

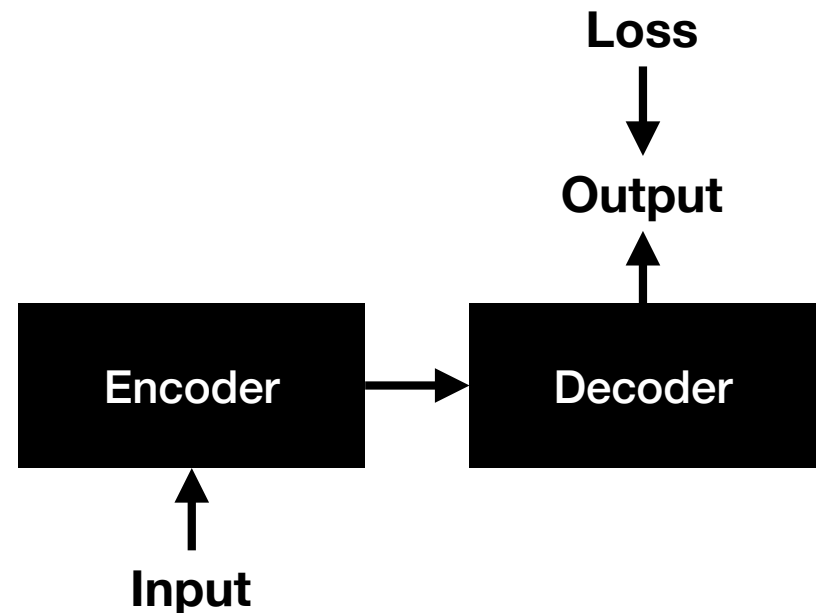
# Raw Data

## Pre-training Encoder-Decoders

Cornell CS 5740: Natural Language Processing  
Yoav Artzi, Spring 2023

# Encoder-decoder

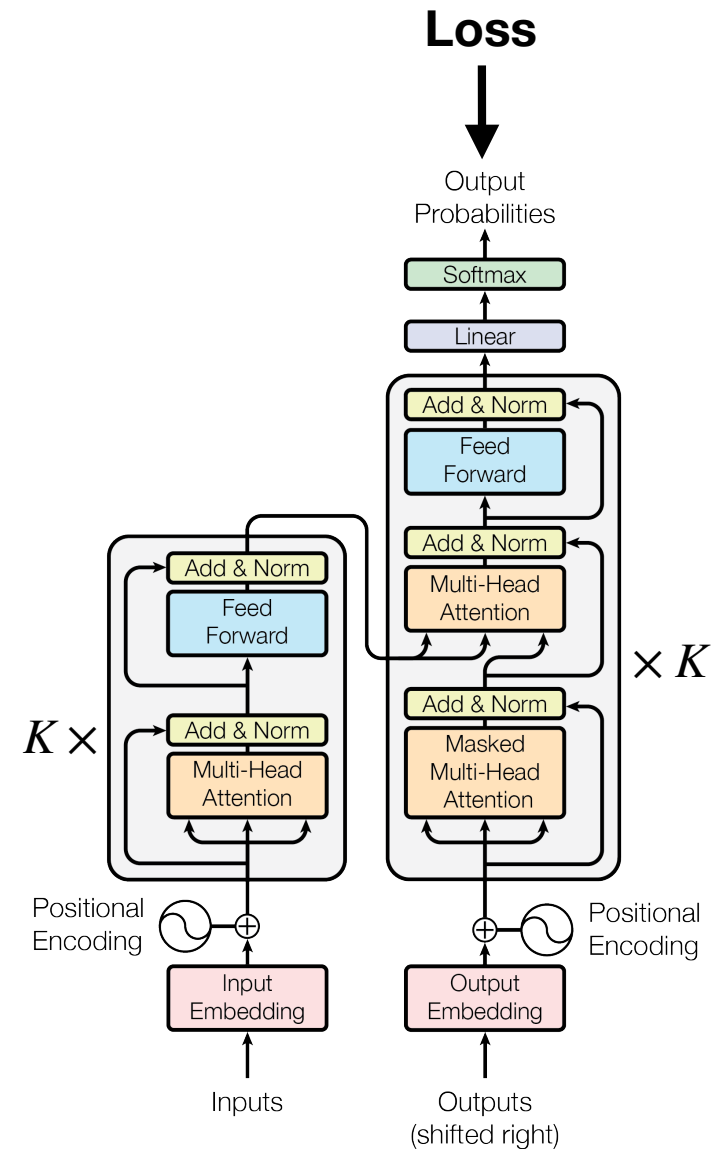
- The model is composed of two components
- Bidirectional **encoder** to process the input
- Autoregressive **decoder** to generate output
- Training is usually done with loss on the output
  - Propagates into the decoder and through it to the encoder



