

Scaling Laws

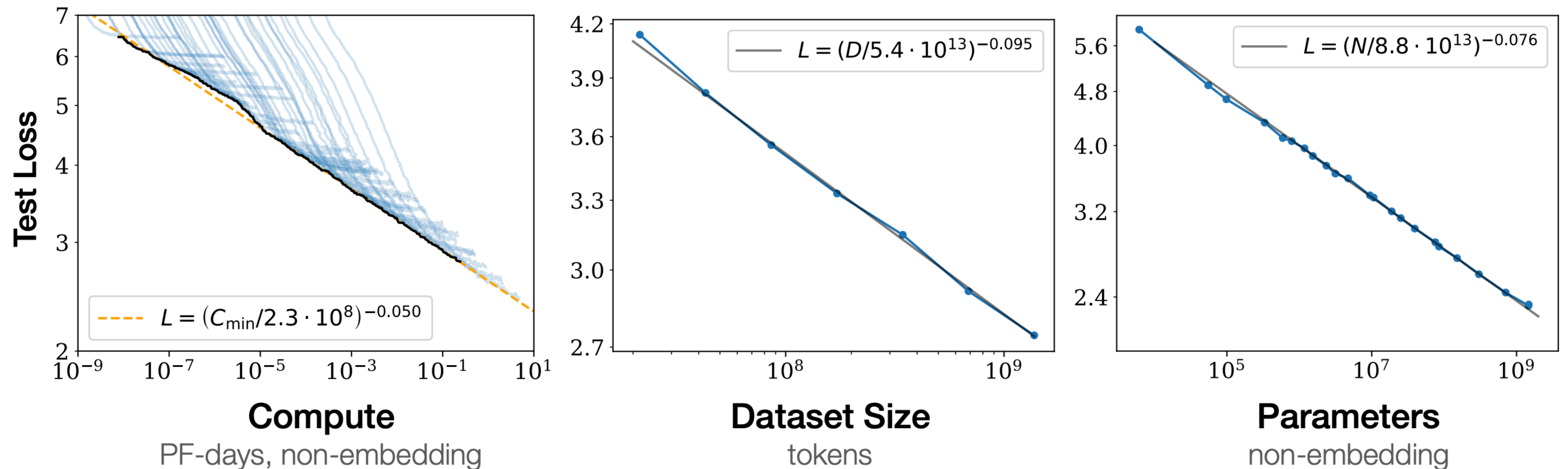


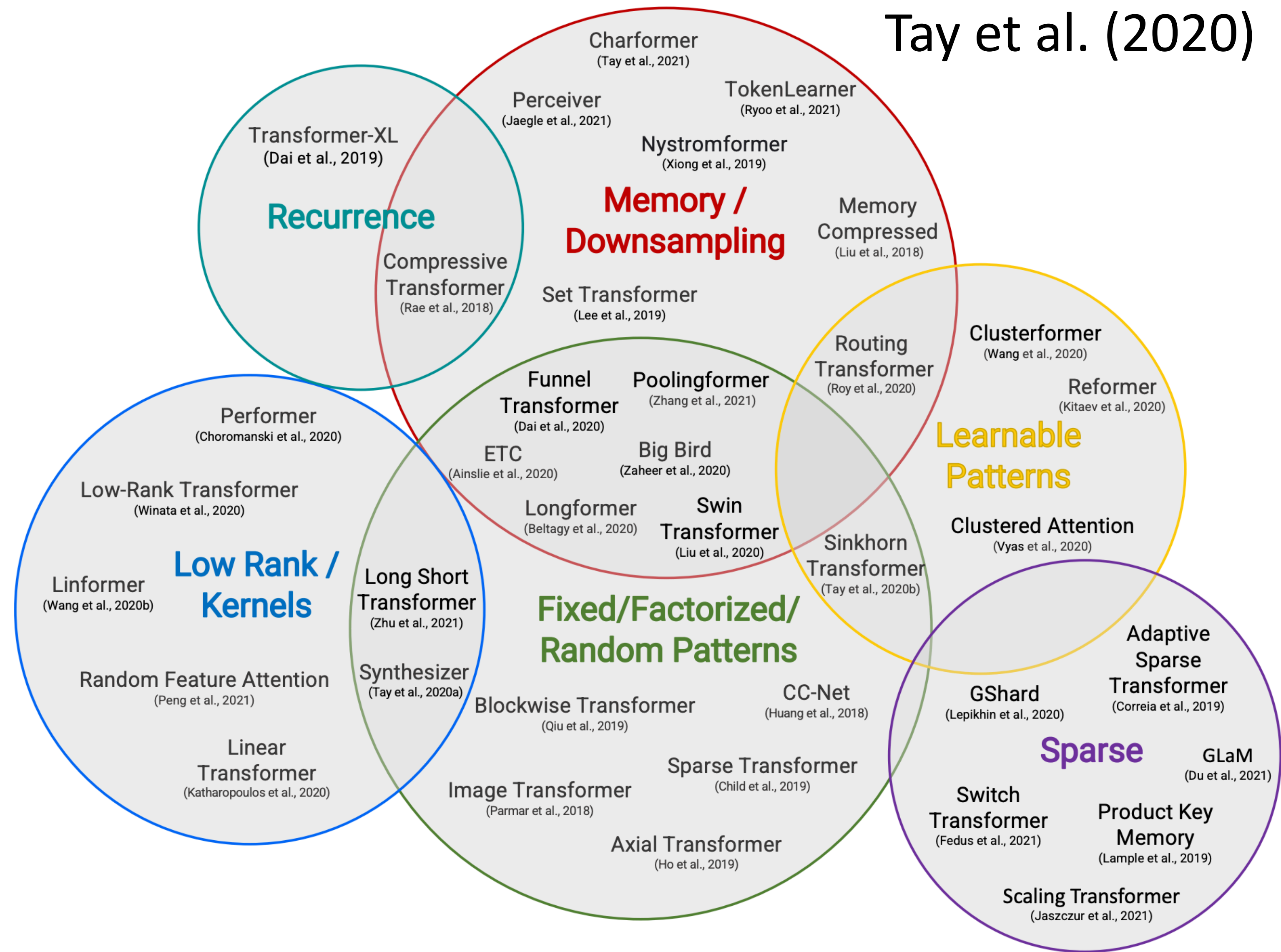
Figure 1 Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute² used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

