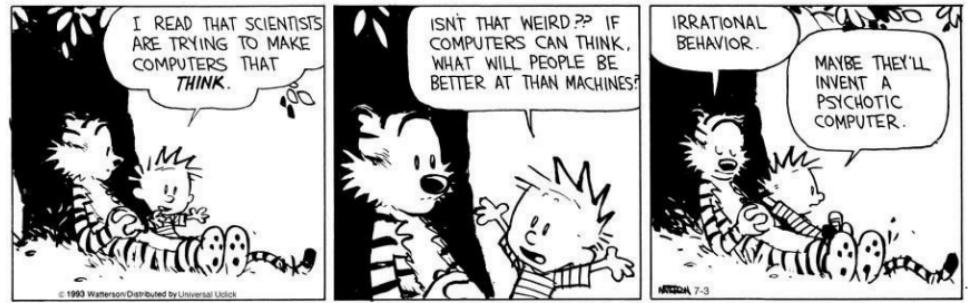


# BERTology review: What Does BERT Look At? An Analysis of BERTs attention

Kevin Clark, Urvashi Khandelwal, Omer Levy, Christopher Manning (2019)

Sualeh Asif

MIT Computer Science Artificial Intelligence Laboratory  
Massachusetts Institute of Technology



## ABOUT ME:

SUALEH ASIF

XG

MIT COMPUTER SCIENCE  
MATH  
THEATRE

### INTERESTS :

MACHINE LEARNING  
REPRESENTATIONS  
EMERGENCE  
VISUAL  
NLP

PERFORMANCE ENGINEERING  
CPUs + GPUs

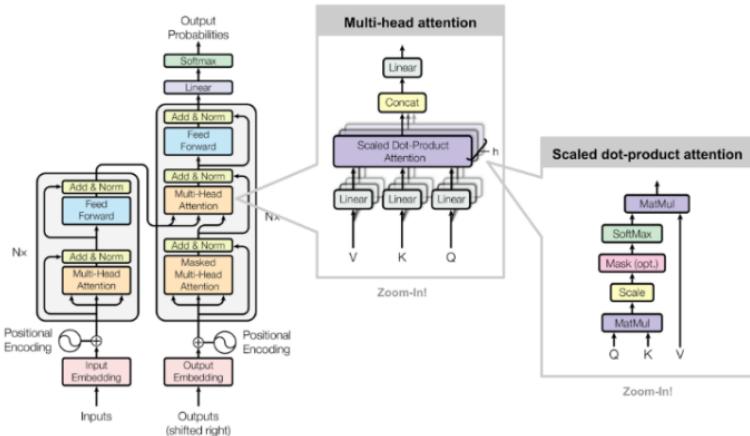
COMPUTATIONAL NUMBER THEORY



[Sualeh@mit.edu](mailto:Sualeh@mit.edu)  
617) 599 5080

# BERTology ....

X

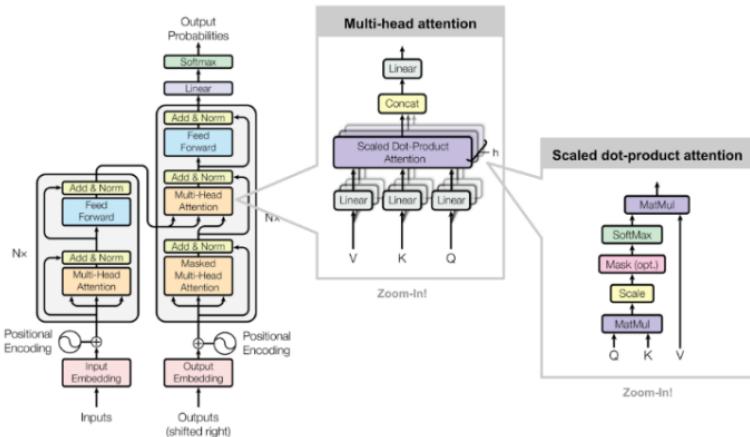


# BERTology

....

But ... why?

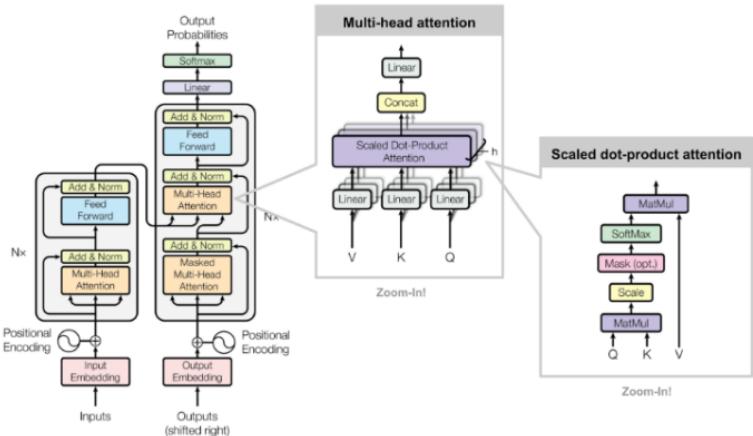
X



# BERTology ....

But ... why?

X

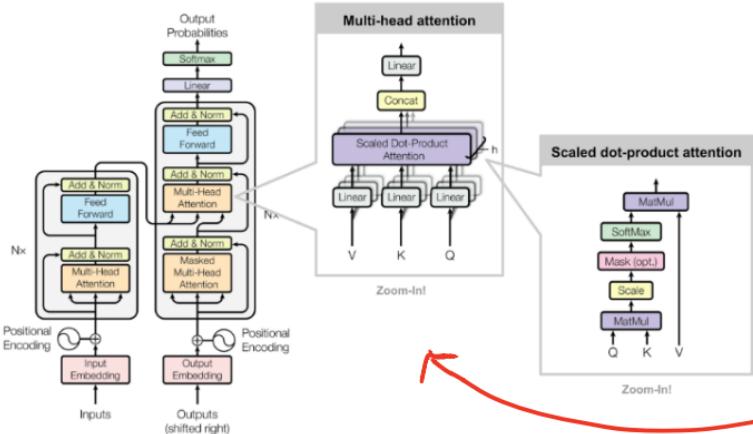


\* what is BERT  
and how did  
we get there?  
*(short history)*

# BERTology ....

But ... why?

X



\* What is BERT and how did we get there?

\* What does BERT teach us?

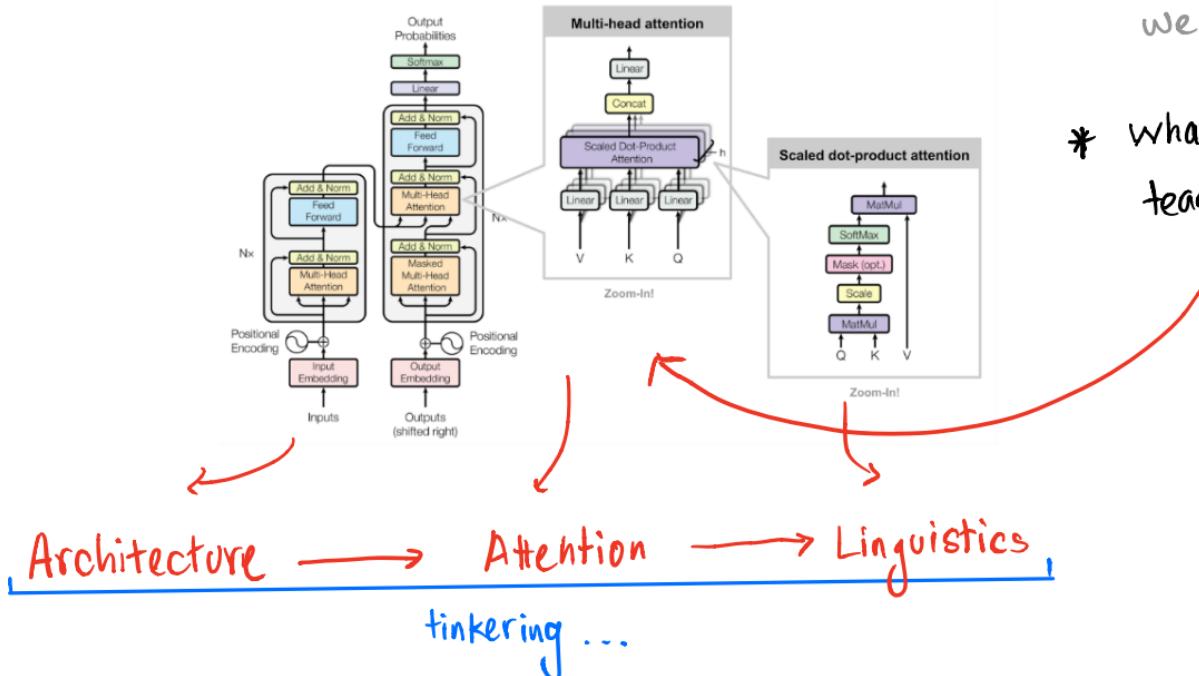
# BERTology ....

