

# Word Embedding and Question Answering

---

Cecilia Xifei Ni z5173159

December 12, 2019

## INTRODUCTION

"Who is the first Prime Minister of Australia?"

Figure 1 is a screen shot of search result returning by Google. This is a question a Question Answering (QA) system should be able to respond to. QA system evaluates texts across the web or database to find answer of a particular question to return in a form of short text.

Current machine algorithms extracts answer from a short paragraph instead of a long content (e.g. an entire wiki page/a news article). When applying modern machine learning algorithms on a long content, the result can be complicated and lengthy.

This essay aims to firstly gives a mathematical heavy explanation on nowadays natural language processing (NLP) algorithms on word embedding, then propose my own heuristics based on them to predict answers from a long content based on the question being asked.

My algorithm will be tested with Google Natural Questions which contains its own private testing dataset.

## GENERAL WORD EMBEDDING MODELS

**Word Embedding** is the collective name for techniques in NLP where words or phrases are mapped into vectors of real numbers.

---

Prime Minister of Australia (1)

---

## Edmund Barton



Australia's first prime minister, **Edmund Barton** at the central table in the House of Representatives in 1901.

[Prime Minister of Australia - Wikipedia](https://en.wikipedia.org/wiki/Prime_Minister_of_Australia)

[https://en.wikipedia.org › wiki › Prime\\_Minister\\_of\\_Australia](https://en.wikipedia.org/wiki/Prime_Minister_of_Australia)

Figure 1: A result returned by Google to answer "Who is the first Prime Minister of Australia?"

As per all machine learning algorithms, the general philosophy is to find the minimum of a converged function. When trying to convert words to number, its original utf code will not work well. What we need is a hard-coded map between the word itself, and its semantics. As human, we learn these semantics from daily experience which machine has no way to access, therefore, we need a mapping from words and its semantic in order for machine to decode.

### FIRST ATTEMPT: HOT ONE ENCODING

Given a collection of  $N$  unique words, each word is of size  $N$ . The vector is very sparse such that only one index is 1 and others are 0. This is to say each word is a **dimension**.

Figure 2 shows an example of hot encoding.

However, this model has some significant drawbacks. Firstly, it probably works fine for English as it has roughly 8000 frequent words. However, for languages like Chinese or Japanese, they are character-based, and each character will generate tons of combinations. This model will suffer from dramatic increase in dimensionalities. Secondly, all vectors in the matrix are independent to each other, while we do want them to learn the inter-connections between them.

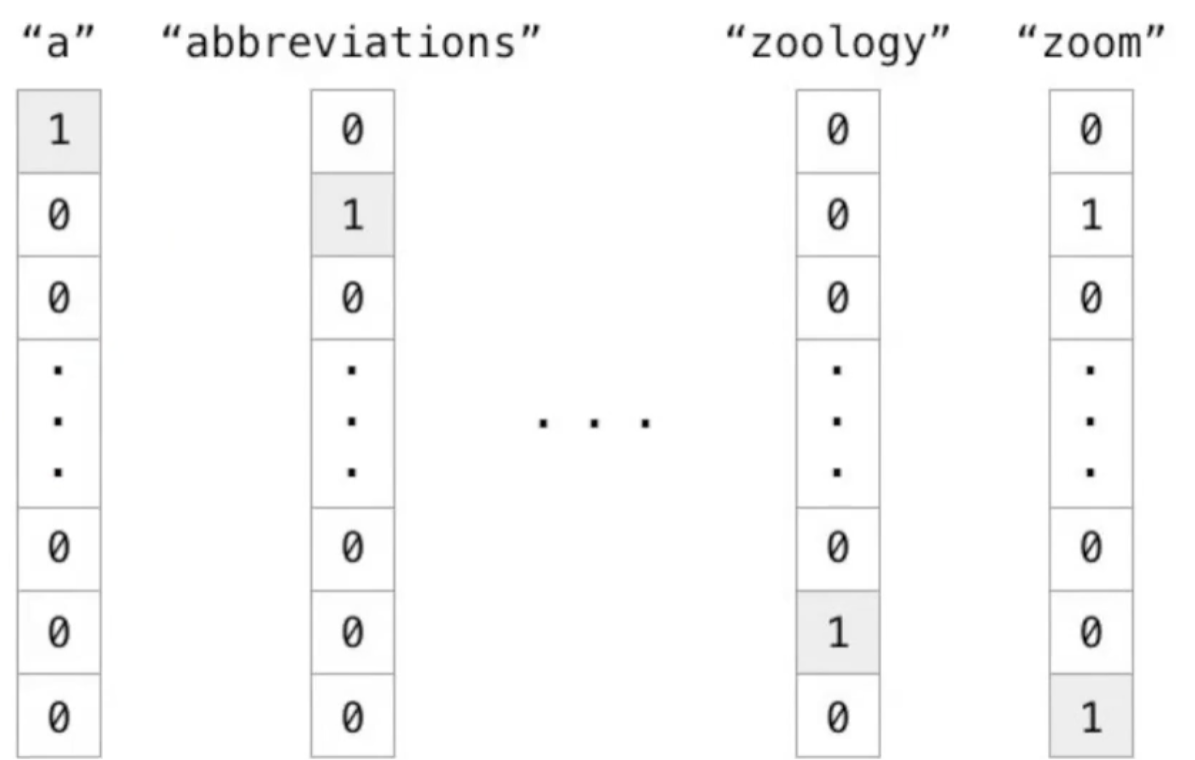


Figure 2: One Hot Encoding Example.

