

Practical Natural Language Processing

Julia Proskurnia
me@juliapro.xyz

July, 2019

About me, myself, I

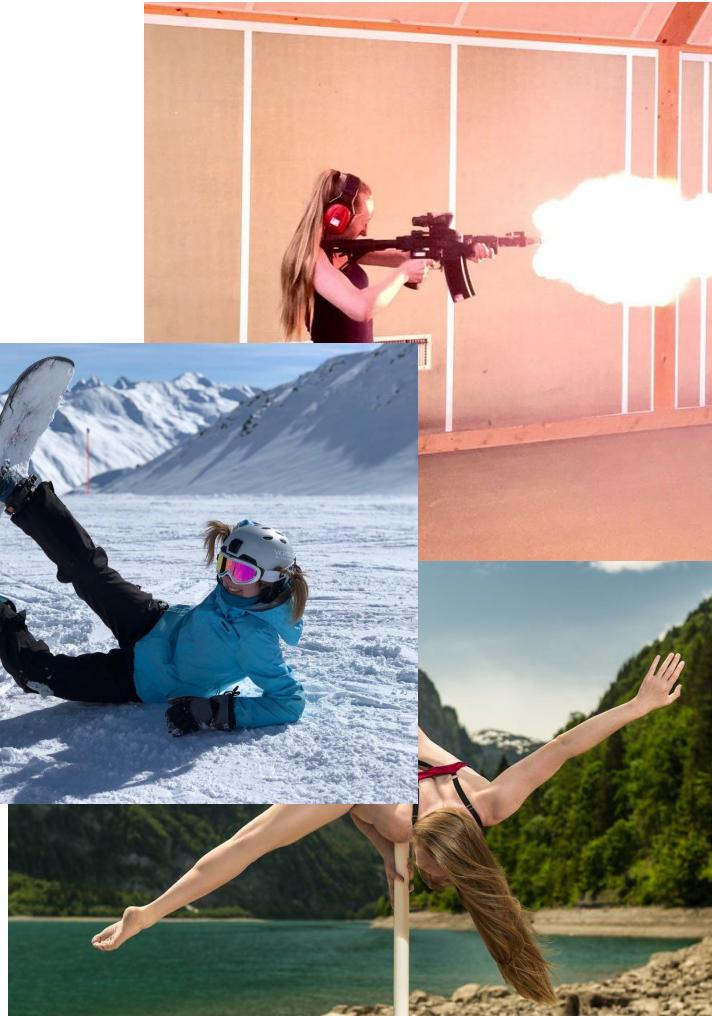
[BA, MS] NTUU KPI (Optimization, Neural Networks,
Time Series Prediction)

[MS] EMDC, Barcelona, Stockholm (Distributed
Systems)

[PhD] EPFL, Lausanne (Profiling, Modeling and
Facilitating Online Activism)

Google, Zurich (Gmail Smartcompose boilerplate
Personalization, Gmail Subject Smartcompose, etc.)

Google, Zurich (Next Generation Assistant)



Course Expectations?

- Would be more useful for people less familiar with NLP.
- Familiarity with Machine Learning concepts would help, but we will cover that briefly as well.

What this course is not about!

- Going over some fancy math for FEW particular algorithms
- Coding scalable solutions
- Go into the wild and write DNN architectures left and right :)
- Not solely Deep Learning in NLP

But...

- You would get some practical starting points on many of the NLP directions
- Lots of reference points
- Buzzwords :)

Overview for what to expect

- Intro to NLP problems
- Intro to the course setup
- Text preprocessing and annotations
- N-Gram extraction, TF-IDF
- Document classification
- Word Embeddings
- Document classification with embeddings
- Topic modeling, clustering
- Intro to entity extraction
- Intro to sequence modeling, deep learning (CNN, LSTM)

Introduction to NLP problems

Structure of the Introduction

- What is Natural Language
- Examples, applications, a few demos
- Typical NLP tasks
- Evaluation and metrics
- The source of data or data collection

Natural Language Processing (or NLP)

Main objective:

program **computers** to **process and analyze**
large amount of natural language data.

What is Natural Language?

"Natural" because evolved naturally, in contrast with Formal Languages.

Features of natural language:

1. Use of shared knowledge
 - Time *flies like* an arrow.
 - Play *rolling stones* on youtube
2. Concise, but ambiguous
 - Remove the stones from the cherries and put *them* in the pie.
 - The hunter shot the tiger; *his* wife too.
3. Unlimited expressive power
 - She was eating a fish *with* (bones, anger, some friends, a fork)
 - Congratulations! Great job!

Challenges for NLP

- Implicit -> allows conciseness (but potentially ambiguous)
- Unlimited expressive power -> flexible interpretation rules
 - (meanings != word as it is written)
- Why do we still understand each other? -> A lot of shared knowledge!

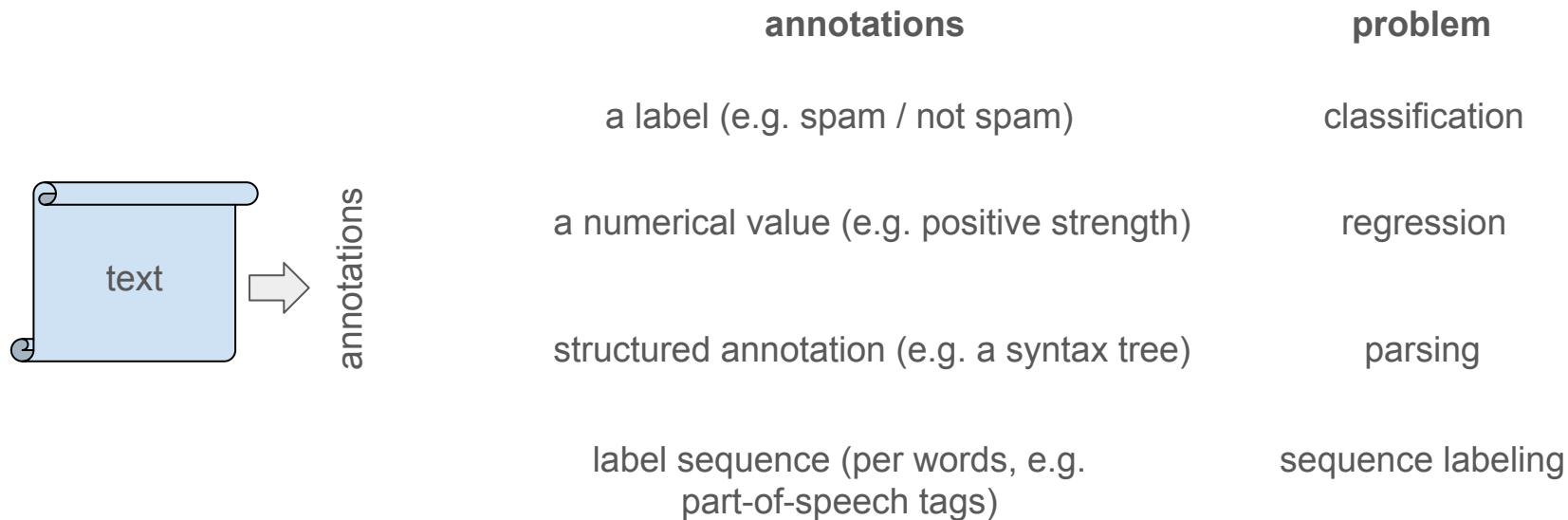
Moreover...

- Existing words (or linguistic units) change over time.
- New words are added.
- Language contains errors.
- Linguistic units denote an idea or even several.
- We have lots of common sense knowledge and understand hierarchical structures in language.

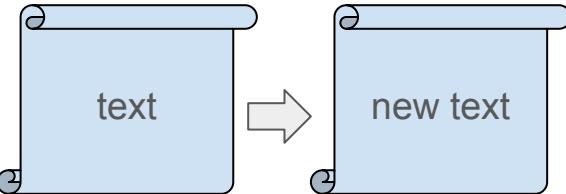
Challenges for NLP

- Full understanding of Natural Language is difficult
 - Similar to understanding / modeling the full world knowledge, needs full-capacity AI
- Natural languages is highly dimensional and sparse
 - There are millions of words and most of them are very rare
- Internationalization requires prior within-domain knowledge
 - Many languages with scarce resources

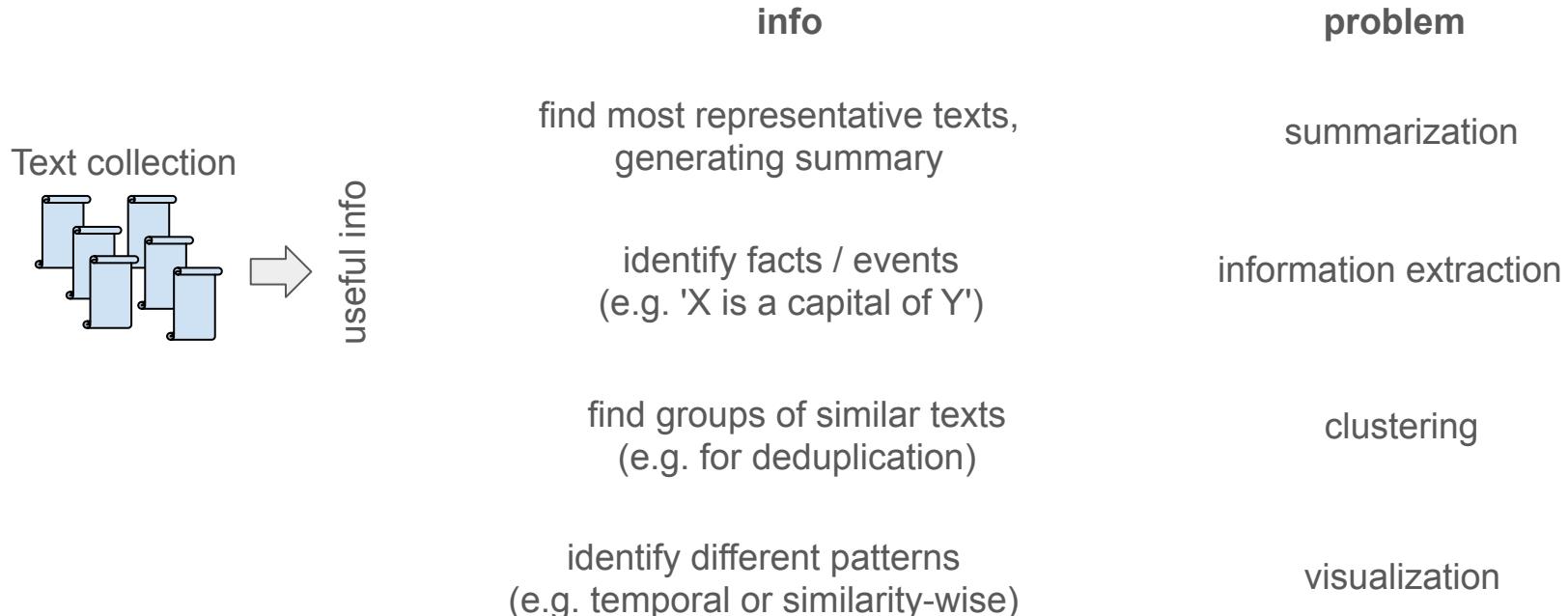
Text to annotations



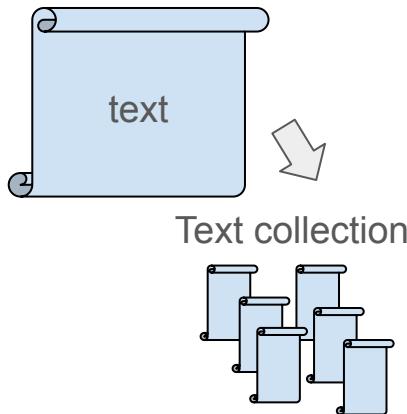
Text to another text

	new text	problem
	from one language to another	translation
	from question to answer	question-answering
	from text to its shortened version	text summarization
	text continuation (e.g. SmartCompose)	language modeling
	generating text rewrites	paraphrase generation

Text collection to useful info



Text to text collection



output

all relevant documents
(e.g. search results for a query)

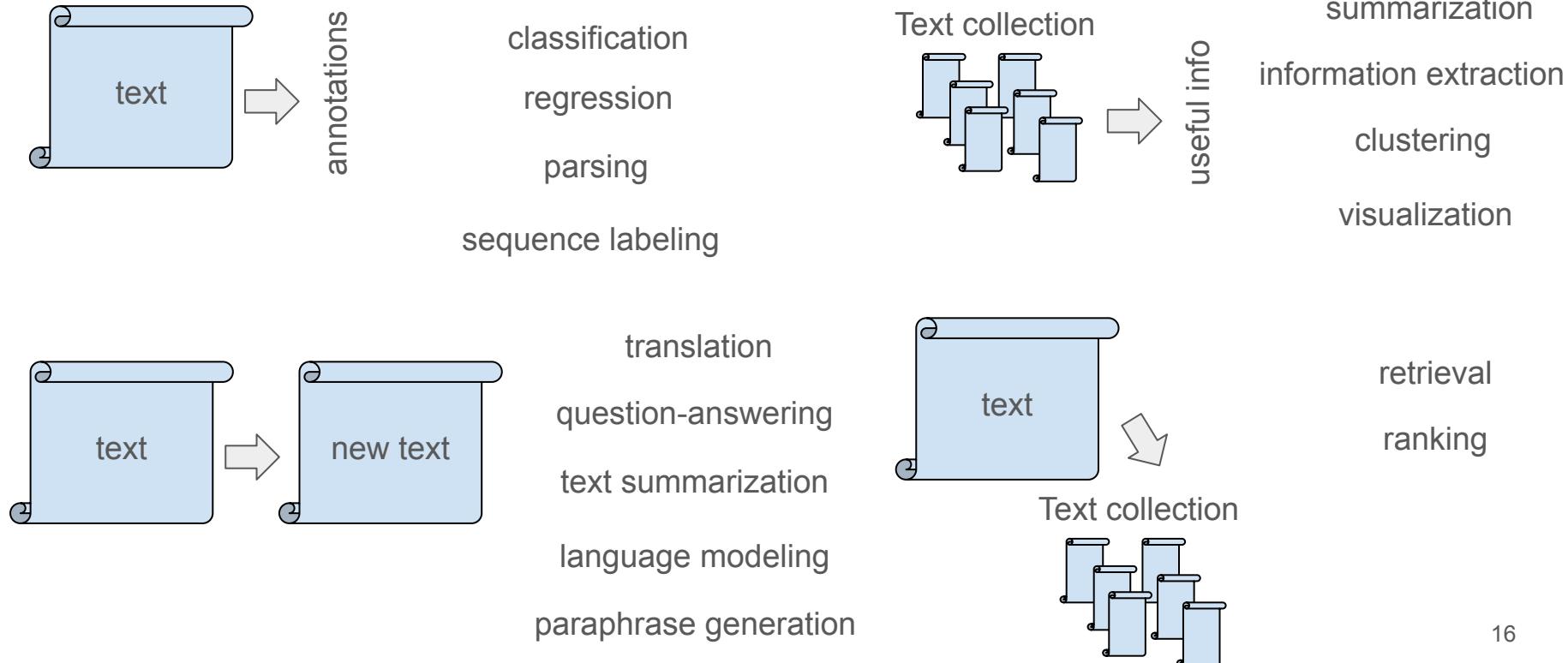
problem

retrieval

select a top doc from a collection
(e.g. knowledge panel)

ranking

Space of NLP problems (non-exhaustive!)



Applications and Examples

More examples in the colab...

Typical NLP tasks

- *Classification*: Classifying documents into particular categories (giving labels).
- *Regression*: Predict numerical values
- *Clustering*: Separating documents into non-overlapping subsets.
- *Ranking*: For a given input, rank documents according to some metric.
- *Association rule mining*: Infer likely associations patterns in data.
- *Structured output*: Parse trees, etc.
- *Sequence-to-sequence*: For a given input/source text, generate output annotations or another text.

Typical NLP Solution Flow

- 1 Formulate a problem statement
- 2 Collect the data
- 3 Get 'golden' or 'silver' answers
- 4 Search for prior resources
- 5 Train the model
- 6 Evaluate the quality
- 7 Iterate

Example: Machine Translation

sequence-to-sequence generation

within-language texts, aligned text pairs

golden: human-rated or generated,
silver: automatically aligned from Wikipedia

pre-trained text representations,
ready tools for text annotation

(many options)

accuracy, readability, ...

find better representations, more features,
better data, etc.

Data sources

Raw text data (Unlabeled/Non-annotated):

- Pay for it :)
- Crawl it from the Web.
 - Examples: [Scrappy](#), [wiki/google crawler](#)
 - IP addresses, proxies, throttling rates etc.
- Twitter (historical or livestream data (< 1% is free))
 - 1% historical tweets from [archive.org](#).
 - 1% via Twitter API.
 - Keywords (up to 2K) that would match tweets in real time 1% total stream of messages.
- News media ([News Archive](#), [GDELT](#))
- Wikipedia ([Wikipedia dumps](#))

Annotation of the data and pitfalls

- Annotate and/or even generate your data using CrowdSourcing
 - a. [Crowdflower](#)
 - b. [MechanicalTurk](#), etc.
- Annotate using auxiliary lexical resources
 - a. [LIWC](#) a tool that analyses your textual input on the presence of various
 - i. "shades" - informal speech; syntactic structures; affect, social words; cognitive, perpetual, biological processes; relativity; personal concerns, etc.

Biases and pitfalls

- Biases
- Not representative
- Noisy data
- Incomplete data
- Incorrect data
- Missing data, etc.

Course practicum setup

- We will use colab.research.com.
- We will focus on the open-source libraries for NLP to give you easier explanations.

Actions:

- Either
 - Open colabs etc.
 - OR open the copy of the colab with jupyter notebook [RECOMMENDED if you want to try out dependency parsing]

Preprocessing and data cleaning

Structure of the Preprocessing and Data Cleaning

- Tokenization + Sentence Splitting
- Normalization / Text cleaning
 - Stemming/lemmatization, stop-words, punctuations, emojis, domain specifics, etc.
- Text Annotation
 - Part-Of-Speech Tagging
 - Parse Trees

Data description: Tweets about natural disasters

DataFrame name: **course_data**

- ‘label’ - whether the tweet is relevant to natural disasters or not
- ‘text’ - raw tweet strings
- ‘text_tokenized’ tokenized input texts
- ‘text_tokenized_cleaned’ tokenized and normalized text

document_embeddings - Text representation with word embedding

document_term_matrix - Bag-of-words representation

corpus_tokens - List with list of cleaned tokens of each document

POS exercise in colab

Stemming and Lemmatization

- **Stemming** is typically faster, more naive word shortening. (crude affix chopping)
It is language dependent.

affect	amus	close	© Lightweight NLP
affect	amuse	close	close
affection	amused	closed	closed
affected	amusement	closely	closely
affecting	amusements	closing	closing
affection	amusing	grate	grate
affections			grateful
affects			gratefully

- **Lemmatization** take deeper approach to extracting lemmas of the words (part of speech, language specific dictionary, etc.).

Lemmatization implies doing “proper” reduction to dictionary headword form

am, are, is → be

car, cars, car's, cars' → car

Porter Stemmer

Step 1a

sses	→ ss	caresses	→ caress
ies	→ i	ponies	→ poni
ss	→ ss	caress	→ caress
s	→ Ø	cats	→ cat

Step 1b

(*v*)ing	→ Ø	walking	→ walk
		sing	→ sing
(*v*)ed	→ Ø	plastered	→ plaster

...

Step 2 (for long stems)

ational	→ ate	relational	→ relate
izer	→ ize	digitizer	→ digitize
ator	→ ate	operator	→ operate
...			

Step 3 (for longer stems)

al	→ Ø	revival	→ reviv
able	→ Ø	adjustable	→ adjust
ate	→ Ø	activate	→ activ
...			

Sample stemmings

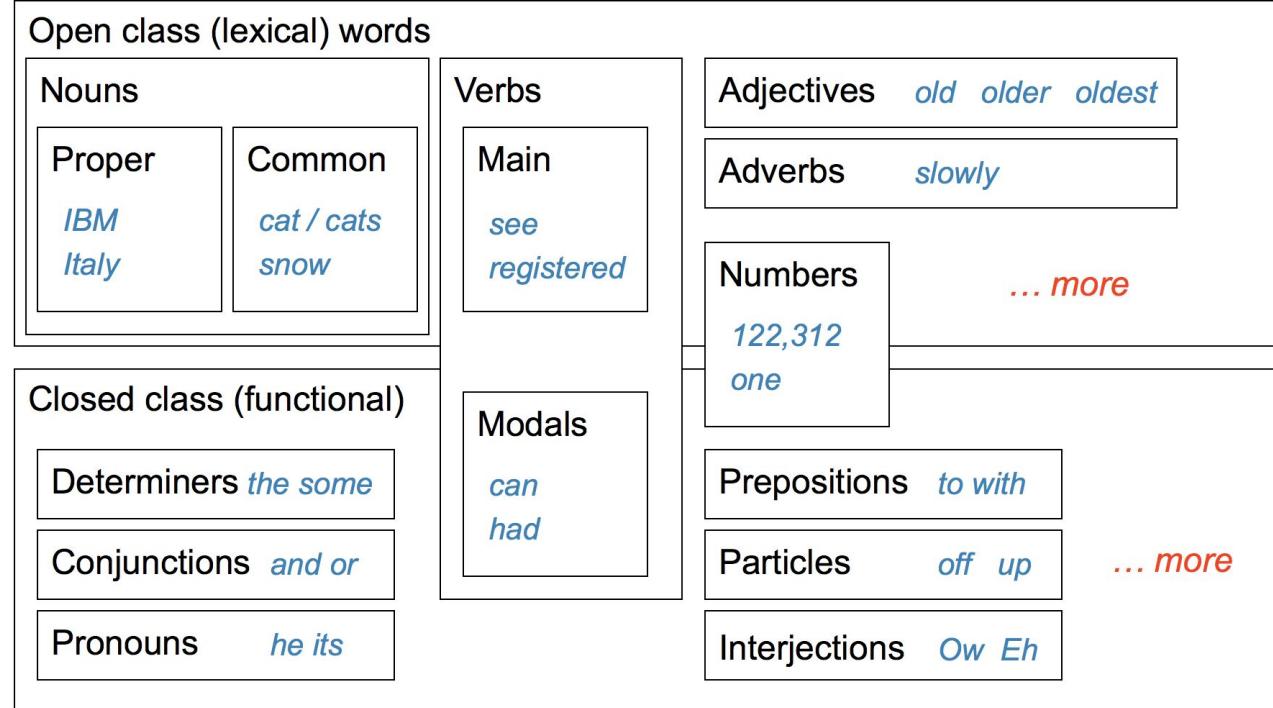
Sample text: Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Lovins stemmer: such an analys can reve featur that ar not eas vis from th vari in th individu gen and can lead to a pictur of expres that is mor biolog transpar and acces to interpre

Porter stemmer: such an analysi can reveal featur that ar not easili visibl from the variat in the individu gene and can lead to a pictur of express that is more biolog transpar and access to interpret

Paice stemmer: such an analys can rev feat that are not easy vis from the vary in the individ gen and can lead to a pict of express that is mor biolog transp and access to interpret

Part-of-Speech Tagging



Part-of-Speech

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	+,%,&
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	\$
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	#
PDT	Predeterminer	<i>all, both</i>	"	Left quote	' or "
POS	Possessive ending	's	"	Right quote	' or "
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	[, (, {, <
PRP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis],), }, >
RB	Adverb	<i>quickly, never</i>	,	Comma	,
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	. ! ?
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	: ; ... --
RP	Particle	<i>up, off</i>			

POS Tagging

- Words often have more than one POS: *back*
 - The back door = JJ
 - On my back = NN
 - Win the voters back = RB
 - Promised to back the bill = VB
- The POS tagging problem is to determine the POS tag for a particular instance of a word.
 - I know *that* he is honest = IN
 - Yes, *that* play was nice = DT
 - You can't go *that* far = RB

40% of the word tokens are ambiguous

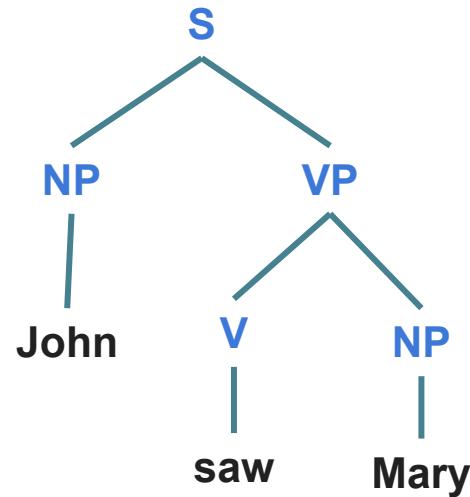
POS Tagging

- What are the main sources of information for POS tagging?
 - Knowledge of neighboring words
 - Bill saw that man yesterday
 - NNP NN DT NN NN
 - VB VB(D) IN VB NN
 - Knowledge of word probabilities
 - *man* is rarely used as a verb....
- The latter proves the most useful, but the former also helps
 - Word the: the → DT
 - Lowercased word Importantly: importantly → RB
 - Prefixes unfathomable: un- → JJ
 - Suffixes Importantly: -ly → RB
 - Capitalization Meridian: CAP → NNP
 - Word shapes 35-year: d-x → JJ

Parse Trees

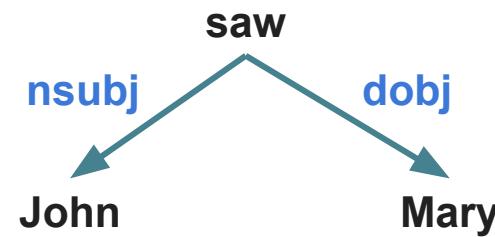
Constituency trees:

structure of nested phrases/constituents



Dependency trees:

actual dependencies between words

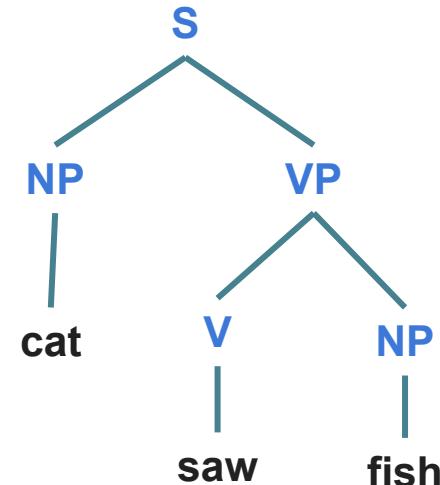


Parse Trees: Why to use?

- To find the root / main intent of the query / sentence.
- To find the modifiers / arguments of known frames (e.g. set alarm WHEN).
- Build event-like feature representations: nsubj('<someone>', 'saw').
- Extract smarter n-grams based on parse trees.

Constituency Parse Trees: Build with Probabilistic Context-Free Grammars

$S \rightarrow NP VP$	1.0	$N \rightarrow$	cat	0.4
$NP \rightarrow NP NP$	0.1	$N \rightarrow$	fish	0.4
$NP \rightarrow NP PP$	0.3	$N \rightarrow$	saw	0.2
$NP \rightarrow N$	0.6	$V \rightarrow$	meets	0.4
$PP \rightarrow P NP$	1.0	$V \rightarrow$	saw	0.2
$VP \rightarrow V NP$	0.7	$V \rightarrow$	greets	0.4
$VP \rightarrow VP PP$	0.3			



Text representation

Text representations

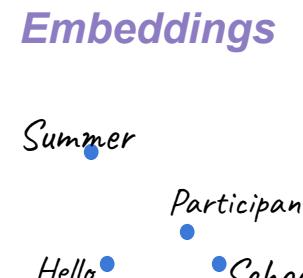
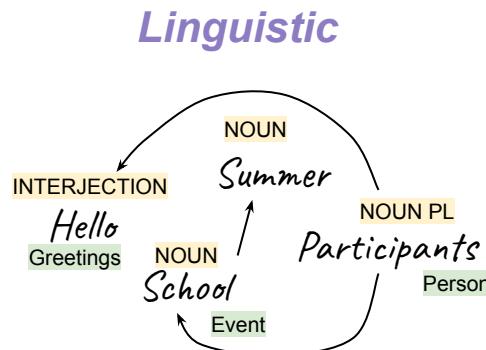
Idea: Represent the text so that the machine can perform its processing.

The text representations are usually inputs for machine-learning models.

Different approaches*:

Lexical

Hello Summer
School Participants



* (that we will discuss)

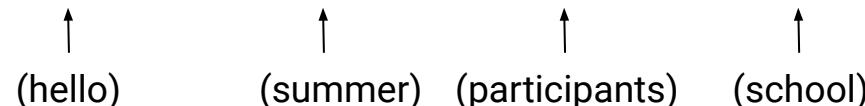
Bag-of-words (BOW) representation

Idea: Split the text into tokens (~words), extract their counts, and build their vector.

Input text: "Hello, Summer school participants!"

Set of normalized words: hello, summer, school, participants

As a **count vector:** [0 0 0 1 0 0 0 ... 0 1 0 0 ... 0 0 1 0 ... 0 0 0 1 0 ... 0]



TF-IDF, Conditional Probability, etc.

Weighting tokens/features:

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right) \quad P(\text{word}_i | \text{document}_j) / P(\text{word}_i | \text{all_documents})$$

$tf_{i,j}$ = number of occurrences of i in j

df_i = number of documents containing i

N = total number of documents

Removing noise:

- Remove tokens that are very frequent
- Remove tokens that have very low TF-IDF

Compare frequency distributions of 2 independent datasets and diminish the intersection.

Bag-of-words: Pros & Cons

Pros:

- Very simple, no need for special annotations

Cons:

- Can have too many low-occurring words
- Breaking collocations, e.g. "happy hour", "new york"
- Breaking syntax structure, e.g. "John saw Mary" or "Mary saw John"
- Purely lexical, no notion of semantics:
 - Different representations for similar concepts, e.g. synonyms ("phone" vs. "telephone")
 - Same representation of ambiguous words, e.g. "buy jaguar car" vs. "jaguar in nature"

N-grams representation

Idea: Take all longer consecutive token sequences as a vector.

