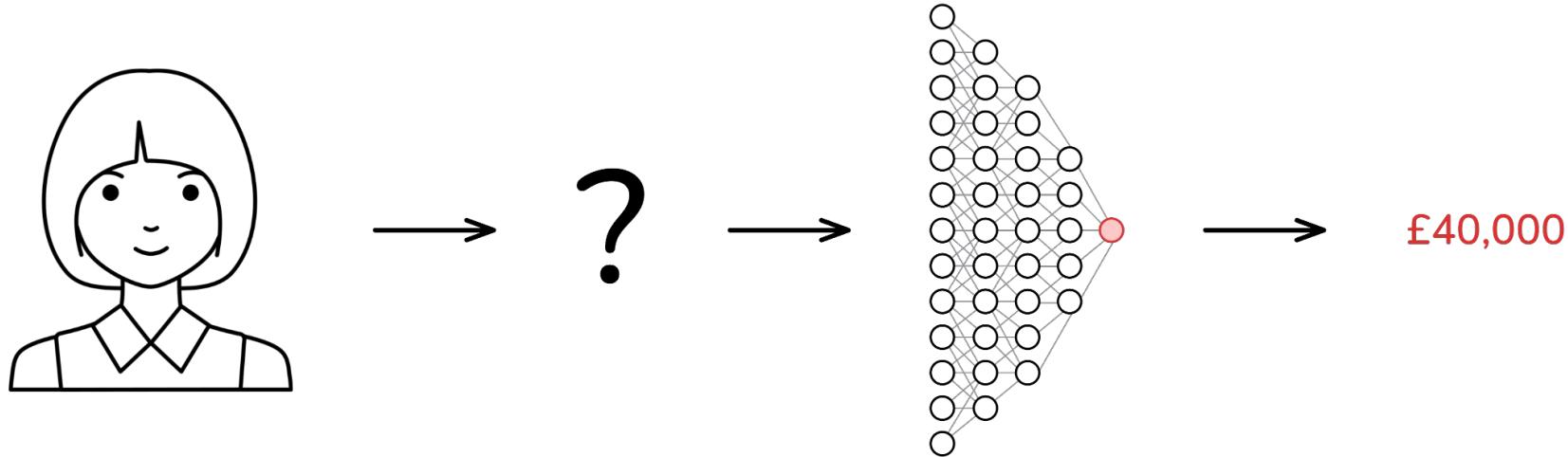


# word2vec

## Word features

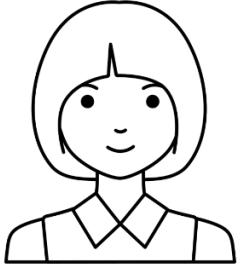
# Income prediction model

# Income prediction model



Convert the user's key characteristics into meaningful numbers the model can process

# Meaningful numbers



- Gender? → Yes, "0" is male, 1 is "female"
- Age? → Yes, nice integer, between 0 and 100
- Height? → Yep, use centimetres, between [20, 250]
- Weight? → Yep, maybe just kilograms, [2, 150]
- Education? → Hmm ...
- City? → Oops ...
- Job? → Oops ...

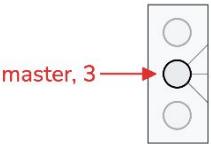
1
28
170
62
?
X
X

# How to encode education?

## Ordinal

0 - none  
 1 - high-school  
 2 - bachelor  
 3 - master  
 4 - PhD

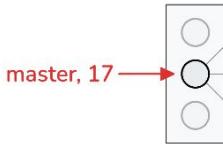
treat "PhD" as 4 times "none"



## Years of Schooling

0 - none  
 12 - high-school  
 15 - bachelor  
 17 - master  
 21 - PhD

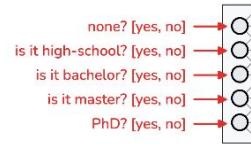
requires some domain knowledge



## One-Hot

[1, 0, 0, 0, 0] - none  
 [0, 1, 0, 0, 0] - high-school  
 [0, 0, 1, 0, 0] - bachelor  
 [0, 0, 0, 1, 0] - master  
 [0, 0, 0, 0, 1] - PhD

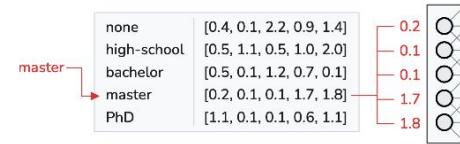
no assumptions, use more neurons



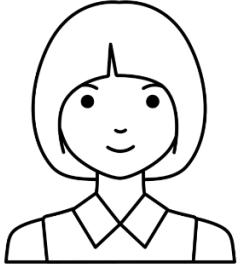
## Learned Embeddings

[0.4, 0.1, 2.2, 0.9, 1.4] - none  
 [0.5, 1.1, 0.5, 1.0, 2.0] - high-school  
 [0.5, 0.1, 1.2, 0.7, 0.1] - bachelor  
 [0.2, 0.1, 0.1, 1.7, 1.8] - master  
 [1.1, 0.1, 0.1, 0.6, 1.1] - PhD

1 - come up with some random numbers i.e. initial features  
 2 - treat those features as trainable parameters  
 3 - more parameters, more data but best accuracy



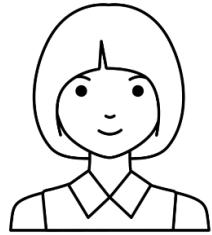
# What about city?



