

# Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer: The T5 Model

Victoria Graf and Abhishek Panigrahi

# Transfer Learning

- Pre-training!
- Start with **unlabeled** data (unlike computer vision)
- General-purpose “English” knowledge

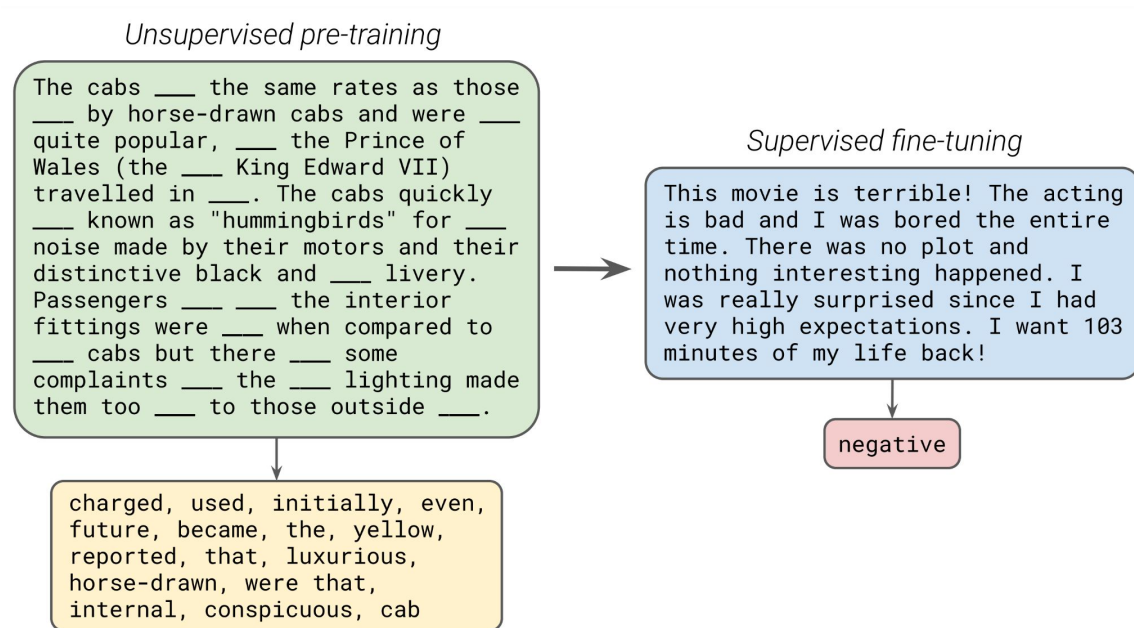
# Transfer Learning



# Transfer Learning



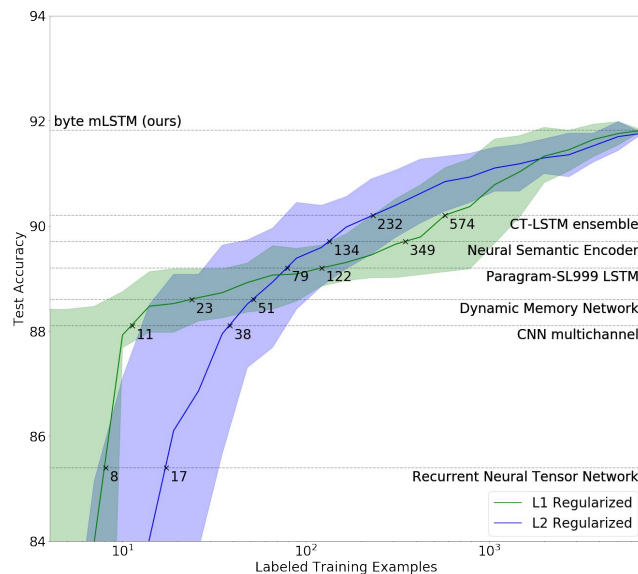
# Transfer Learning



Slide adapted from Colin Raffel

# A Very Brief Context

- 2017: Attention Is All You Need, Unsupervised sentiment neuron
- 2018: ELMo, GPT-1, BERT
  - Bidirectionality
  - Transformers
- 2019: RoBERTa, SpanBERT, ALBERT
- 2020: T5!



<https://openai.com/blog/unsupervised-sentiment-neuron/>

# Transfer Learning: Comparisons?

**Lots of research, so many...**

- Pre-training objectives
- Unlabeled data sets
- Fine-tuning methods
- Model architectures/scales

**... so how do we compare benchmarks?**

# Transfer Learning Comparisons

- Model A has 1B parameters and uses 100M pre-training tokens from BooksCorpus



# Transfer Learning Comparisons

- Model A has 1B parameters and uses 100M pre-training tokens from BooksCorpus
- Model B is made by Google and their deep pockets! It has **2B** parameters and uses **200M** pre-training tokens from Wikipedia

# Transfer Learning Comparisons

- Model A has 1B parameters and uses 100M pre-training tokens from BooksCorpus
- Model B is made by Google and their deep pockets! It has **2B** parameters and uses **200M** pre-training tokens from Wikipedia
- Model B has better performance on SuperGLUE than Model A

# Transfer Learning Comparisons

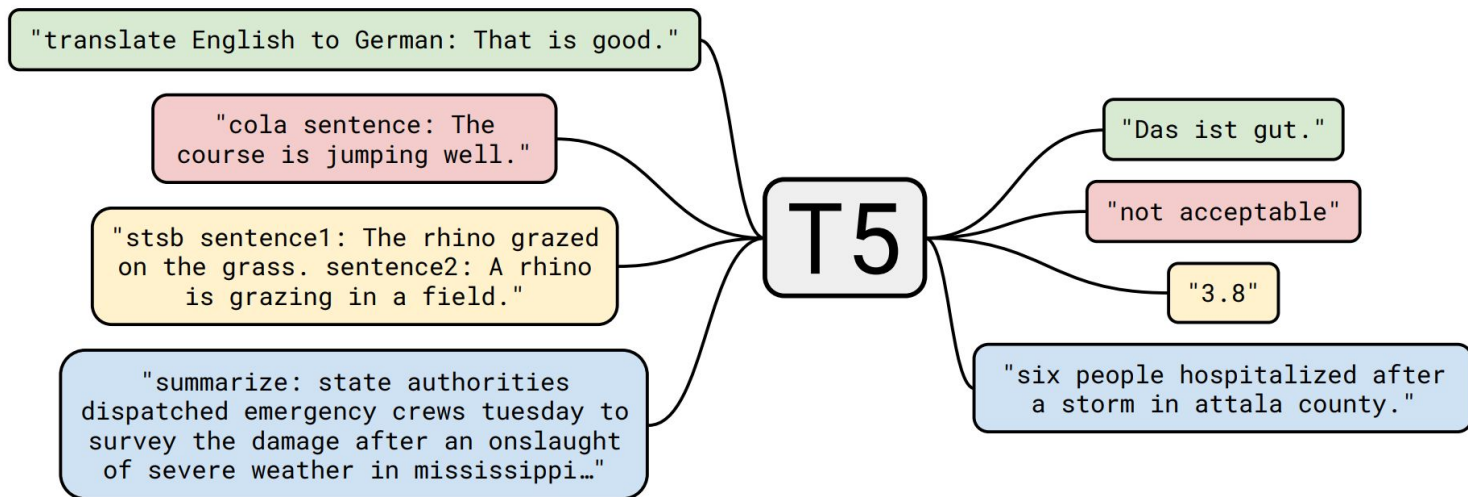
- Model A has 1B parameters and uses 100M pre-training tokens from BooksCorpus
- Model B is made by Google and their deep pockets! It has **2B** parameters and uses **200M** pre-training tokens from Wikipedia
- Model B has better performance on SuperGLUE than Model A

**Is Wikipedia better for pre-training than BooksCorpus?**

# T5: The Basic Idea

- **Text-to-Text Transfer Transformer**
- Every task, one format!
  - Previous attempts included:
    - Question answering
    - Language modeling
    - Span extraction
  - ... but had limitations
- “[Task-specific prefix]: [Input text]” -> “[output text]”

# T5: The Basic Idea



# T5: The Basic Idea

- GLUE and SuperGLUE **classification**; CNN/Daily Mail **abstractive summarization**; SQuAD **question answering**; and WMT English to German, French, and Romanian **translation**
  - GLUE/SuperGLUE: Sentence acceptability judgment, sentiment analysis, paraphrasing/sentence similarity, natural language inference, coreference resolution, sentence completion, word sense disambiguation, question answering
  - French: high resource, Romanian: low resource
- Separate fine-tuning for each task

# Some tasks

Recall: SQuAD, GLUE benchmarks

- CoLA (GLUE): Sentence acceptability
  - **Input:** sentence, **output:** labels “acceptable” or “not acceptable”
  - Ex: “The course is jumping well.” -> not acceptable
- STS-B (GLUE): Sentence similarity
  - **Input:** pair of sentences, **output:** similarity score [1,5]
  - Ex: “sentence1: The rhino grazed. sentence2: A rhino is grazing.” -> 3.8

# Some tasks

- COPA (SuperGLUE): Causal reasoning
  - **Input:** premise and 2 alternatives, **output:** alternative1 or alternative2
  - Ex: “Premise: I tipped the bottle. What happened as a RESULT?  
Alternative 1: The liquid in the bottle froze.  
Alternative 2: The liquid in the bottle poured out.”  
-> alternative2
- ReCoRD/MultiRC (SuperGLUE): Question answering/Reading comprehension



# Input/Output

[Task-specific prefix]: [Input text]

- EnDe (Translation):  
“translate English to German: That is good” -> “Das ist gut”
- CNNDM (Summarization):  
“summarize: state authorities dispatched...” -> “six people hospitalized after storm”

# Input/Output

[Task-specific prefix]: [Input text]

- CoLA (GLUE; Classification):  
“cola sentence: The course is jumping well.” -> “not acceptable”
- STS-B (GLUE; Regression):  
“stsb sentence1: The rhino grazed. sentence2: A rhino is grazing.” -> “3.8”

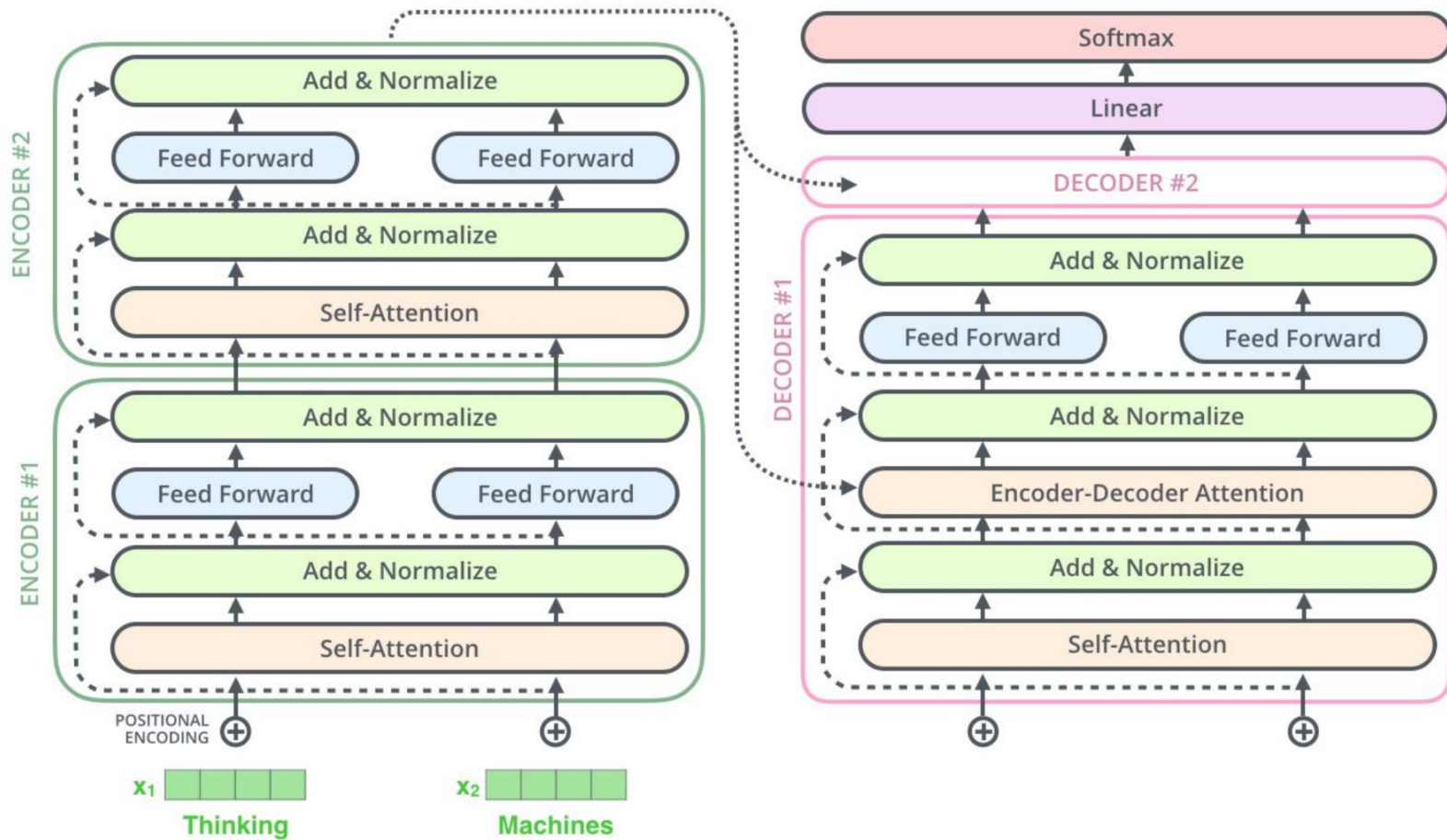
# Input/Output

- CoLA (GLUE; Classification):  
“cola sentence: The course is jumping well.” -> “hamburger”

**“Hamburger” is not a valid CoLA output, so this is a fail!**

# T5 Model

- Encoder-decoder model
  - Baseline size: two stacks of size BERT<sub>BASE</sub>
- Architecture from “Attention Is All You Need”
  - Different position embedding scheme



# C4: The Data

- **Colossal Clean Crawled Corpus**
- Web-extracted text
- English language only (langdetect)
- 750GB

20TB to 750GB? Where did everything go?

