

Sampling Strategies for Generative Models; Transformer Architecture Continued; Pretrained Models; BERT

1. Some Strategies for Sampling in Generative Models

- a. In Generative Models, “Sampling” refers to the choice of the next token using the probability distribution output by the network.
- b. Let’s look at a really simple example:
 - i. Suppose the generative model is creating sentences from the vocabulary $V = [\text{'<s>'}, \text{'</s>'}, \text{'NLP'}, \text{'early'}, \text{'fun'}, \text{'green'}, \text{'hard'}, \text{'interesting'}, \text{'is'}]$
 - ii. We have generated three tokens so far: [$\text{'<s>'}, \text{'NLP'}, \text{'is'}$] and the softmax output (probability distribution) for the next token is:

<u>Token</u>	<u>Probability</u>
'<s>'	0.0
'</s>'	0.01
'NLP'	0.03
'early'	0.12
'fun'	0.12
'green'	0.0
'hard'	0.2
'interesting'	0.51
'is'	0.01

- c. There are five common strategies for sampling [$\text{'<s>'}, \text{'NLP'}, \text{'is'}, ???$]
 - i. Sampling using the unmodified distribution

```
1 print('\nsample: ',sample_choice(outcomes,distribution,option=0))
<s>      0.0
</s>      0.01
NLP       0.03
early     0.12
fun       0.12
green     0.0
hard      0.2
interesting 0.51
is        0.01

sample: fun
```

- ii. Greedy Sampling: Just choose the single most probable

```
1 print('\nsample: ',sample_choice(outcomes,distribution,option=1))
<s>      0.0
</s>      0.0
NLP       0.0
early     0.0
fun       0.0
green     0.0
hard      0.0
interesting 1.0
is        0.0

sample: interesting
```

iii. Softmax with Temperature

```

1 print('\nsample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.01))
<s>          7.095474162284459e-23
</s>          1.928749847963851e-22
NLP          1.4251640827408859e-21
early        1.1548224173015387e-17
fun          1.1548224173015387e-17
green        7.095474162284459e-23
hard          3.4424771084698575e-14
interesting   0.99999999999999655
is           1.928749847963851e-22

sample: interesting

: 1 print('\nsample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.1))
<s>          0.005446284988308871
</s>          0.0060190757806309345
NLP          0.0073517157600377
early        0.018082302955730302
fun          0.018082302955730302
green        0.005446284988308871
hard          0.040242905309378116
interesting   0.893310051481244
is           0.0060190757806309345

sample: interesting

1 print('\nsample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.5))
<s>          0.08396476946896886
</s>          0.08566097032727966
NLP          0.08915686084440715
early        0.10674014184435929
fun          0.10674014184435929
green        0.08396476946896886
hard          0.12526071682556345
interesting   0.23285065904881377
is           0.08566097032727966

sample: </s>

1 print('\nsample: ',sample_choice(outcomes,distribution,option=2,temperature = 10.0))
<s>          0.10986989484510404
</s>          0.10997981969321279
NLP          0.11019999943895172
early        0.11119627595335972
fun          0.11119627595335972
green        0.10986989484510404
hard          0.11208941394957388
interesting   0.11561860562812123
is           0.10997981969321279

sample: green

```

iv. Top-K Sampling: From top K most probable, choose proportionally:

```

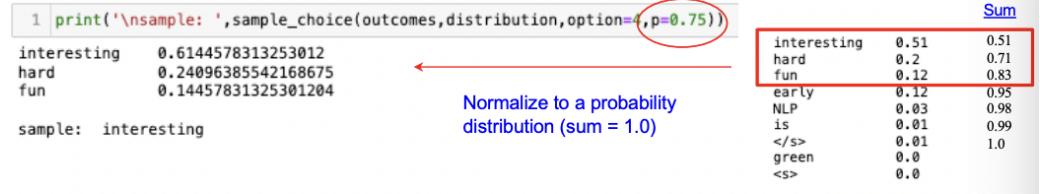
1 print('\nsample: ',sample_choice(outcomes,distribution,option=3,K=4))
interesting   0.5368421052631579
hard          0.2105263157894737
fun           0.12631578947368421
early         0.12631578947368421
sample: fun

```

Normalize to a probability distribution (sum = 1.0)

interesting	0.51
hard	0.2
fun	0.12
early	0.12
NLP	0.03
is	0.01
</s>	0.01
green	0.0
<s>	0.0

v. Top-P: From top choices with sum of probability p, choose proportionally:



