

# Introduction to T5

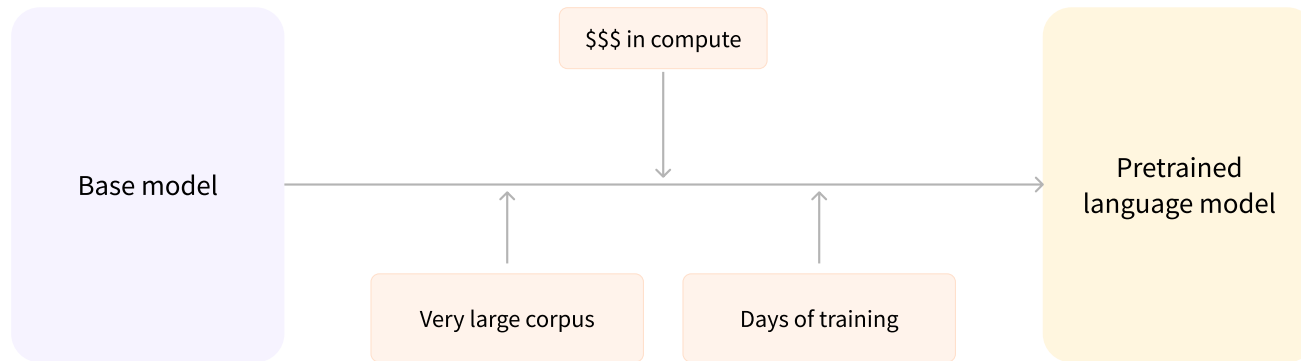
## Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer

Presented by Zengzhi Wang

# Background

- ✓ NLP aims to teach machine general language understanding ability.
- ✓ A historically common approach is to word vectors to map a word to a continuous representation.
- ✓ Recently, pre-training then fine-tuning becomes increasingly popular.

## *Pre-training*



Upstream

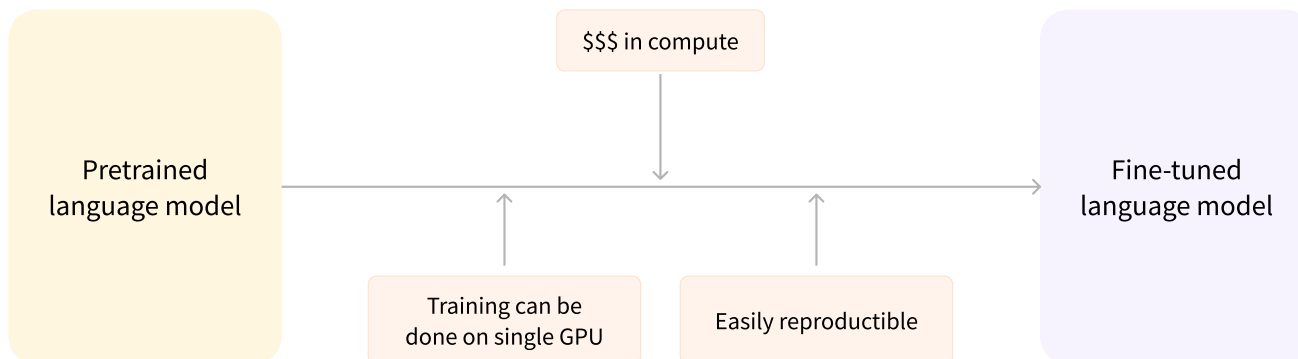


Transfer Learning



Downstream

## *Fine-tuning*



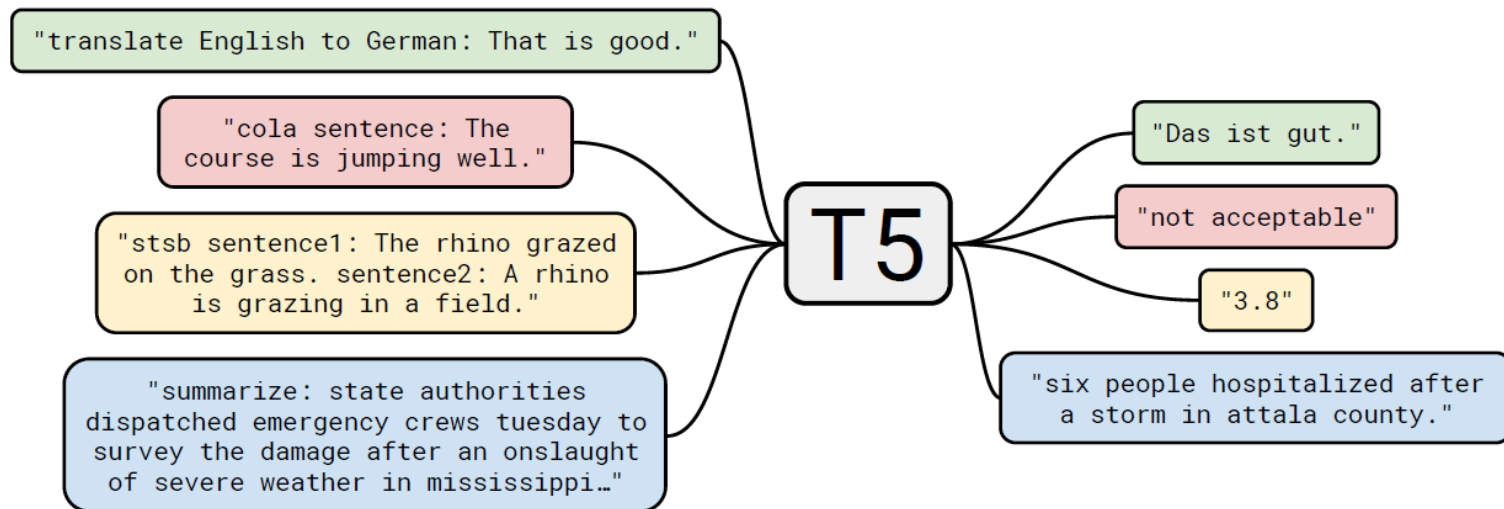
# Background

- Recent works on transfer learning have produced a wide landscape of pre-training objectives, unlabeled datasets, benchmarks, fine-tuning methods, and more.
- The rapid progress make it difficult to compare different algorithms and understand the space of existing methods for transfer learning.

# Motivation

- ✓ Explore the limit of transfer learning.
- ✓ Specifically, leverage a unified text-to-text transformer to systematically study different approaches (pre-training objectives, unlabeled datasets, and other factors) and push the current limits of the field by scaling up models and datasets.

*A comprehensive study on transfer learning for NLP*



# Outline

- Baseline Design
  - Model Architecture
  - Pre-training Corpus
  - Tokenizer/Vocabulary
  - Pre-training Unsupervised Objective
  - Downstream Tasks & IO Format
  - Baseline Performance
- Effect of Model Architectures
- Effect of Unsupervised Objectives
- Effect of Pre-training Dataset
- Effect of Training Strategies
- Effect of Scaling
- Putting it all together
- Takeaways & Outlook

