

Natural Language Processing



github.com/adarsh0806/ODSC

Natural Language Processing

- Introduction to NLP
- Text Processing
- Basic NLP
- Advanced NLP

Introduction to NLP

Introduction to NLP

- Communication and Cognition
- Structured Languages
- Unstructured Text
- Applications and Challenges

Communication and Cognition

Language is...

- a medium of communication
- a vehicle for thinking and reasoning

Structured Languages

- Natural language lacks precisely defined structure

Structured Languages

- Mathematics:

$$y = 2x + 5$$

Structured Languages

- Formal Logic:

$$\text{Parent}(x, y) \wedge \text{Parent}(x, z) \rightarrow \text{Sibling}(y, z)$$

Structured Languages

- SQL:

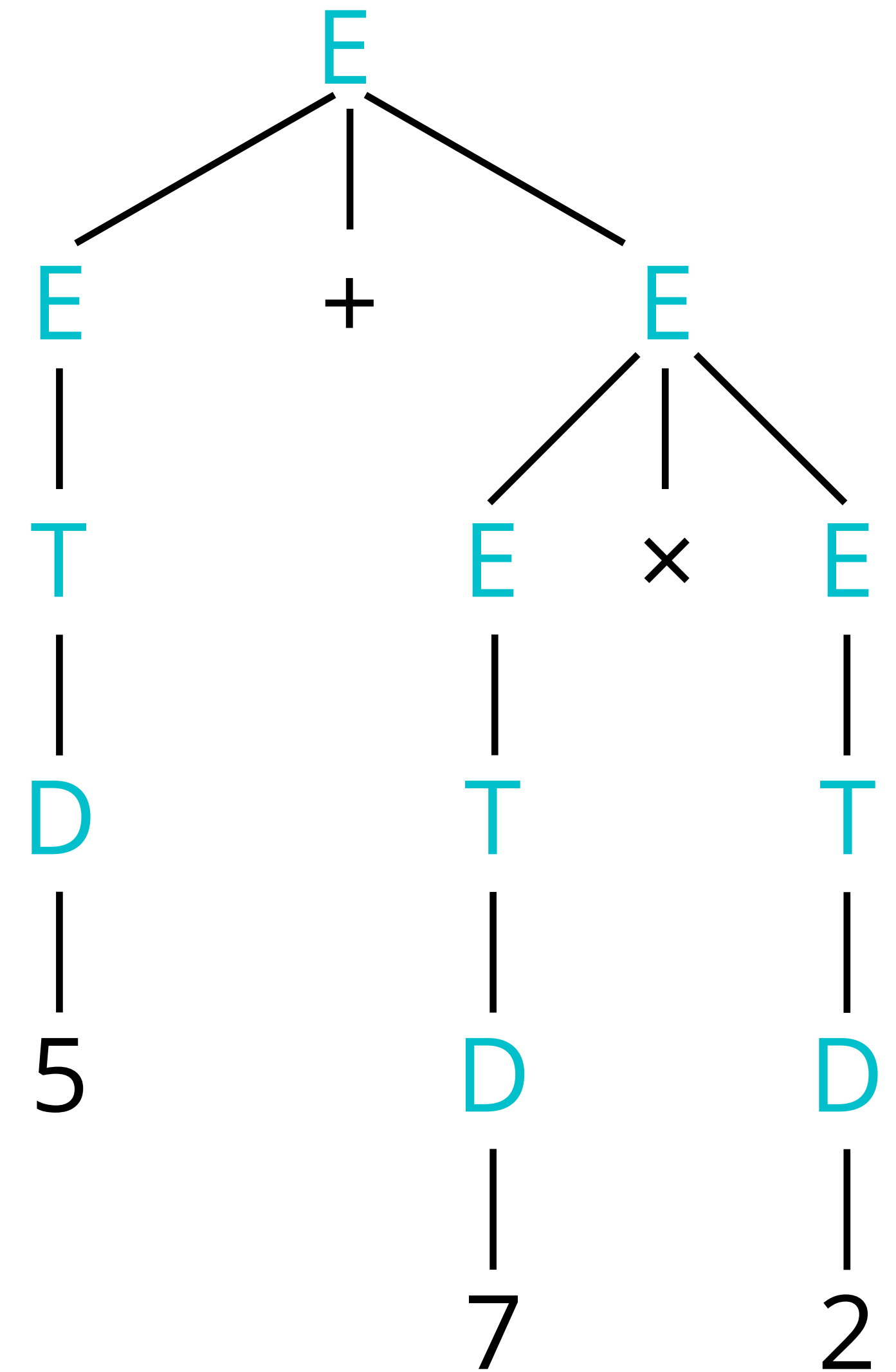
```
SELECT name, email  
FROM users  
WHERE name LIKE 'A%';
```

Grammar

- Arithmetic (single digit):

$E \rightarrow E + E \mid E - E \mid E \times E \mid E \div E \mid (E) \mid D$

$D \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$



Grammar

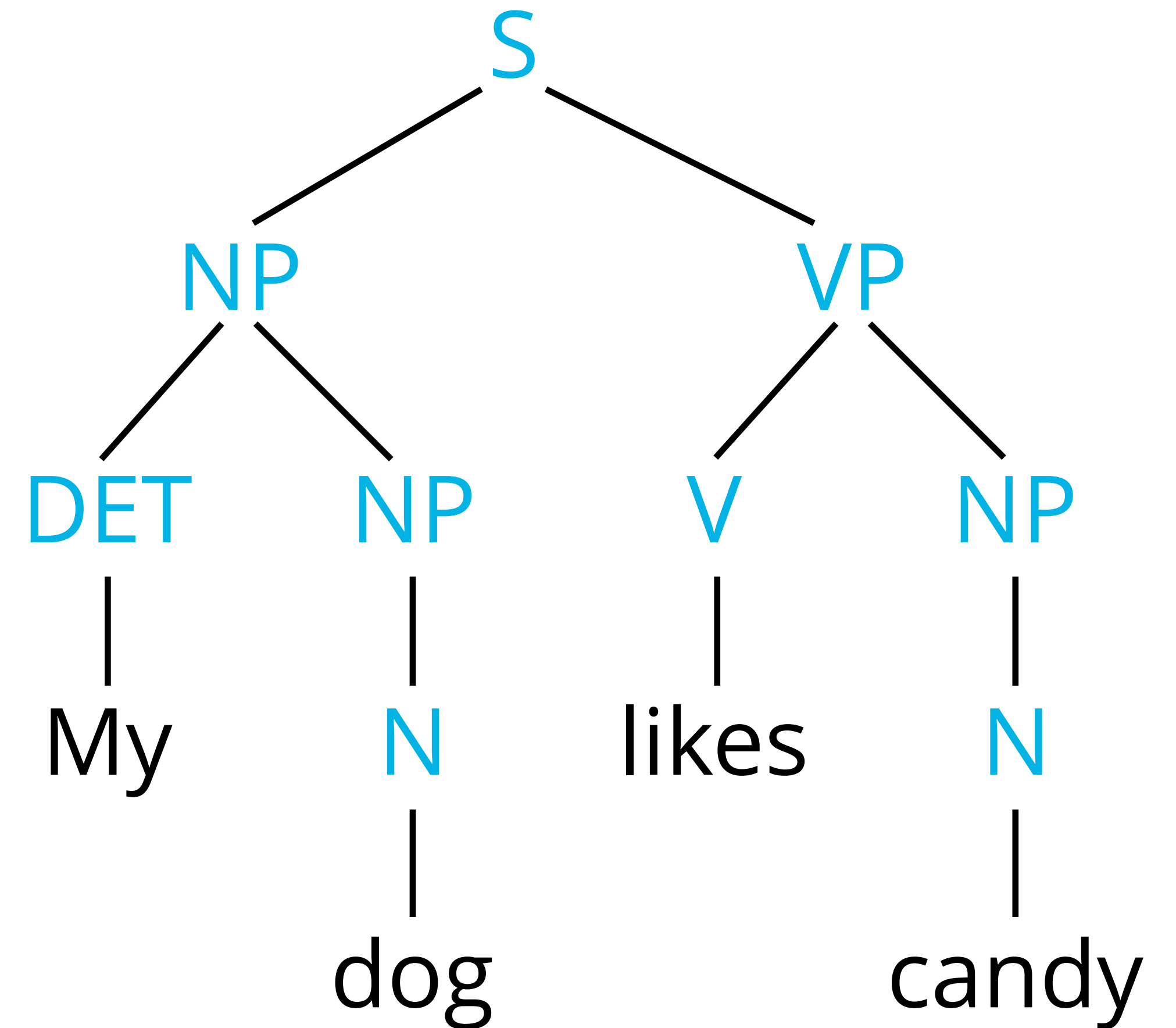
- English sentences (limited):

$S \rightarrow NP VP$

$NP \rightarrow N \mid DET NP \mid ADJ NP$

$VP \rightarrow V \mid V NP$

...



