# Machine Learning

## Representation learning

Dan Goldwasser

dgoldwas@purdue.edu

# Language Models

- **A language model over a given vocabulary V assigns probabilities to strings drawn from V\***

  Can we actually do it?

$$P_{ngram}(w_1...w_i) := P(w_1)P(w_2|w_1)...P(\underbrace{w_i}_{nth\ word} | \underbrace{w_{i-n-1}...w_{i-1}}_{prev.\ n-1\ words})$$

| | |
|---|---|
| **Unigram** | $P(w_1)P(w_2)...P(w_i)$ |
| **Bigram** | $P(w_1)P(w_2|w_1)...P(w_i|w_{i-1})$ |
| **Trigram** | $P(w_1)P(w_2|w_1)...P(w_i|w_{i-2}\ w_{i-1})$ |

# Example: Trigram language model

- Consider the sentence:

*Mr. Smith goes*

$$p(\text{Mr. Smith goes STOP}) = p(\text{Mr.}|*, *))$$
$$p(\text{Smith}|*, \text{Mr.}) \ p(\text{goes}|\text{Mr., Smith}) \ p(\text{STOP}|\text{Smith, goes})$$

# Language Models

| Unigram | $P(w_1)P(w_2)...P(w_i)$ |
|---------|--------------------------|
| **Bigram** | $P(w_1)P(w_2|w_1)...P(w_i|w_{i-1})$ |
| **Trigram** | $P(w_1)P(w_2|w_1)...P(w_i|w_{i-2}\ w_{i-1})$ |

Our goal is to assess whether

**P(Private Customer… Be Toad)**

**>?<**

**P(Private Customer… Be Towed)**

What would be the answer if we use –
**(1)** a *Unigram* model? **(2)** a *Bigram* model?

# Language Models

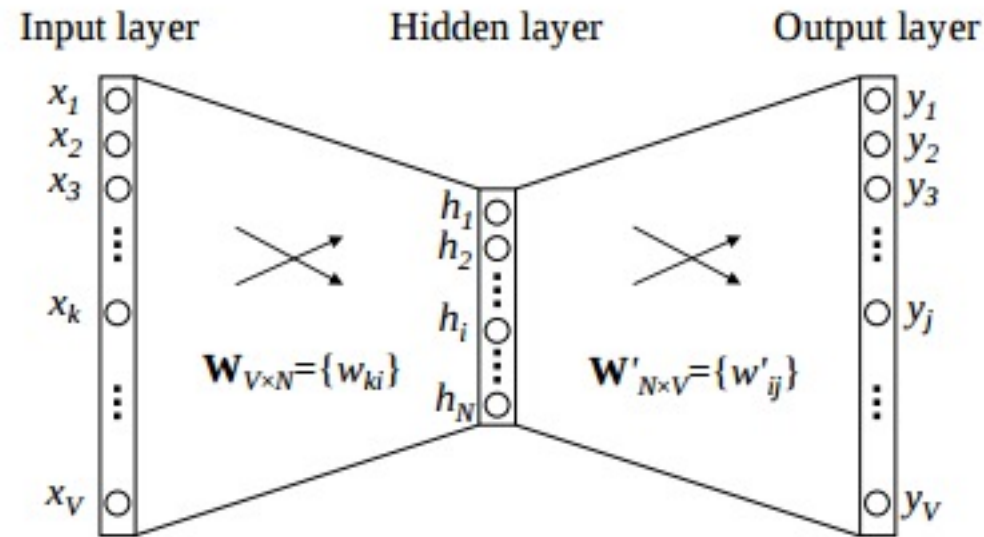| | |
|---|---|
| **Unigram** | • To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have<br>• Every enter now severally so, let<br>• Hill he late speaks; or! a more to leg less first you enter<br>• Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like |
| **Bigram** | • What means, sir. I confess she? then all sorts, he is trim, captain.<br>•Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.<br>•What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?<br>•Enter Menenius, if it so many good direction found'st thou art a strong upon command of fear not a liberal largess given away, Falstaff! Exeunt |
| **Trigram** | • Sweet prince, Falstaff shall die. Harry of Monmouth's grave.<br>• This shall forbid it should be branded, if renown made it empty.<br>• Indeed the duke; and had a very good friend.<br>• Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done. |
| **Quadrigram** | • King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;<br>• Will you not tell me who I am?<br>• It cannot be but so.<br>• Indeed the short and the long. Marry, 'tis a noble Lepidus. |

# So how did we get here?

After my fiancé died, my mother told me to "get out there again." She wanted me to go to a singles bar. I told her I'd rather go to the dentist.

"Just once," she said. "Just to see what it's like."

