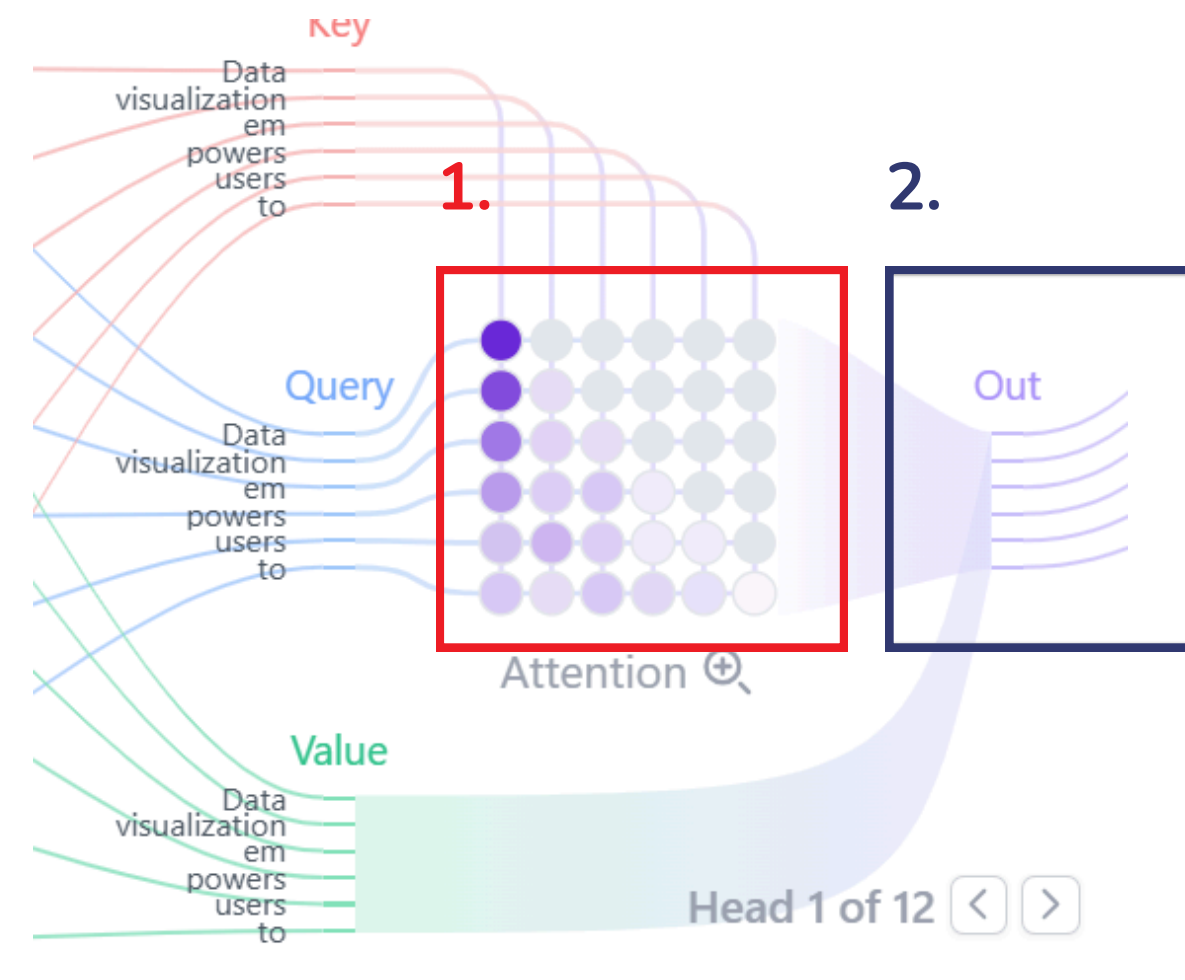
- ✓ Self-attention allows **global context**.
- ✓ All tokens go through the model simultaneously (**parallelization**)

2. TRANSFORMER ARCHITECTURE: SELF-ATTENTION

1. Compare query and key \rightarrow Attention scores.
2. Apply scores to value matrix \rightarrow Each token become context-aware.