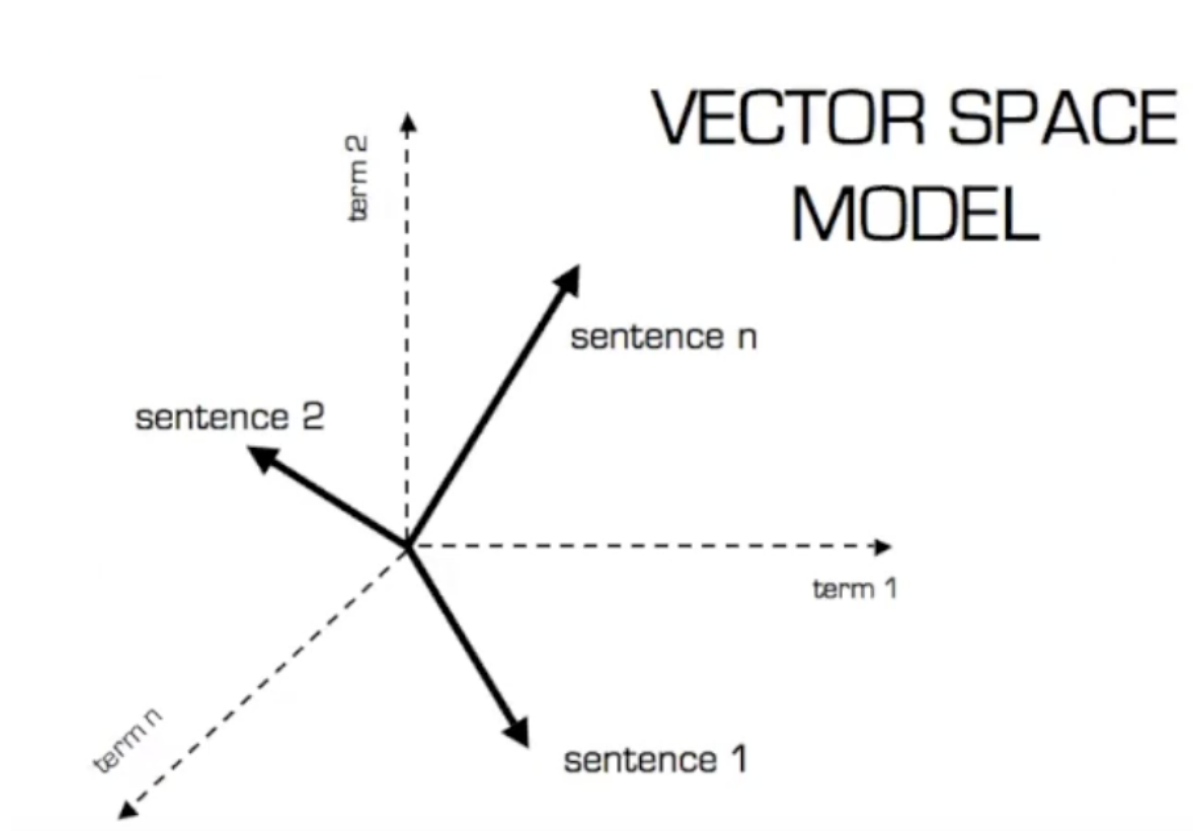


Figure 3: Bag of Words Vector Space Representation

## SECOND ATTEMPT: BAG OF WORDS (BOW) MODEL

On the second attempt, we focus more on extracting information from the perspective of the whole document. Under this model, we represent words in a document as a bag (multiset) of words– we discard order and grammar, and only keep its multiplicity.

Here is an example:

Document 1: "George" "likes" "to" "play" "video" "games", "Mary" "likes" "video" "games" "too".

Document 2: "Mary" "also" "likes" "movies".

The BoW representation of the two above documentations is

BoW1: {"George" : 1, "likes" : 2, "to" : 1, "play" : 1, "video" : 2, "games" : 2, "Mary" : 1, "too": 1}

BoW2: {"Mary" : 1, "also" : 1, "likes": 1, "movie" : 1}

Here, each key is the word, and the value corresponds to occurrence.

This approach could be potentially used to compare similarities between two documents. When comparing them, we construct two vectors with each index corresponding to each unique words. Each be the count of the word appeared in that specific document. The

	0	1
gender	male	female
age	child	adult

Table 1: Table 1: Represent Words in Fewer Dimensions

similarities of two documents can be measured by its Euclidean distance or its cosine similarities.

Intuitively, two documents are similar if their vector points to the similar direction. We define the cosine similarities of two documents as:

$$\text{cosim}(U_k, U_i) = \frac{\langle U_k, U_i \rangle}{\|U_k\| \cdot \|U_i\|} \quad (1)$$

where  $\langle U_k, U_i \rangle$  is the scalar product of  $U_i$  and  $U_k$ .

### THIRD ATTEMPT: NEURAL NETWORK BASED (NNLM)

One intuitive question to ask is, Bag of Words and One Hot Encoding both treat each words as a separate dimension, but do we need so many dimensions?

The answer is probably no. Because a lot of words are similar or related. e.g.

Nouns: dog, cat, pet.

Verbs: fish, fished, fishing.

Adjective: very, great, significant.

Context: play guitar, piano, games, tennis

To not feel overwhelmed by the math at first, let's start with a easy example:

Words to be processed: boy, girl, woman, man

We can simply divide these words into two dimensions: gender and age as per Table 1 shows. This method is called **Distributed Representation**.

However, human language are so complex that we are not able to pre-define all dimensions for all words. Therefore we want to train a function  $\mathbf{f}$  to get the word embedding for us. Figure 4 shows the function mapping.

To the best of my knowledge, we can do this in two ways.

### AUTOENCODER (GAN)

An autoencoder is a type of neural network that learns to label dimensions in an unsupervised manner. Figure 5 shows that an autoencoder is composed of two parts: the input