One day, early last year, I found myself driving to a singles bar in winter snow. I sat in my car for 15 minutes, then drove away. The next day, I went back and sat in my car for another 15 minutes. I did this for a couple of weeks, until I finally mustered up the nerve to walk in.
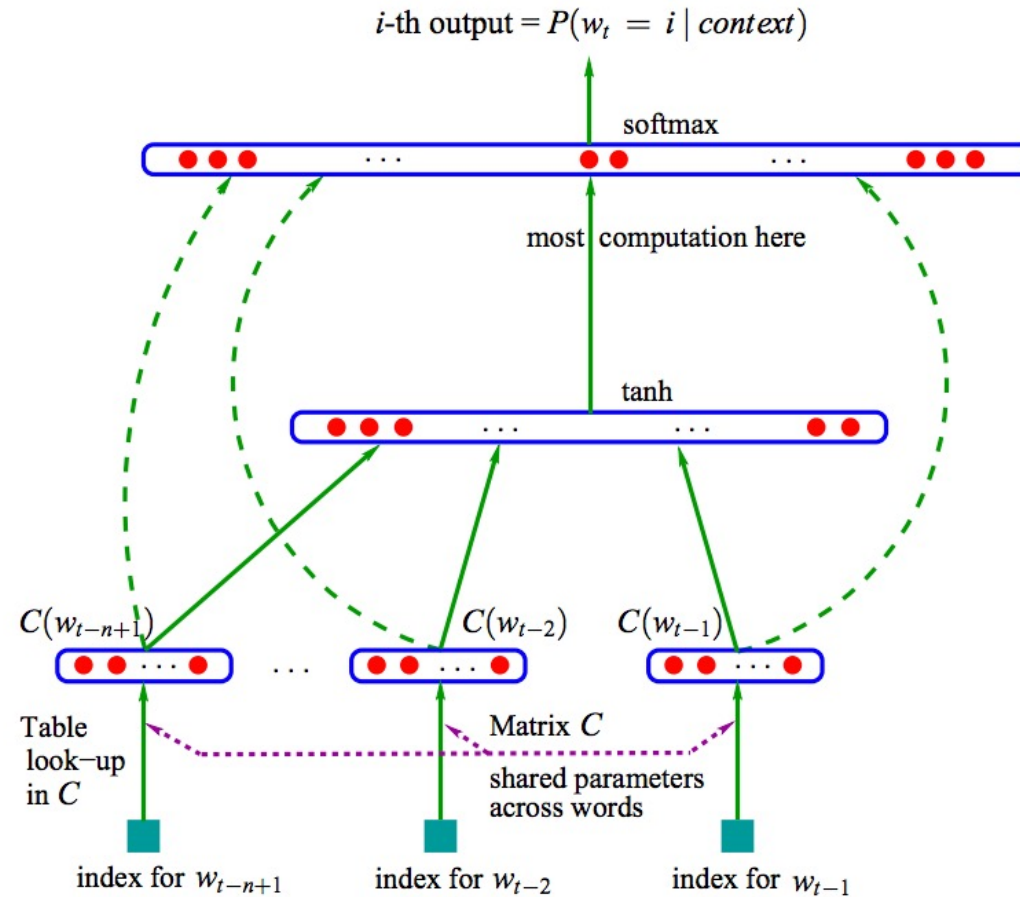
# Word Embedding (no context)
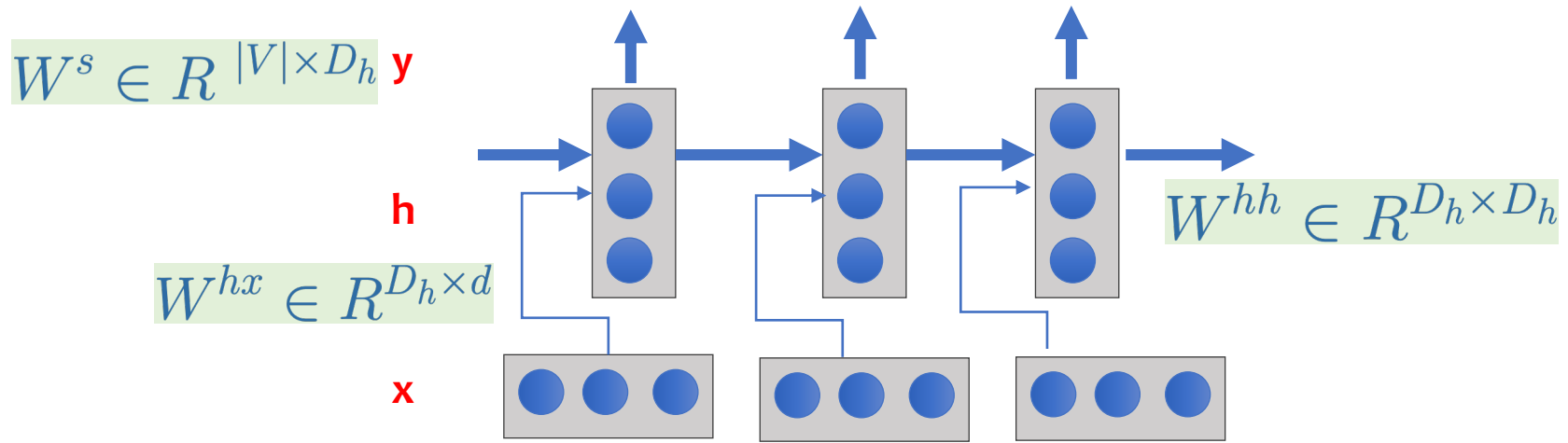
- Is this the same as an Auto-Encoder?



