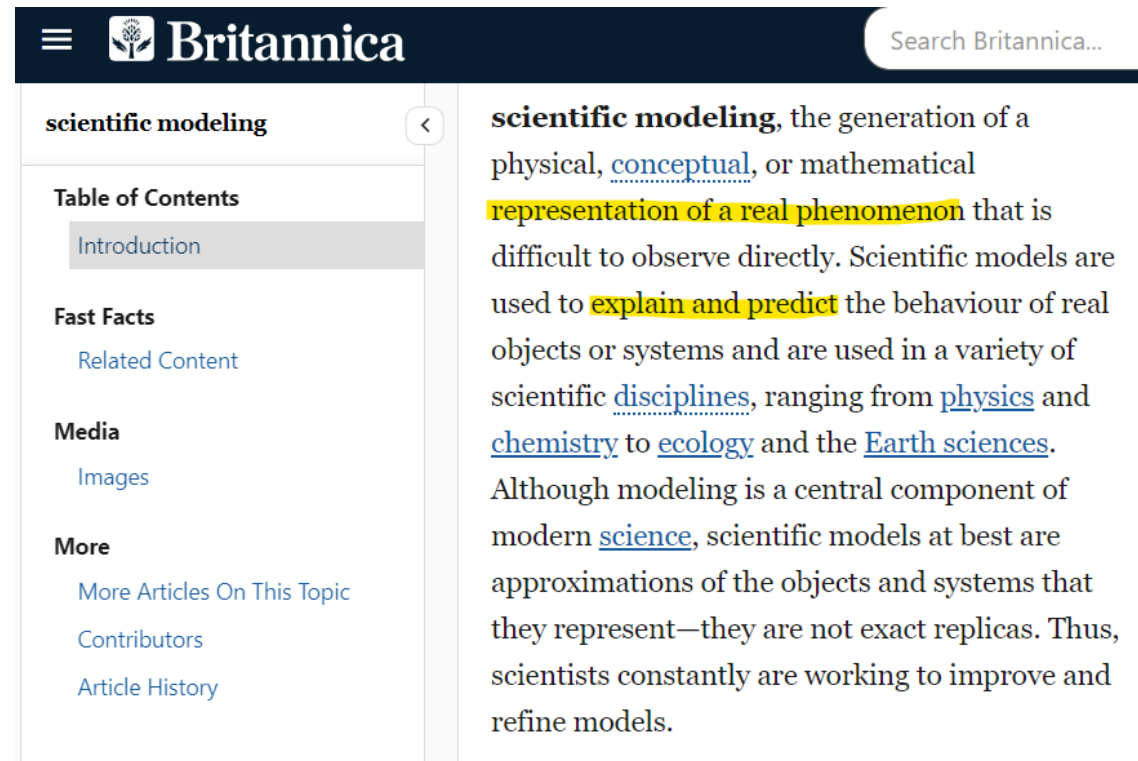


# IS883: Synthesizing Digital Efforts

Mohannad Elhamod

# Language Modeling

# What does a model mean?



The screenshot shows the Britannica website interface for the article "scientific modeling". The left sidebar contains a "Table of Contents" with "Introduction" selected, and sections for "Fast Facts", "Media", and "More" with links to "Related Content", "Images", "More Articles On This Topic", "Contributors", and "Article History". The main content area defines scientific modeling as the generation of a physical, conceptual, or mathematical representation of a real phenomenon that is difficult to observe directly. It states that scientific models are used to explain and predict the behavior of real objects or systems and are used in a variety of scientific disciplines, ranging from physics and chemistry to ecology and the Earth sciences. It also notes that although modeling is a central component of modern science, scientific models at best are approximations of the objects and systems that they represent—they are not exact replicas. Thus, scientists constantly are working to improve and refine models.

**scientific modeling**, the generation of a physical, [conceptual](#), or mathematical [representation of a real phenomenon](#) that is difficult to observe directly. Scientific models are used to [explain and predict](#) the behaviour of real objects or systems and are used in a variety of scientific [disciplines](#), ranging from [physics](#) and [chemistry](#) to [ecology](#) and the [Earth sciences](#). Although modeling is a central component of modern [science](#), scientific models at best are approximations of the objects and systems that they represent—they are not exact replicas. Thus, scientists constantly are working to improve and refine models.