Because he was so small, Stuart was often hard to
find around the house."

verb ↗

↘ *noun*

– *Stuart Little*, E.B. White

Unstructured Text

the quick brown fox jumps over the lazy dog

Unstructured Text

jumps

the

fox

brown

over

dog

quick

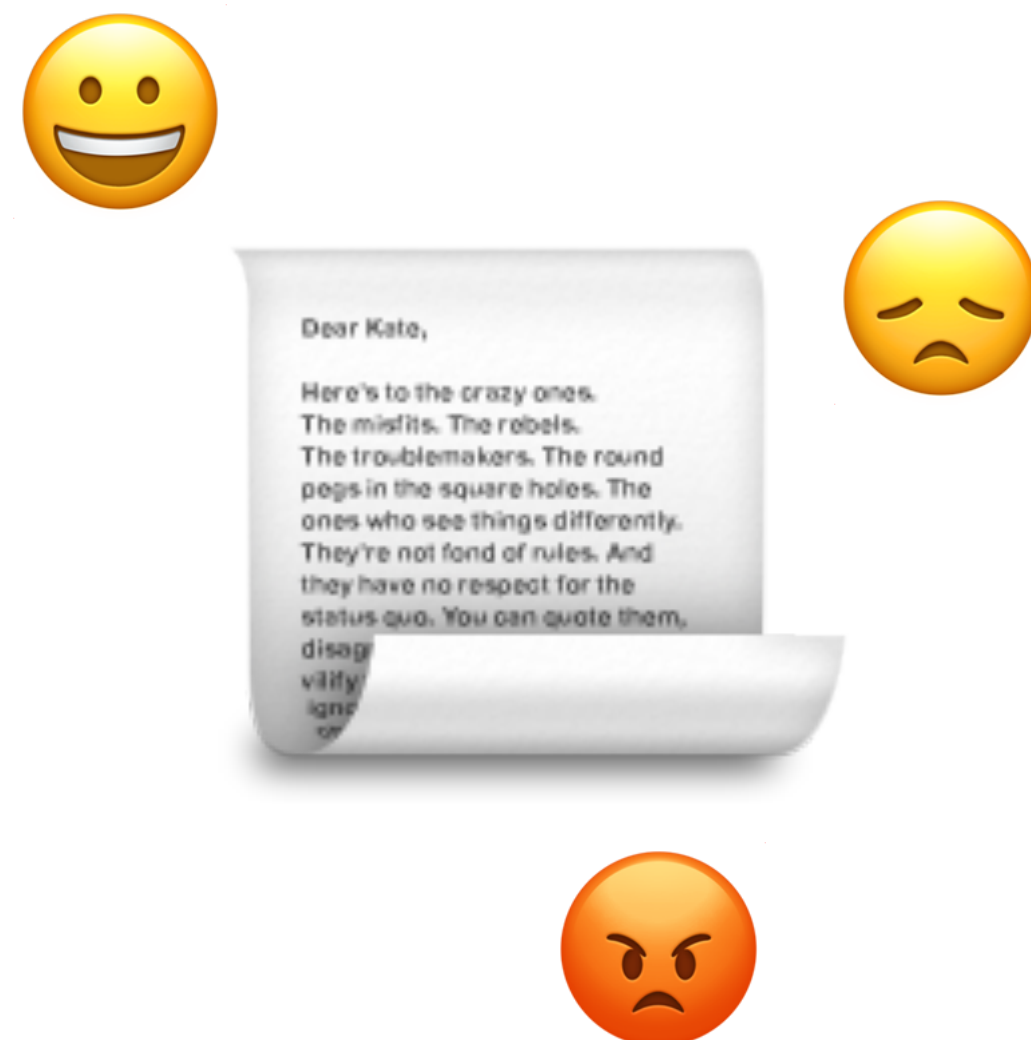
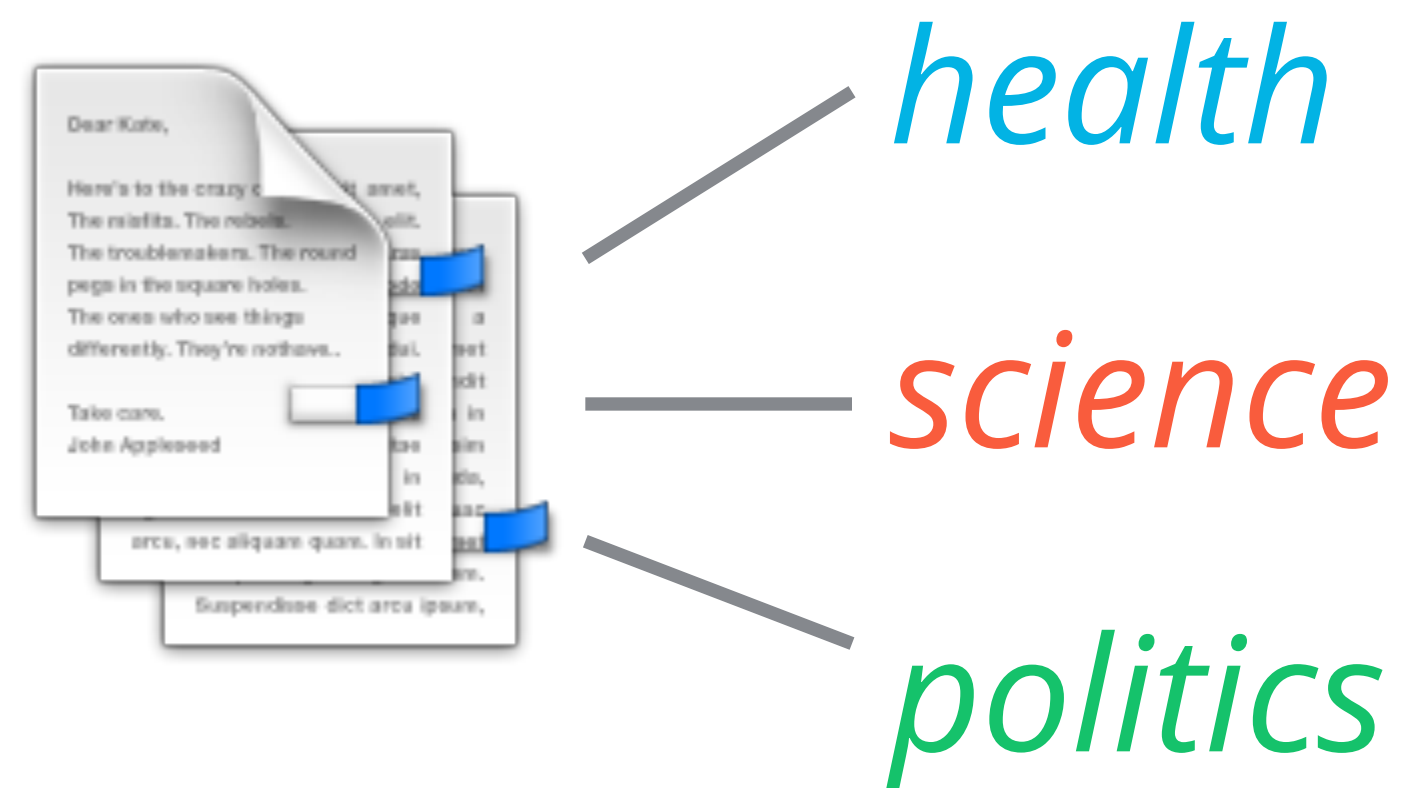
the

