

CS11-747 Neural Networks for NLP

Sentence and Contextualised Word Representations

Graham Neubig



Carnegie Mellon University

Language Technologies Institute

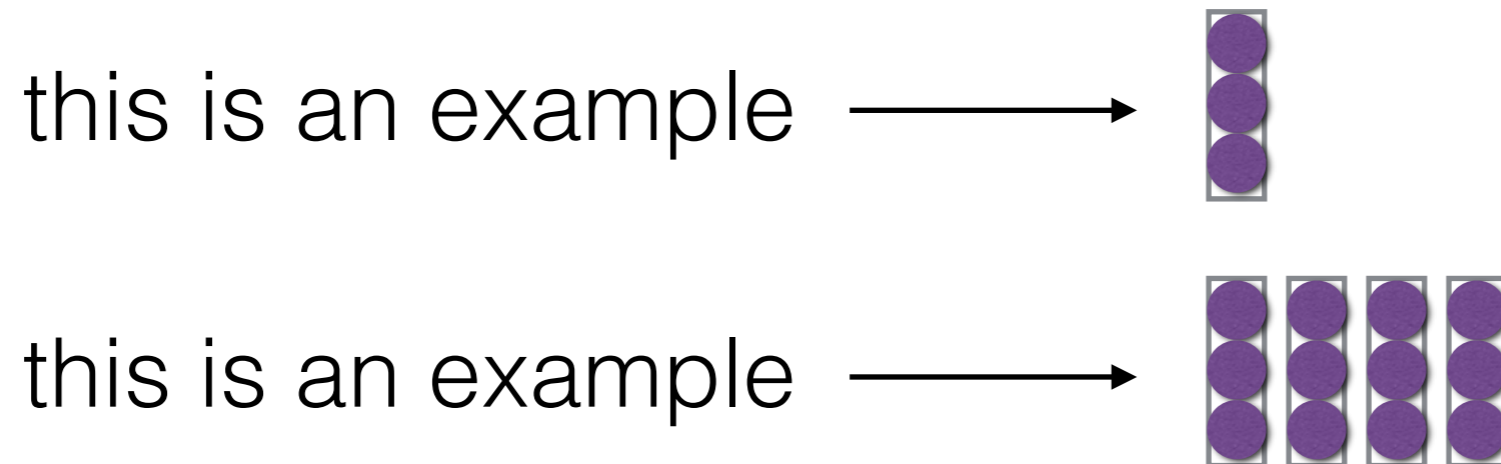
Site

<https://phontron.com/class/nn4nlp2019/>

(w/ slides by Antonis Anastasopoulos)

Sentence Representations

- We can create a vector or sequence of vectors from a sentence



Obligatory Quote!

“You can’t cram the meaning of a whole %&!\$ing sentence into a single \$&!*ing vector!”

— Ray Mooney

Goal for Today

- Briefly Introduce **tasks, datasets** and **methods**
- Introduce different **training objectives**
- Talk about **multitask/transfer learning**

Tasks Using Sentence Representations

Where would we need/use Sentence Representations?

- Sentence Classification
- Paraphrase Identification
- Semantic Similarity
- Entailment
- Retrieval

Sentence Classification

- Classify sentences according to various traits
- Topic, sentiment, subjectivity/objectivity, etc.

I hate this movie

The diagram shows the sentence "I hate this movie" on the left. An arrow points from the end of the sentence to a vertical list of sentiment labels on the right. The labels are: "very good" (green), "good" (green), "neutral" (black), "bad" (red), and "very bad" (red). The "very bad" label is highlighted with a red background, indicating it is the correct classification for the sentence.

very good
good
neutral
bad
very bad

I love this movie

The diagram shows the sentence "I love this movie" on the left. An arrow points from the end of the sentence to a vertical list of sentiment labels on the right. The labels are: "very good" (green), "good" (green), "neutral" (black), "bad" (red), and "very bad" (red). The "very good" label is highlighted with a green background, indicating it is the correct classification for the sentence.

very good
good
neutral
bad
very bad

Paraphrase Identification

(Dolan and Brockett 2005)

- Identify whether A and B mean the same thing

Charles O. Prince, 53, was named as Mr. Weill's successor.



Mr. Weill's longtime confidant, Charles O. Prince, 53, was named as his successor.

- **Note:** *exactly* the same thing is too restrictive, so use a loose sense of similarity

Semantic Similarity/Relatedness

(Marelli et al. 2014)

- Do two sentences mean something similar?

Relatedness score	Example
1.6	A: <i>“A man is jumping into an empty pool”</i> B: <i>“There is no biker jumping in the air”</i>
2.9	A: <i>“Two children are lying in the snow and are making snow angels”</i> B: <i>“Two angels are making snow on the lying children”</i>
3.6	A: <i>“The young boys are playing outdoors and the man is smiling nearby”</i> B: <i>“There is no boy playing outdoors and there is no man smiling”</i>
4.9	A: <i>“A person in a black jacket is doing tricks on a motorbike”</i> B: <i>“A man in a black jacket is doing tricks on a motorbike”</i>

- Like paraphrase identification, but with shades of gray.