word2vec model architecture

# Neural Language Model – fixed context



A Neural Probabilistic Language Model. Bengio et-al 2003

# Recurrent Neural Networks (one/bi *directional* unbounded context)

$$W^s \in R^{|V| \times D_h} \quad \text{y}$$

$$\text{h}$$

$$W^{hh} \in R^{D_h \times D_h}$$

$$W^{hx} \in R^{D_h \times d}$$

$$\text{x}$$

**Input is a word (vectors) sequence:** $\quad x_1, ..., x_{i-1}, x_i, x_{i+1}, ..., x_n$

**At any given time step *i* :** $\qquad h_i = \sigma \left( W^{hh} h_{i-1} + W^{hx} x_i \right)$

$$\hat{y} = \text{softmax}(W^s \ h_i)$$

$$P(x_{i+1} = v_j | x_i, ..., x_1) = \hat{y}_{i,j}$$

# RNN Extensions

- **Key issue**: *long range dependencies between inputs.*
  - "how can we know which word is important to keep around, when predicting the i+1 word?"

- **Solution idea**: *complex hidden units that implement a "memory"*
  - Maintain "old memories" representing relevant long range dependencies
  - Error updates can be back-propagated at different strengths.

# Gated Recurrent Units (GRU)

$$h_i = \sigma \left( W^{hh} h_{i-1} + W^{hx} x_i \right)$$ ⬅ **Original RNN**

**GRU:**

**Update Gate:** $$z_i = \sigma \left( W^z x_i + U^z h_{i-1} \right)$$

**Reset Gate:** $$r_i = \sigma \left( W^r x_i + U^r h_{i-1} \right)$$

**New memory** $$\tilde{h}_i = \tanh \left( W x_i + r_i \circ U h_{i-1} \right)$$

**Final memory (aka *hidden Layer*)** $$h_i = z_t \circ h_{i-1} + \left( 1 - z_i \right) \circ \tilde{h}_i$$

# (Very) Modern Love

After my fiancé died, my mother told me to "get out there again." She wanted me to go to a singles bar. I told her I'd rather go to the dentist.
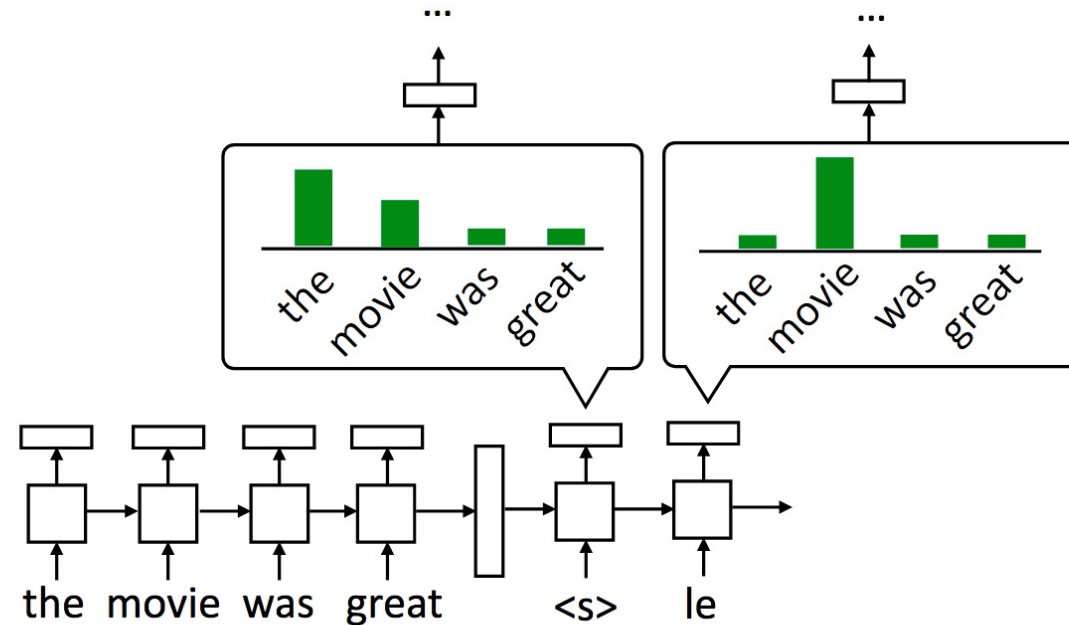
"Just once," she said. "Just to see what it's like."

