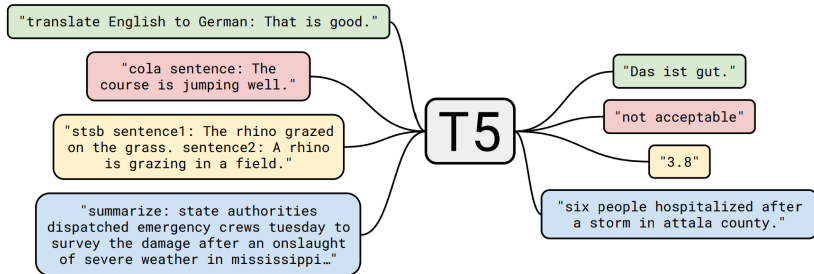# Using the Transformer

# T5 (Raffel et al., 2019)



**Learning goals**

- Understand the improvements over BERT
- Dynamic Masking

**T5: Text-to-Text Transfer Transformer:**

- A complete encoder-decoder Transformer architecture
- All tasks reformulated as text-to-text tasks
- From BERT-size up to 11 Billion parameters



Source: Raffel et al. (2019)

# THE COLOSSAL CLEAN CRAWLED CORPUS (C4)

- Effort to measure the effect of quality, characteristics & size of the pre-training resources
- Common Crawl as basis, careful cleaning and filtering for English language
- Orders of magnitude larger (750GB) compared to commonly used corpora

## Experiments (with respect to C4)

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

| 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 |

Source: Raffel et al. (2019)

# T5 - EXHAUSTIVE EXPERIMENTS

**Performed experiments with respect to ..**

- .. architecture, size & objective
- .. details of the Denoising objective
- .. fine-tuning methods & multi-taks learning strategies

**Conclusions**

- Encoder-decoder architecture works best in this "text-to-text" setting
- More data, larger models & ensembling all boost the performance
  - Larger models trained for fewer steps better than smaller models on more data
  - Ensembling: Using same base pre-trained models worse than complete separate model ensembles
- Different denoising objectives work similarly well
- Updating *all* model parameters during fine-tuning works best (but expensive)