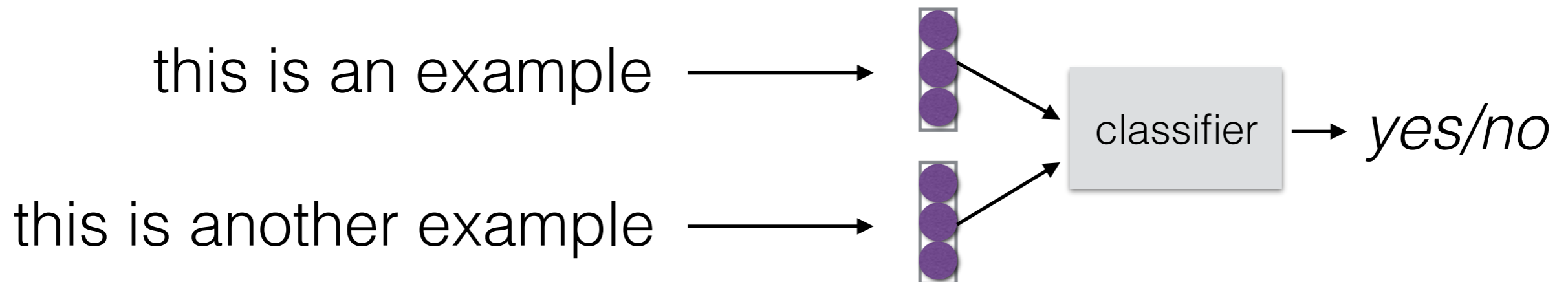
Textual Entailment

(Dagan et al. 2006, Marelli et al. 2014)

- **Entailment:** if A is true, then B is true (c.f. paraphrase, where opposite is also true)
 - The woman bought a sandwich for lunch
→ The woman bought lunch
- **Contradiction:** if A is true, then B is not true
 - The woman bought a sandwich for lunch
→ The woman did not buy a sandwich
- **Neutral:** cannot say either of the above
 - The woman bought a sandwich for lunch
→ The woman bought a sandwich for dinner

Model for Sentence Pair Processing

- Calculate vector representation
- Feed vector representation into classifier



How do we get such a representation?

Multi-task Learning Overview

Types of Learning

- **Multi-task learning** is a general term for training on multiple tasks
- **Transfer learning** is a type of multi-task learning where we only really care about one of the tasks
- **Domain adaptation** is a type of transfer learning, where the output is the same, but we want to handle different topics or genres, etc.

Plethora of Tasks in NLP

- In NLP, there are a plethora of tasks, each requiring different varieties of data
 - **Only text:** e.g. language modeling
 - **Naturally occurring data:** e.g. machine translation
 - **Hand-labeled data:** e.g. most analysis tasks
- And each in many languages, many domains!

Rule of Thumb 1: Multitask to Increase Data

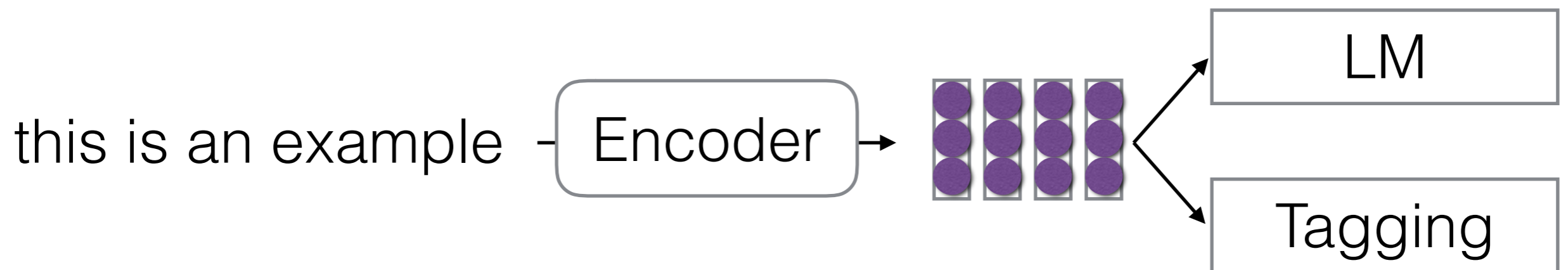
- Perform multi-tasking when one of your two tasks has many fewer data
- **General domain → specific domain**
(e.g. web text → medical text)
- **High-resourced language → low-resourced language**
(e.g. English → Telugu)
- **Plain text → labeled text**
(e.g. LM → parser)

Rule of Thumb 2:

- Perform multi-tasking when your **tasks are related**
- e.g. predicting eye gaze and summarization (Klerke et al. 2016)

Standard Multi-task Learning

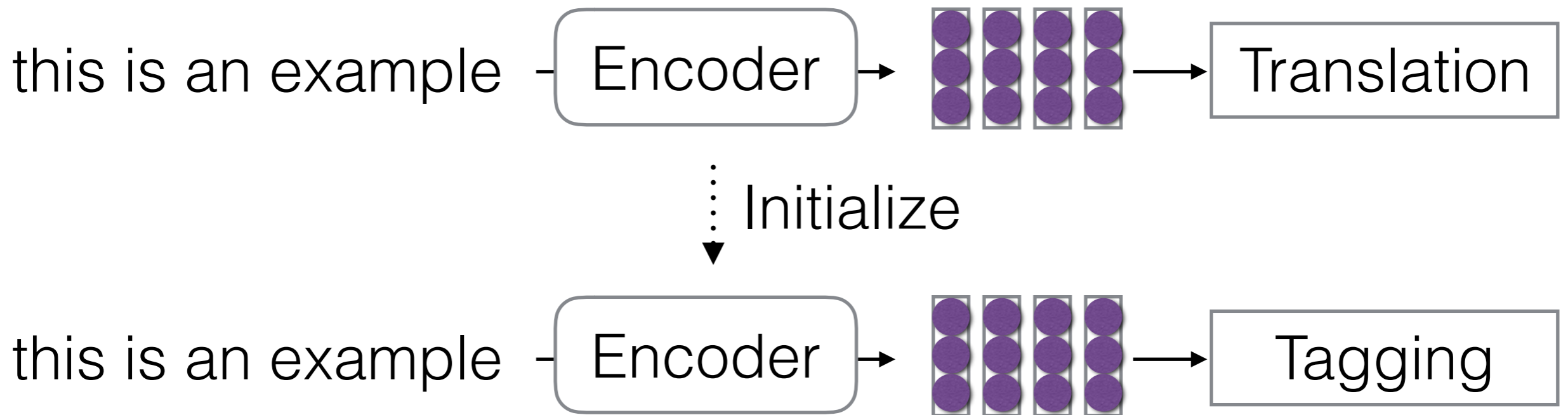
- Train representations to do well on multiple tasks at once