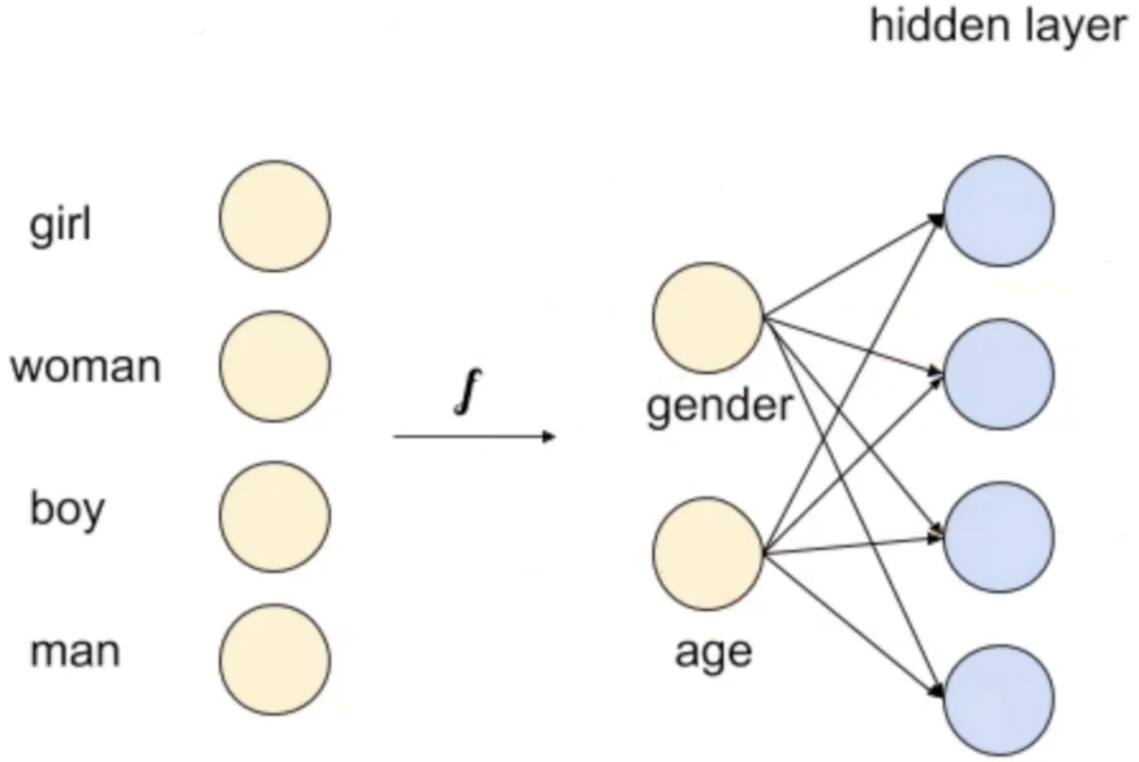


Figure 4: Figure 4: Neural Network Based Mapping

side, marked as  $X$ , and the reconstruction (output) side, marked as  $X'$ . The  $X$  tries to learn a distributed representation of the original data that aims to reduce the dimensions, and the  $X'$  side is where the autoencoder tries to generate from the reduced representation as close as possible to its original input, hence its name, reconstruction.

**BASIC ARCHITECTURE OF AUTOENCODER** An autoencoder consists of two parts, the encoder and the decoder, which can be defined as transitions  $\psi$  and  $\phi$ , such that:

$$\begin{aligned}\psi &: \chi \mapsto \mathfrak{F} \\ \phi &: \mathfrak{F} \mapsto \chi \\ \psi, \phi &= \arg_{\psi, \phi} \min |\chi - (\psi * \phi)\chi|\end{aligned}$$

In the simplest form, we assume there is only one hidden layer. The encoder stage of the autoencoder takes input  $x \in \mathbb{R}^d = \chi$  and maps it to  $\mathbf{h} \in \mathbb{R}^p = \mathfrak{F}$ .

$$\mathbf{h} = \delta(\mathbf{W}x + \mathbf{b})$$

This image  $\mathbf{h}$  is the code section shown in Figure 5.  $\delta$  is normally a sigmoid function defined as following:

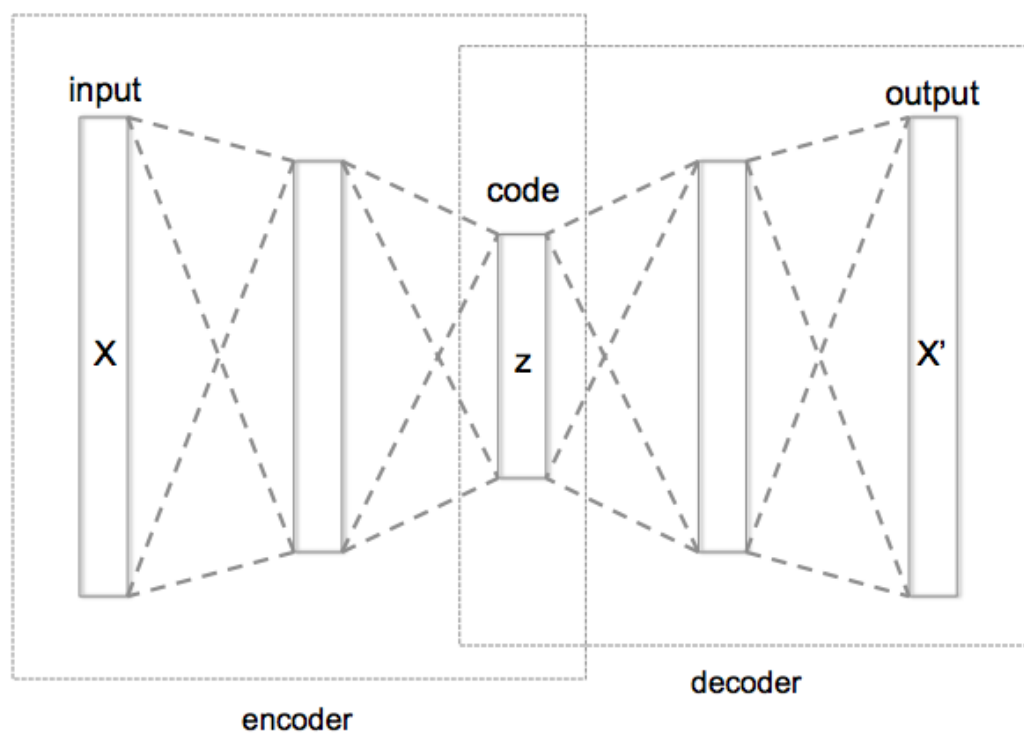


Figure 5: Schematic structure of an autoencoder with 3 fully connected hidden layers. The code ( $z$ , or  $h$  for reference in the text) is the most internal layer.

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

$\mathbf{W}$  is a weight matrix and  $\mathbf{b}$  is a bias vector. The initializations of these two variables are random. They will be learnt iteratively through **backpropagation**.

$$\mathbf{h}' = \sigma(\mathbf{W}'x + \mathbf{b}')$$

The backpropagation computes weight space  $\mathbf{W}$  and bias vector  $\mathbf{b}$  with respect of a loss function  $\mathbb{C}$ .

$x$  : input (vector of features).

$x'$  : target output, the closer to  $x$ , the better.

$\mathbb{C}$ : loss function, which intuitively associate with the "cost" of certain event or representation. The aim of training is to minimize the loss function.

$L$ : the number of layers. In the above example,  $L$  equals to 1. However, here we aims to provide a generalized form.

$W^l = W_{jk}^l$ : the weights between layer  $l-1$  and  $l$ , where  $W_{jk}^l$  is the weight between the  $k$ -th node in layer  $l-1$  and the  $j$ -th node in layer  $l$ .

$f^l$ : sigmoid function at layer  $l$ .

The whole neural network can be represented by this function:

$$g(x) = f^L(W^L f^{L-1}(W^{L-1} \dots f^1(W^1 * x) \dots))$$

For each input  $x_i$ , there will be an output  $x'_i$  corresponding with it. the loss of the model on that pair is the cost of the difference between the predicted output  $g(x_i)$  and the target output  $x'_i$ :

$$\mathbb{C}(y_i, f^L(W^L f^{L-1}(W^{L-1} \dots f^1(W^1 * x) \dots))$$

For each layer, we compute the minimum of the loss function by calculating its derivatives.

$$\frac{\partial \mathbb{C}}{\partial W_{jk}^l}$$

Where the loss function is defined as:

$$\mathbb{C}(x, x') = |x - x'|^2 = |x - \sigma'(\mathbf{W}'(\sigma(\mathbf{W}x + \mathbf{b})) + \mathbf{b}')|$$

where  $x$  is the average across all inputs.