- Transformers scale really well!

Kaplan et al. (2020)

Transformer Runtime

- ▶ Even though most parameters and FLOPs are in feedforward layers, Transformers are still limited by quadratic complexity of self-attention
- ▶ Many ways proposed to handle this



Performers

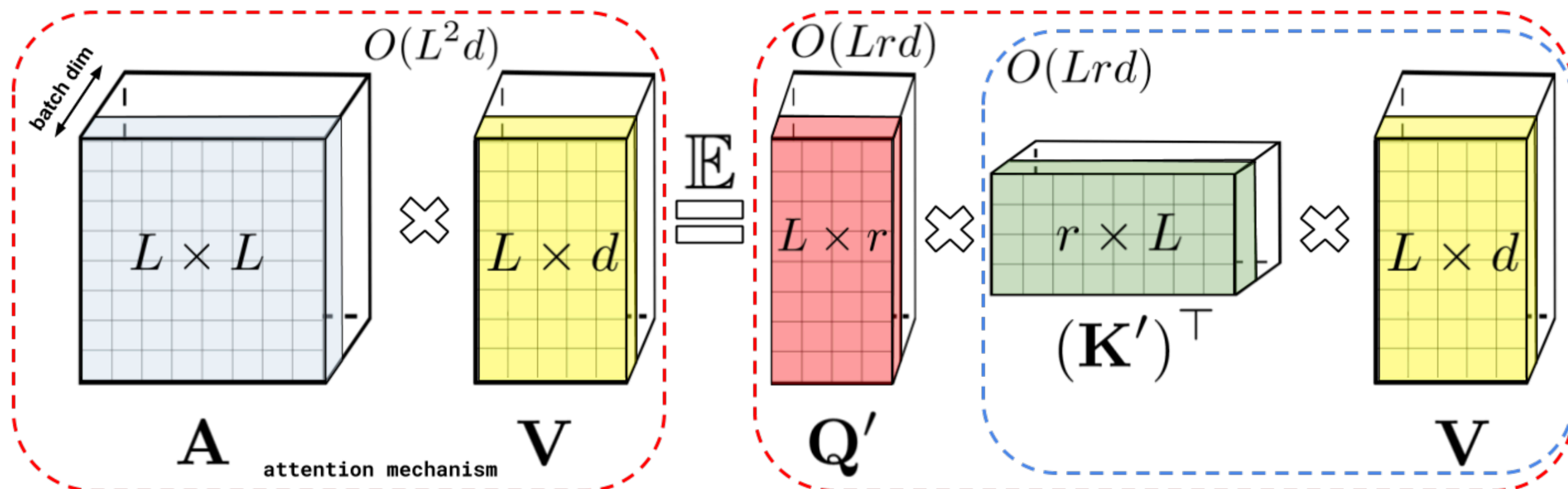
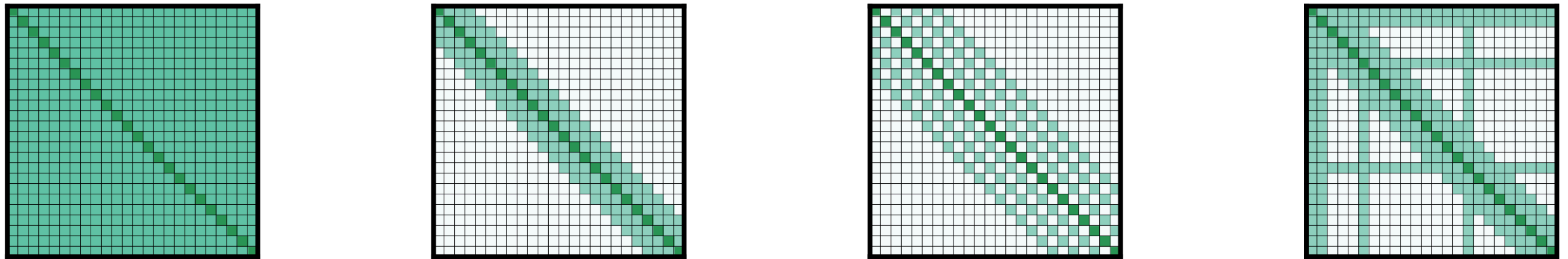