- Gender? → Yes, "0" is male, 1 is "female"
- Age? → Yes, nice integer, between 0 and 100
- Height? → Yep, use centimetres, between [20, 250]
- Weight? → Yep, maybe just kilograms, [2, 150]
- Education? → 3, master
- City? → Oops ...
- Job? → Oops ...

1
28
170
62
3
X
X

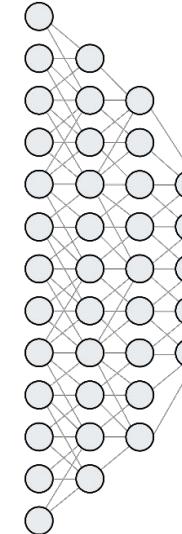
# Recursive Approach



Gender? → 1  
Age? → 28  
Height? → 170  
Weight? → 62  
Education? → 3  
City?

- Population size?
- Average age?
- Average income?
- Crime rate?

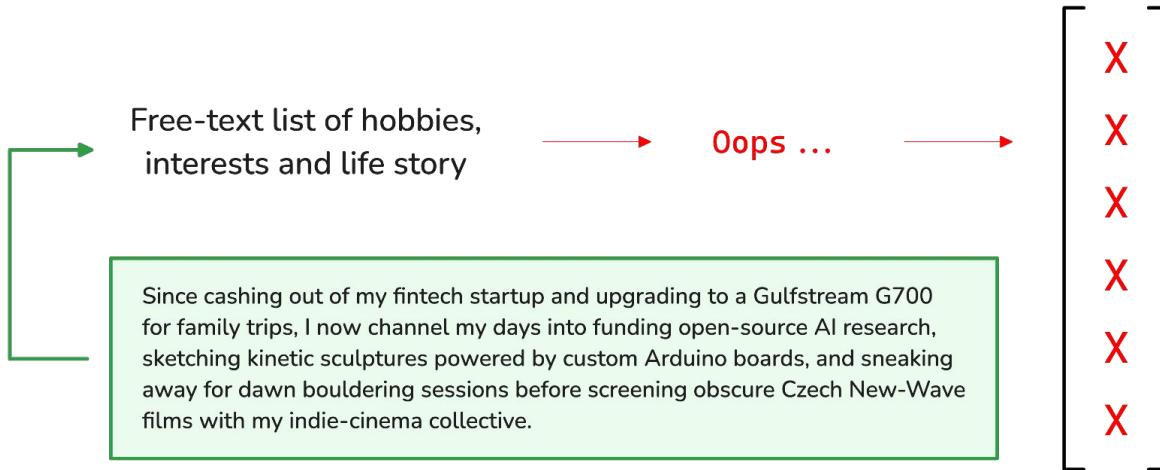
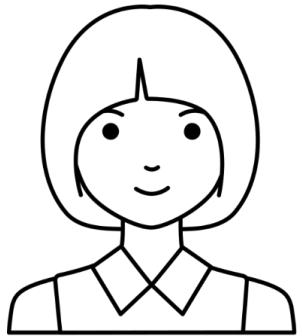
1  
28  
170  
62  
3  
1M  
31  
22k  
4



£40,000

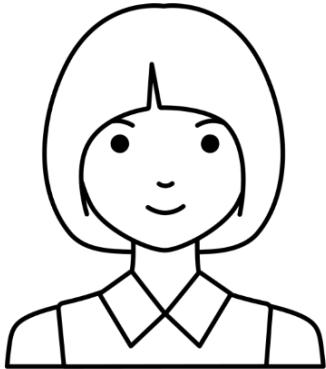
# One more feature...

# One more feature...

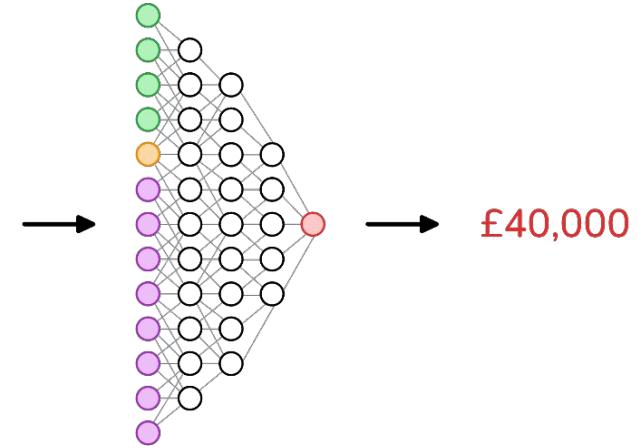


How can we produce meaningful numbers from free text?

# One more feature...



- Gender? → Yes, "0" is male, 1 is "female"
- Age? → Yes, nice integer, between 0 and 100
- Height? → Yep, use centimetres, between [20, 250]
- Weight? → Yep, maybe just kilograms, [2, 150]
- Education? → 3, master
- Description → Since cashing out of my fintech startup...



How can we produce meaningful numbers from free text?

# Distribution hypothesis

What is the meaning of **bardiwac**?

- He handed her her glass of **bardiwiac**.
- Beef dishes are made to complement the **bardiwiacs**.
- Nigel staggered to his feet, face flushed from too much **bardiwiac**.
- I dined off bread and cheese and this excellent **bardiwiac**.



2013



**Efficient Estimation of Word Representations in Vector Space**

---

Tomas Mikolov  
Google Inc., Mountain View, CA  
tmikolov@google.com

Kai Chen  
Google Inc., Mountain View, CA  
kaichen@google.com

Greg Corrado  
Google Inc., Mountain View, CA  
gcorrado@google.com

Jeffrey Dean  
Google Inc., Mountain View, CA  
jeff@google.com

---

**Abstract**

We propose two novel model architectures for computing continuous vector representations of words from very large data sets. The quality of these representations is measured in a word similarity task, and the results are compared to the previously best performing techniques based on different types of neural networks. We observe significant improvements in accuracy at much lower computational cost, i.e. it takes less than a day to train huge word vectors from a 1.6 billion word data set. Furthermore, we show that these vectors provide state-of-the-art performance on our test set for measuring syntactic and semantic word similarities.