## NEURAL NETWORK LANGUAGE MODEL (NNLM)

However, the above model does not fully utilize the context of words, as in words in this category:

Context: play guitar, piano, games, tennis

More specifically, autoencoder does generate vectors that significantly reduce the dimensions compared to previous attempts but it takes into account less than 1 or two words, hence it is not so great at predicting the next word based on the previous context.

Bengio proposed the following model that makes use of the word order in a document. It can be presented by the conditional probability of the next word given all the previous ones:

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$

where  $w_t$  is the  $t$ -th word.

In more details:

The training set is  $w_1, w_2, \dots, w_t \in V$ , where  $V$  is the vocabulary set that is roughly of size  $10^5$  depending on language.

We need to learn  $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$ , in the sense it gives highest likelihood of the next word.

This is subject to constraints that:

For any choice of  $w_1^{t-1}$

$$\begin{aligned} \sum_{i=1}^{|V|} f(i, w_t, \dots, w_{t-n+1}) &= 1 \\ f &> 0 \end{aligned}$$

The function  $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$  has two parts:

1.  $\mathbb{C} : w_i \rightarrow \mathbb{C}(w_i) \in \mathbb{R}^m$  where  $m$  ranges from 30 to 100, which is significant smaller than the size of  $V$  and  $w_i \in V$ .

2. A probability function over words:

A function  $g$  that maps an input sequence of words  $(\mathbb{C}(w_{t-n+1}), \dots, \mathbb{C}(w_{t-1}))$  to the next word  $w_t \in V$ .

$$f(i, w_t, \dots, w_{t-n+1}) = g(i, \mathbb{C}(w_{t-1}), \dots, \mathbb{C}(w_{t-n+1}))$$

$\mathbb{C}$  is being shared for the whole neural network.  $\mathbb{C}$  simply maps a word  $w_t \in V$  to a *feature vector* of  $\mathbb{R}^m$  shown as **figure 7**.

The function  $g$  might be implemented by a feed-forward or recurrent neural network with parameter  $\omega$ . The overall parameter set is  $\theta = (\mathbb{C}, \omega)$

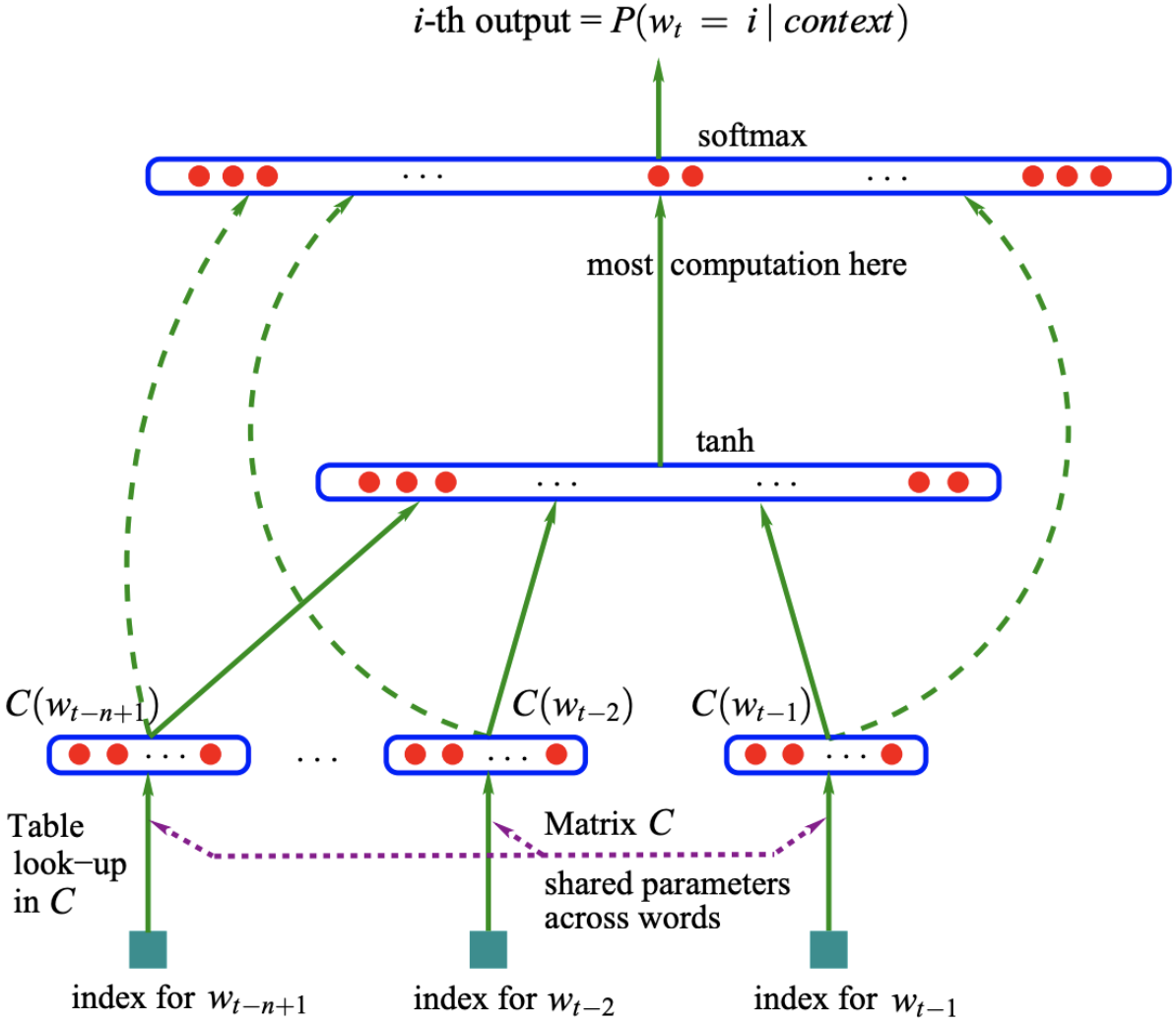


Figure 6: Neural Architecture

Training aims to maximize  $\omega$  the log-likelihood.

$$L = \frac{1}{T} \sum_t \log f(w_t, w_{t-1}, \dots, w_{t-n+1}; \theta) + R(\theta)$$

where  $R$  is a regularization form.

The *softmax* output layer in figure 7 grants the positive probabilities summing up to 1:

$$\hat{P}(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y w_t}}{\sum_i e^{y_i}}$$

The  $y_i$  is the log probability for each possible output word  $w_i \in V$ , defined as follows:

$$y = b + Wx + U \tanh(d + Hx)$$

as figure 7 shows, *tanh* is applied on each word.  $W$  is the weight matrix we described in the previous attempt.  $x$  is the word feature activation vector that is a concatenation from  $\mathbb{C}(w_i)$

$$x = \mathbb{C}(w_{t-1}), \mathbb{C}(w_{t-1}), \dots, \mathbb{C}(w_{t-n+1})$$

$b$  and  $d$  are biases which are additive parameter in neural networks that only have outgoing edges but no incoming edges. Its existence is mainly to not let function passing the origin. Here is a intuitive example of training results:

cat and dog plays similar role in linguistic, syntactically and semantically. backyard and frontyard also plays similar role. running and walking is also a pair. is and was are also quite similar.

