# Representations for a word

我们已经学到

- word2vec
- GloVe
- FastText

在2014年，普遍人为pre-training比random的wordvec对下游任务更有效。

**目前wordvec存在的问题**

- 同一个word type("人")向量相同，没有考虑word token("人次")的context
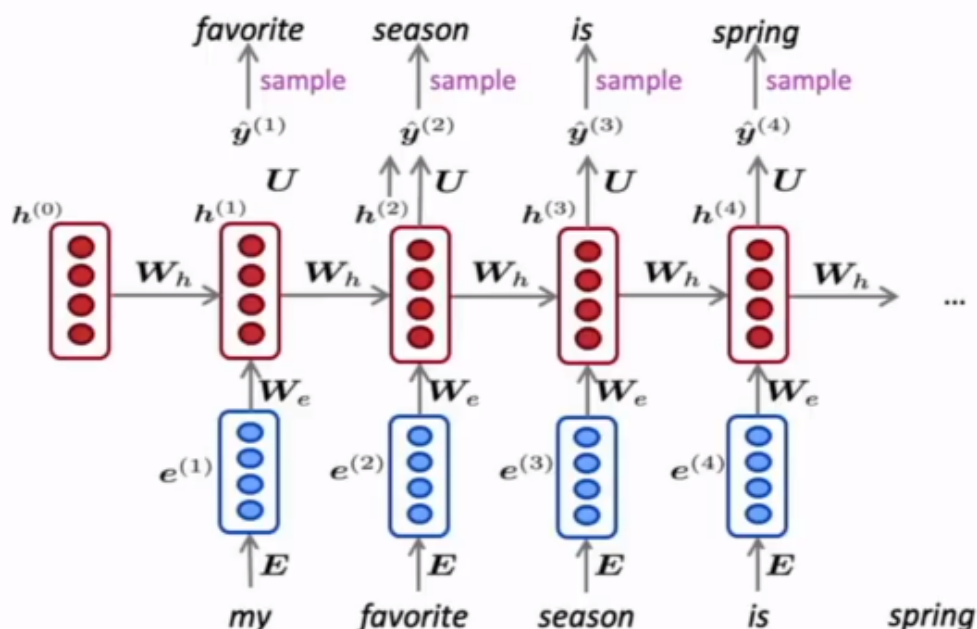- 同义词/派生词也有不同的适用场景，要根据context加以区分

## Tips for unknown words with word vectors

1. char-level
2. Dhingra <mark>为什么有unsupervised wordvec?#因为word2vec可以无监督学习corpus中的word embedding，但最后下游任务的vocab不包含corpus中所有单词，出现unk。</mark>

- 2. Try these tips (from Dhingra, Liu, Salakhutdinov, Cohen 2017)
  - a. If the <UNK> word at test time appears in your unsupervised word embeddings, use that vector as is at test time.
  - b. Additionally, for other words, just assign them a random vector, adding them to your vocabulary
- a. definitely helps a lot; b. may help a little more