# Encoder-decoder

## With Transformers

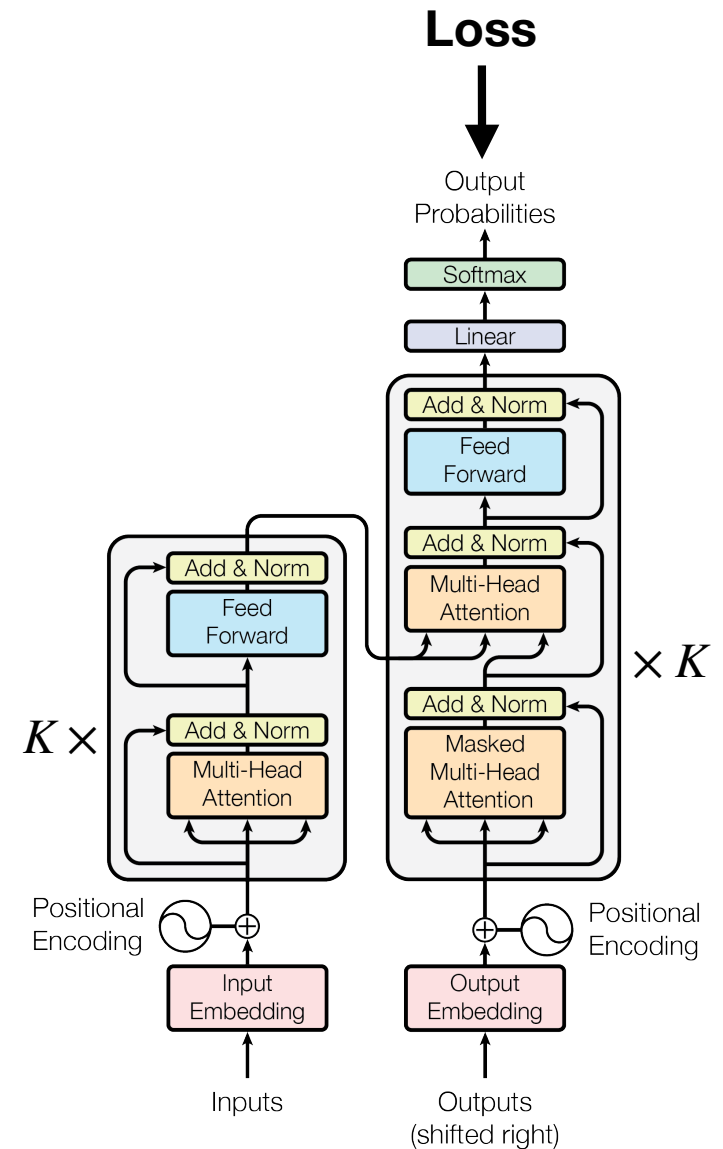
- The model is composed of two components
- Bidirectional **encoder** to process the input
- Autoregressive **decoder** to generate output
- Training is usually done with loss on the output
  - Propagates into the decoder and through it to the encoder



# Encoder-decoder

## With Transformers

- Bidirectional **encoder** to process the input
- Autoregressive **decoder** to generate output
- Why does this structure make sense?



# Encoder-decoder Pre-training

## The BART Recipe

- An encoder-decoder (sequence-to-sequence) pre-trained model
- Extends the BERT approach to encoder-decoder

# BART

## BERT Reminder

- Encoder-only Transformer
- Trained on raw data
- Two self-supervised objects:
  - Masked LM
  - Next-sentence prediction
- Transformed the NLP task landscape — if you have enough data, fine-tuned BERT works really well