2. Some Refinements for Sampling in Generative Models

```

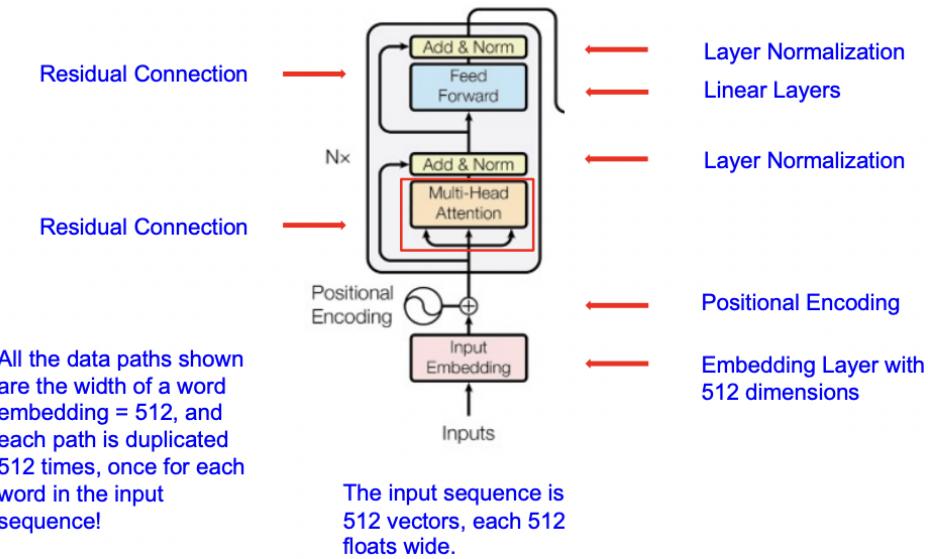
1 # Options
2
3 #   0 = No modification (default)
4 #   1 = Greedy: return one-hot of maximum prob, equivalent to softmax with 0 temperature
5 #   2 = Softmax with temperature (temperature of 1.0 is ordinary softmax)
6 #   3 = Top-K: Choose from top K most probable options (no other modification)
7 #   4 = Nucleus (Top-p): p is a cutoff, choose from the top choices whose cumulative prob just
8 #                     exceeds p (so in general, cumulative prob is slightly more than p)
9
10
11 import math
12 import numpy as np
13 from numpy.random import choice
14
15 def sample_choice(outcomes,distribution,option=3,temperature=0.3,K=5,p=0.25):
16
17     if type(outcomes) != np.ndarray:
18         outcomes = np.array(outcomes)
19     if type(distribution) != np.ndarray:
20         distribution = np.array(distribution)
21
22
23     if option == 0:
24         return choice(a=outcomes,p=distribution/sum(distribution))
25     elif option == 1:
26         oneHot = np.array([0.0]*len(distribution))
27         oneHot[np.argmax(distribution)] = 1.0
28         return outcomes[np.argmax(distribution)]
29     elif option == 2:
30         sum_exp = sum(math.exp(x/temperature) for x in distribution)
31         lst = [math.exp(x/temperature)/sum_exp for x in distribution]
32         return choice(a=outcomes,p=[math.exp(x/temperature)/sum_exp for x in distribution])
33     elif option == 3:
34         indices = np.argsort(distribution)[-K:]
35         return choice(a=outcomes[indices],p=distribution[indices]/sum(distribution[indices]))
36     elif option == 4:
37         indices = np.argsort(distribution)
38         total = 0
39         for C in range(len(distribution)-1,-1,-1):
40             total += distribution[indices[C]]
41             if total > p:
42                 break
43         indices = indices[C:]
44         return choice(a=outcomes[indices],p=distribution[indices]/sum(distribution[indices]))
45

```

a.

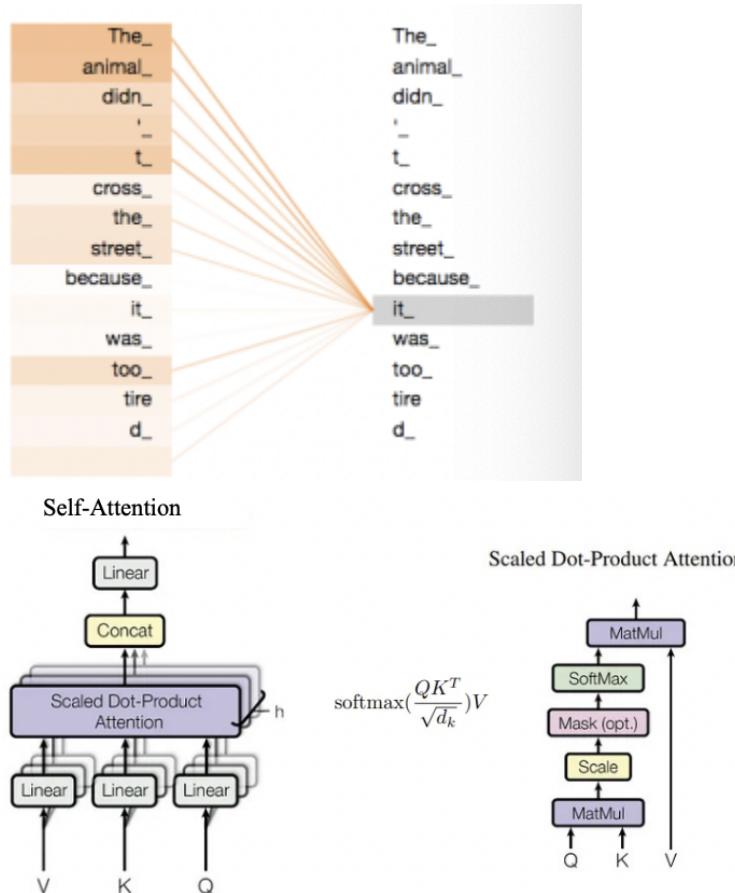
3. Transformer: Multi-Head Attention

- a. The center of the design is the Attention mechanism, here Multi-Head Attention....