# 1. Baseline Design

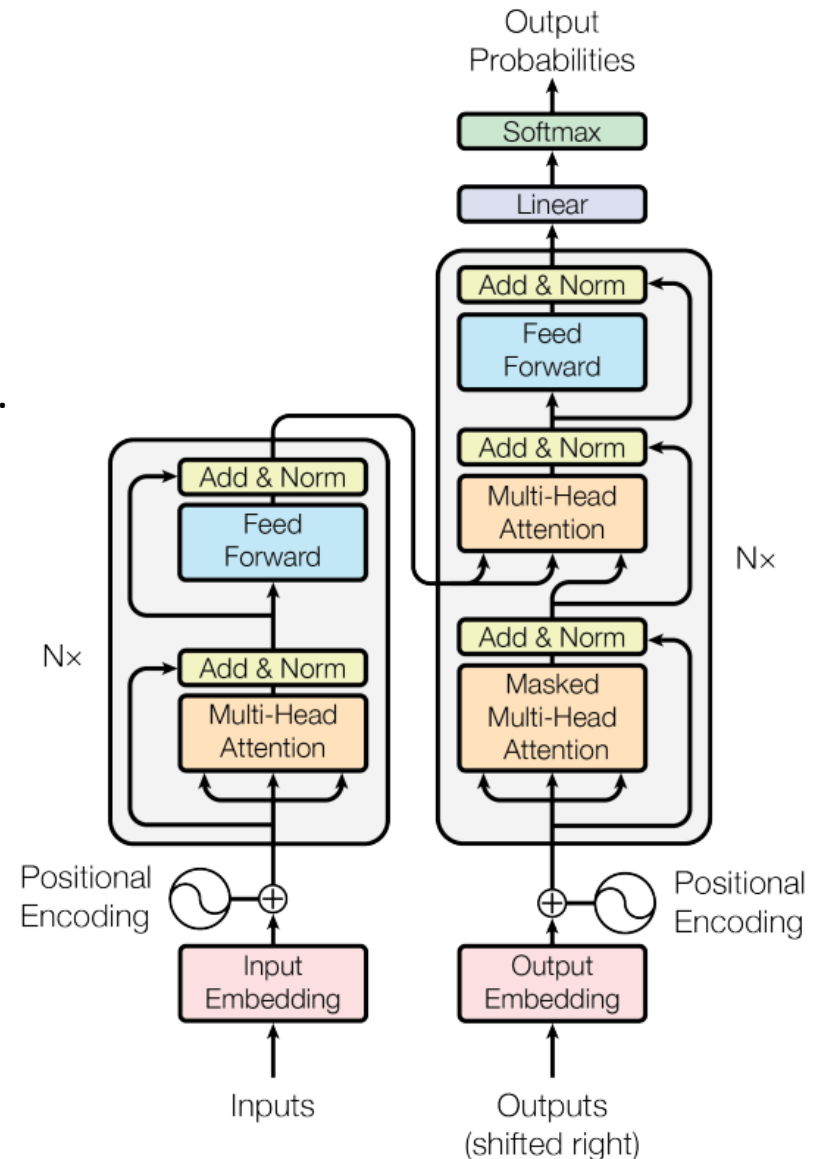
# Model Architecture

Three Modifications:

- ✓ Removing the Layer Norm bias.
- ✓ Placing the layer normalization outside the residual path.
- ✓ Using a different position embedding scheme.
  - A relative position embedding;
  - A scalar that is add to the corresponding logit;
  - 32 embeddings. And assign all relative positions beyond 128 to the same embedding.

Specifically, (T5-Base)

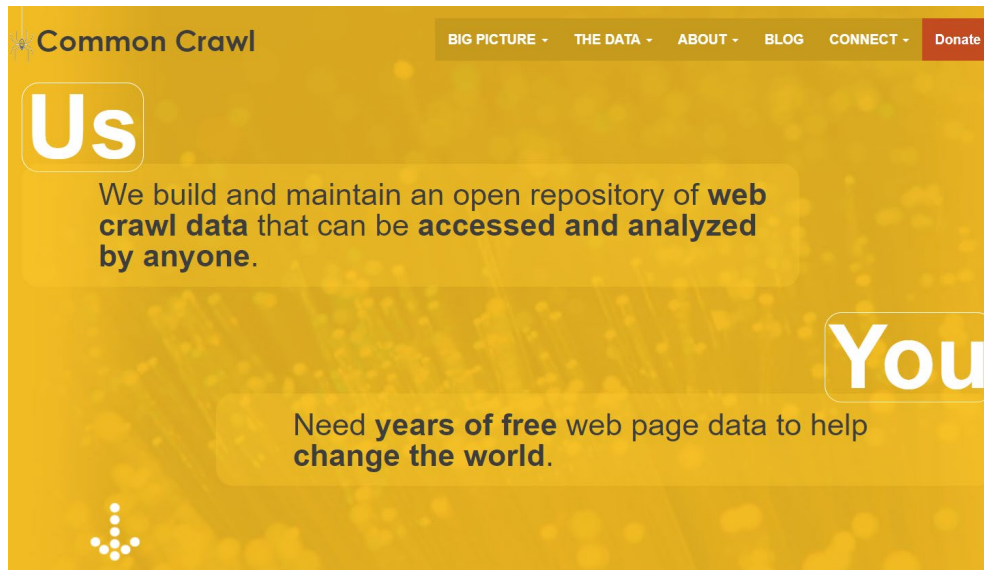
- 12 blocks for Encoder, 12 blocks for Decoder;
- 3072 for FFN hidden-state dimension;
- 768 for hidden-state dimension;
- 12 heads
- 220M parameters



Standard Transformer

# Colossal Clean Crawled Corpus (C4)

- To measure the effect of the quality, characteristics, and size of unlabeled data.
- Common Crawl is a publicly-available web archive that provides “web extracted text” by removing markup and other non-text content from HTML files.



*Unfortunately, the majority of the resulting text is not natural language.*

- Authors use some heuristics rules for cleaning up it, naming the resulting dataset “Colossal Clean Crawled Corpus” (more clean and natural, about 750GB)



# Vocabulary

- ✓ SentencePiece
- ✓ Use a vocabulary of 32000 wordpieces.
- ✓ Multi-lingual corpus, English: German: French: Romanian = 10:1:1:1

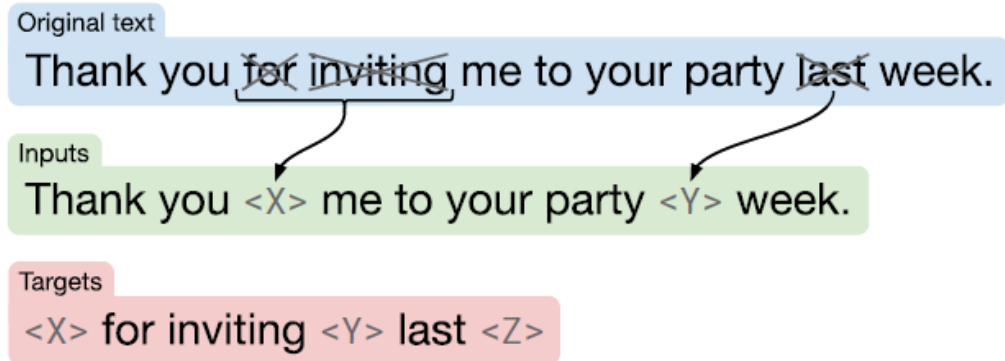
## Intro to SentencePiece:

- ✓ It considers the text as a sequence of Unicode characters, and replaces spaces with a special character “\_”.
- ✓ It does not require a pre-tokenization step, which is very useful for languages where the space character is not used (like Chinese and Japanese).
- ✓ It is a reversible tokenization: there is no special treatment of spaces.

```
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained('t5-base')
tokenizer("This is the Hugging Face course.").tokens()
```

```
>>> ['_This', '_is', '_the', '_Hug', 'ging', '_Face', '_course', '.', '</s>']
```

# Pre-training Objective



- ✓ Randomly samples and then drops out 15% of tokens.
- ✓ All consecutive spans are replaced by a single sentinel token.
- ✓ Instead of [MASK], each sentinel token (e.g., <extra\_id\_0>) is unique to the sequence;
- ✓ Only predict the corresponding sentinel token followed by dropped-out tokens;
- ✓ Final sentinel token is used to mark the end of target sequence.