One day, early last year, I found myself driving to a singles bar in winter snow. I sat in my car for 15 minutes, then drove away. The next day, I went back and sat in my car for another 15 minutes. I did this for a couple of weeks, until I finally mustered up the nerve to walk in.

# *Attention based "reading"*

- Many of the challenges introduced by machine comprehension can be addressed using a general solution based on **attention**

- A general tool, currently used in **all** NLP tasks
  - Essentially, learn meaningful associations between inputs and outputs which can represent structural dependencies.

# Attention

- **Attention**: at each decoder state computes a **distribution over the source inputs** based on the current decoder state
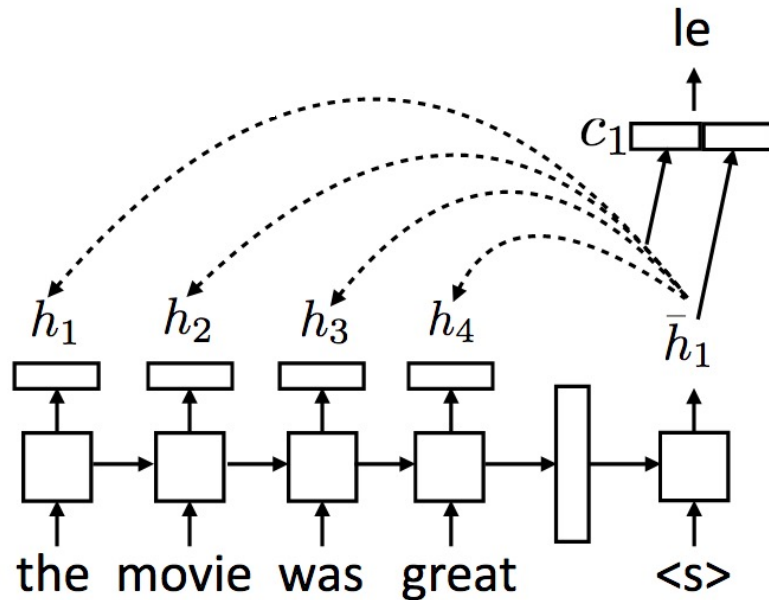
# Attention

- For each decoder state compute the weighted sum of input states

$$P(y_i | \mathbf{x}, y_1, \ldots, y_{i-1}) = \text{softmax}(W[c_i; \bar{h}_i])$$

le

$c_1$

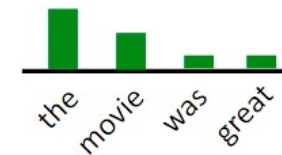$h_1$ $h_2$ $h_3$ $h_4$ $\bar{h}_1$

the movie was great  &lt;s&gt;

$$c_i = \sum_j \alpha_{ij} h_j$$

▸ Weighted sum of input hidden states (vector)

the movie was great

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

$$e_{ij} = f(\bar{h}_i, h_j)$$
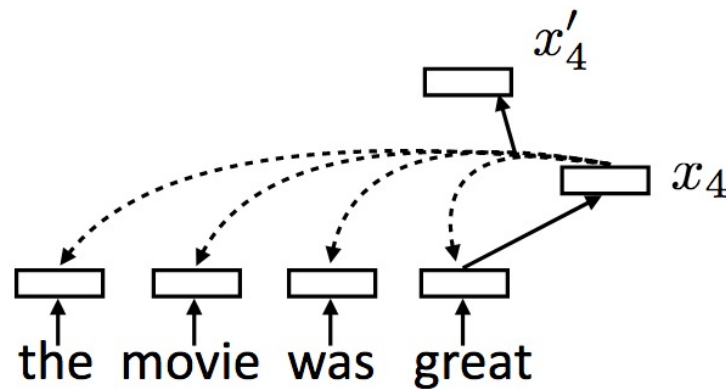
▸ Unnormalized scalar weight

# Attention

We can identify the source word context for output predictions

# Self-attention

- **A new way to represent structure**
  - *Each word forms a query which computes attention over each word*