# What is language modeling anyway?...

Web search engine / ...

I saw a cat|

[Lena-volta](#)

Send

▼

To

Cc

Add a subject

Greetings,

I would like to

Tab

Mohannad Elhamod

Clinical Assistant Professor

Boston University | Questrom School of Business

[elhamod@bu.edu](mailto:elhamod@bu.edu)

QUESTROM

MEANS

BUSINESS.

# What is language modeling anyway?...

I grabbed the **branch** and broke it.

I went to the **branch** and deposited some money.

**Context matters!**

# What is language modeling anyway?...

- I went to \_\_\_\_.
- I woke up at 7 am and went to \_\_\_\_.
- I woke up at 7 am, packed my book and notebook, and went to \_\_\_\_.

**The more context, the more certain**

# What is language modeling anyway?...

I went to the **branch** and deposited some money.

I went to the **bank** and deposited some money.

I went to the **ATM** and deposited some money.

Words which frequently appear in **similar contexts** have **similar meaning**.

[Lena Voita](#)

# What is language modeling anyway?...

I sat at the bank and ... { ... watched the water flow.  
... waited for my turn.

## We process language sequentially\*.

\*We will talk about exceptions later...



# Natural Language Processing (NLP)

Includes text generation:

- Text completion.
- Text summarization.
- Question answering.

But there are also many other tasks such as Text classification: (e.g., Sentiment analysis, Reviews, Fake news) or word classification.

# Formalizing our thoughts

- So, language modeling is the chaining of word probabilities.
- How do we calculate these probabilities?

$P(\text{I saw a cat on } \dots) =$

$P(\text{I}) \cdot P(\text{saw}|\text{I}) \cdot P(\text{a}|\text{I saw}) \cdot P(\text{cat}|\text{I saw a}) \cdot P(\text{on}|\text{I saw a cat}) \cdot \dots$

Probability of I saw a cat on

[Lena-volta](#)

counting...

$$P(\text{cat}) = \frac{N(\text{"cat" in corpus})}{N(\text{all words in corpus})}$$

$$P(\text{cat} | \text{my}) = \frac{N(\text{"my cat" in corpus})}{N(\text{"my" in corpus})}$$

Can you foresee any problem with this calculation?...

# N-grams

Instead, let's just use a context of *fixed-length*.

$P(\text{I saw a cat on a mat}) =$

- $P(\text{I})$
- $P(\text{saw} \mid \text{I})$
- $P(\text{a} \mid \text{I saw})$
- $P(\text{cat} \mid \text{I saw a})$
- $P(\text{on} \mid \text{I saw a cat})$
- $P(\text{a} \mid \text{I saw a cat on})$
- $P(\text{mat} \mid \text{I saw a cat on a})$

- $n=3$  (trigram model):  $P(y_t \mid y_1, \dots, y_{t-1}) = P(y_t \mid y_{t-2}, y_{t-1})$ ,
- $n=2$  (bigram model):  $P(y_t \mid y_1, \dots, y_{t-1}) = P(y_t \mid y_{t-1})$ ,
- $n=1$  (unigram model):  $P(y_t \mid y_1, \dots, y_{t-1}) = P(y_t)$ .

[Lena-volta](#)

# N-grams

Context is like a sliding window into the past.