But ... why?

X

\* What is BERT  
and how did  
we get there ?

\* What does BERT  
teach us ?



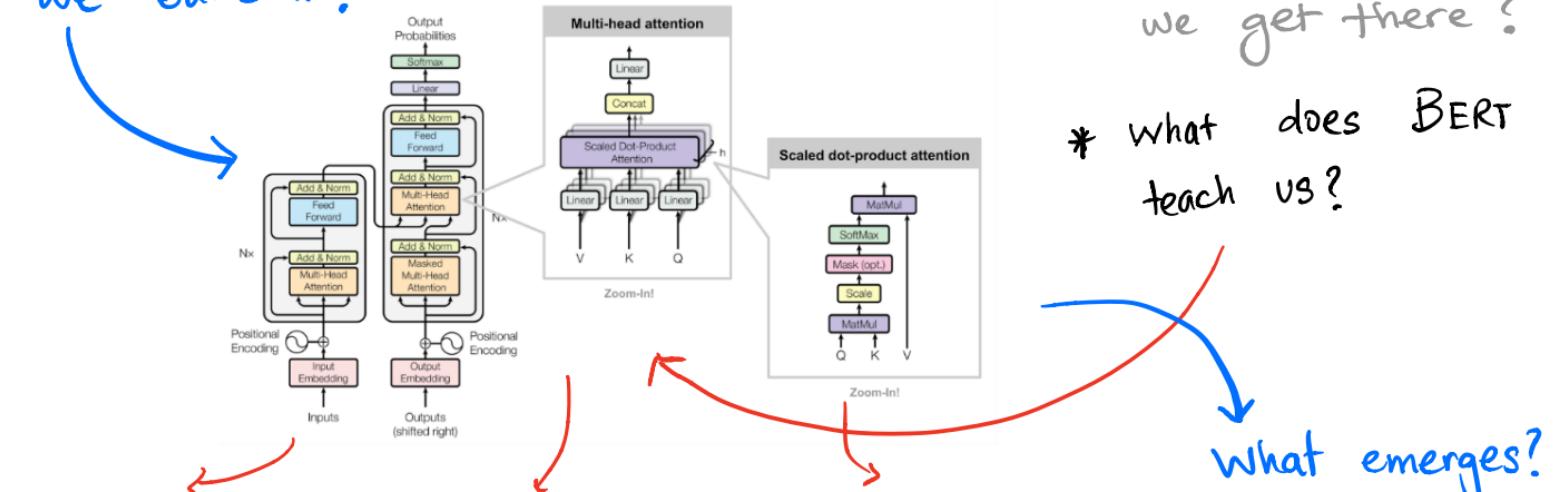
# BERTology ....

X

But ... why?

## Framework:

What we bake in?



Architecture → Attention → Linguistics,  
tinkering ...

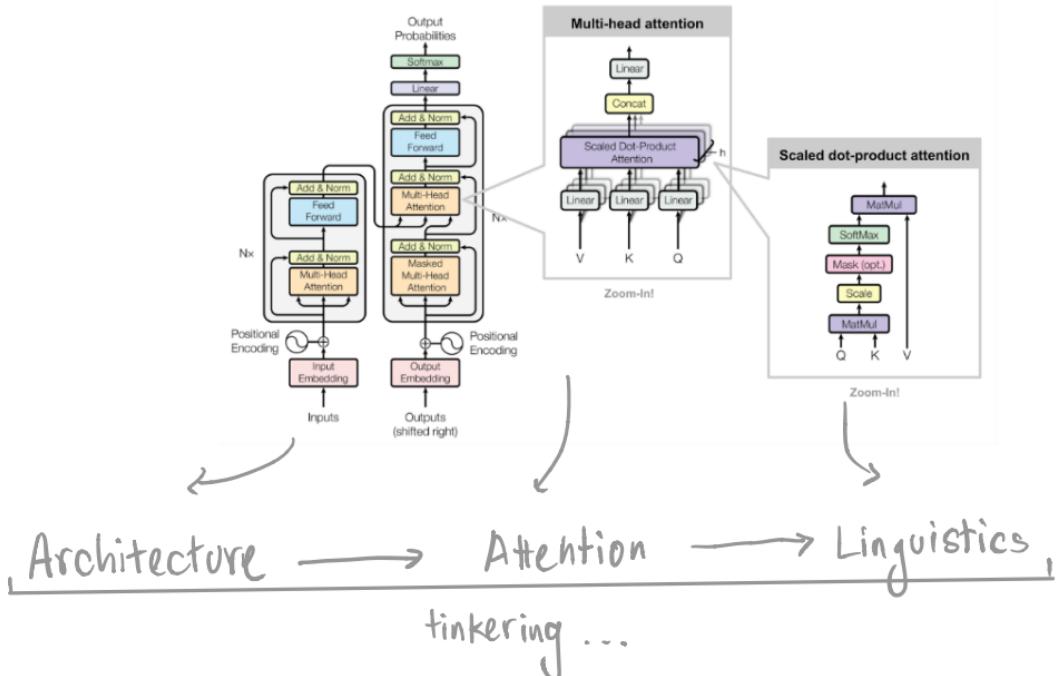
\* What is BERT  
and how did  
we get there?

\* what does BERT  
teach us?

what emerges?

# BERTology ....

But ... why? 



\* What is BERT and how did we get there?

\* What does BERT teach us?

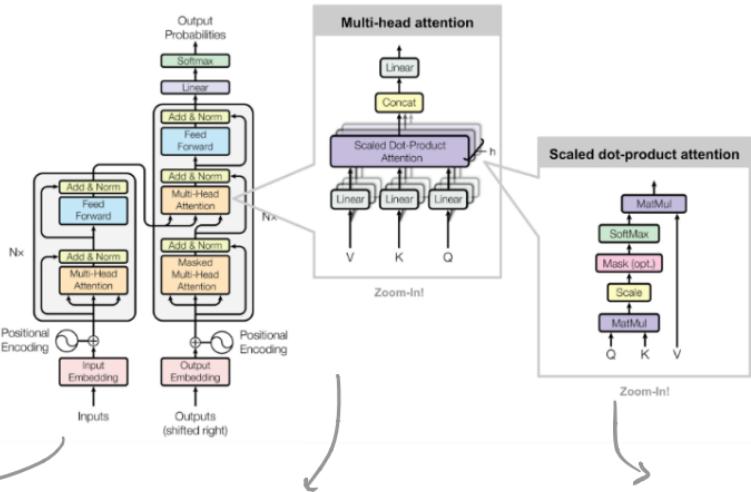
\* What could come next?