From the probability mass of them, we can simply transform:

A cat is walking at frontyard

to A dog is walking at the frontyard

to A dog is running at the frontyard

to A dog was running at the frontyard

to A dog was running at the backyard

#### FOURTH ATTEMPT: WORD2VEC

The above attempt only considers "context" before certain word. However, in 2013, a paper published by Mikolov and Jeff Dean proposed a famous model, Word2Vec, that considers bidirectional context, i.e. context before and after the predicting word.

It has a very similar structure as the attempt we elaborated above, but it is more focusing on **Word Embedding**, i.e. representing a word in vector form, instead of predicting the word based on preceding context.

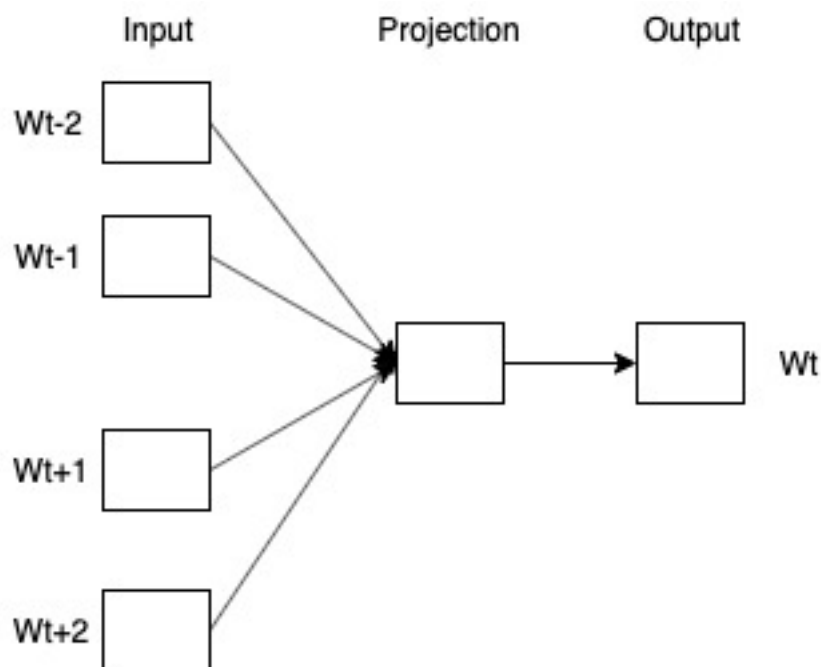


Figure 7: Continuous Bag-of-Words

Word2Vec is a two-layer neural net that takes an input word and outputs it as a feature vector. The core purpose is to gather similar words into the same vector space by calculating the probabilities – the likelihood that they will co-occur.

It has the following two approaches.

#### CONTINUOUS BAG-OF-WORD (CBOW)

The aim of CBOW is to predict the middle word given the words surrounding it.

For example,  $\{I ?? video games\} \rightarrow \{play\}$

As figure 7 shows, this is quite similar to the feed-forward neural network described above. It removes all the hidden layers, and its projection layer is shared among all words, where the projection layer is either the sum or average of all the vectors of the words. Note that despite its name being "Bag-of-words", it bears no relevance to the first attempt we mentioned above.

#### CONTINUOUS SKIP-GRAM (SKIP GRAM)

Opposite to the CBOW model, the objective of *skip gram* is to predict surrounding words from a single word, as shown in figure 8.

Given a sequence of words  $w_1, w_2, w_3, \dots, w_n$ , similar to the maximization of log-likelihood described in the third attempt, the goal of the Skip-gram model is to maximize the average

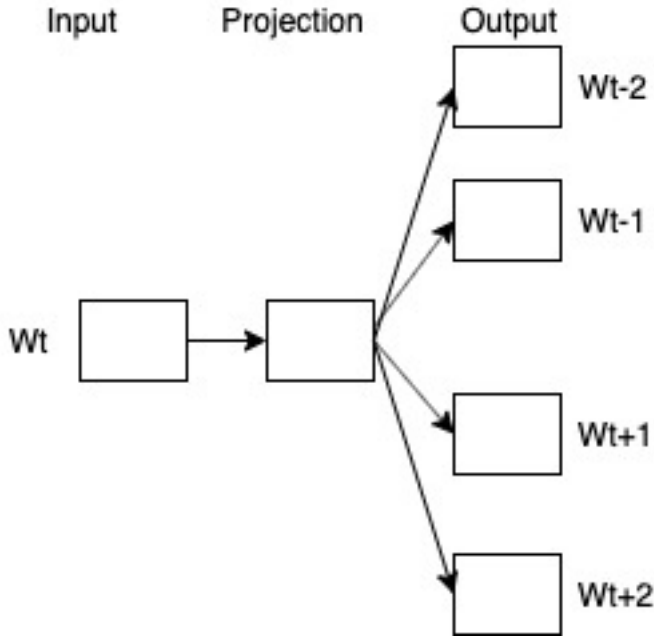


Figure 8: Continuous Skip-Gram

log probability:

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

where  $c$  is the training size. Obviously, larger  $c$  will lead to higher accuracy, in the cost of trading run time. We define  $p(w_{t+j} | w_t)$  using the *softmax* function.

$$p(w_o | w_l) = \frac{\exp(x_{w_o}^T v_{w_l})}{\sum_{w=1}^W \exp(x_{w_o}^T v_{w_l})}$$

where  $x_w$  and  $v_w$  are input and output shown in **figure 8**.  $W$  is the size of vocabulary in the language we trained in.

#### REFINEMENT OF WORD2VEC

**Word2Vec** is a quite expensive algorithm just because of the significant time it takes to run it.

The training complexity of the architecture is proportional to

$$Q = C \times (D + D \times \log_2(V))$$

where  $C$  is the maximum distance of the words.

However, there are several ways we can improve its run-time.