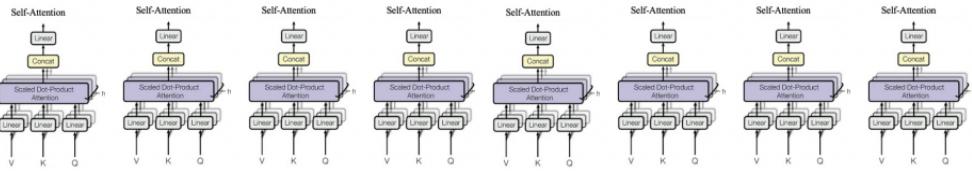


b.

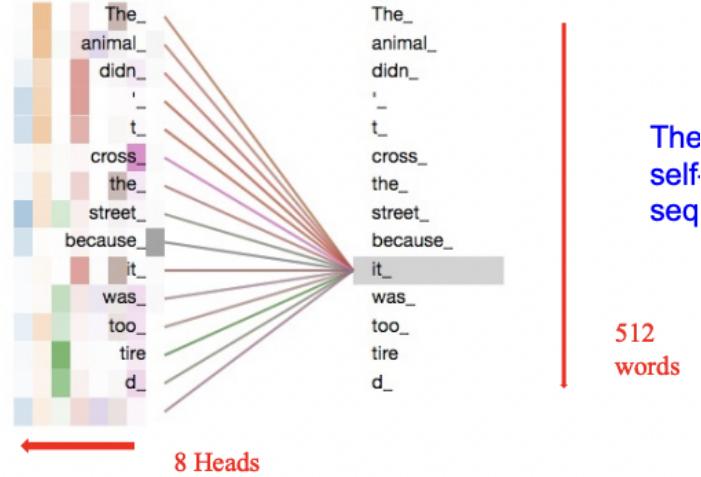
- c. The basic mechanism at work here is self-attention, where the input is processed to determine the dependencies between different words.
- d. Self-attention is implemented by a series of linear transformations, scaled, and then softmaxed to produce the probability distribution which tells us how much each word depends on other words in the same sequence.



- e. Multi-Head Attention applies this self-attention mechanism as 8 (or more) self-attention Heads:



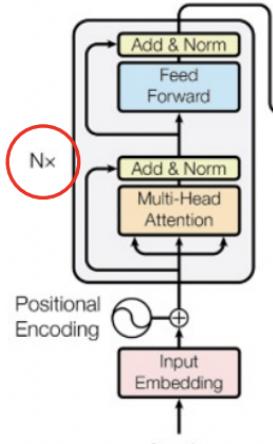
- f. This allows the encoder to try to understand 8 different kinds of dependencies among the words in the input.



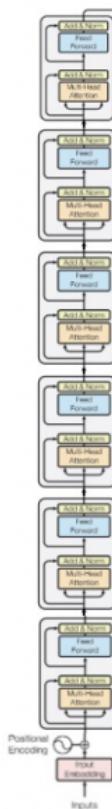
- i. The original design applied 8 self-attention heads to sequences of length 512.

4. Transformer: Encoder

- a. Then, this encoder layer is stacked 6, 8, or more layers deep!



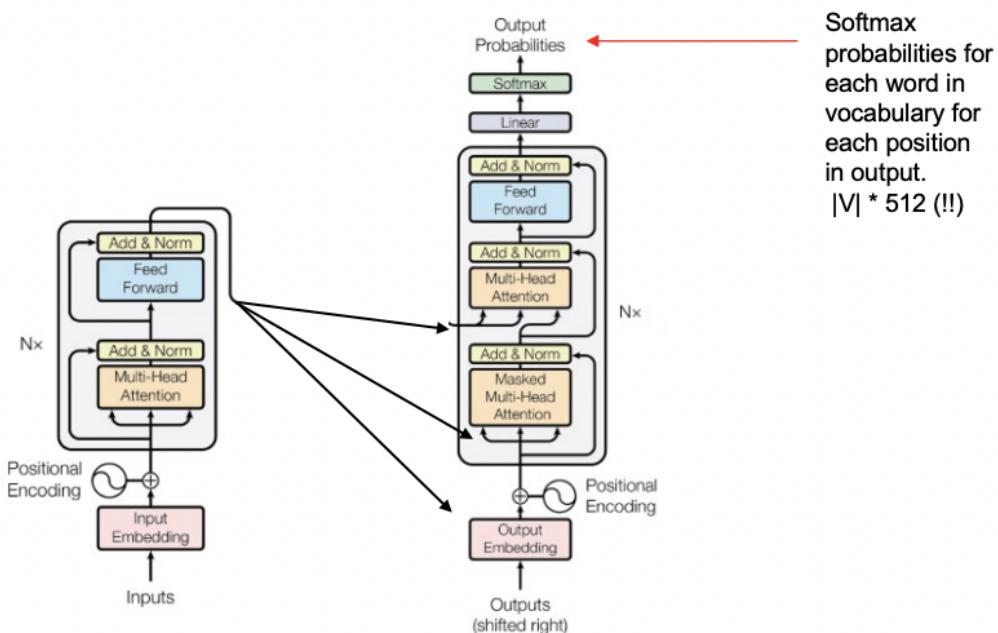
b.



