# Problem: Attention is blind to position

**Solution 1:** (Shaw et al., 2018)*              *simplified

- Encode **relative positions** via distance factors in the weights
- Weights:    $A = \text{Softmax}(X \cdot W^A \cdot X^\top + D)$      $D_{i,j} = d(i - j) \in R$

$$D \in R^{N \times N}$$

**Solution 2:** (Vaswani et al., 2017)

- Encode **absolute positions** via positional embeddings
- Input:     $X^P = X + E_P = x_1 + e_1, \dots, x_n + e_n$

- $E_P \in R^{N \times d}$ can be a decomposition of $D$      $E_P \cdot E_P^\top \approx D$