- In general, as simple as randomly choosing minibatch from one of multiple tasks
- Many many examples, starting with Collobert and Weston (2011)

Pre-training

- First train on one task, then train on another



- Widely used in word embeddings (Turian et al. 2010)
- Also pre-training sentence encoders or contextualized word representations (Dai et al. 2015, Melamud et al. 2016)

Thinking about Multi-tasking, and Pre-trained Representations

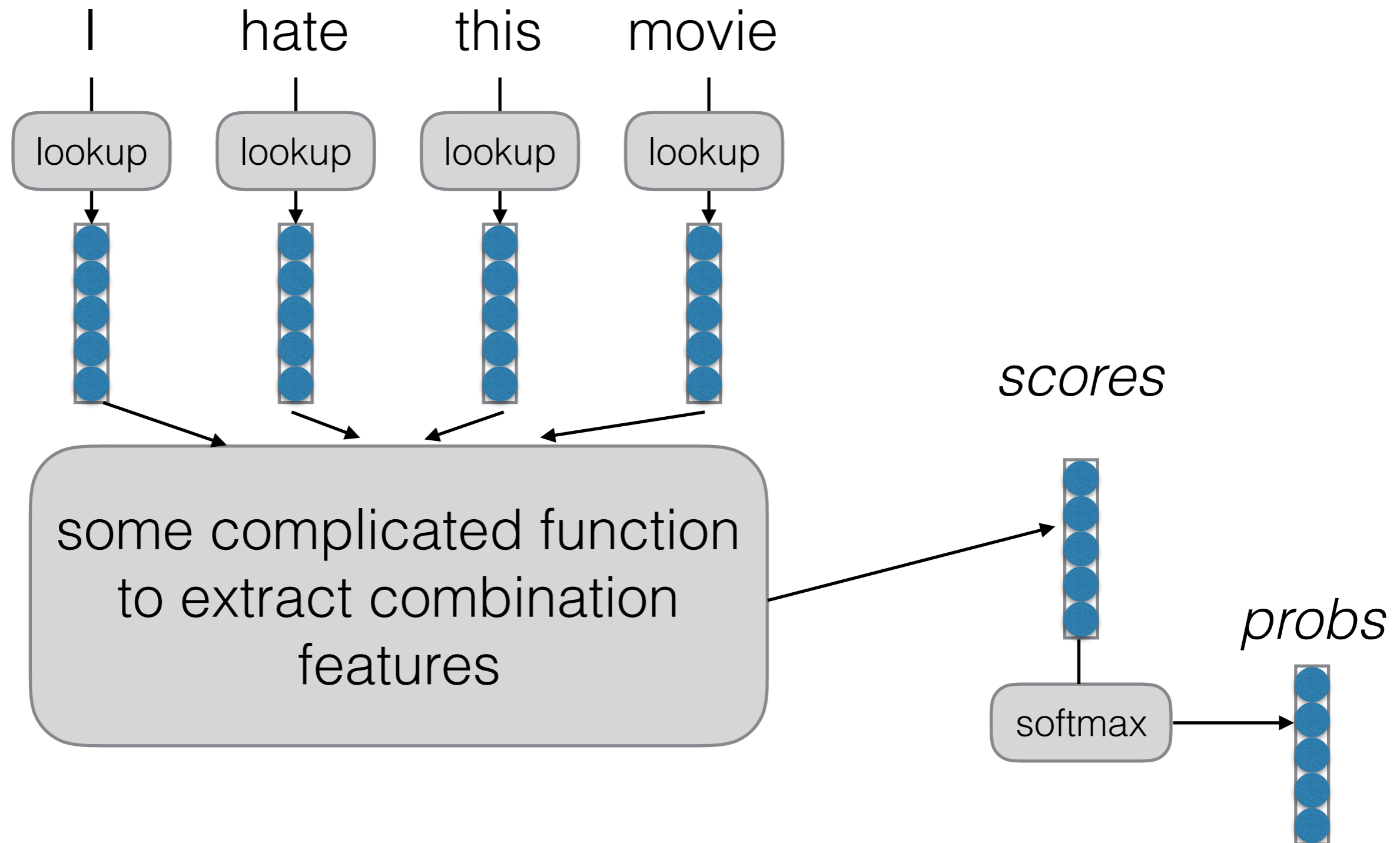
- Many methods have names like SkipThought, ParaNMT, CoVe, ELMo, BERT along with pre-trained models
- These often refer to a combination of
 - **Model:** The underlying neural network architecture
 - **Training Objective:** What objective is used to pre-train
 - **Data:** What data the authors chose to use to train the model
- Remember that these are often conflated (and don't need to be)!

End-to-end vs. Pre-training

- For any model, we can always use an end-to-end training objective
 - **Problem:** paucity of training data
 - **Problem:** weak feedback from end of sentence only for text classification, etc.
- Often better to pre-train sentence embeddings on other task, then use or fine tune on target task

Training Sentence Representations

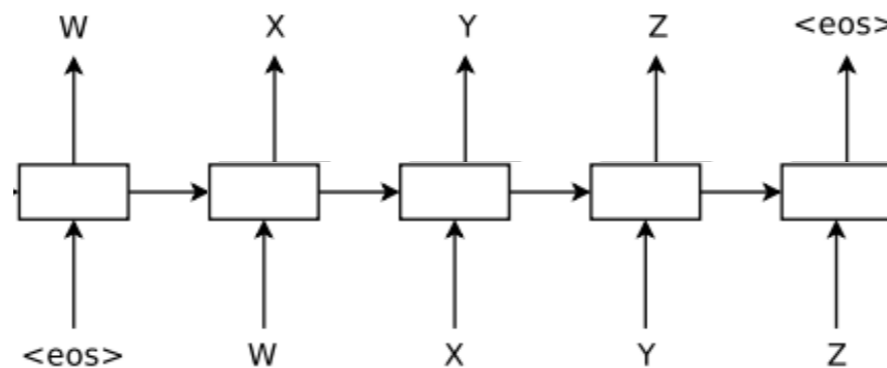
General Model Overview



Language Model Transfer

(Dai and Le 2015)

