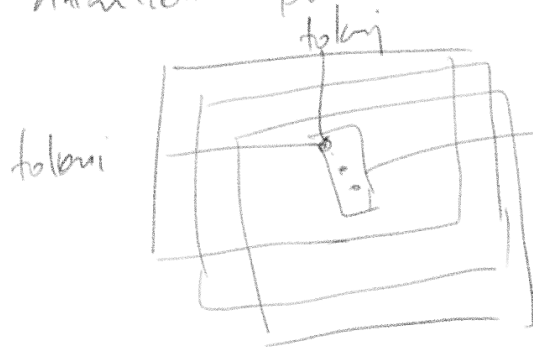# Exploring Geometry of BERT

Can understanding the geometry of internal representations tells us more about BERT, or help us improve BERT (optimize something that helps reduce the computation - maybe

- Does attention matrices encode syntactic features

- Dependency grammar relations

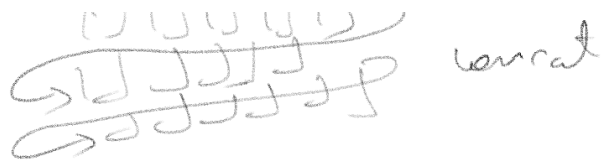- Manning analyzed context encodings They are using attention matrices.

- Attention probe -
  token



token

$\rightarrow$ Concat these to obtain attn vector

$\cancel{e_{ij}}$ $a_{ij}$

- every layer, every head

$a_{ij} - \cancel{12 \times 12}$

overlay 

12

concat

single $a_{ij}$ = [ ... boxes ... ]

$$12 \times 12 = 144$$

Goal → classify a given relation b/w 2 tokens

    ↳ If good ⟹ modelwise attn vector ~~achieves~~ ~~the~~ encodes that relation.

— 2 L2 linear models

① ↳ [ 1 1 1 4 1 ] → dependency b/w i &j
    attn vector     or not [binary classifier]

       ↘ [86.8%]

② multiclass classifier

$a_{ij}$ [ 1 1 1 1 1 1 ] → which ~~relat~~ dependency
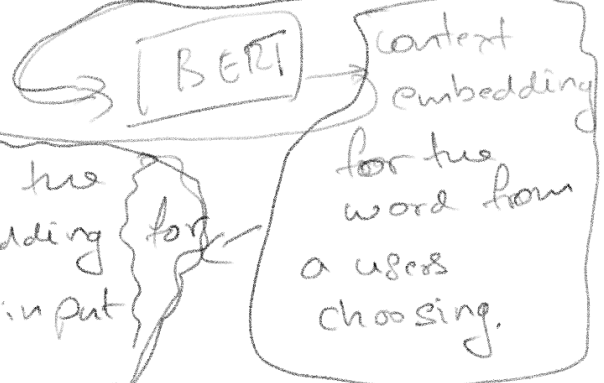         exists b/w tokens i &j.

        ↘ [71.9%]

3.2] Geometry of parse tree embeddings

parse tree distance = (euclidean distance)$^2$

4.1

(W)

[WD sent] { → W'
          W,
          W'

[BERT] → Context embedding for the word from a user's choosing.
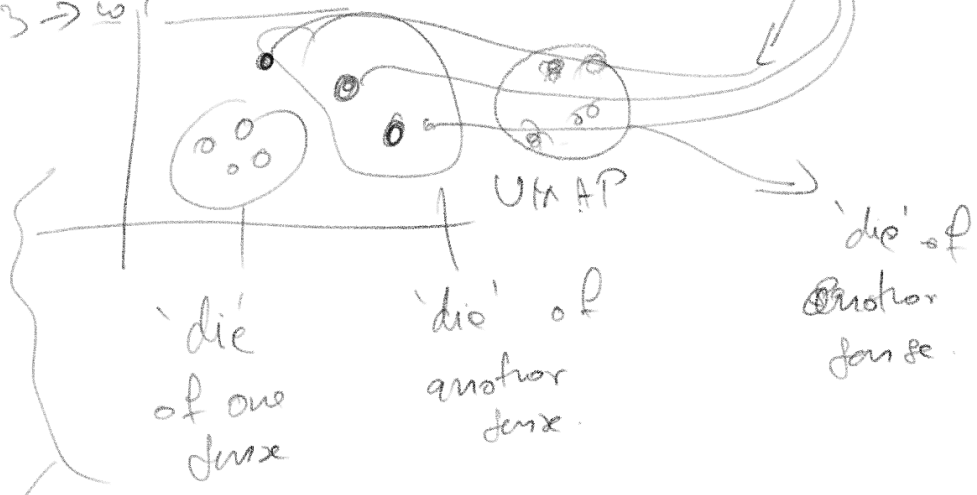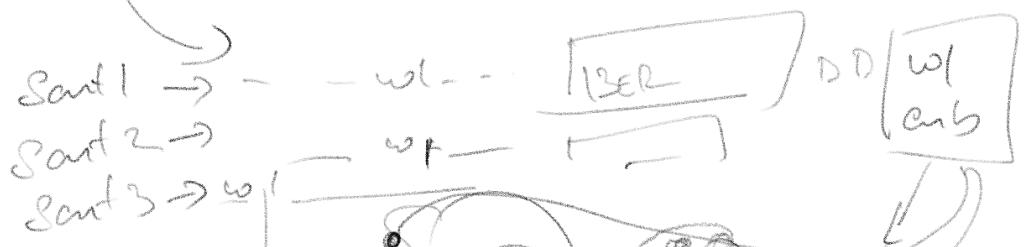
Q How to know the context embedding for a particular input word

context embedding (one for each of i/p tokens)

Sent 1 →  --- W' ---  [BER]  DD [W/ emb]
Sent 2 →  --- W' ---  [        ]
Sent 3 → W'

UMAP

`die` of one sense

`die` of another sense

`die` of another sense

→ This technique even does WSD with
   F1 = 71.1% (Nearest Neigh)
   "... ... CAPTURING WORD SENSE" - note

(ONTEXT EMBEDDINGS ARE [illegible]
Sense due to "attention" - captures the content botter-

○

(ONTEXT EMBEDDINGS ARE [illegible]
Sense due to "attention" - captures the content botter-