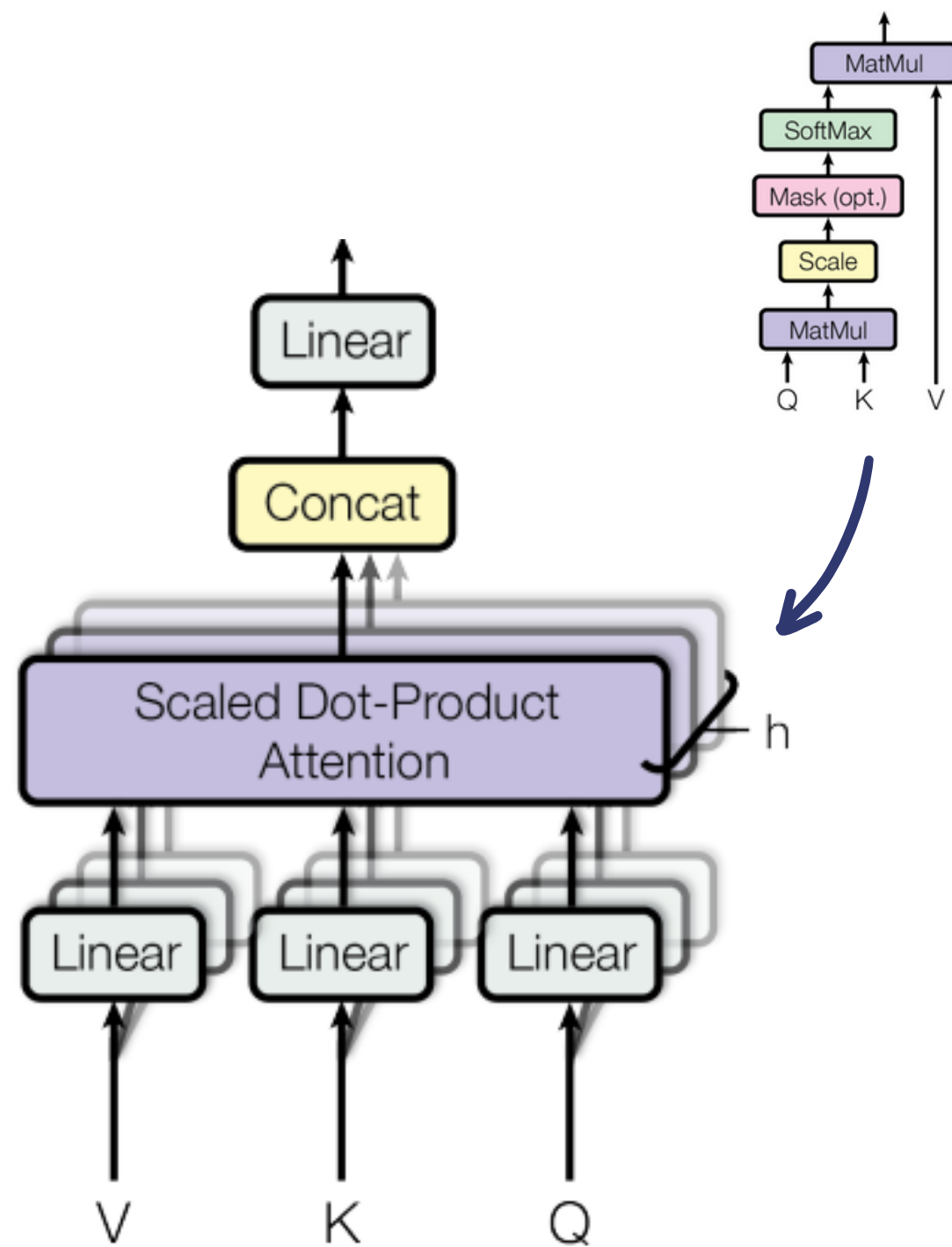


Attention Is All You Need (2017)