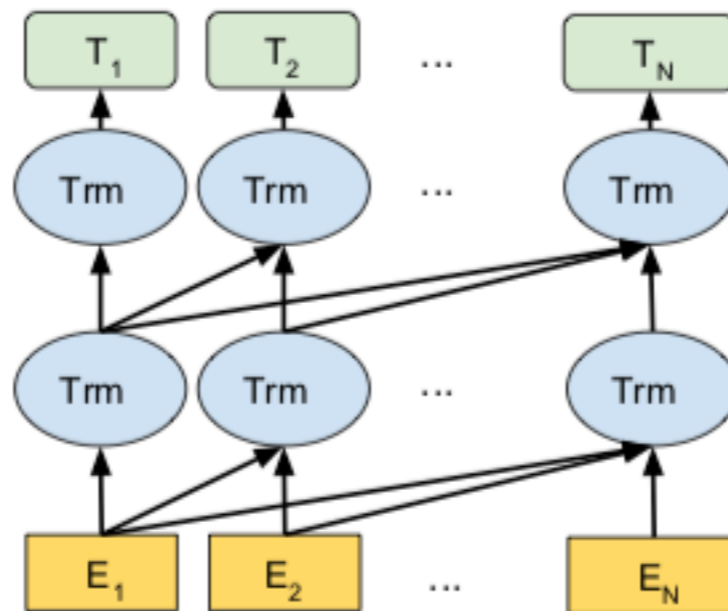
- **Model:** LSTM
- **Objective:** Language modeling objective
- **Data:** Classification data itself, or Amazon reviews



- **Downstream:** On text classification, initialize weights and continue training

Unidirectional Training + Transformer (OpenAI GPT) (Radford et al. 2018)

- **Model:** Masked self-attention
- **Objective:** Predict the next word left->right
- **Data:** BooksCorpus



Downstream: Some task fine-tuning, other tasks additional multi-sentence training

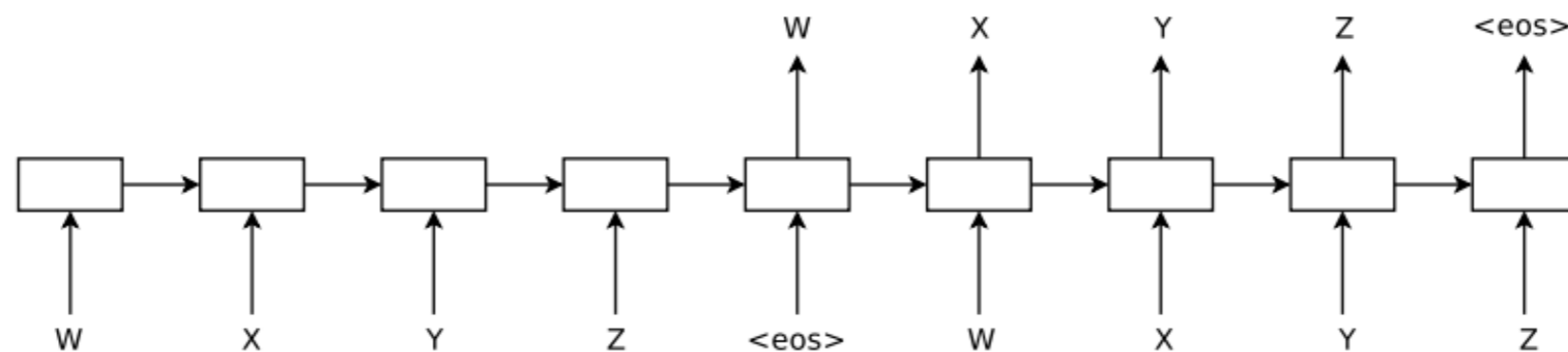
Auto-encoder Transfer

(Dai and Le 2015)

- **Model:** LSTM

- **Objective:** From single sentence vector, re-construct the sentence

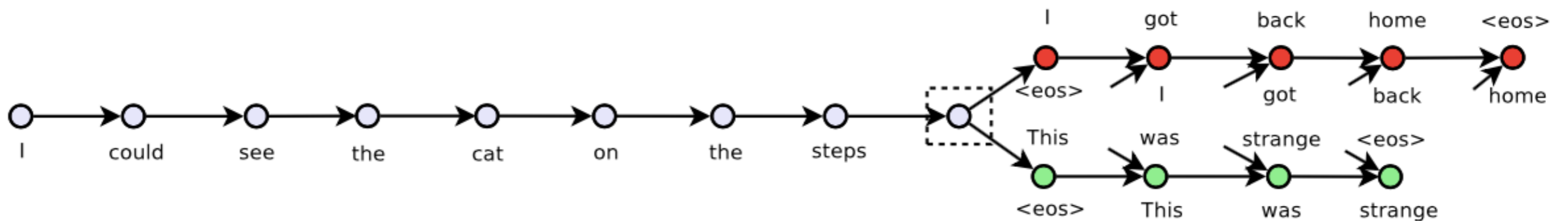
- **Data:** Classification data itself, or Amazon reviews



- **Downstream:** On text classification, initialize weights and continue training

Context Prediction Transfer (Skip-thought Vectors) (Kiros et al. 2015)

- **Model:** LSTM
- **Objective:** Predict the surrounding sentences
- **Data:** Books, important because of context



- **Downstream Usage:** Train logistic regression on $[|u-v|; u*v]$ (component-wise)

Paraphrase ID Transfer (Wieting et al. 2015)

- **Model:** Try many different ones
- **Objective:** Predict whether two phrases are paraphrases or not from
- **Data:** Paraphrase database (<http://paraphrase.org>), created from bilingual data
- **Downstream Usage:** Sentence similarity, classification, etc.
- **Result:** Interestingly, LSTMs work well on in-domain data, but word averaging generalizes better

Large Scale Paraphrase Data (ParaNMT-50MT) (Wieting and Gimpel 2018)

