seq2seq pretrain    $w_1, w_2 w_3, w_4$    reconstruct the input (self-supervised)

[model] → [Model]
↑ ↑ ↑ ↑
$w_1, w_2, w_3, w_4$ (corrupted)

怎么样 corrupt input? MASS (Masked seq to seq pre-training)
          or BART (Bidirectional and autoregressive transformers)

MASS ─ 随机的 mask token    还原的 token
       or delete ─         还原 token
    permutation token      还原的 token
BART < rotation token      还原 token
                                         还原 token e.g: A □ B [SEP] □ E
BART 一般是比较好的 Text infilling } (to token or 多个 token)           ↓
                                                                    A B [SEP] C P E (true)

Unified ← encoder     BERT (bidirectional LM)
         decoder      GPT (left-n-right LM).  需要好好设计 attention 每个token
         seq2seq      BART/MASS (seq2seq LM)  以防看到后面的token
                                              Segment 1    segment 2
                                                  ↓            ↓
                                              encoder      decoder
                                              (可相互看)    (只能往前看)

ELECTRA (efficiently learning an encoder that classifiers token replacements accurately)
NO NO YES NO NO
□ □ □ □ □           predicting yes/no is easier than reconstruction.
↑ ↑ ↑ ↑ ↑
[  Model  ]         every output position is used. (BERT 只 predict masked)
↑ ↑ ↑ ↑ ↑
the chef ate the meal
       因字被replaced

怎么样置换，让语义上假看起来是对的呢。─ 用 smaller BERT 希望能骗到就好.

$\frac{1}{4}$ size of XLNet 就能达到 similar GLUE score

Predict next token

LSTM → ELMO (双向：但是它向、逆向=分前没有分式，已向时没有看到右向的 token)

self-attention → GPT, Megatron, Turing NLG

masking input (whole word { word
                        $w_2$
                                    pharase-level / entity level — ERNIE

BERT                        span (prob)                    __SpanBERT__
                                    span length                    1

                random mask

(Bow.            $\frac{1}{2}W(t-2)$                sum $w(t)$
(using content    $w(t-1)$            →    □ → W
to predict the    $=(W(t+1)$
missing token)    $w(t+1)$

                window 20左右

SpanBERT — Span Boundary Objective
                    $w_6$                            (SBO)
                    SBO
                    2

                $w_1$ $w_2$ $w_3$ □ □ □ □ $w_8$

                    同于 coreference

XLNet

Transformer-XL

            渐  $\stackrel{?}{2}$  $\stackrel{?}{3}$  可4

language model 的话

        可4  $\stackrel{?}{3}$  $\stackrel{?}{2}$  渐1

BERT 的话
                                    random
                                    or
        渐1  $\stackrel{?}{2}$  MASK 可4            渐1  $\stackrel{?}{2}$  MASK  可4

        XLNet - 每次 positional word 不给 content 信息，所以预测可以进行预测。 看 paper

BERT cannot talk?  Word 是 sequence 生成之是左右而右生成的。BERT 也可以 talk)
    Given partial sequence, predict the next token (但 BERT training之是通过把 两边预测
    ∴BERT 不善言语 ...  不适用于做 seq2seq pre-trained model (只能是做 encoder / decoder 部分
    但 MASS/BART 是 pre-trained seq2seq model 用 by self-supervised learning        pretrain 的)

BERT and its family

- Pretrain model
  word embedding. ⎰ word2vec
  ⎱ glove
  ⎱ fasttext · character-level. 可 handle 新词
  每个 "狗" 词义不一样
  contextualized word embedding. ⎰ ELMO
  看过句子 才用给 embedding. ⎱ BERT

Smaller BERT                         Network Compression (see refer for reference)
  distill BERT                         net                    video
  Tiny BERT
  Mobile BERT
  Q8BERT
  ALBERT (最知名. 小且结果 improve)
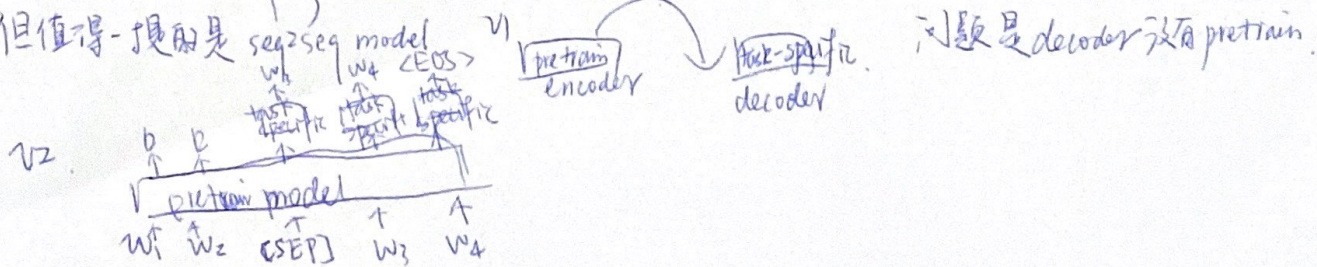
Network Architecture (希望读非常大的 sequence). BERT - 512 token
  Transformer-XL : segment-level recurrence with state reuse

  reformer
  longformer  → 当 sequence 很长 self-attention takes $O(N^2)$ 这两者自来减少 complexity of self-attention

- How to fine-tune ⎰ ① fine tun task-specific
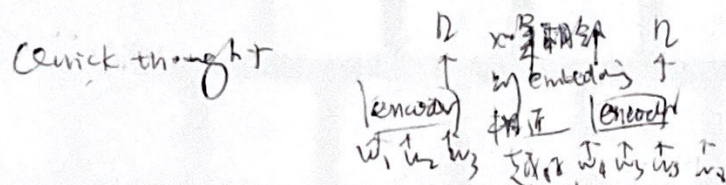  ② fine-tune 入都好 pre-trained & task-specific
  BERT 里 ML 设计了 ③ adaptor: 大多同 pretrain 小心 参数 adaptor & task-specific
                          attention
  但值得一提的是 seq2seq model   v1  [pretrain encoder] → [task-specific decoder]  问题是 decoder 没有 pretrain
                    w3  w4 <EOS>
                   ↑ task ↑ task ↑ task
                  specific specific specific
  v2.   b  b        [pretrain model]
        ↑  ↑  ↑  ↑
        w1 w2 [SEP] w3 w4

Sentence Level    pretrain

(representation for sentence   next sentence

skip thought    [encoder] → ⬚ → [decoder]    表示如果下一个 sentence 是相关的话.
                ↑ ↑ ↑              ↑ ↑ ↑ ↑       把 sentence embedding 找出以入
                $w_1$ $w_2$ $w_3$   $w_4$ $w_5$ $w_6$ $w_7$   (但生成的难度应该会比较大)

(Quick thought    ⬚    x是相邻    ⬚
                  ↑    的 embedding  ↑
            [encoder] 拉近 [encoder]
            $w_1$ $w_2$ $w_3$   就像 $w_1$ $w_2$ $w_3$

BERT train [CLS] embedding 用一个 global 的网络: next sentence prediction

Robustly optimized BERT approach. (ROBERTa)        但结果不如 US
                                    ↙              可能因为很闹事
SOP: sentence order prediction.    (两个是相连 则说 yes.
    ↓ 更用在往的 task.: 两个句子很相近 如果 调顺序 则说 No.)
  ↳ ALBERT.

    StructBERT (Alibaba): 结合了 next sentence prediction 和 SOP.

T5 - ~~comparison~~ 是一个 comparison

another    BERT + knowledge → ERNIE
story ↙                    (不是百度的 ERNIE)

    Audio BERT