Hugging Face is a startup based in New York City and Paris

$p(\text{word})$

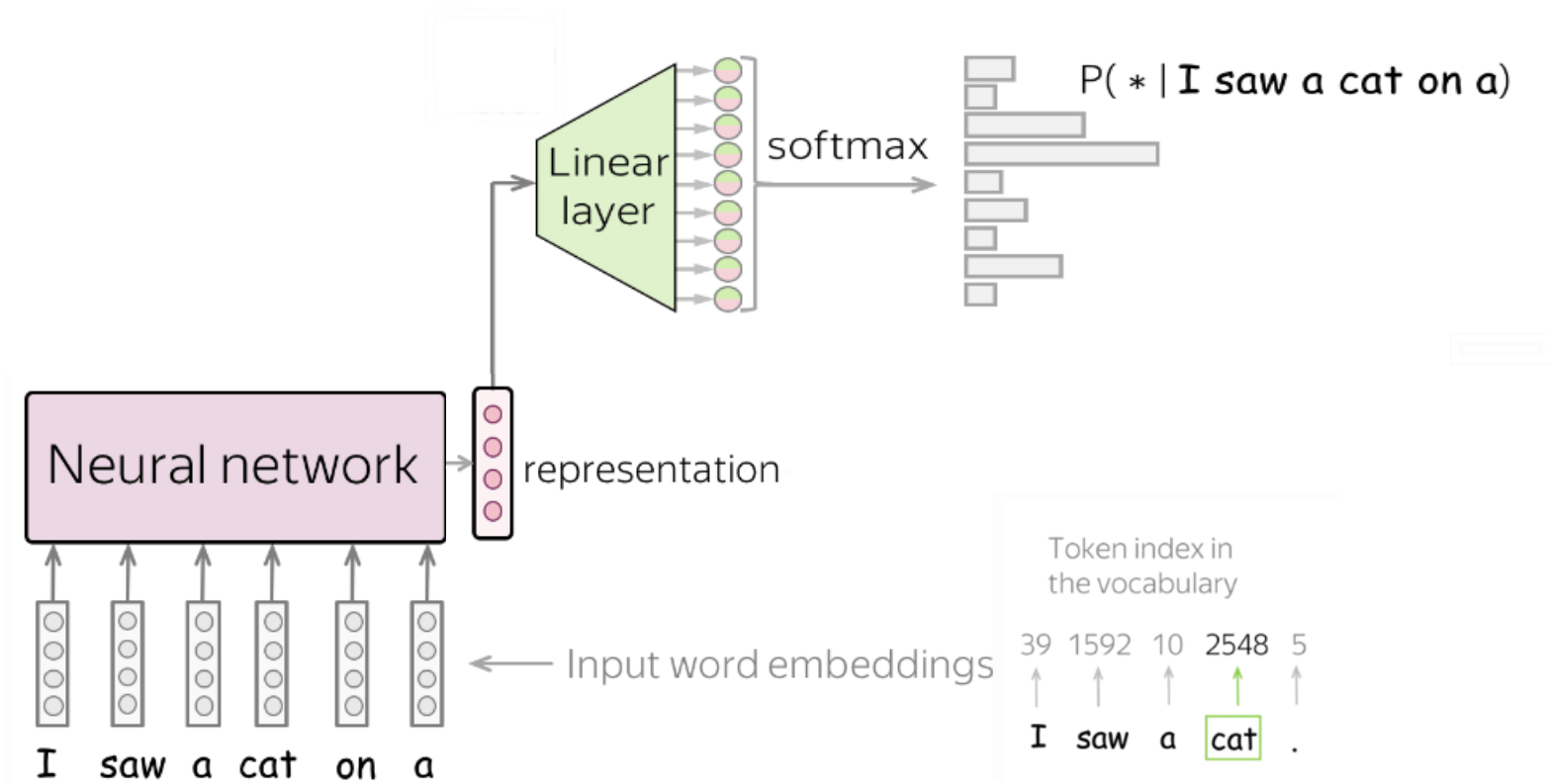
[Huggingface](#)

# Context size

- I went to the beach. My wife sat next to me. She was replying to some emails, and the bird stole our sandwich. Then it started raining suddenly and \_\_\_\_.
- Longer context: predictable outcome.
- Shorter context: Too unpredictable.

# Neural networks for language modeling

# General Model Architecture



[Lena-volta](#)

Can you see any issue with inputting words in an NN?

# What is an embedding?

- embeddings = representation = features = latent space.
- It is a representation of your input.
- [Example for images.](#)