## Did we solve the problem?

在lstm中，每层输出h作为预测单词的contextual word embedding。

**核心思想**

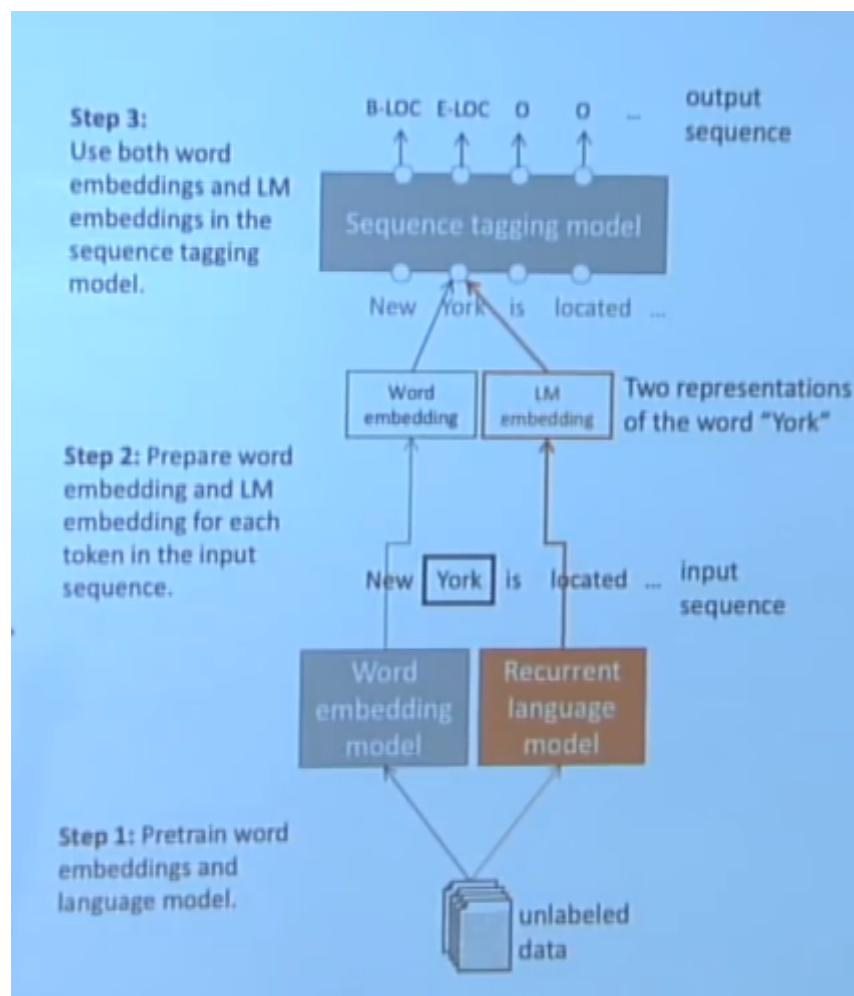可以利用学到的contextual word embedding！

# Contextual word embeddings

## TagLM

用word embedding+LM embedding

https://arxiv.org/pdf/1705.00108.pdf
* Idea: Want meaning of word in context, but standardly learn task RNN only on small task-labeled data (e.g., NER)
* Why don't we do semi-supervised approach where we train NLM on large unlabeled corpus, rather than just word vectors?

如何无监督地pretrain wordvec和LM？#word2vec就是无监督的。LM可以通过"预测下一个单词"任务来完成，这是无监督的。通常pretrain的参数不会在train时改变。

# ELMo (Embedding from Language Models)

ELMo是对TagLM的改进，包括

- 两 层bi-LSTM
- combine每层的特征而不是最后一层

- First run biLM to get representations for each word
- Then let (whatever) end-task model use them
  - Freeze weights of ELMo for purposes of supervised model
  - Concatenate ELMo weights into task-specific model
    - Details depend on task
      - Concatenating into intermediate layer as for TagLM is typical
      - Can provide ELMo representations again when producing outputs, as in a question answering system

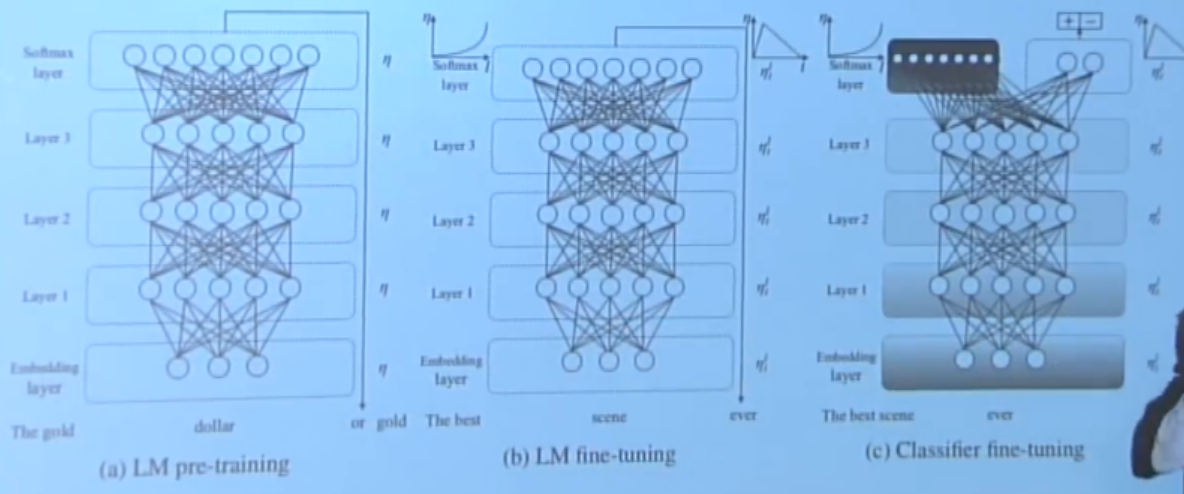ELMo的最大贡献是，这种word embedding可以被用在不同下游任务上。

# ULMfit (Universal Language Model Fine-tuning for Text Classification)

用transfer learning的思想，把一个任务/数据集的embedding迁移到另一个任务/数据集。

流程:

1. 在大的corpus上训练embedding
2. 在小的vocab上fine tune
3. 训练下游分类器

(a) LM pre-training  (b) LM fine-tuning  (c) Classifier fine-tuning

ULMfit提供了一个有效思路: 在大量数据上预训练，在少量数据上fine-tune。

由此出现了大公司的算力竞争，在巨大数据量的corpus上预训练LM：



Let's scale it up!