... do our lessons generalise?

- NLP
- vision
- learning systems

# BERTology ....

But ... why? 



Architecture → Attention → Linguistics,  
tinkering ...

\* BECAUSE WE ARE  
CURIOS HUMANS!

# RUNNING EXAMPLE

(AND MY PROGRAMMING TASK)

## Restaurant-reviews.txt

The	O	Food	B-ASP
bread	B-ASP	is	O
is	O	always	O
top	B-OP	fresh	B-OP
notch	I-OP	and	O
as	O	hot	B-OP
well	O	-	O
.	O	ready	O
		to	O
		eat	O
	!	O	

\* what is BERT  
and how did  
we get there?

\* what does BERT  
teach us?

\* what could come  
next?

... do our lessons  
generalise?

# RUNNING EXAMPLE

(AND MY PROGRAMMING TASK)

Restraunt-reviews.txt

The	O	Food	B-ASP
bread	B-ASP	is	O
is	O	always	O
top	B-OP	fresh	B-OP
notch	I-OP	and	O
as	O	hot	B-OP
well	O	-	O
.	O	ready	O
		to	O
		eat	O
		!	O

1485 sentences

tokenized, embedded  
with bert

(almost) all data  
in the presentation  
comes from this  
dataset!

\* what is BERT  
and how did  
we get there?

\* what does BERT  
teach us?

\* what could come  
next?  
... do our lessons  
generalise?

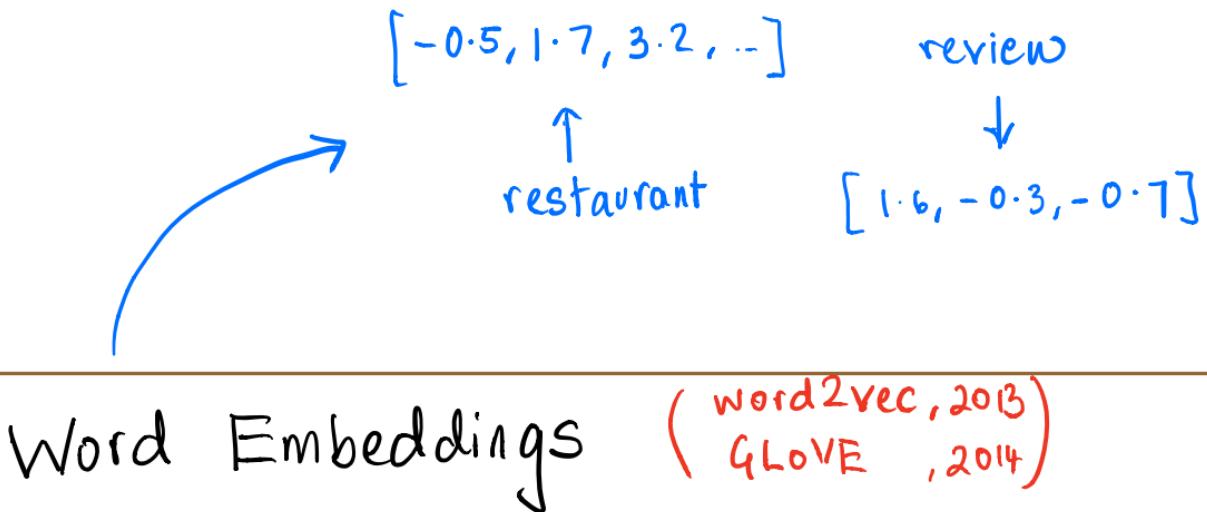
X

- \* what is BERT and how did we get there ?
- \* what does BERT teach us ?
- \* what could come next ?  
... do our lessons generalise ?

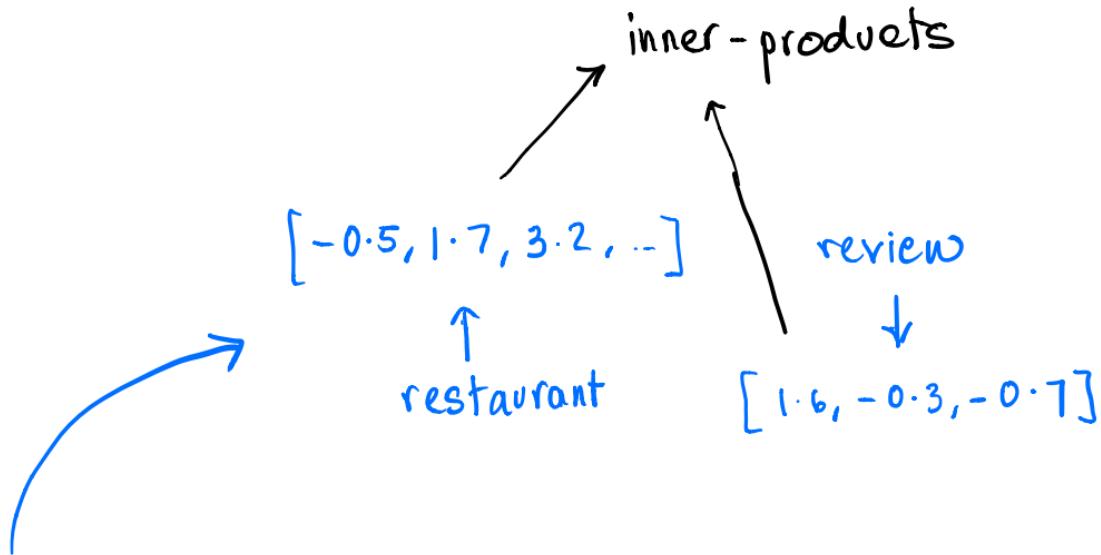
# A BERT RECAP

# A BERT RECAP

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?



# A BERT RECAP



Word Embeddings (word2vec, 2013)  
GLOVE , 2014) + Pre-training

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# A BERT RECAP

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

## Context & Representations (2010-2014)

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

# A BERT RECAP

"Certainly, not the best sushi" vs "BEST tuna roll"

$[0.3, 0.2, \dots]$   $\neq$   $[0.4, 0.5, \dots]$

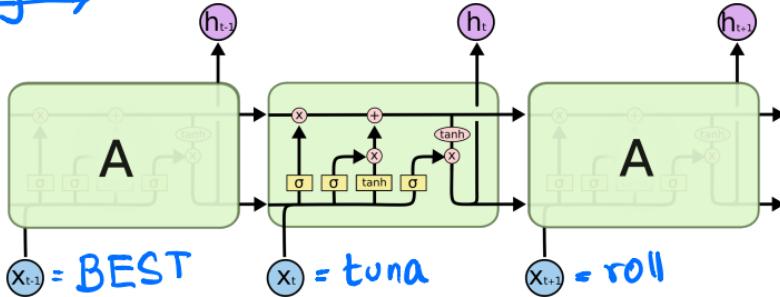
- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

## Context & Representations (2010-2014)

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

# A BERT RECAP

directionality →



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

---

"Large-scale" Supervision & LSTMs (Google, 2015)

---

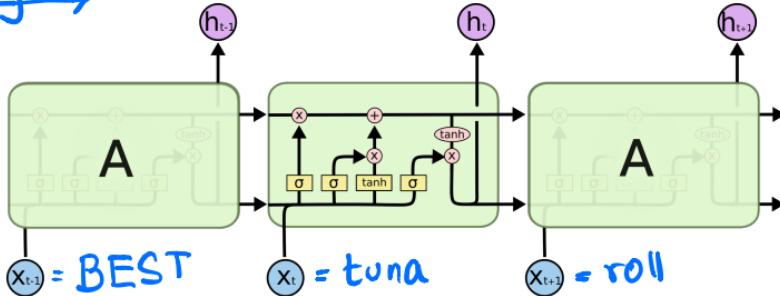
Context & Representations (2010 - 2014)