- **Automatic construction of large paraphrase DB**
 - Get large parallel corpus (English-Czech)
 - Translate the Czech side using a SOTA NMT system
 - Get automated score and annotate a sample
- Corpus is **huge but includes noise**, 50M sentences (about 30M are high quality)
- Trained representations work quite well and generalize

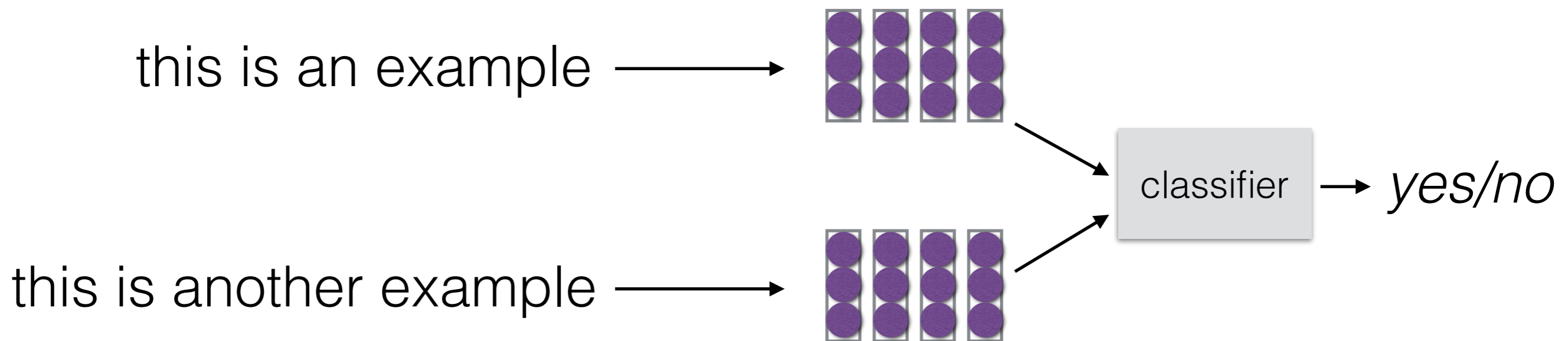
Entailment Transfer (InferSent) (Conneau et al. 2017)

- Previous objectives use no human labels, but what if:
- **Objective:** supervised training for a task such as entailment learn generalizable embeddings?
 - Task is more difficult and requires capturing nuance → yes?, or data is much smaller → no?
- **Model:** Bi-LSTM + max pooling
- **Data:** Stanford NLI, MultiNLI
- **Results:** Tends to be better than unsupervised objectives such as SkipThought

Contextualized Word Representations

Contextualized Word Representations

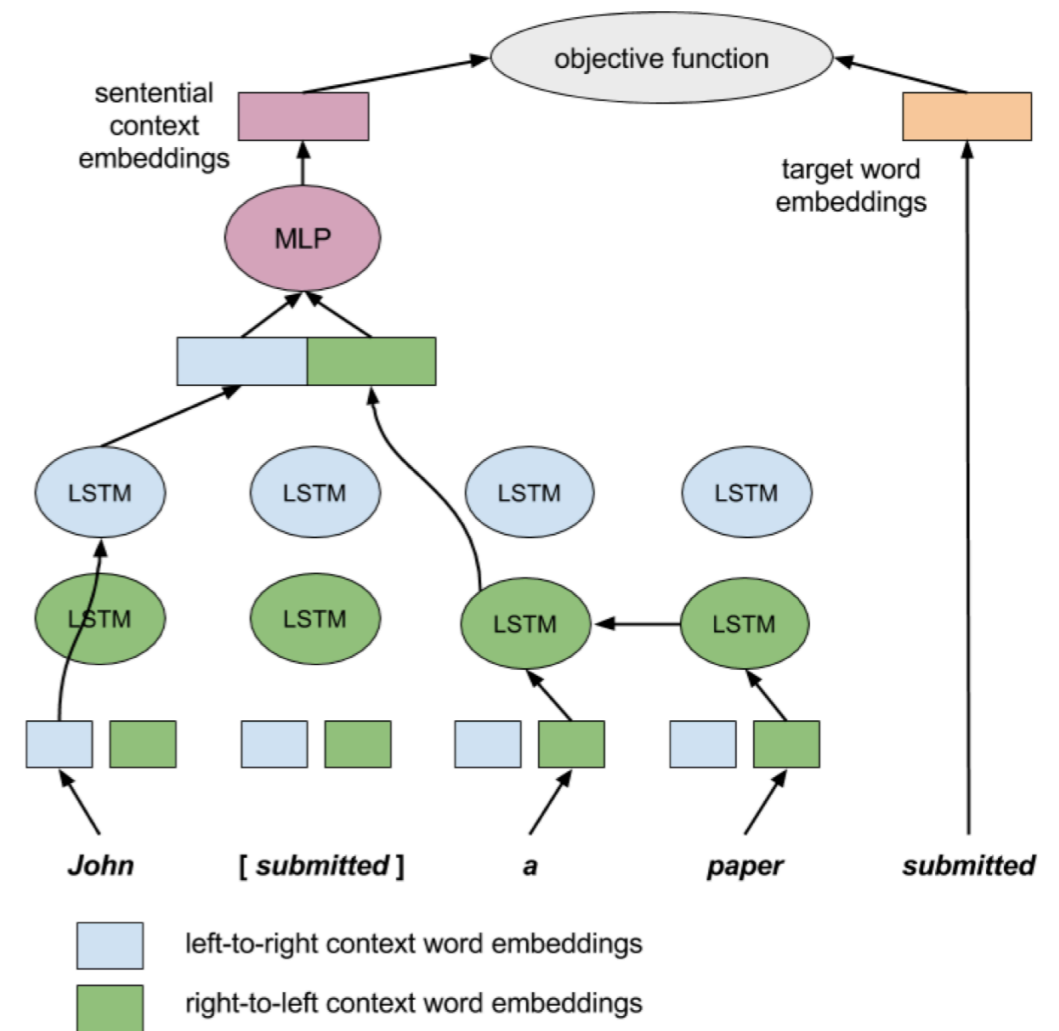
- Instead of one vector per sentence, one vector per word!