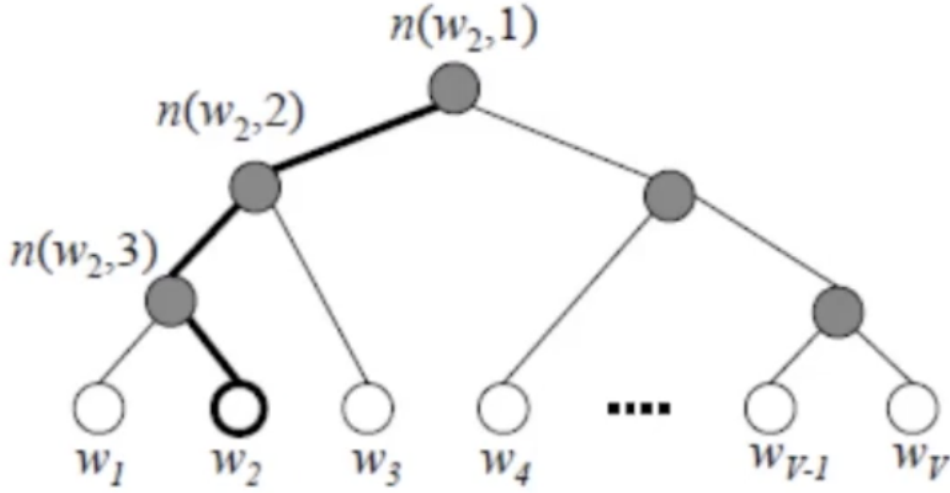


Figure 9: Hierarchical Softmax Model

**HIERARCHICAL SOFTMAX** This model in natural probabilistic language was firstly proposed by Morin and Bengio. The main advantage is that this model improves searching  $W$  to  $\log_2 W$ .

The **Hierarchical Softmax** uses a binary tree representation that represents  $W$  words. For each node, its child nodes' probabilities are explicitly represented.

More intuitively, each node can be reached by traversing from the root of the tree. Nodes that are further down the tree have less probability of being selected as output.

Let  $n(w, j)$  be the  $j$ -th node on the path from root  $w$ , and let  $L(w)$  be the length of the path. therefore,

$$\begin{aligned} n(w, 1) &= \text{root} \\ n(w, L(w)) &= w \end{aligned}$$

For any inner node  $n$ , let  $mid(n)$  be an arbitrary fixed child of  $n$ , and define  $\|x\|$  as

$$\begin{cases} 0 & x \text{ is false} \\ 1 & x \text{ is true} \end{cases}$$

Then we define **Hierarchical Softmax**  $p(w_O|w_I)$  as follows

$$p(w_O|w_I) = \prod_{j=1}^{L(w)-1} \sigma(\|n(w, j+1) = mid(n(w, j))\| x_{n(w, j)}^T v_{w_j})$$

where  $\sigma$  is defined as

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

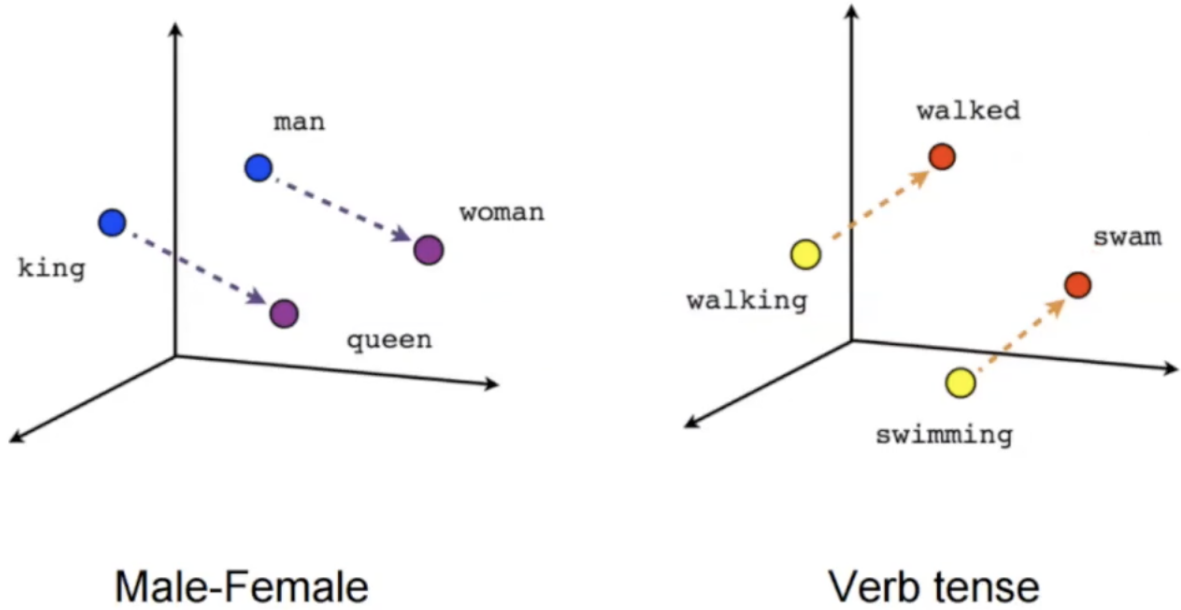


Figure 10: Vector Space of Word2Vec

Also

$$\sum_{w=1}^W p(w|w_I) = 1$$

The above fomulas indicates that the computational cost is proportional to  $L(w_O)$ , which on average is smaller than  $\log W$ .

**NEGATIVE SAMPLING** This model has another name called **Noise Contrastive Estimation (NCE)**, and it was introduced by Gutmann and Hyvarinen.

NCE is roughly a maximum of log probability of the softmax that defines as follows:

$$\log \sigma(x_{w_O}^T v_{w_I}) + \sum_{i=1}^k \mathbb{E}_{w_i} P_n(w) [\log \sigma(-x_{w_O}^T v_{w_I})]$$

$P_n(w)$  is a noise distribution serving as a free parameter in the above formula. Mikolov has observed a number of choices for  $P_n(w)$  and found that the unigram distribution  $U(w)$  raised to the 3/4rd power (i.e.,  $U(w)^{3/4}/Z$ ) outperformed significantly the unigram and the uniform distributions.