Input text: "Hello, Summer school participants!"

Normalized text: "hello gwe school participants"

Extracted n-grams:

- unigrams: hello, summer, school, participants
- bigrams: hello summer, summer school, school participants
- trigrams: hello summer school, summer school participants
- 4-grams: hello summer school participants

N-grams: Pros & Cons

Pros:

- Can capture some structure and connotations
- Quite powerful

Cons:

- Aggravated data sparsity
 - Longer n-grams appear even less frequently than particular unigrams
- Correlated features
 - BUT: still works in practice
- Plus other problems for lexical-level representation

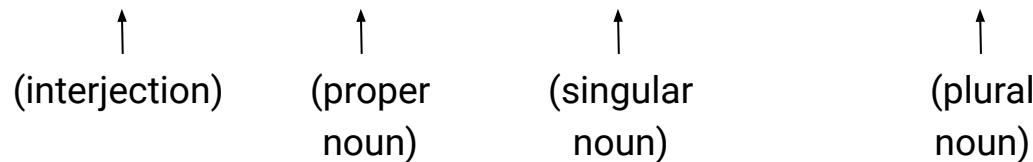
Linguistic representation: POS (Part-of-Speech) Tags

Idea: Capture better word meanings by assigning a POS tag to each word.

Why? Some words can have several POS tags: back (JJ/adjective) vs. back (NN/noun)

Input text: "Hello, Summer school participants!"

Annotated text: Hello/UH Summer/NNP schjool/NN participants/NNS

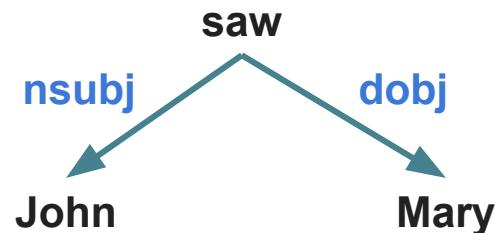


Linguistic representation: Syntax parse trees

Idea: Capture the syntactic structure of the sentence in a tree format.

Dependency trees:

actual dependencies between words



Why?

- Can identify sentence root
- Can extract smarter lexical features:
 - hello participants from
"Hello Summer school participants"

Linguistic representation: Semantic role labels

Why? Get more semantic text understanding in order for machine to act on it.

Input query: "play Rolling Stones on youtube"

Role labels annotations: "play/ACTION rolling stones/ARTIST on youtube/APP"

Current state: Application-specific, usually focusing on particular label sets.

Linguistic representations: Pros & Cons

Pros:

- Can link to more formal language representation (e.g. language grammar)
- Captures more meaning (better disambiguation & connotation structure)

Cons:

- Each layer of annotations adds errors and computational complexity

Entity representation

Idea: Link the lexical forms of entities mentioned in the text to their actual identities (entities defined by non-ambiguous ids).

Examples:

"buy jaguar"  /wordnet/jaguar_car_brand

"jaguar in the nature"  /wordnet/jaguar_animal

Final text representation: a vector of detected entities.

Entity representation: Pros & Cons

Pros:

- Can disambiguate lexical forms into entities, e.g. "jaguar" in context
- Can match different lexical forms into the same entity, e.g. "phone" and "telephone"
- Has information about entity relations, e.g. "/knowledgebase/iphone" is produced by
"/knowledgebase/apple"

Cons:

- Introduces the room for interpretation mistakes
- Only entities! Does not represent (yet) intents,
common adjectives, nouns, etc.

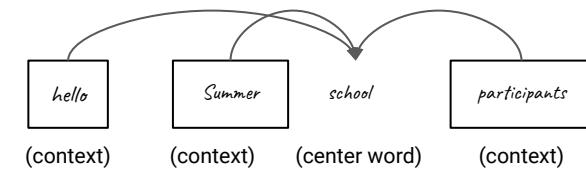
*Hello
Summer
School
Participants*

Word embeddings: Dense semantic representation

Idea: Each word corresponds to a vector in multi-dimensional space.

Desired property:

Similar words are close to each other in that space, while unrelated words are further apart.



Sentence representation: Represent an input text as one multi-dimensional vector.

- Average embedding of all input words
- Train models to represent sequences (Deep Learning)

Embeddings representation: Pros & Cons

Pros:

- Allow direct vector operations (dense representation)
- Can capture text semantics
- Really powerful at data generalizations

Cons:

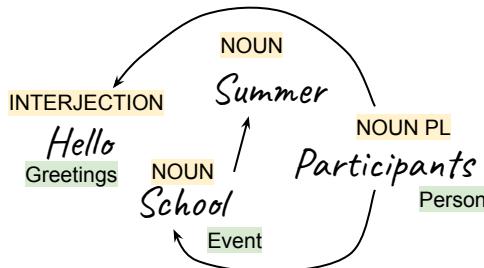
- Less interpretable
- Require tons of data for training
 - BUT: Pre-trained models are available!
- Computationally expensive
 - BUT: Often can be precomputed.

Text representations: Summary

Lexical

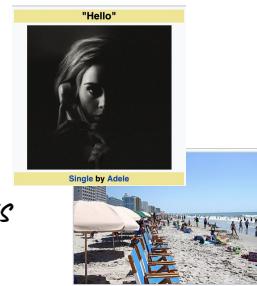
Hello Summer
School Participants

Linguistic

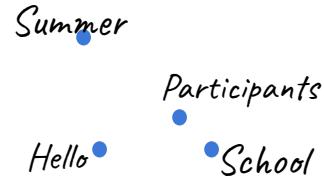


Entities

Hello Summer
School Participants



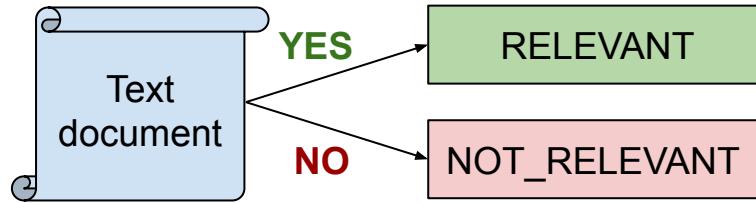
Embeddings



These representations are only the input for the models. The models can build even better representations of the full text by further processing them.

Real text problems: Document classification

Document Classification

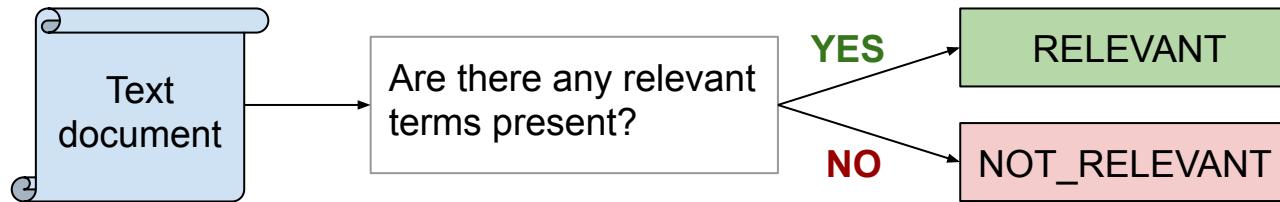


Given a set of labels/classes, need to assign labels to each document.

Examples:

- Spam filtering,
- sentiment analysis,
- documents retrieval (finding relevant documents),
- new topic classification

Rule-based classification

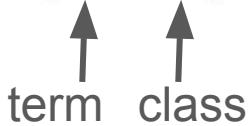


How to find which terms/words are relevant and which are not?



Building lexicons

Pointwise Mutual Information (PMI)

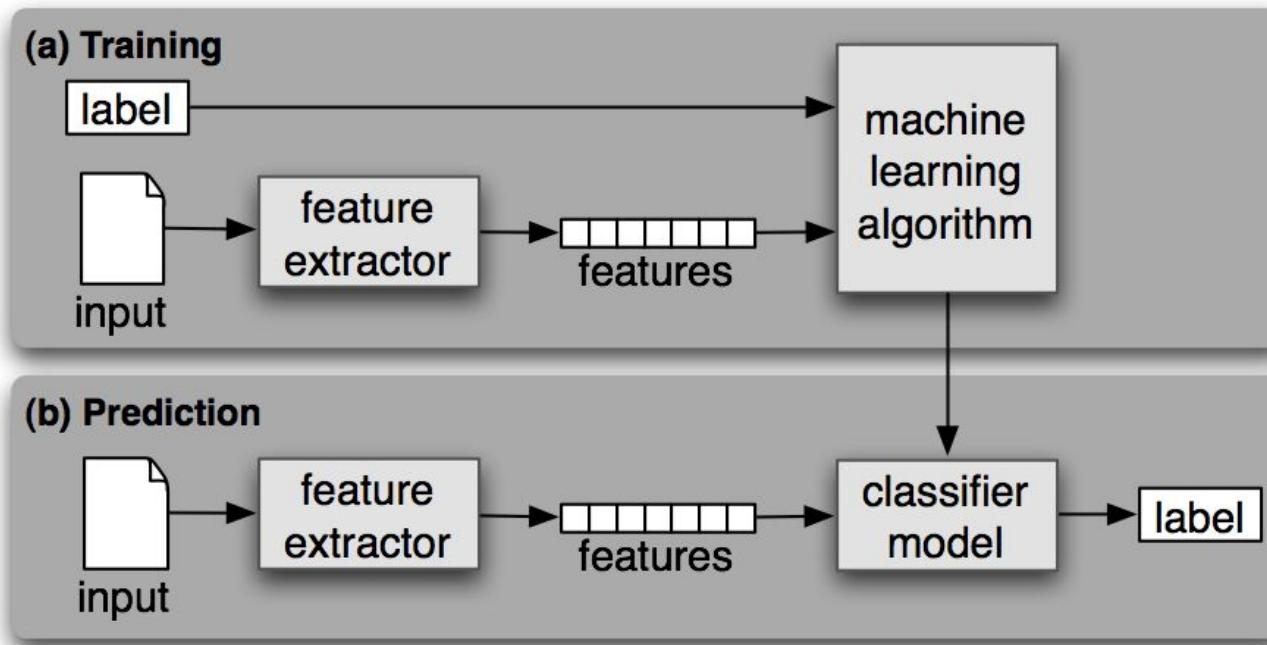
$$PMI(t, c) = \log_2 \frac{P(t, c)}{P(t)P(c)} = \log_2 \frac{P(t|c)}{P(t)}$$


The diagram shows two arrows pointing upwards from the words "term" and "class" to the variables "t" and "c" in the PMI formula. The arrow from "term" points to the first "t", and the arrow from "class" points to the first "c".

$$relatedness_{PMI}(t) = PMI(t, related) - PMI(t, not\ related) =$$

Document Classification

Supervised document classification



Supervised document classification: Common models

Generative models

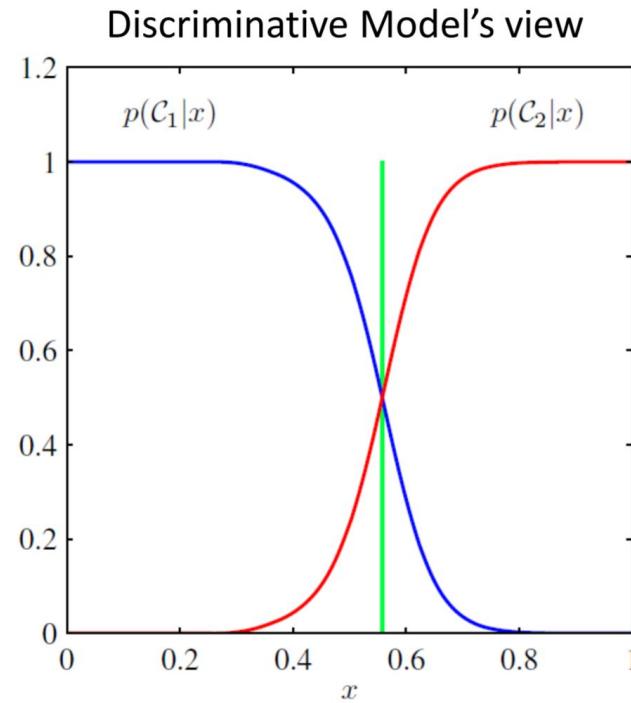
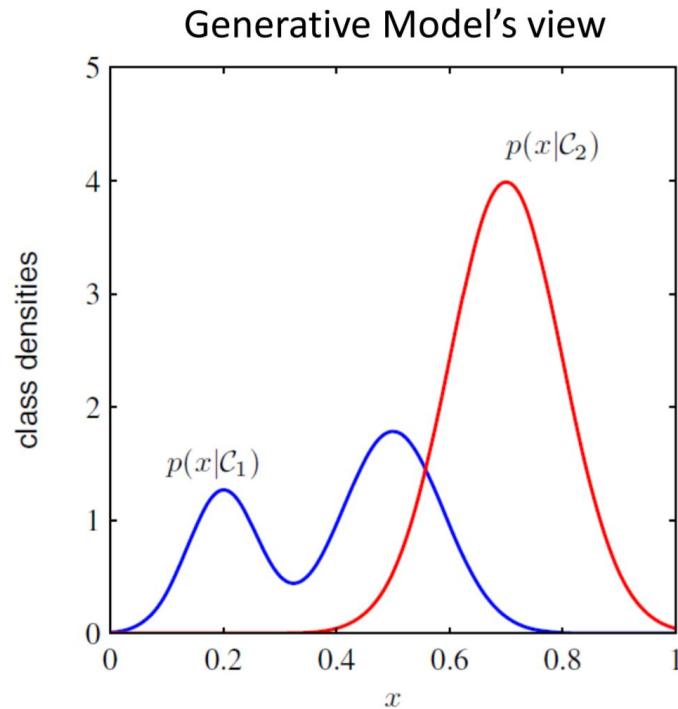
- Try to generate the observed data from the hidden structure (labels y)
- Model joint probability of $p(X, y)$ and try to maximize joint likelihood
- Example: Naive Bayes

Discriminative models

- Directly estimate a decision rule/boundary
- Model conditional probability $p(y | X)$ and try to maximize conditional likelihood (given the data X)
- Example: Logistic Regression

Generative vs. Discriminative

- Binary classification as an example



© Whole set of slides here:

http://www.cs.virginia.edu/~hw5x/Course/TextMining-2018Spring/_site/docs/PDFs/TextCategorization.pdf

Naive Bayes

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

MAP is “maximum a posteriori” = most likely class

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

Dropping the denominator

Multinomial Naïve Bayes

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

Document d
represented as
features
 $x_1..x_n$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

Multinomial Naive Bayes

In such scenario, each document is assumed to be generated as follows.

1. Class y is generated according to $P(y)$.
2. X is generated by sequentially picking words from V (vocabulary) with replacement.
Each word is picked with probability $P(w_j|y)$.

For example, the probability to generate a document $x = w_{j1} \dots w_{jL}$ is

$$P(x|y) = \prod_{l=1}^L P(w_{jl}|y) = \prod_{j=1}^J P(w_j|y)^{n_j(x)}$$

Linear Regression VS Logistic Regression

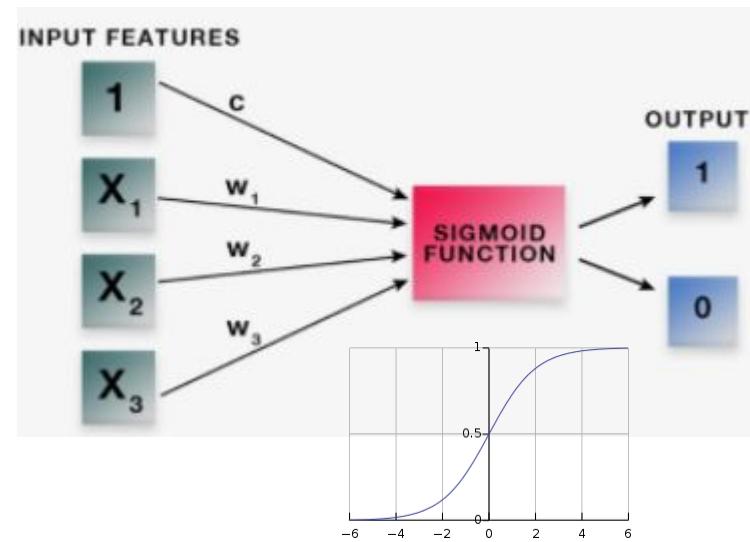
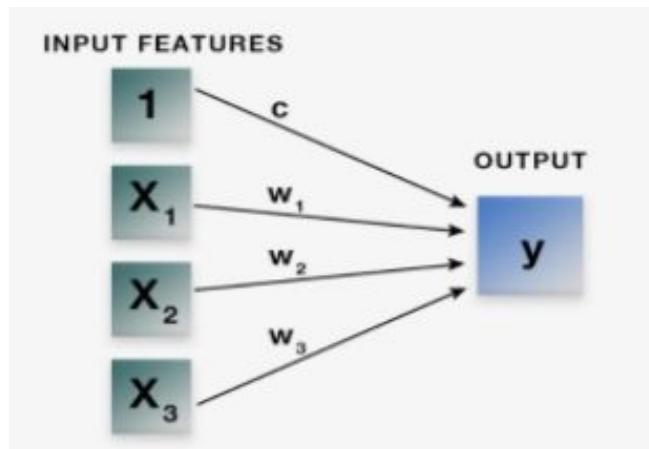


Image source: <https://blog.quantinsti.com/machine-learning-logistic-regression-python/>

Classification: Evaluation Metrics

	Actually correct	Actually incorrect
Selected by the algo	True positives	False positives
Not selected by the algo	False negatives	True negatives

Accuracy: % of items classified correctly = $(tp + tn) / \text{all}$

Precision: % of selected items that are correct = $tp / (tp+fp)$

Recall: % of correct items that are selected = $tp / (tp + fn)$

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

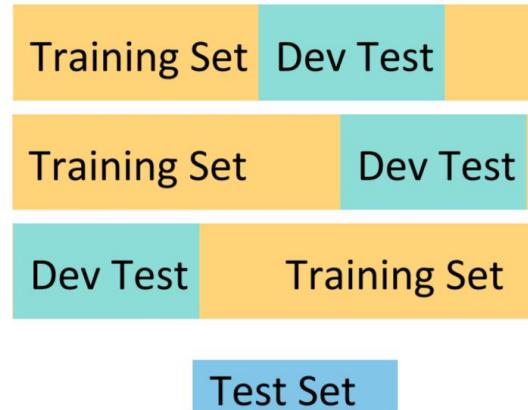
Classification: Common evaluation setup

Training set

Development Test Set

Test Set

- Metric: P/R/F1 or Accuracy
- Unseen test set
 - avoid overfitting ('tuning to the test set')
 - more conservative estimate of performance
- Cross-validation over multiple splits
 - Handle sampling errors from different datasets
 - Pool results over each split
 - Compute pooled dev set performance



Typical approaches to classification

Type	Input	Idea
Rule-based	Explicit linguistic patterns, lexicons	Apply given rules
Supervised	High-quality annotated data	Learn machine-learning classifier based on annotated examples
Semi-supervised	Annotated & non-annotated data	Use more data to find better feature representation, learn more patterns, or mine new training examples
Distant supervision	Pseudo-annotated data based on some heuristics	Use heuristics to get more training data (often of lower quality or biased)

Text Representation: Word Embeddings

Word embeddings is a set of language modeling and feature learning techniques.

Key idea - dimensionality reduction that captures semantic/morphological/contextual/hierarchical/etc.

1. preserving the morphological structure (subword information, etc.);
2. word context representation;
3. global corpus statistics;
4. word hierarchy as in WordNet;
5. relationship between documents and the terms they contain.

Typical Word Represent with one-hot vectors

In traditional NLP, we regard words as discrete symbols:

hotel, conference, motel

Means one 1, the rest 0s

Words can be represented by one-hot vectors:

motel = [0 0 0 0 0 0 0 0 0 1 0 0 0]

hotel = [0 0 0 0 0 0 1 0 0 0 0 0 0]

Vector dimension = number of words in vocabulary (e.g. 500,000)

Drawbacks of the one-hot representation

Example: in web search, if user searches for “Seattle motel”, we would like to match documents containing “Seattle hotel”.

But:

```
motel = [0 0 0 0 0 0 0 0 1 0 0 0 0]  
hotel = [0 0 0 0 0 0 1 0 0 0 0 0 0]
```

These two vectors are orthogonal.

There is no natural notion of **similarity** for one-hot vectors!

Solution:

- Could rely on WordNet’s list of synonyms to get similarity?

How can we introduce the meaning to the machine?

Common solution: Use e.g. **WordNet**, a resource containing lists of **synonym sets** and **hypercnyms** ("is a" relationships).

e.g. synonym sets containing "good":

```
from nltk.corpus import wordnet as wn
for synset in wn.synsets("good"):
    print "%s" % synset.pos(),
    print ", ".join([l.name() for l in synset.lemmas()])
```

```
(adj) full, good
(adj) estimable, good, honorable, respectable
(adj) beneficial, good
(adj) good, just, upright
(adj) adept, expert, good, practiced,
proficient, skillful
(adj) dear, good, near
(adj) good, right, ripe
...
(adv) well, good
(adv) thoroughly, soundly, good
(n) good, goodness
(n) commodity, trade good, good
```

e.g. hypernyms of "panda":

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")
hyper = lambda s: s.hypernyms()
list(pandaclosure(hyper))
```

```
[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]
```

Potential issues with the ontology, knowledge base solutions

- Great as a resource but missing nuance
 - e.g. “proficient” is listed as a synonym for “good”. This is only correct in some contexts.
- Missing new meanings of words
 - e.g. wicked, badass, nifty, wizard, genius, ninja, bombest
 - Impossible to keep up-to-date!
- Subjective
- Requires human labor to create and adapt
- Hard to compute accurate word similarity →

Maybe represent words by their contexts?

- Core idea: A word's meaning is given by the words that frequently appear close-by
 - “*You shall know a word by the company it keeps*” (J. R. Firth 1957: 11)
 - One of the most successful ideas of modern statistical NLP!
- When a word w appears in a text, its **context** is the set of words that appear nearby (within a fixed-size window).
- Use the many contexts of w to build up a representation of w

...government debt problems turning into **banking** crises as happened in 2009...

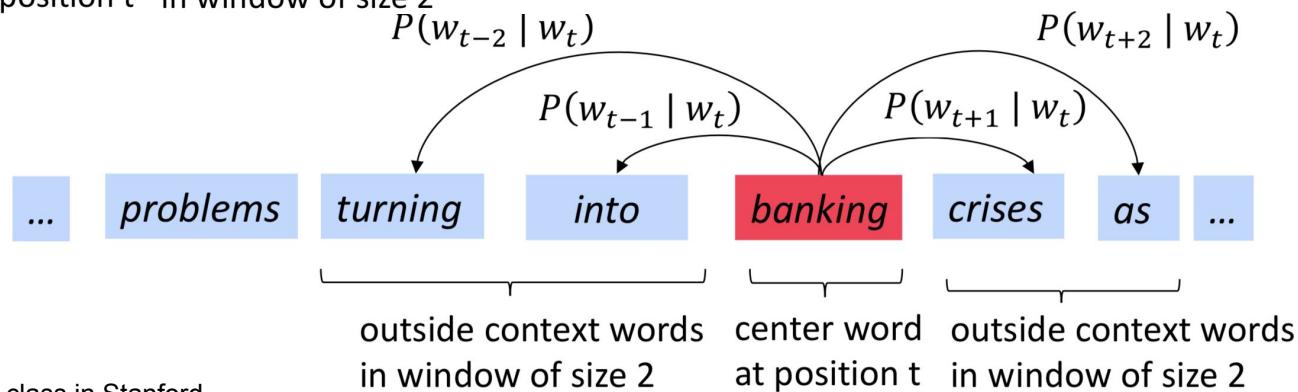
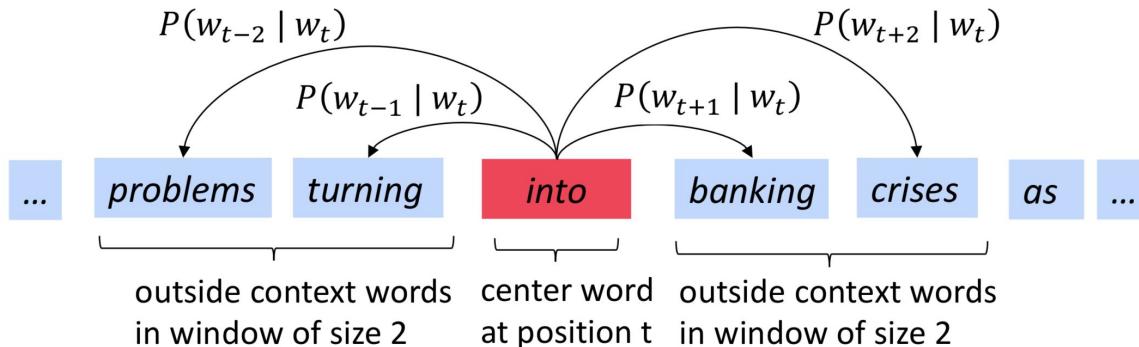
...saying that Europe needs unified **banking** regulation to replace the hodgepodge...

...India has just given its **banking** system a shot in the arm...

Word2Vec (Mikolov et al. 2015)

- We have a large corpus of text
- Every word in a fixed vocabulary is represented by a vector
- Go through each position t in the text, which has a center word c and context (“outside”) words o
- Use the similarity of the word vectors for c and o to calculate the probability of o given c (or vice versa)
- Keep adjusting the word vectors to maximize this probability

Given a center words compute $P(\dots)$ of its context



Word2Vec: Objective Function

For each position $t = 1, \dots, T$, predict context words within a window of fixed size m , given center word w_j .

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

θ is all variables
to be optimized

The **objective function** $J(\theta)$ is the (average) negative log likelihood:

$$\text{minimize } J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

Word2Vec: How to calculate $p(\text{neighbors} \mid \text{center})$?

We will *use two vectors per word w:*

- v_w when w is a center word
- u_w when w is a context word

Then for a center word c and a context word o :

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Dot product compares similarity of o and c .
Larger dot product = larger probability

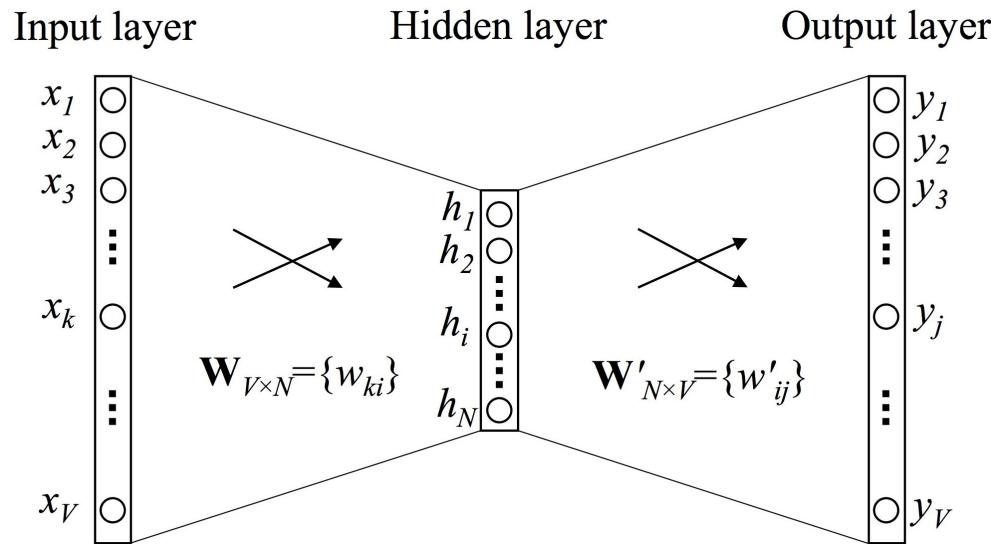
After taking exponent,
normalize over entire vocabulary

Word2Vec: Parameters of the model

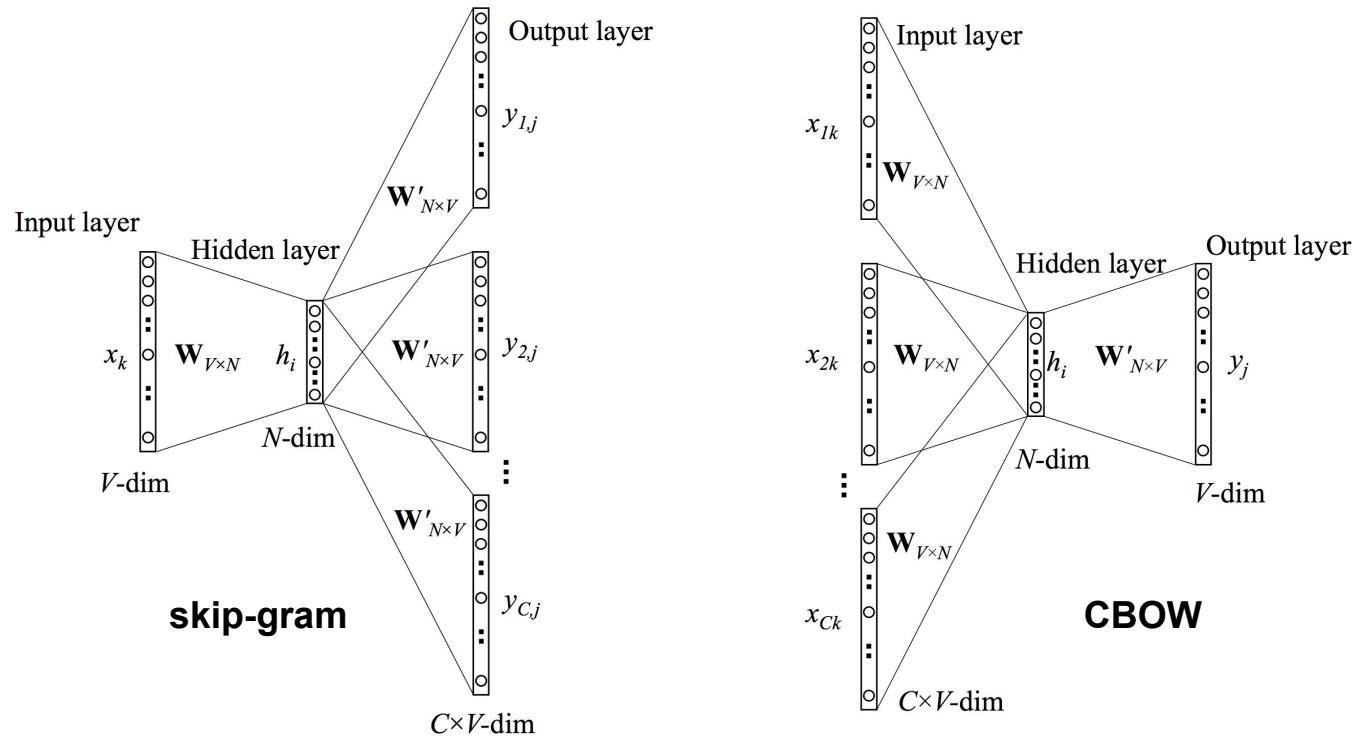
In our case with d -dimensional vectors and V -many words:

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

Word2vec One Context Word



Word Embeddings



© Taken from “word2vec Parameter Learning Explained”
<https://arxiv.org/pdf/1411.2738.pdf>

Word Embeddings Visualization

<https://projector.tensorflow.org/>

Document Similarity

Topic Modeling

Overview: Similarity, Topicality

Find a document that is more similar to a given input text?

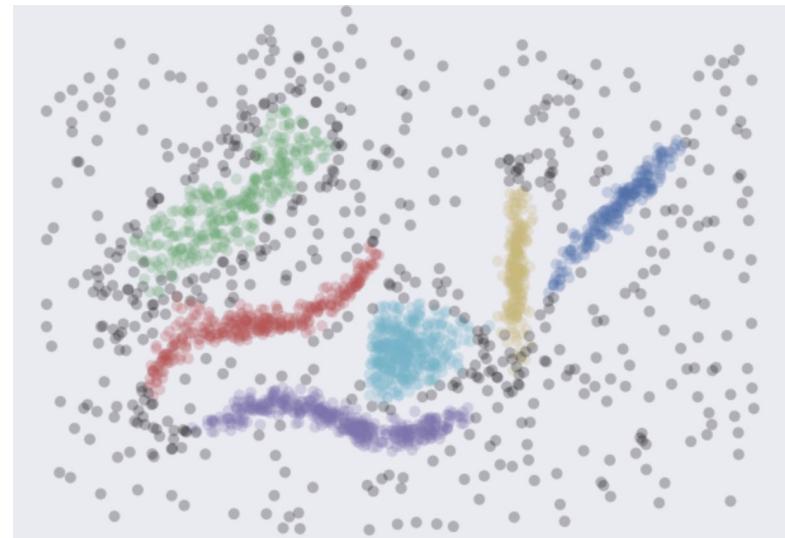
- Similarity measures
 - Text based (edit distance)
 - Vector space representation based (euclidean distance, cosine similarity, WMD)
- Document similarity
 - Tf-idf representation
 - W2V representation

Clustering

Group input documents into known or unknown number of groups based on similarity between them.

Why?

- Reduce input size
- Discover hidden patterns in your data without any supervision
- Visualize these patterns



© taken from lmcinnes github

<http://nbviewer.jupyter.org/github/lmcinnes/hdbscan/blob/master/notebooks/Comparing%20Clustering%20Algorithms.ipynb>

LDA

Topics

