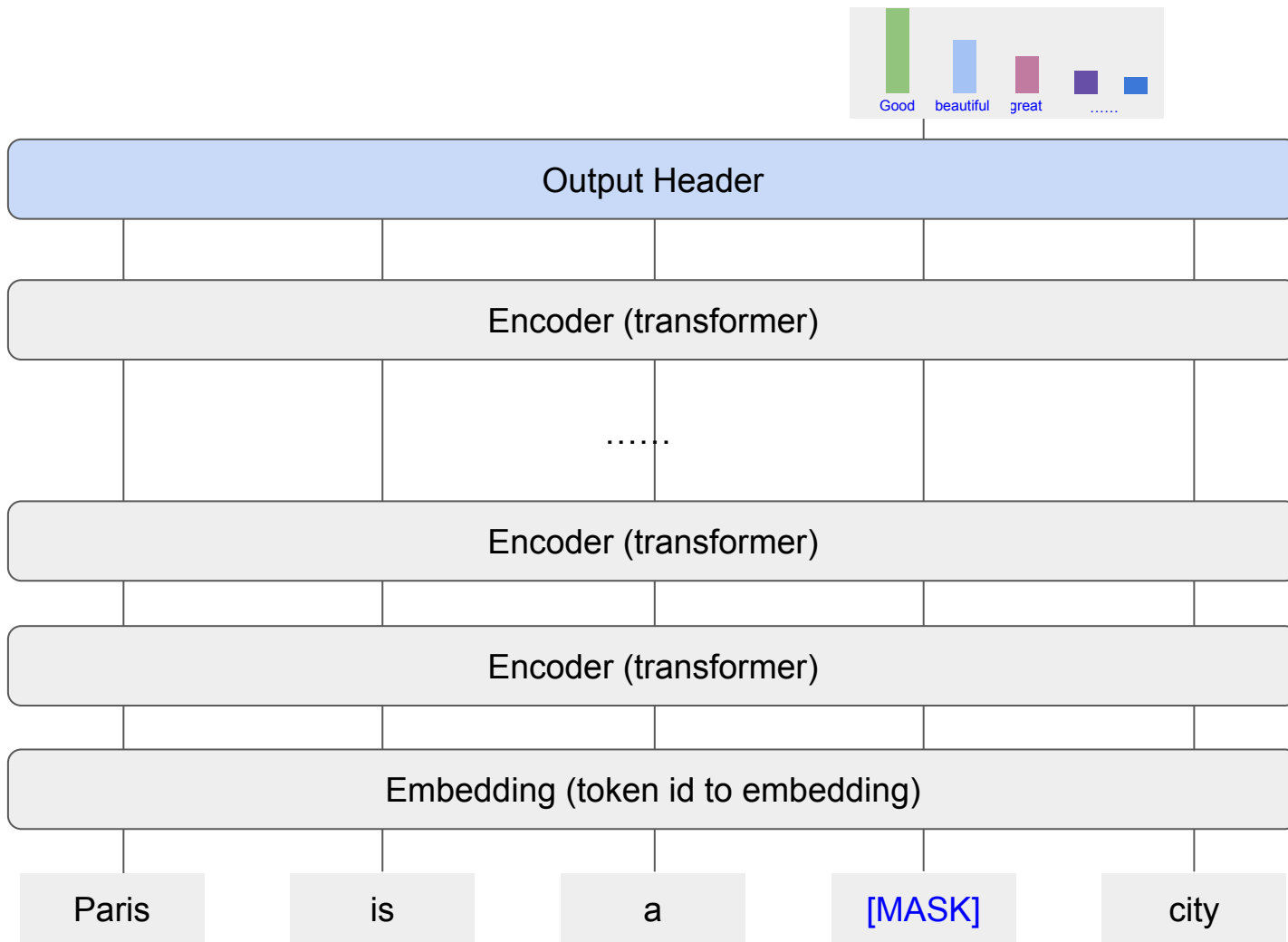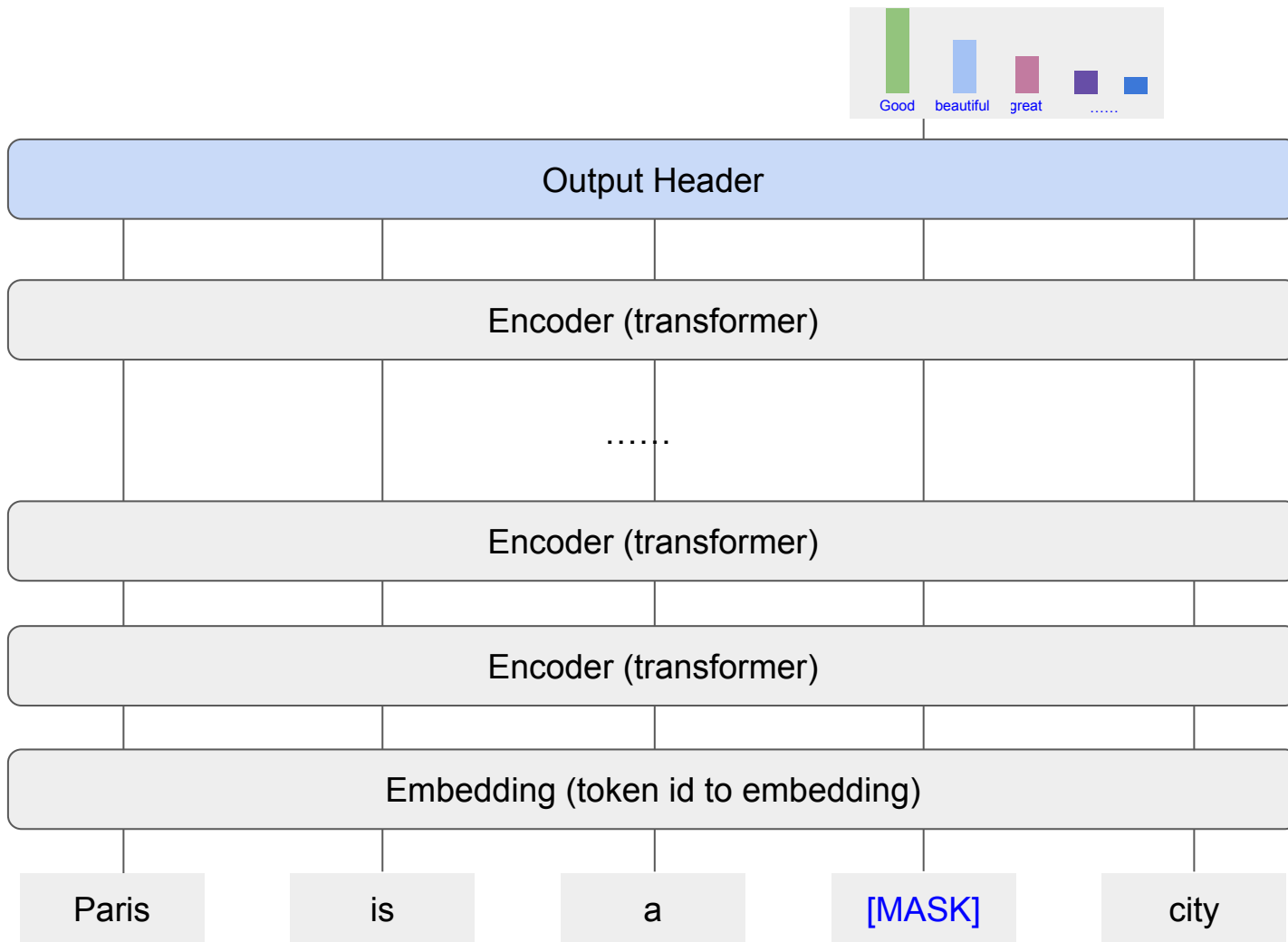BERT Masked Token Prediction