How to train this representation?

Central Word Prediction Objective (context2vec) (Melamud et al. 2016)

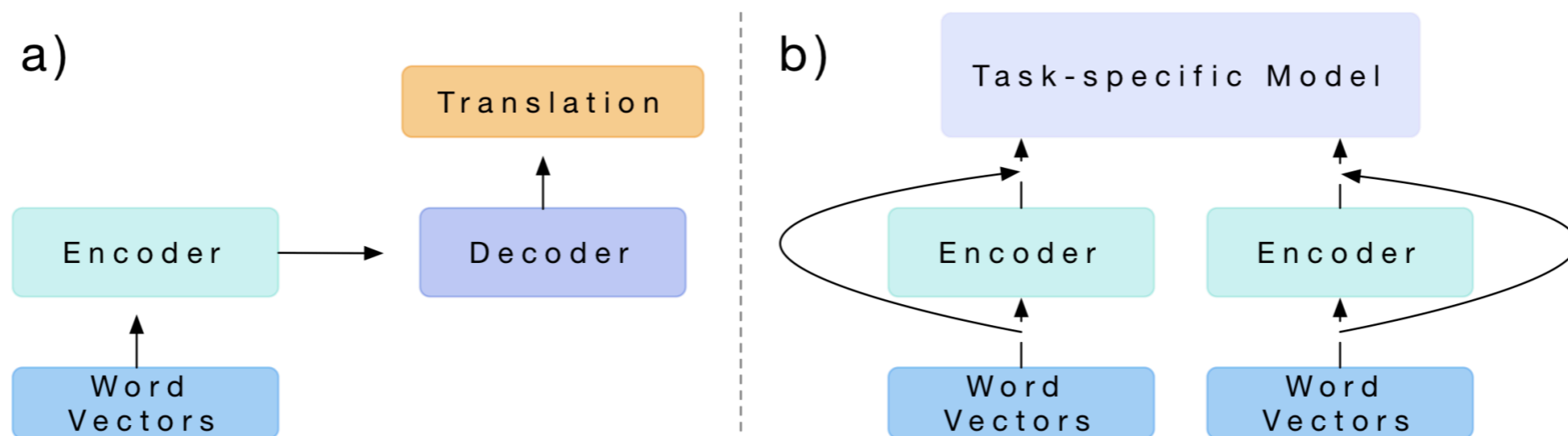
- **Model:** Bi-directional LSTM
- **Objective:** Predict the word given context
- **Data:** 2B word ukWaC corpus
- **Downstream:** use vectors for sentence completion, word sense disambiguation, etc.



Machine Translation Objective (CoVe)

(McMann et al. 2017)

- **Model:** Multi-layer bi-directional LSTM
- **Objective:** Train attentional encoder-decoder
- **Data:** 7M English-German sentence pairs

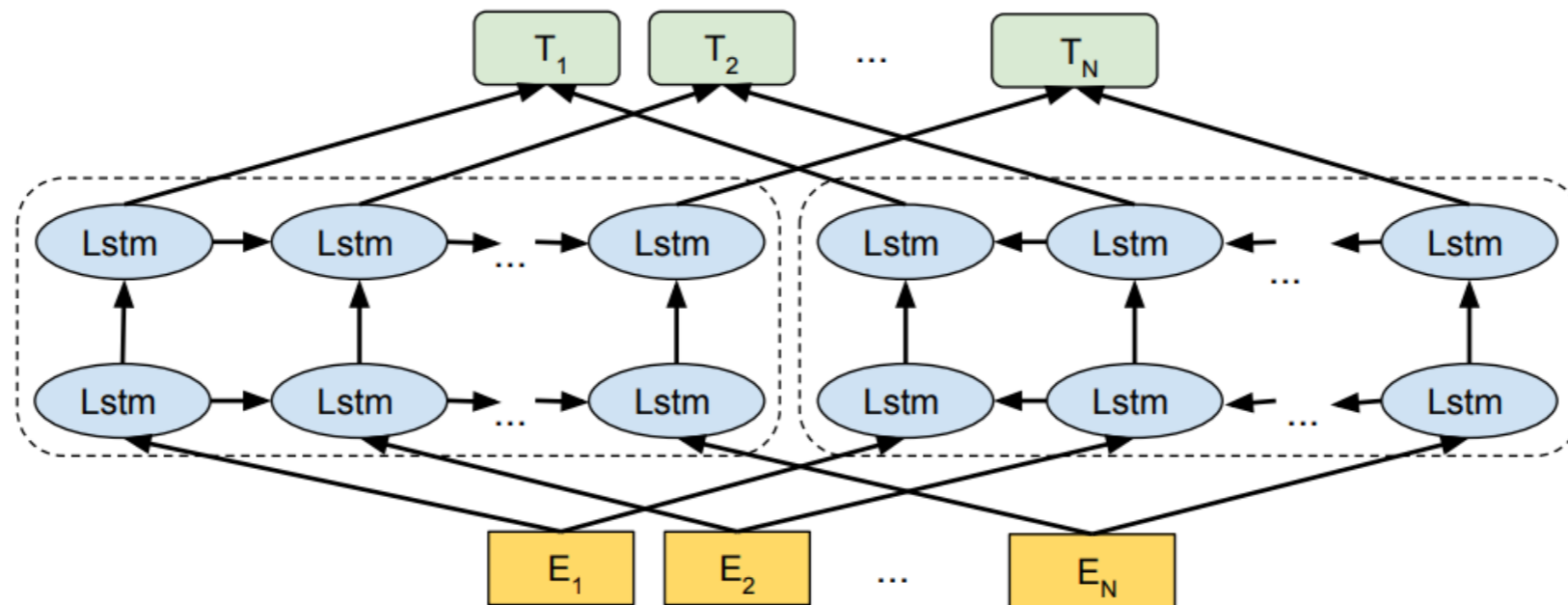


Downstream: Use bi-attention network over sentence pairs for classification

Bi-directional Language Modeling Objective (ELMo)