# Training

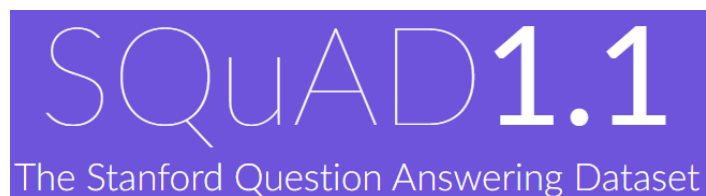
## Details:

- ✓ Number of pre-training steps:  $2^{19}$ ;
- ✓ Max sequence length: 512 ( $2^9$ );
- ✓ Batch size: 128 ( $2^7$ );
- ✓ In total, these corresponds to pre-training on  $2^{35}$  (34B) tokens;
- ✓ Note that  $2^{35}$  only covers a fraction of the entire C4 dataset (w/o any repeat).

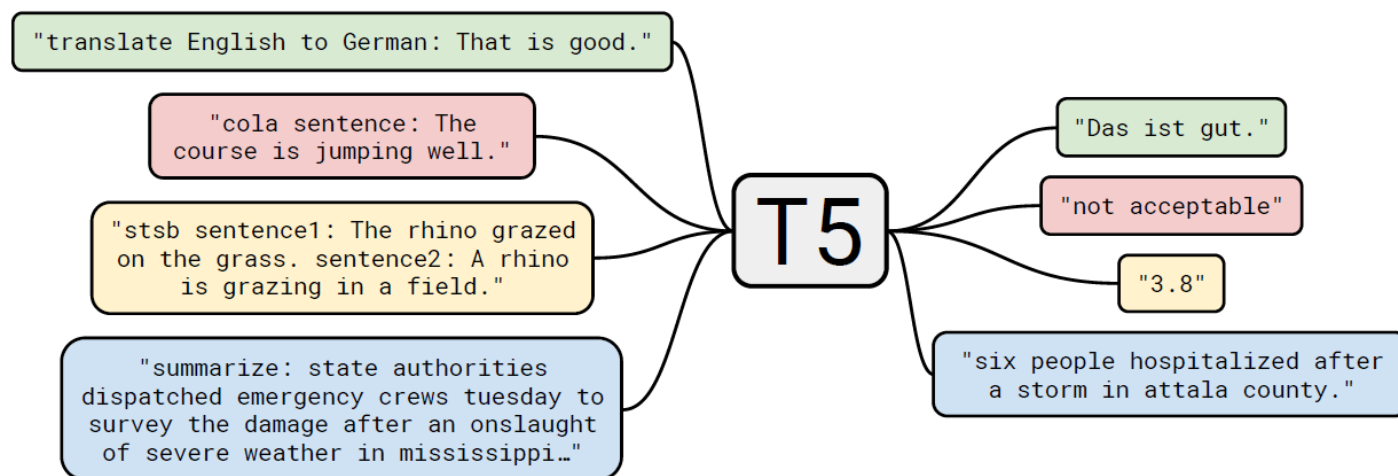
## Comparison:

- ✓ BERT used roughly 137B tokens;
- ✓ RoBERTa used 2.2T tokens.
- ✓ Wow, this baseline only need a reasonable computational budget.

# Downstream Tasks & IO Format



CNN/Daily Mail  
WMT English to German, French,  
and Romanian translation.



*Count the output as wrong when it does not fall into predefined label texts, though authors never observed this behavior.*

# Baseline Performance

	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

Table 1: Average and standard deviation of scores achieved by our baseline model and training procedure. For comparison, we also report performance when training on each task from scratch (i.e. without any pre-training) for the same number of steps used to fine-tune the baseline model. All scores in this table (and every table in our paper except Table 14) are reported on the validation sets of each data set.

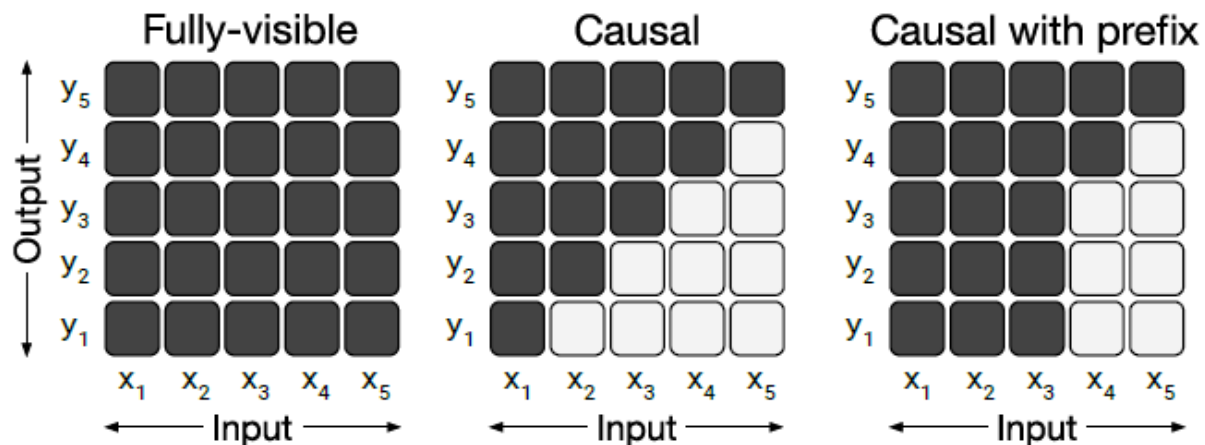
- ✓ Baseline vs. no pre-training variant.
- ✓ Baseline yields comparable performance to existing models with similar size.
- ✓ Can not directly compare this baseline to BERT-Base because it is an encoder-decoder model and was pre-trained for roughly  $\frac{1}{4}$  as many steps.

## 2. Effect of Architectures

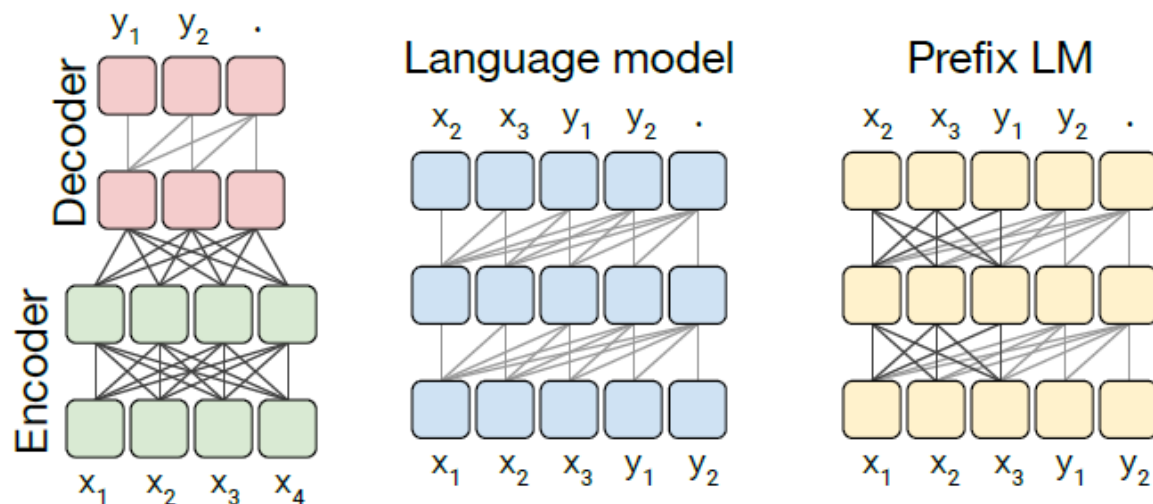
# Recall

*A core factor for different architectures is the “mask” used by different attention mechanisms in the model.*

Attention  
Matrix



Transformer  
Architecture



# Comparison among Different Model Structure

Model	# Layers	# Param.	# FLOPs
BERT-Base	L	P	\
Enc.-Dec.	L for Enc. & L for Dec.	2P	M
Enc.-Dec.	L for Enc. & L for Dec.	P (sharing)	M
Enc-Dec.	L/2 for Enc. & L/2 for Dec.	P	M/2
Dec.-only	L	P	M
Dec.-only w/ Prefix	L	P	M

- ✓ For the computational cost, L layers in the decoder-only model must be applied to both the input and output sequence, while the encoder is only applied to the input sequence and the decoder is only applied to the output sequence. They are approximately equivalent.