$$\alpha_{i,j} = \text{softmax}(x_i^\top x_j) \quad \text{scalar}$$

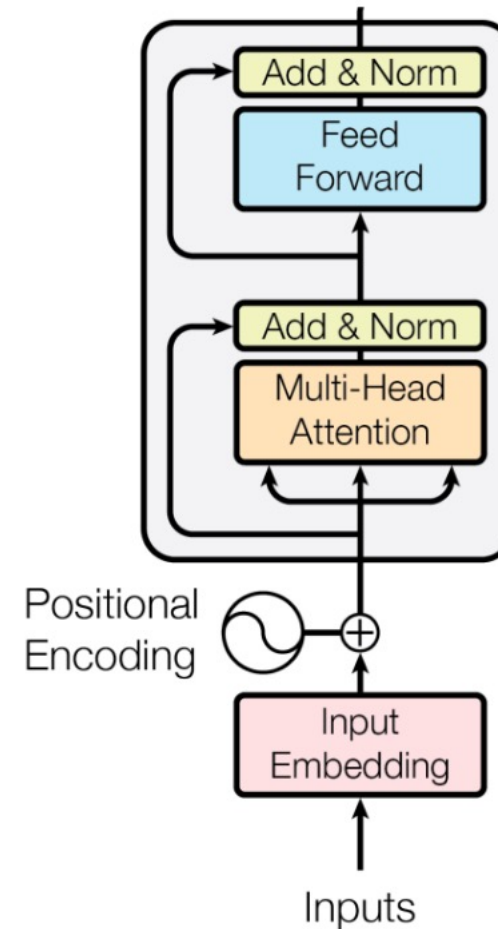$$x_i' = \sum_{j=1}^{n} \alpha_{i,j} x_j \quad \text{vector = sum of scalar * vector}$$

The representation of each word is a function of its neighbors. Does that sound familiar?

Vaswani et al. (2017)

# Transformers

- The idea of self attention was extremely influential in NLP
  - **No fixed position representation as in LSTM instead structure is represented the attention assignments.**
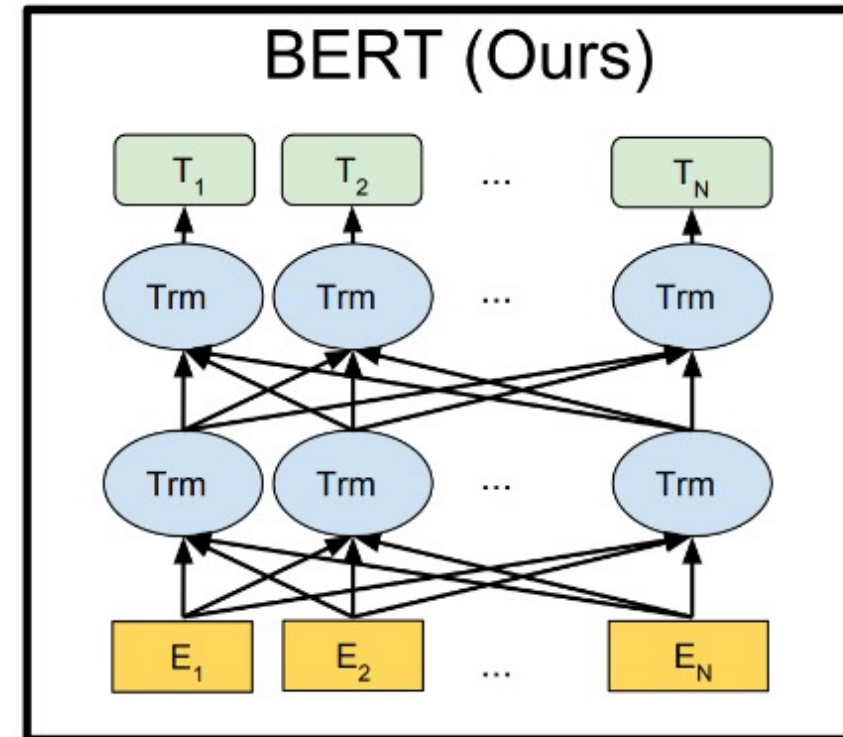
  <span style="color:red">**I like bananas but not carrots.**
  **Vs.**
  **I like carrots but not bananas**</span>

- In reality, position information is needed, but it is used differently compared to an LSTM, by encoding it as part of the input
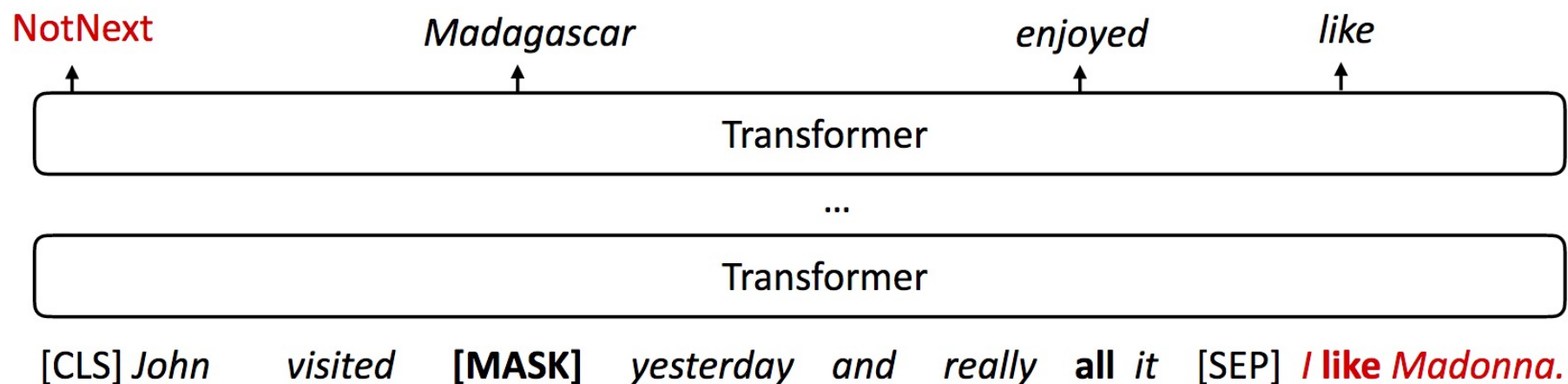
# BERT

