

트랜스 포머의 Attn은 $O(n^2)$ 의 복잡도를 갖음

⇒ 계산량, 시간, 메모리 소요↑

* 긴문장이 주어졌을때 Bert와 같은경우, Truncated_Seq를 사용하게 되는데 일반 텍스트 하의 개개의 문 라인이 와.

⇒ 하의 개개의 제한과, 문맥적 때문에 문맥 파악이 어려울 수 있다.

Model	attention matrix	char-LM	other tasks	pretrain
Transformer-XL (2019)	ltr	yes	no	no
Adaptive Span (2019)	ltr	yes	no	no
Compressive (2020)	ltr	yes	no	no
Reformer (2020)	sparse	yes	no	no
Sparse (2019)	sparse	yes	no	no
Routing (2020)	sparse	yes	no	no
BP-Transformer (2019)	sparse	yes	MT	no
Blockwise (2019)	sparse	no	QA	yes
Our Longformer	sparse	yes	multiple	yes

Table 1: Summary of prior work on adapting Transformers for long documents. ltr: left-to-right.

< 기존 연구 >

long Document는 : ltr : 긴 텍스트를 left → right 방향으로 임베딩. / Bidirectional 방법을 사용하는 것은 약함.
 Sparse Attn 방식이 있다 - General한 방식 / Attn 복잡도 낮춤.
 ⇒ Longformer.

복잡도를 낮추기 위한 방법들.

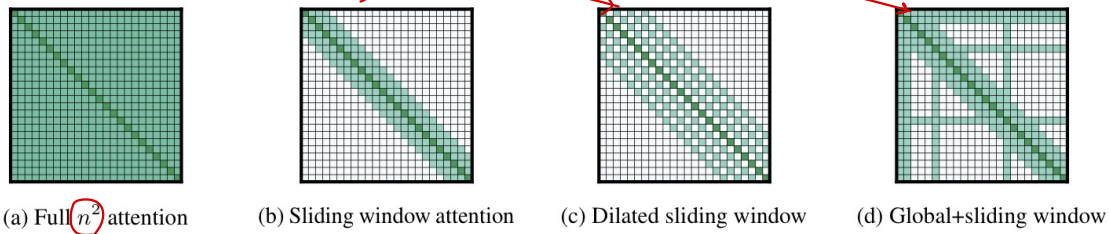
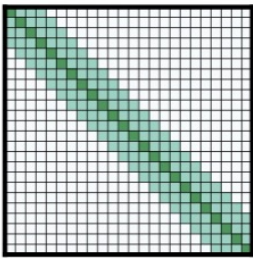


Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.

LongFormer는 문장이 길더라도 레코드를 고려한 whole Contextual representation을 활용하면서 모델 아키텍처에 의해서 양의 성능을 낼 수 있다.

1) Sliding window.



(b) Sliding window attention

: 둘은 근처에 있는 토큰들끼리만 Attn을 계산. 멀리 떨어진 토큰끼리는 계산하지 x.

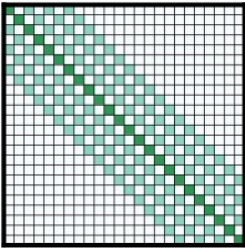
W : fixed window size.

⇒ Each Token attends $\frac{1}{2}W$ tokens on each side.

Computational Cost: $O(n \times w)$.

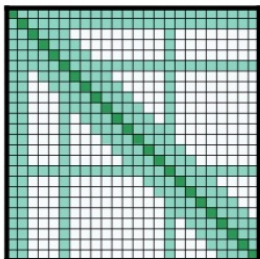
Receptive field: $L \times W$

2).



(c) Dilated sliding window

- 위의 슬라이딩 방법보다 좀더 효율적이고, 초기 컨텍스트를 좀더 확장하며 만든 방식 (더 넓은 범위)
- Receptive field: $L \times D \times W$.
- 특히 Multi-head Attn에서 서로 다른 dilated를 사용해서 Attn Score를 계산
 - Dilated가 낮은 애들은 좀더 지향적인 애들을
 - Dilated가 높은 애들은 긴 문맥 정보를 처리함



(d) Global+sliding window

- * 1). 이는 local Context 만 바라보기 때문에 컨텍스트에는 부족한 부분이 생김.
- * 하지만의 모델은, 문제들을 풀며, [sep] 같이 special 토큰을 생성하게 하거나 [cls] 토큰을 전체적인 문맥에 의해서 에니션하듯 하도록 유도하게 함. (?)
따라서 글로벌 에니션 방식에서는 미리 선택해준 special 토큰에 대해서 전체 토큰에 대해서 에니션 할 수 있도록 하였다.

special 토큰의 개수가 문맥길이에 비해 상당히 작기 때문에 몇몇 토큰에 대해 글로벌 에니션을 적용하였을 개수에 있어서 무언이 안된다.