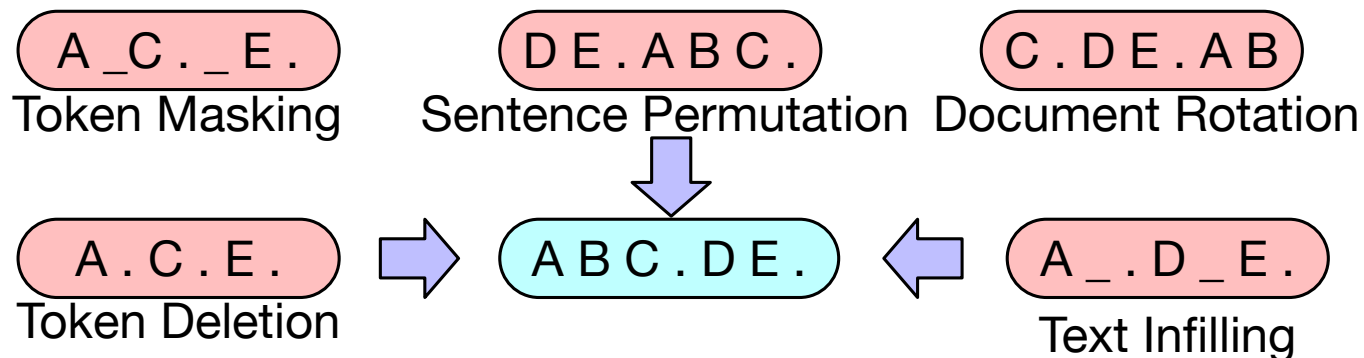
# BART

- How does the BERT learning approach adapt to an encoder-decoder architecture?
  - Output is generated by decoder, and the loss is on the output
  - Input is a sequence of tokens

# BART

## Denoising Self-supervised Objective

- Corrupt the input following five different recipes



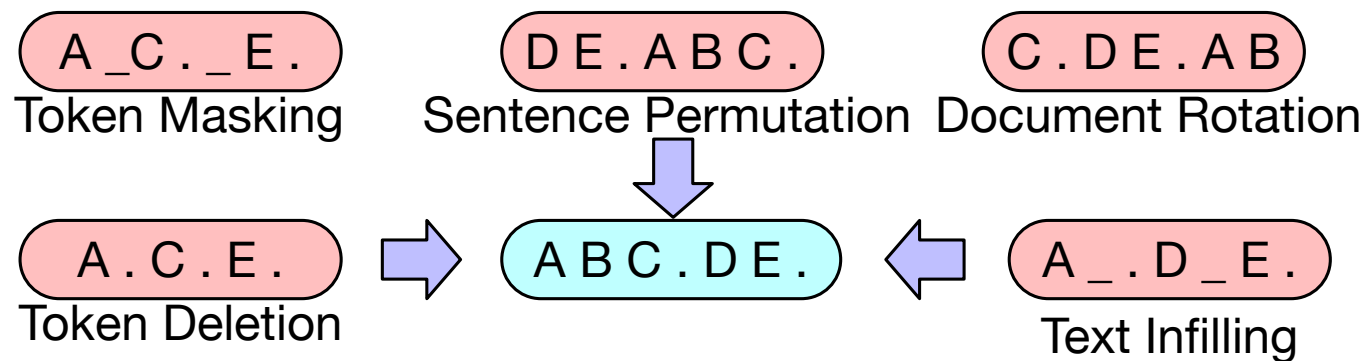
- Try to recover the pre-corrupted input by generating it using the decoder
- Train on a lot of raw text data, just like with BERT
- How to compute the loss? Loss can be computed using “teacher forcing”



# BART

## Denoising Self-supervised Objective

- Corrupt the input following five different recipes



- Try to recover the pre-corrupted input by generating it using the decoder
- How to compute the loss? Loss can be computed using “teacher forcing”

# BART

## What Do We Get?

- BERT: a pre-trained encoder
- BART: pre-trained decoder and encoder
  - Can use both
  - Or can use only the decoder wherever we would use BERT