lazy

Unstructured Text

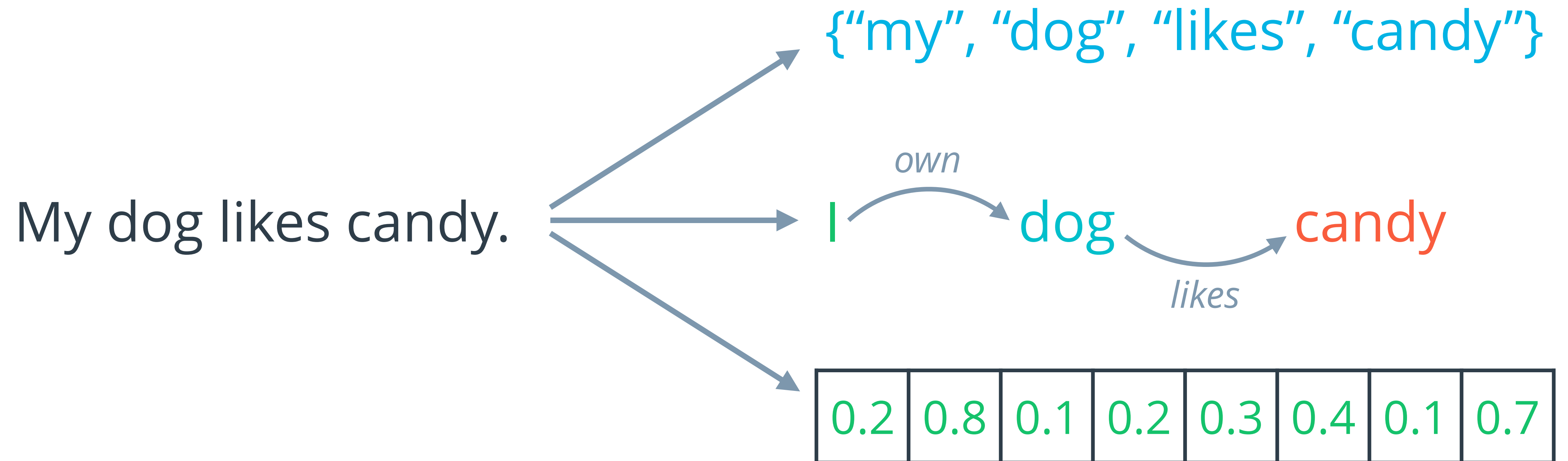
jumps
the fox
brown
over dog
+ quick the lazy -

Applications



what time is it?
¿que hora es?

Challenges: Representation



Challenges: Temporal Sequence

I want to buy a gallon of milk

water

petrol

Challenges: Context

The old Welshman came home toward daylight, spattered with candle-grease, smeared with clay, and almost worn out. He found Huck still in the bed that had been provided for him, and delirious with fever. The physicians were all at the cave, so the Widow Douglas came and took charge of the patient.

The diagram illustrates the flow of context in the text. It features several colored boxes and arrows: a light green box around 'Welshman', a light green box around 'He', an orange box around 'Huck', an orange box around 'him', an orange box around 'patient.', and a pink box around 'Widow Douglas'. An orange box also surrounds the phrase 'delirious with fever'. Arrows indicate the flow of context: one from 'He' to 'Huck', one from 'him' to 'Huck', and one from 'patient.' to 'Widow Douglas'.

—*The Adventures of Tom Sawyer*, Mark Twain

"Mary went back home. ..."



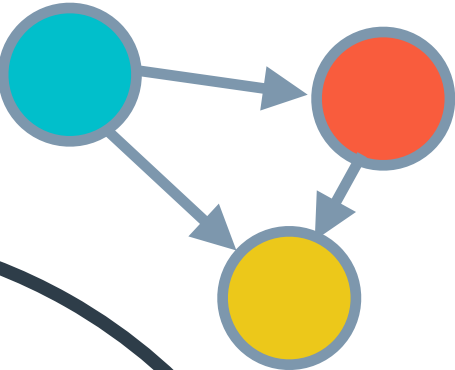
Process

{<"mary", "go", "home">, ... }

Transform

{<0.4, 0.8, 0.3, 0.1, 0.7>, ... }

Analyze



Predict



Present

Text Processing

Text Processing

- Tokenization
- Stop Word Removal
- Stemming and Lemmatization

Tokenization

“Jack and Jill went up the hill” → <“jack”, “and”, “jill”,
“went”, “up”, “the”, “hill”>

Tokenization

"No, she didn't do it."

<"no,", "she", "didn",
"", "t", "do" "it", ".">

<"no", "she", "didnt",
"do", "it">



Tokenization

Big money behind big special effects tends to suggest a big story. Nope, not here. Instead this huge edifice is like one of those over huge luxury condos that're empty in every American town, pretending as if there's a local economy huge enough to support such.

—Rotten Tomatoes

<"big", "money", "behind", "big", "special", "effects", "tends", "to", "suggest", "big", "story", "nope", "not", "here", "instead", "this", "huge", "edifice", "is", "like", "one", "of", "those", "over", "huge", "luxury", "condos", "that", "re", "empty", "in", "every", "american", "town", "pretending", "as", "if", "there", "local", "economy", "huge", "enough", "to", "support", "such">

<"big", "money", "behind", "big", "special", "effects", "tends", "to", "suggest", "big", "story">

<"nope", "not", "here">

<"instead", "this", "huge", "edifice", "is", "like", "one", "of", "those", "over", "huge", "luxury", "condos", "that", "re", "empty", "in", "every", "american", "town", "pretending", "as", "if", "there", "local", "economy", "huge", "enough", "to", "support", "such">



Stop Word Removal

wristwatch invented 1904 Louis Cartier.