- **Transformer-based approach** *instead of an LSTM-based like ELMo.*
  - Transformer vs. LSTM
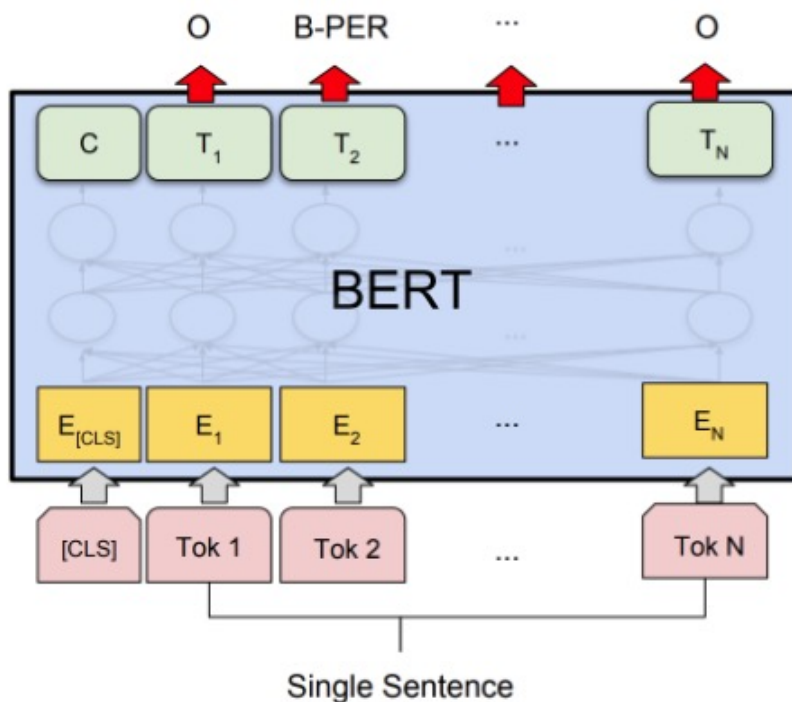  - Masked language objective instead of usual LM
  - Fine-tuned at test time

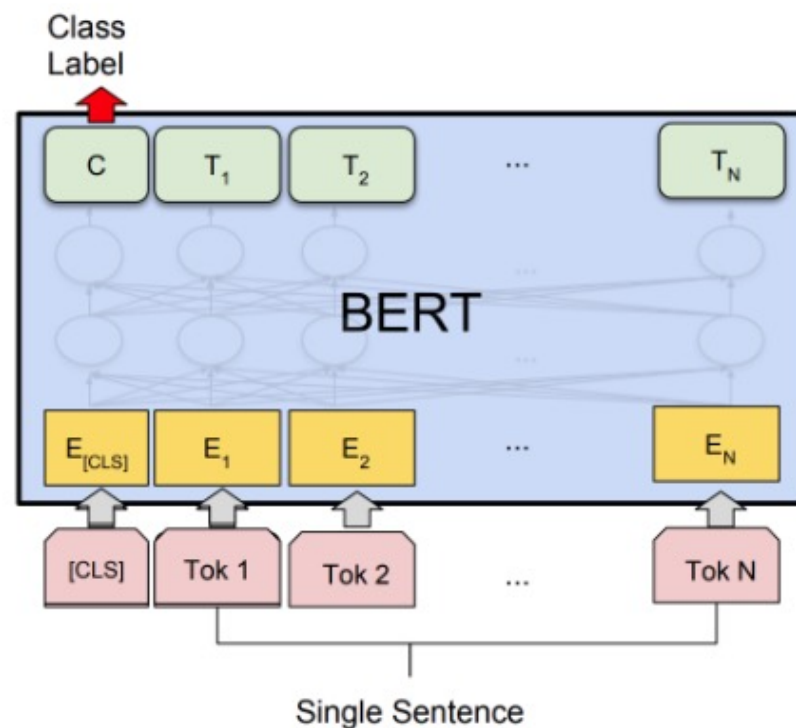# Next sentence

- **BERT objective**: masked LM + next sentence



NotNext    Madagascar    enjoyed    like

Transformer

...