**1 Introduction**

Many current NLP systems and techniques treat words as atomic units - there is no notion of similarity between words, as these are represented as indices in a vocabulary. This choice has several good reasons - simplicity, robustness and the observation that simple models trained on huge amounts of data outperform complex systems trained on less data. An example is the popular N-gram model used for statistical language modeling - today, it is possible to train N-grams on virtually all available data (trillions of words) [1].

However, these techniques are at their limits in many tasks. For example, the amount of relevant in-domain data for automatic speech recognition is limited - the performance is usually dominated by the size of high quality transcribed speech data (often just millions of words). In machine translation, the existing corpora for many languages contain only a few billions of words at best. Thus, there are situations where even the basic techniques will not result in any significant progress, and we have to focus on more advanced techniques.

With progress of machine learning techniques in recent years, it has become possible to train more complex models on much larger data set, and they typically outperform the simple models. Probably the most successful concept is to use distributed representations of words [10]. For example, neural network based language models significantly outperform N-gram models [1, 27, 17].

**1.1 Goals of the Paper**

The main goal of this paper is to introduce techniques that can be used for learning high-quality word vectors from huge data sets with billions of words, and with millions of words in the vocabulary. As far as we know, none of the previously proposed architectures has been successfully trained on more

1

arXiv:1301.3781v3 [cs.CL] 7 Sep 2013

**Distributed Representations of Words and Phrases and their Compositionality**

---

Tomas Mikolov  
Google Inc., Mountain View, CA  
tmikolov@google.com

Ilya Sutskever  
Google Inc., Mountain View, CA  
ilyasut@google.com

Kai Chen  
Google Inc., Mountain View, CA  
kai@google.com

Greg Corrado  
Google Inc., Mountain View, CA  
gcorrado@google.com

Jeffrey Dean  
Google Inc., Mountain View, CA  
jeff@google.com

---

**Abstract**

The recently introduced continuous skip-gram model is an efficient method for learning high-quality distributed vector representations that capture a large number of precise syntactic and semantic relationships in text. In this paper we present several extensions that improve both the quality of the vectors and the training speed. By subsampling of the frequent words we obtain significant speedup and also learn more regular word representations. We also describe a simple alternative to the hierarchical softmax called negative sampling.

An important limitation of word embeddings is their lack of disambiguity to word order and their inability to represent idiomatic phrases. For example, the meanings of "Canada" and "Air" cannot be easily combined to obtain "Air Canada". Motivated by this example, we present a simple method for finding phrases in text, and show that learning good vector representations for millions of phrases is possible.

**1 Introduction**

Distributed representations of words in a vector space help learning algorithms to achieve better performance in natural language processing tasks by grouping similar words. One of the earliest use of word representations dates back to 1986 due to Rumelhart, Hinton, and Williams [13]. This idea has since been applied to standard language modeling with considerable success [1]. The followings are some applications to automatic speech recognition and machine translation [14, 7], and a wide range of NLP tasks [2, 20, 15, 3, 18, 19, 9].

Recently, Mikolov et al. [8] introduced the Skip-gram model, an efficient method for learning high-quality vector representations of words from large amounts of unstructured text data. Unlike most of the previously used neural network architectures for learning word vectors, training of the Skip-gram model (Figure 1) does not involve dense matrix multiplications. This makes the training extremely efficient as optimized single-machine implementation can train on more than 100 billion words in a day.

The word representations computed using neural networks are very interesting because the learned vectors explicitly encode many linguistic regularities and patterns. Somewhat surprisingly, many of these patterns can be represented as linear translations. For example, the result of a vector calculation  $\text{vec}(\text{"Madrid"}) - \text{vec}(\text{"Spain"}) + \text{vec}(\text{"France"})$  is closer to  $\text{vec}(\text{"Paris"})$  than to any other word vector [9, 8].

1

arXiv:1310.4546v1 [cs.CL] 16 Oct 2013

Mikolov et al. 2013 (January)

Mikolov et al. 2013 (October)

# word2vec - recipe

**1**

Treat the entire corpus as one long string.  
 Slide a window of five tokens across it; for each window, mask the middle token and use the remaining four to predict the "?"

For dinner we served dark red bardiwac with steak and mushrooms.

For dinner we served dark red bardiwac with steak and mushrooms.

dark red ? with steak

? =  $f(\text{dark, red, with, steak})$

**2**

Because any given context can sensibly be completed by several different words, we don't force the model to output one "correct" token. Instead, we teach it to return a probability distribution over the whole vocabulary

$$p_{\theta}(\cdot | \text{"dark"}, \text{"red"}, \text{"with"}, \text{"steak"})$$

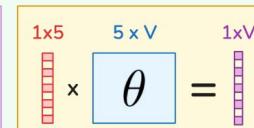
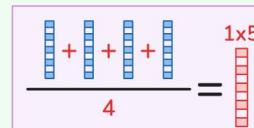