# BART

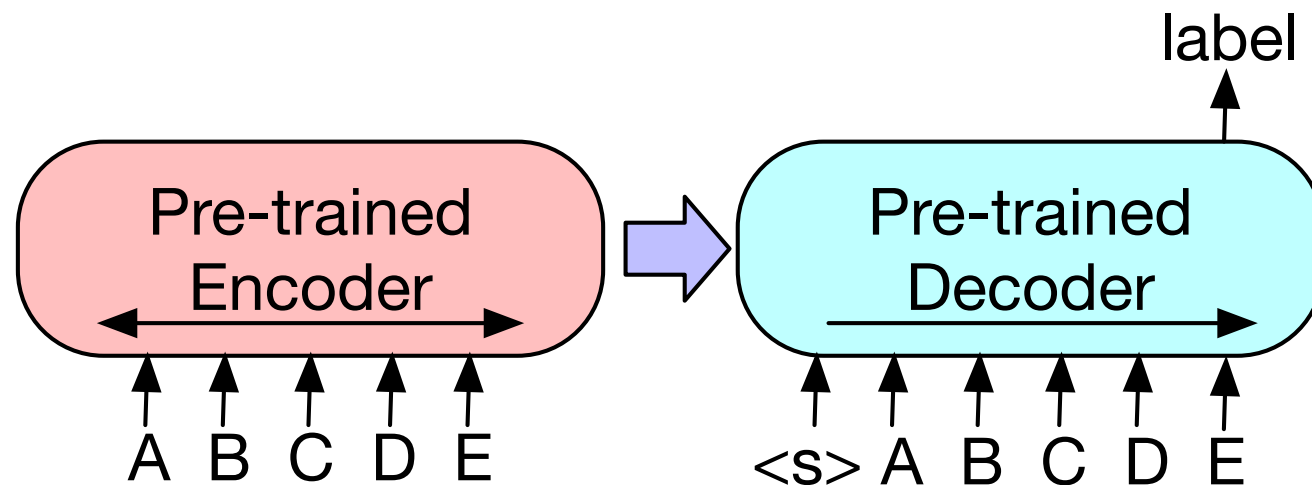
## How to Use?

- Similar to BERT: fine-tune for the end task
- Very natural for summarization
  - Because input and output vocabulary are the same
- How can we use for classification
- What about machine translation?

# BART

## Classification

- The input is given the encoder
- The same input is forced decoded in the decoder via “teacher forcing”
- The representation from the final decoder hidden state is given to a classification head



# BART

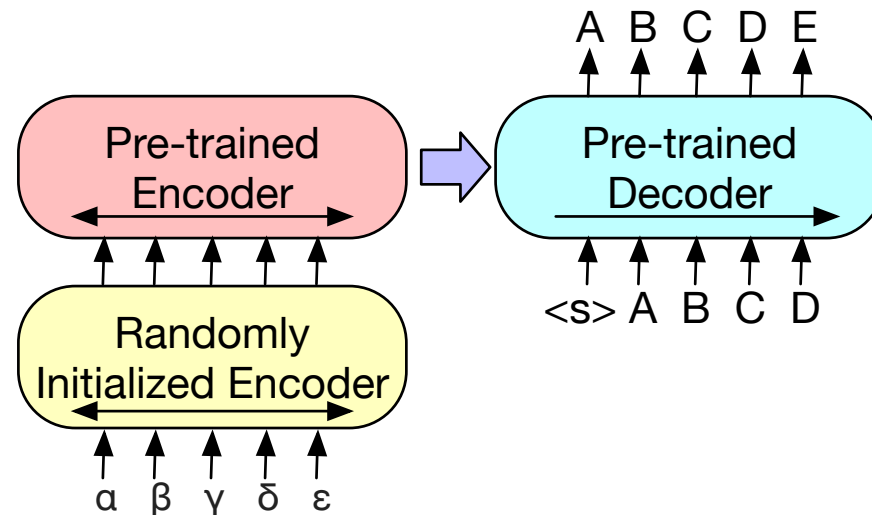
## Machine Translation

- In MT, the input and output vocabularies are different
- When is that a problem with BART?
- How can we solve it?

# BART

## Machine Translation

- In MT, the input and output vocabularies are different
- When is that a problem with BART?
- How can we solve it?
- Add a small pre-encoder encoder to replace the BART input embeddings with computed embeddings