# C4: The Data

- Retain:
  - Sentences with terminal punctuation marks
  - Pages with at least 5 sentences, sentences with at least 3 words
- Deduplicate three sentence spans
- Remove:
  - References to Javascript
  - Lorem ipsum text
  - Code

# C4: The Data

Menu

Lemon

Introduction

The lemon, *Citrus Limon* (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae.

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

Article

The origin of the lemon is unknown, though lemons are thought to have first grown in Assam (a region in northeast India), northern Burma or China.

A genomic study of the lemon indicated it was a hybrid between bitter orange (sour orange) and citron.

Please enable JavaScript to use our site.

Home

Products

Shipping

Contact

FAQ

Dried Lemons, \$3.59/pound

Organic dried lemons from our farm in California.

Lemons are harvested and sun-dried for maximum flavor.

Good in soups and on popcorn.

The lemon, *Citrus Limon* (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae.

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Curabitur in tempus quam. In mollis et ante at consectetur.

Aliquam erat volutpat.

Donec at lacinia est.

Duis semper, magna tempor interdum suscipit, ante elit molestie urna, eget efficitur risus nunc ac elit.

Fusce quis blandit lectus.

Mauris at mauris a turpis tristique lacinia at nec ante.

Aenean in scelerisque tellus, a efficitur ipsum.

Integer justo enim, ornare vitae sem non, mollis fermentum lectus.

Mauris ultrices nisl at libero porta sodales in ac orci.

```
function Ball(r) {
  this.radius = r;
  this.area = pi * r ** 2;
  this.show = function(){
    drawCircle(r);
  }
}
```



# C4: The Data

Menu

Lemon

Introduction

The lemon, *Citrus Limon* (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae.

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

Article

The origin of the lemon is unknown, though lemons are thought to have first grown in Assam (a region in northeast India), northern Burma or China.

A genomic study of the lemon indicated it was a hybrid between bitter orange (sour orange) and citron.

Please enable JavaScript to use our site.

Home

Products

Shipping

Contact

FAQ

Dried Lemons, \$3.59/pound

**Organic dried lemons from our farm in California.**

**Lemons are harvested and sun-dried for maximum flavor.**

**Good in soups and on popcorn.**

The lemon, *Citrus Limon* (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae.

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Curabitur in tempus quam. In mollis et ante at consectetur.

Aliquam erat volutpat.

Donec at lacinia est.

Duis semper, magna tempor interdum suscipit, ante elit molestie urna, eget efficitur risus nunc ac elit.

Fusce quis blandit lectus.

Mauris at mauris a turpis tristique lacinia at nec ante.

Aenean in scelerisque tellus, a efficitur ipsum.

Integer justo enim, ornare vitae sem non, mollis fermentum lectus.

Mauris ultrices nisl at libero porta sodales in ac orci.

```
function Ball(r) {
  this.radius = r;
  this.area = pi * r ** 2;
  this.show = function(){
    drawCircle(r);
  }
}
```

# C4: The Data

750GB? What does that mean?

Data set	Size
★ C4	745GB
C4, unfiltered	6.1TB
RealNews-like	35GB
WebText-like	17GB
Wikipedia	16GB
Wikipedia + TBC	20GB

# Vocabulary

- 32,000 wordpieces shared across input and output
- Pre-training is English, but fine-tuning includes German, French, and Romanian
- Trained SentencePiece model 10:1:1:1 English : German : French : Romanian
  - Can handle fixed set of languages

# mT5

- mC4: Common Crawl dataset covering 101 languages!
  - Only line length filter, no punctuation filter
  - How do you sample across languages?
    - “Boosting” the probability of training on low-resource languages without overfitting
- Similar architecture to T5
- 6 tasks from the XTREME multilingual benchmark
  - Entailment, reading comprehension, NER, paraphrase identification
- Illegal predictions (XQuAD)

# Experiments

- Even Google has a budget...
  - NOT combinatorial
  - Standard deviation only found for baseline
  - $\sim 2^{35}$  or 34B pre-training tokens (much less than BERT!)
- Inverse square-root learning rate schedule with warm-up
- Results reported on validation sets

# Baseline Objective

Original text

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

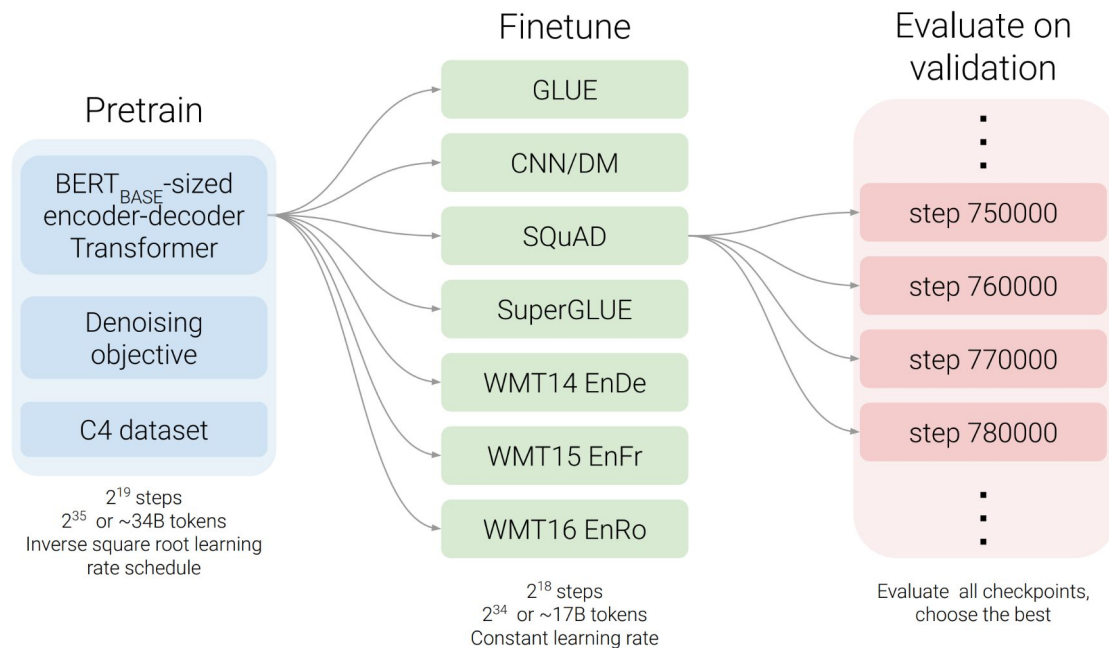
Inputs

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

Targets

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

# Workflow



Slide adapted from Colin Raffel

# Baseline Performance

**Bold** scores are within two standard deviations of the best score in a given experiment

	GLUE	CNNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline average	<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>
Baseline standard deviation	0.235	0.065	0.343	0.416	0.112	0.090	0.108
No pre-training	66.22	17.60	50.31	53.04	25.86	<b>39.77</b>	24.04

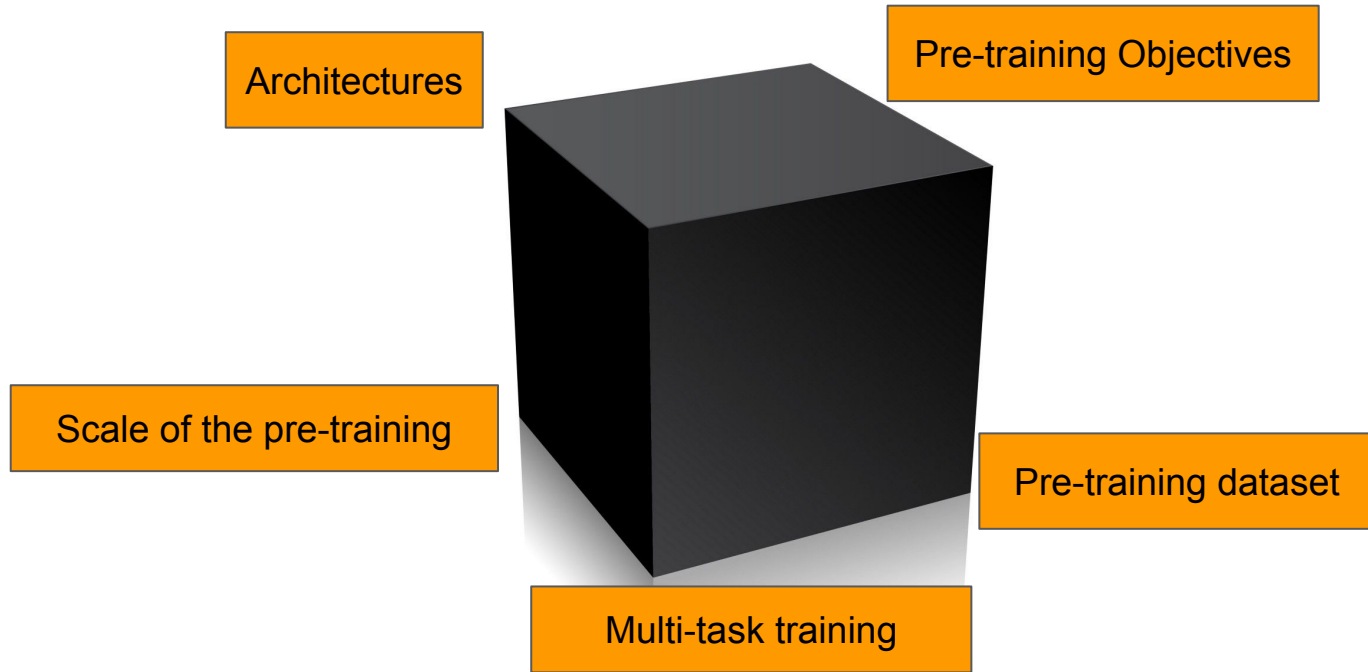


# Baseline Performance

- GLUE/SuperGLUE are sets of tasks including CoLA, STS-B, etc.
- CNNDM is a summarization task
- EnDe/EnFr/EnRo are translation tasks

	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline average	<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>
Baseline standard deviation	0.235	0.065	0.343	0.416	0.112	0.090	0.108
No pre-training	66.22	17.60	50.31	53.04	25.86	<b>39.77</b>	24.04

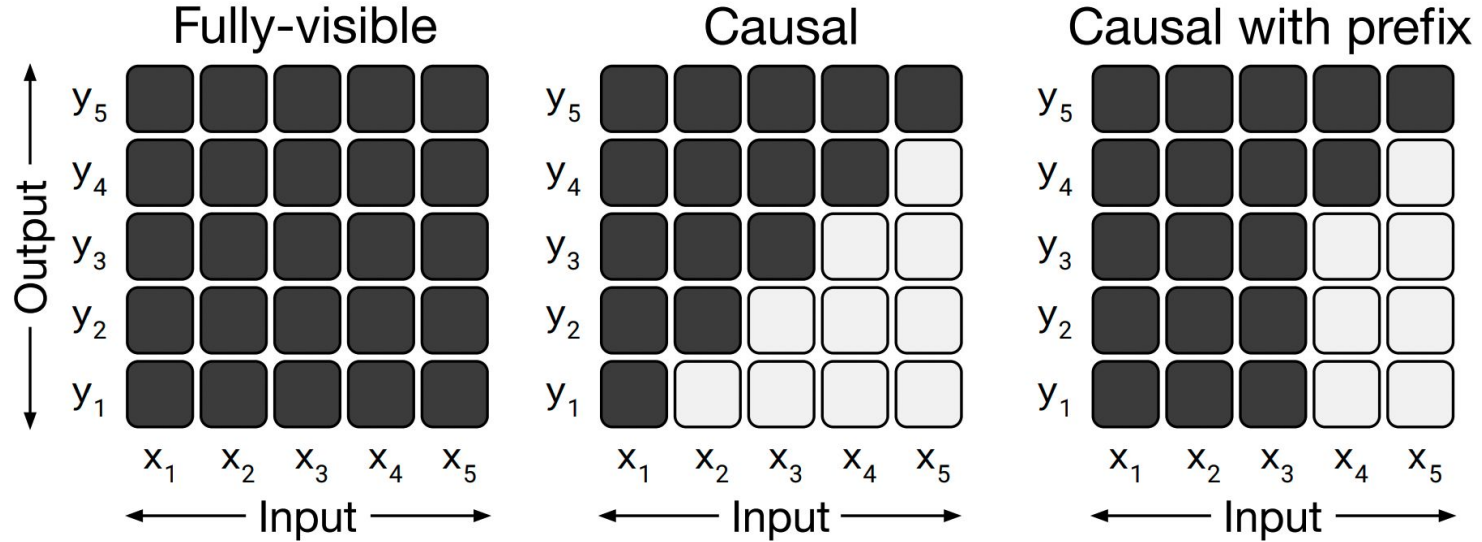
# Axis of Decisions for Pre-training and Fine-tuning