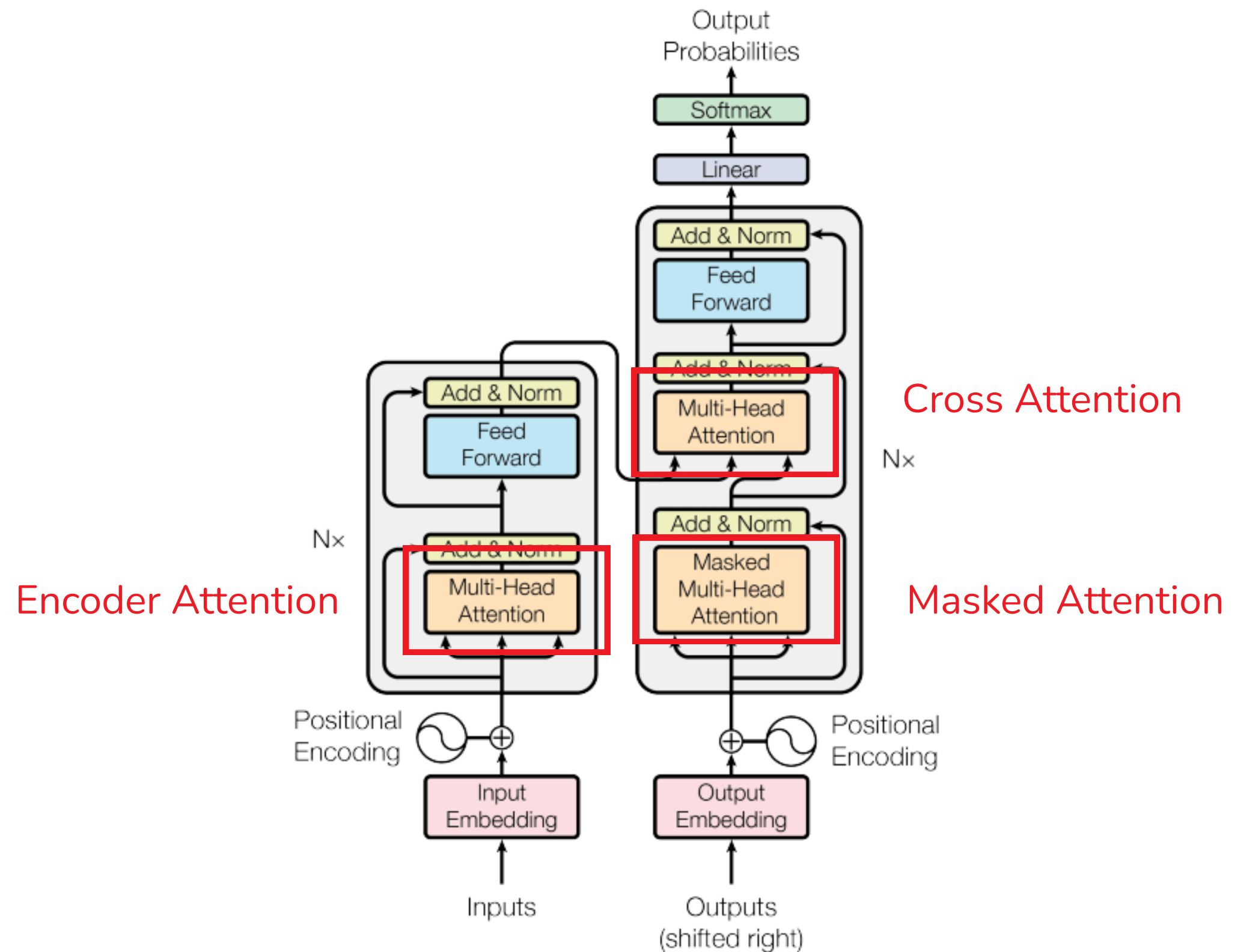


<https://poloclub.github.io/transformer-explainer/>

2. TRANSFORMER ARCHITECTURE: MULTI-HEAD ATTENTION



Multi-head Attention



3. IMPLEMENTATION

- I built the whole transformer **from scratch** using PyTorch, following the paper and other's projects on GitHub.



```
Attention
Generate Code Markdown

# nn.MultiheadAttention can be used instead but for learning purposes, I built one from scratch
# https://pytorch.org/docs/stable/generated/torch.nn.MultiheadAttention.html

class MultiHeadAttention(nn.Module):
    def __init__(self, d_model, num_heads, drop_out=0.1):
        """
        d_model: the dimension of the embedding vector for EACH token
        num_heads: the number of attention heads
        drop_out: the dropout rate
        """
        super().__init__()
        assert d_model % num_heads == 0, "MultiHeadAttention.__init__(): d_model must be divisible by num_heads"
