## Modify attention with block mask $M$

query $\boldsymbol{x}_i$, the entries of $M$ are set to minus $\infty$ all tokens outside $\mathcal{N}_i$ except the **[cls] token**, and **zero otherwise**

**cls** token



## block mask $M$



## Local attention: fast + inductive bias

$$\mathrm{Attention}(Q, K, V) = \mathrm{softmax}\left(\frac{QK^\top}{\sqrt{D}} + M\right)V$$