

# Using the Transformer

## BigBird (Zaheer et al., 2020)



### Learning goals

- Understand subtleties of Self-Attention
- BigBird architecture using patterns

# INTRODUCING PATTERNS

## Reasoning:

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

## Example:

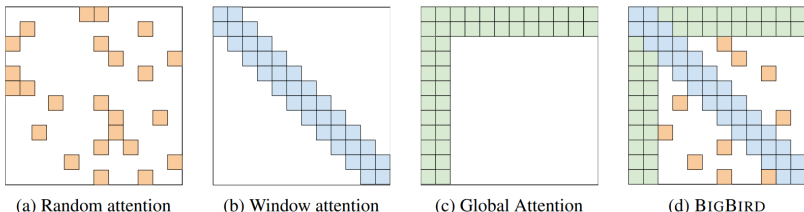


Figure 1: Building blocks of the attention mechanism used in BIGBIRD. White color indicates absence of attention. (a) random attention with  $r = 2$ , (b) sliding window attention with  $w = 3$  (c) global attention with  $g = 2$ . (d) the combined BIGBIRD model.

Source: Zaheer et al. (2020)