```
gene      0.04
dna       0.02
genetic   0.01
...

```

```
life      0.02
evolve   0.01
organism 0.01
...

```

```
brain     0.04
neuron   0.02
nerve    0.01
...

```

```
data      0.02
number   0.02
computer 0.01
...

```

Documents

Seeking Life's Bare (Genetic) Necessities

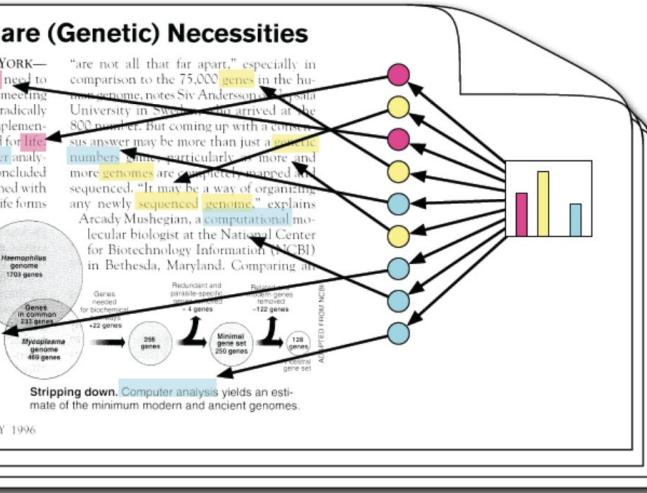
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genomic meeting here,¹⁰ two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing. Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

LDA

Topics



Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism ~~need~~ to survive? Last week at the genome meeting here,¹ two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all

the genomes, he says, will reveal which genes

are needed for biochemical processes.

Redundant and gene-redundant genes are removed.

Relaxed and minimal genomes are identified.

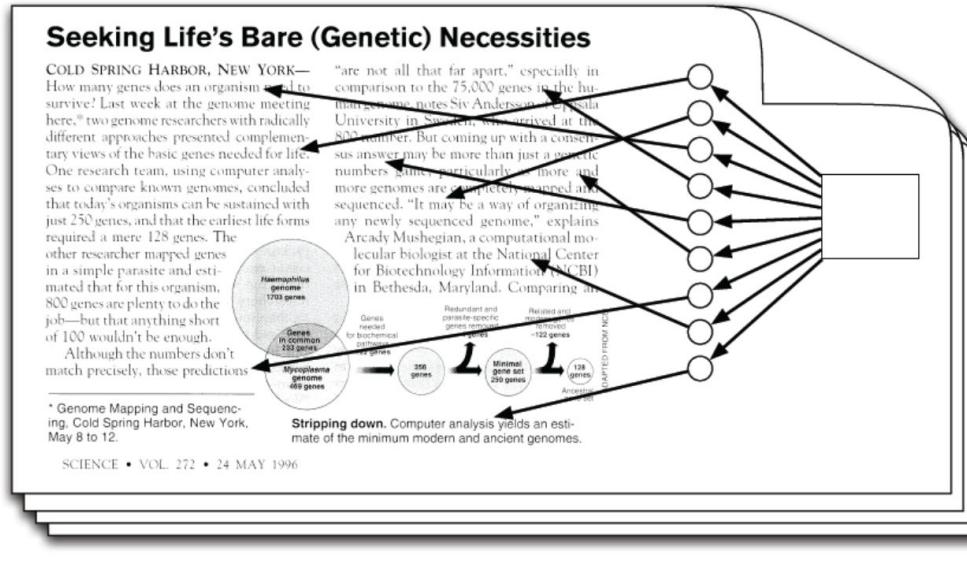
128 genes are identified.

128 genes are identified.

128 genes are identified.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

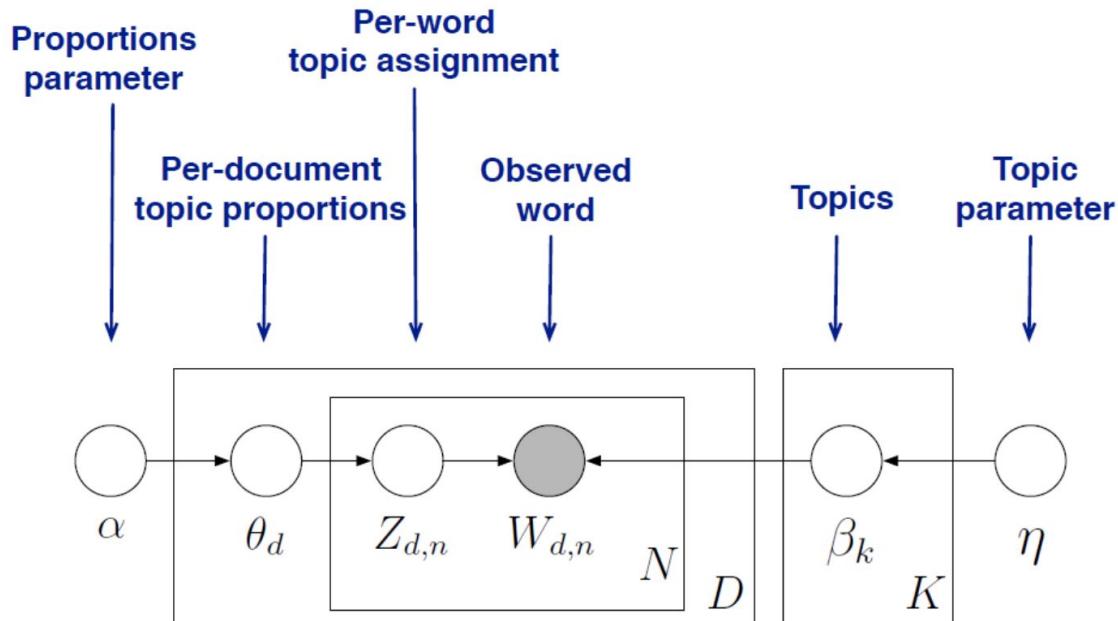
Topic proportions and assignments



- Our goal is to **infer** the hidden variables
- I.e., compute their distribution conditioned on the documents

$$p(\text{topics, proportions, assignments} | \text{documents})$$

LDA



- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; unshaded nodes are hidden.
- Plates indicate replicated variables.

General learning schema

To learn the topic representation of each document and the words associated to each topic:

- Go over the documents and **assign randomly** each word in the topic to one of the K topics
- for each document and each word in the document
 - for each topic t ,
 - (1) $p(topic_k|document_n)$ - the portion of words in document n that are currently assigned to the topic k ;
 - (2) $p(word_i|topic_k)$ - portion of assignments to topic k over all documents that come from this word i .
 - Reassign word i to a new topic, where we choose a topic with probability
 $p(topic_k|document_n) \cdot p(word_i|topic_k)$ - the probability that topic k generated word i .
- Stop once steady

Advanced NLP 1: Linking text to semantics

Extracting entities: Linking text to semantics

What are entities?

entity

/ɛnˈtɪti/ 

noun

a thing with distinct and independent existence.

What are **named** entities?

From [wiki](#): **Named entity** is a real-world object, such as persons, locations, organizations, products, etc., that can be denoted with a proper name. It can be abstract or have a physical existence.

Examples: Donald Trump, Paris, Google.

Detecting entities: Subtasks

Named Entity Recognition: Identify where named entities appear in the text, with particular types (e.g. Organization, Location, Person, Time)

Example: flight to Paris/LOC, set alarm at 8am/Time.

Entity Linking (or disambiguation): Match the identity of the entity mentioned in the text to actual entity (defined by non-ambiguous id)

Example: Donald [Trump](#) and [Trump](#) card.

Coreference resolution: Find all expressions that refer to the same entity in a text.

Example: The food was salty so guests did not enjoy [it](#).

Databases of semantic relations

- ***Is-a***: the relationship between two entities in which one entity inherits from the other.
 - a. ***Hypernyms*** are entities with broader/more general meaning,
 - b. ***Hyponyms*** are entities with narrower/more specific meaning.

Example: iPhone is a phone, phone is an electronic device.
- ***Meronymy*** denotes a constituent part of, or a member of something. Meronym is a part of a whole!

Example: battery is a part of a smartphone.

Databases: [WordNet](#)

Databases of semantic relations

- **Synsets** denotes synonyms that are interchangeable in some context w/o changing the true value of the statement;
Example: there were done **minor** or **small** mistakes.
- **Metonymy** substitution of the name of an attributes for the thing that was meant.
Example: Business executives are in town VS Suits are in town.
- **Anything**
Example: Find who is married to who? What causes what?

Knowledge Databases and Ontologies

Wikipedia: knowledge about the world; not explicitly structured, but often used as a source of entities in the literature.

DBPedia: mined structured content from Wikipedia

WordNet: lexical database for English, representing different word senses as 'synsets' and also containing synset relations.

OMCS: Commonsense knowledge database, representing human knowledge "The sky is blue". Not easy to mine on the web!

YAGO: is an open source knowledge base developed at the Max Planck Institute.

Knowledge Graph

A database of entities and facts
(as relations between entities)

Goal: understand/answer **factual**
questions about **the world**