        self.d_model = d_model
        self.num_heads = num_heads
        self.head_dim = d_model // num_heads

        self.q_linear = nn.Linear(d_model, d_model)
        self.v_linear = nn.Linear(d_model, d_model)
        self.k_linear = nn.Linear(d_model, d_model)

        # Final linear projection after concatenating heads (Learns how to merge all heads back into a single representation)
        self.out_proj = nn.Linear(d_model, d_model)

        self.dropout = nn.Dropout(drop_out)
```

Work is still in progress...

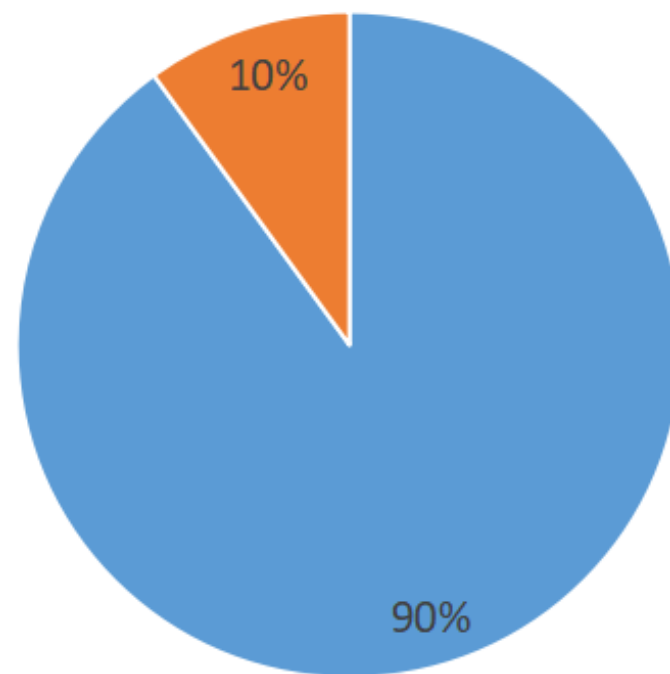
3. IMPLEMENTATION: TOKENIZER



<https://github.com/google/sentencepiece>

Dataset: 100,000 sequences ([AI Hub](#))

Train - Validation Data Ratio



■ Train ■ Validation

- Sequences are **reversible** (white spaces are treated as symbols).

```
Hello_World.
```

Then, this text is segmented into small pieces, for example:

```
[Hello] [_Wor] [ld] [.]
```