(Peters et al. 2018)

- **Model:** Multi-layer bi-directional LSTM
- **Objective:** Predict the next word left->right, next word right->left independently
- **Data:** 1B word benchmark LM dataset

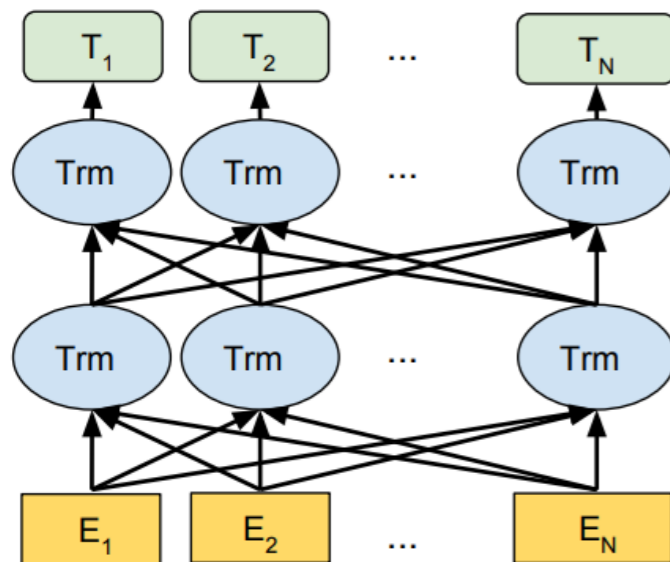


Downstream: Finetune the weights of the linear combination of layers on the downstream task

Masked Word Prediction (BERT)

(Devlin et al. 2018)

- **Model:** Multi-layer self-attention. Input sentence or pair, w/ [CLS] token, subword representation



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{##ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

- **Objective:** Masked word prediction + next-sentence prediction
- **Data:** BooksCorpus + English Wikipedia

Masked Word Prediction

(Devlin et al. 2018)

1. predict a masked word
 - 80%: substitute input word with [MASK]
 - 10%: substitute input word with random word
 - 10%: no change
- Like context2vec, but **better suited for multi-layer self attention**

Consecutive Sentence Prediction

(Devlin et al. 2018)

1. classify two sentences as consecutive or not:
 - 50% of training data (from OpenBooks) is "consecutive"

Input = [CLS] the man [MASK] to the store [SEP]
penguin [MASK] are flight ##less birds [SEP]

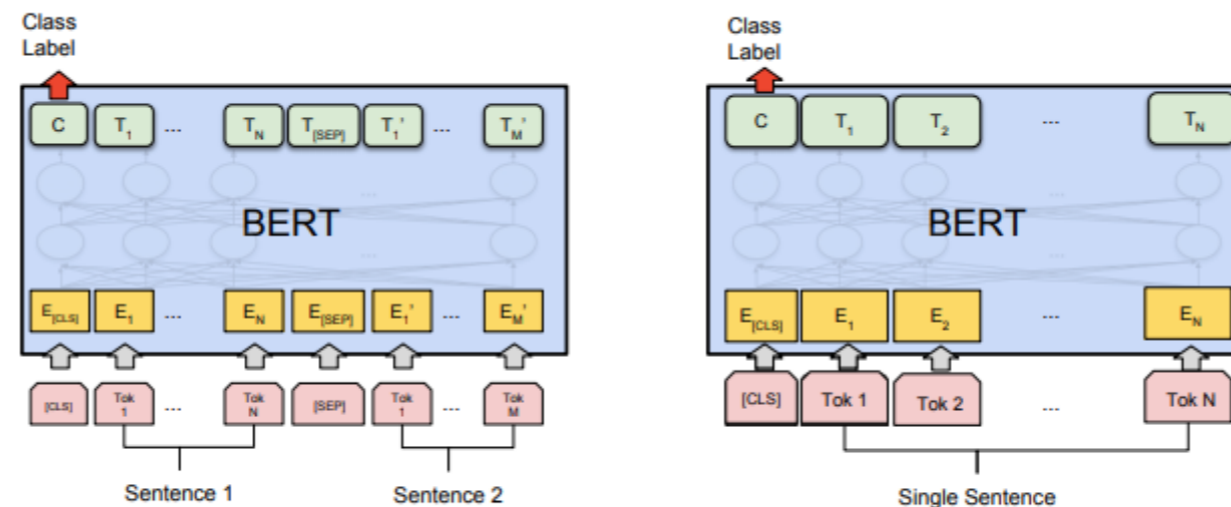
Label = NotNext

Input = [CLS] the man went to [MASK] store [SEP]
he bought a gallon [MASK] milk [SEP]