Barack Obama

44th U.S. President



barackobama.com

Barack Hussein Obama II is an American politician who served as the 44th President of the United States from January 20, 2009, to January 20, 2017. A member of the Democratic Party, he was the first African American to be elected to the presidency and previously served as a United States Senator from Illinois. [Wikipedia](#)

Born: August 4, 1961 (age 57 years), [Kapi'olani Medical Center for Women and Children, Honolulu, Hawaii, United States](#)

Height: 1.85 m

Vice president: [Joe Biden \(2009–2017\)](#)

Education: [Harvard Law School \(1988–1991\), MORE](#)

Parents: [Ann Dunham, Barack Obama Sr.](#)

Entity Linking: Common approach

1. Annotate text
2. Find entity candidates by their surface forms (= n-grams of how an entity can appear in the text)
3. Disambiguate to particular entity based on the context (e.g. co-occurring words)

Advanced NLP 2: Sequence Modeling

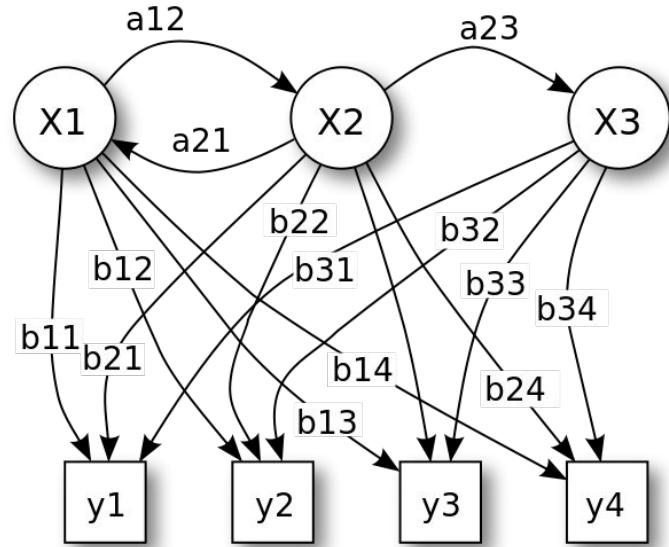
Sequence Modeling (HMM, CRF)

Hidden Markov Models

$$p(l, s) = p(l_1) \prod_i p(l_i | l_{i-1}) p(w_i | l_i)$$

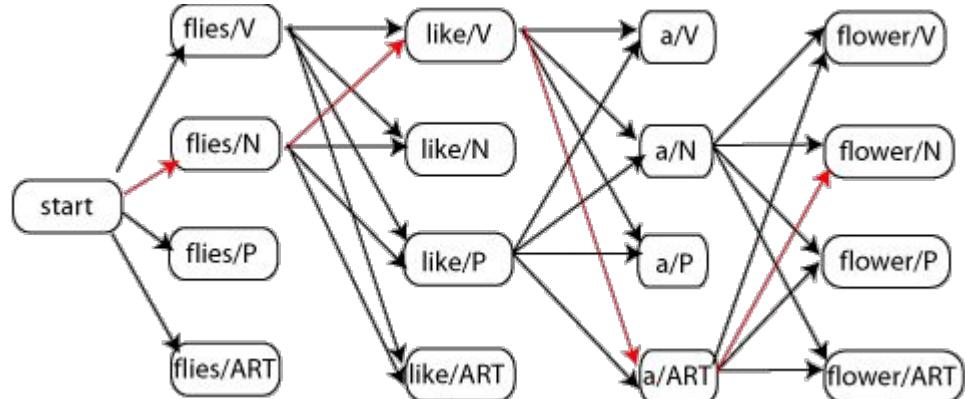
where

- $p(l_i | l_{i-1})$ are transition probabilities
- $p(w_i | l_i)$ are emission probabilities (ϵ)



Hidden Markov Models

POS example:



Conditional Random Fields

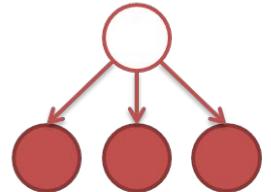
Feature function is a function that takes in as input:

- a sentence s
- the position i of a word in the sentence
- the label l_i of the current word
- the label l_{i-1} of the previous word

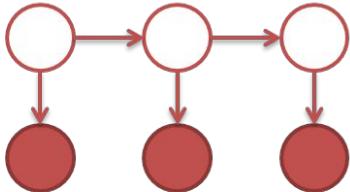
So given a sequence s we can now score it as follows $score(l|s) = \sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l_i, l_{i-1})$

$$p(l|s) = \frac{\exp[score(l|s)]}{\sum_{l'} \exp[score(l'|s)]} = \frac{\exp[\sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l_i, l_{i-1})]}{\sum_{l'} \exp[\sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l'_i, l'_{i-1})]}$$

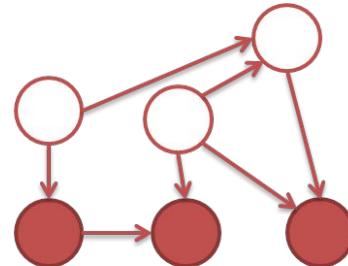
Arbitrary features can be also added on top :)



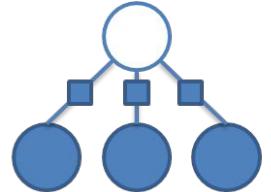
Naïve Bayes



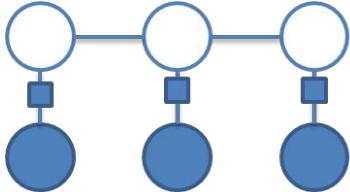
Markov models



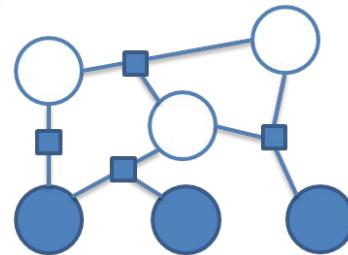
Directional Models



Logistic Regression



Linear-chain CRF



CRF

Sequences

Graphs

Generative

Discriminative

Adapted from C. Sutton, A. McCallum, "An Introduction to Conditional Random Fields", ArXiv, November 2010

Sequence Modeling

Motivation: Why Sequence Models?

Recurrent Neural Networks: Overview

Recurrent Neural Network Architectures

Types of RNN Cells

Outline

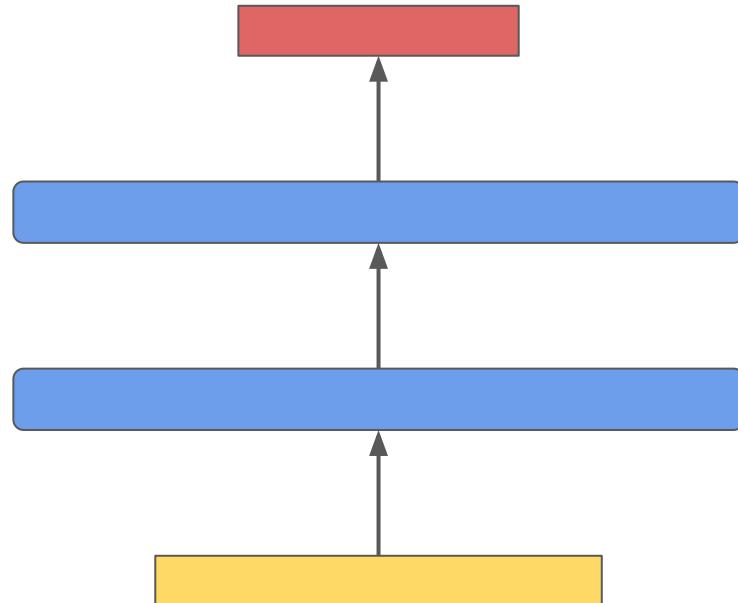
Motivation: Why Sequence Models?

Recurrent Neural Networks: Overview

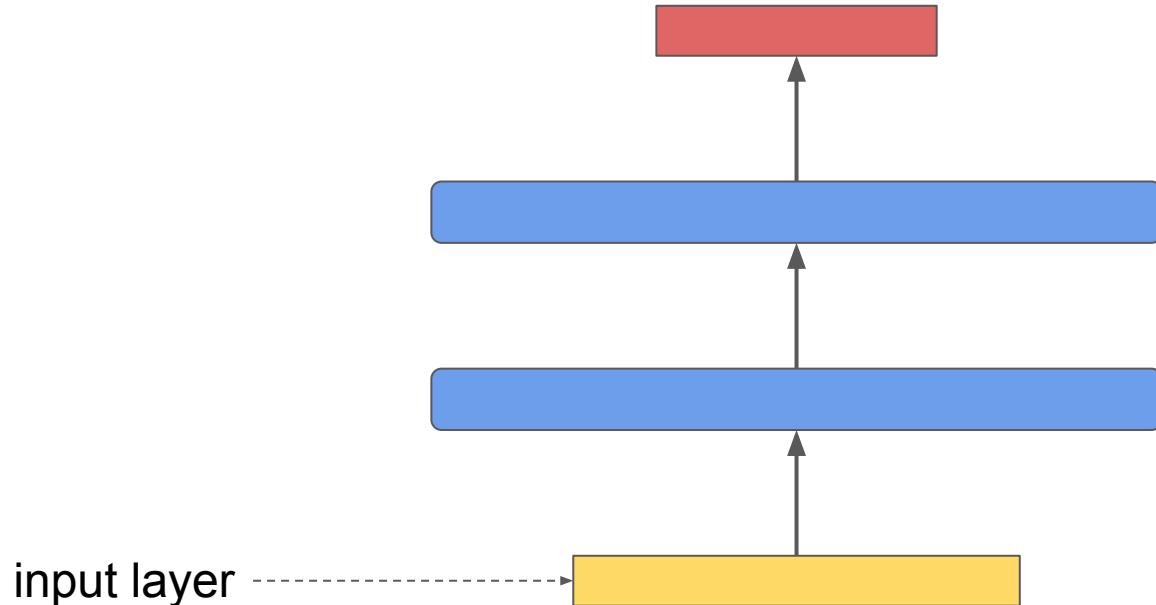
Recurrent Neural Network Architectures

Types of RNN Cells

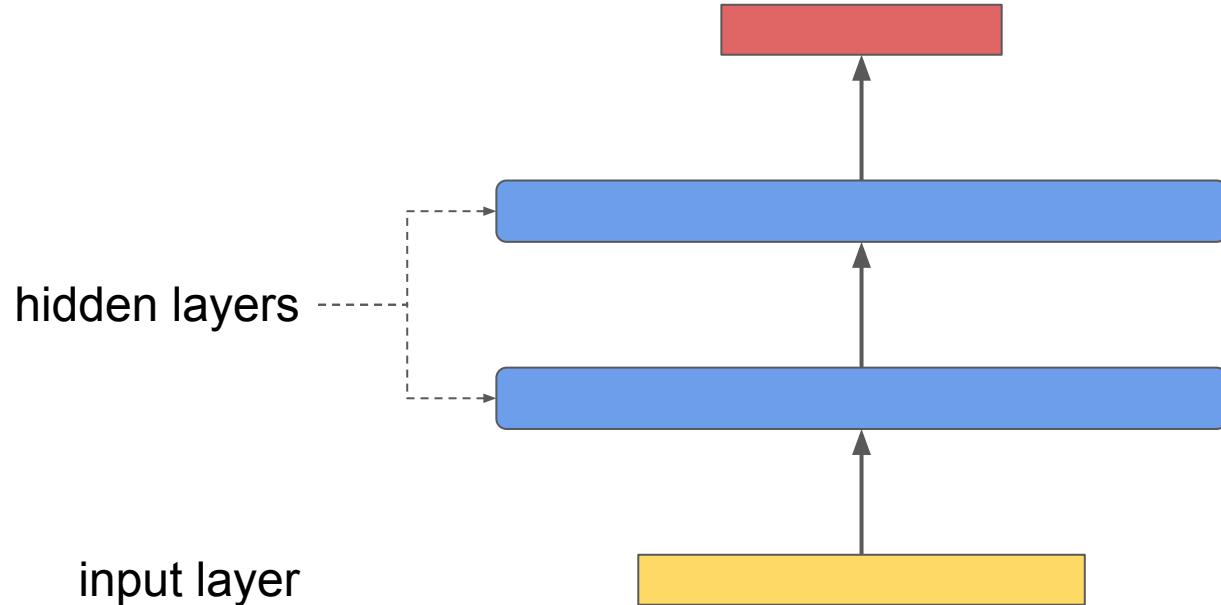
Feed Forward Networks



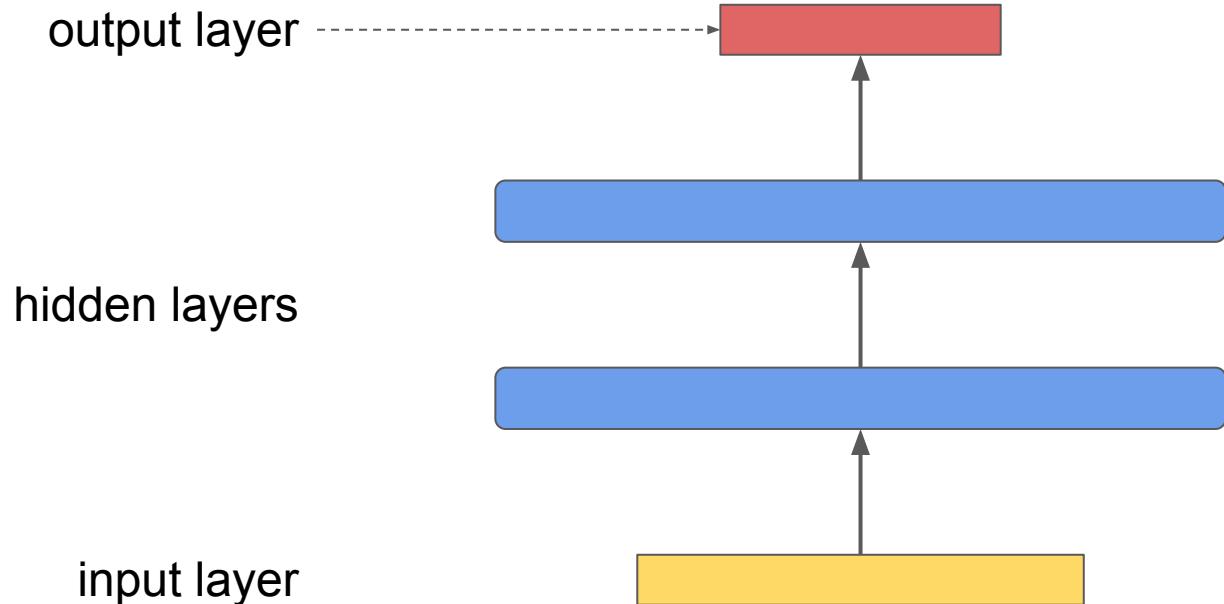
Feed Forward Networks



Feed Forward Networks

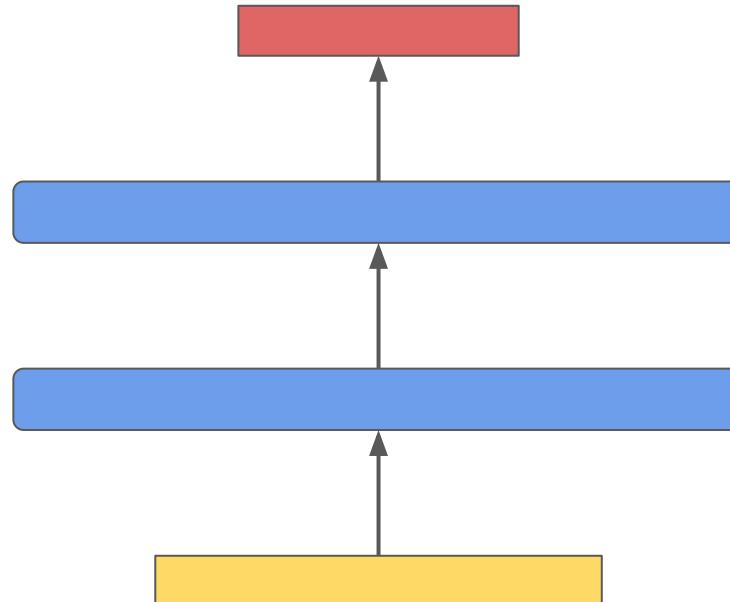


Feed Forward Networks



Feed Forward Networks

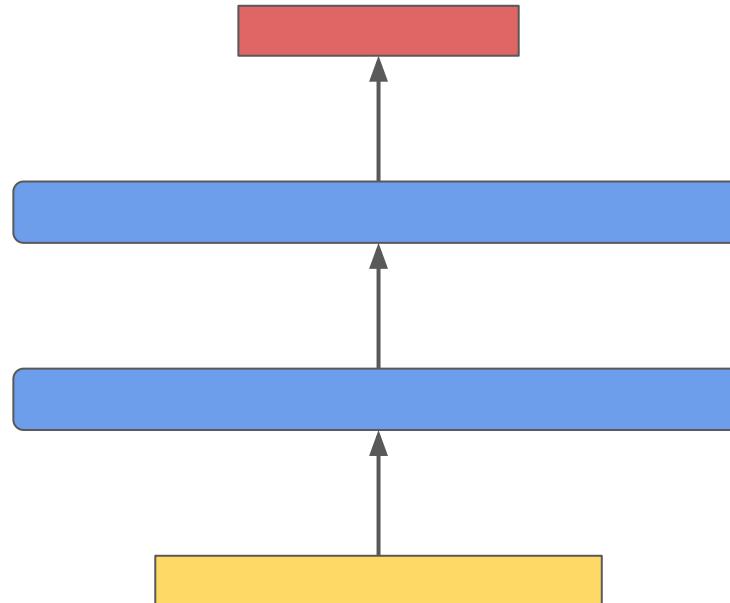
Fixed size layers



Feed Forward Networks

Fixed size layers

Inference is stateless

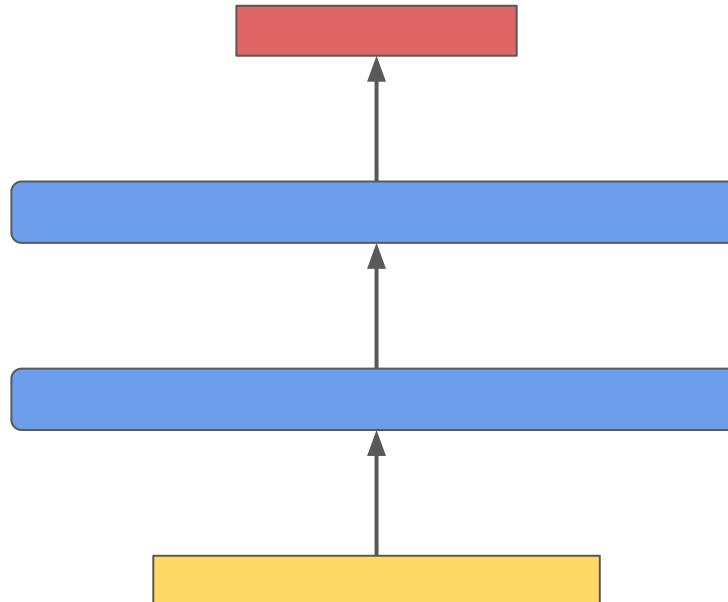


Feed Forward Networks

Fixed size layers

Inference is stateless

Nodes are unordered



Language as Input

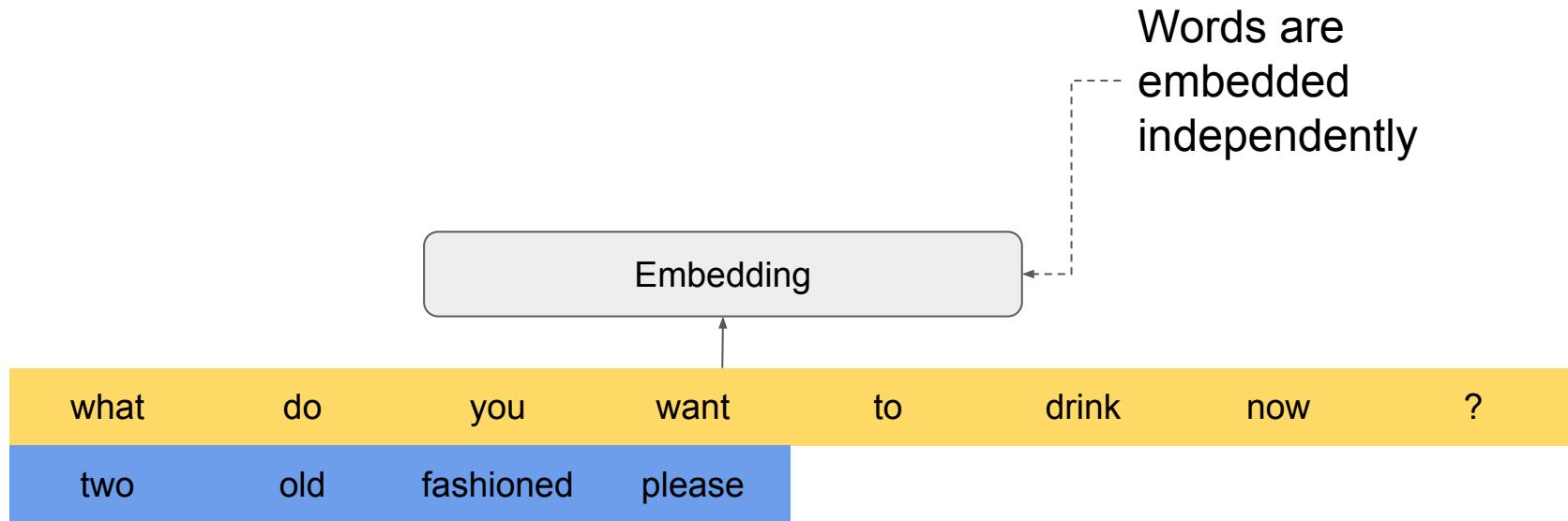
Input to a language model can have variable length. For example,

what	do	you	want	to	drink	now	?
two	old	fashioned	please				

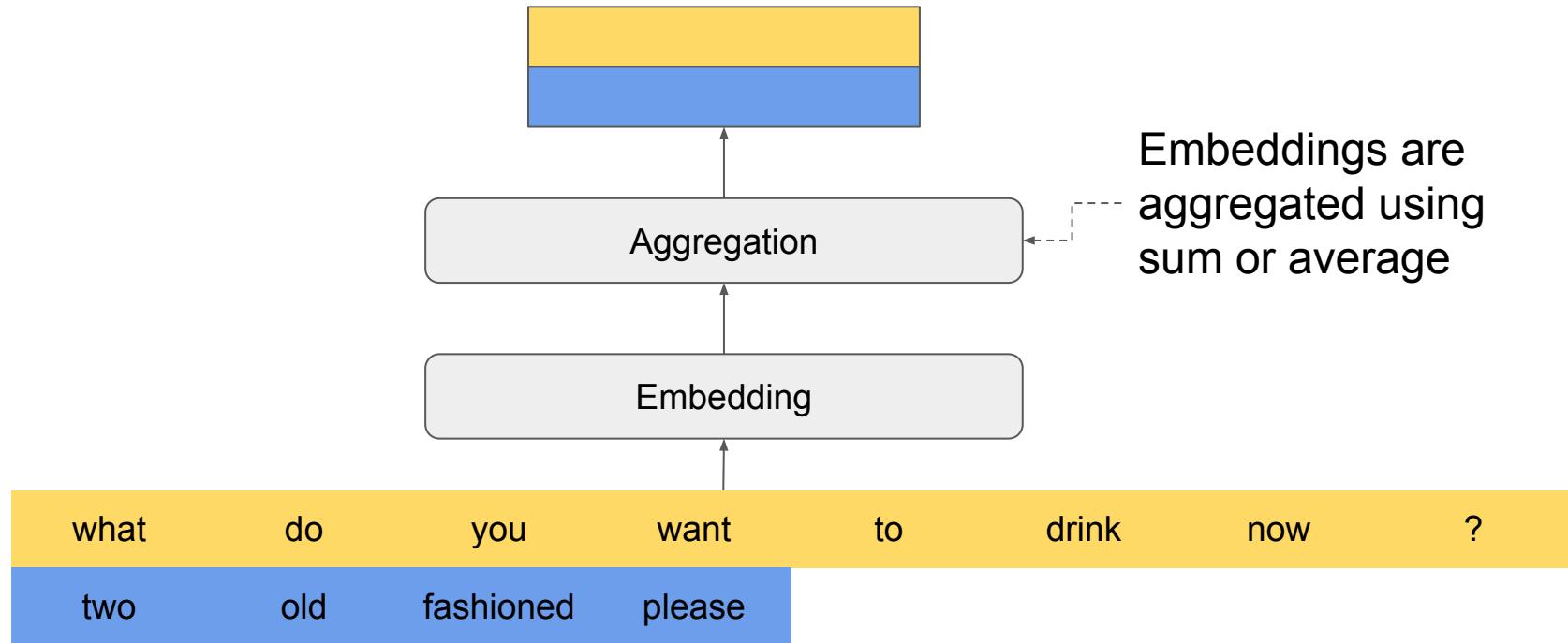
Language as Input: the “Typical” Approach

what do you want to drink now ?
two old fashioned please

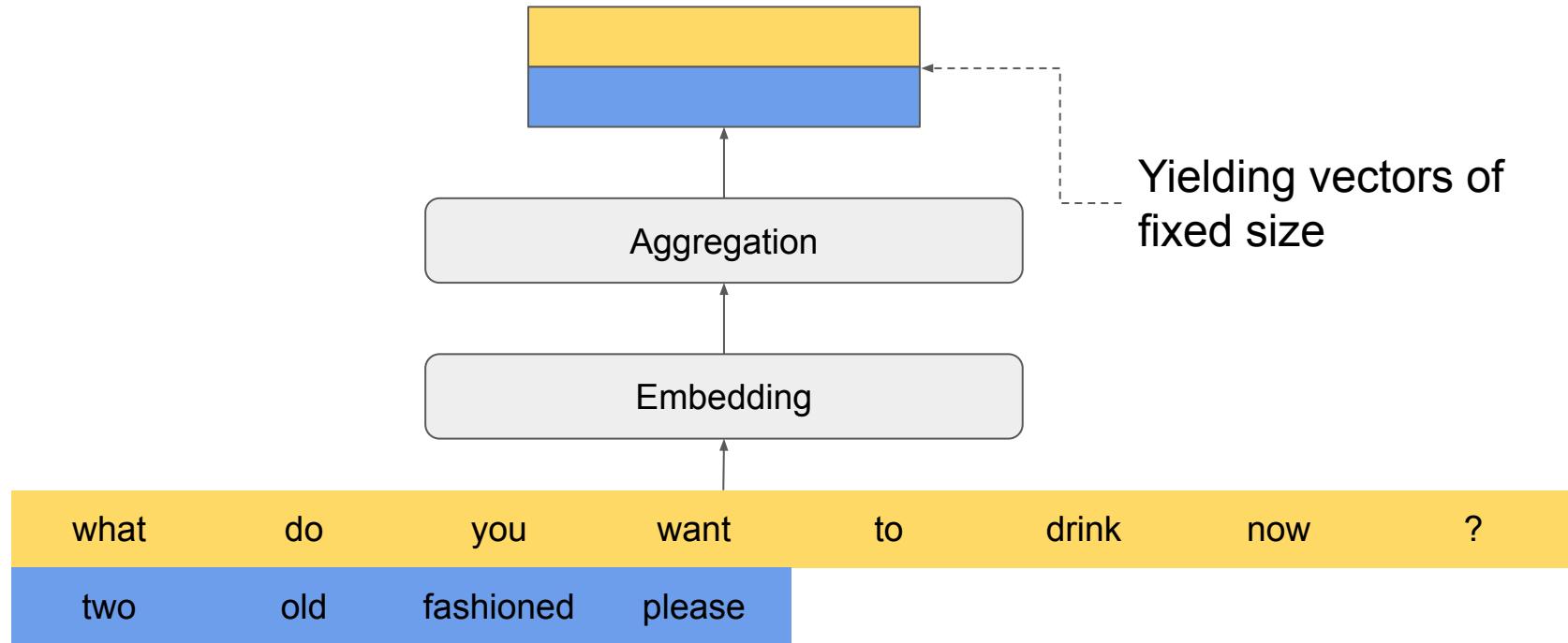
Language as Input: the “Typical” Approach



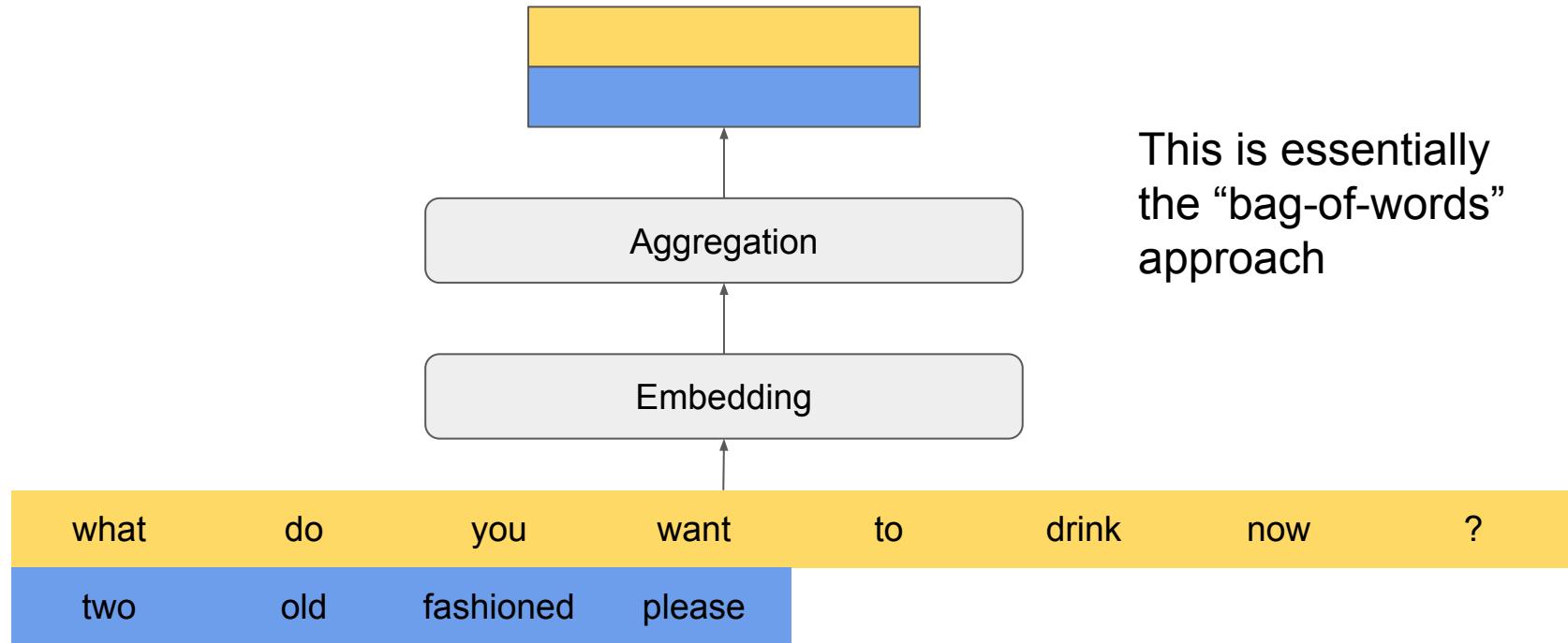
Language as Input: the “Typical” Approach



Language as Input: the “Typical” Approach



Language as Input: the “Typical” Approach



Bag of Words: Pros and Cons

Pro: Allows variable length data to be fed into a typical feed-forward network.

Pro: Variable length features can easily be combined with fixed length features.

Pro: Widespread empirical success.

Bag of Words: Pros and Cons

Pro: Allows variable length data to be fed into a typical feed-forward network.

Pro: Variable length features can easily be combined with fixed length features.

Pro: Widespread empirical success.

Con: Destroys the structure of the input:

The cat sat on the mat != {cat, sat, the, the, mat, on}

Structure and Ordering is Important

For example:

Humor or sarcasm

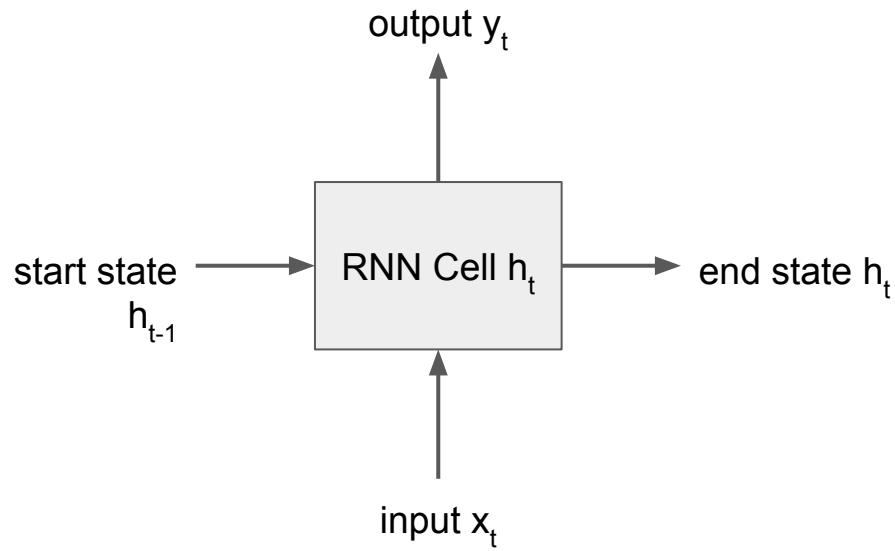
Modifiers such as **somewhat**, **extremely** or **not**

Multi-word phrases such as **tiger mosquito** or **zurich city bridge**

Big Idea of Recurrent Neural Networks



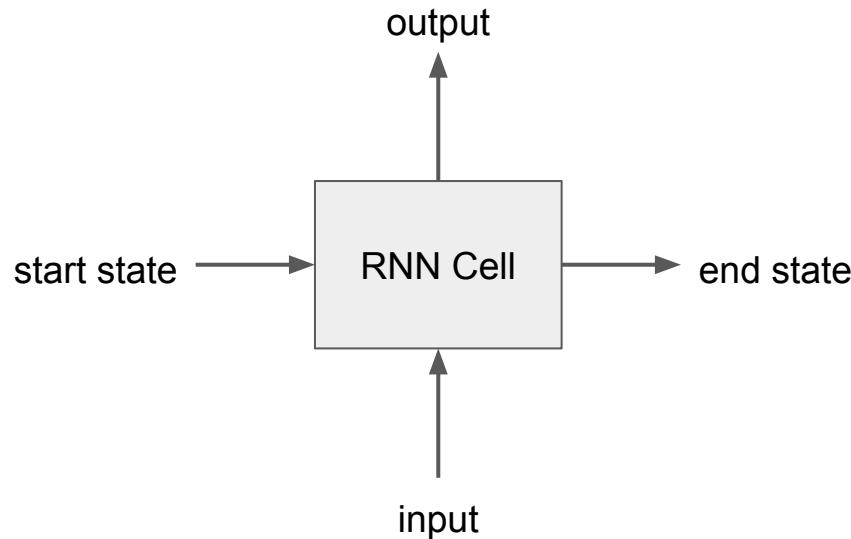
Recurrent Neural Nets: A Cell...