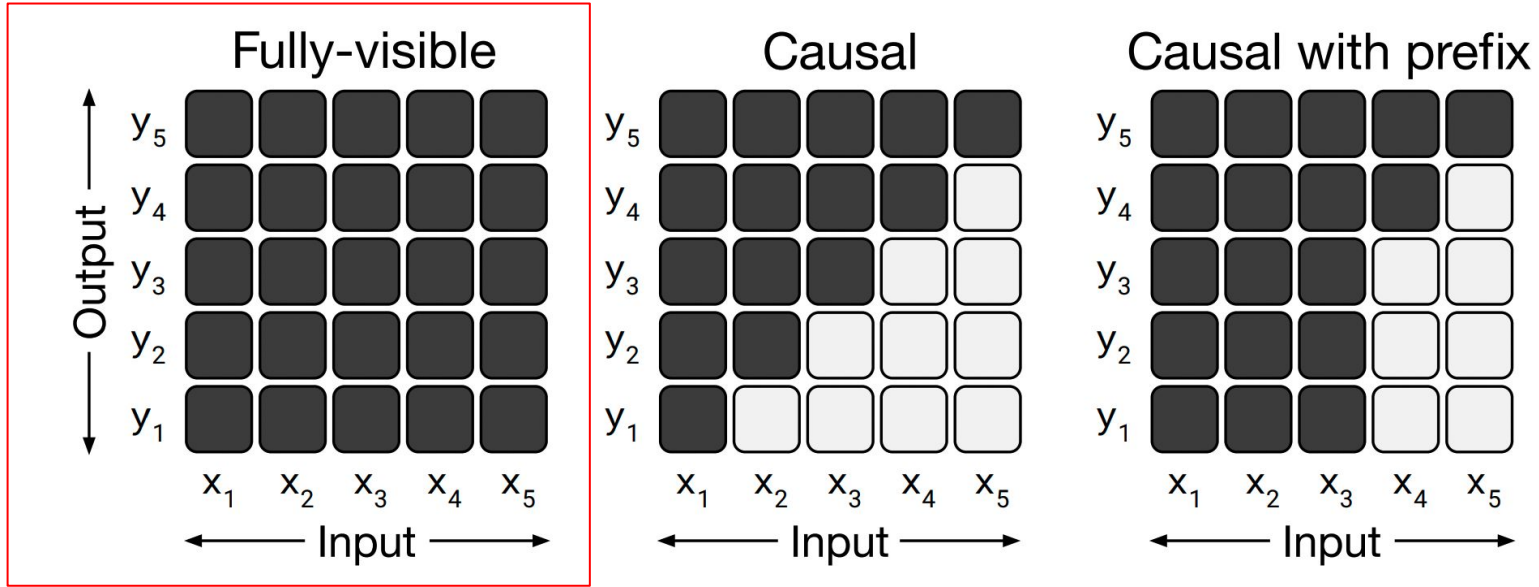
Understand the **first** order effect of each aspect by altering it while keeping other aspects of pre-training fixed.

# Architectures

# Different Attention Mask Patterns

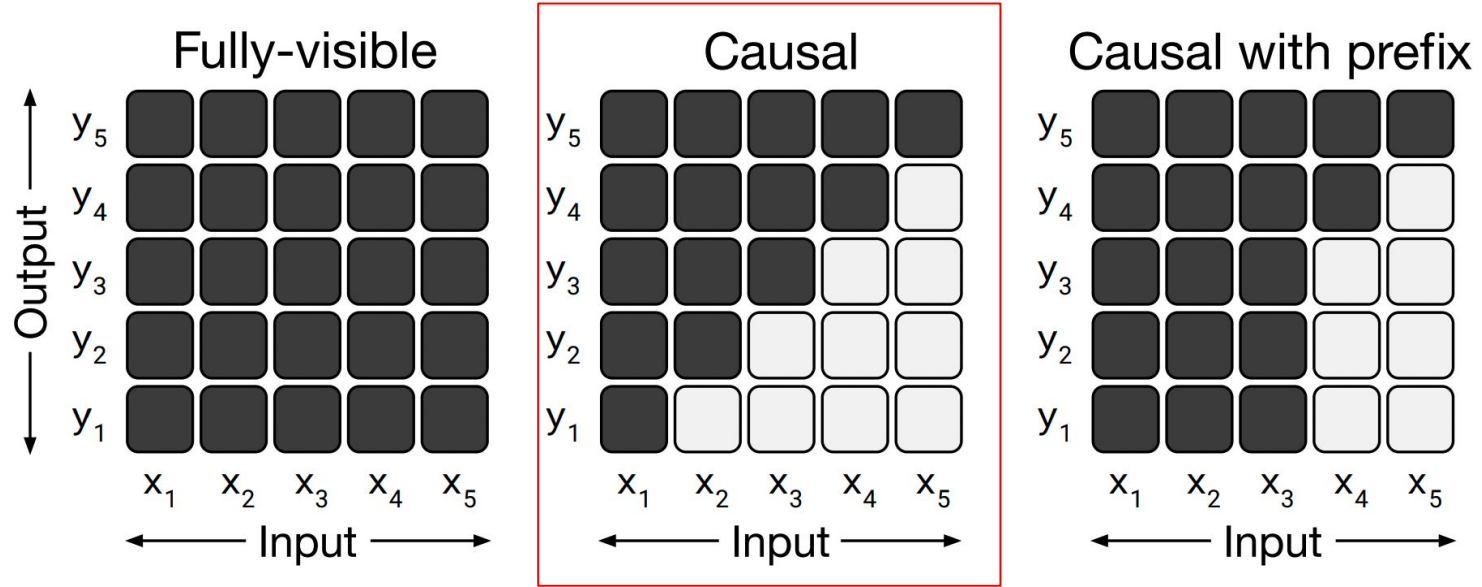


# Different Attention Mask Patterns



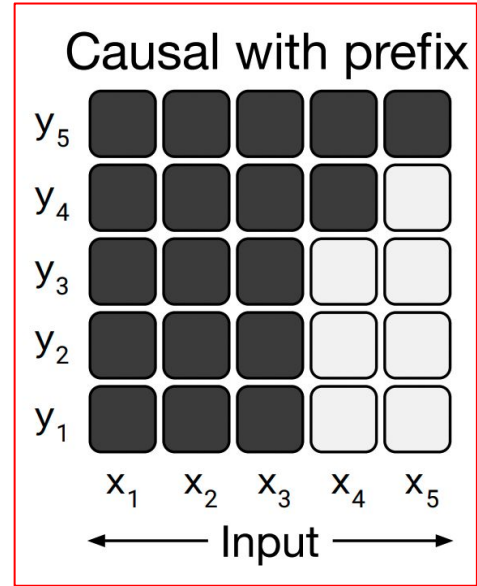
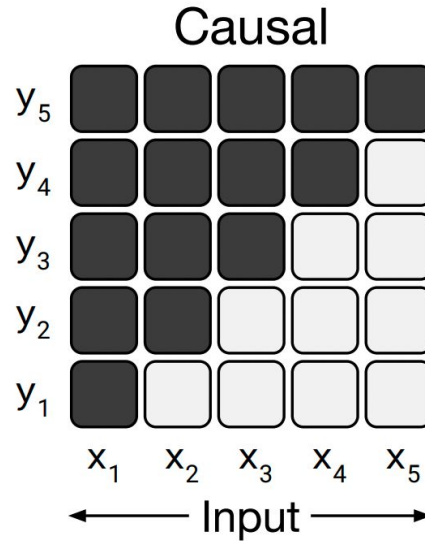
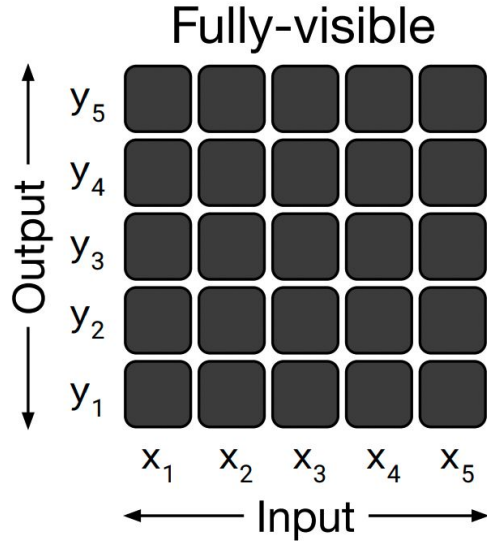
Fully visible mask allows the self attention mechanism to attend to the full input.

# Different Attention Mask Patterns



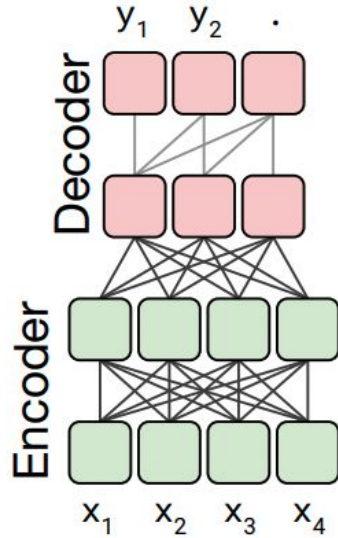
A causal mask doesn't allow output elements to look into the future.

# Different Attention Mask Patterns

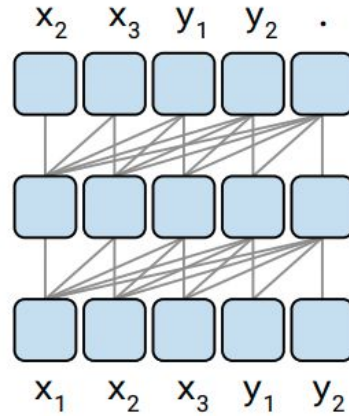


Causal mask with prefix allows to fully-visible masking on a portion of input.