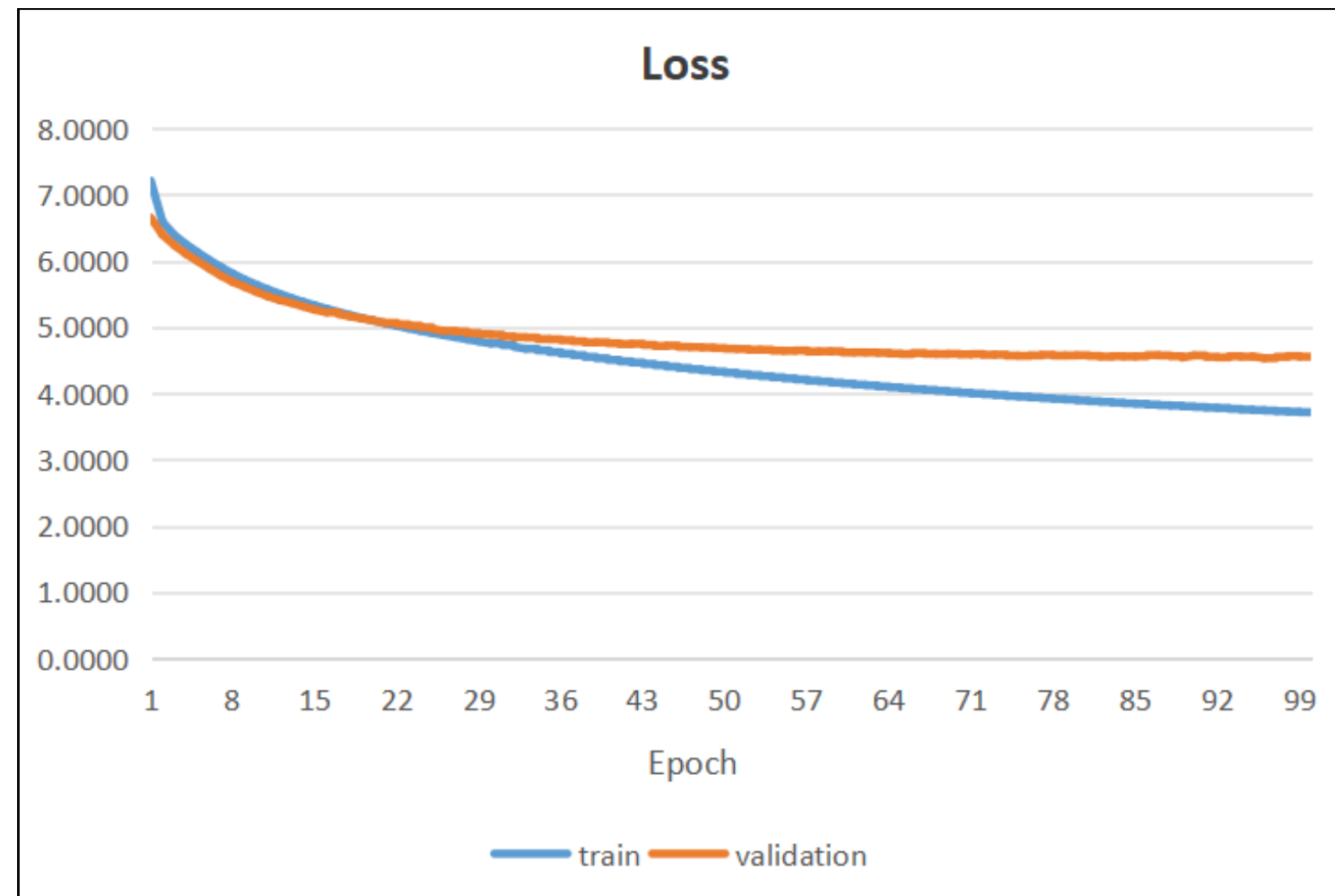
- sentencepiece*'s unigram model is **superior to morphology**.
 - E.g. The word **obviously** can be tokenized into **_obvious** + **ly**

```
3963  _obvious  -11.3057
```

```
243   ly      -7.58103
```

This is an actual example of using *sentencepiece* on English

3. IMPLEMENTATION: RESULTS



Loss curve plateaus for the majority of the training process, suggesting that the model is **underfitting**.