```
output, end_state =  
rnn_cell(input, start_state)
```

$$\begin{aligned} \mathbf{h}_t &= g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c}) \\ \hat{\mathbf{y}}_t &= \mathbf{W}\mathbf{h}_t + \mathbf{b} \end{aligned}$$

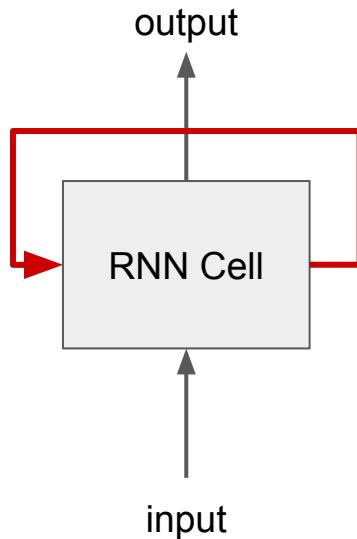
Recurrent Neural Nets: A Cell...



The only rule:

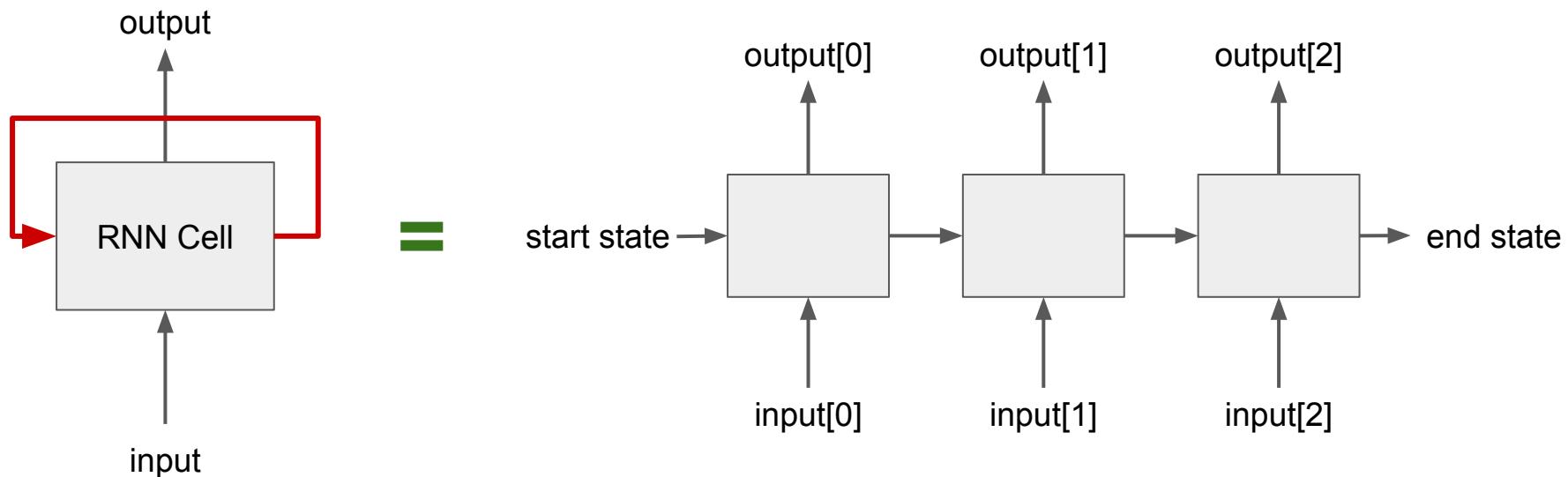
start state and **end state** must have the same dimension.

Recurrent Neural Nets: A Cell... Wrapped in a For Loop



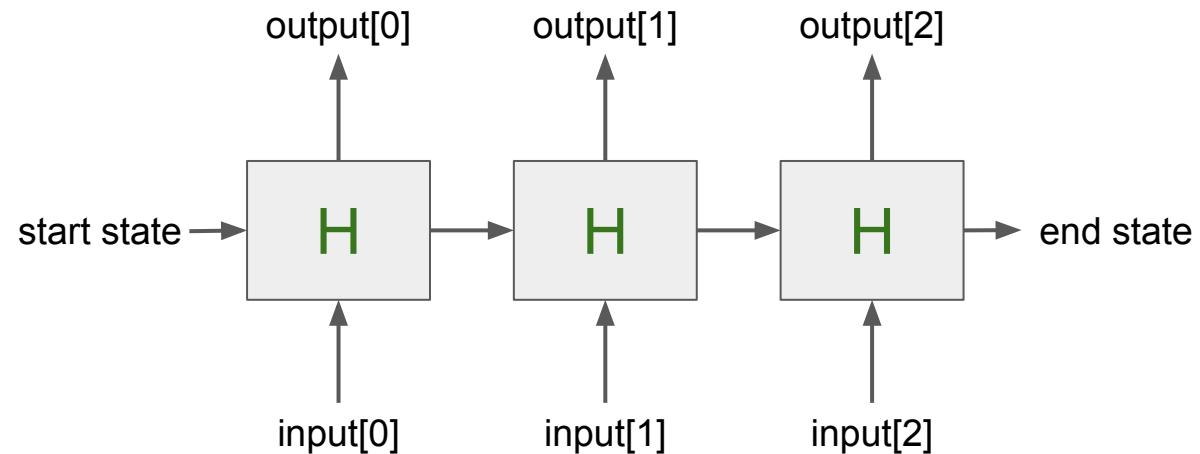
```
for i in range(sequence_length):  
    output[i], end_state =  
        rnn_cell(input[i], start_state)  
    start_state = end_state
```

Recurrent Neural Nets: Unrolling the For Loop



Recurrent Neural Nets: Unrolling the For Loop

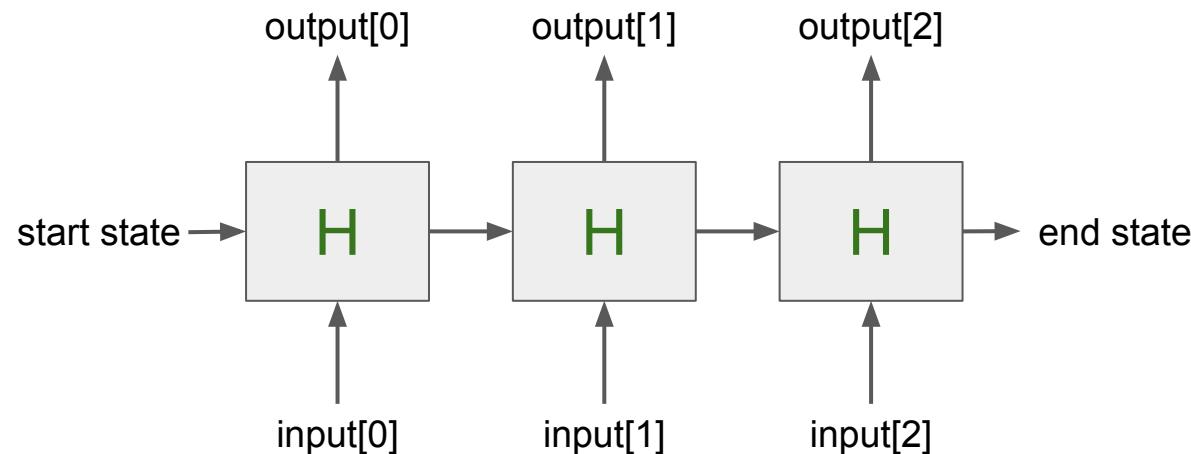
Each cell is **identical** and they **share weights**.



Recurrent Neural Nets: Unrolling the For Loop

Each cell is **identical** and they **share weights**.

When we train an RNN, we are really training a **single cell** that is used repeatedly.

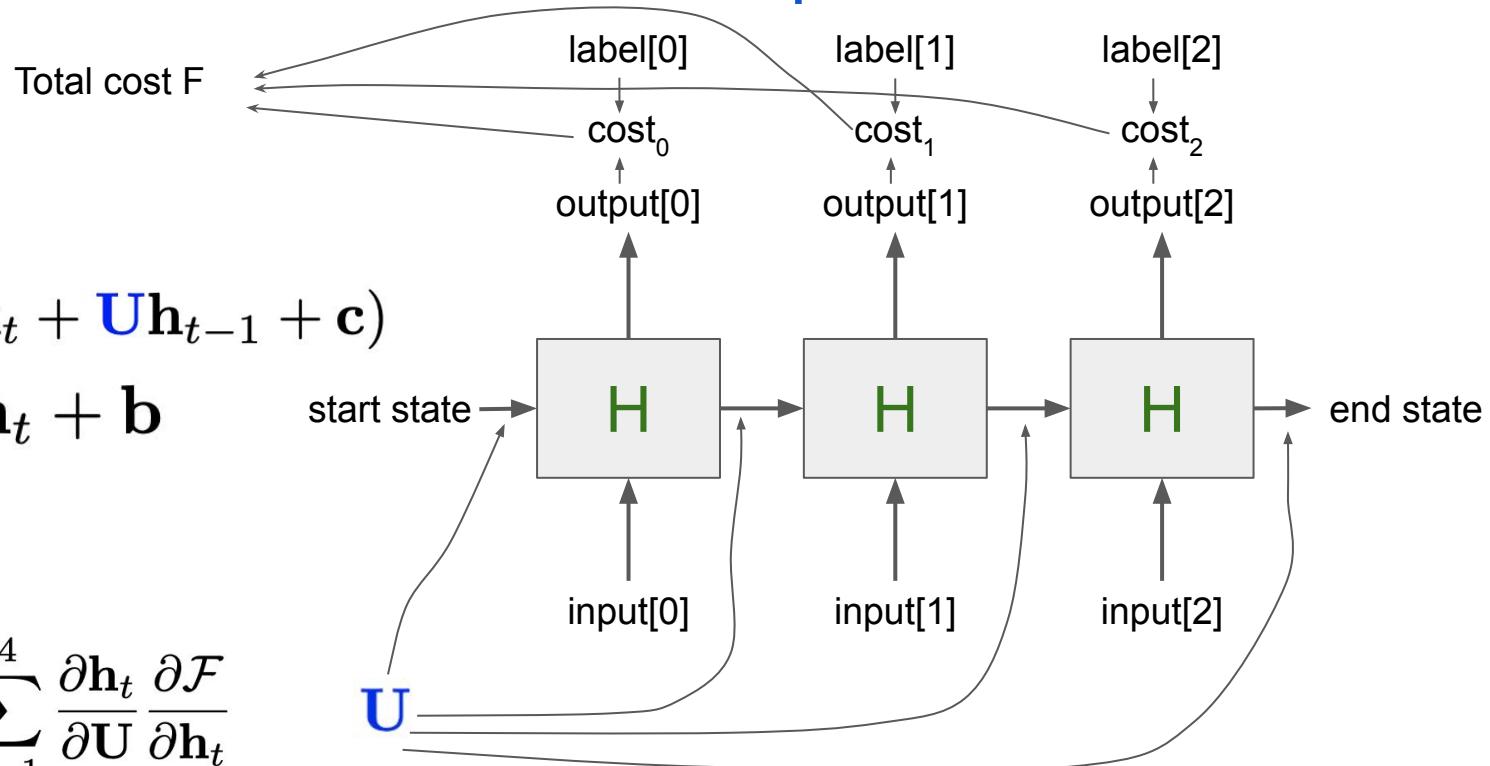


Recurrent Neural Nets: Error computation

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$

$$\frac{\partial \mathcal{F}}{\partial \mathbf{U}} = \sum_{t=1}^4 \frac{\partial \mathbf{h}_t}{\partial \mathbf{U}} \frac{\partial \mathcal{F}}{\partial \mathbf{h}_t}$$



Outline

Motivation: Why Sequence Models?

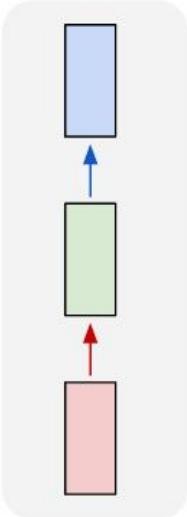
Recurrent Neural Networks: Overview

Recurrent Neural Network Architectures

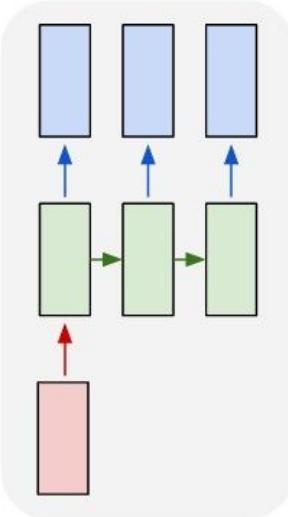
Types of RNN Cells

Types Of NNs

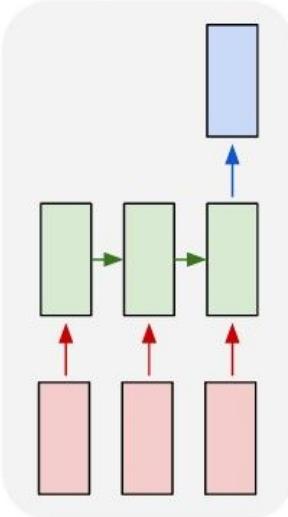
one to one



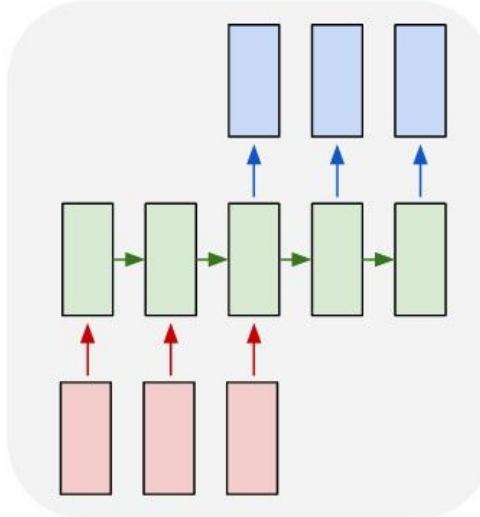
one to many



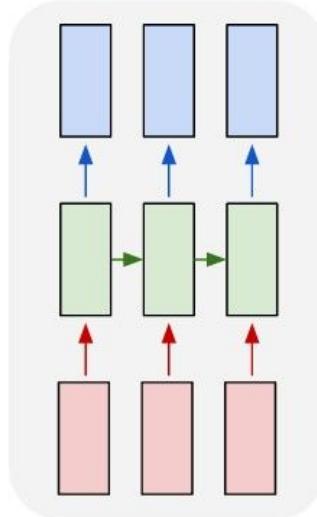
many to one



many to many

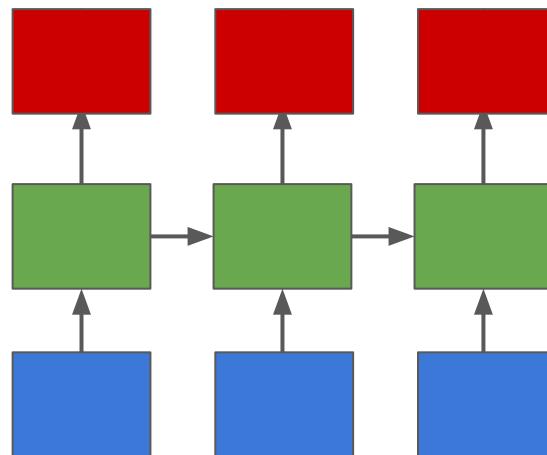


many to many



Taken from [2].

RNN Architectures: Many to Many

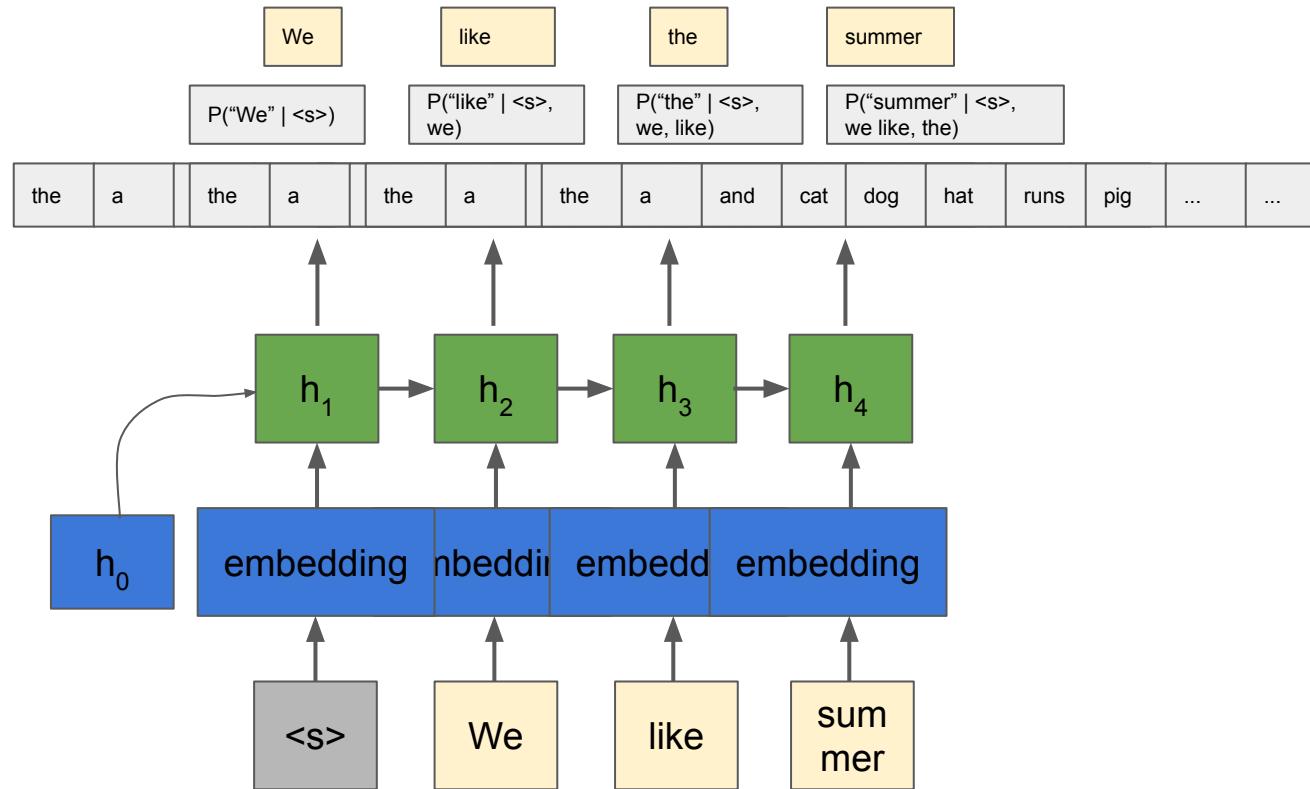


Example: speech recognition

Input: windowed audio

Output: sequence of letters

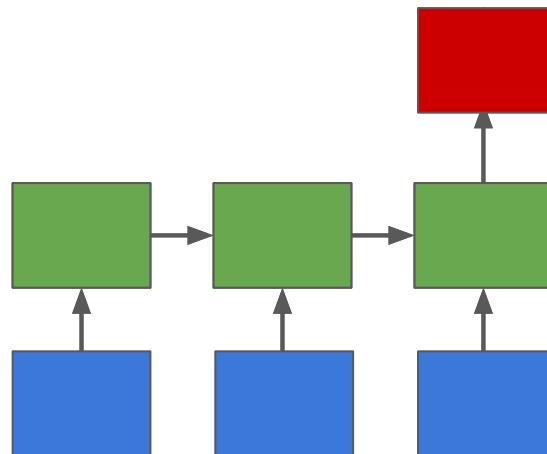
RNN Architectures: Example



$$\mathbf{u} = \mathbf{W}\mathbf{h} + \mathbf{b}$$

$$p_i = \frac{\exp u_i}{\sum_j \exp u_j}$$

RNN Architectures: Many to One

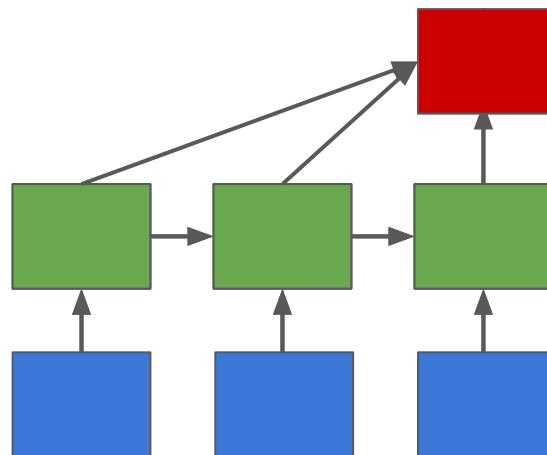


Example: **spam filtering**

Input: **words in an email**

Output: **spam or not spam**

RNN Architectures: Many to One

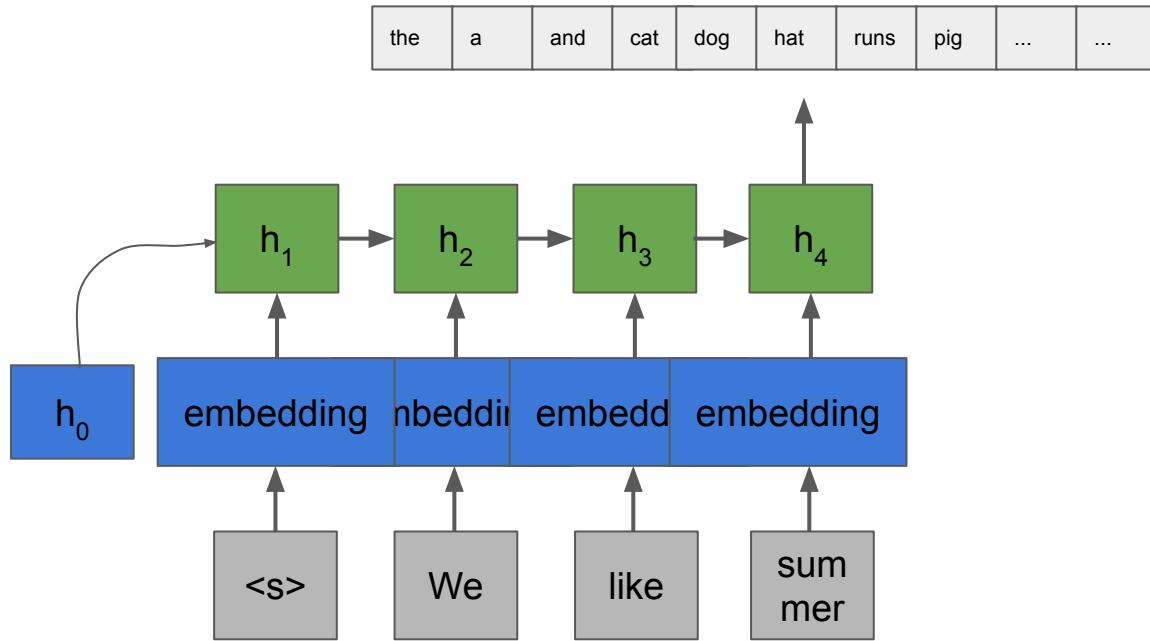


Example: **spam filtering**

Input: **words in an email**

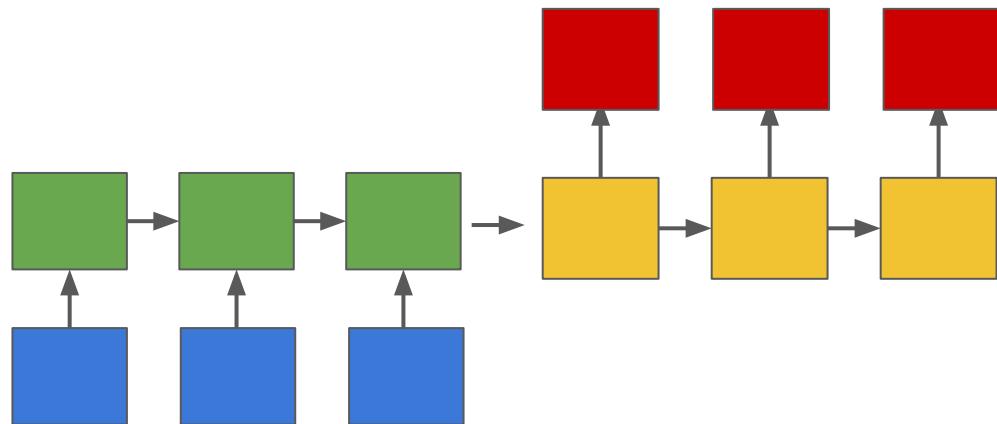
Output: **spam or not spam**

RNN Architectures: Example



$$\mathbf{u} = \mathbf{W}\mathbf{h} + \mathbf{b}$$
$$p_i = \frac{\exp u_i}{\sum_j \exp u_j}$$

RNN Architectures: Sequence to Sequence

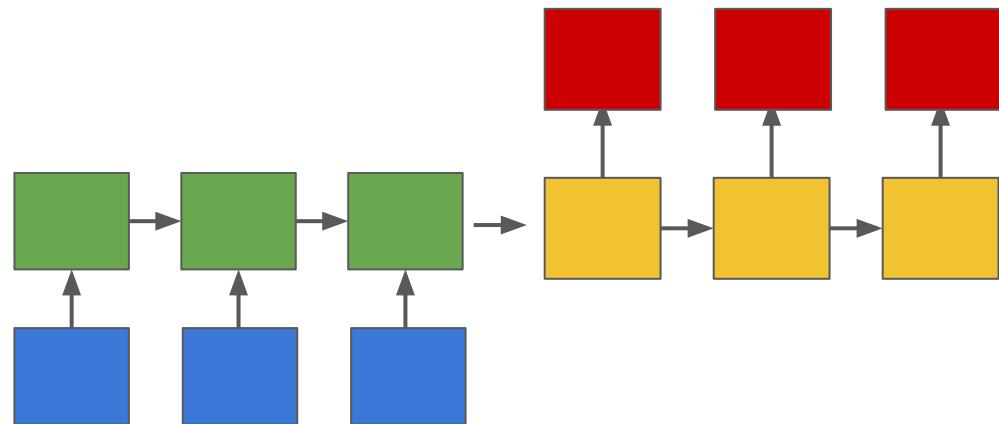


Example: **translation**

Input: **Spanish**

Output: **Mandarin**

RNN Architectures: Sequence to Sequence

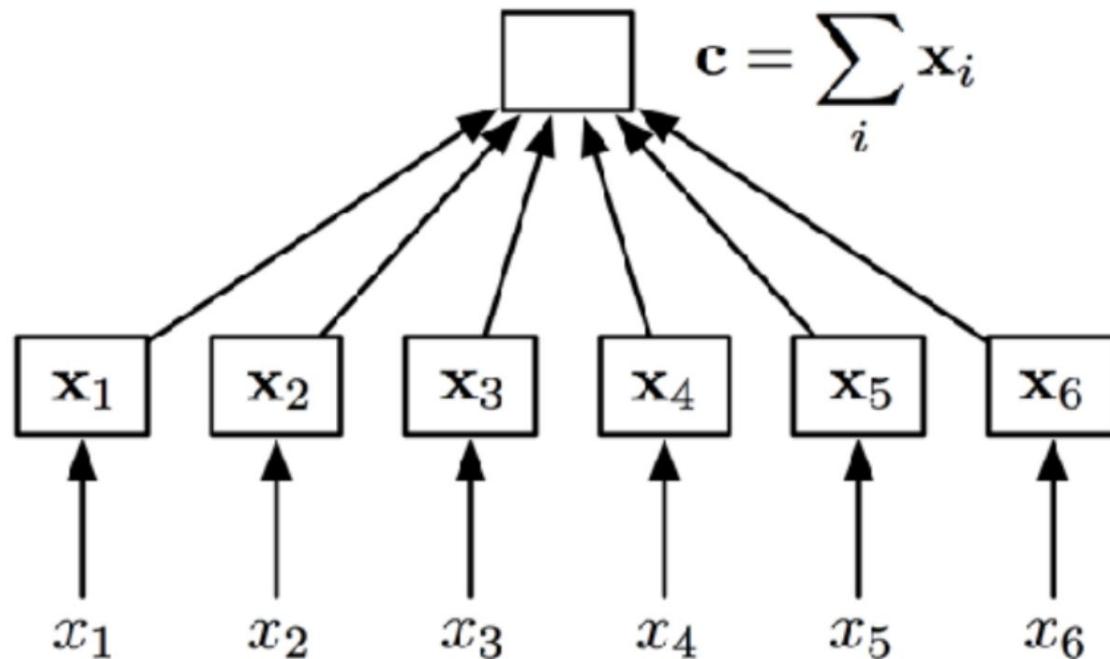