# Performance Comparison

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
Encoder-decoder	LM	$2P$	$M$	79.56	18.59	76.02	64.29	26.27	39.17	26.86
Enc-dec, shared	LM	$P$	$M$	79.60	18.13	76.35	63.50	26.62	39.17	27.05
Enc-dec, 6 layers	LM	$P$	$M/2$	78.67	18.26	75.32	64.06	26.13	38.42	26.89
Language model	LM	$P$	$M$	73.78	17.54	53.81	56.51	25.23	34.31	25.38
Prefix LM	LM	$P$	$M$	79.68	17.84	76.87	64.86	26.28	37.51	26.76

➤ How to perform *denoising* objective?

- Enc.-Dec. /Prefix-LM: Thank you <X> me to your party last week. => <X> for inviting <Y>
- Causal LM: concatenate the input and target.



➤ How to perform *language modeling* objective?

- Enc.-Dec./Prefix-LM: Thank you for inviting me to your party last week.
- Causal LM: generate from beginning to end in an autoregressive manner.

# Performance Comparison

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
Encoder-decoder	LM	$2P$	$M$	79.56	18.59	76.02	64.29	26.27	39.17	26.86
Enc-dec, shared	LM	$P$	$M$	79.60	18.13	76.35	63.50	26.62	39.17	27.05
Enc-dec, 6 layers	LM	$P$	$M/2$	78.67	18.26	75.32	64.06	26.13	38.42	26.89
Language model	LM	$P$	$M$	73.78	17.54	53.81	56.51	25.23	34.31	25.38
Prefix LM	LM	$P$	$M$	79.68	17.84	76.87	64.86	26.28	37.51	26.76

- ✓ The Enc.-Dec. architecture with the denoising objective performs best.
- ✓ The parameter-sharing variant performs nearly as well.
- ✓ Halving the number of layers significantly hurt performance.
- ✓ Enc.-Dec. outperforms prefix-LM, suggesting enc.-dec. attention is beneficial.
- ✓ Denoising objective always performs better than LM objective on downstream tasks.

### 3. Effect of Unsupervised Objectives<sub>(Based on Enc.-Dec.)</sub>

# High-Level Comparison

Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style <a href="#">Devlin et al. (2018)</a>	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>



Objective	GLUE	CNNDM	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 ( <a href="#">Devlin et al., 2018</a> )	<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

# Exploration on BERT-style Objective

Recall BERT-style objective:

- ✓ mask 15% tokens, where 90% of tokens are replaced with [MASK], and 10% are replaced with a random token.

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 .	(original text)
Deshuffling	party me for your to . last fun you inviting week Thank	(original text)
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
Random spans	Thank you <X> to <Y> week .	<X> for inviting me <Y> your party last <Z>

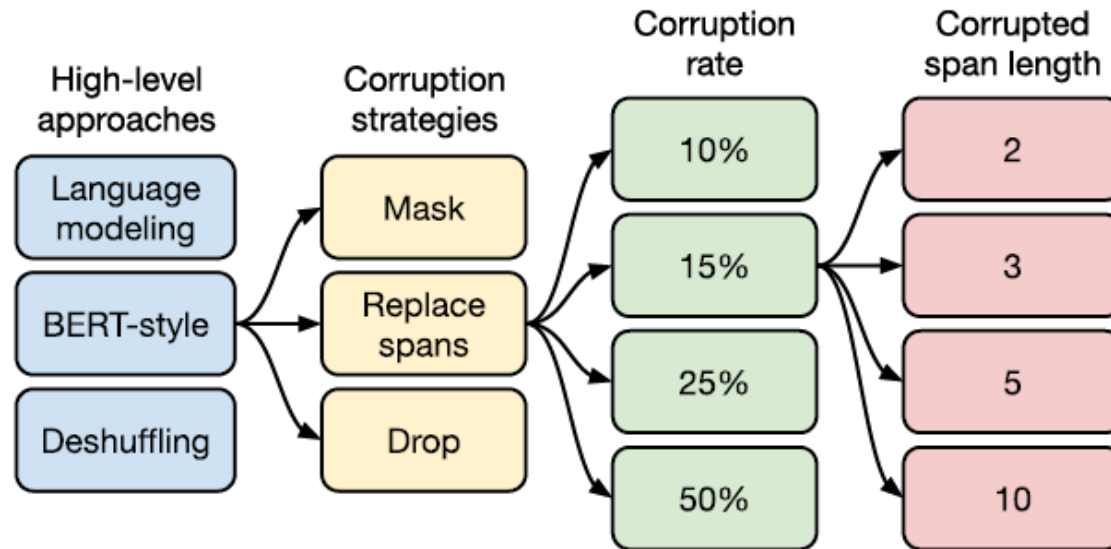
Objective	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
BERT-style (Devlin et al., 2018)	82.96	19.17	<b>80.65</b>	69.85	26.78	<b>40.03</b>	27.41
MASS-style (Song et al., 2019)	82.32	19.16	80.10	69.28	26.79	<b>39.89</b>	27.55
★ Replace corrupted spans	83.28	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>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

# Exploration on Replacing Corrupted Spans

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>

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>

# Overview of Explorations on Objectives



- ✓ Denoising objectives significantly outperform language modeling and de-shuffling.
- ✓ There is not a remarkable performance across variants of denoising objectives.

## 4. Effect of Pre-training Dataset (based on baseline)



# Different Potential Datasets

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>

- ✓ Removing the heuristic filtering from C4 degrades performance.
- ✓ In some cases, it's better to pretrain on a specific domain than the diverse C4 dataset.
- ✓ For example, pretraining on Wikipedia + TBC boosts the model performance on SQuAD and SuperGLUE.

# Different Dataset Size (based on C4)

*With the same pre-training computational budget ( $2^{35}$  tokens)*

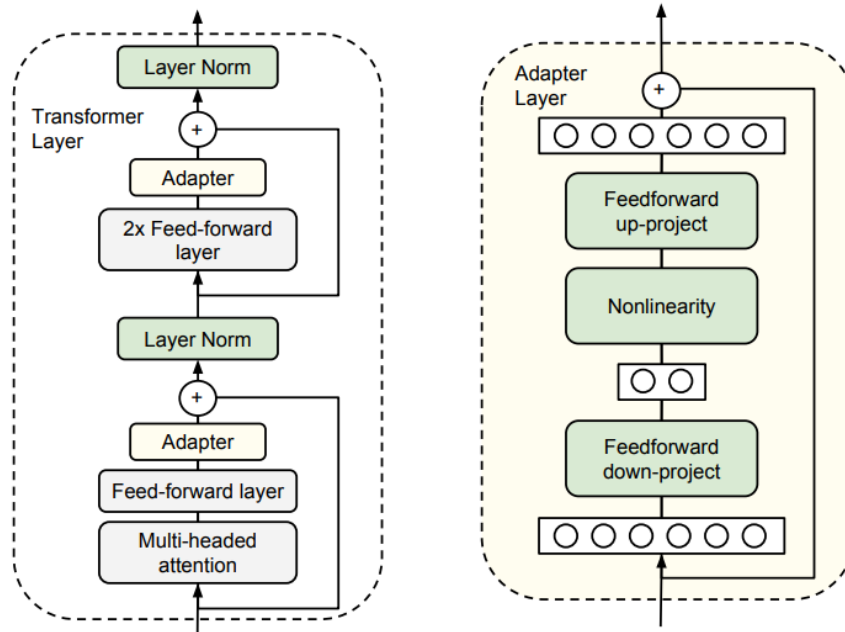
Number of tokens	Repeats	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Full data set	0	83.28	19.24	80.88	71.36	26.98	39.82	27.65
$2^{29}$	64	82.87	19.19	80.97	72.03	26.83	39.74	27.63
$2^{27}$	256	82.62	19.20	79.78	69.97	27.02	39.71	27.33
$2^{25}$	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
$2^{23}$	4,096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

Table 9: Measuring the effect of repeating data during pre-training. In these experiments, we only use the first  $N$  tokens from C4 (with varying values of  $N$  shown in the first column) but still pre-train over  $2^{35}$  tokens. This results in the data set being repeated over the course of pre-training (with the number of repeats for each experiment shown in the second column), which may result in memorization (see Figure 6).