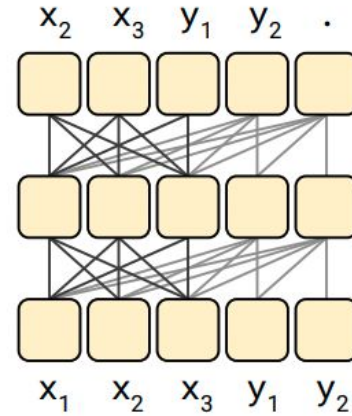
# Transformer Architecture Variants



Language model

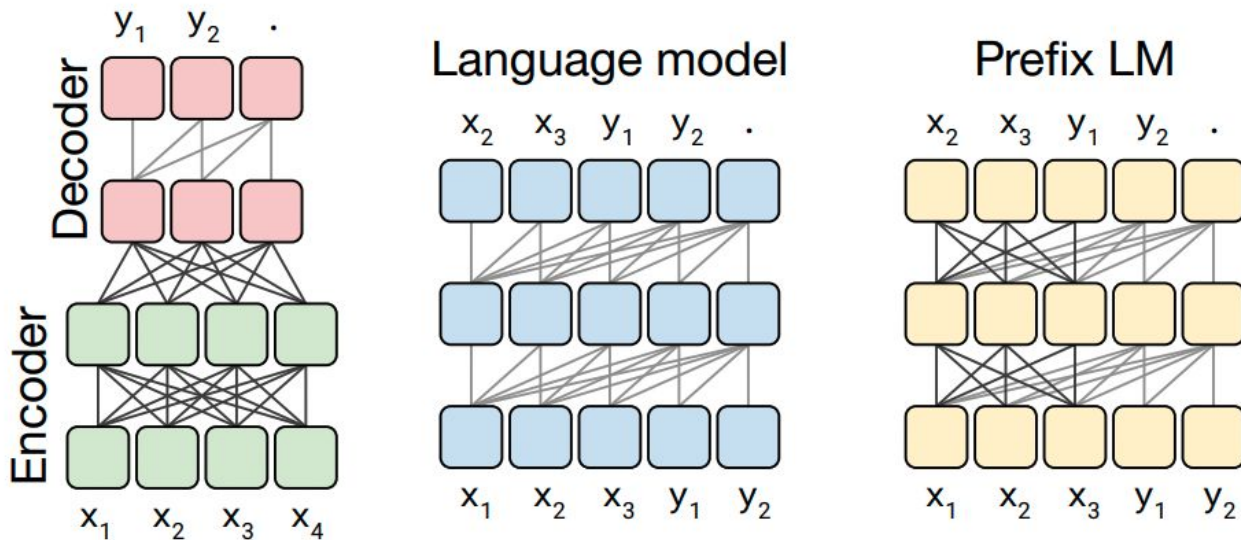


Prefix LM



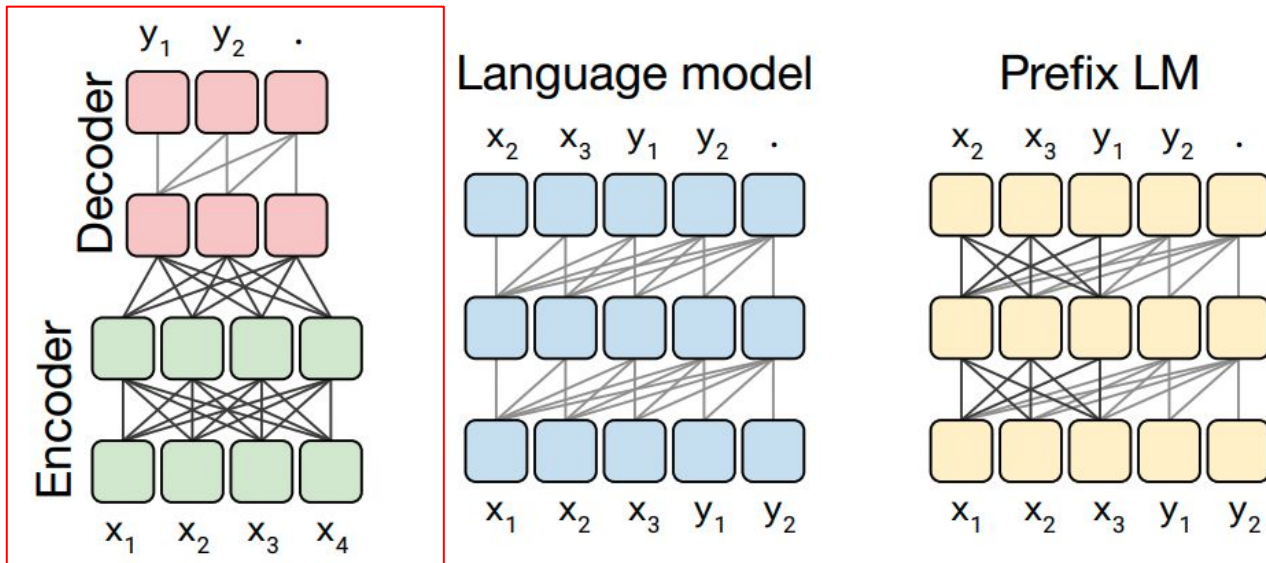


# Transformer Architecture Variants



Translation: That is good -> Das ist gut.

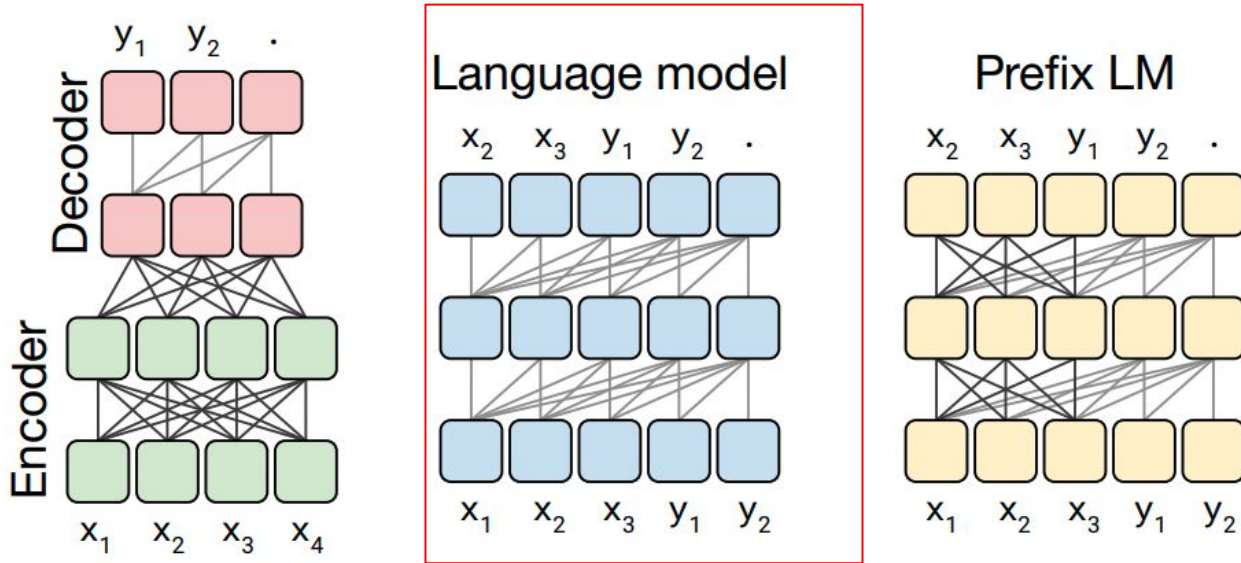
# Transformer Architecture Variants



Translation: That is good -> Das ist gut.

Translate English to German: That is good. Target: Das is gut.

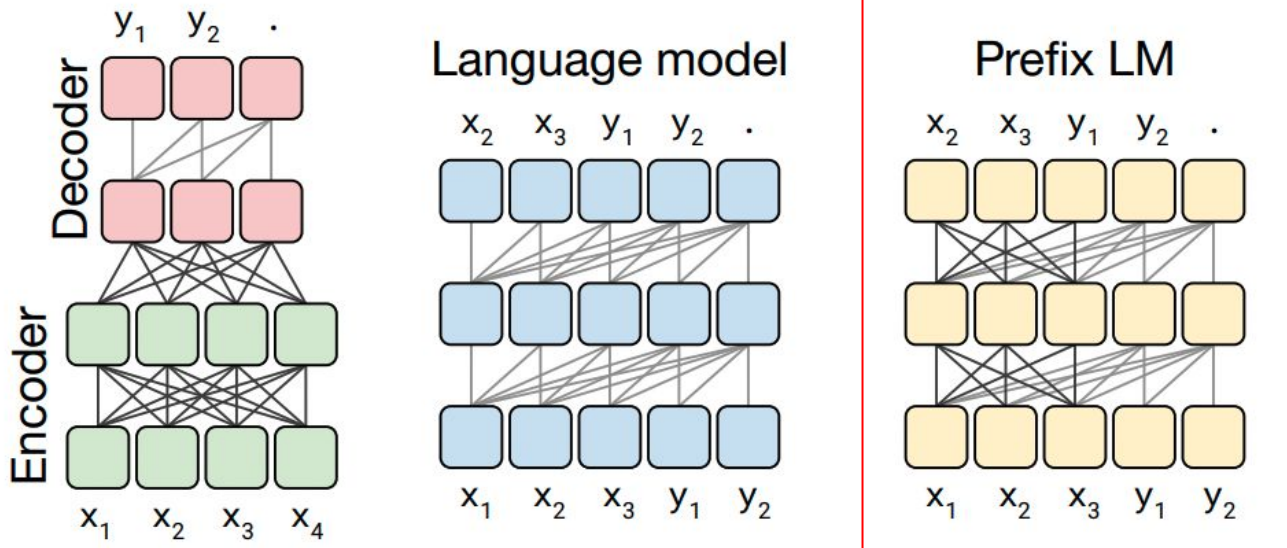
# Transformer Architecture Variants



Translation: That is good -> Das ist gut.

- Translate English to German: That is good. Target: Das ist gut.
  - “Good” representation can only look at “Translate English to German: That is”.

# Transformer Architecture Variants



Translation: That is good -> Das ist gut.

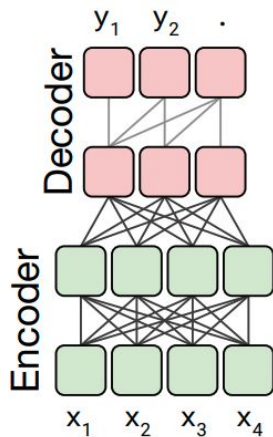
- Translate English to German: That is good. Target: Das is gut.
  - “Good” representation can look at “Translate English to German: That is. Target:”.

# Performance of different Architectural Variants

---

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>

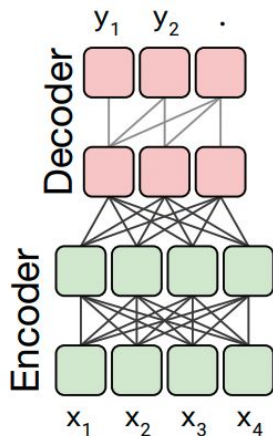
---



# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>

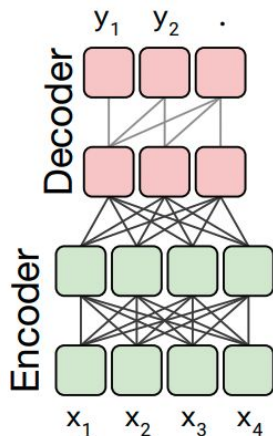
Input: Thank you for <X> me to your party <Y>.  
Target: <X> inviting <Y> last week.



# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>

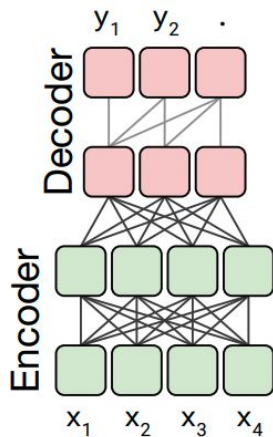
Number of parameters



# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>

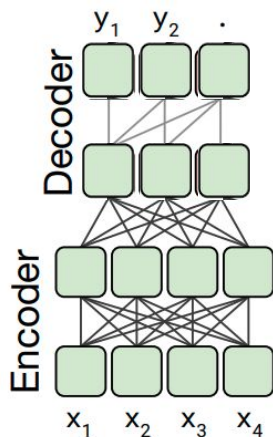
Number of flops





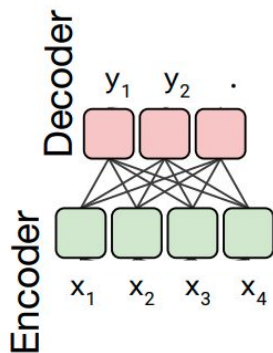
# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>



# Performance of different Architectural Variants

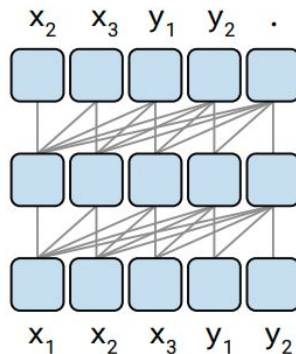
Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95



# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86

Language model

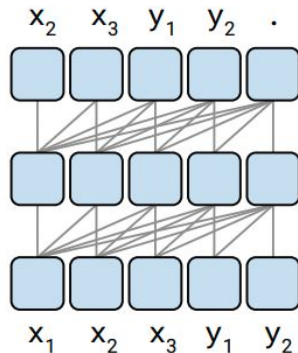


# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86

Language model is decoder-only

Language model

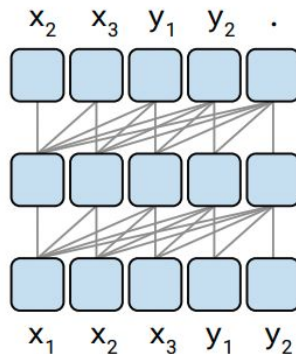


# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86

LM looks at both input and target, while encoder only looks at input sequence and decoder looks at output sequence.