BERT Masked Token Prediction

Output Header

Output Size: vocab_size
(each word in vocab has a
probability, sum is 100%)

Encoder (transformer)

......

Encoder (transformer)

Encoder (transformer)

Output Size:
sequence_len * hidden_size

Embedding (token id to embedding)

Output Size:
sequence_len * hidden_size

Paris    is    a    [MASK]    city

Output Size:  sequence_len

Good    beautiful    great    ......

**BERT Embedding**

Embedding = pooling (*e.g average on sequence_len dimension*) of transformer output

Embedding size: *(sequence_len * hidden_size) / sequence_len = hidden_size*

Encoder (transformer)

……

Encoder (transformer)

Encoder (transformer)

Output Size:
sequence_len * hidden_size

Embedding (token id to embedding)

Output Size:
sequence_len * hidden_size

| Paris | is | a | [MASK] | city |

Output Size: sequence_len

GPT Next Token Prediction

Paris    is    a

Good  beautiful  great  ......

Output Header

Decoder (transformer)

......

Decoder  (transformer)

Decoder (transformer)

Embedding (token id to embedding)

[START]    Paris    is    a

| Paris | is | a | | | |
|-------|-----|-----|--|--|--|


Good   beautiful   great   ......

**GPT Next Token Prediction**

## Output Header

Output Size: vocab_size
(each word in vocab has a
probability, sum is 100%)

## Decoder (transformer)

......