- My BLEU score (ENG-KOR): **11.35**
- OPUS-MT project's BLEU score:
 - ENG-KOR: 13.3
 - KOR-ENG: 41.3

Translation:

ENG: I am going home tomorrow

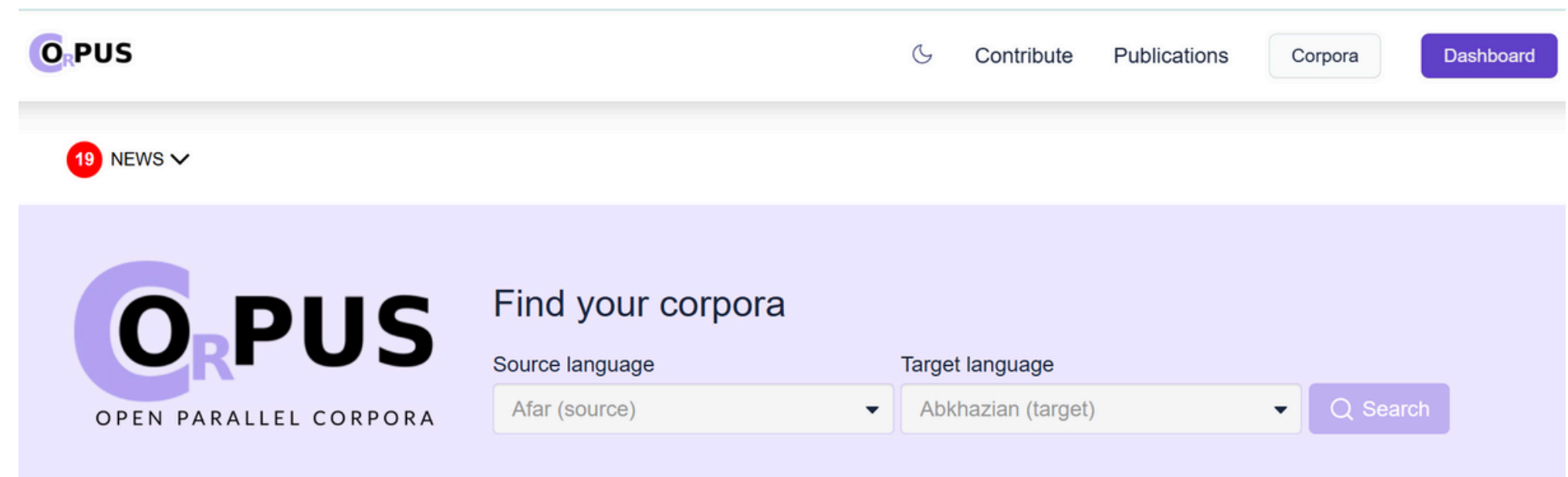
KOR: 나는 내일 집에 가는 당신과 있을 것이에요.

(I will be with you on your way home tomorrow.)