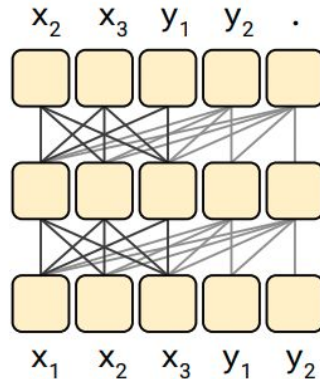
Language model



# Performance of different Architectural Variants

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	$P$	$M$	81.82	18.61	78.94	68.11	26.43	37.98	27.39

Prefix LM



# Performance of different Architectural Variants

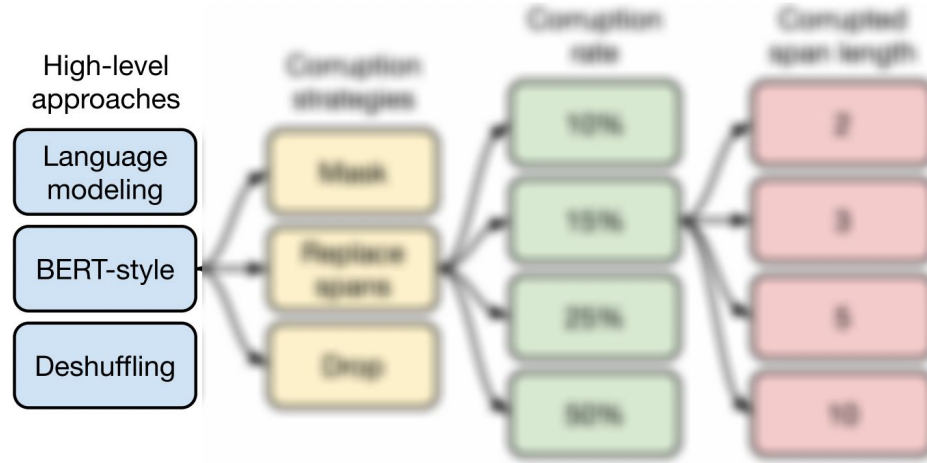
Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	$P$	$M$	81.82	18.61	78.94	68.11	26.43	37.98	27.39

1. Sharing parameters in encoder and decoder models perform nearly as well as the baseline.
2. Halving the number of layers in encoder and decoder hurts the performance.
3. Performance of Encoder and Decoder with shared parameters is better than decoder only LM and prefix LM.

# Objectives

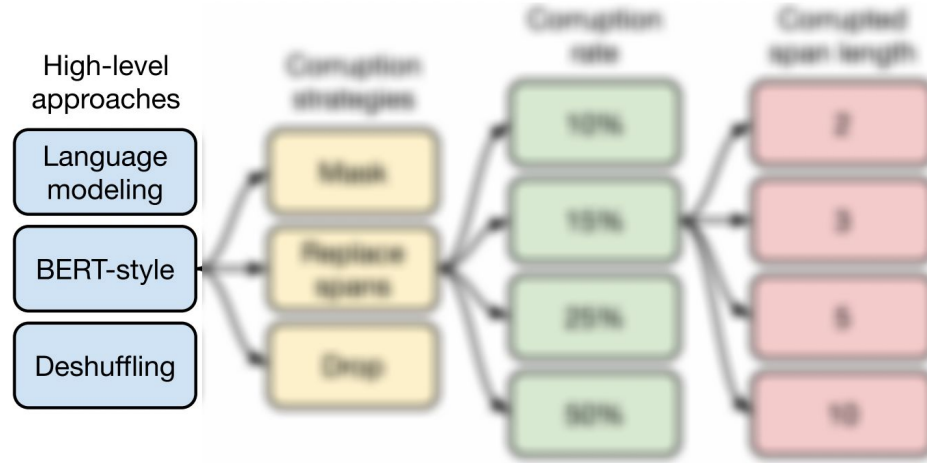


# Different Unsupervised Objectives



Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style Devlin et al. (2018)	Thank you <M> <M> me to your party apple week .	<i>(original text)</i>
Deshuffling	party me for your to . last fun you inviting week Thank	<i>(original text)</i>

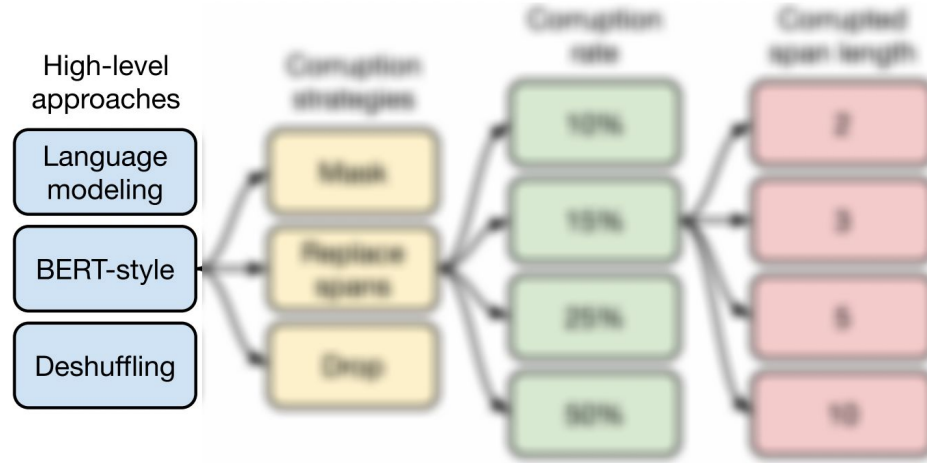
# Different Unsupervised Objectives



Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style Devlin et al. (2018)	Thank you <M> <M> me to your party apple week .	<i>(original text)</i>
Deshuffling	party me for your to . last fun you inviting week Thank	<i>(original text)</i>

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

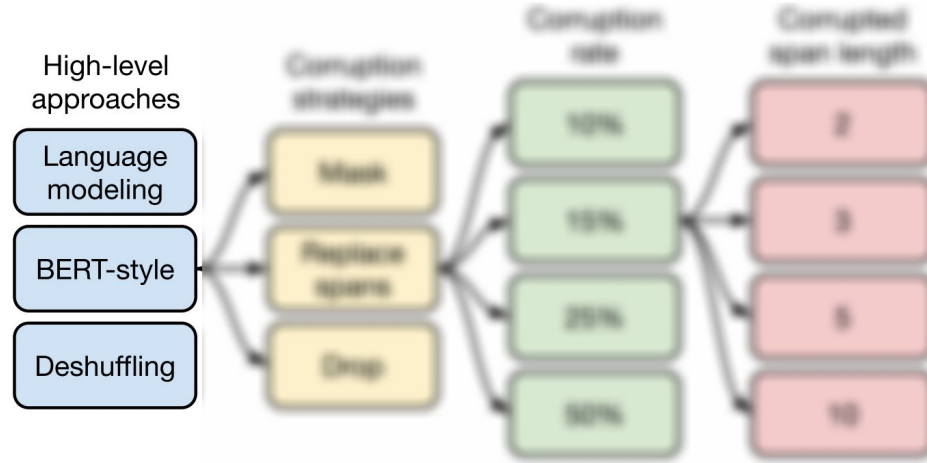
# Different Unsupervised Objectives



Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style Devlin et al. (2018)	Thank you <M> <M> me to your party apple week .	<i>(original text)</i>
Deshuffling	party me for your to . last fun you inviting week Thank	<i>(original text)</i>

- Thank you <M> <M> me to your party apple week . Thank you for inviting me to your party last week.

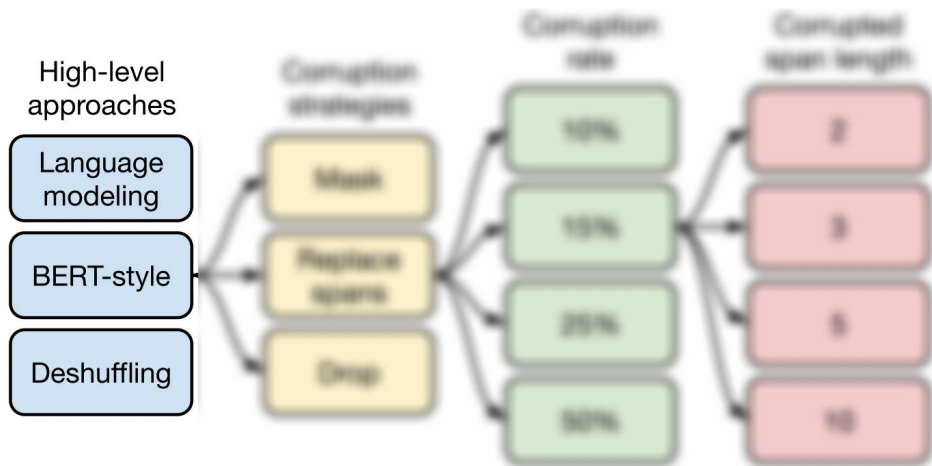
# Different Unsupervised Objectives



Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style Devlin et al. (2018)	Thank you <M> <M> me to your party apple week .	<i>(original text)</i>
Deshuffling	party me for your to . last fun you inviting week Thank	<i>(original text)</i>

- party me for your to . last fun you inviting week Thank .
Thank you for inviting me to your party last week.

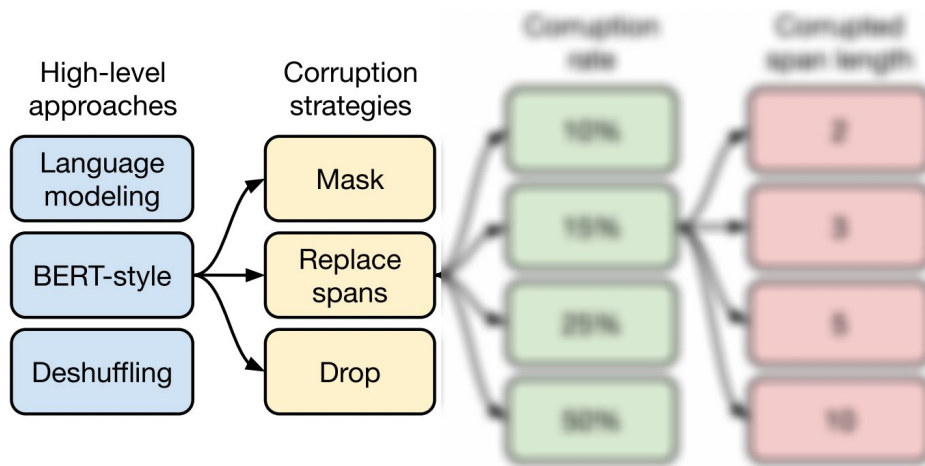
# Performance of the three disparate pre-training objectives



1. BERT-style objective performs best.
2. Prefix LM works well on translation tasks.
3. Deshuffling objective is significantly worse.

Objective	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
Prefix language modeling	80.69	18.94	77.99	65.27	<b>26.86</b>	39.73	<b>27.49</b>
BERT-style (Devlin et al., 2018)	<b>82.96</b>	<b>19.17</b>	<b>80.65</b>	<b>69.85</b>	<b>26.78</b>	<b>40.03</b>	<b>27.41</b>
Deshuffling	73.17	18.59	67.61	58.47	26.11	39.30	25.62

# Different BERT-style Unsupervised Objectives



Objective

Inputs

Targets

MASS-style Song et al. (2019)

I.i.d. noise, replace spans

I.i.d. noise, drop tokens

Thank you <M> <M> me to your party <M> week .

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

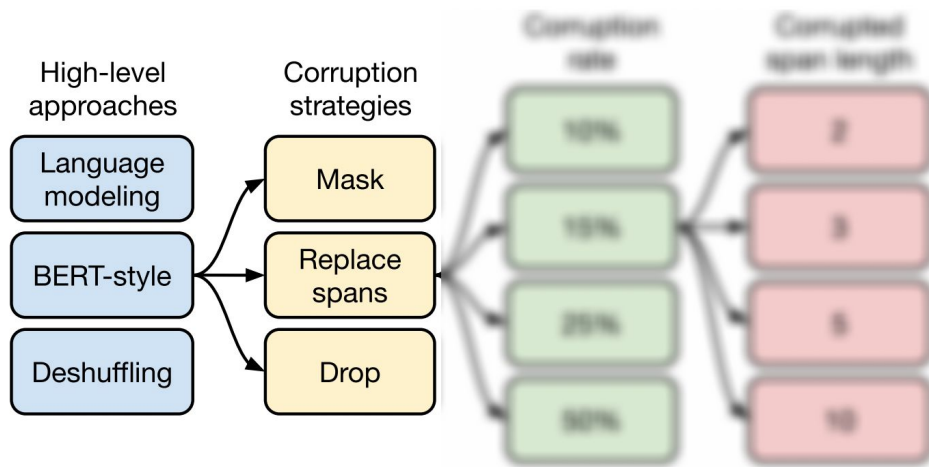
Thank you me to your party week .

*(original text)*

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

for inviting last

# Different BERT-style Unsupervised Objectives



- Thank you <M> <M> me to your party <M> week . Thank you for inviting me to your party last week

MASS-style Song et al. (2019)

Thank you <M> <M> me to your party <M> week .

*(original text)*

I.i.d. noise, replace spans

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

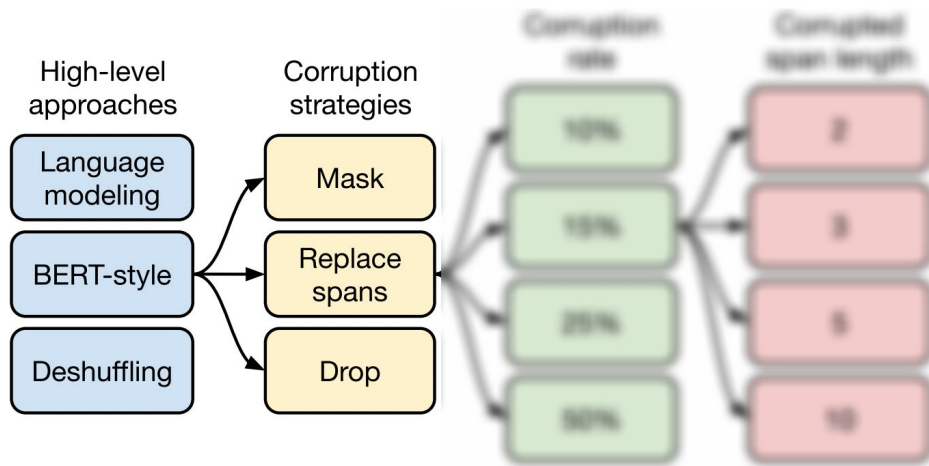
<X> for inviting <Y> last <Z>

I.i.d. noise, drop tokens

Thank you me to your party week .

for inviting last

# Different BERT-style Unsupervised Objectives



- Thank you <X> me to your party <Y> week . <X> for inviting <Y> last <Z>

MASS-style Song et al. (2019)

Thank you <M> <M> me to your party <M> week .