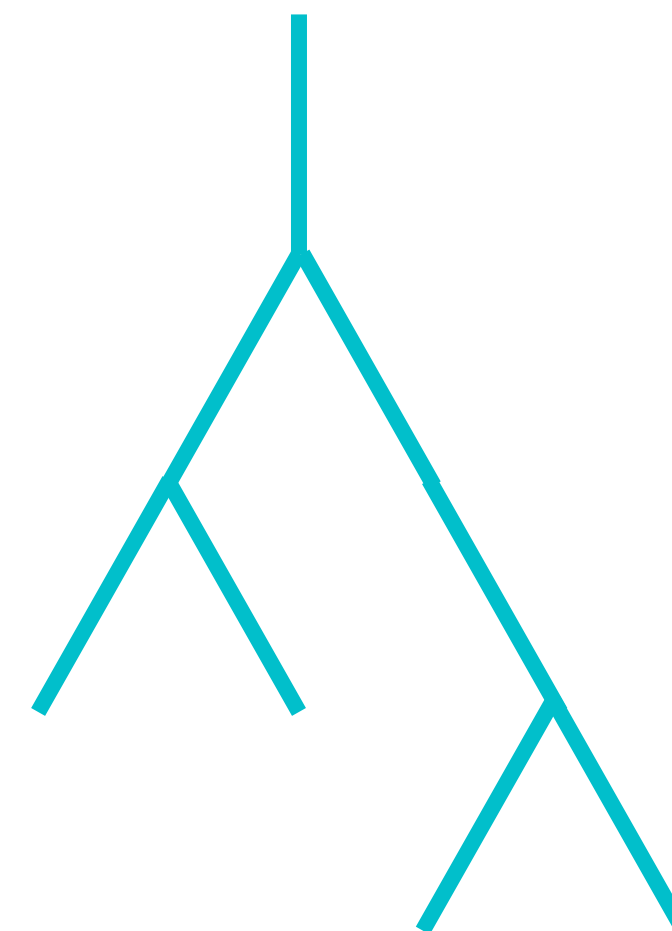
Stemming

branching

branched

branches

branch

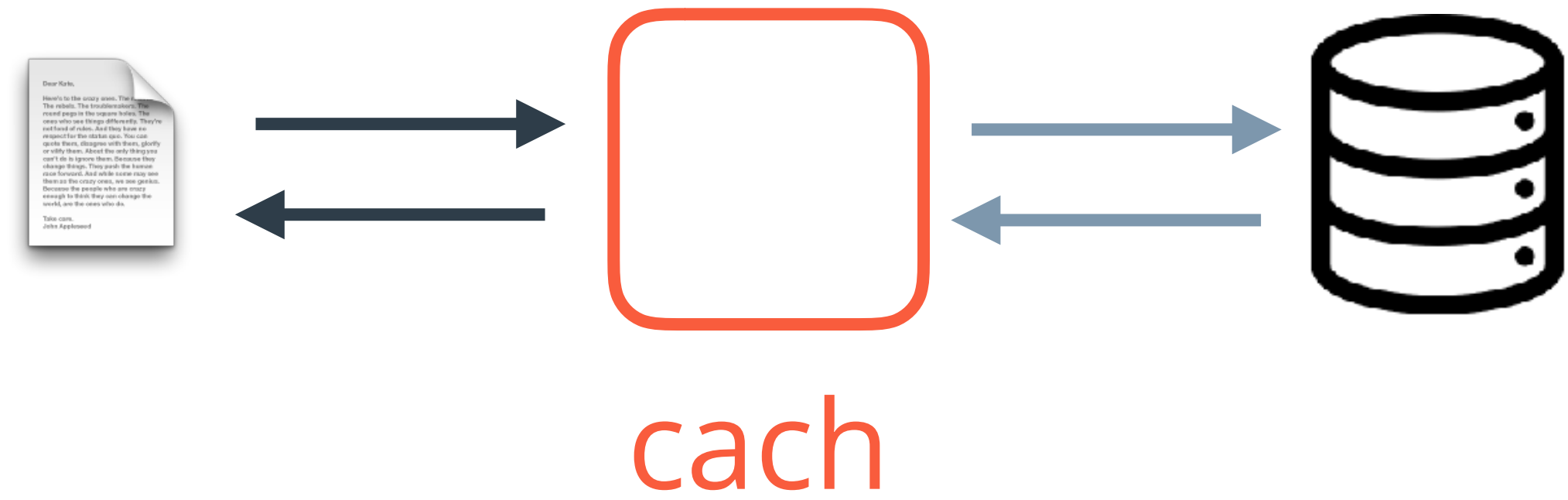


Stemming

~~caching~~

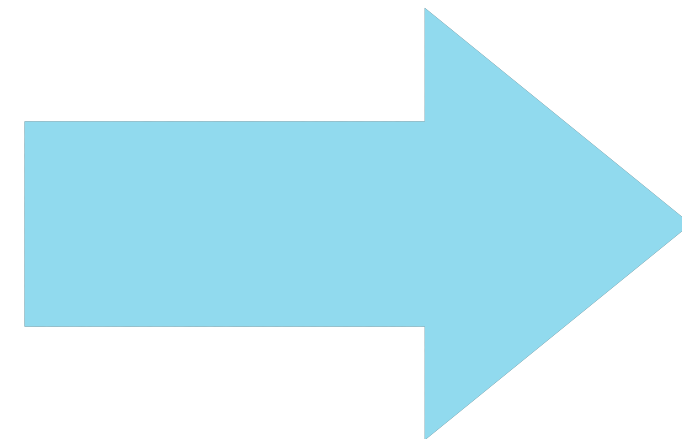
~~cached~~

~~caches~~



Lemmatization

is
was
were

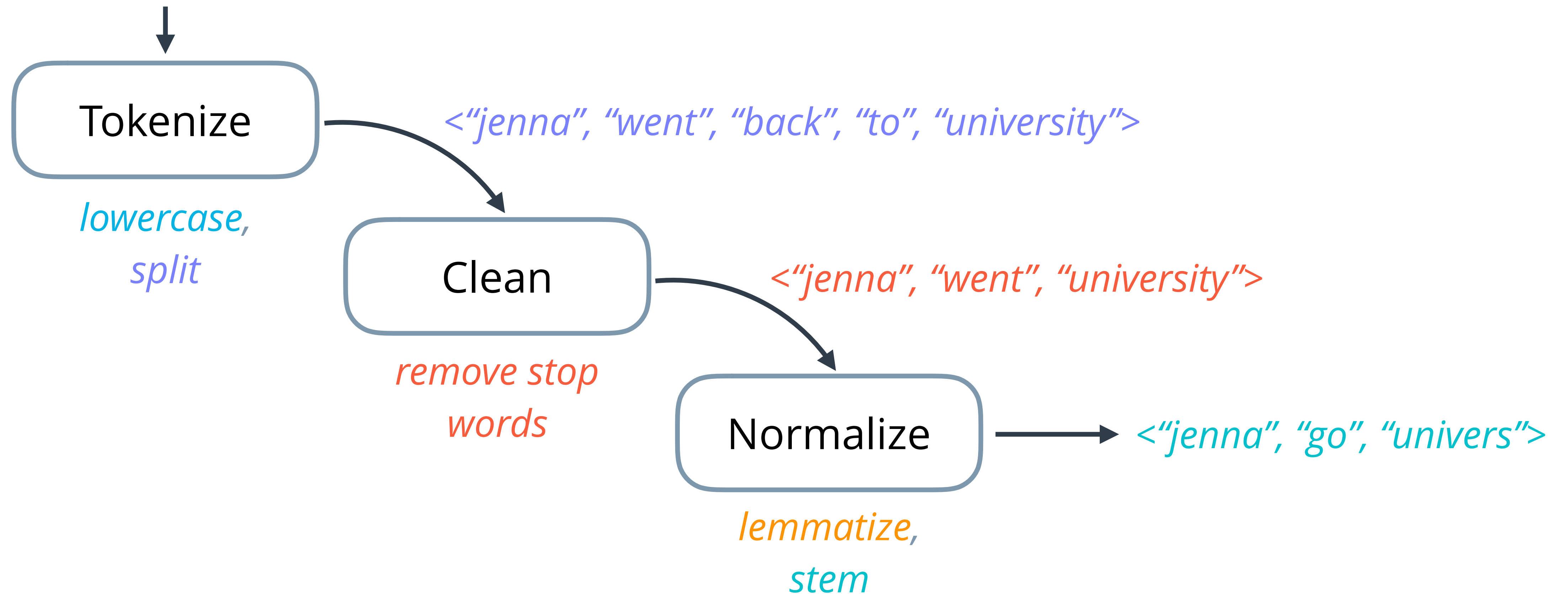


be



Text Processing Summary

"Jenna went back to University."