We often use different cells for
encoding and decoding.

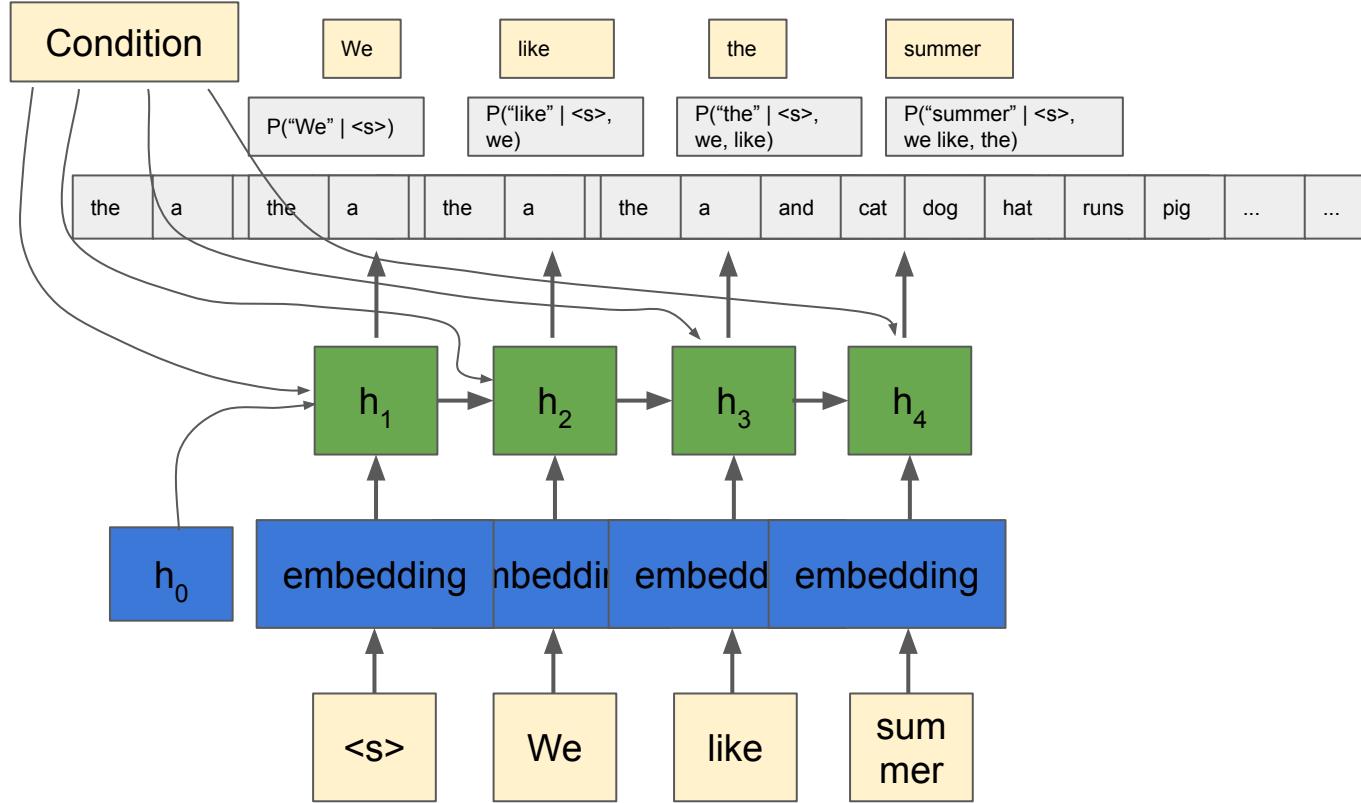
Conditional LM

x "input"	w "text output"
An author	A document written by that author
A topic label	An article about that topic
{SPAM, NOT_SPAM}	An email
A sentence in French	Its English translation
A sentence in English	Its French translation
A sentence in English	Its Chinese translation
An image	A text description of the image
A document	Its summary
A document	Its translation
Meteorological measurements	A weather report
Acoustic signal	Transcription of speech
Conversational history + database	Dialogue system response
A question + a document	Its answer
A question + an image	Its answer

Encoder



RNN Architectures: Example



RNN Architectures: Beam Search

E.g., for $b=2$:

$\textcolor{red}{x} = \textit{Bier trinke ich}$
beer drink I

$\langle s \rangle$

logprob=0

w_0

w_1

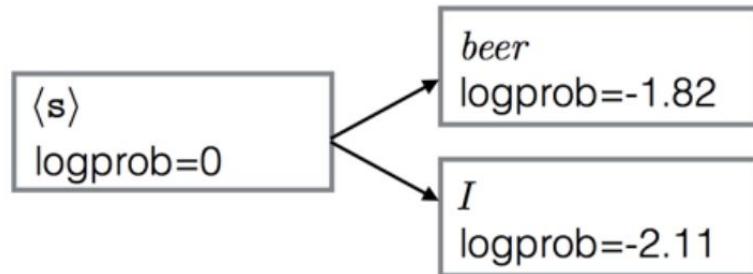
w_2

w_3

RNN Architectures: Beam Search

E.g., for $b=2$:

$x = \text{Bier trinke ich}$
beer drink I



w_0

w_1

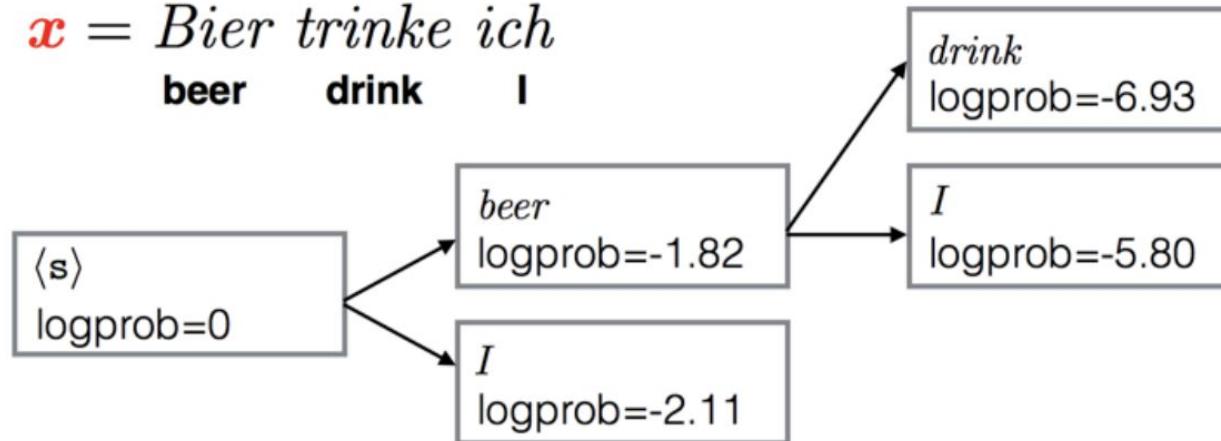
w_2

w_3

RNN Architectures: Beam Search

E.g., for $b=2$:

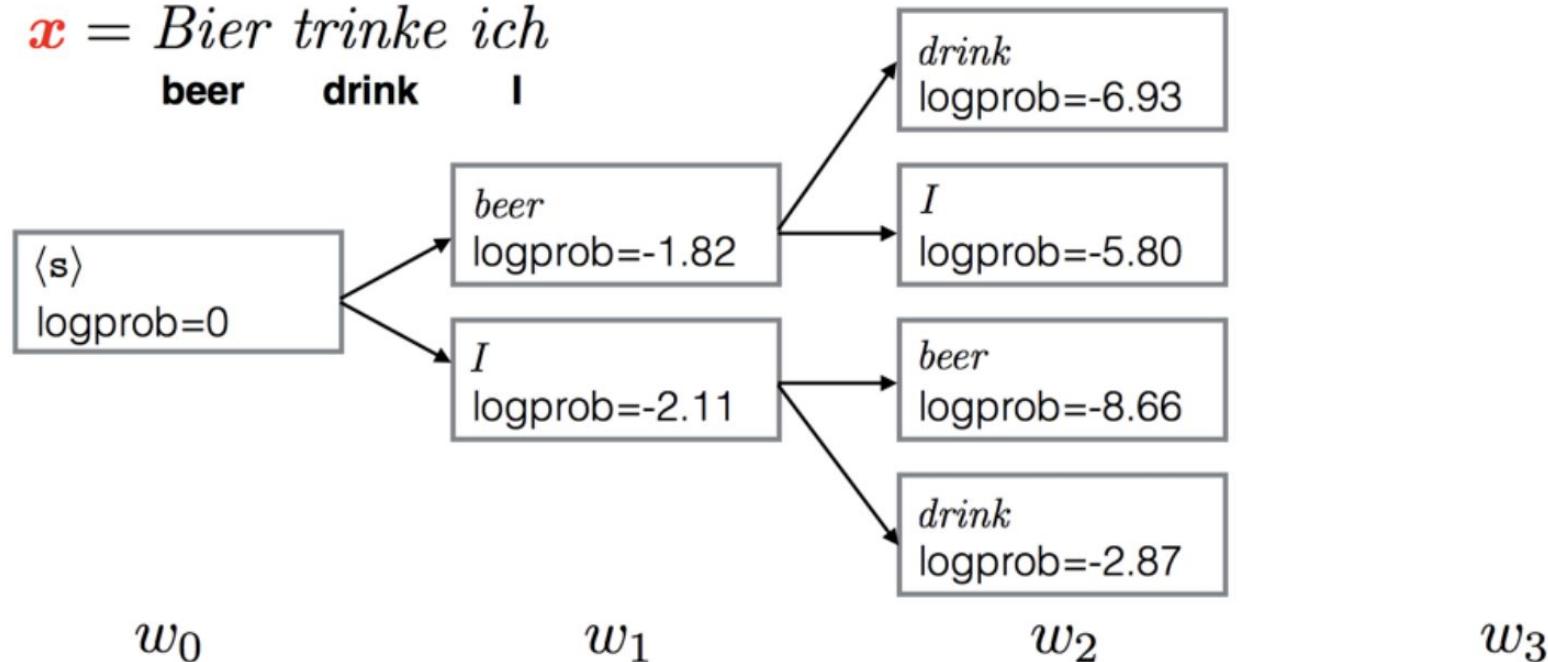
$\textcolor{red}{x} = \textit{Bier trinke ich}$
beer drink I



RNN Architectures: Beam Search

E.g., for $b=2$:

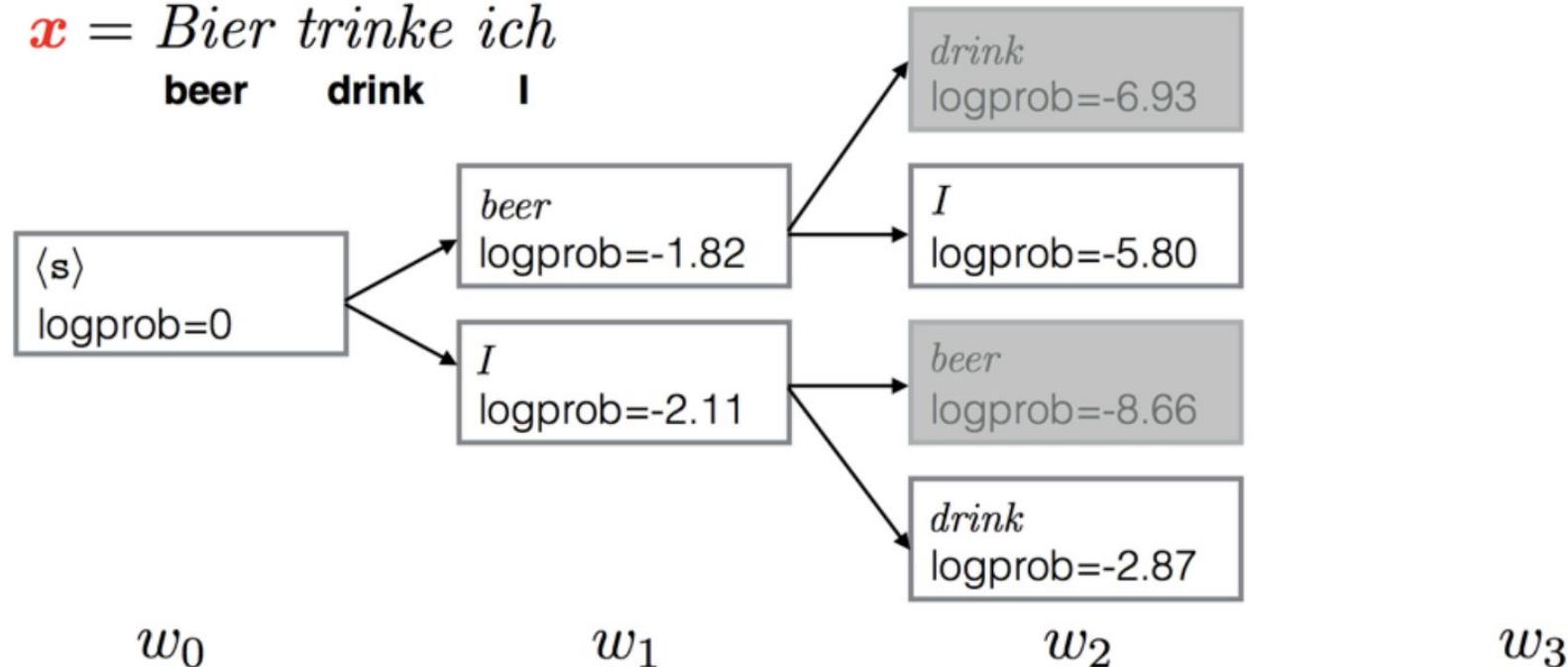
$x = \text{Bier trinke ich}$
beer **drink** **I**



RNN Architectures: Beam Search

E.g., for $b=2$:

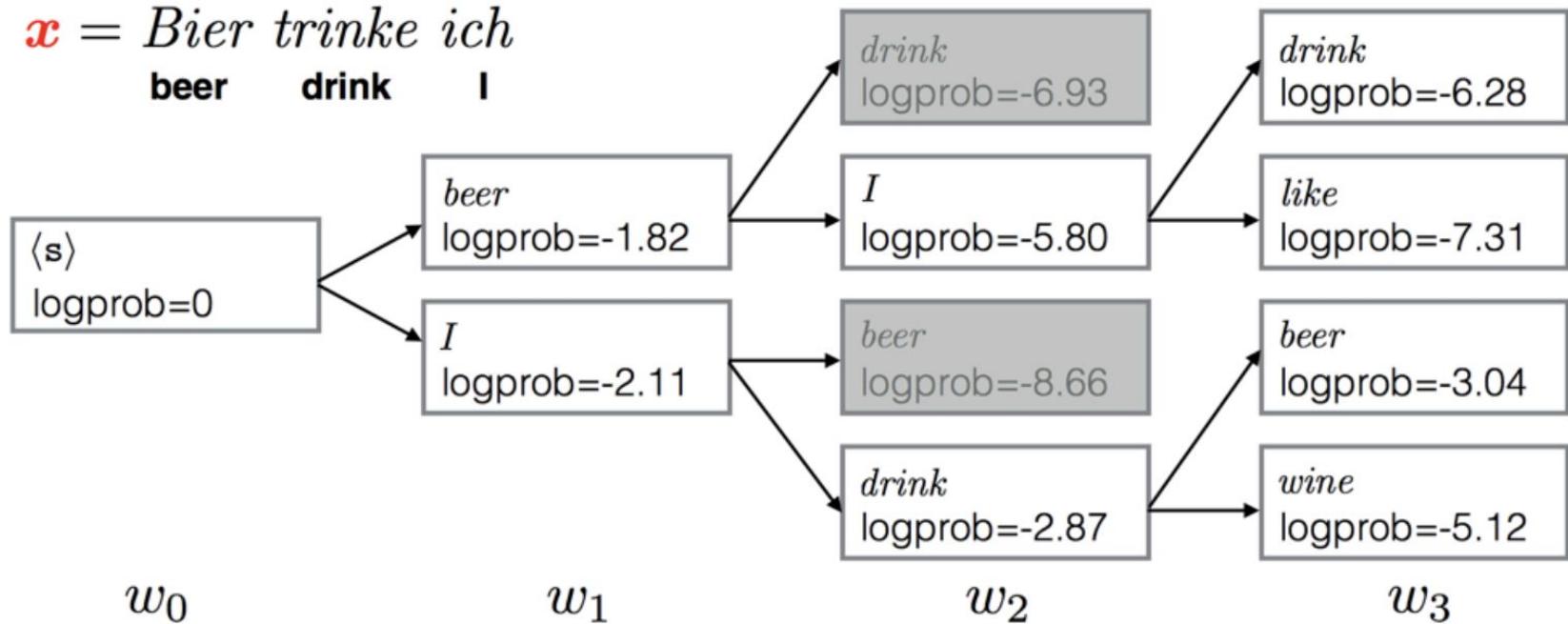
$\textcolor{red}{x} = \textit{Bier trinke ich}$
beer **drink** **I**



RNN Architectures: Beam Search

E.g., for $b=2$:

$x = Bier trinke ich$
beer drink I

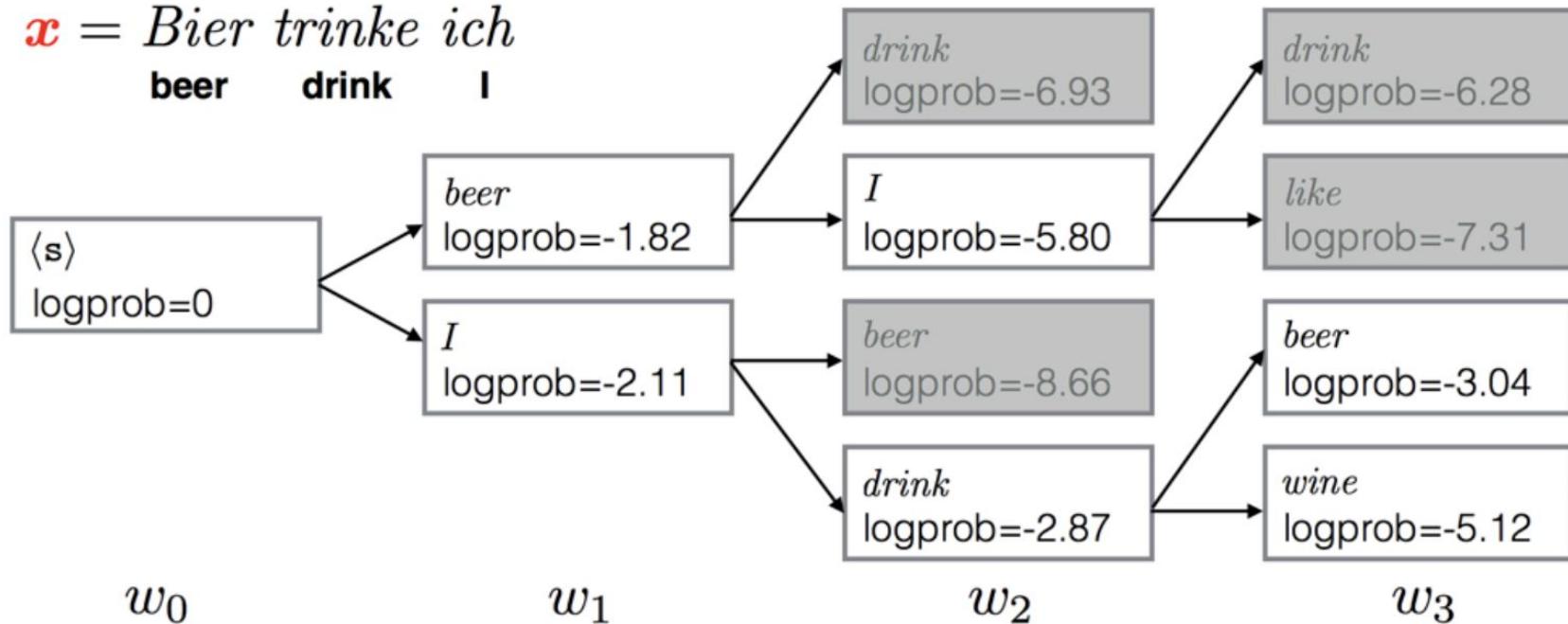


RNN Architectures: Beam Search

E.g., for $b=2$:

$x = \text{Bier trinke ich}$

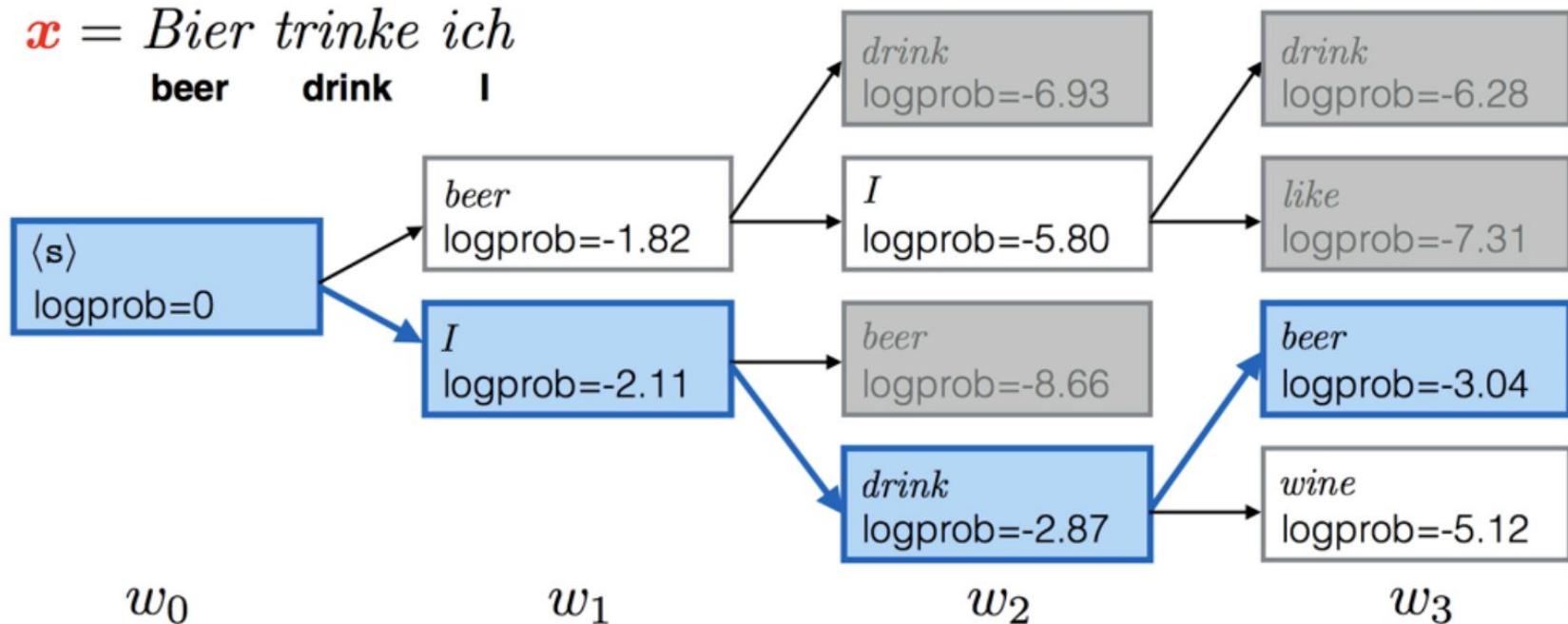
beer drink I



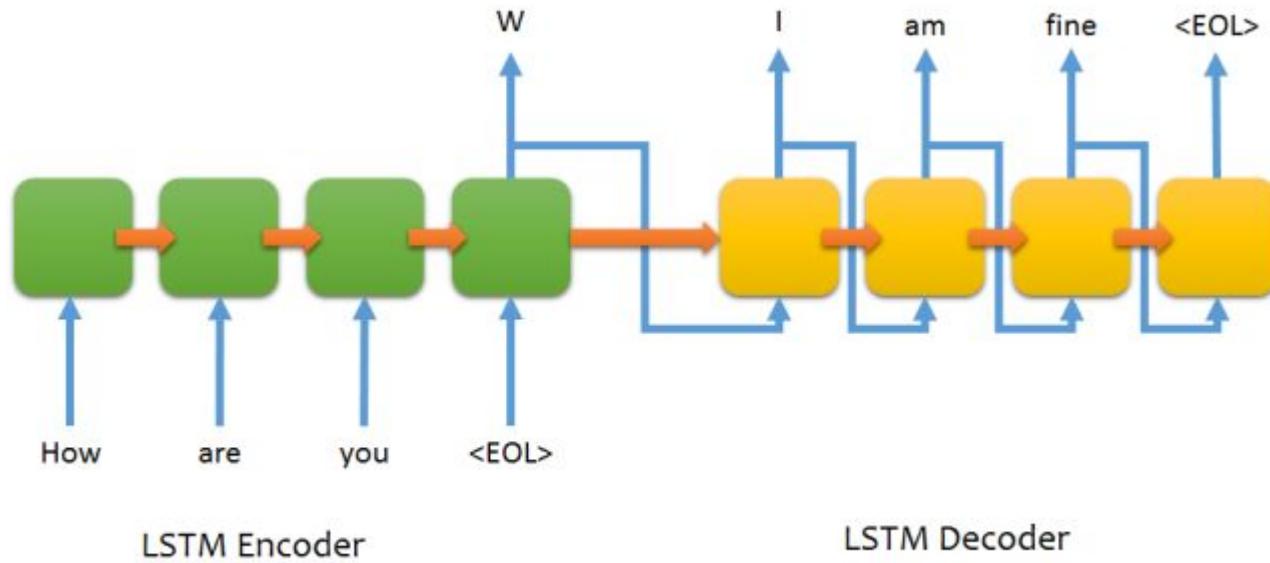
RNN Architectures: Beam Search

E.g., for $b=2$:

$\mathbf{x} = \text{Bier trinke ich}$
beer drink I

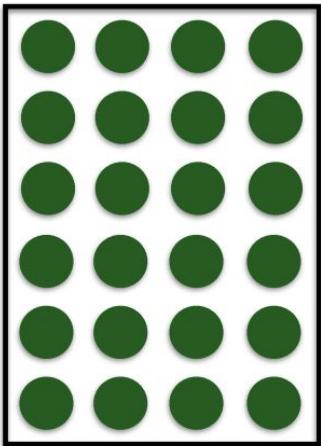


RNN Architectures: Sequence to Sequence

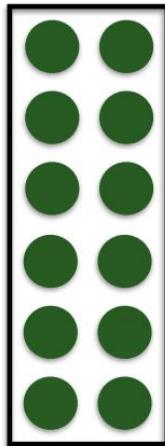


[Source](#)

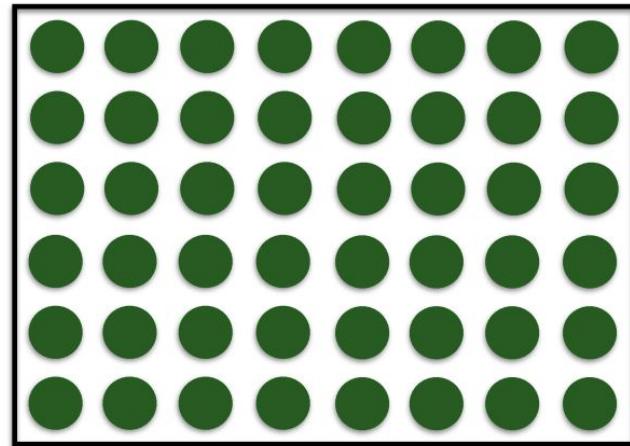
Sentences as matrices



Ich möchte ein Bier



Mach's gut

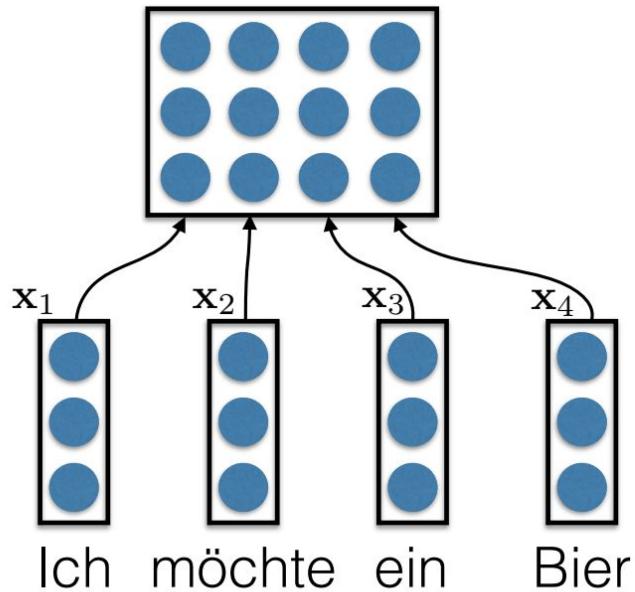


Die Wahrheiten der Menschen sind die unwiderlegbaren Irrtümer

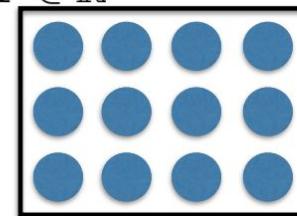
Question: How do we build these matrices?

Sentences as matrices

$$\mathbf{f}_i = \mathbf{x}_i$$

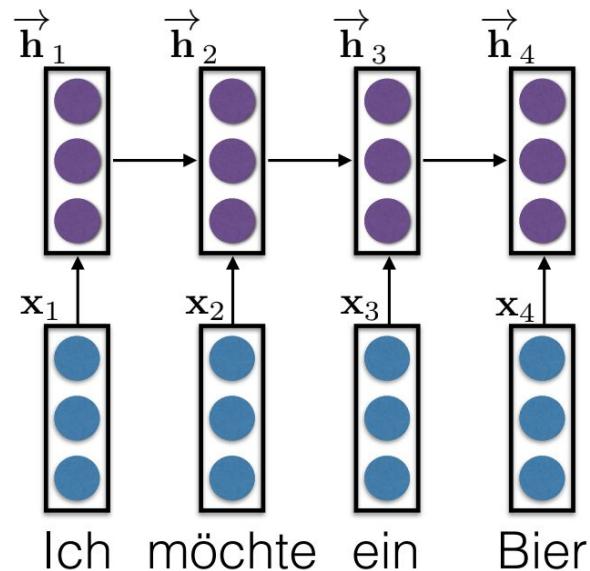


$$\mathbf{F} \in \mathbb{R}^{n \times |\mathbf{f}|}$$

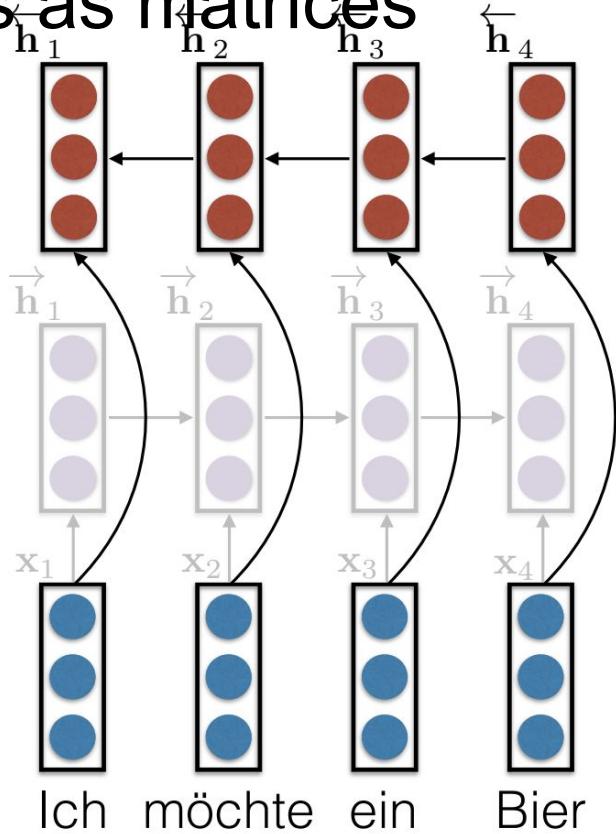


Ich möchte ein Bier

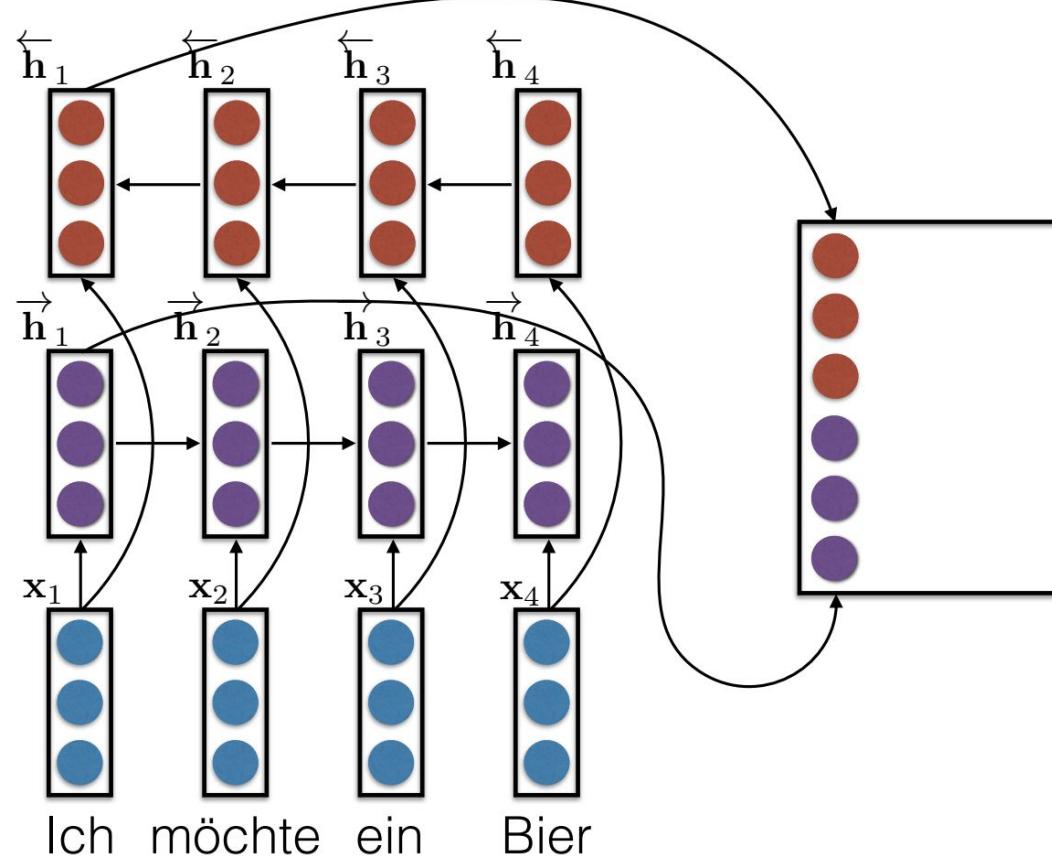
Sentences as matrices



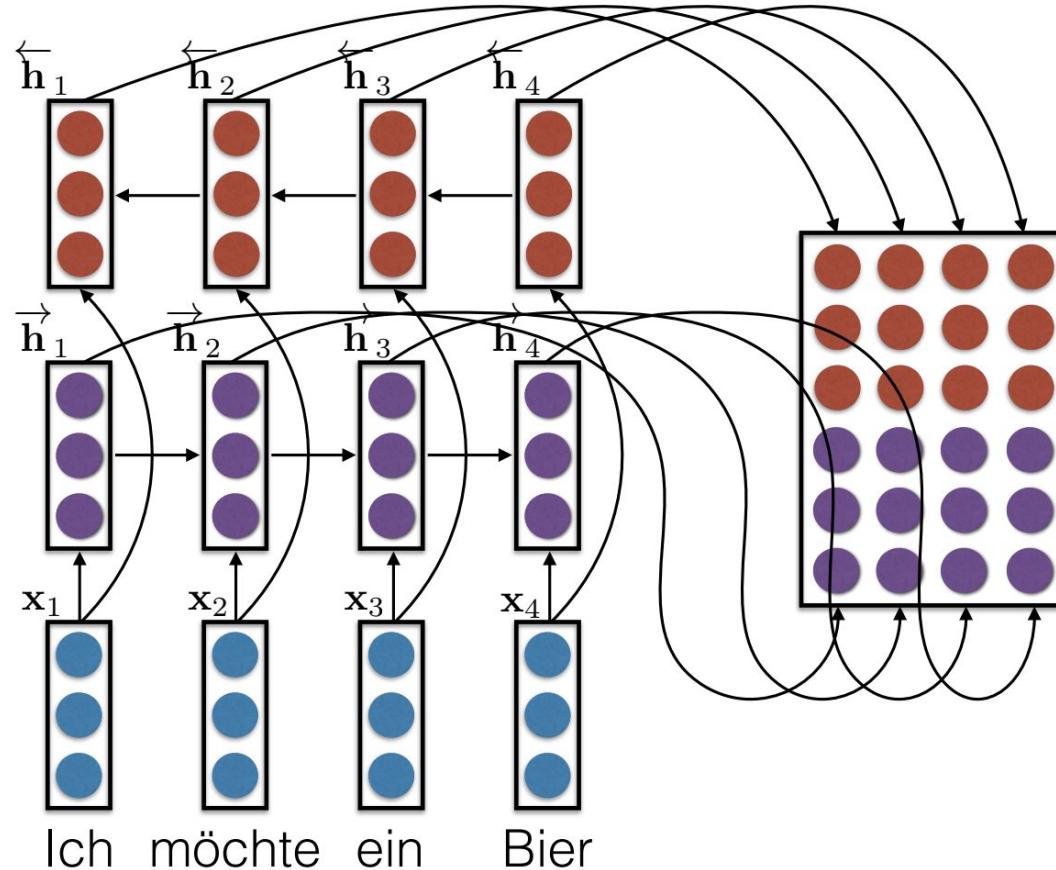
Sentences as matrices



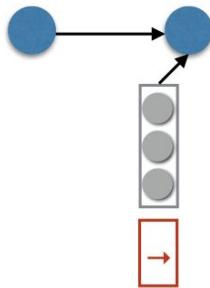
$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$



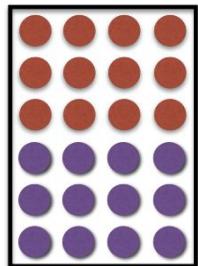