---

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

# A BERT RECAP

directionality →



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

---

"Large-scale" Supervision & LSTMs (Google, 2015)

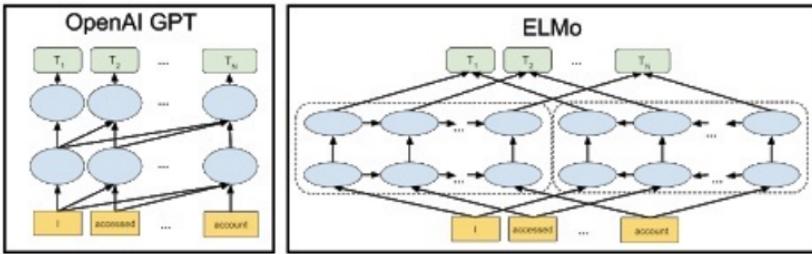
---

Context & Representations (2010 - 2014)

---

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

# A BERT RECAP



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

Bidirection "L  $\longleftrightarrow$  R" ( $ELM_o$ ) + Generative Pre-Training

"Large-scale" Supervision & LSTMs (Google, 2015)

Context & Representations (2010 - 2014)

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

# A BERT RECAP

Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

"Transformers"

→

Bidirectional "L  $\longleftrightarrow$  R" (ELMo) + Generative Pre-Training

"Large-scale" Supervision & LSTMs (Google, 2015)

Context & Representations (2010 - 2014)

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?
  - ... do our lessons generalise?

# A BERT RECAP

BERT =   
 Bidirectional  
 Encoder  
 Representations from  
 Transformers

Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

Bidirection "L  $\longleftrightarrow$  R" (ELMo) + Generative Pre-Training

"Large-scale" Supervision & LSTMs (Google, 2015)

Context & Representations (2010 - 2014)

Word Embeddings (<sup>word2vec, 2013</sup><sub>GLOVE, 2014</sub>) + Pre-training

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?
  - ... do our lessons generalise?



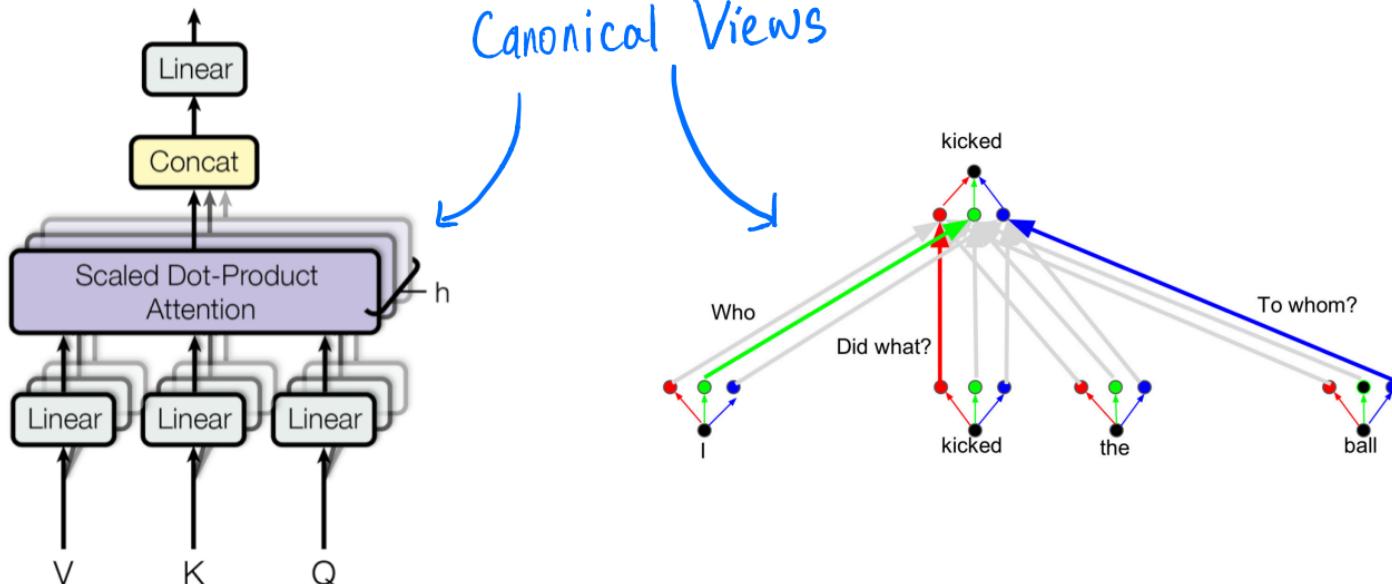
# A BERT RECAP

Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?
  - ... do our lessons generalise?

# A BERT RECAP

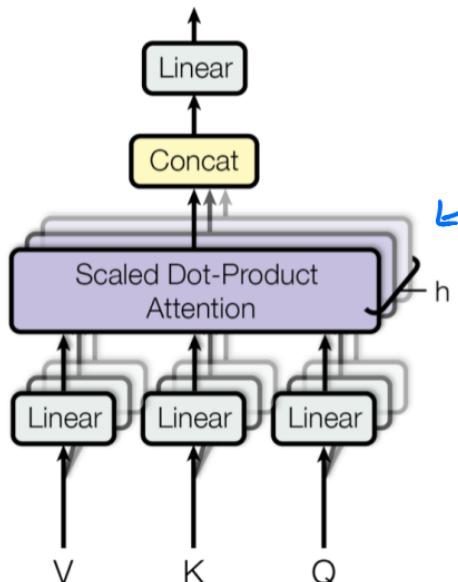
## Multi-headed "Attention" mechanisms (Vaswani et.al 2017)



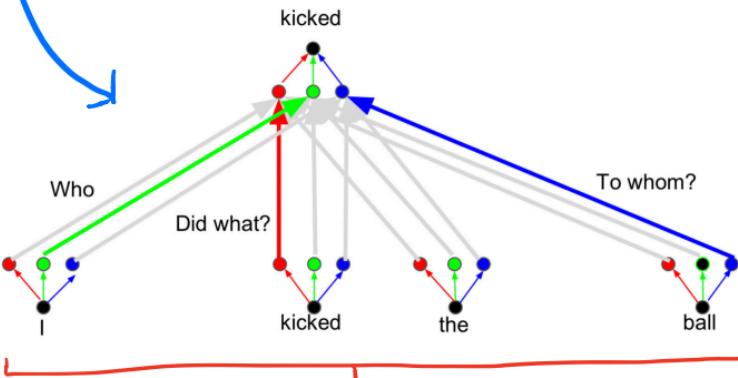
- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# A BERT RECAP