# BART

## Performance

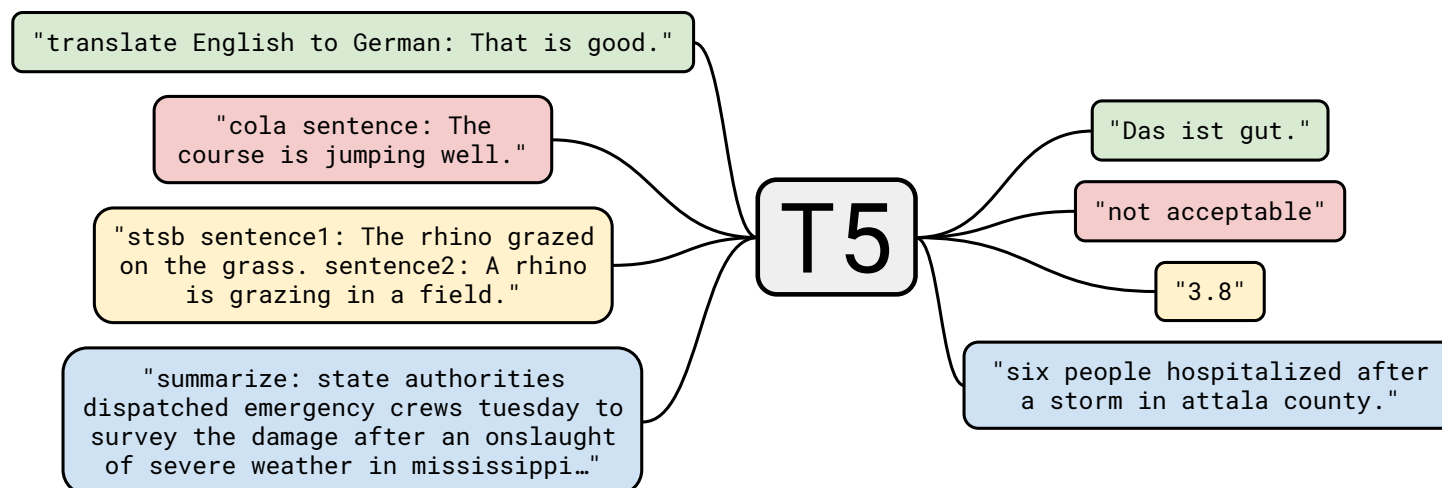
- Can do anything that BERT does
- But can also do generation tasks (e.g., summarization)

Model	SQuAD 1.1 F1	MNLI Acc	ELI5 PPL	XSum PPL	ConvAI2 PPL	CNN/DM PPL
BERT Base (Devlin et al., 2019)	88.5	<b>84.3</b>	-	-	-	-
BART Base						
w/ Token Masking	90.4	84.1	25.05	7.08	11.73	6.10
w/ Token Deletion	90.4	84.1	24.61	6.90	11.46	5.87
w/ Text Infilling	<b>90.8</b>	84.0	24.26	<b>6.61</b>	<b>11.05</b>	5.83
w/ Document Rotation	77.2	75.3	53.69	17.14	19.87	10.59
w/ Sentence Shuffling	85.4	81.5	41.87	10.93	16.67	7.89
w/ Text Infilling + Sentence Shuffling	<b>90.8</b>	83.8	24.17	6.62	11.12	<b>5.41</b>

# Encoder-decoder Pre-training

## The T5 Recipe

- Concurrent and similar to BART
- Also adopted the text-to-text approach for all NLP tasks





# T5

## Pretraining

- Pretraining is similar to the denoising objective of BART:
  - Input: text with gaps (phrases removed)
  - Output: sequence of phrases to fill the gaps

Original text

Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs

Thank you <X> me to your party <Y> week.

Targets

<X> for inviting <Y> last <Z>

# T5

## Results

- T5 was trained on one of the first very large corpora: 750GB of text, with pre-training using  $2^{35}$  tokens
- First to show the impact of data scale
- Why did they repeat the smaller scale datasets?

Number of tokens	Repeats	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Full data set	0	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
$2^{29}$	64	<b>82.87</b>	<b>19.19</b>	<b>80.97</b>	<b>72.03</b>	<b>26.83</b>	<b>39.74</b>	<b>27.63</b>
$2^{27}$	256	82.62	<b>19.20</b>	79.78	69.97	<b>27.02</b>	<b>39.71</b>	27.33
$2^{25}$	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
$2^{23}$	4,096	76.34	18.33	70.92	59.29	26.37	38.84	25.81

# BART and T5

## Takeaways

- BART and T5 are very useful for all sorts of sequence-to-sequence tasks with language
  - T5 comes in different sizes
  - There are various customization (e.g., CodeT5)
- Extended the generalizations conclusions from BERT, and demonstrated the impact of data scale

# Acknowledgements

- Slides are based on slides from UT Austin CS 388 by Greg Durrett