+ 복잡도 $O(n)$.

* 석연 스텔 도리에 대해서만 잘 이해를 한다.

[CLS].[]

↓
스텔 도리에 대한 명확한 명칭을 알려주지 않는다.

* Dilated \Rightarrow 레이어를 층층이 쌓아서, 같은 층의 정보 위로 각층의 쿼리 수를 크게..

CNN 같아?

\Rightarrow 인. 레이어에서는, Window 내에서 fully attention을 하되.

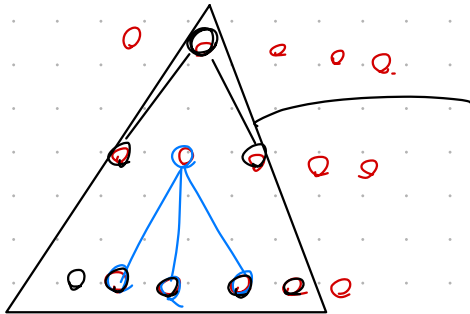
위로 각층 dilation을 준다.

영어 Ver.

* Transformer는 매우 harsh한 제약조건이 있어서 긴 문장은 문장 단위로 끊어서 봐야 함.

⇒ 각 문장은 독립적이 되어 버린다.

*. 여를 살펴봐 (dilated).



let's say kernel size is 3.

원래 차를 더 넓게 정보를 획득 함.

⇒. 낮은 차에서는 local 한 정보를 얻는다. 반면 차를 갈수록 좀더 넓은 정보를 얻는다.

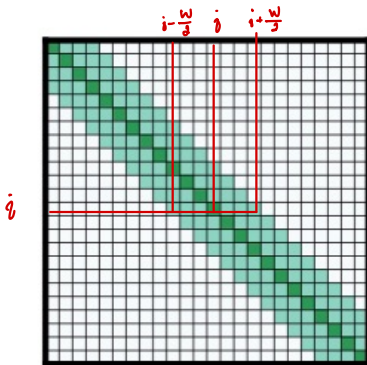
* CNN과 비슷한 느낌인데

생각해보면 make sense 함.

왜냐하면 차의 값이 주변에 있는 애들이 서로에게 중요한 역할을 하겠지.

* CNN에서와 마찬가지로 transformer에서도 함.

* $O(m \cdot w)$
↑
Window size.



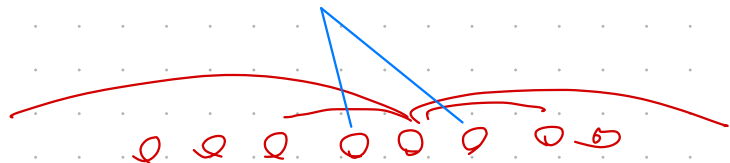
(b) Sliding window attention

위에서 봤던 CNN 같은것은. 1개의 layer를 올라가면, 2w 만큼 Context를 얻는다.

NLP는 주변 정보가 중요하잖아...

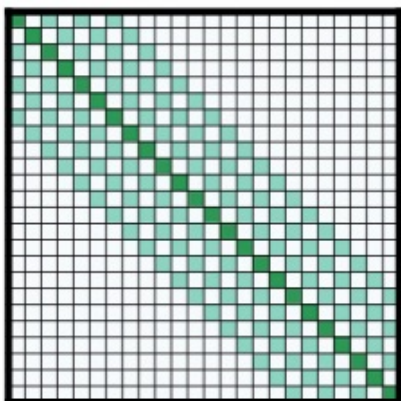


Oh well! 주변 정보 안중요

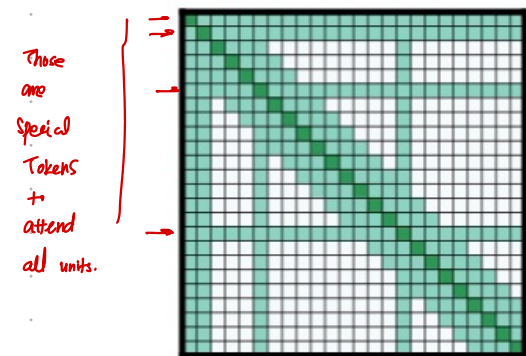


그래서 낮은 layer에서 (b) 쓰고.

높은 차에서 (c) 쓰고.



(c) Dilated sliding window



(d) Global+sliding window

→ 이는 Special Token이 왔다고 함.