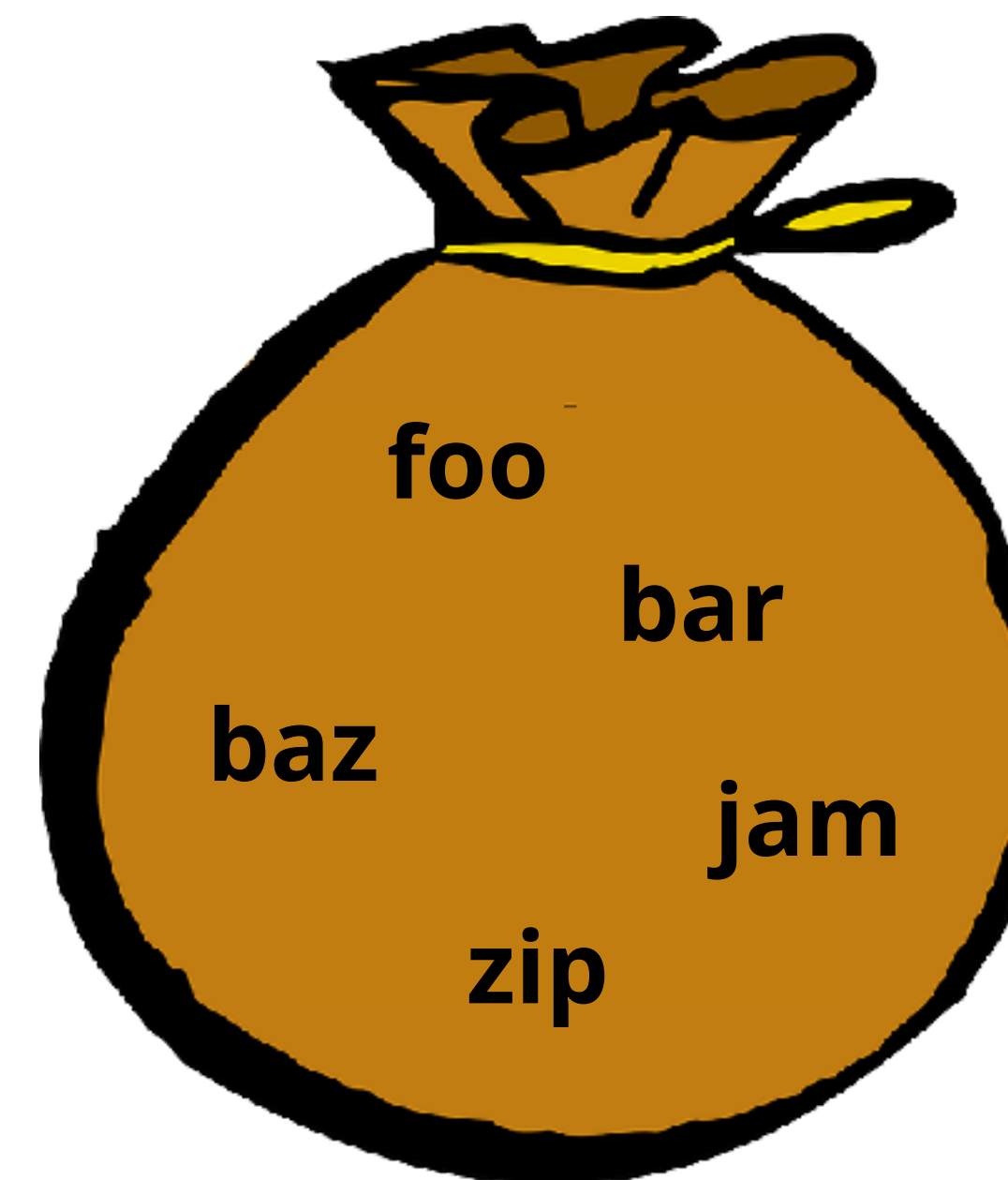
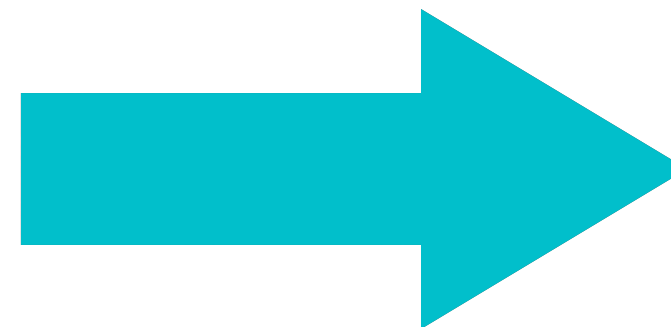


Basic NLP

Basic NLP

- Bag of Words Representation
- Document-Term Matrix
- Task: Document Classification

Bag of Words



Bag of Words

“Little House on the Prairie” → {“littl”, “hous”, “prairi”}

“Mary had a Little Lamb” → {“mari”, “littl”, “lamb”}

“The Silence of the Lambs” → {“silenc”, “lamb”}

“Twinkle Twinkle Little Star” → {“twinkl”, “littl”, “star”} ?

Bag of Words

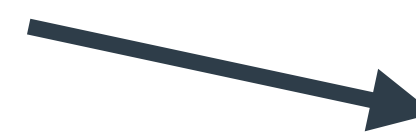
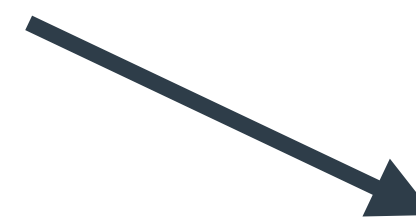
“Little House on the Prairie”

“Mary had a Little Lamb”

“The Silence of the Lambs”

“Twinkle Twinkle Little Star”

corpus (D)



littl hous prairi mari
lamb silenc twinkl star

vocabulary (V)

Bag of Words

“Little House on the Prairie”

“Mary had a Little Lamb”

“The Silence of the Lambs”

“Twinkle Twinkle Little Star”

littl	hous	prairi	mari	lamb	silenc	twinkl	star

Bag of Words - Term Matrix

term frequency

“Little House on the Prairie”

“Mary had a Little Lamb”

“The Silence of the Lambs”

“Twinkle Twinkle Little Star”

	littl	hous	prairi	mary	lamb	silenc	twinkl	star
“Little House on the Prairie”	1	1	1	0	0	0	0	0
“Mary had a Little Lamb”	1	0	0	1	1	0	0	0
“The Silence of the Lambs”	0	0	0	0	1	1	0	0
“Twinkle Twinkle Little Star”	1	0	0	0	0	0	2	1

Document Similarity

a “Little House on the Prairie”

b “Mary had a Little Lamb”

littl	hous	prairi	mari	lamb	silenc	twinkl	star
1	1	1	0	0	0	0	0
1	0	0	1	1	0	0	0

$$\mathbf{a} \cdot \mathbf{b} = \sum a_i b_i = 1 + 0 + 0 + 0 + 0 + 0 + 0 + 0$$

dot product

Document Similarity

a “Little House on the Prairie”

b “Mary had a Little Lamb”

	littl	hous	prairi	mari	lamb	silenc	twinkl	star
a	1	1	1	0	0	0	0	0
b	1	0	0	1	1	0	0	0

$$\mathbf{a} \cdot \mathbf{b} = \sum a_i b_i = 1 \quad \text{dot product}$$

$$\cos(\theta) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \cdot \|\mathbf{b}\|} = \frac{1}{\sqrt{3} \times \sqrt{3}} = \frac{1}{3} \quad \text{cosine similarity}$$

Term Specificity

“Little House on the Prairie”

“Mary had a Little Lamb”

“The Silence of the Lambs”

“Twinkle Twinkle Little Star”

document frequency —

littl	hous	prairi	mari	lamb	silenc	twinkl	star
1/3	1/1	1/1	0/1	0/2	0/1	0/1	0/1
1/3	0/1	0/1	1/1	1/2	0/1	0/1	0/1
0/3	0/1	0/1	0/1	1/2	1/1	0/1	0/1
1/3	0/1	0/1	0/1	0/2	0/1	2/1	1/1
3	1	1	1	2	1	1	1