#### RESULTS OF WORD2VEC

Figure 10 shows that **Word2Vec** can be successful in creating the vector representation of a word, for example, king and queen has a closer Euclidean distance, and a smaller angle

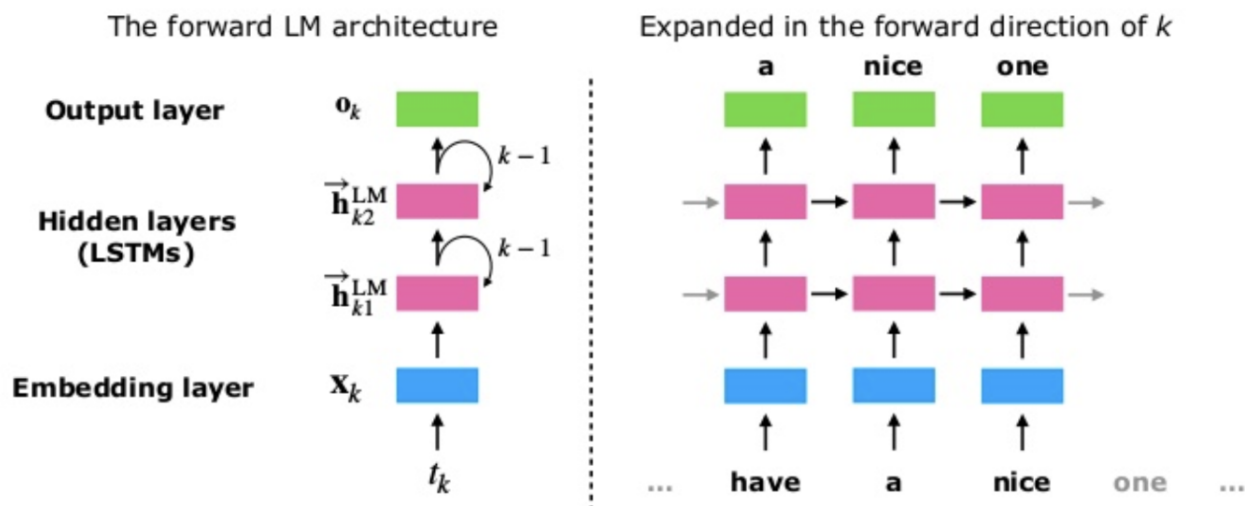


Figure 11: Two-layer Bidirectional LSTM Backbone

in-between.

Moreover, with vector representation of a word, we can write formula based off words. For instance:

France + Paris - Being = China

The result of France + Paris - Being would be equal or very close to the value of Beijing.

Isn't that interesting?

### GLOBAL VECTORS FOR WORD REPRESENTATION (GLOVE)

This model is highly similar to **Word2Vec** introduced above. The only difference is **Word2Vec** looks for limited set of context surrounding the certain word, while **GloVe** calculates based on the whole document. Its training is performed on aggregated global word-word co-occurrence statistics from the corpus. Therefore, **GloVe** can only be offline, which is its main drawback.

This algorithm isn't used that frequently as the relevance between a word and its whole document is quite limited.

### EMBEDDING FROM LANGUAGE MODEL (ELMo)

All the above models does not consider words with multiple meanings. Peters published a paper in 2018 of deep contextualized word representation that models complex characteristics of word use and how these uses vary across linguistic context.

Remarkably, previous models, after training, get constant vectors, while **ELMo** only gets



a pre-trained vector. And the final output will be calculated after fine-tuning. So **ELMo** changes NLP from static to dynamic!

Given a sequence of  $N$  tokens,  $(t_1, t_2, t_3, \dots, t_N)$ , we propose a forward language model that computes the probabilities of the sequence  $t_k$  given all the history token  $(t_1, t_2, t_3, \dots, t_{k-1})$ .

$$p(t_1, t_2, t_3, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, t_2, \dots, t_{k-1})$$

The forward: **ELMo** uses LSTM layer to compute output. At the forward direction, for each position  $k$ , each LSTM layer outputs a context dependent representation  $\vec{\mathbf{h}}_{k,j}^{L,M}$  where  $j = 1, 2, 3, \dots, L$ . The top layer LSTM output  $\vec{\mathbf{h}}_{k,j}^{L,M}$  is used to predict the next token  $t_{k+1}$  with softmax.

The backward: A backward LM is similar to a forward LM, except it runs over the sequence in reverse, predicting the previous token given the future context:

$$p(t_1, t_2, t_3, \dots, t_N) = \prod_{k=1}^N p(t_k | t_{k+1}, t_{k+2}, \dots, t_N)$$

Its implementation is also quite similar: we produce  $\overleftarrow{\mathbf{h}}_{k,j}^{L,M}$  of  $t_k$  given  $(t_{k+1}, t_{k+2}, \dots, t_N)$

The formulation of **ELMo** jointly combines the forward and backward directions that maximize the log likelihood:

$$\sum_{k=1}^N (\log p(t_k | t_1, t_2, \dots, t_{k-1}; \theta_x, \vec{\theta}_{LSTM}, \theta_S) + (\log p(t_k | (t_{k+1}, t_{k+2}, \dots, t_N; \theta_x, \vec{\theta}_{LSTM}, \theta_S)$$

Hence, for each token  $t_k$ , **ELMo** computes a set of  $2L + 1$  representations.

$$\begin{aligned} R_k &= \{x_k^{LM}, \vec{\mathbf{h}}_{k,j}^{L,M}, \overleftarrow{\mathbf{h}}_{k,j}^{L,M} | j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{L,M} | j = 1, \dots, L\} \end{aligned}$$

where  $\mathbf{h}_{k,j}^{L,M}$  is the token layer.

More generally, we compute a task specific weighting of all biLM layers:

$$ELMo_k^{task} = E(R_k; \theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{L,M}$$

where  $s^{task}$  is a softmax-normalized weights and the  $\gamma$  allows task to scale the entire **ELMo** vector. The indices  $k$  and  $j$  correspond to the index of the word and the index of the layer from which the hidden state is being extracted from.

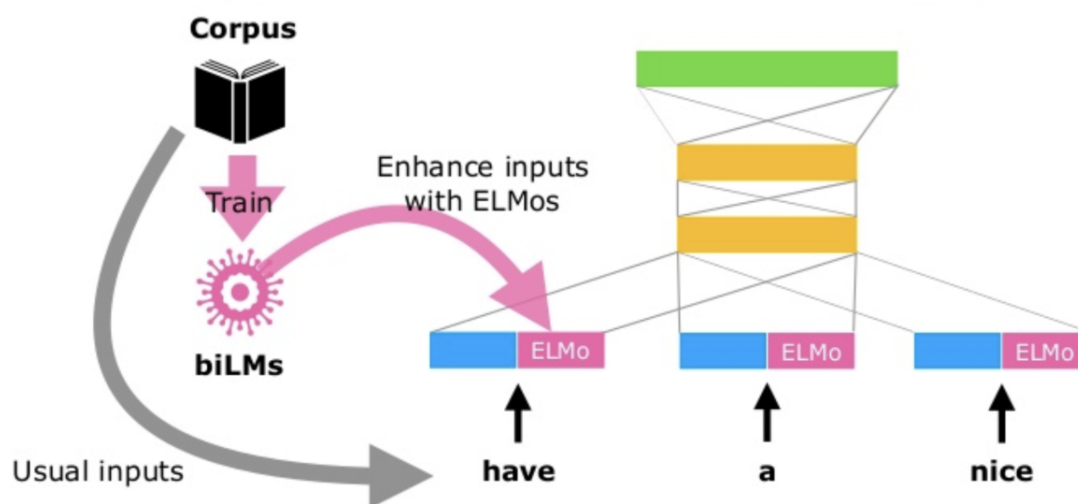


Figure 12: Overview of the procedure on training ELMo

To summarize: • **ELMo** dynamically determines word embedding in downstream task

- ELMo generates three embeddings:

1. word embedding
2. 1st LSTM layer embedding
3. 2nd LSTM layer embedding

- the pre-training step involves getting three embedding  $(v_1, v_2, v_3)$  per word.