4. CONCLUSIONS: FUTURE WORK

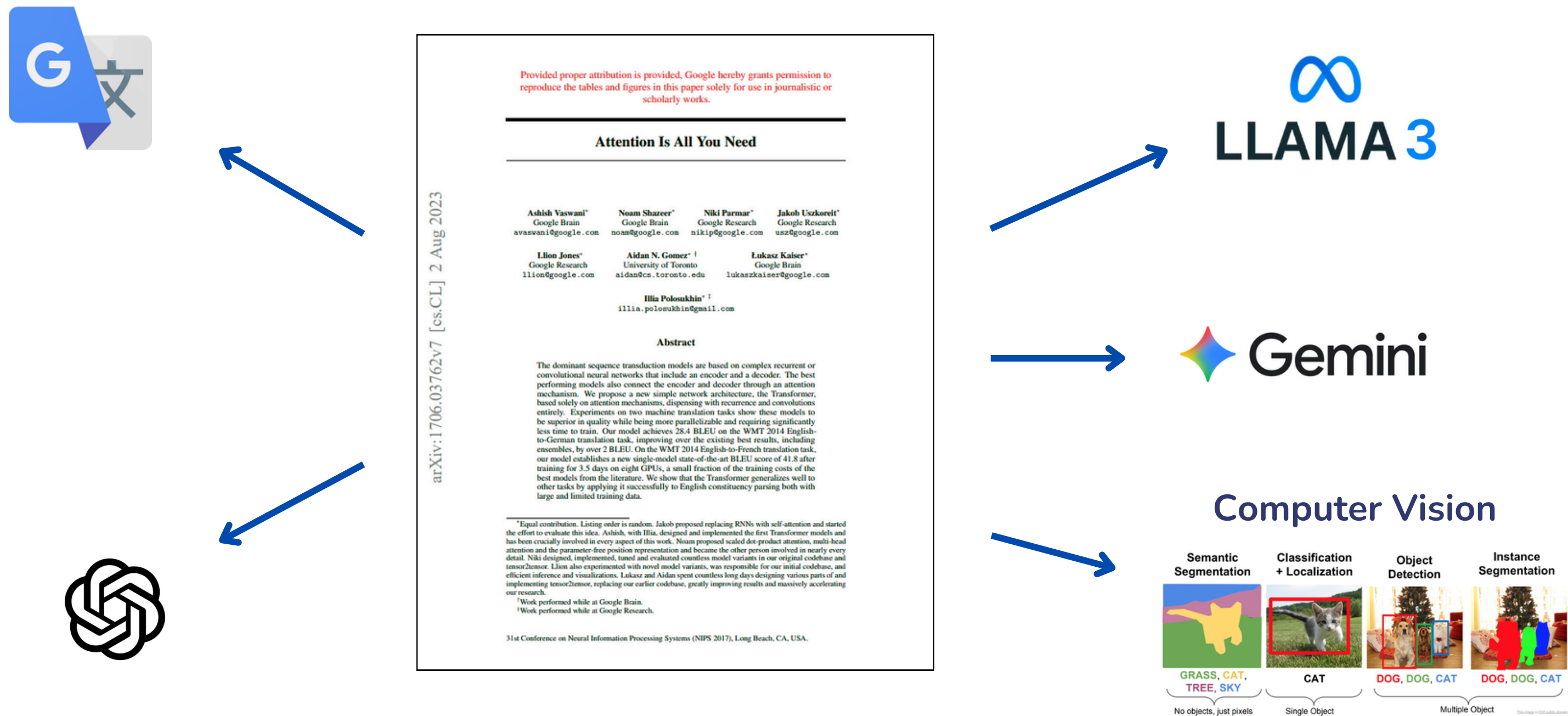
- Apply **Beam Search**.
- **Larger dataset**.
- Fine-tune hyperparameters (embedding size, increase attention heads).

OPUS corpora: <https://opus.nlpl.eu>



The screenshot shows the OPUS website interface. At the top, there is a navigation bar with the OPUS logo on the left and links for 'Contribute', 'Publications', 'Corpora', and 'Dashboard' on the right. Below the navigation bar, there is a section with a red circle containing the number '19' and the text 'NEWS' with a dropdown arrow. The main content area has a light purple background. On the left, there is a large 'OPUS' logo with 'OPEN PARALLEL CORPORA' written below it. To the right of the logo, the text 'Find your corpora' is displayed. Below this text, there are two dropdown menus: 'Source language' with 'Afar (source)' selected, and 'Target language' with 'Abkhazian (target)' selected. To the right of these dropdowns is a purple button with a magnifying glass icon and the text 'Search'.

4. CONCLUSIONS: INFLUENCE OF TRANSFORMER



Attention Is All You Need (2017)

<https://arxiv.org/pdf/1706.03762>

5. REFERENCES

- Attention Is All You Need (2017) <https://arxiv.org/abs/1706.03762>
- Sentencepiece <https://github.com/google/sentencepiece>
- SamLynnEvans/Transformer <https://github.com/SamLynnEvans/Transformer>
- Hufon/pytorch-transformer-kor-eng <https://github.com/Hufon/pytorch-transformer-kor-eng>
- CLIP: Connecting text and images <https://openai.com/index/clip/>
- Transformer Explainer <https://poloclub.github.io/transformer-explainer/>
- OPUS corpora <https://opus.nlpl.eu>

THANK YOU!