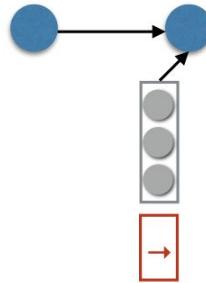
$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$



Attention

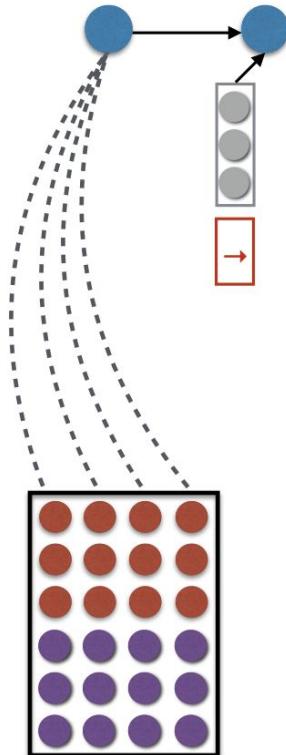


Attention



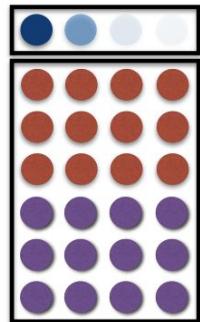
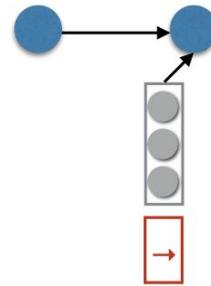
Ich möchte ein Bier

Attention



Ich möchte ein Bier

Attention

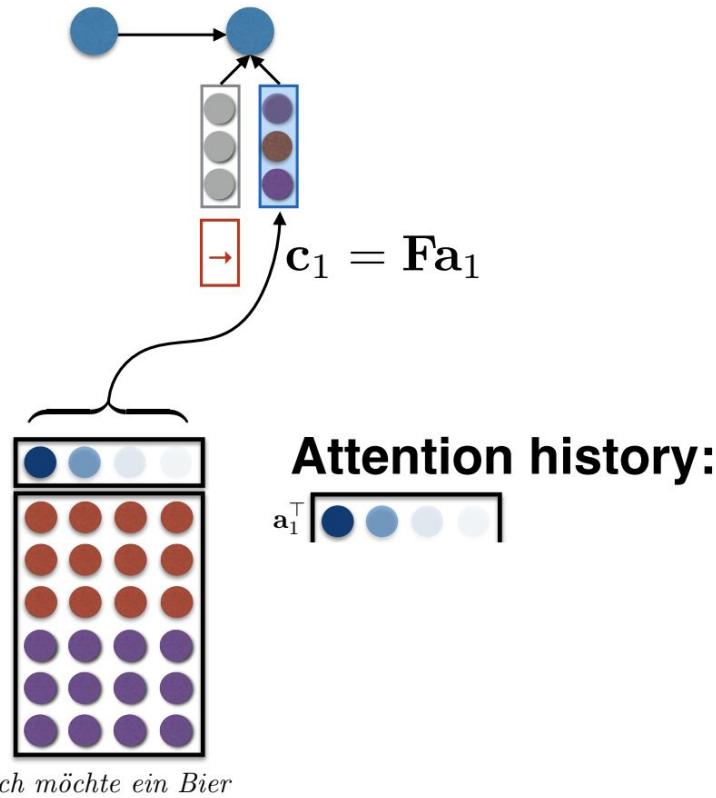


Ich möchte ein Bier

Attention history:

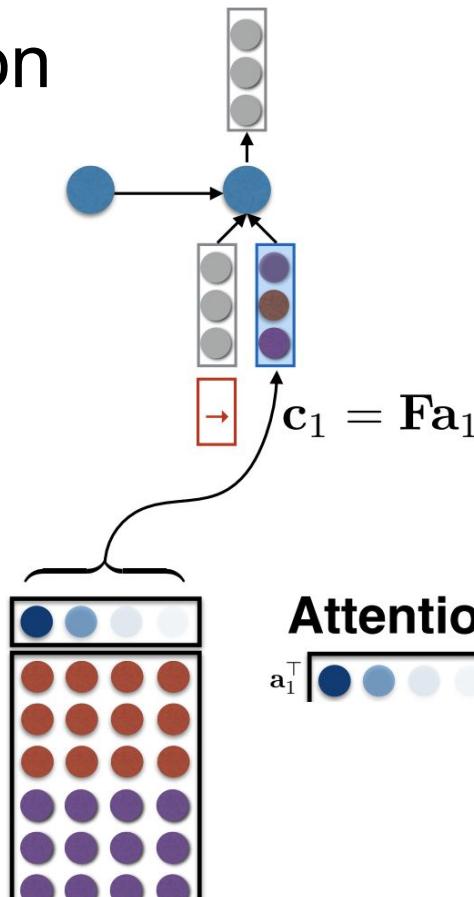
$$a_1^T [\text{dark blue} \quad \text{medium blue} \quad \text{light blue} \quad \text{white}]$$

Attention



Ich möchte ein Bier

Attention

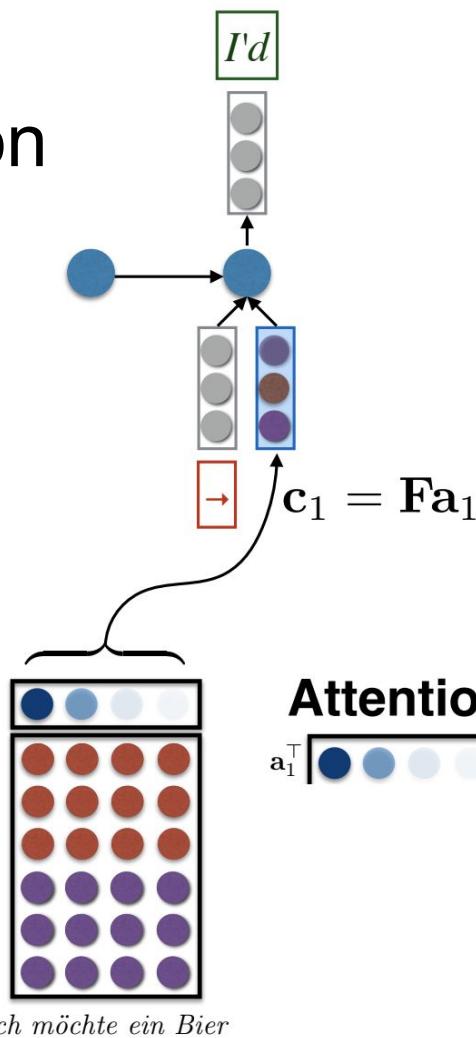


Attention history:

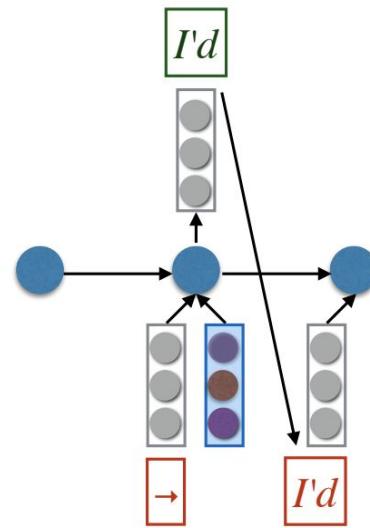
$$a_1^\top [\text{blue circle} \quad \text{blue circle} \quad \text{light blue circle} \quad \text{light blue circle}]$$

Ich möchte ein Bier

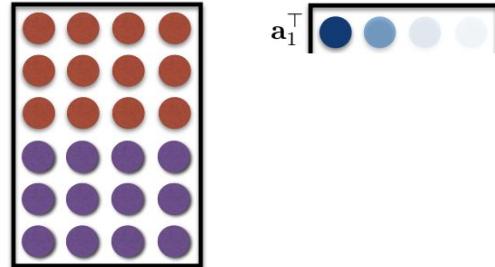
Attention



Attention

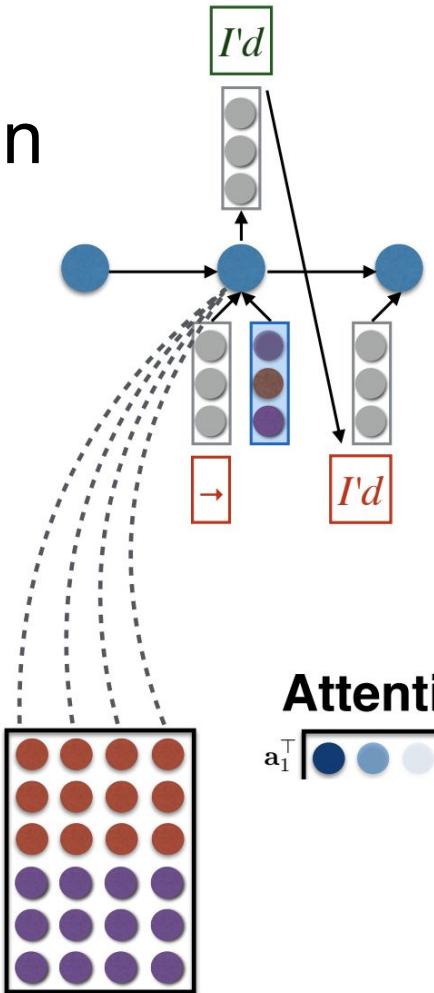


Attention history:



Ich möchte ein Bier

Attention

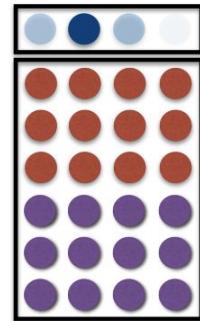
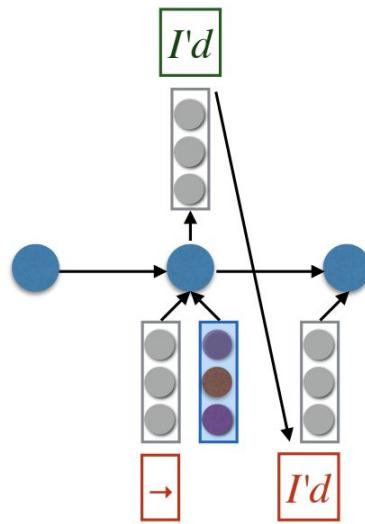


Attention history:

$$a_1^T \quad \boxed{\bullet \bullet \bullet \bullet}$$

Ich möchte ein Bier

Attention

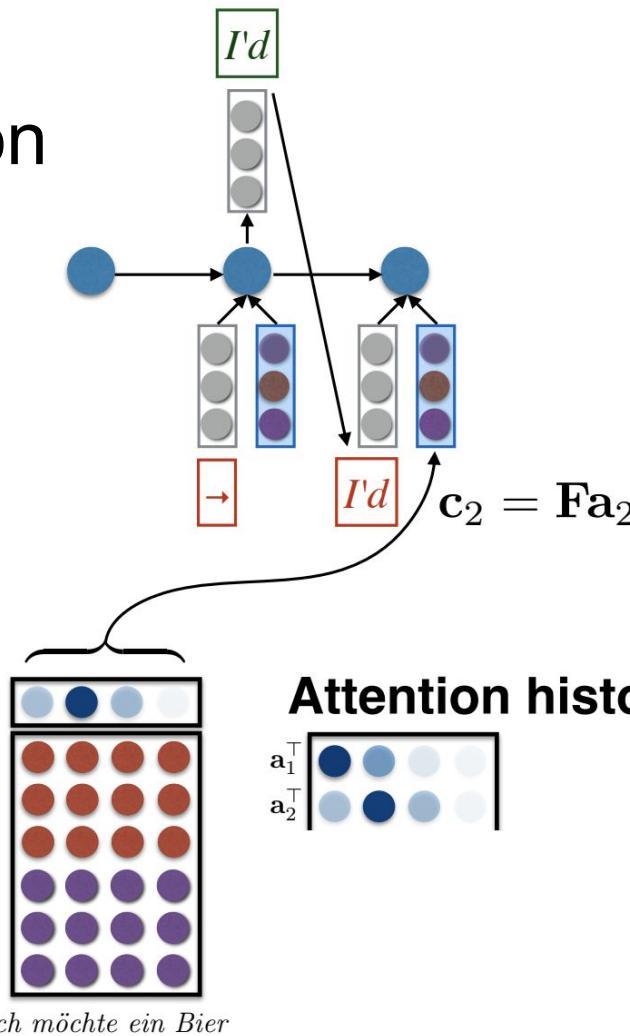


Attention history:

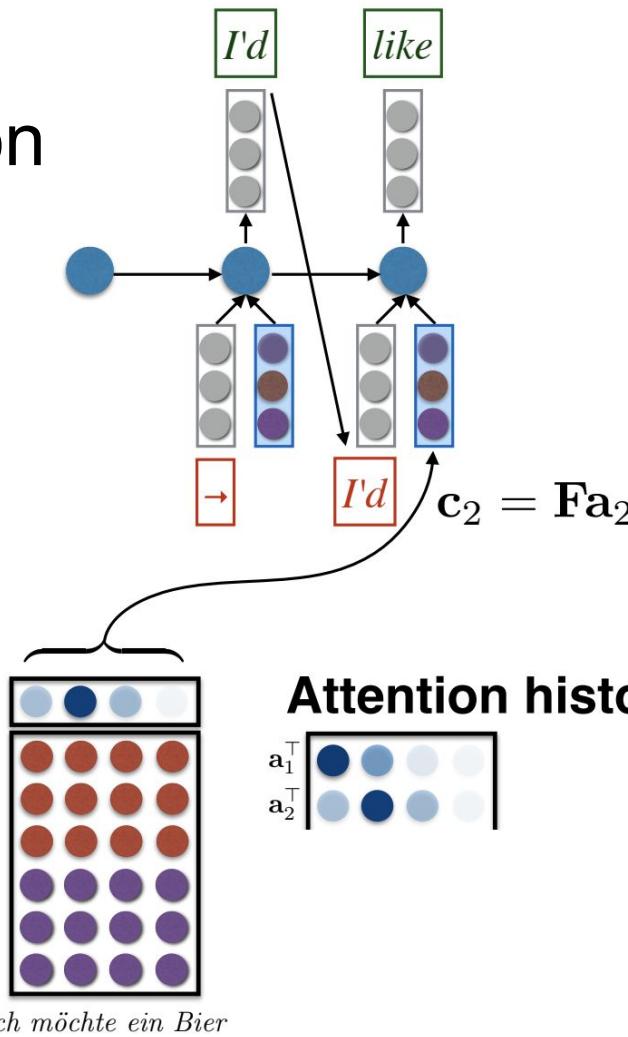
$$\begin{matrix} a_1^\top & | \\ a_2^\top & | \end{matrix} \begin{matrix} \text{dark blue} & \text{light blue} & \text{white} & \text{white} \\ \text{light blue} & \text{dark blue} & \text{white} & \text{white} \end{matrix}$$

Ich möchte ein Bier

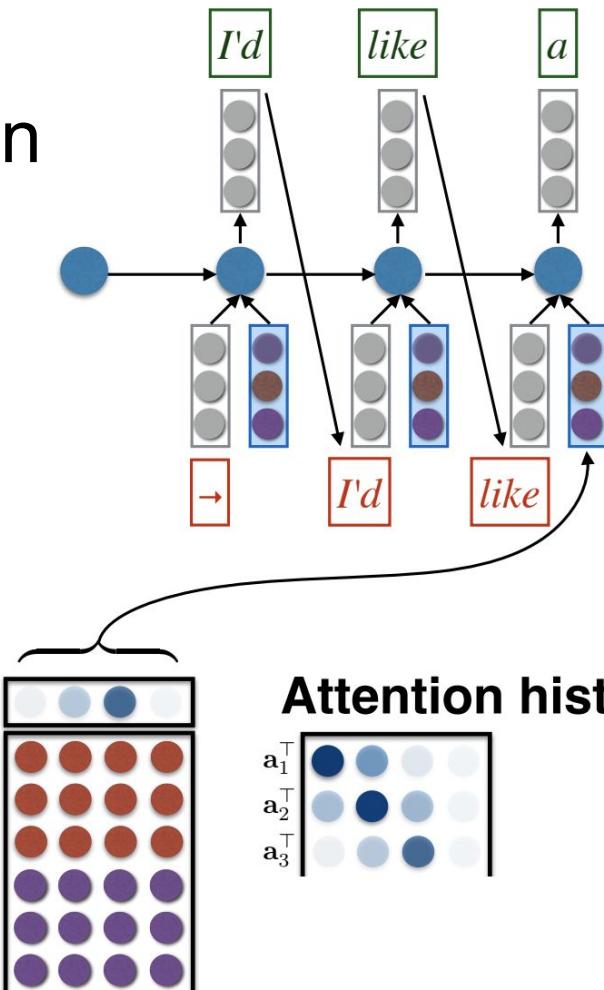
Attention



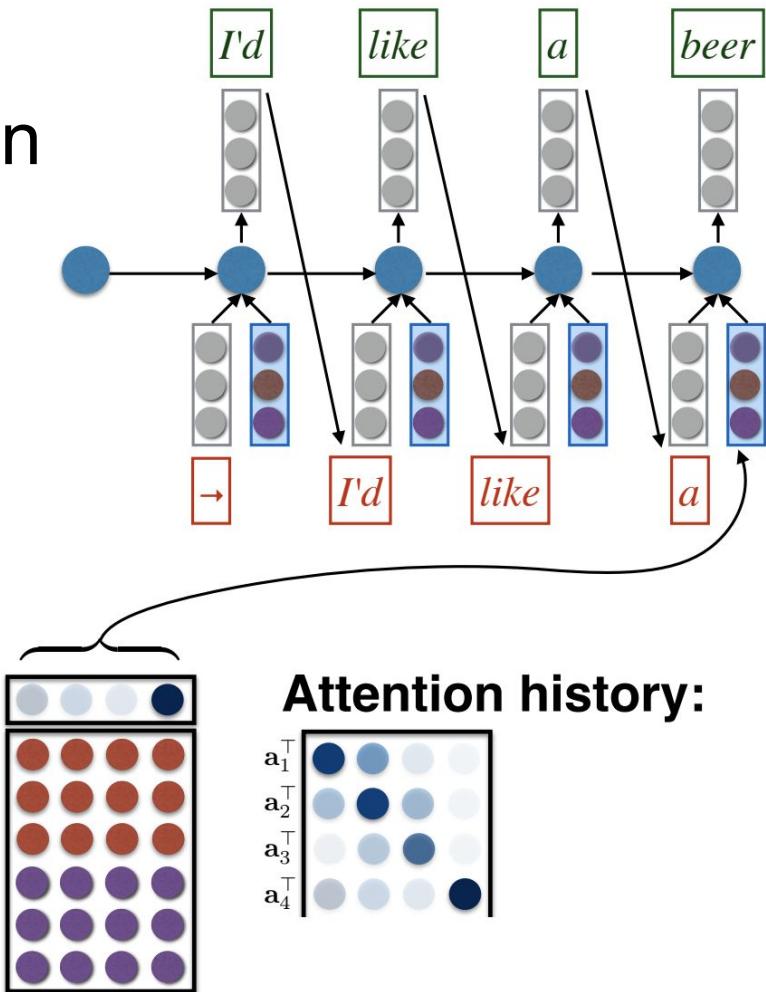
Attention



Attention

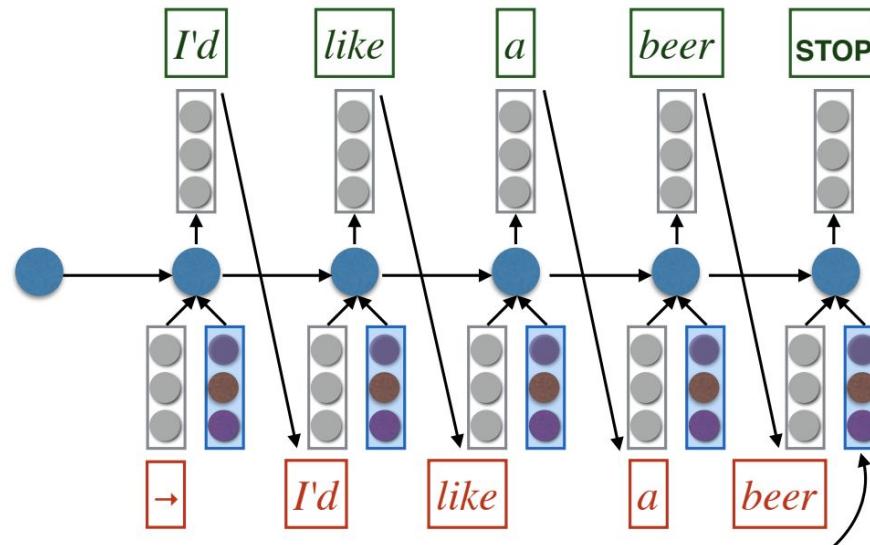


Attention

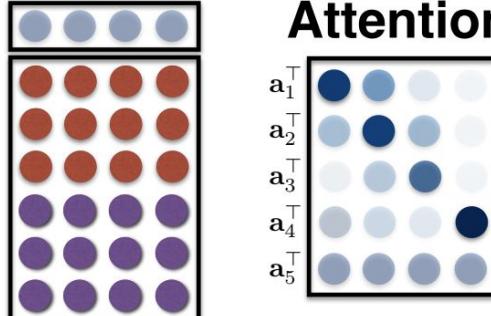


Ich möchte ein Bier

Attention



Attention history:



Ich möchte ein Bier

Outline

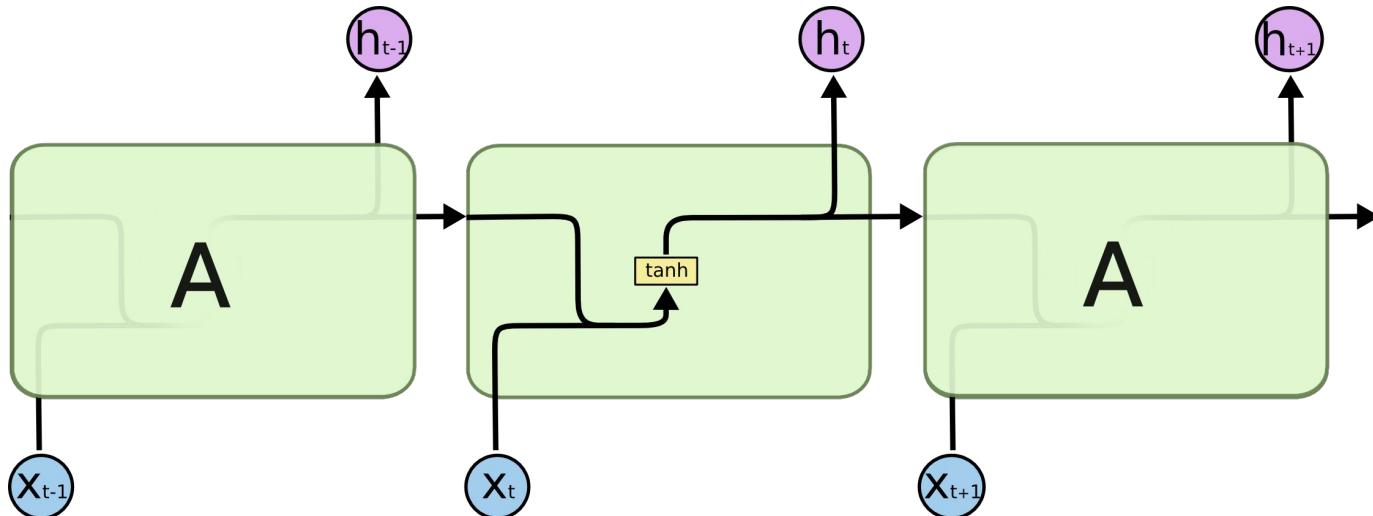
Motivation: Why Sequence Models?

Recurrent Neural Networks: Overview

Recurrent Neural Network Architectures

Types of RNN Cells

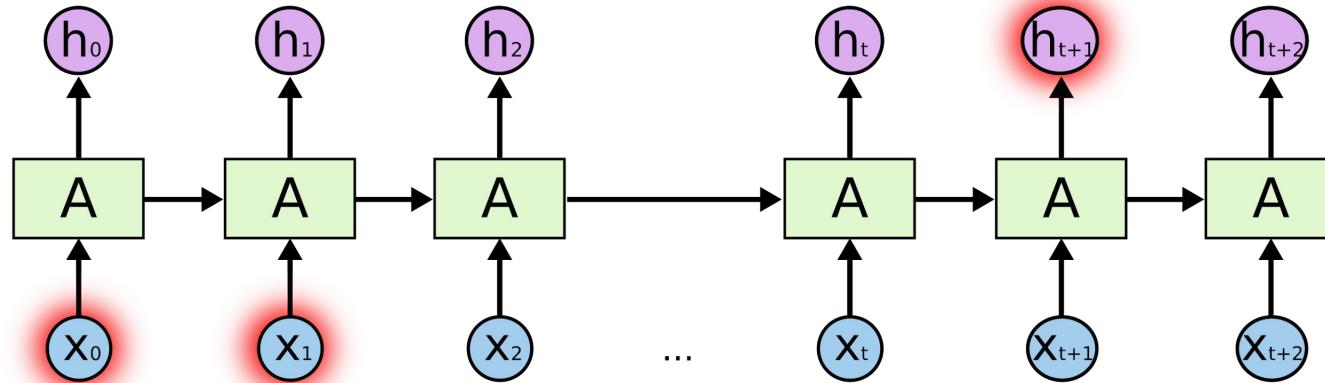
Simple RNN Cell Example



$$h_t = \tanh(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

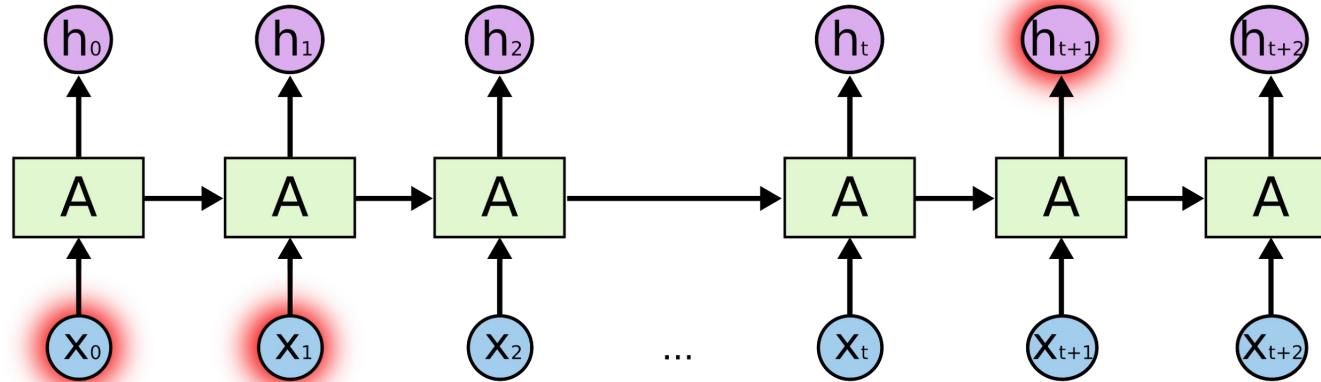
Taken from <http://colah.github.io/posts/2015-08-Understanding-LSTMs/> [3].

Problem: Vanishing Gradients



Taken from [3].

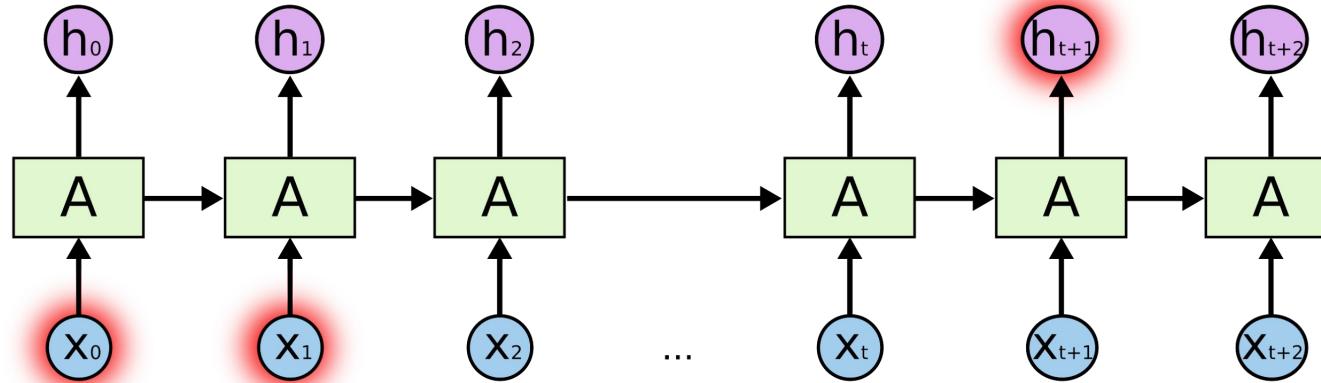
Problem: Vanishing Gradients



Using simple RNN cells, the model has difficulty learning connections between inputs and outputs that are far apart.

Taken from [3].

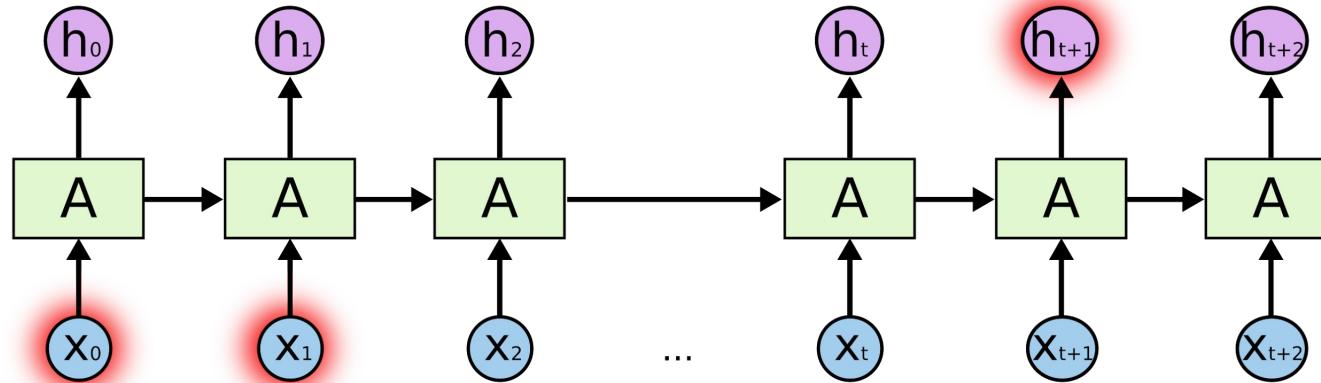
Problem: Vanishing Gradients



In other words, your model has **very small gradients**.

Taken from [3].

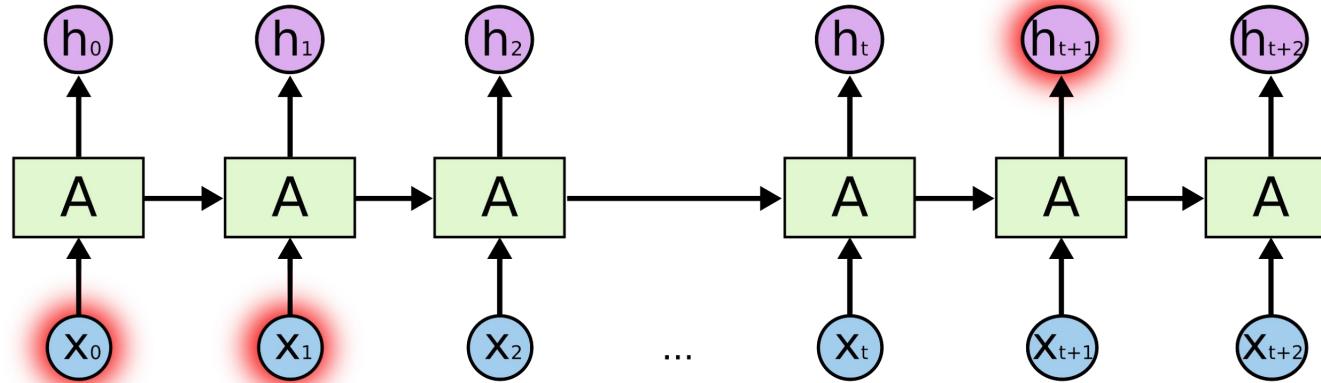
Problem: Vanishing Gradients



Gradients from $t+T$ at time t vanish exponentially with T if
at each time step $|dh(t)/dh(t-1)| < 1.0$

Taken from [3].

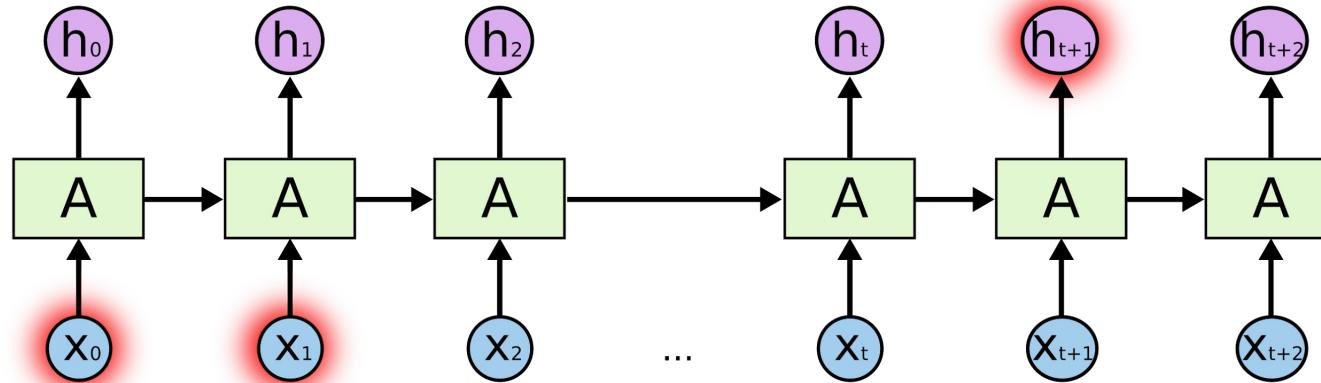