- the fine-tuning step involves freezing the embedding and train weight  $(w_1, w_2, w_3)$  for  $(v_1, v_2, v_3)$  per word.

- The final embedding is  $w_1 v_1 + w_2 v_2 + w_3 v_3$

Why three embeddings only? because normal words have less or equal to three meanings, and this is practically proven to work the best.

## BIDIRECTIONAL ENCODER REPRESENTATION FROM TRANSFORMER (BERT)

**BERT** is a model published by Google in 2018 Nov that basically takes advantage of all the previous models that we introduced.

Similar to **ELMo**, it has two phases: pre-training and fine-tuning, and it also utilizes two-layer bidirectional model.

Similar to **Word2Vec**, it also predicts word from a given context, and transform it to a vector.

It uses Transformer which is a model proposed in the famous paper *Attention is All You*

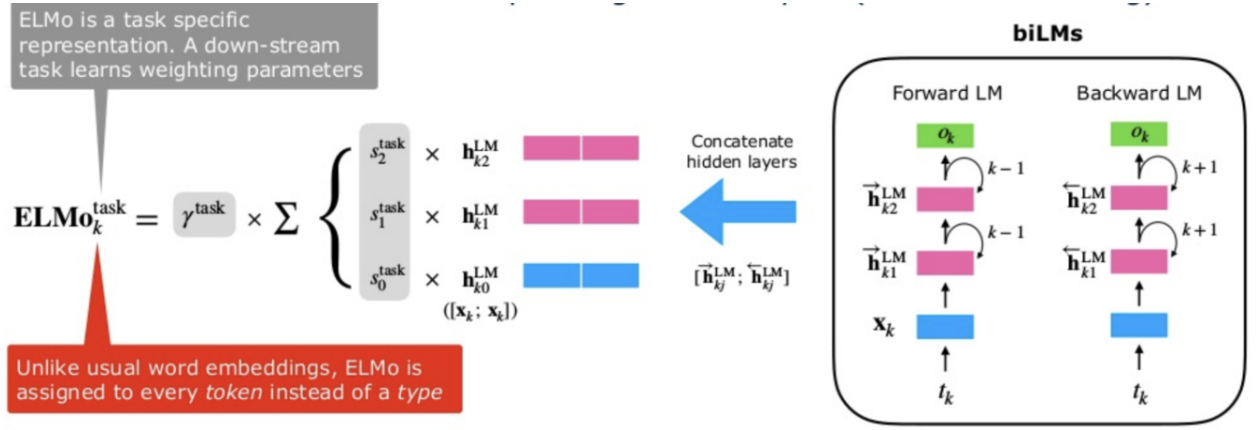


Figure 13: A graphic elaboration on ELMo function

Need instead of RNN-GPT (Generative Pre-training).

In this report, we denote the number of layers, aka, Transformer Blocks, as  $L$ .

The hidden size as  $H$ .

The number of self-attention heads as  $A$ .

#### TRANSFORMER

##### INPUT/OUTPUT REPRESENTATION

The input of **BERT** can be both a sentence or a set of sentences, which is an arbitrary continuous text.

Then **BERT** uses WordPiece embeddings with a 30,000-token library. The first token of a sequence is normally a special classification token ( $[\text{CLS}]$ ). The final hidden state of this token is used as the "aggregate sequence representation for classification tasks". Sentence pairs are packed together into a single sequence. **BERT** differentiate the sentences in two ways.

First, **BERT** separate them with a special token ( $[\text{SEP}]$ ). Second, **BERT** adds a learned embedding to every token implying the sentence it belongs, be it sentence A or sentence B.

As shown in **Figure 14**, **BERT** denotes the input embedding as  $E$ , the final hidden vector of the special  $[\text{CLS}]$  for the  $i$ -th input token as  $T_i \in \mathbb{R}^H$

##### MASKED LANGUAGE MODEL

Intuitively, it is quite obvious that a bi-directional NLP model is better than models that purely considers left-to-right or right-to-left. Sadly, the conditional probability determines that the direction is either left-to-right or right-to-left. A bidirectional model requires the

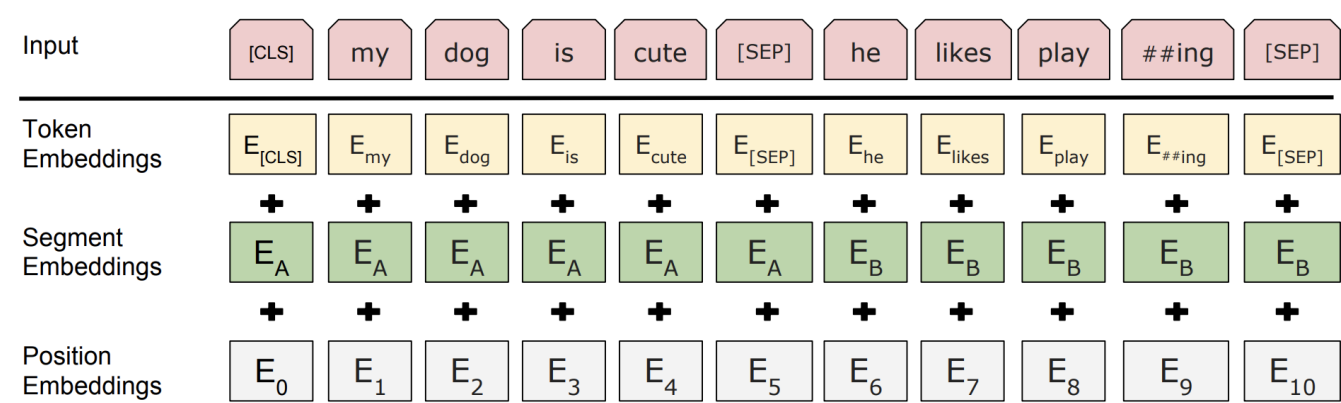


Figure 14: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

words to be predicted to see itself.

How **BERT** solves this problem is using a **Masked Language Model** (Masked LM in short)

QUESTION ANSWERING

## REFERENCES

- [1] Prof. Mike Gal, "Physical waves", Lectures, *UNSW: PHYS1241*
- [2] R Nave, "Work Functions for Photoelectric Effect" *Hyperphysics* (2017), found at <http://hyperphysics.phy-astr.gsu.edu/hbase/Tables/photoelec.html>, accessed 26 Oct. 2017