## Multi-headed "Attention" mechanisms (Vaswani et.al 2017)



Canonical Views

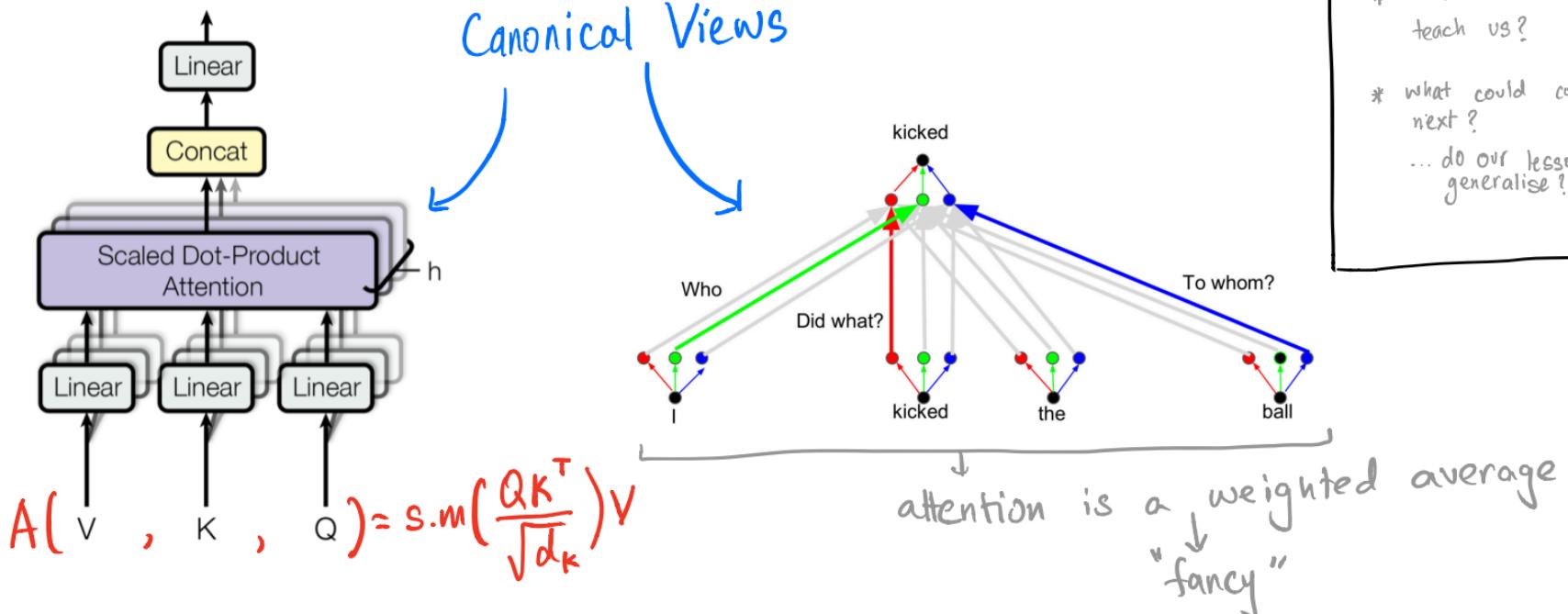


attention is a weighted average  
"fancy"

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# A BERT RECAP

## Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

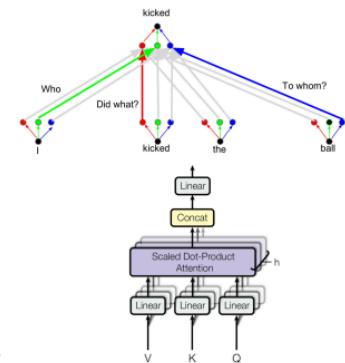


# A BERT RECAP

Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

But what do we bake in? ←

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?



# A BERT RECAP

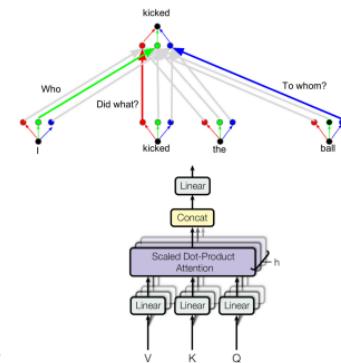
## Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

But what do we bake in? ←

- \* well-constructed "context"
- \* separate learning capacity  
.... allows composition!

but ... warning !!

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?



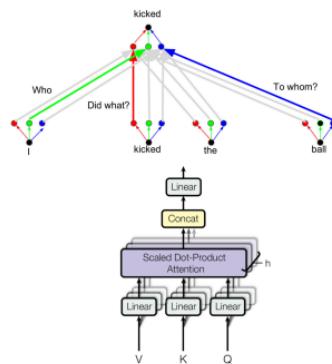
# A BERT RECAP

## Multi-headed "Attention" mechanisms (Vaswani et.al 2017)

But what do we bake in? ←

- \* well-constructed "context"
- \* separate learning capacity  
.... allows composition!
- \* curious to hear what you think??

- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?



X

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERTS ATTENTION

Clark, Khandewar, Levy, Manning

# AN ANALYSIS OF BERT'S ATTENTION

BUT WHY?

- \* Language modeling is the "ultimate" task
  - ... so we want to understand
  - ... how BERT models language.

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

BUT WHY?

- \* Language modeling is the "ultimate" task  
so we want to understand  
... how BERT models language.
- \* Our understanding of attention could better  
guide what we bake in future models.

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

BUT WHY?

- \* Language modeling is the "ultimate" task
  - so we want to understand
  - ... how BERT models language.
- \* Our understanding of attention could better guide what we bake in future models.
- \* Curiosity: we don't fully understand if BERT understands the "structure" of language? What about features?

\* What is BERT and how did we get there?

\* What does BERT teach us?

\* What could come next?

... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

BUT HOW?

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

BUT HOW?

- \* Syntactic trees. (Lin et.al.)
- \* Parts of speech/syntax (Tenney / Liu et.al.)
- \* Subject - Predicate Agreement. (Goldberg et.al.)
- \* Numerology (Wallace et.al.)
- \* Knowledge bases / Reasoning (Petroni et.al.)

spoiler: BERT  
is not great  
at this!

X

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn.  
⇒ ~50% next token, ~50% prev. token

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn?  
⇒ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
⇒ substantial amount of attn. to [SEP], [CLS]



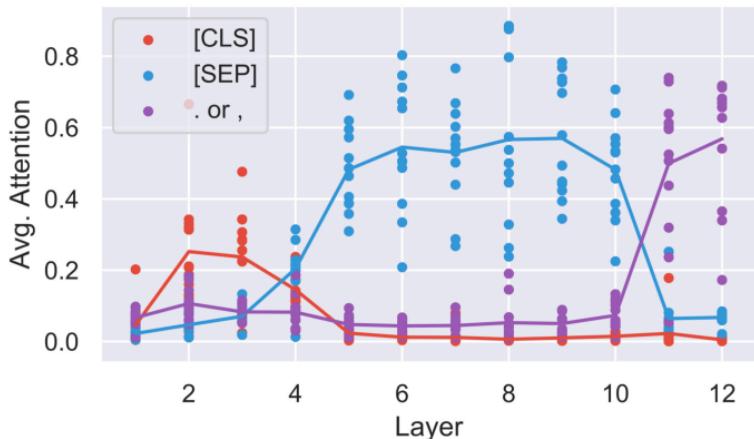
Reminder:  
[CLS] < sentence > [SEP]

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn?  
→ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
→ substantial amount of attn. to [SEP], [CLS]

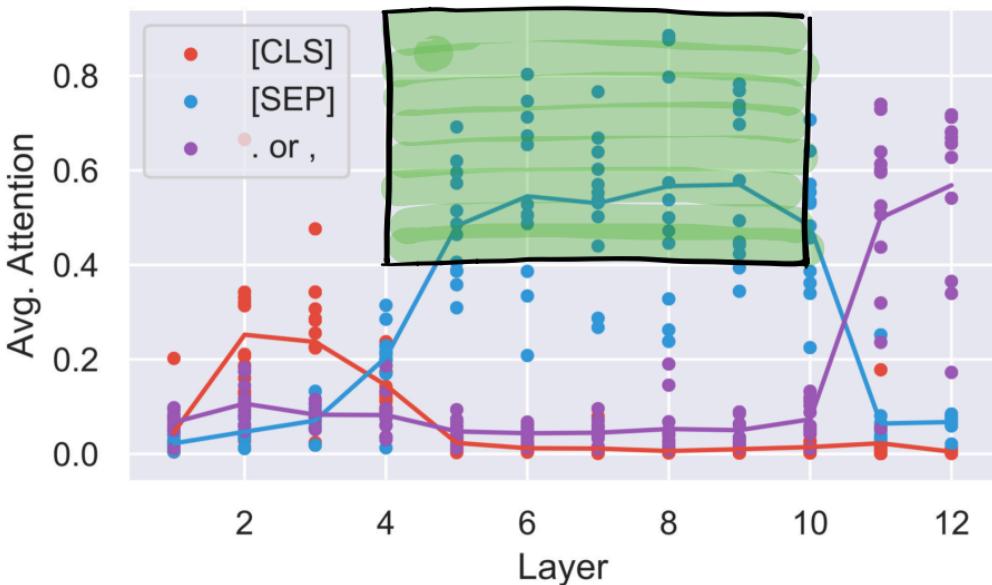


Reminder:  
[CLS] < sentence > [SEP]

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?



- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

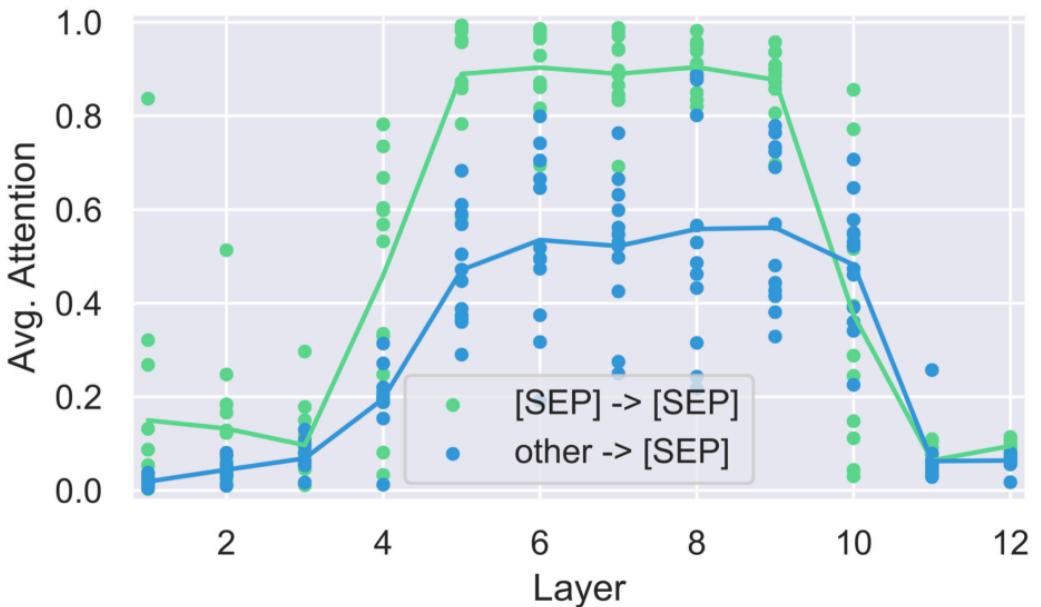
What is on the surface?

- \* Forward attn. vs Backward attn?  
→ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
→ substantial amount of attn. to [SEP], [CLS]  
hmm... why?  
\* Maybe, "aggregate surface-level information"?

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?



- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

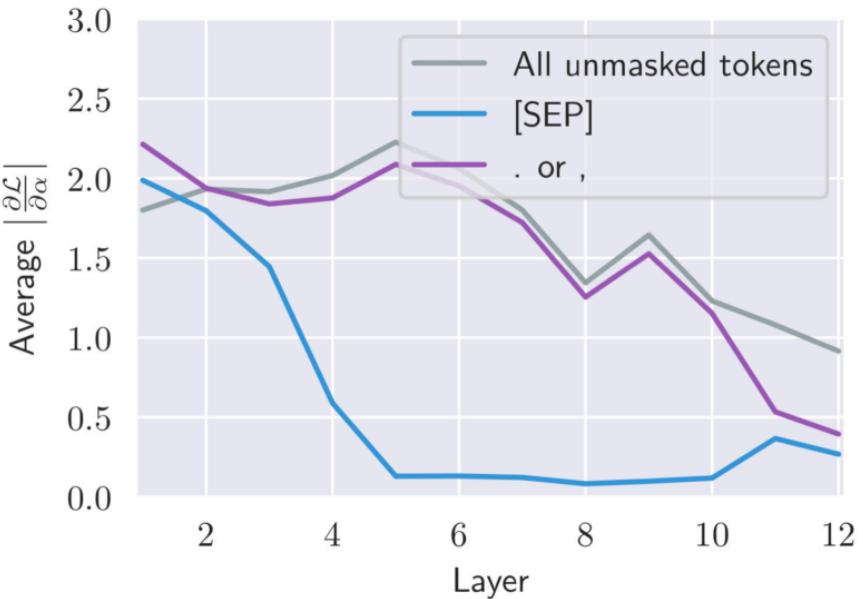
What is on the surface?

- \* Forward attn. vs Backward attn?  
    ⇒ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
    ⇒ substantial amount of attn. to [SEP], [CLS]  
        hmm... why?  
    → Maybe, "aggregate surface-level information"?
- \* Maybe, its just a "NO-OP"

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?



- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn?  
→ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
→ substantial amount of attn. to [SEP], [CLS]  
hmm... why?  
→ \* Maybe, "aggregate surface-level information"?
- ✓ \* Maybe, its just a "NO-OP"

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn?  
    ⇒ ~50% next token, ~50% prev. token
- \* What about specific tokens?  
    ⇒ substantial amount of attn. to [SEP], [CLS]
- \* Focused vs. Broad Attn.

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

- \* Forward attn. vs Backward attn?  
⇒ ~50% next token, ~50% prev. token
  - \* What about specific tokens?  
⇒ substantial amount of attn. to [SEP], [CLS]
  - \* Focused vs. Broad Attn.  
⇒ lower layers : "broad" attention  
higher layers : "focused" attention
- high entropy.

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

(0,0)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

(2,1)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

"prev"  
word  
(2,3)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

(10,5)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

"broad" attention

"next" word

exclusive  
[SEP]  
focus.

\* What is BERT and how did we get there?

\* What does BERT teach us?

\* What could come next?

... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

What is on the surface?

specificity  
increases  
with  
depth!

(0,0)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

(2,1)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

"prev"  
word  
(2,3)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

(10,5)

Food  
is  
always  
fresh  
and  
hot  
-  
ready  
to  
eat  
!  
[SEP]

"broad" attention

"next" word

exclusive  
[SEP]  
focus.

\* What is BERT  
and how did  
we get there?

\* What does BERT  
teach us?

\* What could come  
next?  
... do our lessons  
generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

Idea: A word "predicts" another word if it gets most of its attention.

Do "predictions" learn language features?

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?
  - ... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

Idea: A word "predicts" another word if it gets most of its attention.

Do "predictions" learn language features?

Evaluation: Check accuracy with Penn Tree-bank annotated with Stanford Dependencies!

⇒ no single head does well overall, but they do very well with "composition"

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

common relations:

- prepositions

- noun subjects

- adjective  
modifiers

Relation	Head	Accuracy	Baseline
All	7-6	34.5	26.3 (1)
prep	7-4	66.7	61.8 (-1)
pobj	9-6	<b>76.3</b>	34.6 (-2)
det	8-11	<b>94.3</b>	51.7 (1)
nn	4-10	70.4	70.2 (1)
nsubj	8-2	58.5	45.5 (1)
amod	4-10	75.6	68.3 (1)
dobj	8-10	<b>86.8</b>	40.0 (-2)
advmod	7-6	48.8	40.2 (1)
aux	4-10	81.1	71.5 (1)
poss	7-6	<b>80.5</b>	47.7 (1)
auxpass	4-10	<b>82.5</b>	40.5 (1)
ccomp	8-1	<b>48.8</b>	12.4 (-2)
mark	8-2	<b>50.7</b>	14.5 (2)
prt	6-7	<b>99.1</b>	91.4 (-1)

Best head!