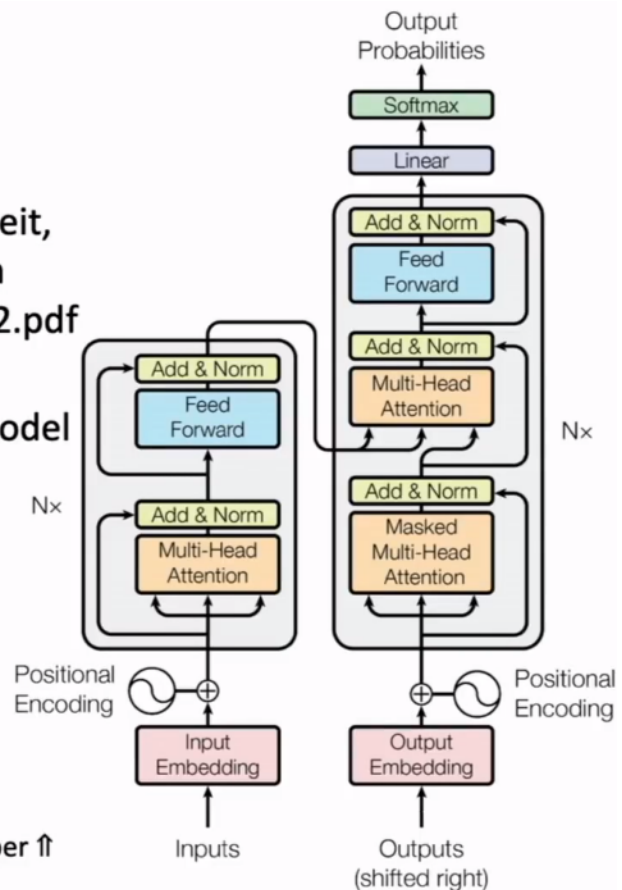| ULMfit | GPT | BERT | GPT-2 |
| Jan 2018 | June 2018 | Oct 2018 | Feb 2019 |
| Training: | Training | Training | Training |
| 1 GPU day | 240 GPU days | 256 TPU days | ~2048 TPU v3 days according to a reddit thread |
| | | ~320–560 GPU days | |

fast.ai  OpenAI  Google AI  OpenAI

# Transformers

Attention可以表示各个timestep之间的关系，可以只用attention。

# Transformer Overview

Attention is all you need. 2017.
Aswani, Shazeer, Parmar, Uszkoreit,
Jones, Gomez, Kaiser, Polosukhin
https://arxiv.org/pdf/1706.03762.pdf

- Non-recurrent sequence-to-sequence encoder-decoder model
- Task: machine translation with parallel corpus
- Predict each translated word
- Final cost/error function is standard cross-entropy error on top of a softmax classifier

This and related figures from paper ⇑

38

attention的q, K, V的解释。

# Dot-Product Attention (Extending our previous def.)

- Inputs: a query q and a set of key-value (k-v) pairs to an output
- Query, keys, values, and output are all vectors

- Output is weighted sum of values, where
- Weight of each value is computed by an inner product of query and corresponding key
- Queries and keys have same dimensionality $d_k$ value have $d_v$

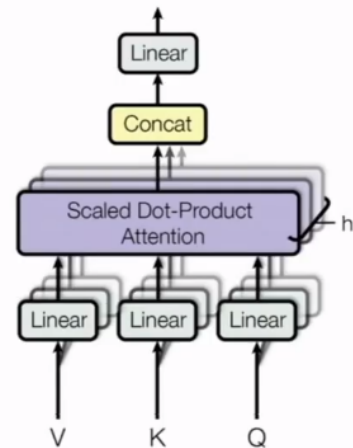$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