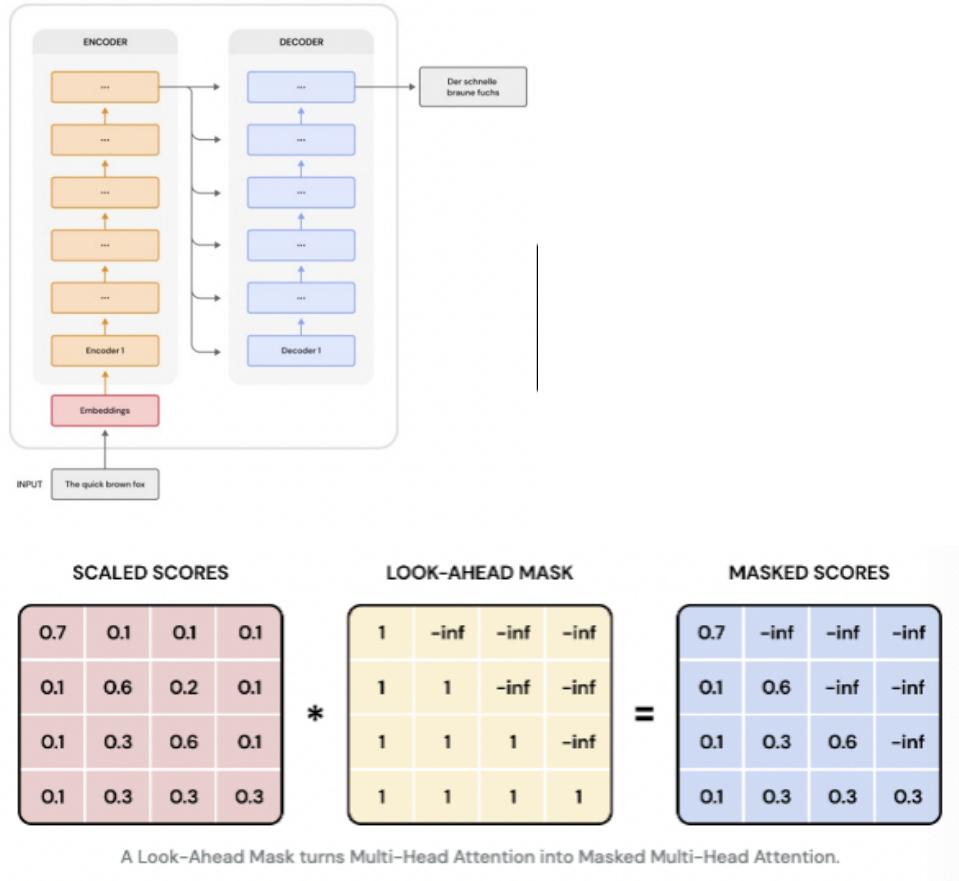
c.

5. Transformer: Decoder

- a. The Decoder stage has essentially the same features as the Encoder, except for a stage called Masked Multi-Head Attention.
- b. Apparently, the classic diagram (from the original paper) is a bit misleading, and there are additional connections between the encoder and decoder:



- c. The Decoder stage has essentially the same features as the Encoder, except for Multi-Headed Attention is masked to simulate left to right generation of output; Each stage of Decoder takes the context vector (output of encoder) as an additional input:

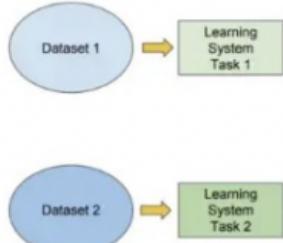


6. Transfer Learning: Leveraging pre-trained models

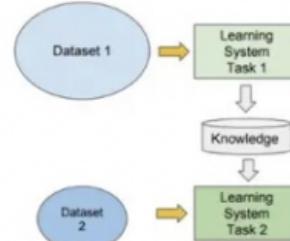
- a. A hugely important development in ML in the last decade is the development of transfer learning systems, consisting of two stages:
 - i. Pre-Training: A (large) model is trained on a (large) corpus to produce generic outputs;
 - ii. Downstream: A NN is added on top of the pre-trained model to build models for a different task (called a "downstream task")

Traditional ML vs Transfer Learning

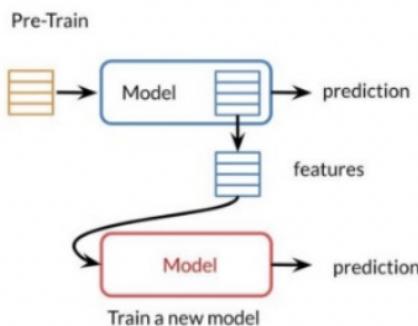
- Isolated, single task learning:
 - Knowledge is not retained or accumulated. Learning is performed w/o considering past learned knowledge in other tasks