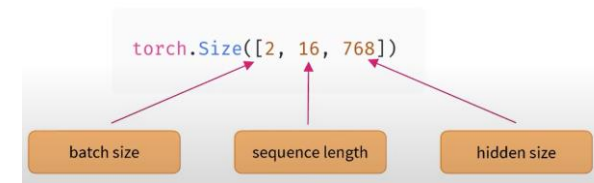
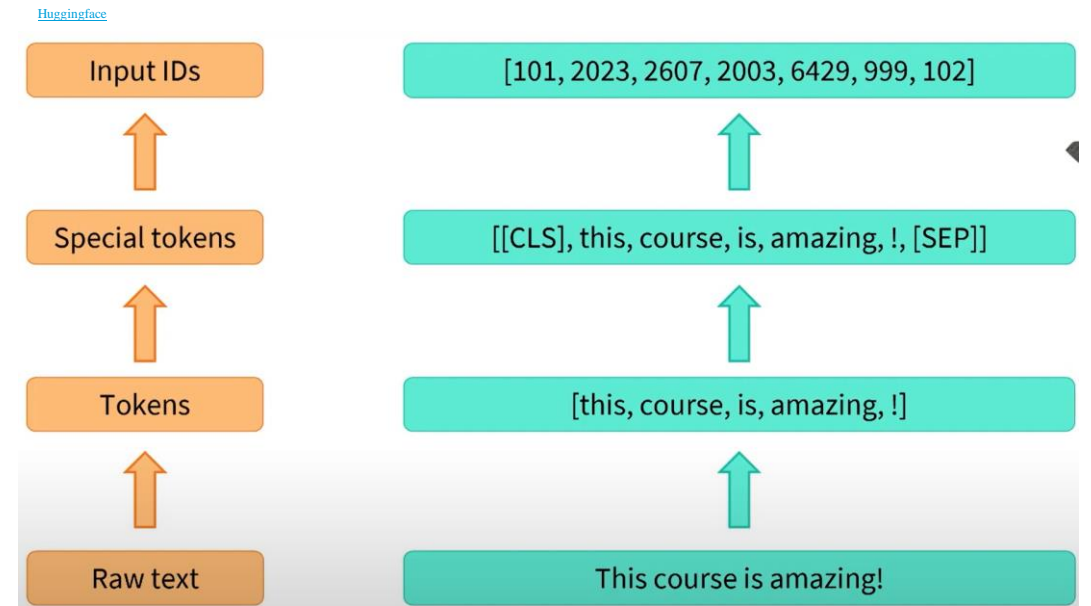
# Word embeddings

- We ideally want words that have similar meanings to have smaller distances.
- [Demo](#)
- Examples:
  1. [Word2Vec \(Google\)](#)
  2. [GloVe \(Stanford\)](#)
  3. [Train your own!](#)

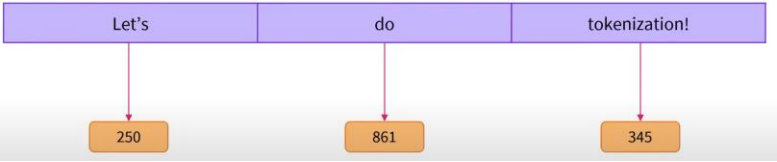
# Tokenization

So, every time we have sentences to generate, we represent them as:

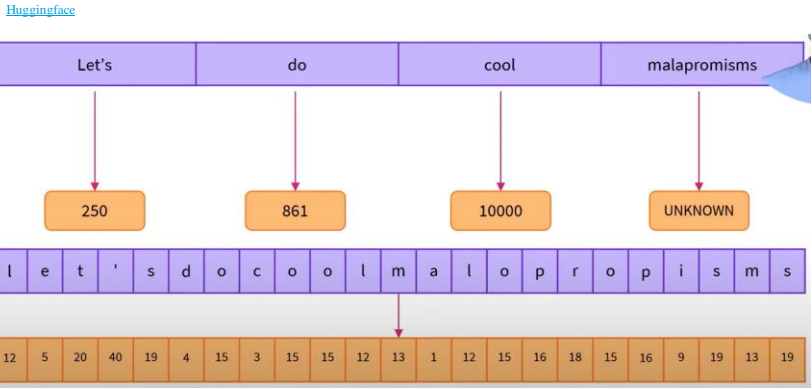
- A batch of sentences (i.e., batches)
- Each sentence is represented as a sequence of tokens (sequence length)
- Each token is represented as a vector (hidden size)



# Why do word level?



Issues?



Issues?