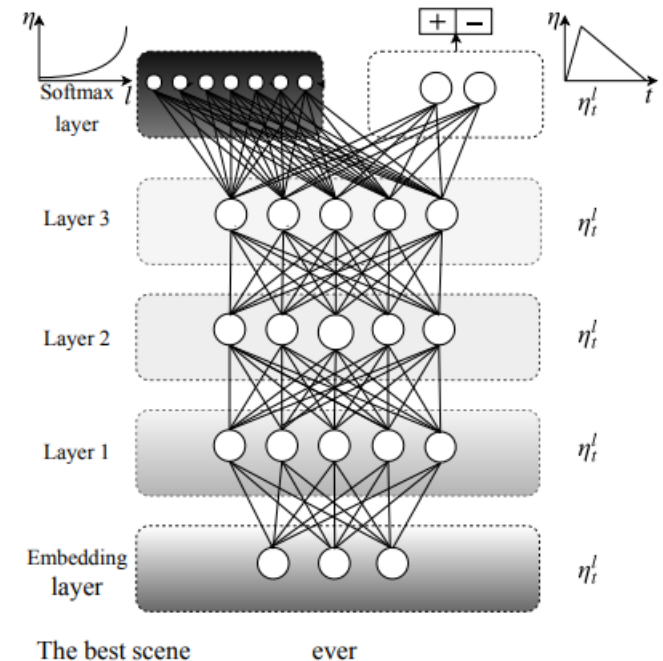
- ✓ Performance degrades as the data size shrinks.
- ✓ Some amount of repetition is ok.
- ✓ Use large pre-training datasets whenever possible.

## 5. Effect of Training Strategy (based on baseline)

# Fine-tuning Methods



## Adapter Tuning



## Gradual Unfreezing

More and more parameters  
are fine-tuned from top to  
down.

Parameter-efficient transfer learning for NLP. 2019 ICML.

Universal language model fine-tuning for text classification. 2018 ArXiv.

# Comparison of Different Fine-tuning Methods

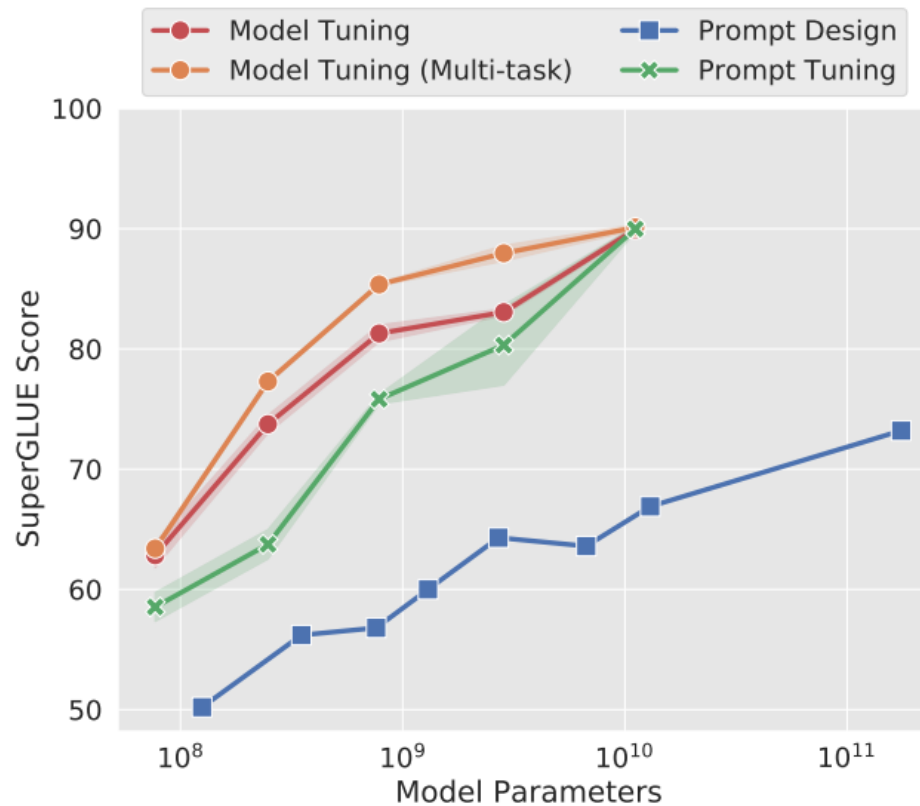
Fine-tuning method	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ All parameters	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
Adapter layers, $d = 32$	80.52	15.08	79.32	60.40	13.84	17.88	15.54
Adapter layers, $d = 128$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Adapter layers, $d = 512$	81.54	17.78	79.18	64.30	23.45	33.98	25.81
Adapter layers, $d = 2048$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Gradual unfreezing	82.50	18.95	79.17	<b>70.79</b>	26.71	39.02	26.93

Table 10: Comparison of different alternative fine-tuning methods that only update a subset of the model's parameters. For adapter layers,  $d$  refers to the inner dimensionality of the adapters.

*Note: we gradually unfreeze layers in the encoder and decoder in parallel, starting from the top in both cases. An additional layer are fine-tuned after every  $2^{18/12}$  steps.*

- ✓ Lower-resource tasks, such as SQuAD, work well with a small value of  $d$
- ✓ Higher-resource tasks require a large dimensionality.
- ✓ Gradual unfreezing degrades the performance across all tasks.

# Extension: The Power of Scale



The Power of Scale for Parameter-Efficient Prompt Tuning. EMNLP 2021

2022/7/25

Introduction to T5

30

# Multi-task Learning

- ✓ In the unified text-to-text framework, MTL corresponds to mixing data sets together.
- ✓ An extremely important factor in MTL is how much data from each task.
- ✓ How to set the proportion of data from each task?
  - data set sizes;
  - difficulty of learning the task;
  - task interference or negative transfer.

# Mixing Strategies for MTL

*A relaxed MTL: allowing to select a different checkpoint for each task.*

- Examples-proportional mixing:
  - Some tasks are so large, e.g. C4 and WMT English to French.
  - $r_m = \min(e_m, K) / \sum \min(e_n, K)$  , where K is the artificial data set size limit.
- Temperature-scaled mixing:
  - Scale the mixing rate with a temperature T, and renormalize it.
  - Set K to a large value  $2^{21}$ .
- Equal mixing:
  - Each example in a batch is randomly sampled from each task.
  - *A suboptimal strategy as model will quickly overfit on low-resource tasks and underfit on high-resource tasks.*



# Comparison of Different Mixing Strategies

Mixing strategy	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline (pre-train/fine-tune)	<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>
Equal	76.13	19.02	76.51	63.37	23.89	34.31	26.78
Examples-proportional, $K = 2^{16}$	80.45	19.04	77.25	69.95	24.35	34.99	27.10
Examples-proportional, $K = 2^{17}$	81.56	19.12	77.00	67.91	24.36	35.00	27.25
Examples-proportional, $K = 2^{18}$	81.67	19.07	78.17	67.94	24.57	35.19	27.39
→ Examples-proportional, $K = 2^{19}$	81.42	<b>19.24</b>	79.78	67.30	25.21	36.30	<b>27.76</b>
Examples-proportional, $K = 2^{20}$	80.80	<b>19.24</b>	<b>80.36</b>	67.38	25.66	36.93	<b>27.68</b>
Examples-proportional, $K = 2^{21}$	79.83	18.79	79.50	65.10	25.82	37.22	27.13
→ Temperature-scaled, $T = 2$	81.90	<b>19.28</b>	79.42	69.92	25.42	36.72	27.20
Temperature-scaled, $T = 4$	80.56	<b>19.22</b>	77.99	69.54	25.04	35.82	27.45
Temperature-scaled, $T = 8$	77.21	19.10	77.14	66.07	24.55	35.35	27.17

- ✓ MTL underperforms pretraining-then-finetuning on most tasks.
- ✓ The "equal mixing" strategy results in dramatically degraded performance.