*(original text)*

I.i.d. noise, replace spans

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

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

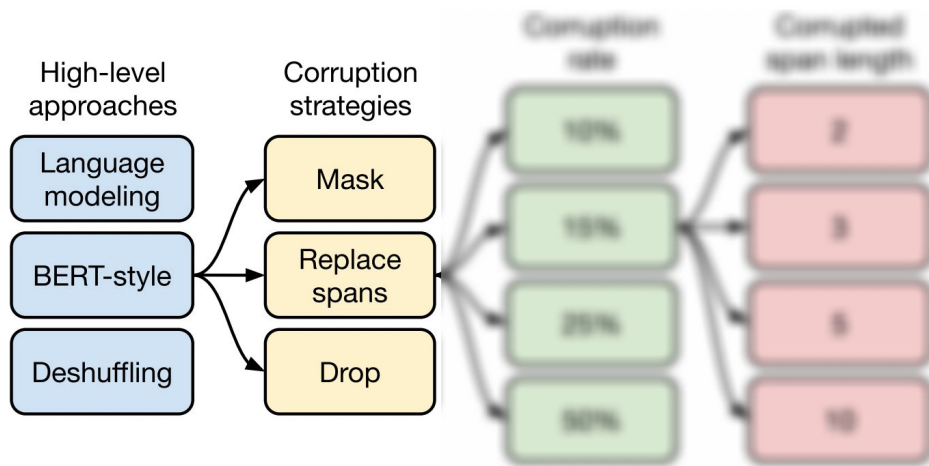
I.i.d. noise, drop tokens

Thank you me to your party week .

for inviting last



# Different BERT-style Unsupervised Objectives



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

MASS-style Song et al. (2019)

I.i.d. noise, replace spans

I.i.d. noise, drop tokens

Thank you <M> <M> me to your party <M> week .

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

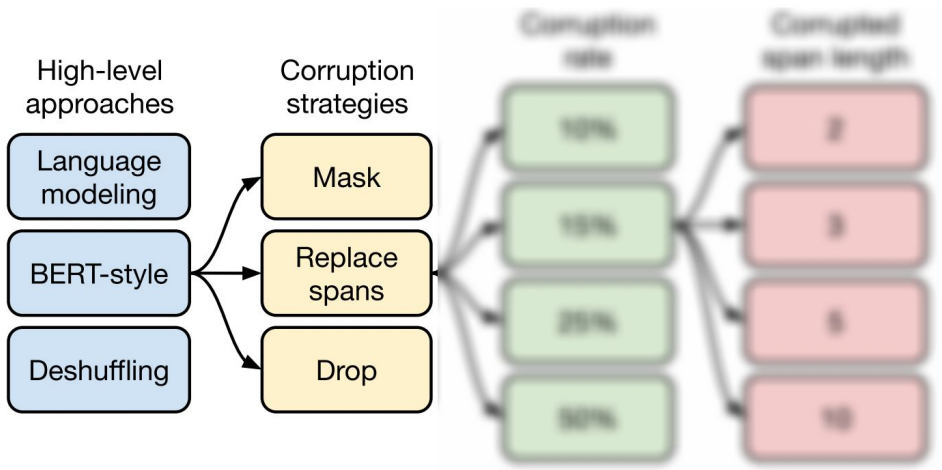
Thank you me to your party week .

*(original text)*

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

for inviting last

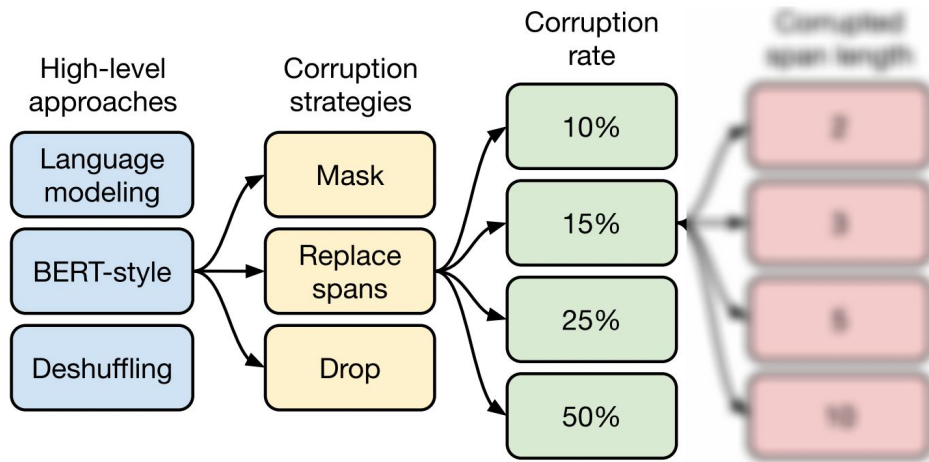
# Comparison of variants of the BERT-style pre-training objective



1. All the variants perform similarly.
2. “Replace corrupted spans” and “Drop corrupted tokens” are more appealing because target sequences are shorter, speeding up training.

Objective	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
BERT-style (Devlin et al., 2018)	82.96	19.17	<b>80.65</b>	69.85	26.78	<b>40.03</b>	27.41
MASS-style (Song et al., 2019)	82.32	19.16	80.10	69.28	26.79	<b>39.89</b>	27.55
★ Replace corrupted spans	83.28	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	<b>27.65</b>
Drop corrupted tokens	<b>84.44</b>	<b>19.31</b>	<b>80.52</b>	68.67	<b>27.07</b>	39.76	<b>27.82</b>

# Different Corruption Rates



Objective

Inputs

Targets

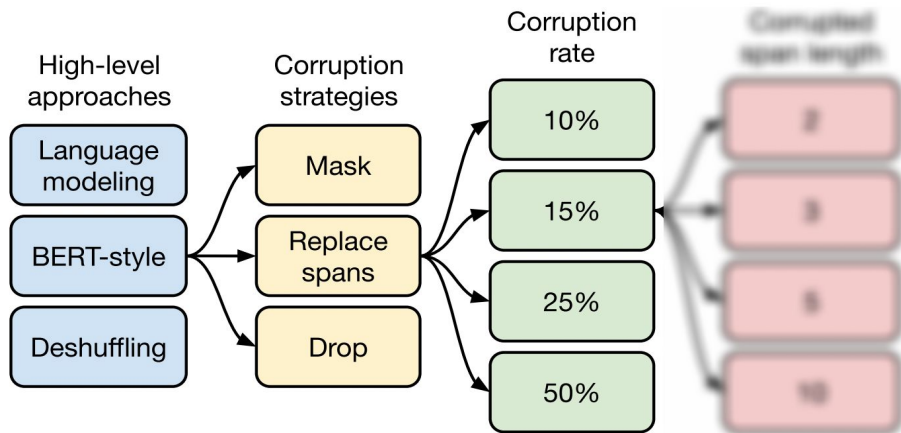
- Thank <X> for inviting me to <Y> party last week. <X> you <Y> your <Z>.
- Thank <X> for <Y> me to your party <Z>. <X> you <Y> inviting <Z> last week.

I.i.d. noise, replace spans

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

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

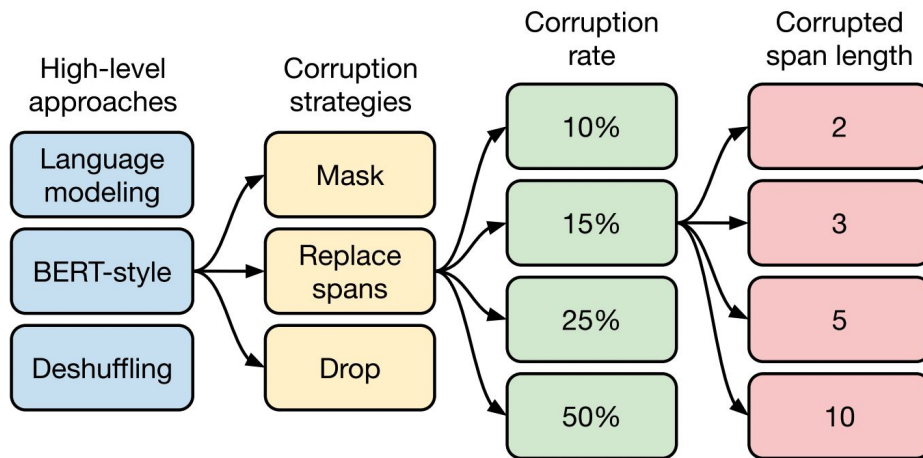
# Performance of the i.i.d. corruption objective with different corruption rates



1. Larger corruption rate leads to downstream performance degradation.
2. Larger corruption rate also leads to longer targets, slowing down training.

Corruption rate	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
10%	<b>82.82</b>	19.00	<b>80.38</b>	69.55	<b>26.87</b>	39.28	<b>27.44</b>
★ 15%	<b>83.28</b>	19.24	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
25%	<b>83.00</b>	<b>19.54</b>	<b>80.96</b>	70.48	<b>27.04</b>	<b>39.83</b>	<b>27.47</b>
50%	81.27	19.32	79.80	70.33	<b>27.01</b>	<b>39.90</b>	<b>27.49</b>

# Different Corruption Rates



---

Objective

Inputs

Targets

---

- Thank <X> for inviting me to <Y> party last <Z>. <X> you <Y> your <Z> week.
  - Thank <X> me to your party week. <X> you for inviting <Z>.

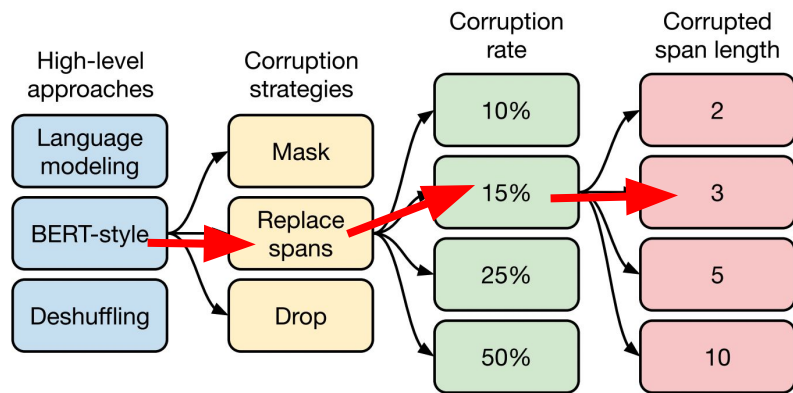
---

Random spans

Thank you <X> to <Y> week .

<X> for inviting me <Y> your party last <Z>

# Performance of the span-corruption objective for different average span lengths



1. Average span length of 3 works well on most non-translation tasks.
2. Span corruption produces shorter target sequences and leads to speedup in training.

Span length	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline (i.i.d.)	<b>83.28</b>	19.24	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
2	<b>83.54</b>	19.39	<b>82.09</b>	<b>72.20</b>	<b>26.76</b>	<b>39.99</b>	<b>27.63</b>
3	<b>83.49</b>	<b>19.62</b>	<b>81.84</b>	<b>72.53</b>	<b>26.86</b>	39.65	<b>27.62</b>
5	<b>83.40</b>	19.24	<b>82.05</b>	<b>72.23</b>	<b>26.88</b>	39.40	<b>27.53</b>
10	82.85	19.33	<b>81.84</b>	70.44	<b>26.79</b>	39.49	<b>27.69</b>

# Pre-training dataset

# Performance from pre-training on different data sets.

Data set	Size	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21
RealNews-like	35GB	<b>83.83</b>	<b>19.23</b>	80.39	②72.38	<b>26.75</b>	<b>39.90</b>	<b>27.48</b>
WebText-like	17GB	<b>84.03</b>	<b>19.31</b>	<b>81.42</b>	71.40	<b>26.80</b>	<b>39.74</b>	<b>27.59</b>
Wikipedia	16GB	①81.85	<b>19.31</b>	81.29	68.01	<b>26.94</b>	39.69	<b>27.67</b>
Wikipedia + TBC	20GB	83.65	<b>19.28</b>	<b>82.08</b>	③73.24	<b>26.77</b>	39.63	<b>27.57</b>

Pre-training on in-domain data tends to help downstream task.

① Much worse on COLA

Check whether a sentence is linguistically correct?