## Decoder  (transformer)

Output Size:
sequence_len * hidden_size

## Decoder (transformer)

## Embedding (token id to embedding)

Output Size:
sequence_len * hidden_size

| [START] | Paris | is | a |
|---------|-------|-----|---|

Output Size:  sequence_len

**GPT Embedding**

Embedding = pooling (*e.g average on sequence_len dimension*) of transformer output

Embedding size: *(sequence_len * hidden_size) / sequence_len = hidden_size*

Decoder (transformer)

……

Decoder  (transformer)

Decoder (transformer)

Embedding (token id to embedding)

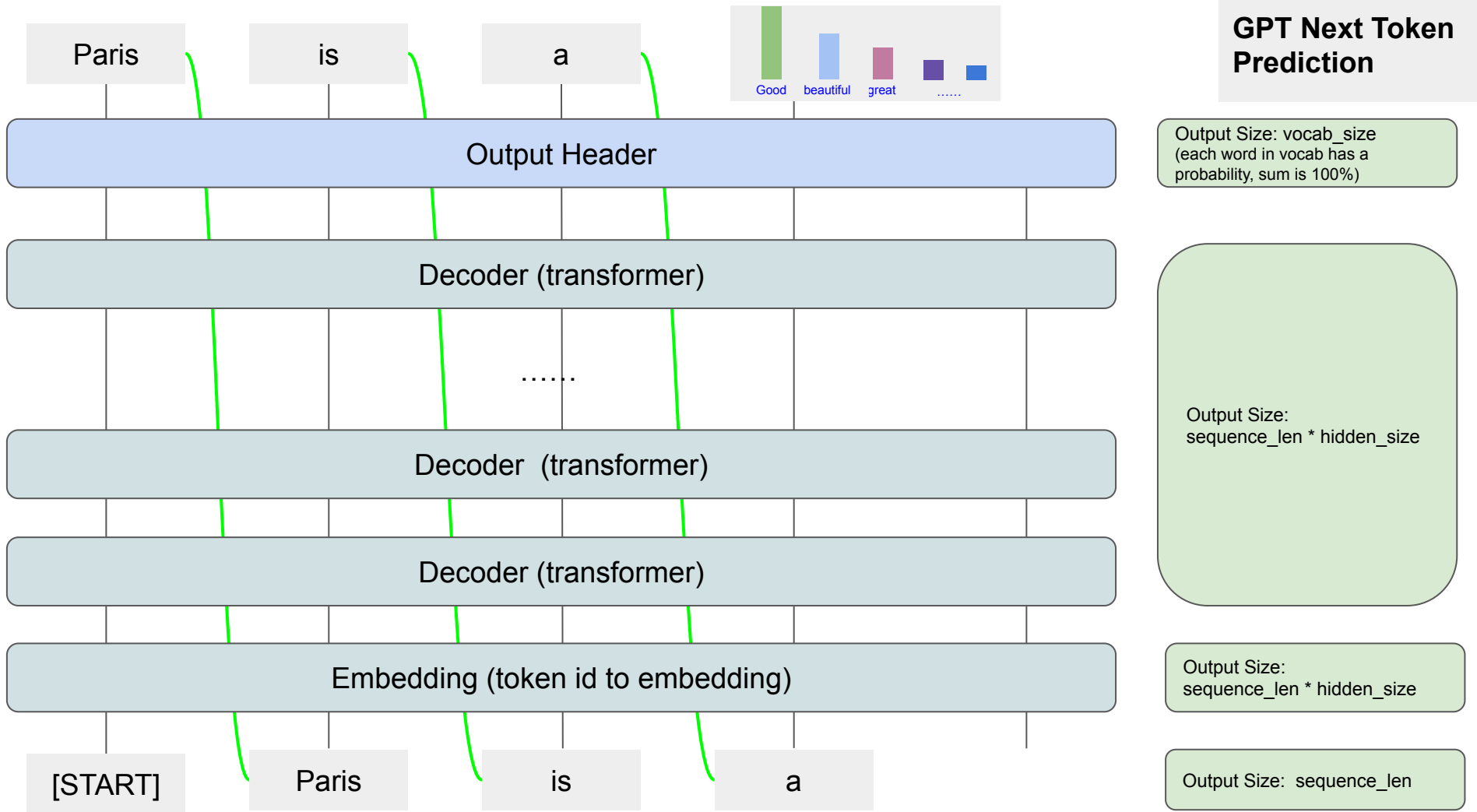[START]     Paris     is     a

Output Size:
sequence_len * hidden_size

Output Size:
sequence_len * hidden_size

Output Size:  sequence_len