Integrating Exact Attention Locally

- Algorithm: Given Q, K and V, let Q'=arphi(Q) and K'=arphi(K) and compute $\mathsf{LT}(Q'\cdot (K')^\mathsf{T})\cdot V$

• Instead of using $\mathsf{LT}(Q'\cdot (K')^\mathsf{T})$ in the diagonal blocks, what if we use $LT((Q \cdot K^{\mathsf{T}})^4)$?

This is exactly what we are trying to approximate

The computation that is needed is of the same order

Improves the model quality as well!

Integrating Exact Attention Locally

- Algorithm: Given Q, K and V, let $Q' = \varphi(Q)$ and $K' = \varphi(K)$ and compute $\mathrm{LT}(Q' \cdot (K')^\mathsf{T}) \cdot V$
- Instead of using $\mathrm{LT}(Q'\cdot (K')^{\mathsf{T}})$ in the diagonal blocks, what if we use $\mathrm{LT}((Q\cdot K^{\mathsf{T}})^4)$?
 - This is exactly what we are trying to approximate
 - The computation that is needed is of the same order
 - Improves the model quality as well!

Model Perplexities

Perplexities on Wiki-40B