② Much better on ReCoRD

Question answering on News dataset

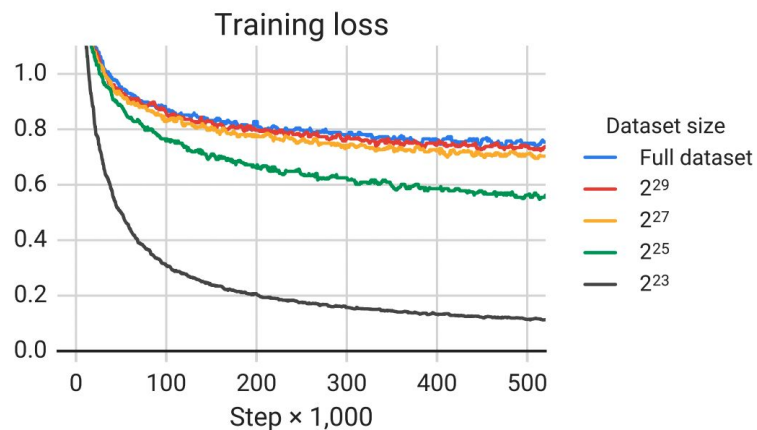
③ Much better on MultiRC

Question answering on Novel dataset



# Effect of repeating data during pre-training

Number of tokens	Repeats	GLUE	CNN3M	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



1. Performance degrades as dataset size shrinks.
2. Model memorizes the pre-training data, with smaller dataset size.

# Scaling

Scaling strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline	83.28	19.24	80.88	71.36	26.98	39.82	27.65
1× size, 4× training steps	85.33	19.33	82.45	74.72	27.08	40.66	27.93
1× size, 4× batch size	84.60	19.42	82.52	74.64	27.07	40.60	27.84
2× size, 2× training steps	<b>86.18</b>	19.66	<b>84.18</b>	77.18	27.52	<b>41.03</b>	28.19
4× size, 1× training steps	<b>85.91</b>	19.73	<b>83.86</b>	<b>78.04</b>	27.47	40.71	28.10
4× ensembled	84.77	<b>20.10</b>	83.09	71.74	<b>28.05</b>	40.53	<b>28.57</b>
4× ensembled, fine-tune only	84.05	19.57	82.36	71.55	27.55	40.22	28.09

1. Advantage in increasing model size compared to simply increasing batch size or number of training steps.
2. Not much of a difference between increasing size + training and increasing size only
  - a. Improving training time and model size are complementary means of improving performance.
3. Ensembling helps, except in SuperGLUE.

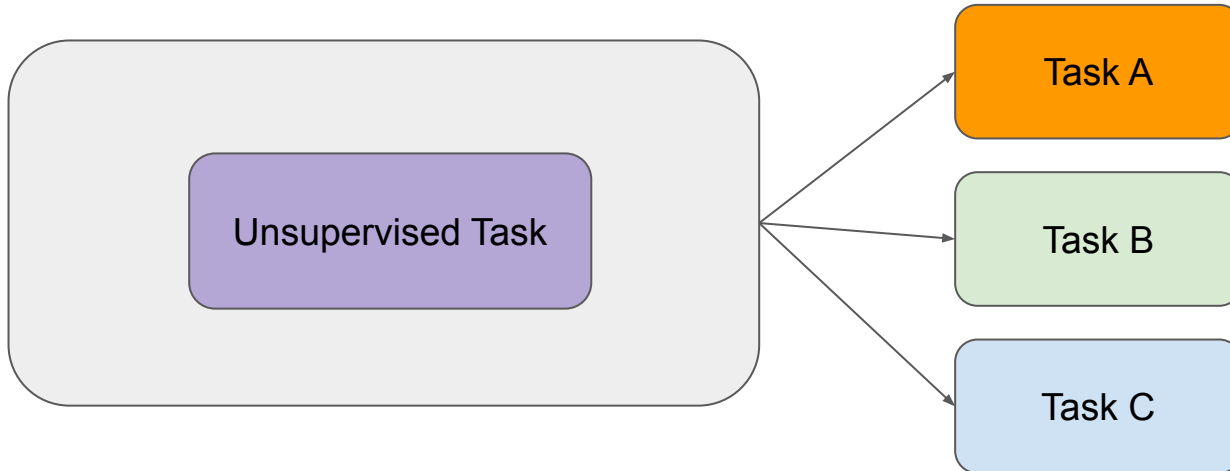
# Multi-task training

# Multi-task

---

Training strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65

---

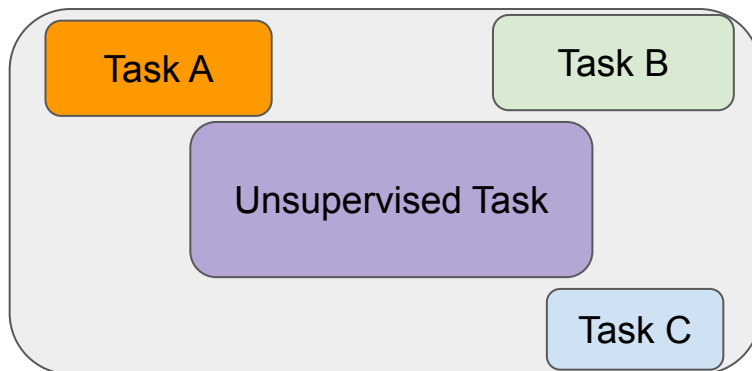


# Multi-task

---

Training strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76

---

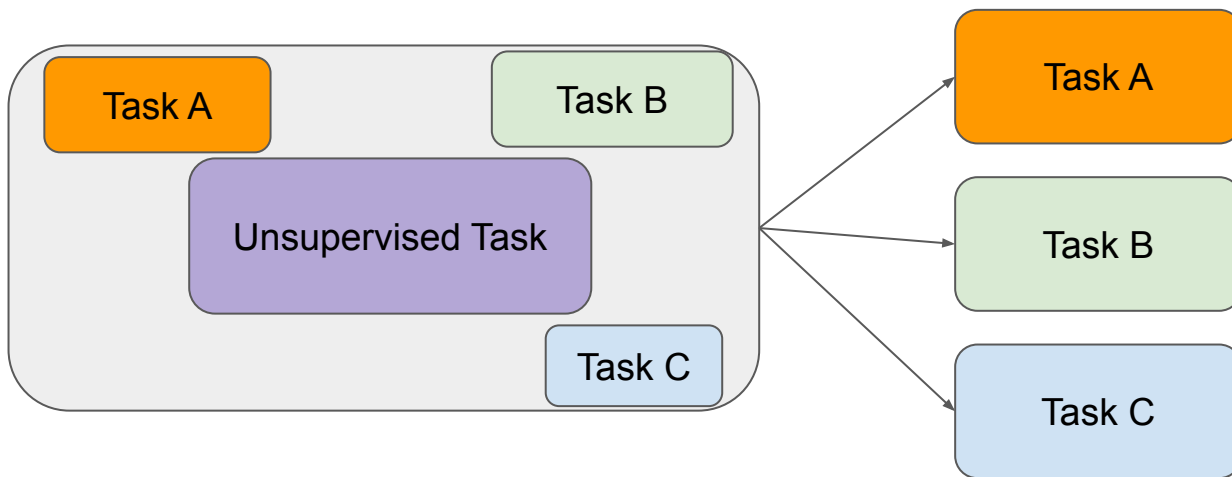


# Multi-task

---

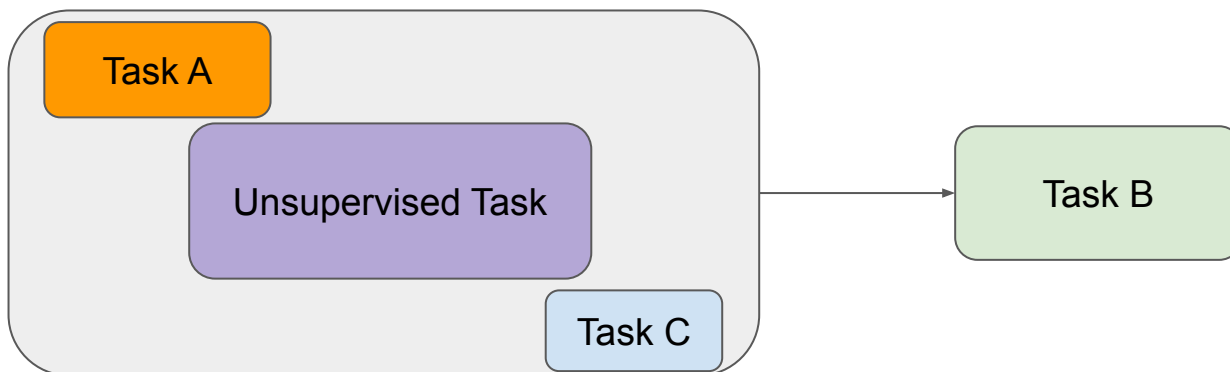
Training strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76
Multi-task pre-training + fine-tuning	<b>83.11</b>	<b>19.12</b>	<b>80.26</b>	<b>71.03</b>	<b>27.08</b>	39.80	<b>28.07</b>

---



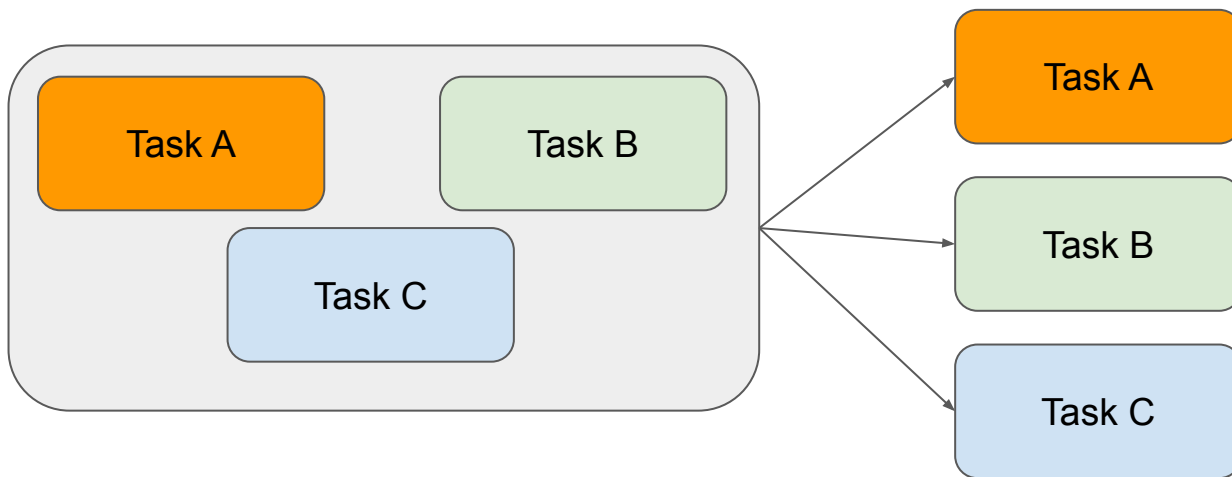
# Multi-task

Training strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76
Multi-task pre-training + fine-tuning	<b>83.11</b>	<b>19.12</b>	<b>80.26</b>	<b>71.03</b>	<b>27.08</b>	39.80	<b>28.07</b>
Leave-one-out multi-task training	81.98	19.05	79.97	<b>71.68</b>	<b>26.93</b>	39.79	<b>27.87</b>



# Multi-task

Training strategy	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76
Multi-task pre-training + fine-tuning	<b>83.11</b>	<b>19.12</b>	<b>80.26</b>	<b>71.03</b>	<b>27.08</b>	39.80	<b>28.07</b>
Leave-one-out multi-task training	81.98	19.05	79.97	<b>71.68</b>	<b>26.93</b>	39.79	<b>27.87</b>
Supervised multi-task pre-training	79.93	18.96	77.38	65.36	26.81	<b>40.13</b>	<b>28.04</b>