**3**

## CBOW

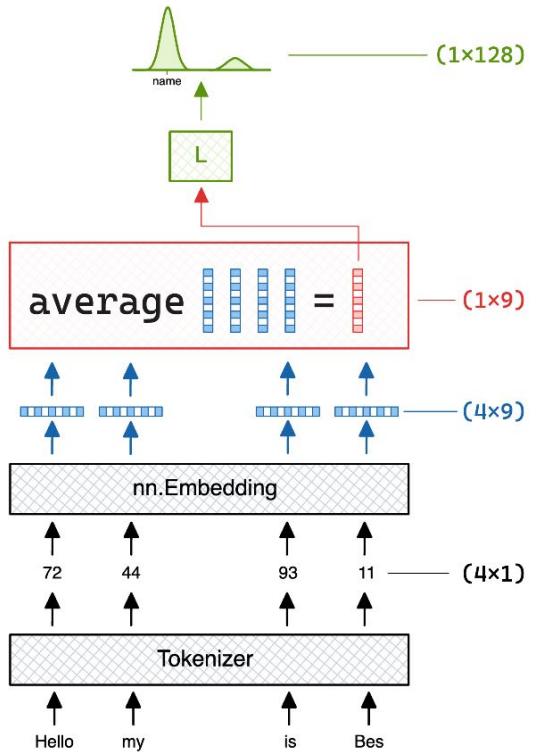
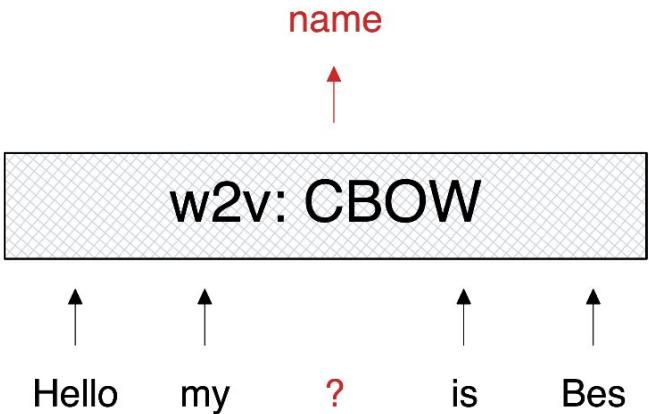
Continuous Bag  
Of Words

dark →  
red →  
with →  
steak →

1	dark	[0.4, 0.1, 2.2, 0.9, 1.4]
2	steak	[0.5, 1.1, 0.5, 1.0, 2.0]
...	...	...
23,451	with	[0.5, 0.1, 1.2, 0.7, 0.1]
23,452	red	[0.1, 1.1, 0.2, 0.5, 0.2]



# Continuous Bag of Words



```

1  #
2  #
3  #
4  import torch
5  import torch.nn.functional as F
6  #
7  #
8  #
9  #
10 vocab = {
11   "Hello": 72,
12   "my": 44,
13   "name": 21,
14   "is": 93,
15   "Bes": 11
16 }
17 #
18 #
19 #
20 sentence = ["Hello", "my", "is", "Bes"]
21 #
22 #
23 #
24 #
25 class CBOW(torch.nn.Module):
26     def __init__(self):
27         super(CBOW, self).__init__()
28         self.emb = torch.nn.Embedding(128, 9)
29         self.linear = torch.nn.Linear(9, 128)
30     def forward(self, inputs):
31         embs = self.emb(inputs)
32         embs = embs.mean(dim=1)
33         out = self.linear(embs)
34         probs = F.log_softmax(out, dim=1)
35         return probs
36
37
38

```

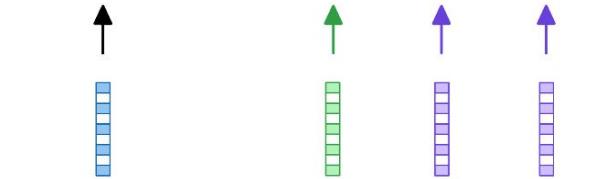
# Skip-Gram

Hello my is Bes

w2v: skip-gram

↑  
name

$$J(\mathbf{v}_c, \mathbf{v}_w, \mathbf{v}_{n_1}, \dots, \mathbf{v}_{n_k}) = -\log(\sigma(\mathbf{v}_w^\top \mathbf{v}_c)) - \sum_{i=1}^k \log(\sigma(-\mathbf{v}_{n_i}^\top \mathbf{v}_c))$$