Figure 1: Approximation of the regular attention mechanism $\mathbf{A}\mathbf{V}$ (before \mathbf{D}^{-1} -renormalization) via (random) feature maps. Dashed-blocks indicate order of computation with corresponding time complexities attached.

- No more len^2 term, but we are fundamentally approximating the self-attention mechanism (cannot form \mathbf{A} and take the softmax)

Choromanski et al. (2020)

Longformer



(a) Full n^2 attention

(b) Sliding window attention

(c) Dilated sliding window

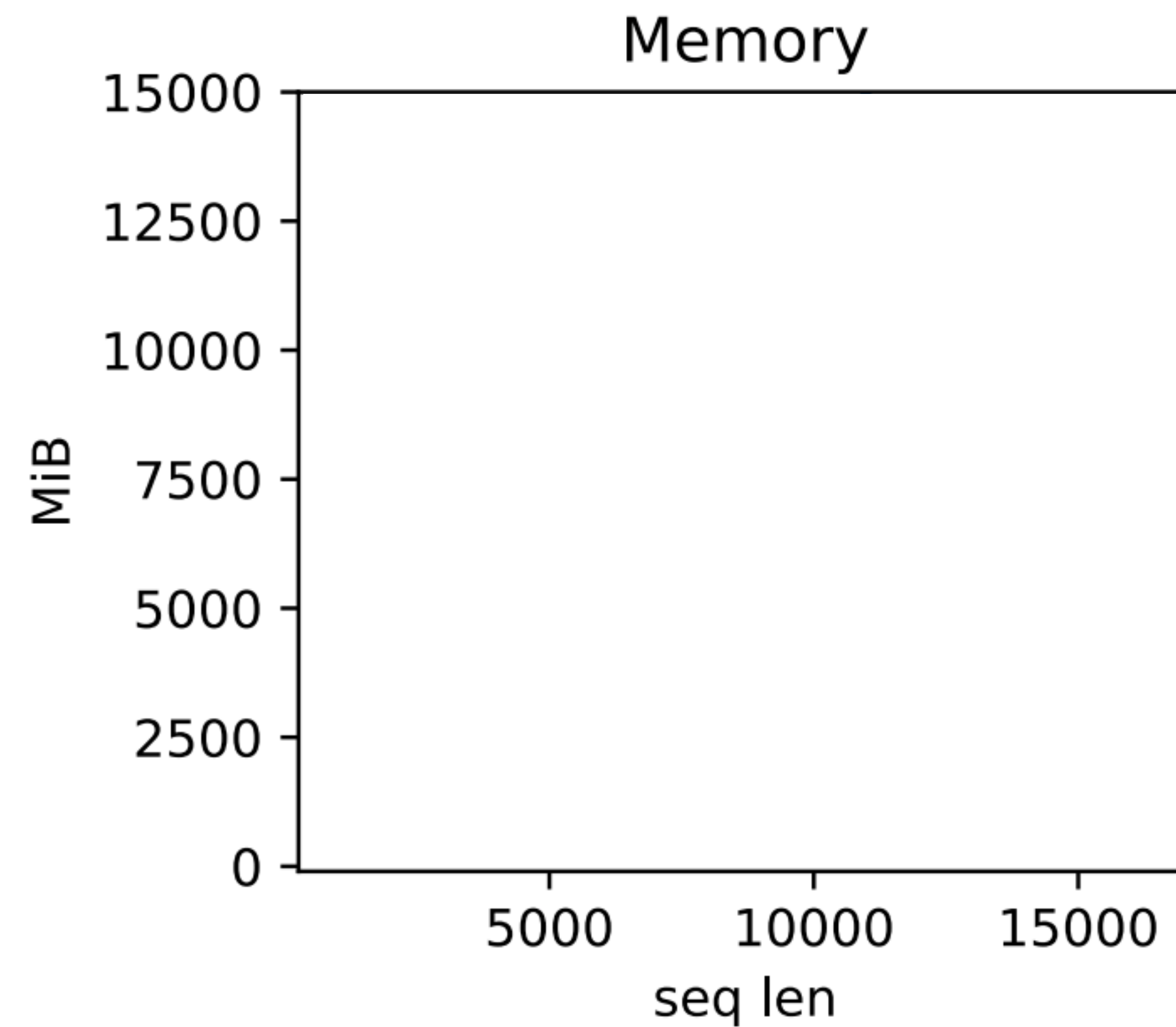
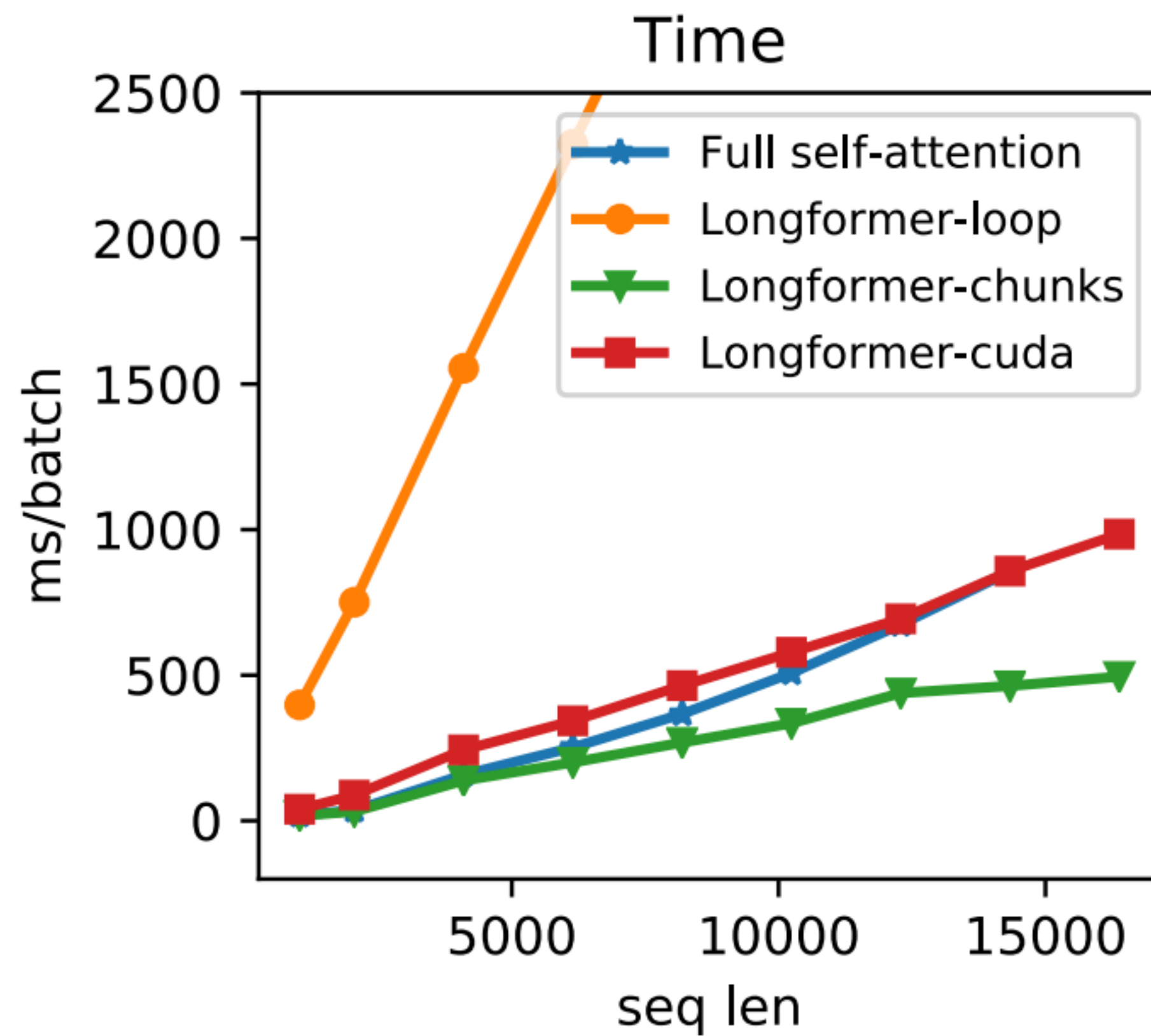
(d) Global+sliding window

Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.

- ▶ Use several pre-specified self-attention patterns that limit the number of operations while still allowing for attention over a reasonable set of things
- ▶ Scales to 4096-length sequences

Beltagy et al. (2021)

Attention Maps



- ▶ Loop = non-vectorized version
- ▶ What will the memory profile look like?

Beltagy et al. (2021)