*How to close the gap between MTL and pretraining-then-finetuning?*

# Combining MTL with Fine-tuning

*The model is pretrained on all tasks at once and then fine-tuned on the individual supervised tasks.*

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>

Table 12: Comparison of unsupervised pre-training, multi-task learning, and various forms of multi-task pre-training.

- ✓ Fine-tuning after multi-task pretraining leads to a comparable performance to baseline.
- ✓ “leave-one-out” multi-task training + fine-tuning is slight worse.
- ✓ Unsupervised pre-training is an important factor in most tasks.

## 6. Scaling

# How to Scale?

Recall the baseline:

- 12 blocks for Encoder, 12 blocks for Decoder;
- 3072 for FFN hidden-state dimension;
- 768 for hidden-state dimension;
- 12 heads
- 220M parameters



A large version (refer to BERT-Large):

- 4096 for FFN hidden-state dim.;
- 1024 for hidden-state dim.;
- 16 heads;

*You were just given  $4\times$  more compute. How should you use it?*

- Larger models: 16 and 32 blocks with a large version (roughly 2x and 4x the computational cost).
- More steps: such as 4x steps.
- Large batch size.
- Ensemble of 4 separately models (with averaging the decoder logits).

# Comparison of Different Scaling Strategy

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

Table 13: Comparison of different methods of scaling up our baseline model. All methods except ensembling fine-tuned models use 4× the computation as the baseline. “Size” refers to the number of parameters in the model and “training time” refers to the number of steps used for both pre-training and fine-tuning.

# 7. Putting It All Together

# Bag of Tricks

*Combining all insights and exploring the limits of transfer learning.*

- Objective: i.i.d. denoising objective --> span corruption (with mean length of 3 and corruption rate of 15%).
- Longer training:  $2^{40}$  tokens
  - 32 times as many as baseline;
  - 8, 2 and 1/2 times the size of BERT, XLNet and RoBERTa;
- Model sizes: Small, Base, Large, 3B(XL) and 11B(XXL).
  - Using small models can be helpful when computational resources are limited.
- Multi-task pre-training:
  - Step1: multi-task pretraining;
  - Step2: fine-tuning on the individual task after Step1;
- Decoding: greedy decoding --> beam search.

# Performance of T5 Variants

Model	GLUE Average	CoLA Matthew's	SST-2 Accuracy	MRPC F1	MRPC Accuracy	STS-B Pearson	STS-B Spearman
Previous best	89.4 <sup>a</sup>	69.2 <sup>b</sup>	97.1 <sup>a</sup>	<b>93.6<sup>b</sup></b>	<b>91.5<sup>b</sup></b>	92.7 <sup>b</sup>	92.3 <sup>b</sup>
T5-Small	77.4	41.0	91.8	89.7	86.6	85.6	85.0
T5-Base	82.7	51.1	95.2	90.7	87.5	89.4	88.6
T5-Large	86.4	61.2	96.3	92.4	89.9	89.9	89.2
T5-3B	88.5	67.1	97.4	92.5	90.0	90.6	89.8
T5-11B	<b>90.3</b>	<b>71.6</b>	<b>97.5</b>	92.8	90.4	<b>93.1</b>	<b>92.8</b>

Model	QQP F1	QQP Accuracy	MNLI-m Accuracy	MNLI-mm Accuracy	QNLI Accuracy	RTE Accuracy	WNLI Accuracy
Previous best	74.8 <sup>c</sup>	<b>90.7<sup>b</sup></b>	91.3 <sup>a</sup>	91.0 <sup>a</sup>	<b>99.2<sup>a</sup></b>	89.2 <sup>a</sup>	91.8 <sup>a</sup>
T5-Small	70.0	88.0	82.4	82.3	90.3	69.9	69.2
T5-Base	72.6	89.4	87.1	86.2	93.7	80.1	78.8
T5-Large	73.9	89.9	89.9	89.6	94.8	87.2	85.6
T5-3B	74.4	89.7	91.4	91.2	96.3	91.1	89.7
T5-11B	<b>75.1</b>	90.6	<b>92.2</b>	<b>91.9</b>	96.9	<b>92.8</b>	<b>94.5</b>

Model	SQuAD EM	SQuAD F1	SuperGLUE Average	BoolQ Accuracy	CB F1	CB Accuracy	COPA Accuracy
Previous best	90.1 <sup>a</sup>	95.5 <sup>a</sup>	84.6 <sup>d</sup>	87.1 <sup>d</sup>	90.5 <sup>d</sup>	95.2 <sup>d</sup>	90.6 <sup>d</sup>
T5-Small	79.10	87.24	63.3	76.4	56.9	81.6	46.0
T5-Base	85.44	92.08	76.2	81.4	86.2	94.0	71.2
T5-Large	86.66	93.79	82.3	85.4	91.6	94.8	83.4
T5-3B	88.53	94.95	86.4	89.9	90.3	94.4	92.0
T5-11B	<b>91.26</b>	<b>96.22</b>	<b>88.9</b>	<b>91.2</b>	<b>93.9</b>	<b>96.8</b>	<b>94.8</b>



# Performance of T5 Variants

Model	MultiRC F1a	MultiRC EM	ReCoRD F1	ReCoRD Accuracy	RTE Accuracy	WiC Accuracy	WSC Accuracy
Previous best	84.4 <sup>d</sup>	52.5 <sup>d</sup>	90.6 <sup>d</sup>	90.0 <sup>d</sup>	88.2 <sup>d</sup>	69.9 <sup>d</sup>	89.0 <sup>d</sup>
T5-Small	69.3	26.3	56.3	55.4	73.3	66.9	70.5
T5-Base	79.7	43.1	75.0	74.2	81.5	68.3	80.8
T5-Large	83.3	50.7	86.8	85.9	87.8	69.3	86.3
T5-3B	86.8	58.3	91.2	90.4	90.7	72.1	90.4
T5-11B	<b>88.1</b>	<b>63.3</b>	<b>94.1</b>	<b>93.4</b>	<b>92.5</b>	<b>76.9</b>	<b>93.8</b>

Model	WMT EnDe BLEU	WMT EnFr BLEU	WMT EnRo BLEU	CNN/DM ROUGE-1	CNN/DM ROUGE-2	CNN/DM ROUGE-L
Previous best	<b>33.8<sup>e</sup></b>	<b>43.8<sup>e</sup></b>	<b>38.5<sup>f</sup></b>	43.47 <sup>g</sup>	20.30 <sup>g</sup>	40.63 <sup>g</sup>
T5-Small	26.7	36.0	26.8	41.12	19.56	38.35
T5-Base	30.9	41.2	28.0	42.05	20.34	39.40
T5-Large	32.0	41.5	28.1	42.50	20.68	39.75
T5-3B	31.8	42.6	28.2	42.72	21.02	39.94
T5-11B	32.1	43.4	28.1	<b>43.52</b>	<b>21.55</b>	<b>40.69</b>

# Scale Is Not the Only Factor

Model	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Baseline-1T	84.80	19.62	83.01	73.90	27.46	40.30	28.34
T5-Base	<b>85.97</b>	<b>20.90</b>	<b>85.44</b>	<b>75.64</b>	<b>28.37</b>	<b>41.37</b>	<b>28.98</b>

Table 15: Performance comparison of T5-Base to our baseline experimental setup used in the rest of the paper. Results are reported on the validation set. “Baseline-1T” refers to the performance achieved by pre-training the baseline model on 1 trillion tokens (the same number used for the T5 model variants) instead of  $2^{35} \approx 34\text{B}$  tokens (as was used for the baseline).