*Word-based tokenization*

- Very large vocabularies
- Large quantity of out-of-vocabulary tokens
- Loss of meaning across very similar words

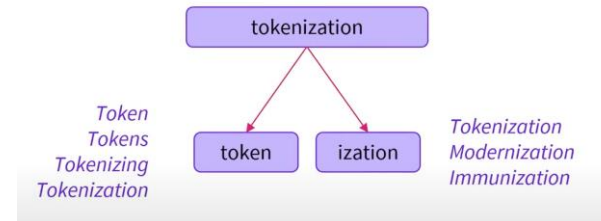
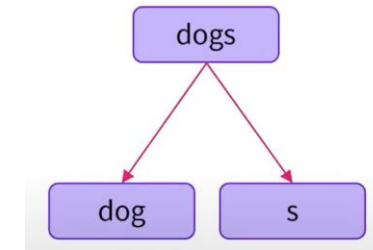
*Character-based tokenization*

- Very long sequences
- Less meaningful individual tokens

# Why do word level?

How about sub-words?

- Preserves word morphology.
- Can represent new words.
- Handles misspelling.
- Examples:



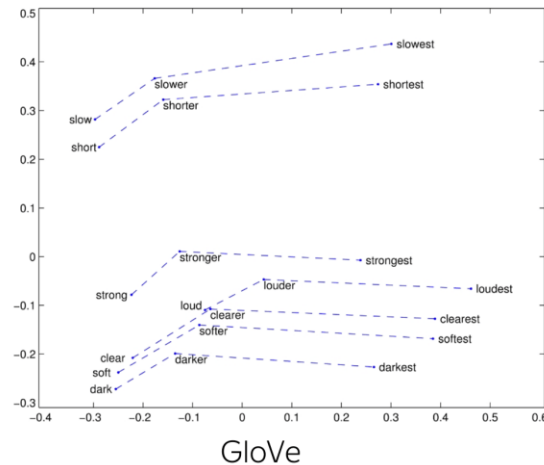
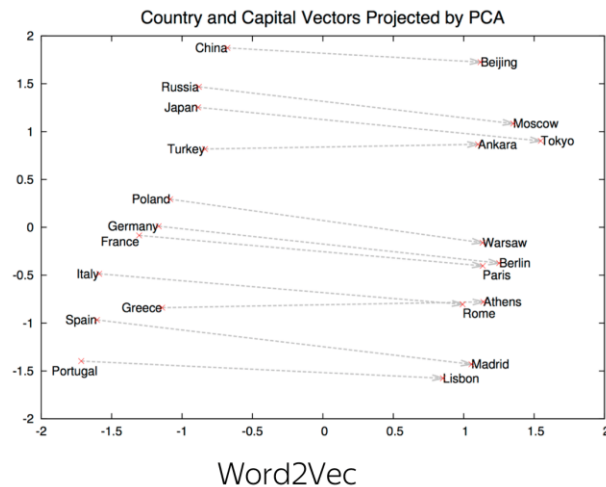
# Word embeddings

Word embeddings can also be used find directionality in the corpus:

- [Demo 1 \(semantics\)](#)
- [Demo 2 \(vector view\)](#)
- [Demo 3 \(dimensionality\)](#)

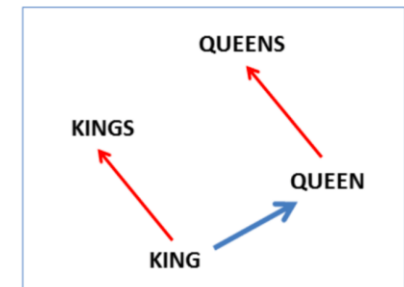
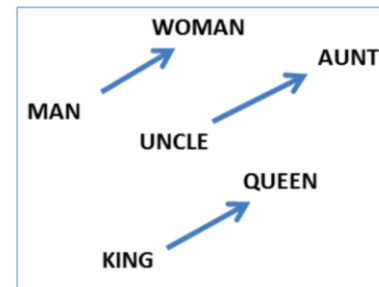
# Word embeddings

Word embeddings can also be used find directionality in the corpus.



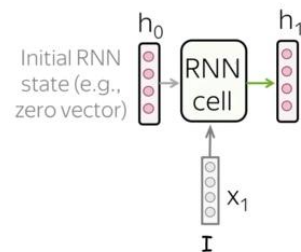
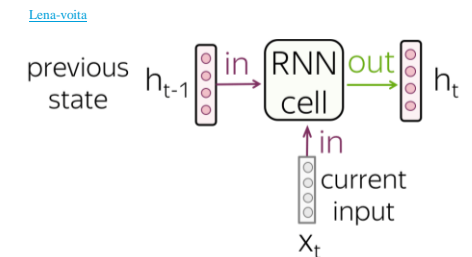
semantic:  $v(\text{king}) - v(\text{man}) + v(\text{woman}) \approx v(\text{queen})$

syntactic:  $v(\text{kings}) - v(\text{king}) + v(\text{queen}) \approx v(\text{queens})$