Term Specificity

“Little House on the Prairie”

“Mary had a Little Lamb”

“The Silence of the Lambs”

“Twinkle Twinkle Little Star”

littl	hous	prairi	mari	lamb	silenc	twinkl	star
1/3	1	1	0	0	0	0	0
1/3	0	0	1	1/2	0	0	0
0	0	0	0	1/2	1	0	0
1/3	0	0	0	0	0	2	1

TF-IDF

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \cdot \text{idf}(t, D)$$

term frequency
 $\text{count}(t, d) / |d|$

inverse document frequency
 $\log(|D| / |\{d \in D : t \in d\}|)$

Task: Document Classification

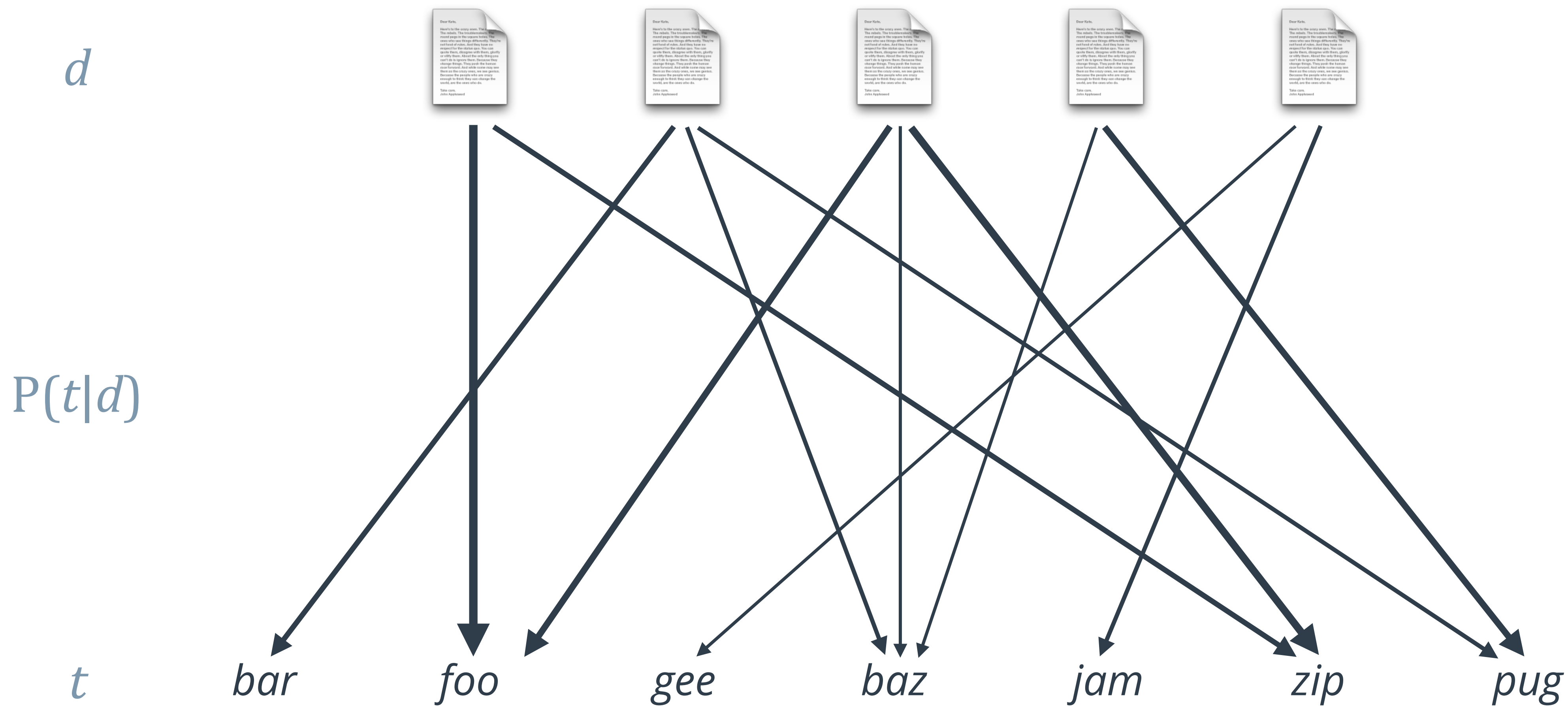
Spam Detection

Advanced NLP

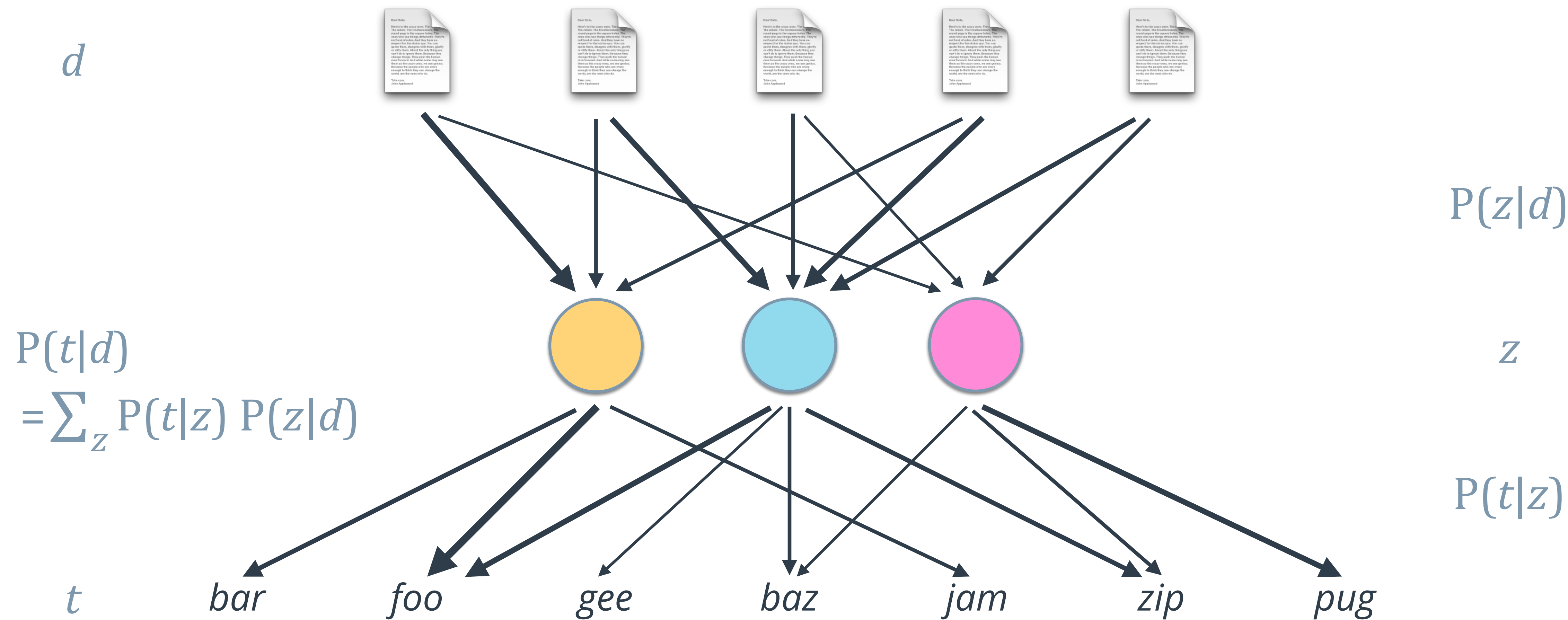
Advanced NLP

- Latent Variables