\* what is BERT  
and how did  
we get there?

\* what does BERT  
teach us?

\* what could come  
next?

... do our lessons  
generalise?

surprising difference.

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

But ... what exactly do these look like ?

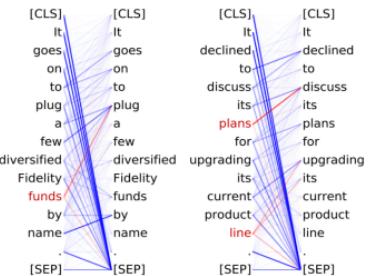
- \* What is BERT and how did we get there ?
- \* What does BERT teach us ?
- \* What could come next ?
  - ... do our lessons generalise ?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

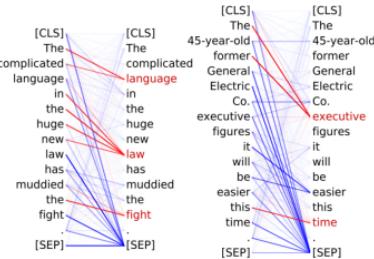
Head 8-10

- Direct objects attend to their verbs
- 86.8% accuracy at the dobj relation



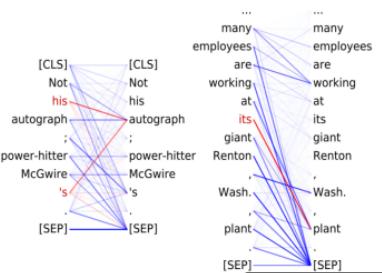
Head 8-11

- Noun modifiers (e.g., determiners) attend to their noun
- 94.3% accuracy at the det relation



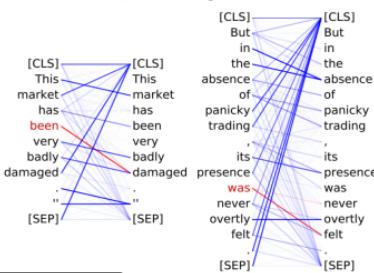
Head 7-6

- Possessive pronouns and apostrophes attend to the head of the corresponding NP
- 80.5% accuracy at the poss relation



Head 4-10

- Passive auxiliary verbs attend to the verb they modify
- 82.5% accuracy at the auxpass relation



\* what is BERT and how did we get there?

\* what does BERT teach us?

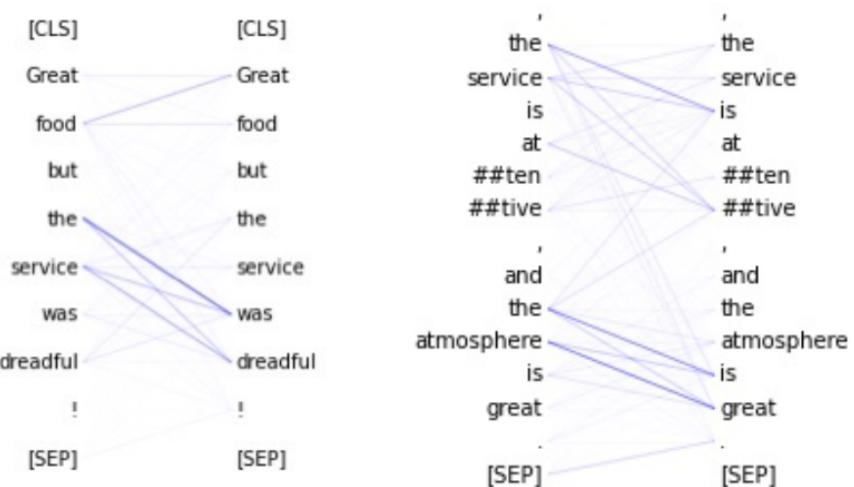
\* what could come next?

... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

amodj: adjectival modifiers



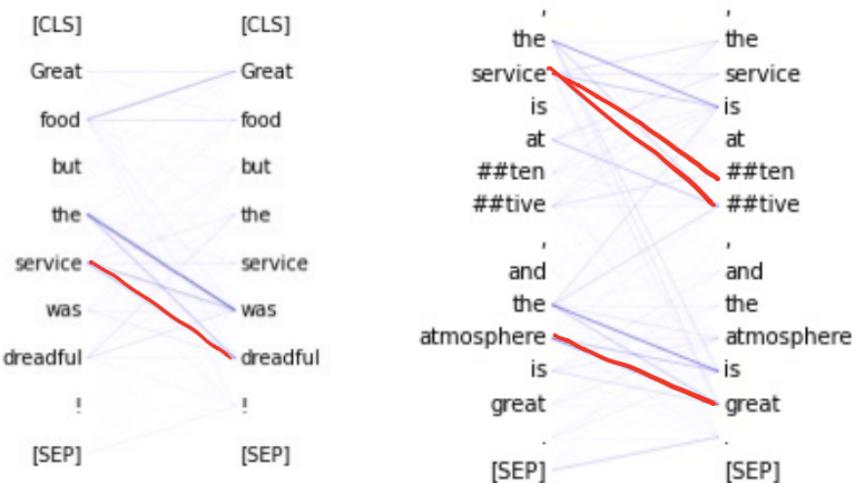
- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

amodj: adjectival modifiers

(8,4)

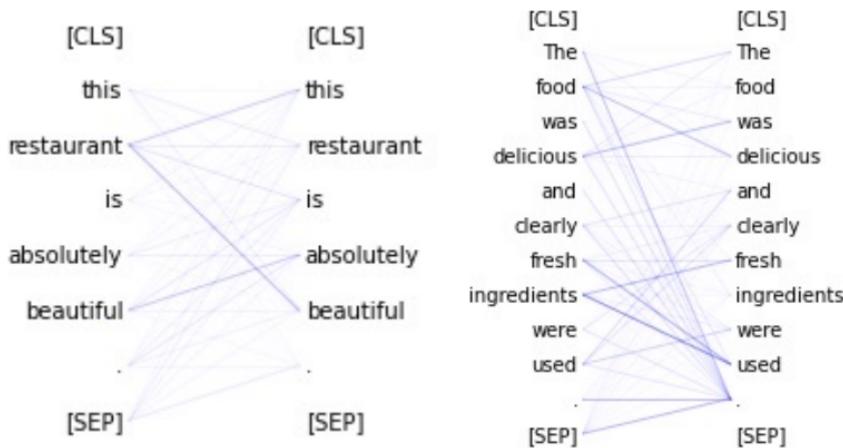


- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

n modj: nominal modifiers



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

n modj: nominal modifiers

(9,7)



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

early "broad" focused heads.

(0,0)



- \* what is BERT and how did we get there?
- \* what does BERT teach us?
- \* what could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## LEARNING FEATURES OF LANGUAGE

How far can we push extracting dependencies?

Idea: probing heads

- \* What is BERT and how did we get there?
  - \* What does BERT teach us?
  - \* What could come next?
    - ... do our lessons generalise?
-

# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

What kind of similarities do heads have?

- \* What is BERT and how did we get there?
  - \* What does BERT teach us?
  - \* What could come next?
    - ... do our lessons generalise?
-

# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

What kind of similarities do heads have?

Idea: Compute distances ...

$$d(H_i, H_j)$$

plot this on  
2d.

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

What kind of similarities do heads have?

Idea: Compute

$$d(H_i, H_j) = JS(H_i, H_j)$$

plot this on  
2d.



$$JS(p||q) = \frac{1}{2} D_{KL}(p || \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q || \frac{p+q}{2})$$

$$D(p||q) = \int p \log \frac{p}{q} dx$$

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

What kind of similarities do heads have?