# Recurrent Neural Nets (RNNs)

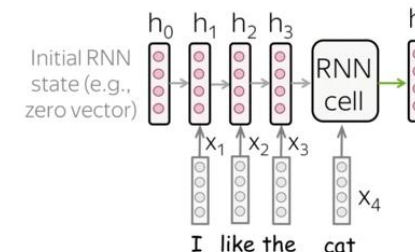
- Combines the embeddings of previous context and current word  
→ gives next word.



Get new state from RNN

Text: I like the cat on a mat <eos>  
↑  
we are here  
not read yet

.....

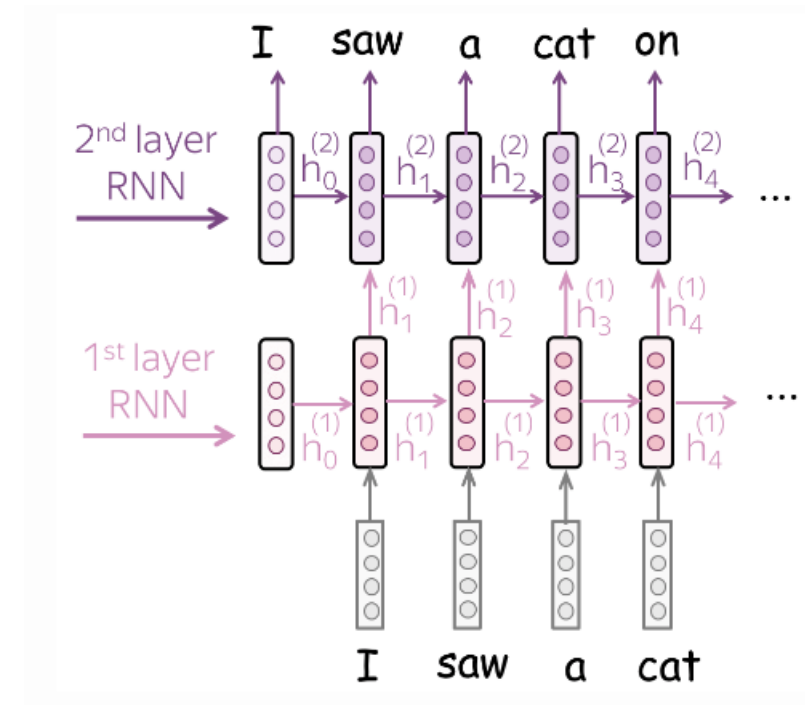


Get new state from RNN

Text: I like the cat on a mat <eos>  
↑  
we are here  
not read yet

# Recurrent Neural Nets (RNNs)

- We can add more layers and units per layer to increase complexity.



[Lena-volta](#)

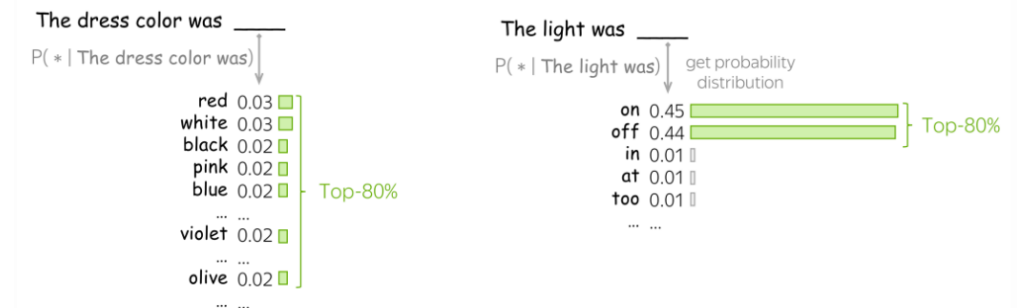
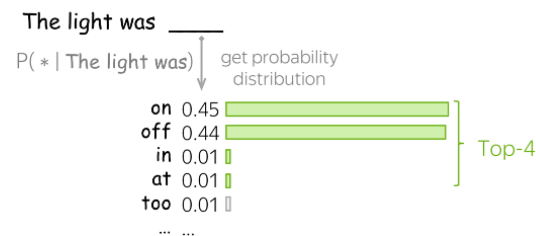
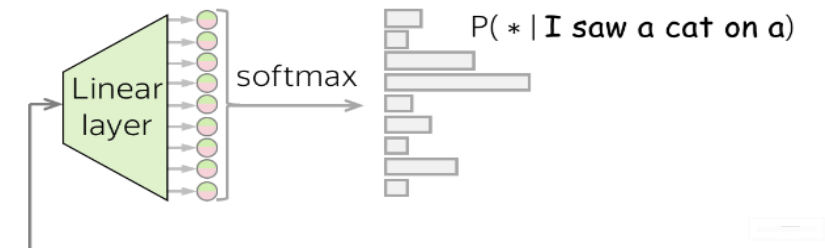