multi-head attention希望$Q_i, K_i, V_i$的不同i注意不同的东西。
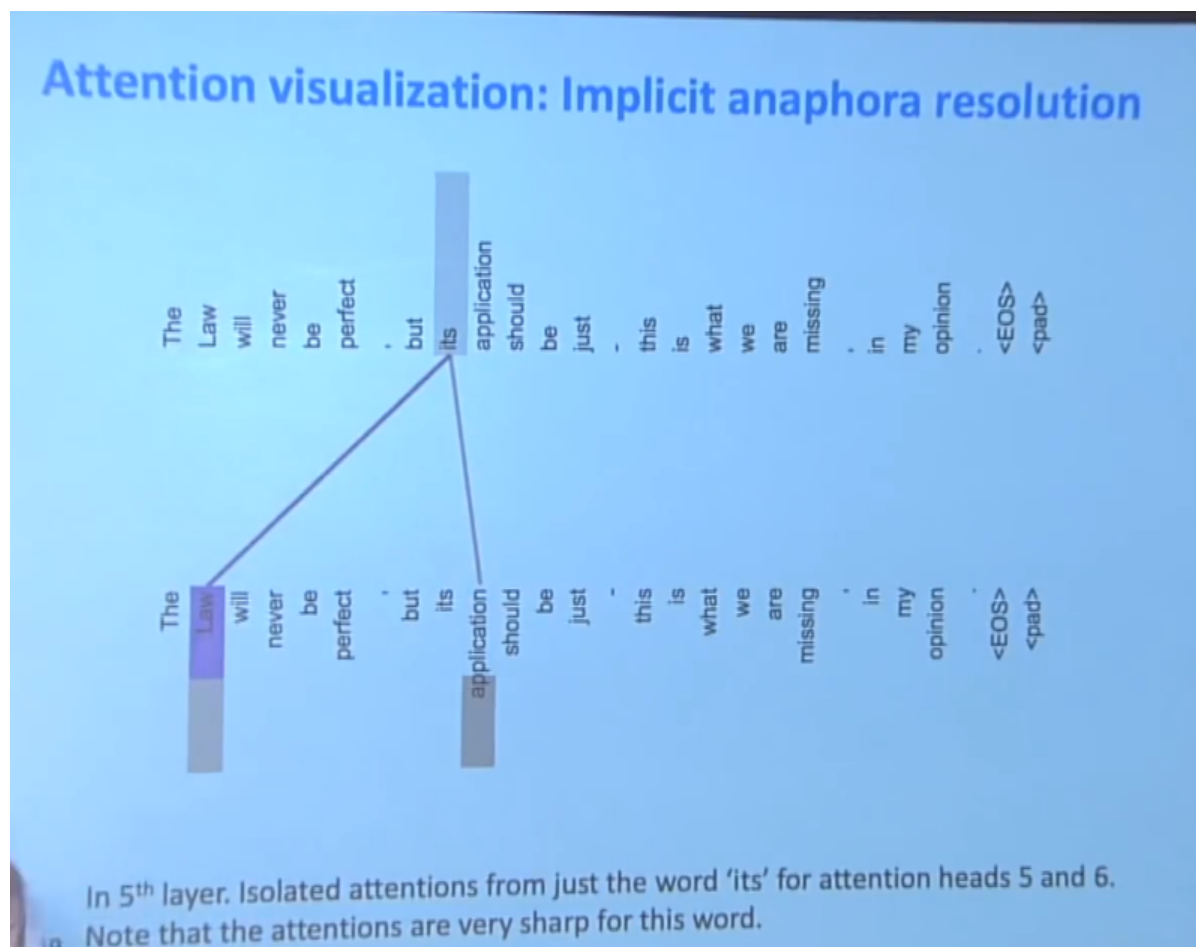
# Multi-head attention

- Problem with simple self-attention:
- Only one way for words to interact with one-another
- Solution: Multi-head attention
- First map Q, K, V into h=8 many lower dimensional spaces via W matrices
- Then apply attention, then concatenate outputs and pipe through linear layer

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



attention可视化



In 5th layer. Isolated attentions from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.

# BERT (Bidirectional Encoder Representations from Transformers)

Language Model也可以是双向的。<mark>seq2seq不是双向的吗?</mark>

单向是因为:

- LM会预测下一个word
- 看到两边的context是一 种作弊

**解决方法**



- **Solution**: Mask out *k*% of the input words, and then predict the masked words
  - They always use *k* = 15%

  store        gallon
   ↑            ↑

  the man went to the [MASK] to buy a [MASK] of milk

- Too little masking: Too expensive to train
- Too much masking: Not enough context