Label = IsNext

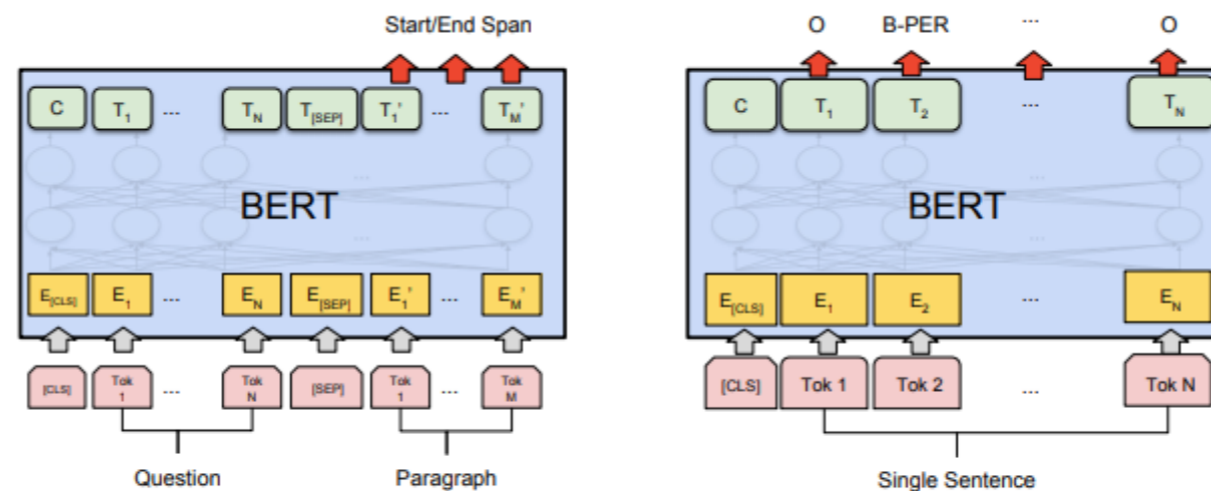
Using BERT with pre-training/finetuning

- Use the pre-trained model as the first “layer” of the final model, then train on the desired task



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA



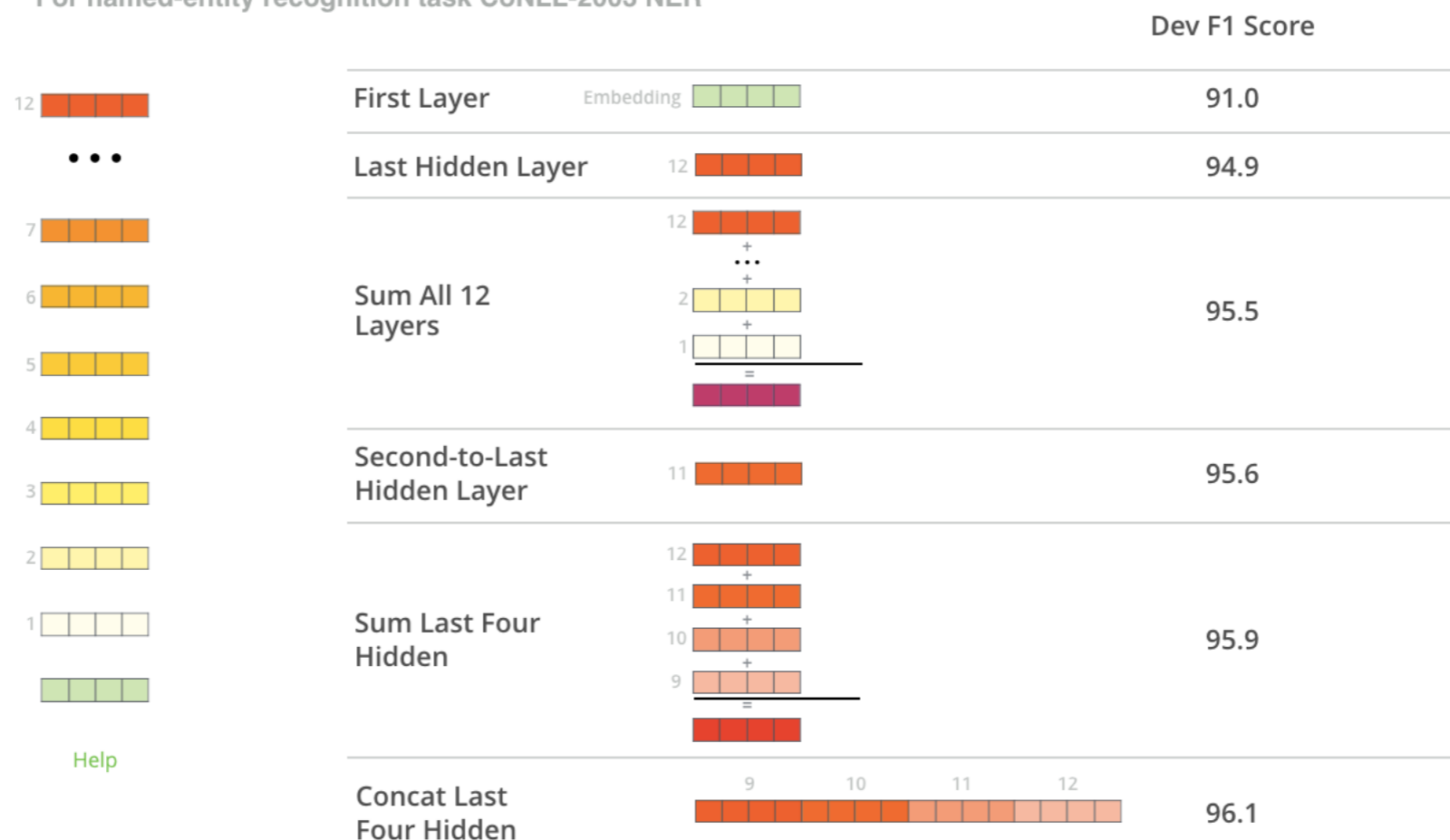
(c) Question Answering Tasks:
SQuAD v1.1

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

Using BERT for Representations

- Use the pre-trained model to obtain contextualised word representations for the input

What is the best contextualized embedding for “Help” in that context?
For named-entity recognition task CoNLL-2003 NER



[visualization from The Illustrated BERT: <https://jalammar.github.io/illustrated-bert/>]

Which Method is Better?

Which Model?

- Not very extensive comparison...
- Wieting et al. (2015) find that simple word averaging is more robust out-of-domain
- Devlin et al. (2018) compare unidirectional and bi-directional transformer, but no comparison to LSTM like ELMo (for performance reasons?)

Which Training Objective?

- Not very extensive comparison...
- Zhang and Bowman (2018) control for training data, and find that bi-directional LM seems better than MT encoder
- Devlin et al. (2018) find next-sentence prediction objective good compliment to LM objective

Which Data?

- Not very extensive comparison...
- Zhang and Bowman (2018) find that more data is probably better, but results preliminary.
- Data with context is probably essential.

Questions?