# Using the Transformer

# Efficient Transformers



**Learning goals**

- Understand the efficiency problems and shortcomings of transformer-based models

- Learn about some strategies to alleviate them

# THE $\mathcal{O}(N^2)$ PROBLEM

**Quadratic time & memory complexity of Self-Attention**

- *Inductive bias of Transformer models:*
  Connect all tokens in a sequence to each other
- **Pro:** Can (theoretically) learn contexts of arbitrary length
- **Con:** Bad scalability limiting (feasible) context size

**Resulting Problems:**

- Several tasks require models to consume longer sequences
- *Efficiency:* Are there more efficient modifications which achieve similar or even better performance?

# EFFICIENT TRANSFORMERS

**Broad overview on so-called "X-formers"** `▶ Tay et al. (2020)`

- Efficient & fast Transformer-based models
  $\rightarrow$ Reduce complexity from $\mathcal{O}(n^2)$ to (up to) $\mathcal{O}(n)$
- Claim on-par (or even) superior performance
- Different techniques used:
    - Fixed/Factorized/Random Patterns
    - Learnable Patterns (extension of the above)
    - Low-Rank approximations or Kernels
    - Recurrence (see e.g. `▶ Transformer-XL (Dai et al., 2019)`)
    - Memory modules

**Side note:**

- Most Benchmark data sets not explicitly designed for evaluating long-range abilities of the models.
- Recently proposed: `▶ Longe Range Arena: A benchmark for efficient transformers` (Tay et al., 2020)

# INTRODUCING PATTERNS

**Reasoning:**

- Making every token attend to every other token might be unnecessary
- Introduce sparsity in the commonly dense attention matrix

**Example:**



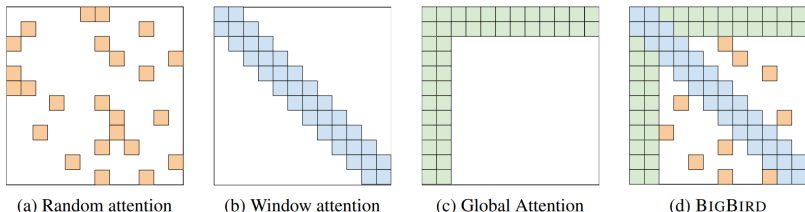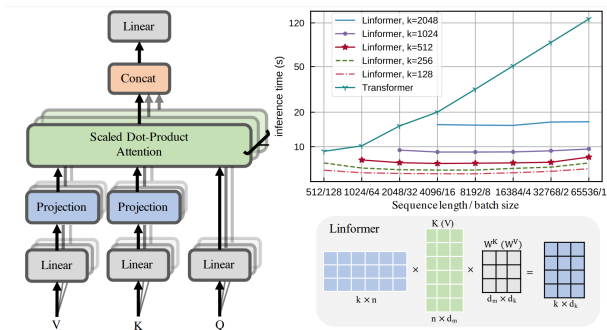(a) Random attention   (b) Window attention   (c) Global Attention   (d) BIGBIRD

Figure 1: Building blocks of the attention mechanism used in BIGBIRD. White color indicates absence of attention. (a) random attention with $r = 2$, (b) sliding window attention with $w = 3$ (c) global attention with $g = 2$. (d) the combined BIGBIRD model.

Source: Zaheer et al. (2020)

# LINEAR SELF-ATTENTION

**Reasoning** ▸ Wang et al. (2020)

- Most information in the Self-Attention mechanism can be recovered from the first few, largest singular values
- Introduce additional *k*-dimensional projection before self-attention



Source: Wang et al. (2020)

## DEBERTA

**Disentangled Attention** ▸ He et al. (2020)

- Each token represented by two vectors for content ($\mathbf{H}_i$) and relative position ($\mathbf{P}_{i|j}$)
- Calculation of the Attention Score:

$$A_{i,j} = \{\boldsymbol{H_i}, \boldsymbol{P_{i|j}}\} \times \{\boldsymbol{H_j}, \boldsymbol{P_{j|i}}\}^\mathsf{T}$$
$$= \boldsymbol{H_i}\boldsymbol{H_j}^\mathsf{T} + \boldsymbol{H_i}\boldsymbol{P_{j|i}}^\mathsf{T} + \boldsymbol{P_{i|j}}\boldsymbol{H_j}^\mathsf{T} + \boldsymbol{P_{i|j}}\boldsymbol{P_{j|i}}^\mathsf{T}$$

- with content-to-content, content-to-position, position-to-content and position-to-position attention

# DISENTANGLED ATTENTION

**Standard (Single-head) Self-Attention:**

$$Q = HW_q, K = HW_k, V = HW_v, A = \frac{QK^\mathsf{T}}{\sqrt{d}}$$

$$H_o = \text{softmax}(A)V$$

**Disentangled Attention*:**

$$Q_c = HW_{q,c}, K_c = HW_{k,c}, V_c = HW_{v,c}, Q_r = PW_{q,r}, K_r = PW_{k,r}$$

$$\tilde{A}_{i,j} = \underbrace{Q_i^c K_j^{c\mathsf{T}}}_{\text{(a) content-to-content}} + \underbrace{Q_i^c K_{\delta(i,j)}^{r}{}^\mathsf{T}}_{\text{(b) content-to-position}} + \underbrace{K_j^c Q_{\delta(j,i)}^{r}{}^\mathsf{T}}_{\text{(c) position-to-content}}$$

$$H_o = \text{softmax}\left(\frac{\tilde{A}}{\sqrt{3d}}\right)V_c$$

*Position-to-position part is removed since it, according to the authors, does not provide much additional information as *relative* position emebeddings are used.