- Task: Topic Modeling
- Word Embeddings

Bag of Words: Graphical Model



Latent Variables



Missing Priors

$$P(t|d) = \sum_z P(t|z) P(z|d)$$

conditional probabilities

$$P(t, d) = ? \quad P(t, z) = ?$$

joint probabilities

$$P(d)$$

$$P(z)$$

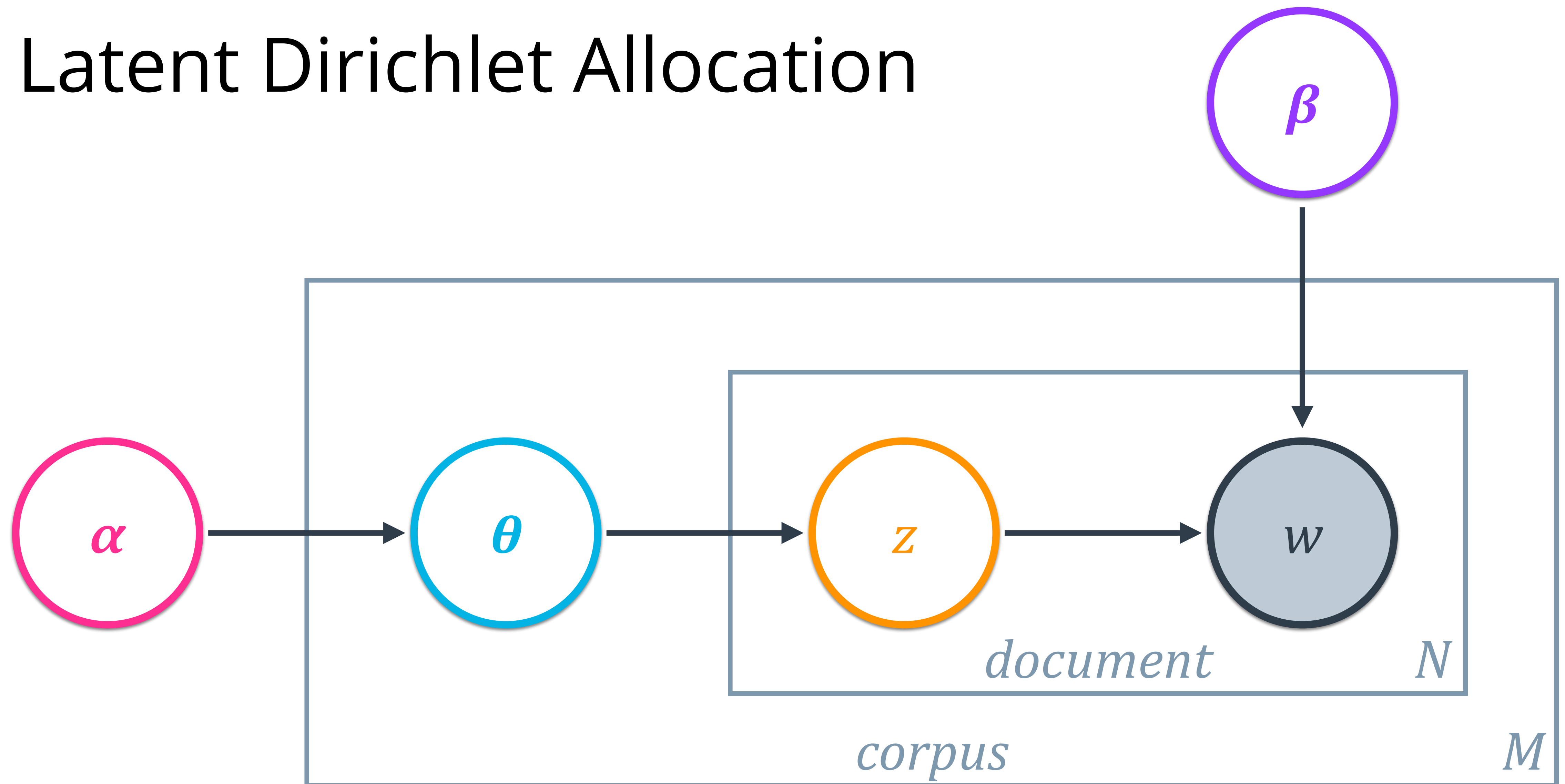
prior probabilities

α

β

Dirichlet distributions

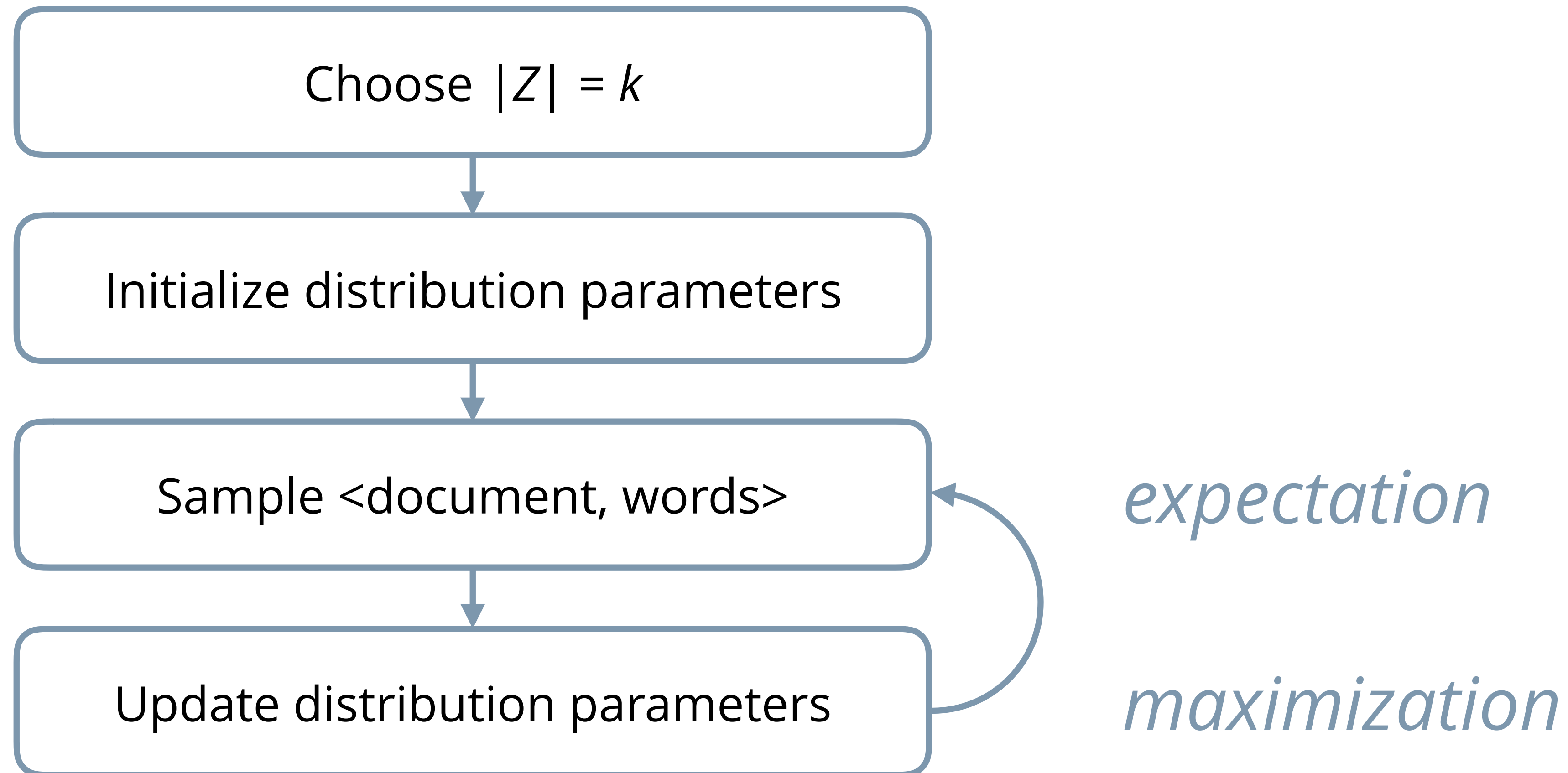
Latent Dirichlet Allocation



LDA: Use Cases

- Topic modeling, document categorization.
- Mixture of topics in a new document: $P(\mathbf{z} \mid \mathbf{w}, \alpha, \beta)$
- Generate collections of words with desired mixture.