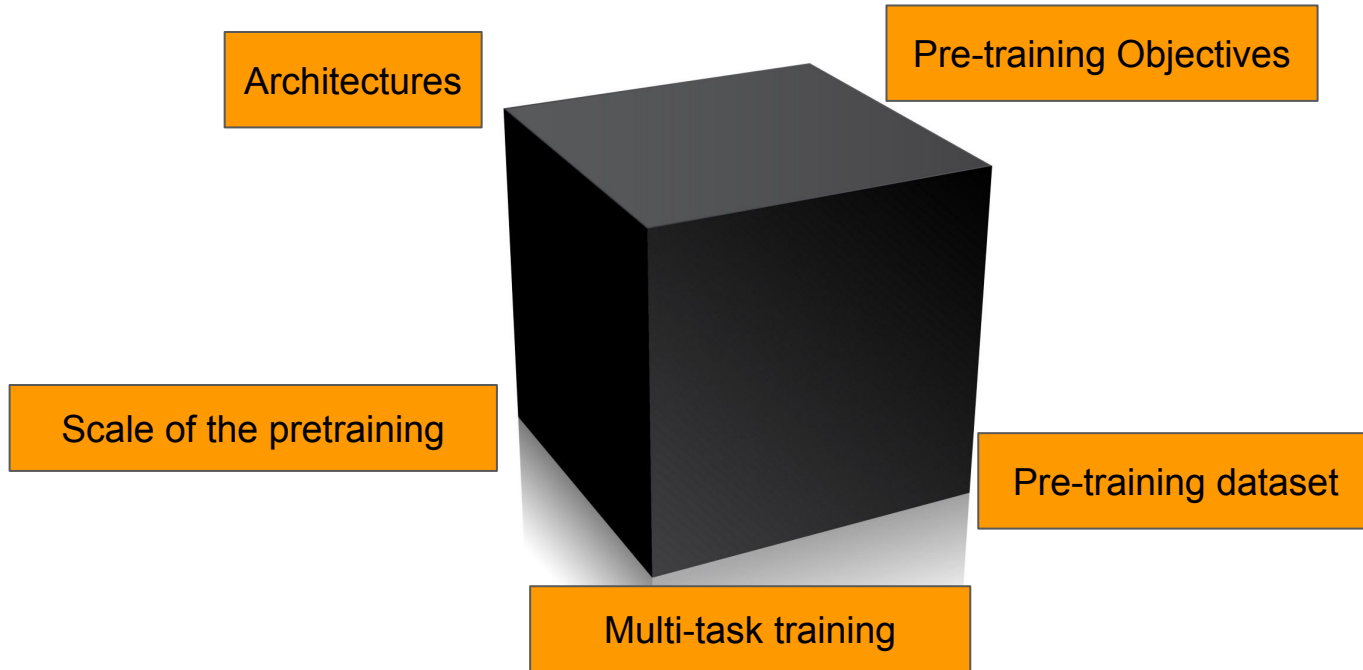


# Multi-task

Training strategy	GLUE	CNN4	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76
Multi-task pre-training + fine-tuning	<b>83.11</b>	<b>19.12</b>	<b>80.26</b>	<b>71.03</b>	<b>27.08</b>	39.80	<b>28.07</b>
Leave-one-out multi-task training	81.98	19.05	79.97	<b>71.68</b>	<b>26.93</b>	39.79	<b>27.87</b>
Supervised multi-task pre-training	79.93	18.96	77.38	65.36	26.81	<b>40.13</b>	<b>28.04</b>

1. Multi-task pre-training + fine-tuning works as well as unsupervised pre-training + fine-tuning.
2. Practical benefit of Multi-task pre-training + fine-tuning is to monitor downstream performance during pre-training.

# Putting it all together



## Encoder-decoder architecture

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<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>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
Enc-dec, 6 layers	Denoising	$P$	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	$P$	$M$	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	$P$	$M$	81.82	18.61	78.94	68.11	26.43	37.98	27.39

## Span prediction objective

	Span length	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline (i.i.d.)		<b>83.28</b>	19.24	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
	2	<b>83.54</b>	19.39	<b>82.09</b>	<b>72.20</b>	<b>26.76</b>	<b>39.99</b>	<b>27.63</b>
	3	<b>83.49</b>	<b>19.62</b>	<b>81.84</b>	<b>72.53</b>	<b>26.86</b>	39.65	<b>27.62</b>
	5	<b>83.40</b>	19.24	<b>82.05</b>	<b>72.23</b>	<b>26.88</b>	39.40	<b>27.53</b>
	10	82.85	19.33	<b>81.84</b>	70.44	<b>26.79</b>	39.49	<b>27.69</b>

## C4 dataset

Data set	Size	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21
RealNews-like	35GB	<b>83.83</b>	<b>19.23</b>	80.39	72.38	<b>26.75</b>	<b>39.90</b>	<b>27.48</b>
WebText-like	17GB	<b>84.03</b>	<b>19.31</b>	<b>81.42</b>	71.40	<b>26.80</b>	<b>39.74</b>	<b>27.59</b>
Wikipedia	16GB	81.85	<b>19.31</b>	81.29	68.01	<b>26.94</b>	39.69	<b>27.67</b>
Wikipedia + TBC	20GB	83.65	<b>19.28</b>	<b>82.08</b>	<b>73.24</b>	<b>26.77</b>	39.63	<b>27.57</b>

## Multi-task pre-training

Training strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Unsupervised pre-training + fine-tuning	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	27.65
Multi-task training	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	27.76
Multi-task pre-training + fine-tuning	<b>83.11</b>	<b>19.12</b>	<b>80.26</b>	<b>71.03</b>	<b>27.08</b>	39.80	<b>28.07</b>
Leave-one-out multi-task training	81.98	19.05	79.97	<b>71.68</b>	<b>26.93</b>	39.79	<b>27.87</b>
Supervised multi-task pre-training	79.93	18.96	77.38	65.36	26.81	<b>40.13</b>	<b>28.04</b>

## Bigger model trained longer

Scaling strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline	83.28	19.24	80.88	71.36	26.98	39.82	27.65
1× size, 4× training steps	85.33	19.33	82.45	74.72	27.08	40.66	27.93
1× size, 4× batch size	84.60	19.42	82.52	74.64	27.07	40.60	27.84
2× size, 2× training steps	<b>86.18</b>	19.66	<b>84.18</b>	77.18	27.52	<b>41.03</b>	28.19
4× size, 1× training steps	<b>85.91</b>	19.73	<b>83.86</b>	<b>78.04</b>	27.47	40.71	28.10
4× ensembled	84.77	<b>20.10</b>	83.09	71.74	<b>28.05</b>	40.53	<b>28.57</b>
4× ensembled, fine-tune only	84.05	19.57	82.36	71.55	27.55	40.22	28.09

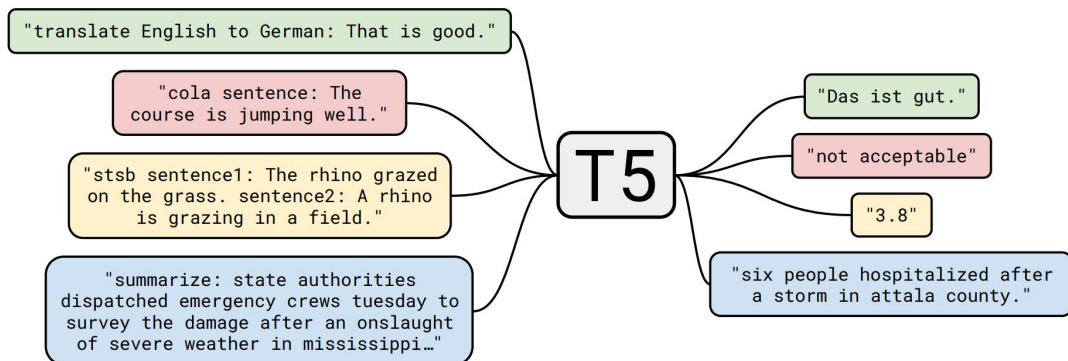
# Model Variants

Model	Parameters	No. of layers	$d_{\text{model}}$	$d_{\text{ff}}$	$d_{\text{kv}}$	No. of heads
Small	60M	6	512	2048	64	8
Base	220M	12	768	3072	64	12
Large	770M	24	1024	4096	64	16
3B	3B	24	1024	16384	128	32
11B	11B	24	1024	65536	128	128

Model	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
Previous best	89.4	20.30	95.5	84.6	<b>33.8</b>	<b>43.8</b>	<b>38.5</b>
T5-Small	77.4	19.56	87.24	63.3	26.7	36.0	26.8
T5-Base	82.7	20.34	92.08	76.2	30.9	41.2	28.0
T5-Large	86.4	20.68	93.79	82.3	32.0	41.5	28.1
T5-3B	88.5	21.02	94.95	86.4	31.8	42.6	28.2
T5-11B	<b>89.7</b>	<b>21.55</b>	<b>95.64</b>	<b>88.9</b>	32.1	43.4	28.1

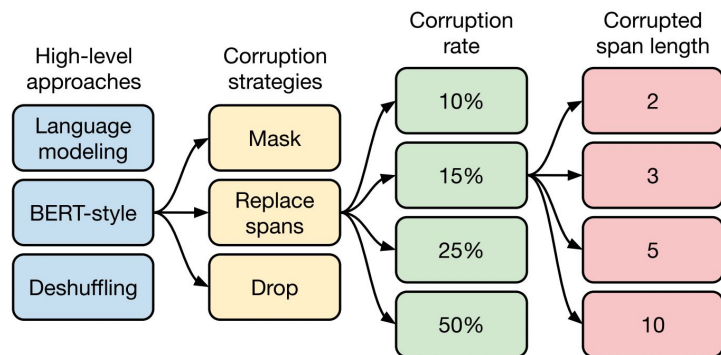
# Let's review!

- **Unified text-to-text framework**
- Supports both discriminative and generative tasks
  - Classification, summarization, translation, etc.
  - Better on GLUE/SuperGLUE, SQuAD, and summarization; less on translation



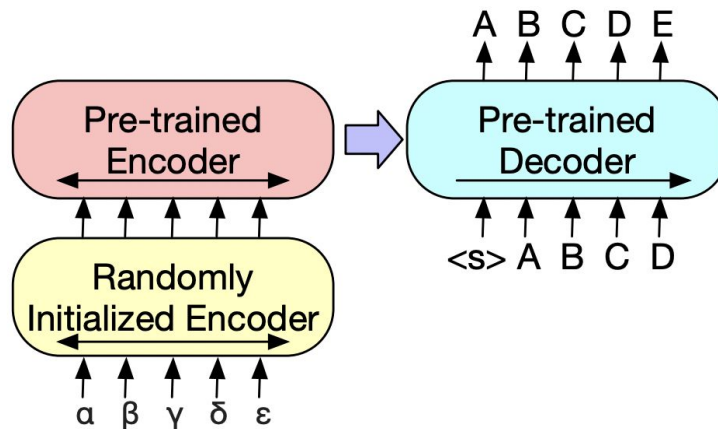
# Let's review!

- “Empirical comparison of existing techniques”
  - Evidence for encoder-decoder models, span masking, multi-task pre-training
- Still no limit on large model improvements?
- C4 as a large, clean corpus



# Other Variants

- BART (Lewis et al. 2020)
  - Similar Architecture as T5.
  - Performs competitive to RoBERTa and XLNet on discriminative tasks.
  - Outperformed existing methods on question answering, and summarization tasks.
  - Improved results on machine translation with fine-tuning on target language.
- mT5 (Xue et al. 2021)
  - Discussed earlier!



# Other Variants

- AlexaTM 20B (Soltan et al. 2022)
  - Larger architecture on multilingual C4 dataset.
  - Can outperform much larger autoregressive models (GPT-3 175B) in zero shot tasks.

Model	BoolQ (acc)	CB (acc)	RTE (acc)	ReCoRD (acc)	WSC (acc)	WiC (acc)	CoPA (acc)	MultiRC (f1a)	Avg
PaLM 540B	<b>88.0</b>	<u>51.8</u>	<b>72.9</b>	<b>92.9</b>	<b>89.1</b>	<b>59.1</b>	<b>93.0</b>	<b>83.5</b>	<b>78.8</b>
GPT3 175B	60.5	46.4	63.5	<u>90.2</u>	65.4	0.0	<u>91.0</u>	<u>72.9</u>	61.2
BLOOM 175B	63.5	33.9	52.0	NA	51.9	50.6	56.0	57.1	NA
GPT3 13B	66.2	19.6	62.8	89.0	64.4	0.0	84.0	71.4	57.2
UL 20B	63.1	41.1	60.7	88.1	<u>79.9</u>	49.8	85.0	36.2	63.0
AlexaTM 20B	<u>69.44</u>	<b>67.9</b>	<u>68.59</u>	88.4	68.27	<u>53.29</u>	78.0	59.57	<u>69.16</u>



# Q1. Describe how T5 is adapted to sentence classification tasks

[Task-specific prefix]: [Input text]

- CoLA (GLUE; Classification):  
“cola sentence: The course is jumping well.” -> “not acceptable”
- STS-B (GLUE; Regression):  
“stsb sentence1: The rhino grazed. sentence2: A rhino is grazing.” -> “3.8”

“cola sentence: The course is jumping well.” -> “hamburger”

**“Hamburger” is not a valid CoLA output, so this is a fail!**

Q2. Can you think of a reason why generating the entire output performs worse than only generating the masked spans?

Objective	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
BERT-style (Devlin et al., 2018)	82.96	19.17	<b>80.65</b>	69.85	26.78	<b>40.03</b>	27.41
MASS-style (Song et al., 2019)	82.32	19.16	80.10	69.28	26.79	<b>39.89</b>	27.55
★ Replace corrupted spans	83.28	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	<b>27.65</b>
Drop corrupted tokens	<b>84.44</b>	<b>19.31</b>	<b>80.52</b>	68.67	<b>27.07</b>	39.76	<b>27.82</b>

Q3. Would you expect a BERT encoder or a T5 encoder to learn richer linguistic features, assuming both were the same size and trained for the same number of steps? How would it change if the average masked span length was increased?