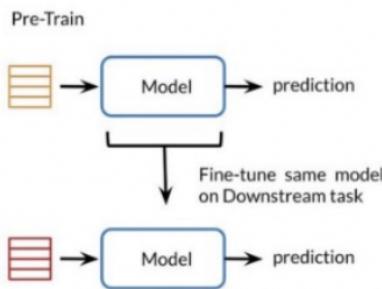
- Learning of a new tasks relies on the previous learned tasks:
 - Learning process can be faster, more accurate and/or need less training data



- b.
- There are two flavors of transfer learning:
 - Feature-based: Use outputs of pre-trained model as "features" input to downstream model; further training only occurs in downstream model:

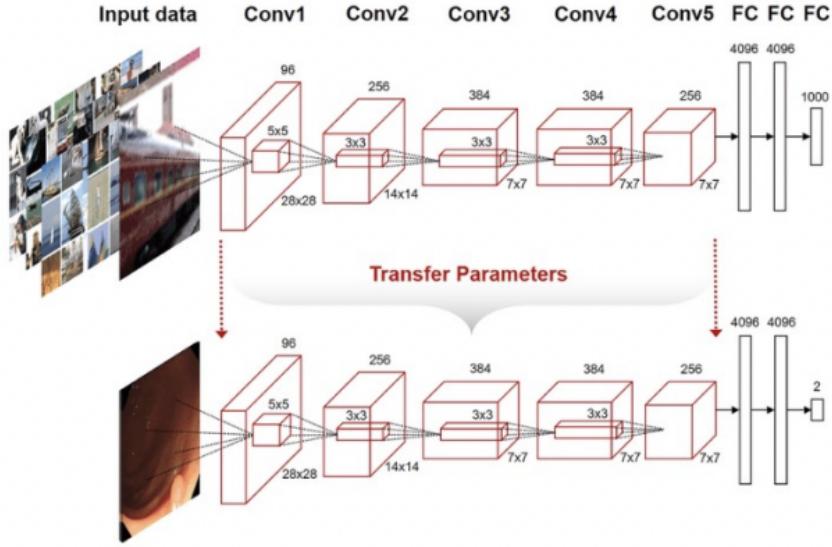


- Example: ELMO
- Fine-Tuning: Outputs of pre-trained model are inputs to downstream model, but BOTH models are trained; typically, changes to pre-trained model are minimal.



- Example: GPT, BERT

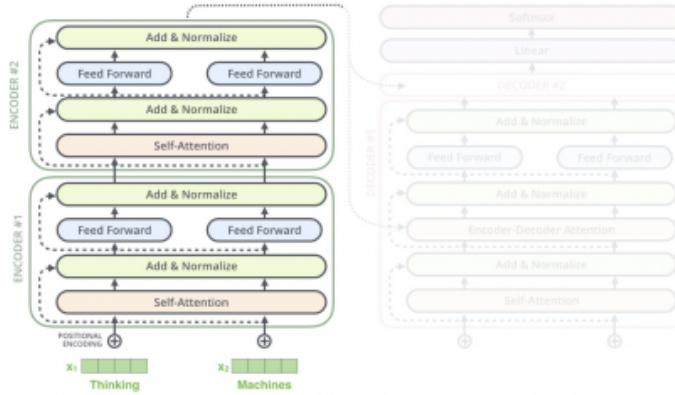
- d. This was first used in image processing (but is now used everywhere):



AlexNet structure of transfer learning from ImageNet database. The parameters are transferred in all layers from Conv1 to Conv5 except FC.

7. Transfer Learning for NLP: GPT, ELMO, and BERT

- a. Pre-trained models are typically very large and trained for a long time on huge amounts of data (so they have a lot of knowledge you can leverage for your specific task).
- b. The most successful systems in NLP are
 - i. GPT (Generative Pre-trained Transformer)
 - ii. ELMO (Embeddings from Language Models)
 - iii. BERT (Bidirectional Encoder Representations from Transformers)
- c. These all
 - i. Are based on transformers, but use ONLY the encoder phase;
 - ii. Produce contextual embeddings;
 - iii. GPT is autoregressive (outputs feed back into inputs)



- d. The most significant difference in the models is how they process the input sequence:
- GPT uses a simplified "causal" transformer model which only uses left-context:

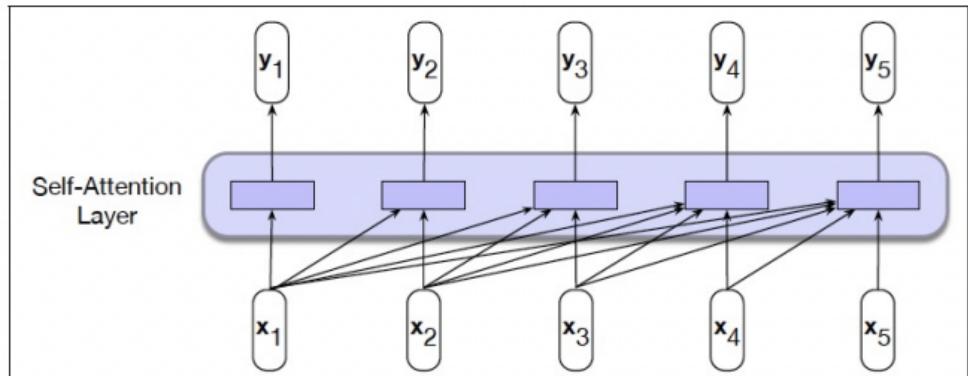


Figure 11.1 A causal, backward looking, transformer model like Chapter 9. Each output is computed independently of the others using only information seen earlier in the context.