- ✓ Baseline vs. Baseline-1T
- ✓ Baseline-1T vs. T5-Base

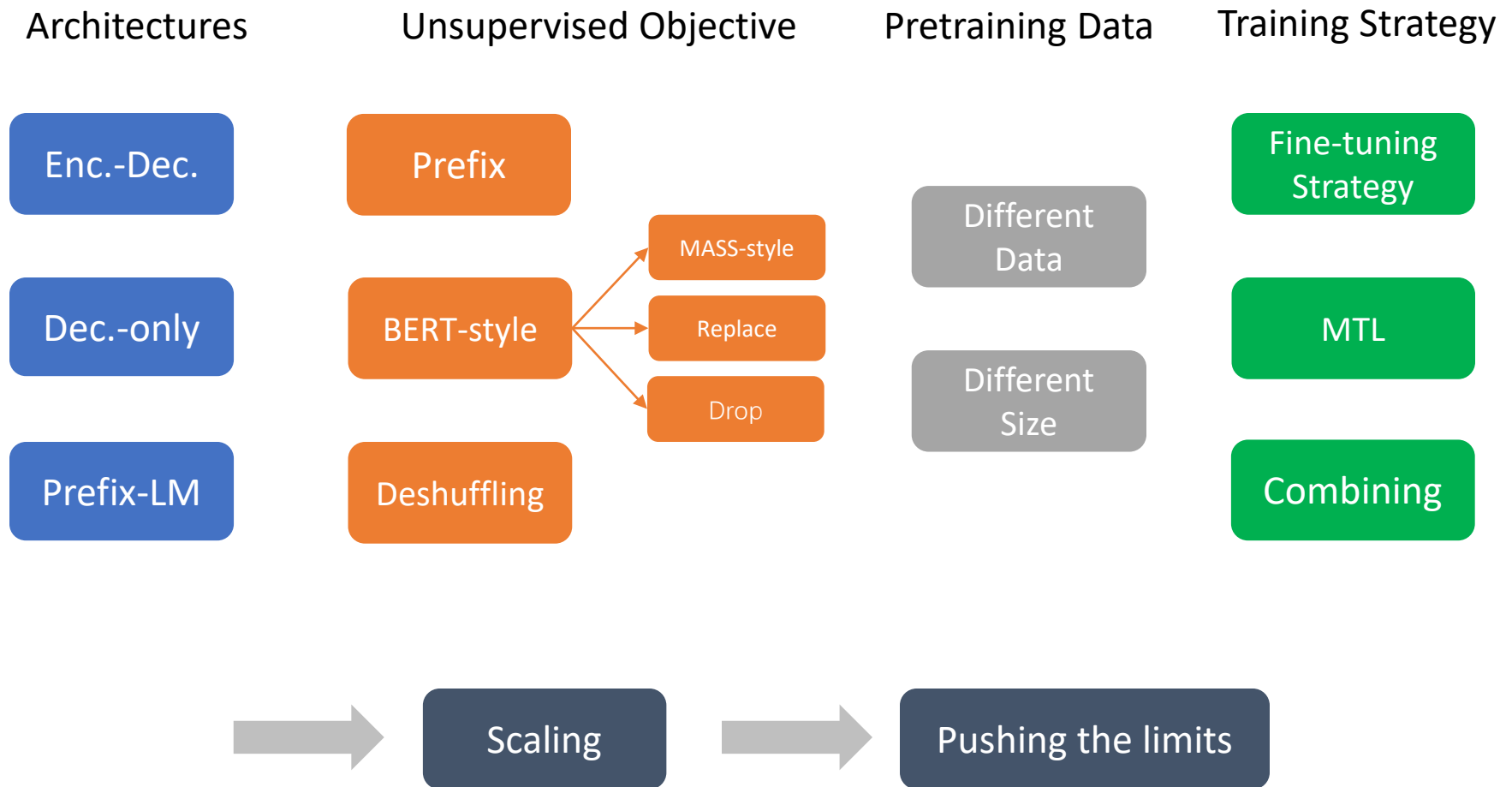
# All Results

Table	Experiment	GLUE													CNN/DM		SQuAD		Score Average	SuperGLUE				WMT										
		Score Average	CoLA MCC	SST-2 Acc	MRPC F1	MRPC Acc	STSBC PCC	STSBC SOC	QQP F1	QQP Acc	MNLI <sub>mm</sub> Acc	MNLI <sub>lm</sub> Acc	QNLI Acc	RTE Acc	R-1-F	R-2-F	R-L-F	EM F1		BoolQ Acc	CB F1	CB Acc	CoPA Acc	MultIRc F1	MultIRc EM	ReCoRD F1	ReCoRD EM	RTE Acc	WiC Acc	WSC Acc	EnDe BLEU	EnFr BLEU	EnRo BLEU	
1	★ Baseline average	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28	41.33	19.24	38.77	80.88	88.81	71.36	76.62	91.22	91.96	66.20	66.13	25.78	69.05	68.16	75.34	68.04	78.56	26.98	39.82	27.65
1	Baseline standard deviation	0.235	1.111	0.569	0.729	1.019	0.374	0.418	0.108	0.070	0.291	0.231	0.361	1.393	0.065	0.065	0.058	0.343	0.226	0.416	0.365	3.237	2.560	2.741	0.716	1.011	0.370	0.379	1.228	0.850	2.029	0.112	0.090	0.108
1	No pre-training	66.22	12.29	80.62	81.42	73.04	72.58	72.97	81.94	86.62	68.02	67.98	75.69	58.84	39.19	17.60	36.69	50.31	61.97	53.04	65.38	71.61	76.79	62.00	59.10	0.84	20.33	17.95	54.15	54.08	65.38	25.86	39.77	24.04
2	★ Ene/dec, denoising	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28	41.33	19.24	38.77	80.88	88.81	71.36	76.62	91.22	91.96	66.20	66.13	25.78	69.05	68.16	75.34	68.04	78.56	26.98	39.82	27.65
2	Ene/dec, shared, denoising	82.81	55.24	91.86	91.58	88.24	87.43	87.58	88.69	91.60	83.88	84.01	90.23	73.65	41.11	18.78	38.48	80.63	88.49	70.73	77.13	95.04	96.43	65.00	66.12	22.98	68.95	68.09	70.76	68.18	75.96	26.72	39.03	27.46
2	Ene/dec, 6 layers, denoising	80.88	46.26	92.09	91.51	87.99	87.01	86.76	87.93	90.87	82.20	82.41	88.83	71.48	40.83	18.97	38.31	77.59	86.07	68.42	73.79	91.70	92.86	67.00	61.02	19.62	61.26	60.33	72.20	65.99	75.00	26.38	38.40	26.95
2	Language model, denoising	74.70	24.50	90.80	86.08	78.92	85.22	85.42	85.40	88.99	76.72	77.05	86.02	64.62	39.49	17.93	36.91	61.14	71.37	55.02	65.47	90.08	71.43	58.00	43.03	2.94	53.35	52.31	53.07	58.62	63.46	25.09	35.28	25.86
2	Prefix LM, denoising	81.82	49.99	92.43	91.43	88.24	87.20	86.98	88.41	91.39	82.32	82.93	88.71	74.01	40.48	18.61	37.90	78.94	87.31	68.11	75.50	93.37	91.07	60.00	63.43	21.20	65.03	64.11	71.48	65.67	73.08	26.43	37.98	27.39
2	Ene/dec, LM	79.56	42.03	91.86	91.84	88.24	87.13	87.00	88.21	91.15	81.68	81.96	88.54	65.70	40.67	18.59	38.13	76.02	84.85	64.29	72.23	95.74	89.29	57.00	60.53	16.26	59.28	58.30	65.34	64.89	70.19	26.27	39.17	26.86
2	Ene/dec, shared, LM	79.60	44.83	92.09	90.20	85.78	86.03	85.87	87.77	91.02	81.74	82.29	89.16	65.34	40.16	18.13	37.59	76.35	84.86	63.50	70.49	91.41	87.50	55.00	60.21	16.89	57.83	56.73	63.54	63.48	70.19	26.62	39.17	27.05
2	Ene/dec, 6 layers, LM	78.67	38.72	91.40	90.40	86.52	86.82	86.49	87.87	91.03	80.99	80.92	88.05	65.70	40.29	18.26	37.70	75.32	84.06	64.06	71.38	85.25	89.29	60.00	57.56	16.79	55.22	54.30	66.79	63.95	71.15	26.13	38.42	26.89
2	Language modeling	73.78	28.53	89.79	85.23	78.68	84.22	84.00	84.88	88.70	74.94	73.78	84.84	58.84	38.97	17.54	36.37	53.81	64.55	56.51	64.22	59.92	71.43	64.00	53.04	1.05	46.81	45.78	58.84	56.74	69.23	25.23	34.31	25.38
2	Prefix LM, LM	79.68	41.26	92.09	90.11	86.27	86.82	86.32	88.35	91.35	81.71	82.02	89.04	68.59	39.66	17.84	37.13	76.87	85.39	64.86	71.47	93.37	91.07	57.00	58.67	16.89	59.25	58.16	64.26	66.30	71.15	26.28	37.51	26.76
4	Language modeling with prefix	80.69	44.22	93.00	91.68	88.48	87.20	87.18	88.39	91.41	82.66	83.09	89.29	68.95	40.71	18.94	38.15	77.99	86.43	65.27	73.55	83.95	87.50	55.00	59.65	18.89	61.78	60.76	68.59	65.67	73.08	26.86	39.73	27.49
4	BERT-style (Devlin et al., 2018)	82.96	52.49	92.55	92.79	89.95	87.68	87.66	88.47	91.44	83.60	84.05	90.33	75.45	41.27	19.17	38.72	80.65	88.24	69.85	76.48	94.37	94.64	61.00	63.29	25.08	66.76	65.85	72.20	69.12	75.00	26.78	39.82	27.41