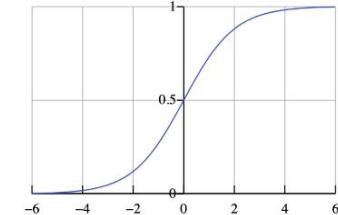


nn.Embedding

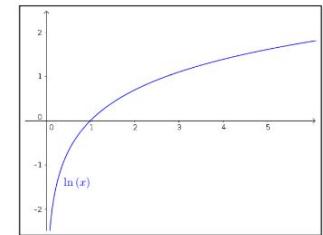
Tokenizer

name Hello sport rocket

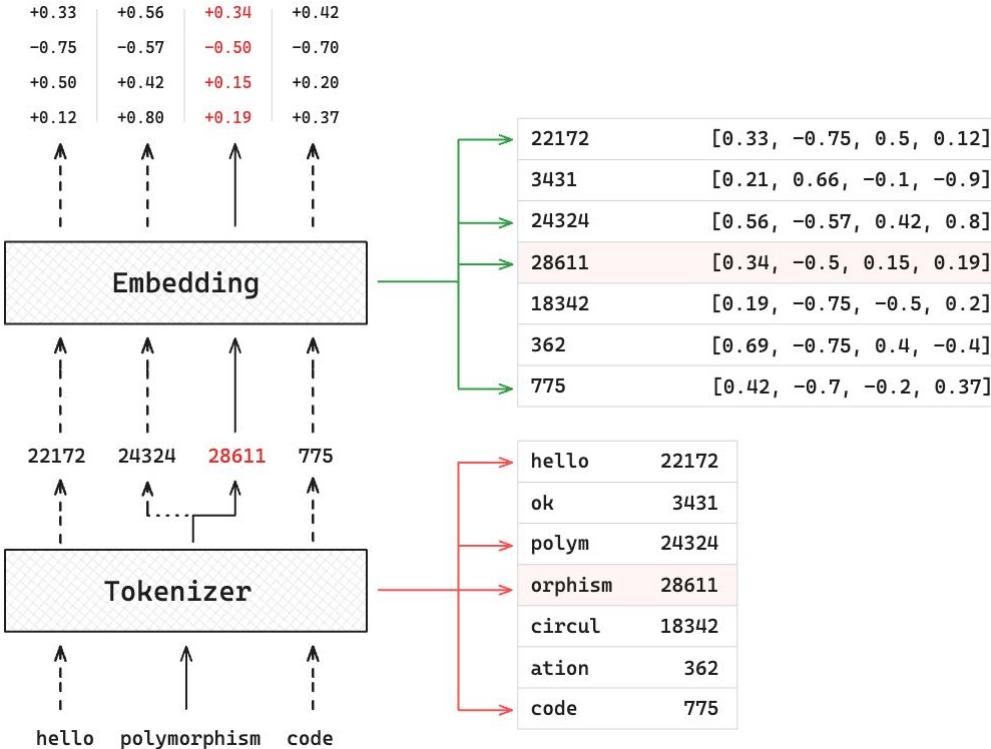
Sigmoid



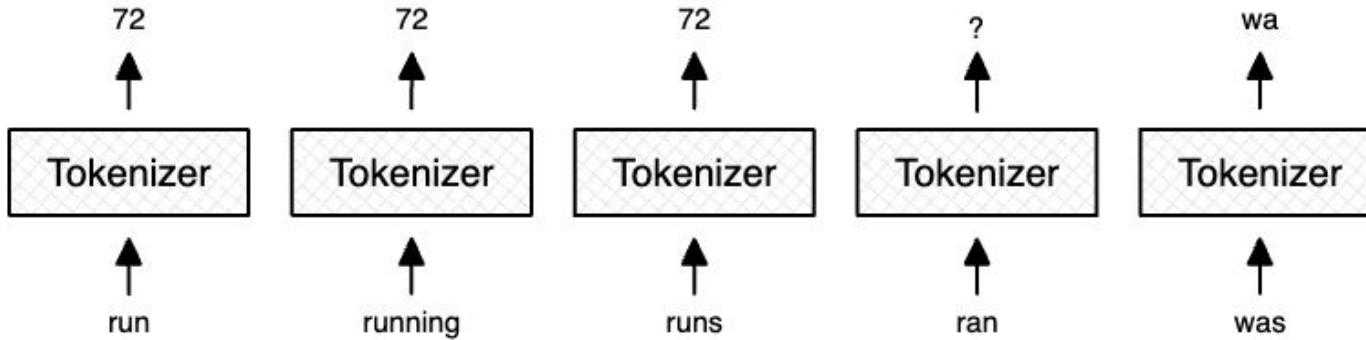
Logarithm