Transformer

[CLS] *John    visited    **[MASK]**    yesterday    and    really    **all** it    [SEP] I **like** Madonna.*

# BERT in practice

**Very flexible, can be used for NLI, classification, tagging, etc.**



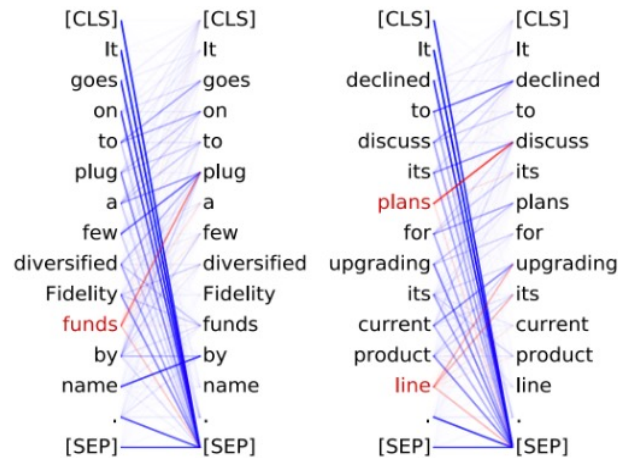(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

(b) Single Sentence Classification Tasks: SST-2, CoLA
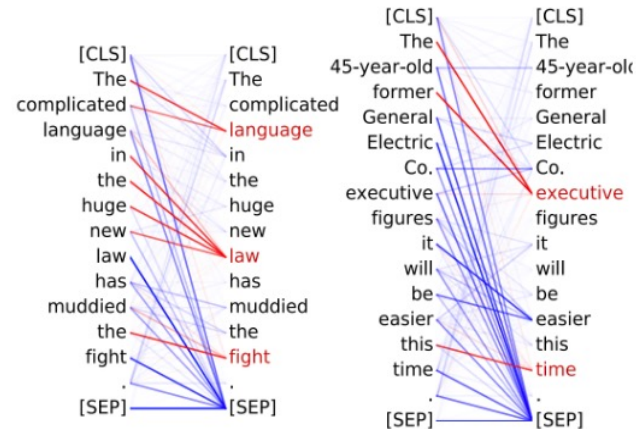
# What does BERT learn?



Head 8-10
- **Direct objects** attend to their verbs
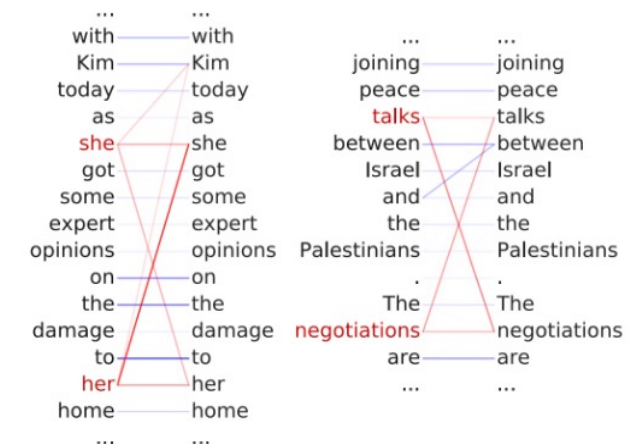- 86.8% accuracy at the dobj relation

Head 8-11
- **Noun modifiers** (e.g., determiners) attend to their noun
- 94.3% accuracy at the det relation

Head 5-4
- **Coreferent** mentions attend to their antecedents
- 65.1% accuracy at linking the head of a coreferent mention to the head of an antecedent

▸ Still way worse than what supervised systems can do, but interesting that this is learned organically

# Discussion

- **Current NLP trend**: *train a **very complex** neural language model using **massive** amounts of data*
- The learned representation should capture "language understanding capability"
  - Word meaning
  - Linguistic structure
  - World knowledge
  - Bad stuff expressed through language.