- ELMO and BERT use a bidirectional transformer architecture:

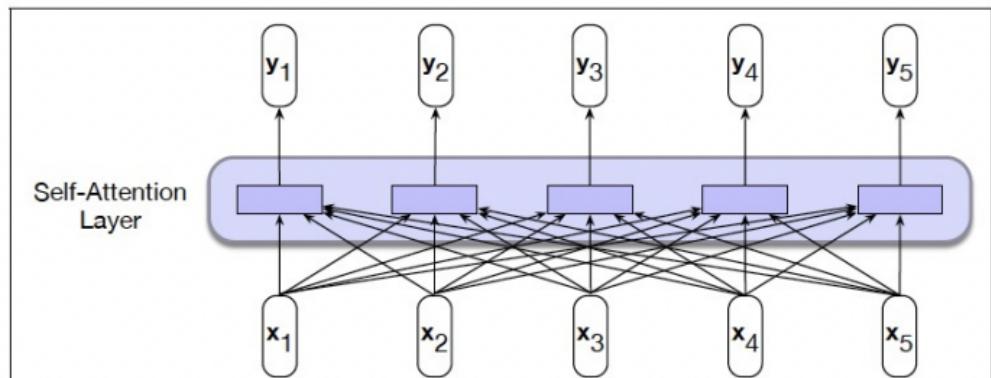


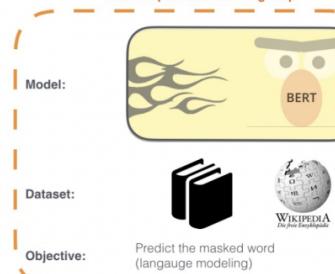
Figure 11.2 Information flow in a bidirectional self-attention model. In processing each element of the sequence, the model attends to all inputs, both before and after the current one.

8. BERT

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

Semi-supervised Learning Step



2 - **Supervised** training on a specific task with a labeled dataset.

Supervised Learning Step



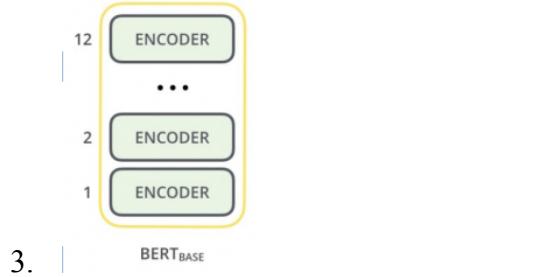
Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Alreides, please find attached...	Not Spam

a.

b. Bert has two implementations:

i. Bert base:

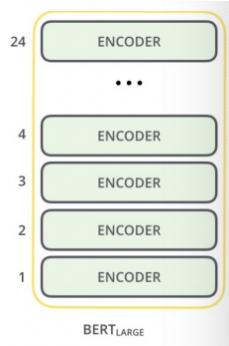
1. 12 layers; 12 attention heads per layer;
2. 768 hidden units; 110 M parameters



3.

ii. Bert large:

1. 16 layers; 16 attention heads per layer;
2. 1024 hidden units; 340 M parameters



3.

iii. BERT added two significant ideas for training which allowed it to achieve SOTA performance on significant tasks:

1. Training using a Masked Language Model (MLM);
2. Focus on a "next sentence prediction" task with two sentences as input.

iv. Note: The first versions of BERT used the 800 M word BooksCorpus and a 2.5 B word English Wikipedia corpus.

c. The Masked Language Model is based on an educational theory/testing paradigm known as the Cloze Task, where students learn a language by filling in blanks in a story or piece of text:

A screenshot of a Cloze Task interface. On the left, there is an illustration of three pink pigs standing next to a straw house. To the right, there is a text area containing a story with several blank spaces for words. On the far right, there is a 'Word vault' window listing various words with checkmarks next to them. The story text is as follows:

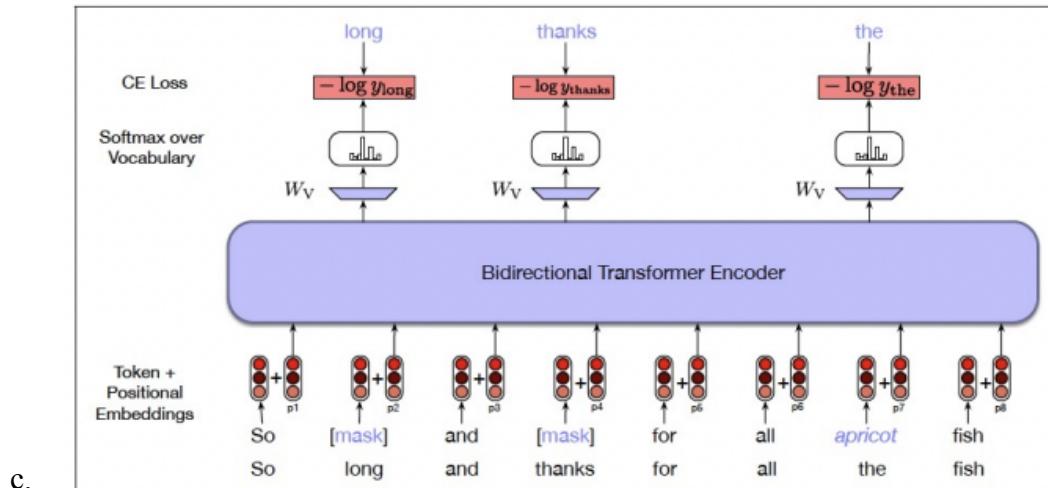
Once upon a time there was an ___ pig with three ___ pigs, and because there was not enough to ___ them, she sent them out to ___ their fortunes. The ___ little pig went off and met a man who had a bundle of ___. He asked the man, "Please give me that ___ to build a ____." The man gave him the ___ and the little pig ___ a house with it.

Word vault contents:

- ✓ old
- ✓ little
- ✓ feed
- ✓ seek
- ✓ test
- ✓ straw
- ✓ straw
- ✓ house
- ✓ straw

9. Training BERT

- a. Masked Language Modeling uses unannotated text from a large corpus. 15% of the words in the corpus are selected for the training phase:
 - i. of these, 80% of replaced with the token [MASK]
 - ii. 10% are replaced with randomly-selected tokens
 - iii. 10% are left unchanged.
- b. The model is trained to predict the missing tokens



- c.
- d. A variation of MLM uses spans (subsequences of the input sequence); all of the words in the span are replaced as before:

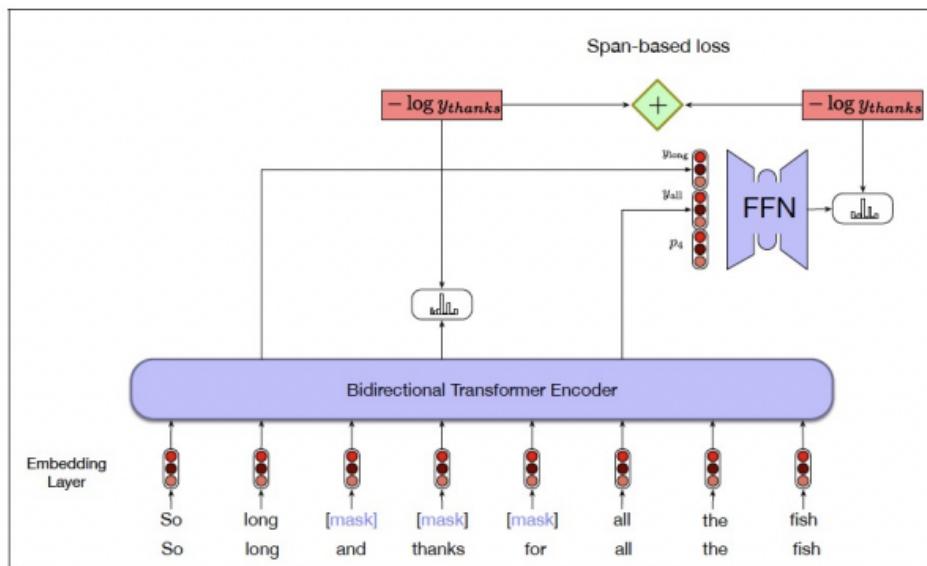


Figure 11.6 Span-based language model training. In this example, a span of length 3 is selected for training and all of the words in the span are masked. The figure illustrates the loss computed for word *thanks*; the loss for the entire span is based on the loss for all three of the words in the span.

- e. The Next Sentence Prediction task is to input TWO sentences starting with the token [CLS] and separated by the token [SEP]. The training set has 50% sentences that are next to each other in the corpus, and 50% random sentences.

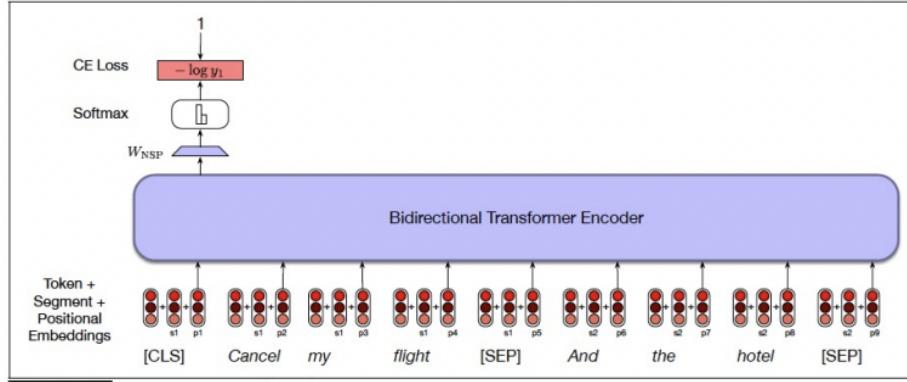


Figure 11.7 An example of the NSP loss calculation.