# Tokenizer

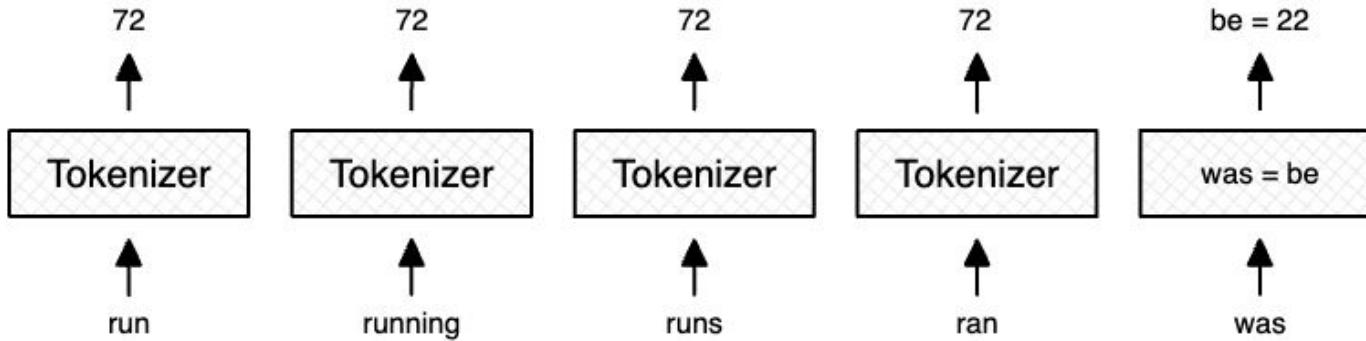


# Stemming

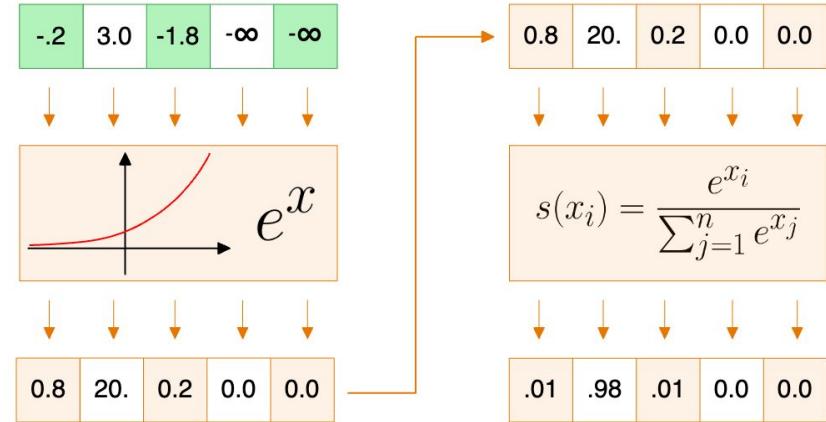
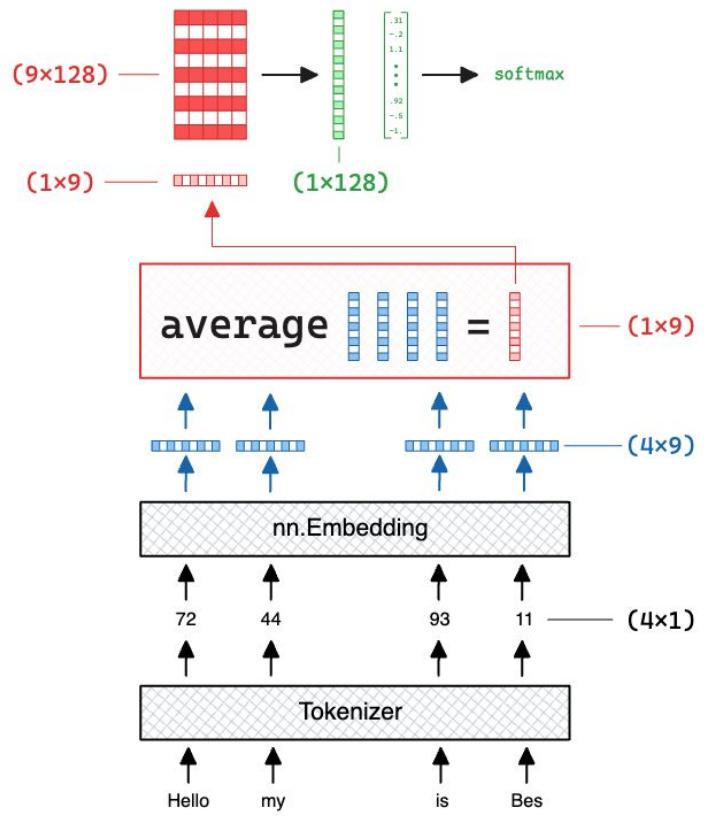


remove ["ing", "s"]

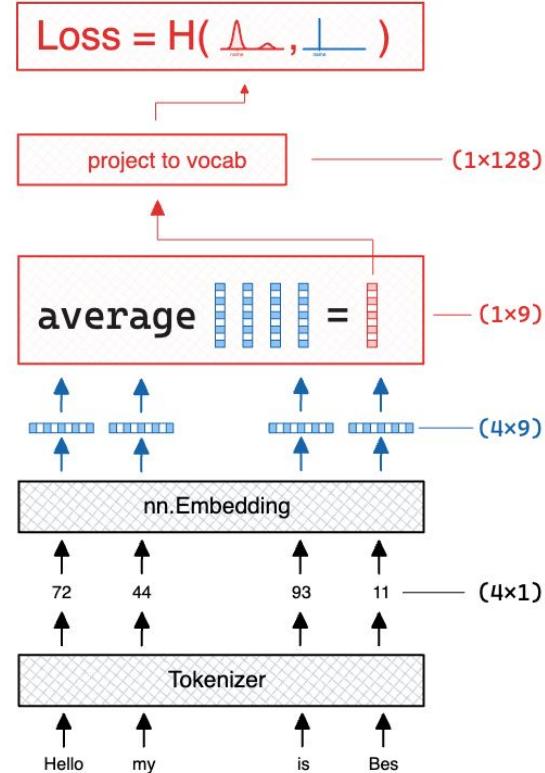
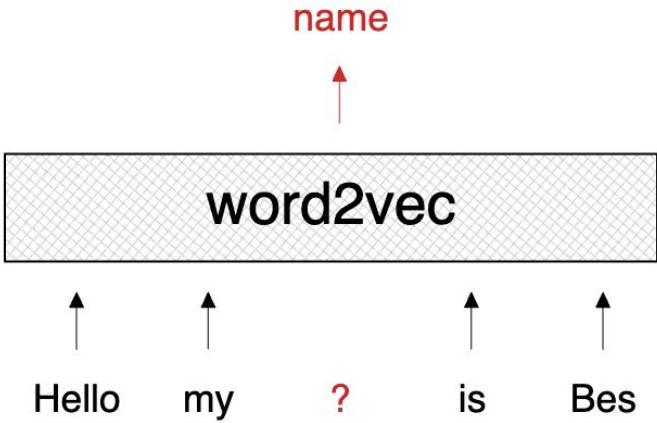
## Lemmatization