Idea: Compute

$$d(H_i, H_j) = \sum_{\text{tok data}} JS(H_i(\text{tok}), H_j(\text{tok}))$$

plot this on  
2d.

$$JS(p||q) = \frac{1}{2} D_{KL}(p || \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q || \frac{p+q}{2})$$

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

What kind of similarities do heads have?

Idea: Compute  $d(H_i, H_j) = \sum_{\text{tok} \in \text{data}} \text{JS}(H_i(\text{tok}), H_j(\text{tok}))$

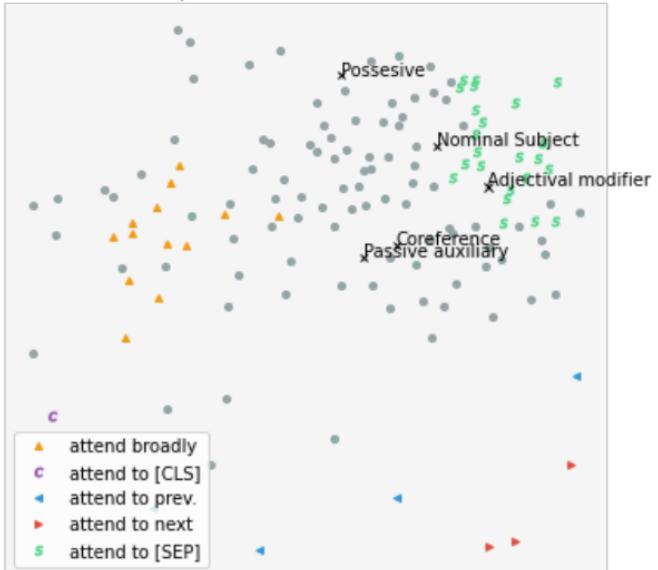
... plot!! (hammer: multi-dimensional scaling)

- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

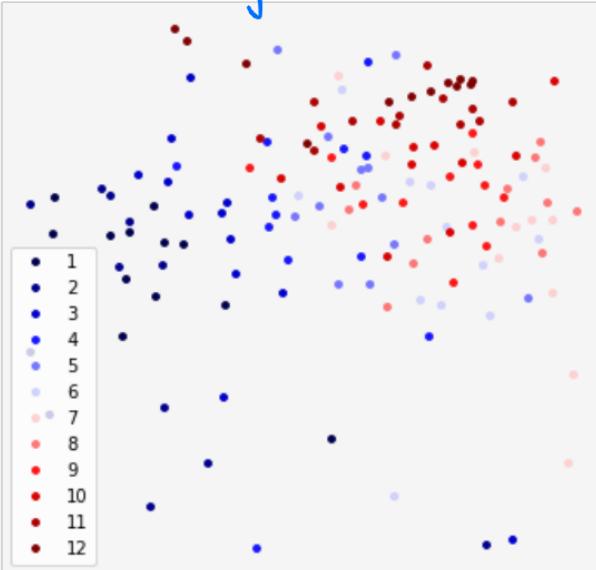
# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

behaviour:



layers:

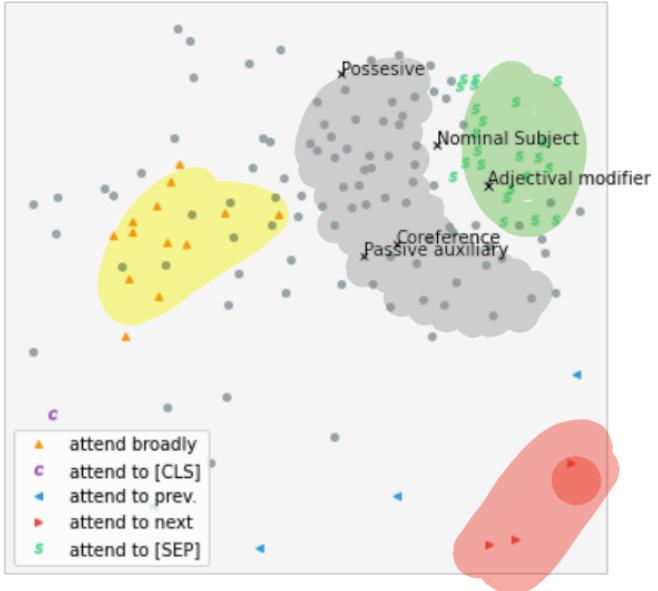


- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

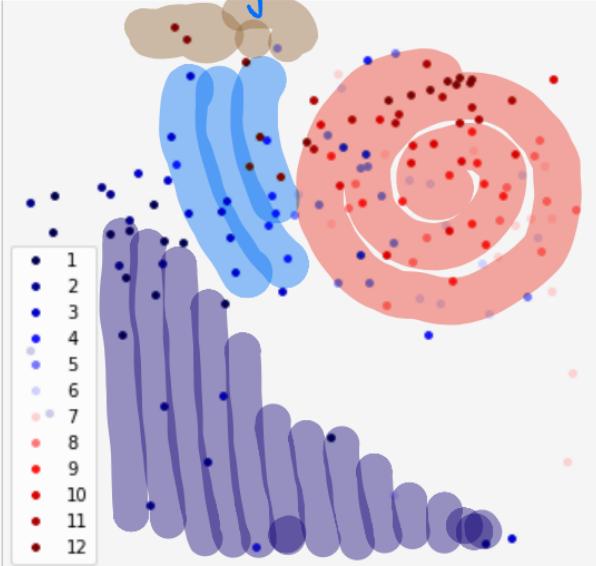
# AN ANALYSIS OF BERT'S ATTENTION

## Similarity and Composition

behaviour:



layers:



- \* What is BERT and how did we get there?
- \* What does BERT teach us?
- \* What could come next?  
... do our lessons generalise?

# WHAT COMES NEXT

X

Yeah ..... so ????

- \* what is BERT and how did we get there ?
- \* what does BERT teach us ?
- \* what could come next ?  
... do our lessons generalise ?

# WHAT COMES NEXT

Yeah..... so ????

\* We have a **Toolbox !!** to study language nets.

H

- \* What is BERT and how did we get there ?
- \* What does BERT teach us ?
- \* What could come next ?
  - ... do our lessons generalise ?

# WHAT COMES NEXT

Yeah.... so ????

- \* We have a Toolbox !! to study language nets.  
*boring?*
- \* "Attention" is more interesting than vectors/neurons.  
i.e. turn "our" attention to higher level abstractions.  
think multi-modal, many agents ....

H

- \* What is BERT and how did we get there ?
- \* What does BERT teach us ?
- \* What could come next ?  
... do our lessons generalise ?

# WHAT COMES NEXT

Yeah.... so ????

- \* We have a Toolbox !! to study language nets.  
→ boring?
- \* "Attention" is more interesting than vectors/neurons.  
i.e. turn "our" attention to higher level abstractions.  
think multi-modal, many agents ....
- \* We shouldn't be afraid of using more computation!  
i.e. there is hope and we should **keep trying** ...

X

- \* What is BERT and how did we get there ?
- \* What does BERT teach us ?
- \* What could come next ?  
... do our lessons generalise ?

# WHAT COMES NEXT

Yeah..... so ????

- \* We have a Toolbox !! to study language nets.  
→ boring?
- \* "Attention" is more interesting than vectors/neurons.  
i.e. turn "our" attention to higher level abstractions.  
think multi-modal, many agents ....
- \* We shouldn't be afraid of using more computation!  
i.e. there is hope and we should keep trying ...  
what do all of you think??

X

- \* What is BERT and how did we get there ?
- \* What does BERT teach us ?
- \* What could come next ?  
... do our lessons generalise ?

# BERTology review: What Does BERT Look At? An Analysis of BERTs attention

Kevin Clark, Urvashi Khandelwal, Omer Levy, Christopher Manning (2019)

Sualeh Asif

MIT Computer Science Artificial Intelligence Laboratory  
Massachusetts Institute of Technology