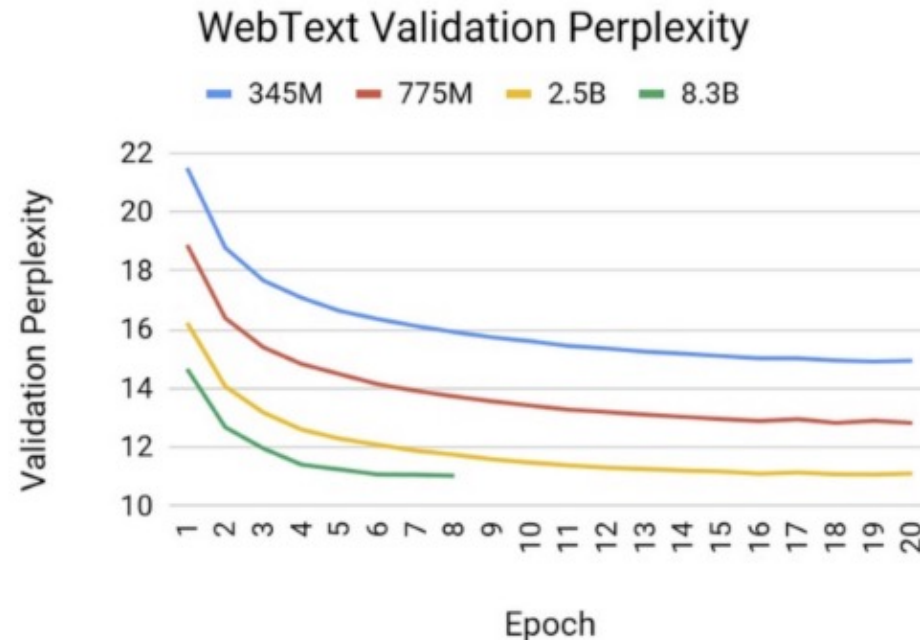# Language Modeling

- *Impossible, **but**.. Bigger models do better and better*
  - Limit not yet reached

▸ NVIDIA: trained 8.3B parameter GPT model (5.6x the size of GPT-2)

▸ Arguable these models are still underfit: larger models still get better held-out perplexities



NVIDIA blog (Narasimhan, August 2019)

# Open Questions

- Language modeling as a pre-training, and transformer-based architectures are now SoTA for all major NLP task.
  - GPT-3 uses 175 Billion parameters.

**That's a lot!**

Is this a sustainable direction? Currently all models developed in industry

Building these models in academia would require a collective effort, might require changing how research is done.

# Open Questions

- **What is really learned by a language model and how should we think about them?**
  - It's clear that our original definition was too simplistic (capture word-patterns for tasks like context sensitive spelling correction)

  - Do we still need symbolic knowledge graphs? Can everything "just" be learned from raw text using a language model?

# Evaluating Learned Representations

Q: What is human life expectancy in the United States?
**A: Human life expectancy in the United States is 78 years.**

Q: Who was president of the United States in 1955?
**A: Dwight D. Eisenhower was president of the United States in 1955.**

Q: What party did he belong to?
**A: He belonged to the Republican Party.**

## GPT 3 Generated Answers

Q: Which is heavier, a mouse or an elephant?
**A: An elephant is heavier than a mouse.**

Q: Which is heavier, a mouse or a toaster?
**A: A toaster is heavier than a mouse.**

Q: Which is heavier, a toaster or a pencil?
**A: A pencil is heavier than a toaster.**

https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html

# Evaluating Machine Comprehension

- The models we discussed are very flexible.
  - Highly expressive models, trained over massive amounts of text.
  - **What do they actually learn?**
  - Memorize facts, capture common-sense judgements, world knowledge, etc.
- Many tasks can be formulated as LM
  - Question answering, discourse relation, etc.

# GPT3: new directions

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1    Translate English to French:        ←——— task description

2    sea otter => loutre de mer          ←——— examples

3    peppermint => menthe poivrée        ←—

4    plush girafe => girafe peluche      ←—

5    cheese =>                           ←——— prompt
```

Brown et al. (2020)

# GPT3: new directions

▶ **Key observation:** few-shot learning only works with the very largest models!



Brown et al. (2020)

# Summary

- Machine learning has gone through a remarkable transformation over the last decade.
  - Learning from "clean" labels → learn from weak/indirect/self supervision
  - Generalization for one task → One shot learning relying on representation
  - Feature engineering → learned representation, at **massive** scale
  - Connect modalities: vision and language, text and speech,…
  - Perform real world tasks..

# Final exam

- Similar in scope to the mid-term.
- Bring pens + calculator.
- Cheatsheet (1 page, 2 sides).

# Final exam

- Short questions ("true/false", draw lines, etc.)
  - **Always explain your answers!**
- **Pick one**: We can use AdaBoost to help reduce bias/variance/both/neither, when comparing the final classifier to the original base classifier.
- **Pick one**: We can use Bagging to help reduce the bias/variance/both/ when comparing the final classifier to the original base classifier.
- **Which of the following..** Will provide a different vectorized representation to the two phrases "I like A but not B" and "I like B but not A" (1) Averaged word embedding (2) LSTM (3) RNN (4) Unigram BoW (5) Bigram BoW

# Final exam

- Theory Question:
  - How many examples would we need to ensure an error no greater than ε, with probability of no more than $\delta$, given that we learn functions of the class monotone conjunctions over N variables?

  - We define a rectangle using two points (x1,y1) and (x2,y2). Any points contained within the rectangle are classified as positive. What would be the number of examples now?

# Final exam

- Algorithm Question
- The AdaBoost algorithm requires base models that can be trained w.r.t a distribution of examples, updated throughout the run of the algorithm.
  - Suggest and justify your choice of a base model.
  - Explain how you would modify the algorithm to take the distribution of the data into account.

# Final Exam

- You are given the following dataset.
  - Write down the initial distribution over the examples
  - Draw the classifier that would be trained by your model
  - What is the distribution at the next step?