# Recurrent Neural Nets (RNNs)

[Demo](#)

# Sampling The Distribution

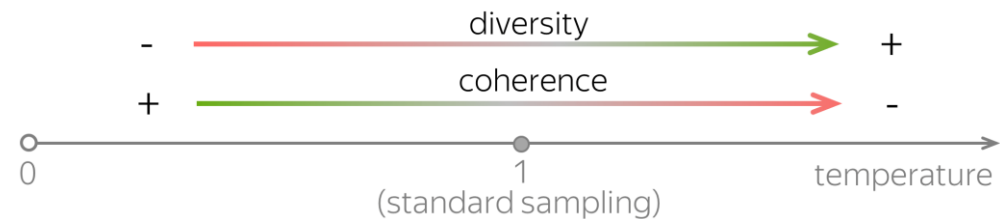
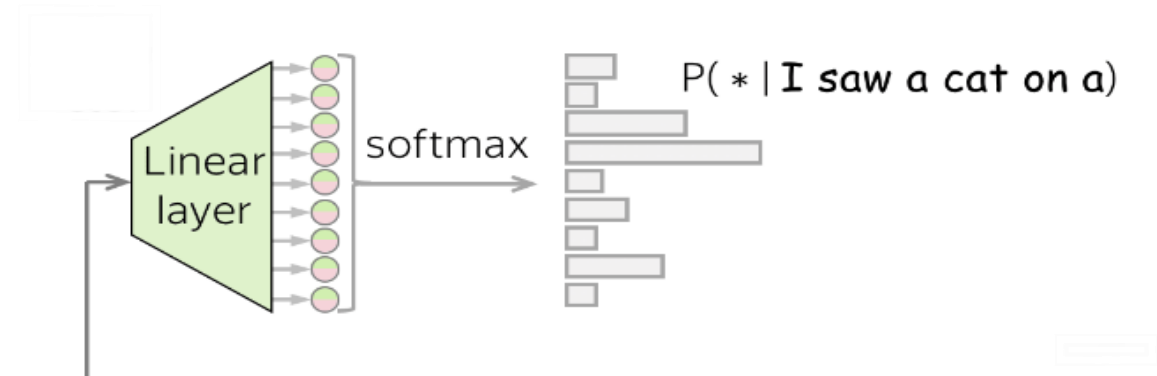
- Always take top probability?
  - That makes the model deterministic (no creativity).
- Alternative?
  - Top-k or top-p.



[Lena-vaita](#)

# Sampling The Distribution

- Some words have way higher probability than others.
- This can be manually tuned through **temperature**.
- [Demo](#)