4	DeShuffling	73.17	22.82	87.16	86.88	81.13	84.03	83.82	86.38	89.90	76.30	76.34	84.18	58.84	40.75	18.59	38.10	67.61	76.76	58.47	69.17	63.70	78.57	56.00	59.85	12.70	45.52	44.36	57.04	64.89	68.27	26.11	39.30	25.62
5	BERT-style (Devlin et al., 2018)	82.96	52.49	92.55	92.79	89.95	87.68	87.66	88.47	91.44	83.60	84.05	90.33	75.45	41.27	19.17	38.72	80.65	88.24	69.85	76.48	94.37	94.64	61.00	63.29	25.08	66.76	65.85	72.20	69.12	75.00	26.78	39.82	27.41
5	MASS-style (Song et al., 2019)	82.32	47.01	91.83	92.53	89.71	88.21	88.18	88.58	91.44	82.96	83.67	90.02	72.26	41.16	19.16	38.55	80.10	88.07	69.28	75.08	84.98	89.29	63.00	64.46	23.50	66.71	65.91	72.20	67.71	78.85	26.79	39.89	27.55
5	★ Replace corrupted spans	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28	41.33	19.24	38.77	80.88	88.81	71.36	76.62	91.22	91.96	66.20	66.13	25.78	69.05	68.16	75.34	68.04	78.56	26.98	39.82	27.65
5	Drop corrupted tokens	84.44	60.04	92.89	92.79	89.95	87.28	86.85	88.56	91.54	83.94	83.92	90.74	79.42	41.27	19.31	38.70	80.52	88.28	69.67	75.90	96.02	94.64	56.00	65.06	23.92	65.54	64.60	71.12	67.40	74.04	27.07	39.76	27.82
6	Corruption rate = 10%	82.82	52.71	92.89	91.55	88.24	88.02	87.94	88.67	91.40	83.59	84.51	90.33	75.45	41.05	19.00	38.75	80.38	88.36	69.55	74.98	92.37	92.86	62.00	66.04	24.66	67.93	67.07	70.76	67.24	75.96	26.87	39.28	27.44
6	★ Corruption rate = 15%	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28	41.33	19.24	38.77	80.88	88.81	71.36	76.62	91.22	91.96	66.20	66.13	25.78	69.05	68.16	75.34	68.04	78.56	26.98	39.82	27.65
6	Corruption rate = 25%	83.00	53.47	93.00	92.44	89.46	87.36	87.36	88.68	91.53	84.44	84.15	90.77	74.01	41.69	19.54	39.14	80.96	86.61	70.48	76.39	93.02	92.86	68.00	65.46	24.66	68.20	67.39	73.65	67.87	72.12	27.04	39.83	27.47
6	Corruption rate = 50%	81.27	46.26	91.63	91.11	87.99	87.87	87.64	88.70	91.57	83.64	84.10	90.24	70.76	41.51	19.32	38.89	79.80	87.76	70.33	75.02	93.05	92.86	68.00	62.97	24.13	64.94	64.13	72.20	68.50	77.88	27.01	39.90	27.42
7	★ Baseline (i.i.d.)	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28	41.33	19.24	38.77	80.88	88.81	71.36	76.62	91.22	91.96	66.20	66.13	25.78	69.05	68.16	75.34	68.04	78.56	26.98	39.82	27.65
7	Average span length = 2	83.54	53.82	92.20	93.05	90.44	87.85	87.71	88.42	91.40	84.28	84.46	90.88	77.62	41.23	19.39	38.69	82.09	89.69	72.20	77.06	94.37	91.07	70.00	66.28	26.13	71.34	70.61	75.34	68.34	78.85	26.76	39.99	27.63
7	Average span length = 3	83.49	53.90	92.43	92.25	89.46	87.49	87.53	88.72	91.51	84.85	84.84	90.99	77.26	41.50	19.62	39.44	81.84	86.69	72.53	76.95	94.37	94.64	70.00	67.64	28.75	70.84	69.90	74.73	67.71	77.88	26.86	39.65	27.62
7	Average span length = 5	83.40	52.12	93.12	92.63	89.71	88.70	88.47	88.84	91.64	84.32	84.29	90.79	76.90	41.39	19.24	38.82	82.05	89.79	72.23	77.06	83.06	89.29	69.00	68.16	30.12	71.38	70.53	75.51	69.91	79.81	26.88	39.40	27.53
7	Average span length = 10	82.85	50.11	92.09	91.95	88.97	88.45	88.22	88.86	91.63	84.34	84.28	91.07	76.17	41.38	19.33	38.80	81.84	89.39	70.44	76.4													

# Takeaways & Outlook

# Takeaways

## Text-to-Text Framework



# Outlook

- The inconvenience of large models:
  - Scaling up may continue to be a promising way to achieve better performance.
  - Authors advocate for research on methods that achieve stronger performance with cheaper models. (e.g. distillation and parameter sharing)
- More efficient knowledge extraction:
  - Obtain general-purpose “knowledge”;
  - Denoising objective in this work;
  - Any more efficient way?
- Formalizing the similarity between tasks:
  - Pretraining on unlabeled in-domain data can improve performance on downstream tasks.
  - Formulating the similarity between the pre-training and downstream tasks could help choose pre-training tasks and unlabeled data.
- Language-agnostic models:
  - English-only pre-training does not achieve SOTA results on translation tasks;
  - Vision: models can perform a given NLP task with good performance regardless of the text’s language.

Thanks