Problem: Vanishing Gradients



In other words, your model has **very small gradients**.
It will learn **extremely slowly**.

Taken from [3].

Problem: Gradients Exploding

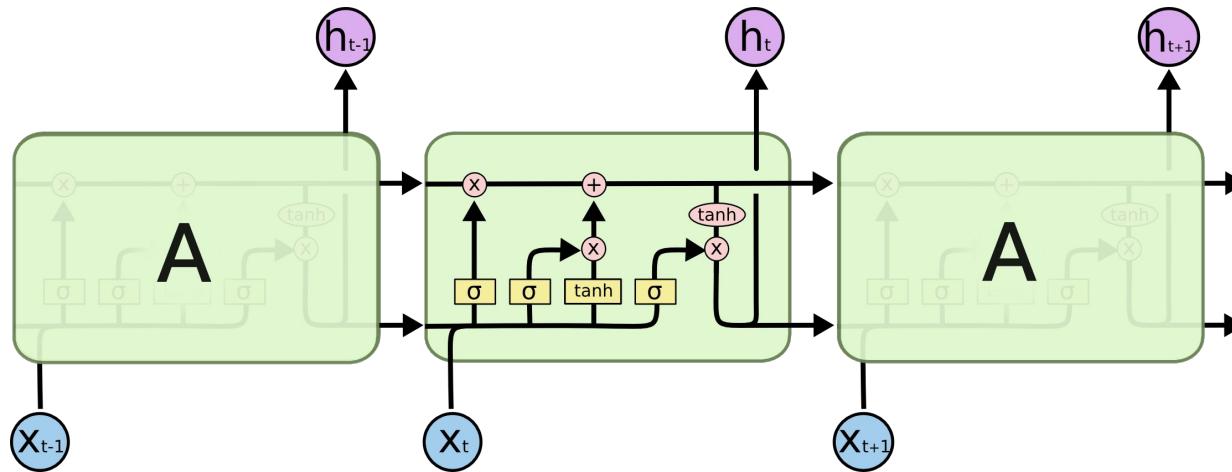


Opposite issue! Start to see NaNs during training

Reference: [On the difficulty of training recurrent neural networks](#)

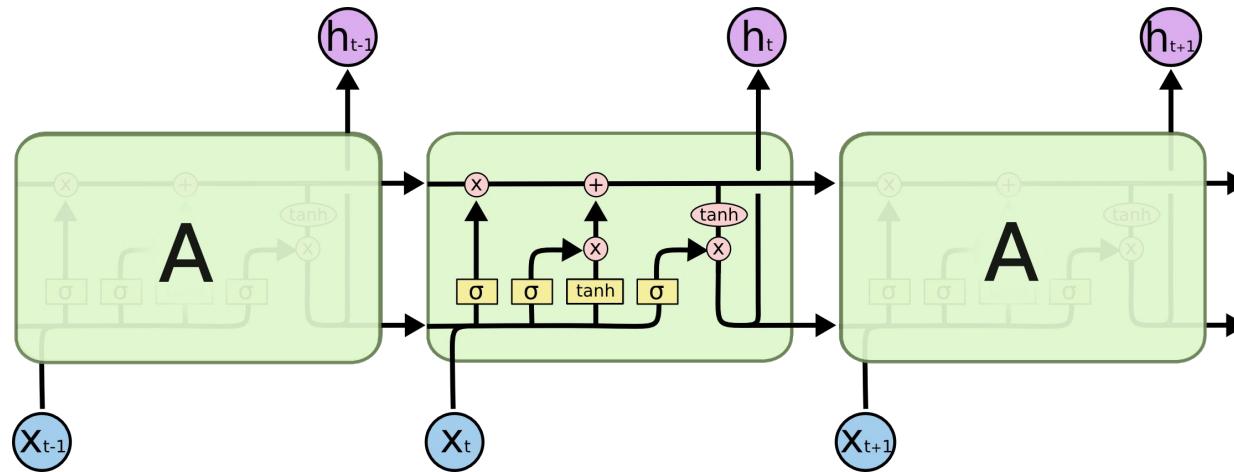
Taken from [3].

A solution: Long Short Term Memory Cells



Taken from [3].

A solution: Long Short Term Memory Cells



Neural Network
Layer



Pointwise
Operation



Vector
Transfer



Concatenate

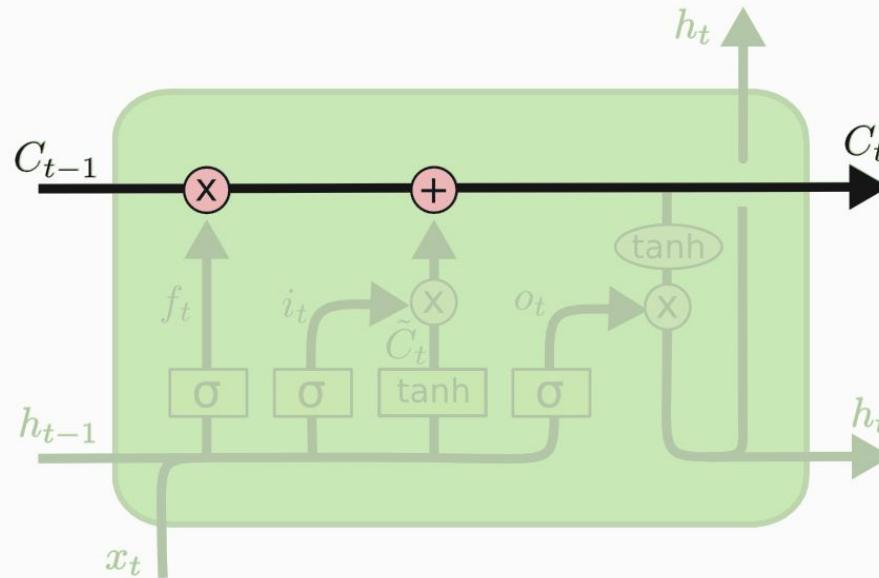


Copy

Taken from [3].

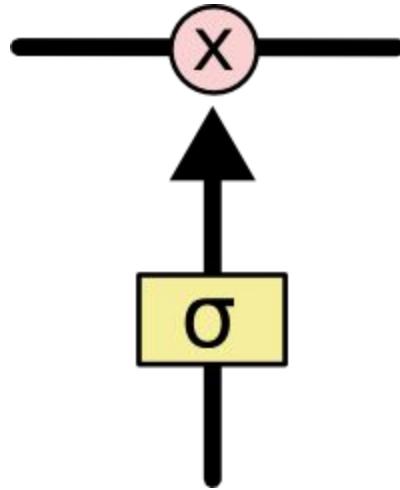
LSTM - Cell State

- “Conveyer belt”
- LSTM can “add” or “remove” information to cell state via *gates*.



Taken from [3].

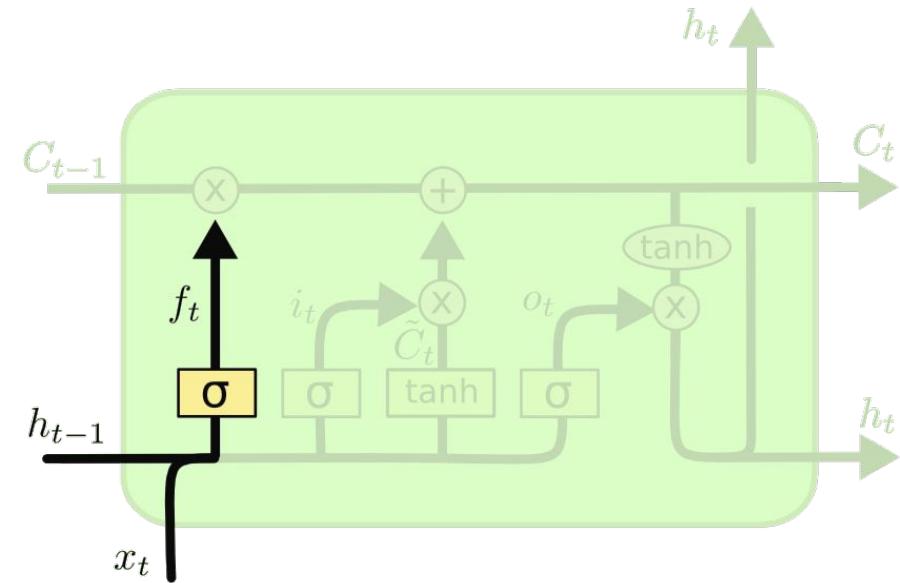
Gates: Optionally Let Information Through



- Element-wise sigmoid and element-wise multiplication
- Differentiable: trainable
- $\sigma(\cdot) \in [0, 1]$

Taken from [3].

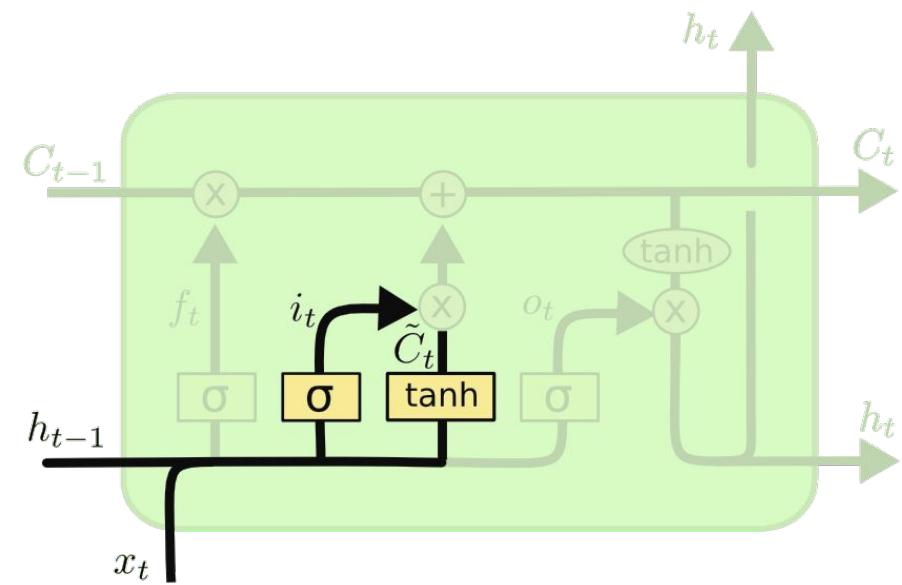
LSTM Cell: Forget Gate



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Taken from [3].

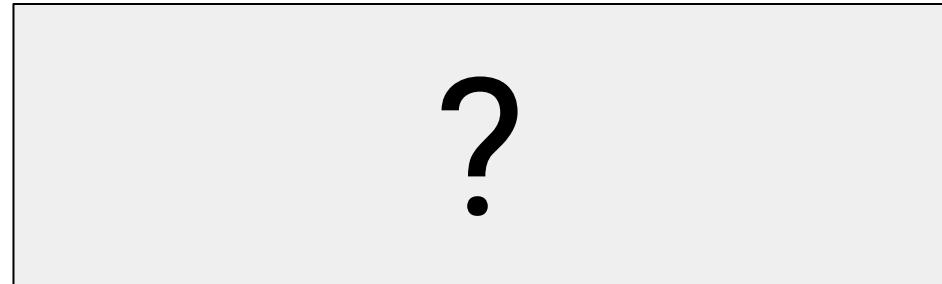
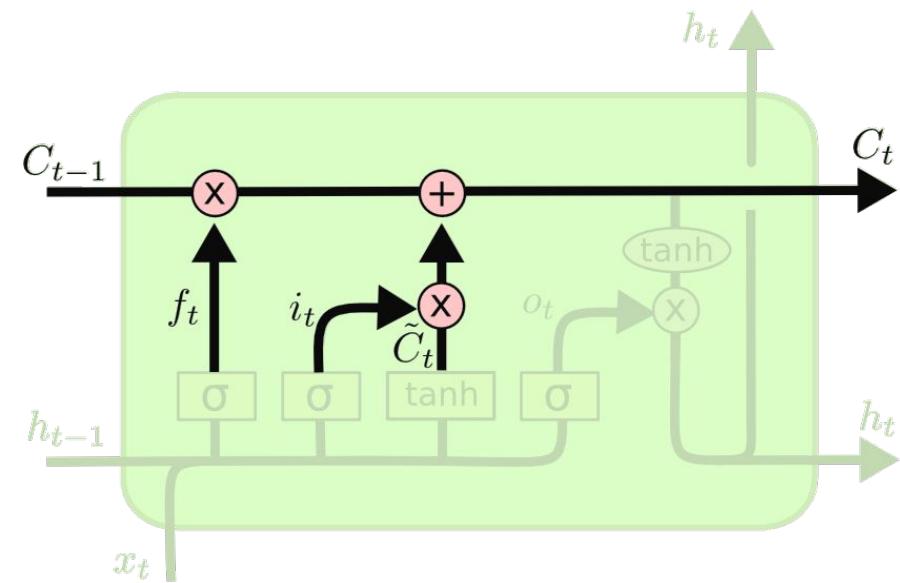
LSTM Cell: Input Gate and Candidate State



?

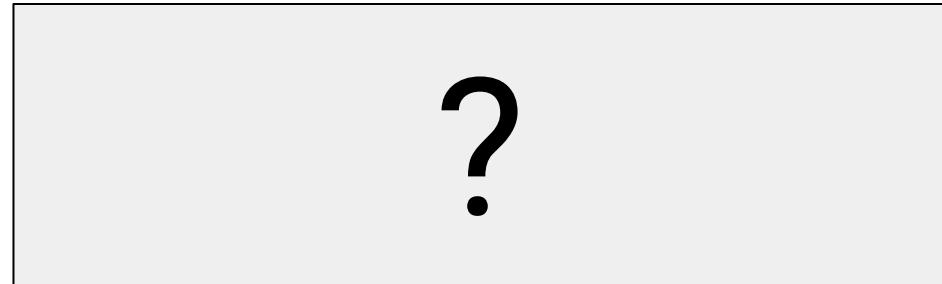
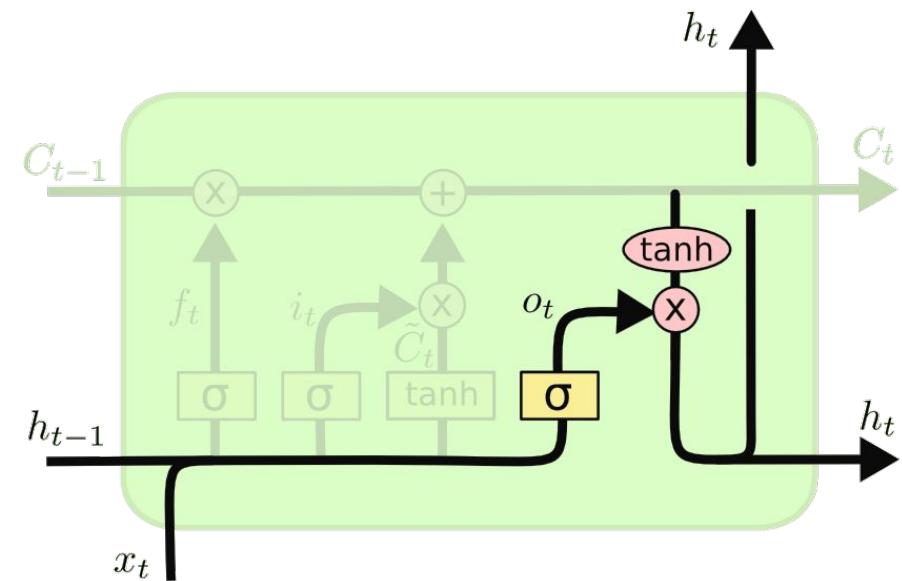
Taken from [3].

LSTM Cell: Update the Cell State



Taken from [3].

LSTM Cell: Output Gate



- Output filtered version of cell state
- $\tanh(\cdot) \in [-1, 1]$

Taken from [3].

LSTM Cell: Take Away

LSTM cells propagate state on a '**conveyor belt.**'

LSTM cells contain a series of **gates** that alter state at each time step.

Some gates **add** to the state while other gates **forget**.

LSTM cells can propagate information over many time steps.

They can also erase state quickly.

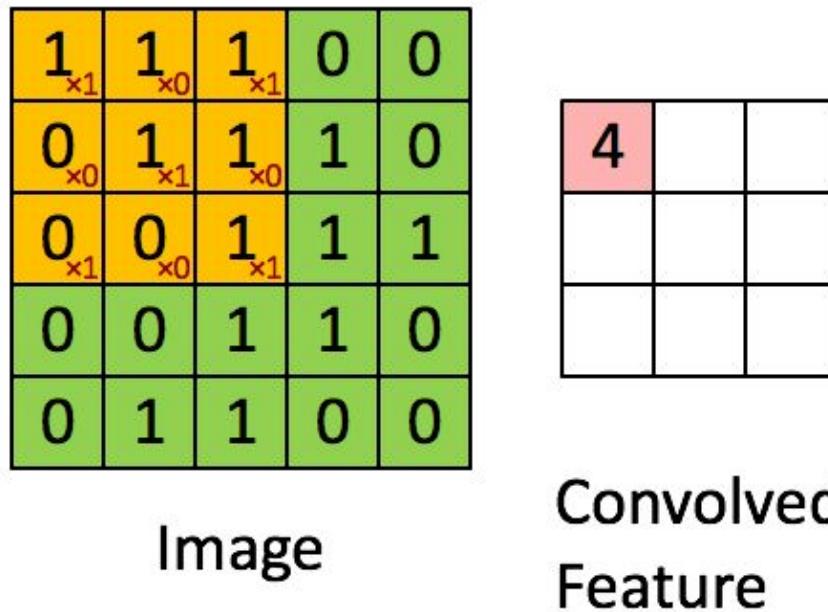
Empirical results

- LSTM is probably the most popular, but:
- People use different variations of LSTM, gated recurrent units (GRU), etc.

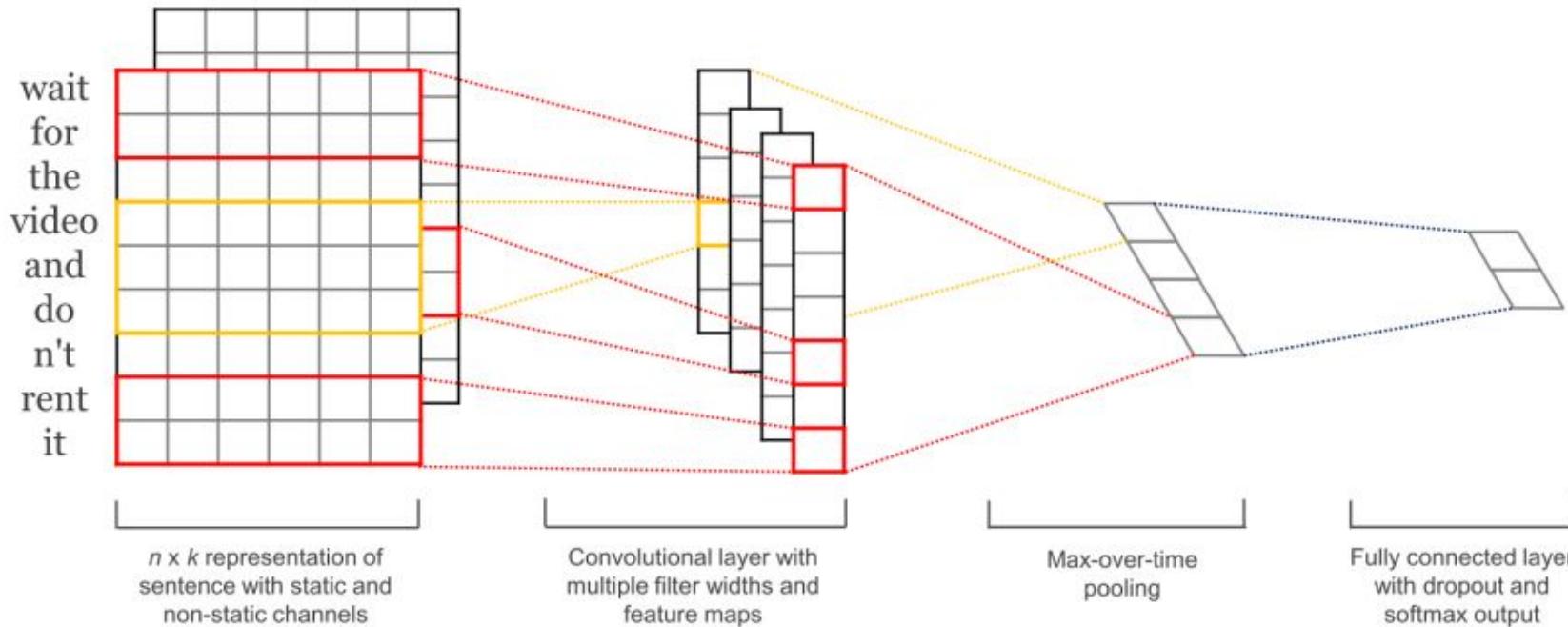
RNN Cells In TensorFlow

<https://www.tensorflow.org/tutorials/recurrent>

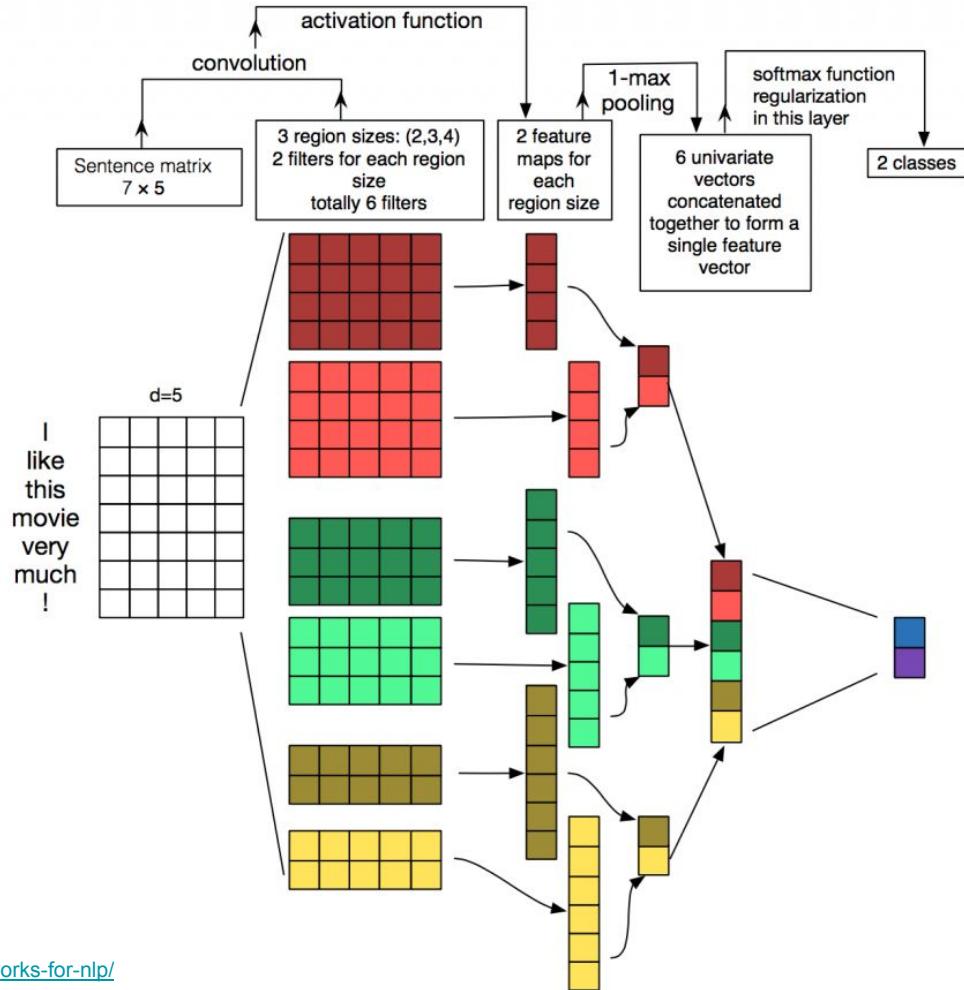
Convolutional Neural Networks



CNN

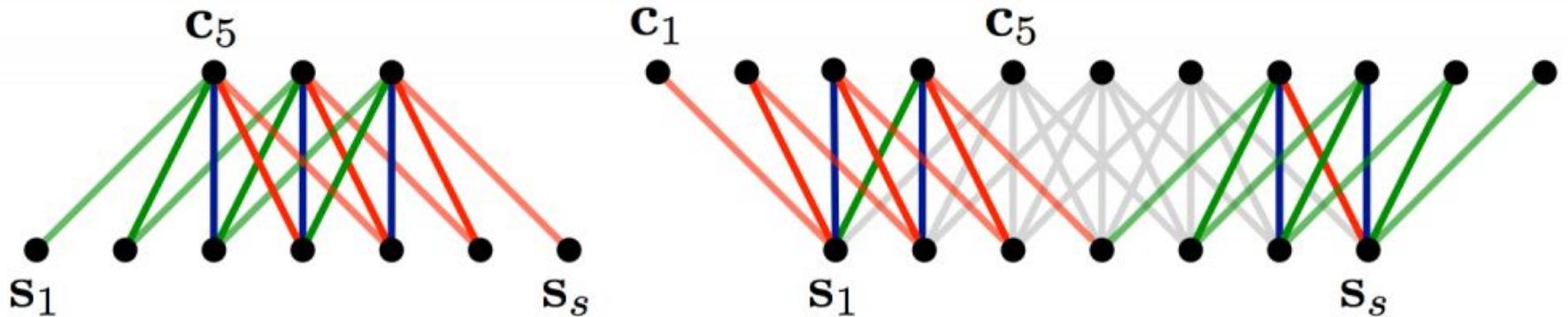


Convolutional Neural Networks on Text



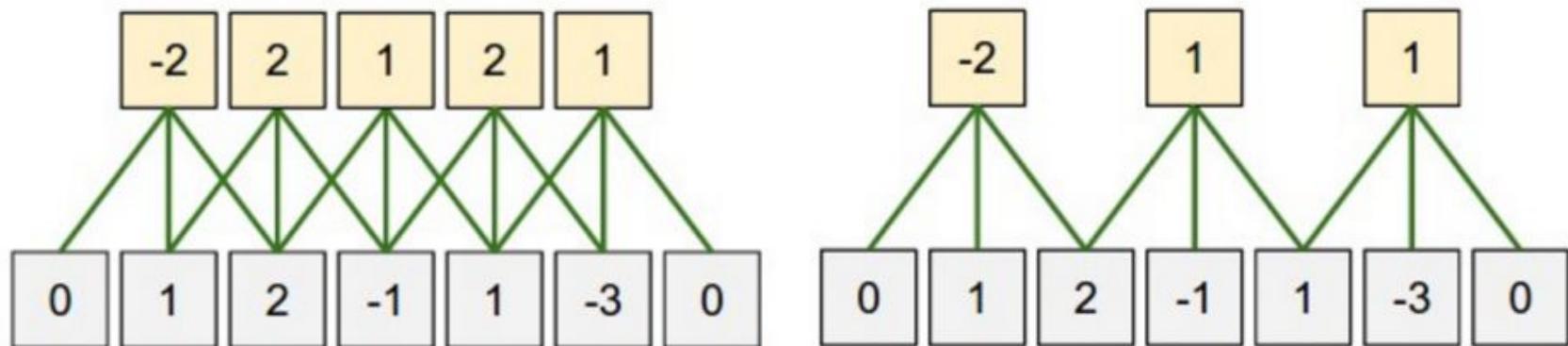
Convolutional Neural Networks on Text

Narrow vs Wide Convolution



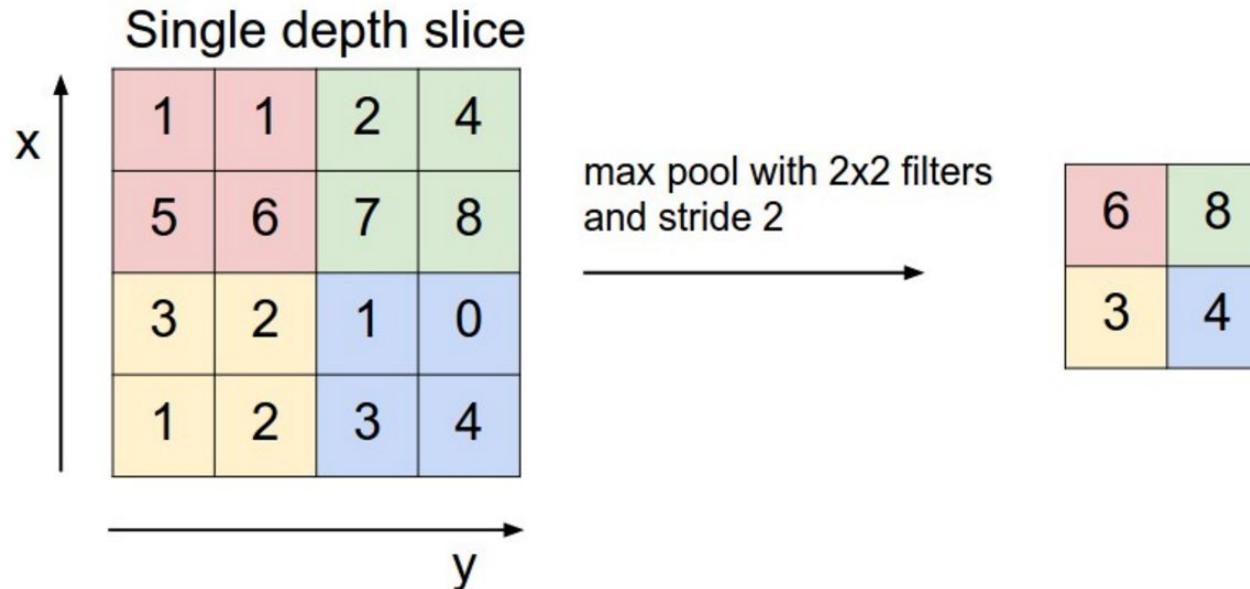
Convolutional Neural Networks on Text

Stride Size



Convolutional Neural Networks on Text

Pooling Layers



Convolutional Neural Networks on Text Channels

RGB for images

Word2Vec, Glove, etc for text

Case Studies

SmartReply

Text Simplification and its evaluation

All the best!

Feel free to ask me anything related to
Google, NLP, ML, DL, PhD, Career,
Online Activism, Sport :)

me@juliapro.xyz