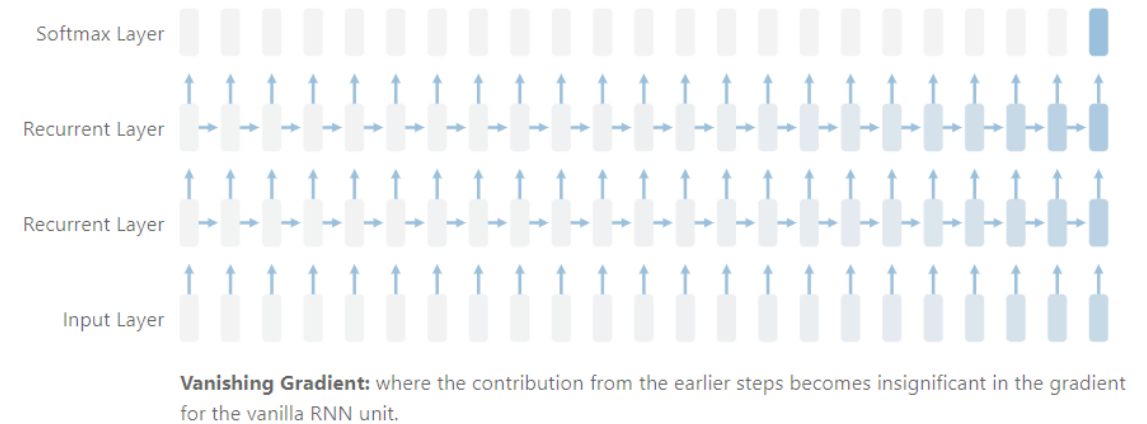
[Lena-volta](#)

# Measuring The Metric

- What are we looking for?
  - A model that is not surprised by the new text it seen.
- We use perplexity.
  - Takes values between 1 and number of possible tokens.
  - Smaller is better.
  - [Demo](#)

# RNNs (issues)

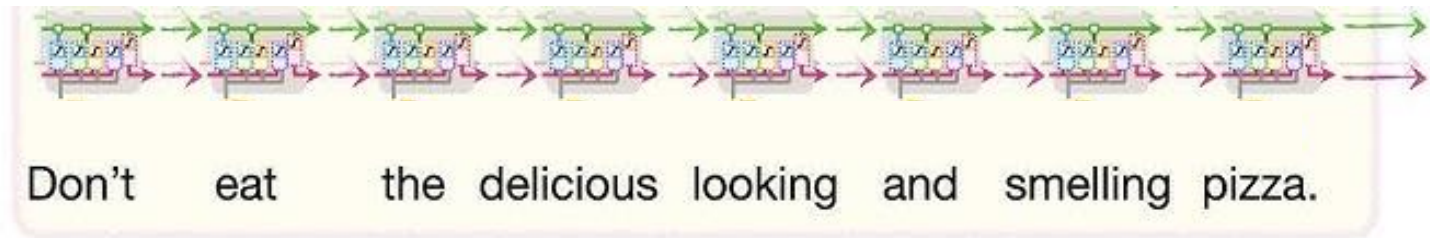
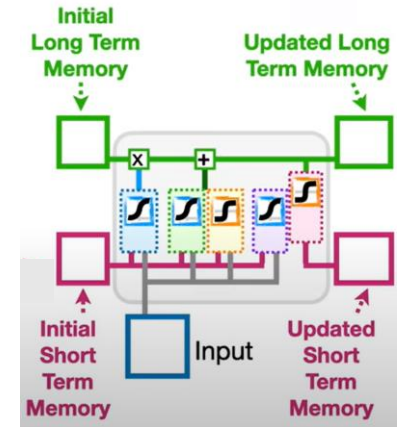
- Gradient becomes insignificant for long contexts
  - The network forgets early words...
  - It is called the “vanishing gradient” or “memorization” problem.
  - RNNs have an issue memorizing long contexts.



[distill.pub](#)

# Attempted solution: LSTM

- Instead of one representation, let's have two!
  - One for short-term memory
  - And one for long-term memory.
- Somewhat of an improvement.

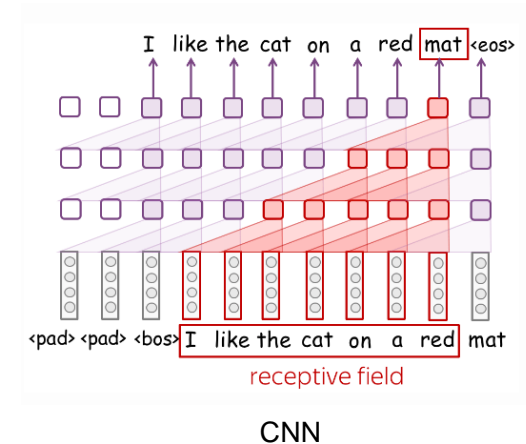


[StatQuest](#)

# There are other similar models

GRUs, CNNs...

But we will not talk about them here.



# Question for next week...

What limitations are inherent in traditional NLP models, and how might they be addressed?



# Assignment