10. Using BERT

- a. Bert can be used for sentence classification if a single sentence is input:

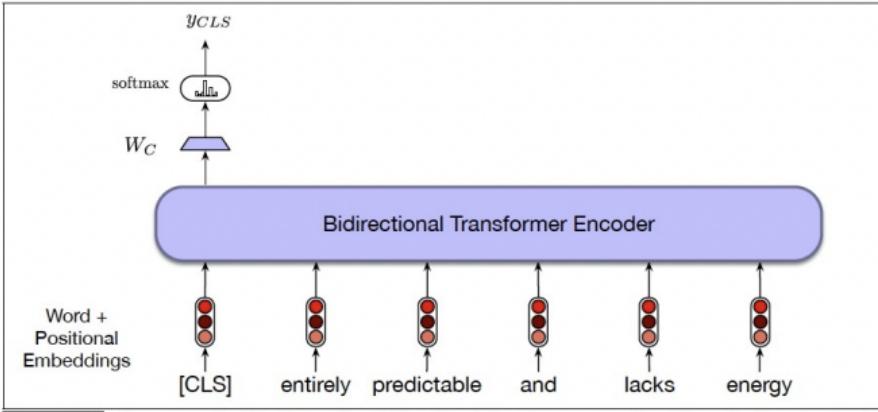


Figure 11.8 Sequence classification with a bidirectional transformer encoder. The output vector for the [CLS] token serves as input to a simple classifier.

- b. BERT can classify the relationship between two sentences:

- Neutral
 - a: Jon walked back to the town to the smithy.
 - b: Jon traveled back to his hometown.
- Contradicts
 - a: Tourist Information offices can be very helpful.
 - b: Tourist Information offices are never of any help.
- Entails
 - a: I'm confused.
 - b: Not all of it is very clear to me.

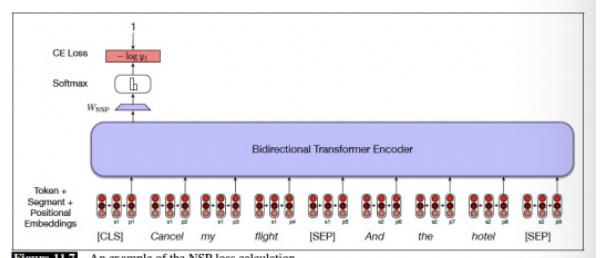
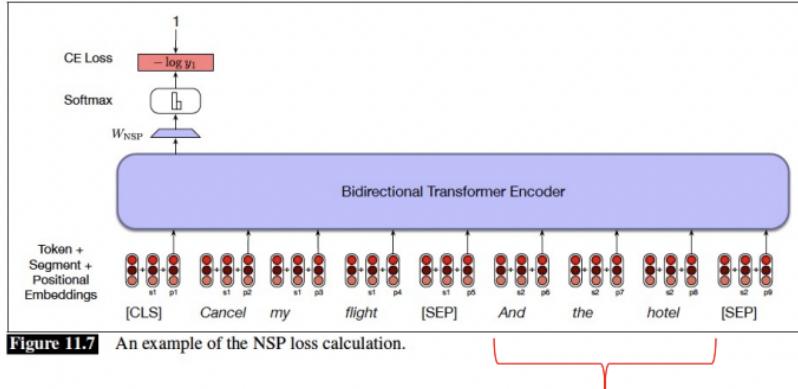


Figure 11.7 An example of the NSP loss calculation.

- c. BERT can generate the most likely sentence to follow a given sentence:



Use Beam Search to
find most likely
sentence to follow.

- d. Bert can be used for sequence labelling if all of the outputs are used:

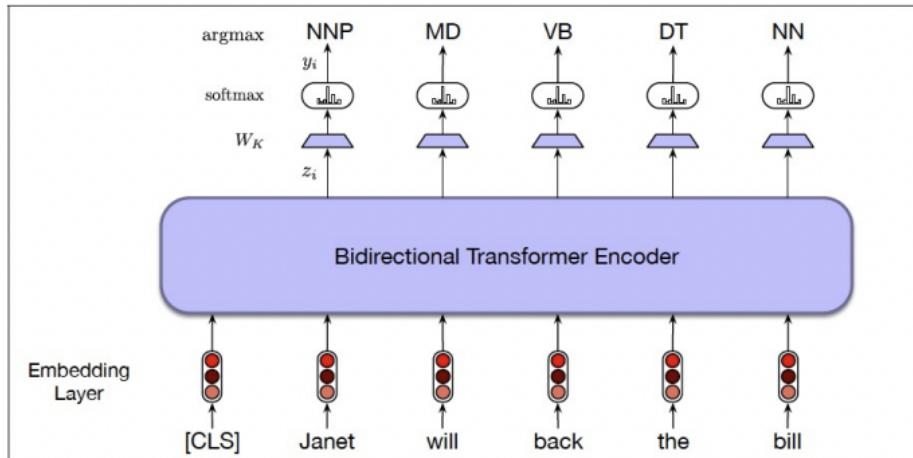


Figure 11.9 Sequence labeling for part-of-speech tagging with a bidirectional transformer encoder. The output vector for each input token is passed to a simple k-way classifier.