vocab mapping



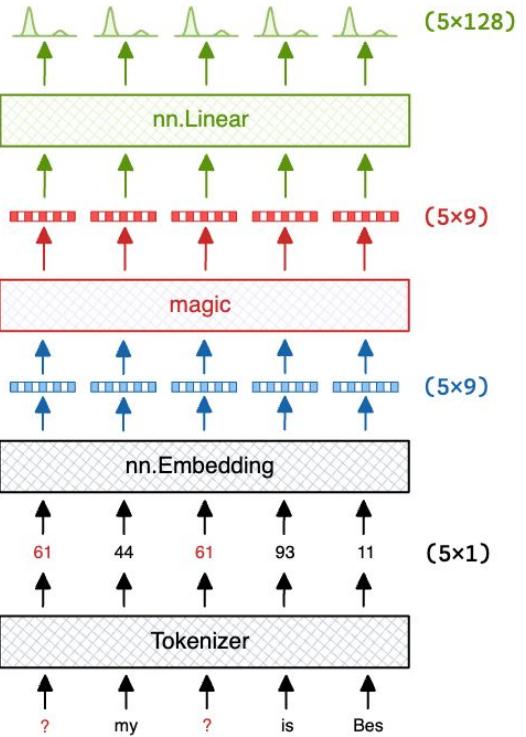
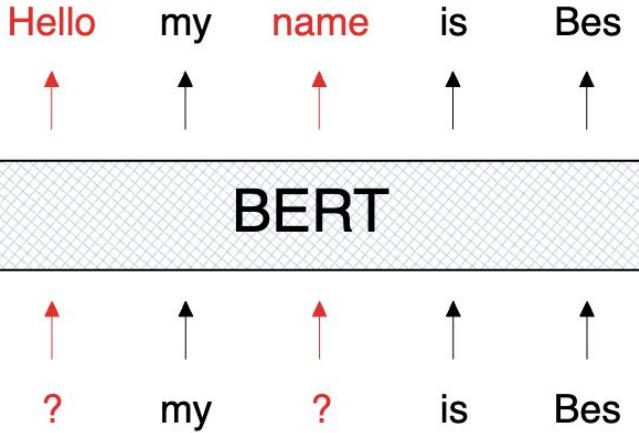
# Overview



```

1  #
2  #
3  #
4  import torch
5  import torch.nn.functional as F
6  #
7  #
8  #
9  #
10 vocab = {
11   "Hello": 72,
12   "my": 44,
13   "name": 21,
14   "is": 93,
15   "Bes": 11
16 }
17 #
18 #
19 #
20 sentence = ["Hello", "my", "is", "Bes"]
21 #
22 #
23 class CBOW(torch.nn.Module):
24     def __init__(self):
25         super(CBOW, self).__init__()
26         self.emb = torch.nn.Embedding(128, 9)
27         self.linear = torch.nn.Linear(9, 128)
28     def forward(self, inputs):
29         embs = self.emb(inputs)
30         embs = embs.mean(dim=1)
31         out = self.linear(embs)
32         out = F.log_softmax(out, dim=1)
33         return out
34 
```

# BERT



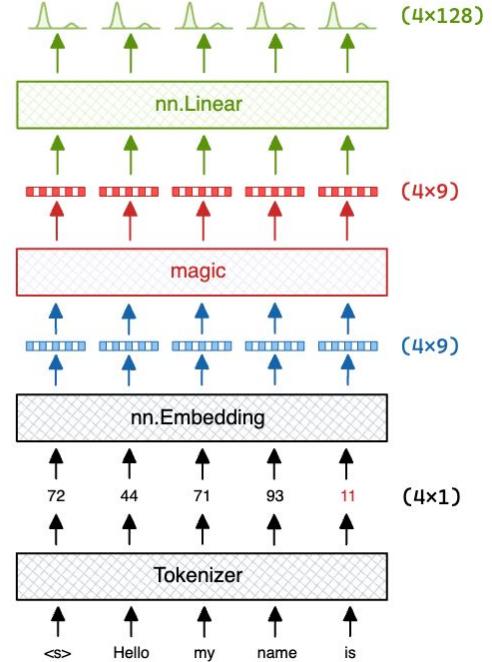
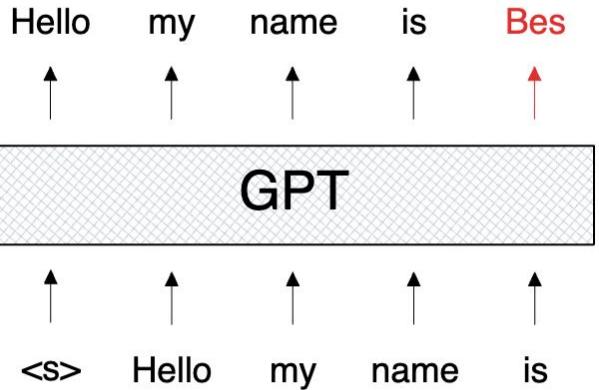
```

1  # 
2  # 
3  import torch
4  #
5  #
6  class Block(torch.nn.Module):
7  ... def __init__(self):
8  ... super().__init__()
9  ... self.W_Q = torch.nn.Linear(9, 9)
10 ... self.W_K = torch.nn.Linear(9, 9)
11 ... self.W_V = torch.nn.Linear(9, 9)
12 ... self.ins = torch.nn.Linear(9, 9)
13 ... self.relu = torch.nn.ReLU()
14 ... self.out = torch.nn.Linear(9, 9)
15 ...
16 ... def forward(self, x):
17 ...     Q = self.W_Q(x)
18 ...     K = self.W_K(x)
19 ...     V = self.W_V(x)
20 ...     x = Q @ K.T
21 ...     x = x @ V
22 ...     x = self.out(x)
23 ...     x = self.relu(x)
24 ...     x = self.ins(x)
25 ...
26 ...     return x
27 ...
28 ...
29 ...
30 class BERT(torch.nn.Module):
31 ... def __init__(self):
32 ... super().__init__()
33 ... self.embed = torch.nn.Embedding(53, 9)
34 ... self.blocks = [Block() for _ in range(5)]
35 ... self.linear = torch.nn.Linear(9, 53)
36 ...
37 ... def forward(self, x):
38 ...     x = self.embed(x)
39 ...     x = x.squeeze(0)
40 ...     for block in self.blocks: x = block(x)
41 ...     x = self.linear(x)
42 ...
43 ...     return x
44

```

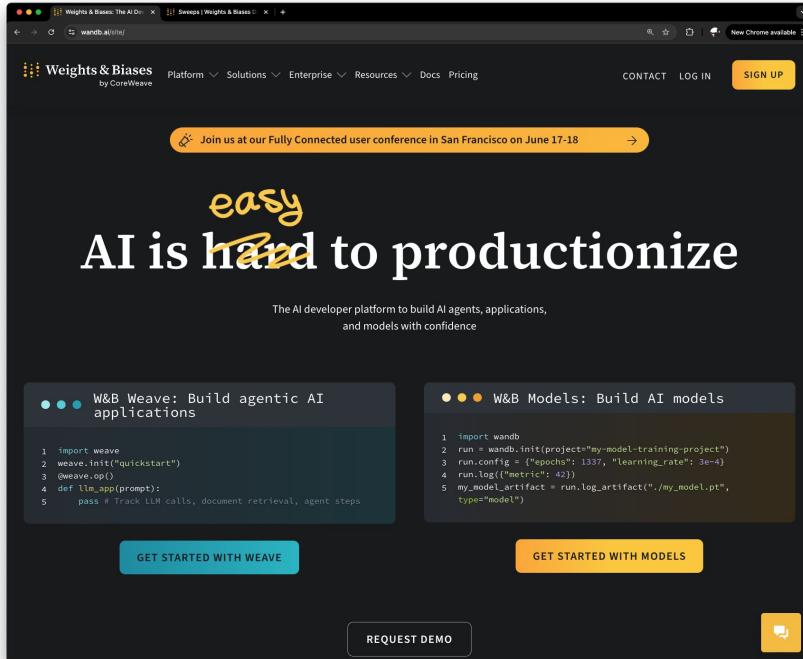


# Overview

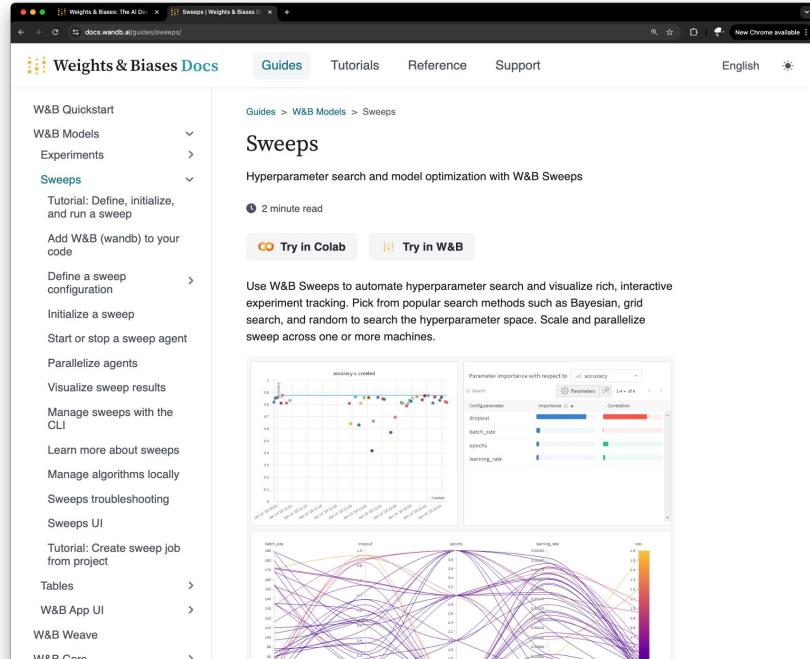


```
1 import torch
2 import torch.nn.functional as F
3
4
5 vocab = {
6     "I": 61,
7     "Hello": 72,
8     "my": 44,
9     "name": 21,
10    "is": 93,
11    "Bes": 11
12 }
13
14
15 sentence = ["?", "my", "?", "is", "Bes"]
16
17
18 class Magic(torch.nn.Module):
19     def __init__(self):
20         super(Magic, self).__init__()
21         ...
22     def forward(self, inputs):
23         ...
24
25
26 class BERT(torch.nn.Module):
27     def __init__(self):
28         super(BERT, self).__init__()
29         self.emb = torch.nn.Embedding(128, 9)
30         self.magics = torch.nn.ModuleList([Magic() for _ in range(3)])
31         self.linear = torch.nn.Linear(9, 128)
32
33     def forward(self, inputs):
34         embs = self.emb(inputs)
35         for magic in self.magics:
36             embs = magic(embs)
37         probs = F.log_softmax(embs, dim=1)
38
39         return probs
40
```

# Weights & Biases



The homepage features a large dark banner with the text "easy AI is hard to productionize". Below the banner, there's a callout for the "Fully Connected" conference and two sections: "W&B Weave: Build agentic AI applications" and "W&B Models: Build AI models". Each section includes a snippet of code and a "GET STARTED" button.



The documentation page for "Sweeps" shows the navigation bar "Weights & Biases Docs" and "Sweeps". It includes a sidebar with links like "W&B Quickstart", "W&B Models", "Experiments", "Sweeps", and "Tables". The main content area has a "Try in Colab" button and a "Try in W&B" button. It explains what W&B Sweeps are and provides a screenshot of the interface showing a scatter plot of accuracy vs. created and parameter importance.

# Good luck!