LDA: Parameter Estimation



Task: Topic Modeling

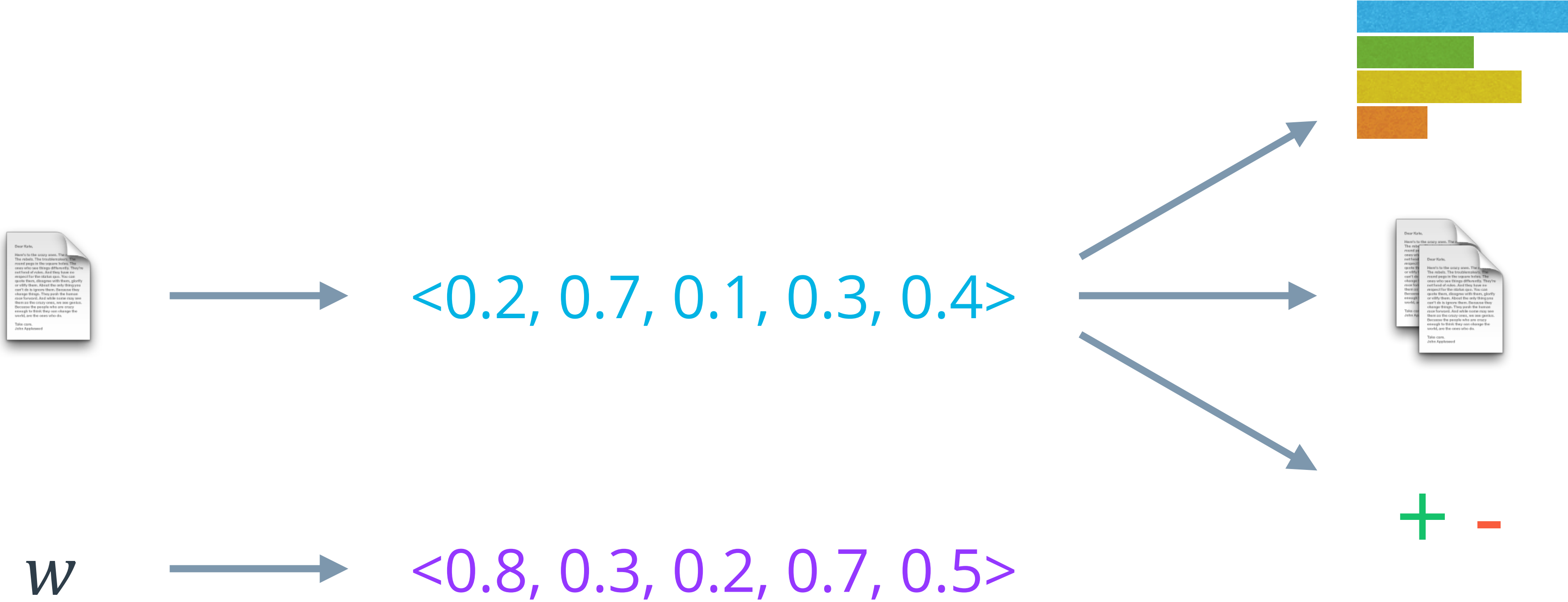
Categorize Newsgroups Data

LDA: Further Reading

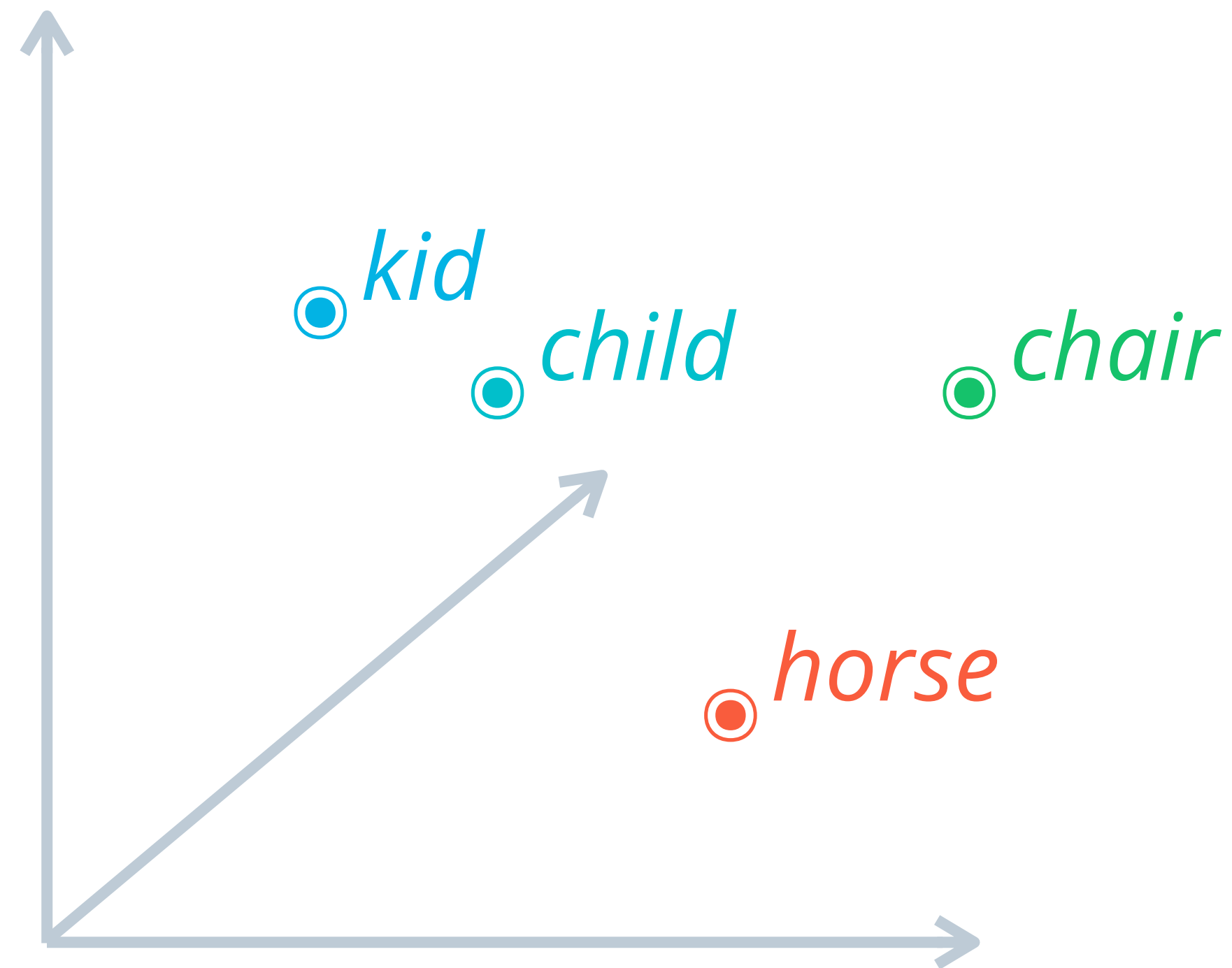
David Blei, Andrew Ng, Michael Jordan, 2003. [Latent Dirichlet Allocation](#),
In *Journal of Machine Learning Research*, vol. 3, pp. 993-102.

Thomas Boggs, 2014. [Visualizing Dirichlet Distributions with matplotlib](#).

Document vs. Word Representations

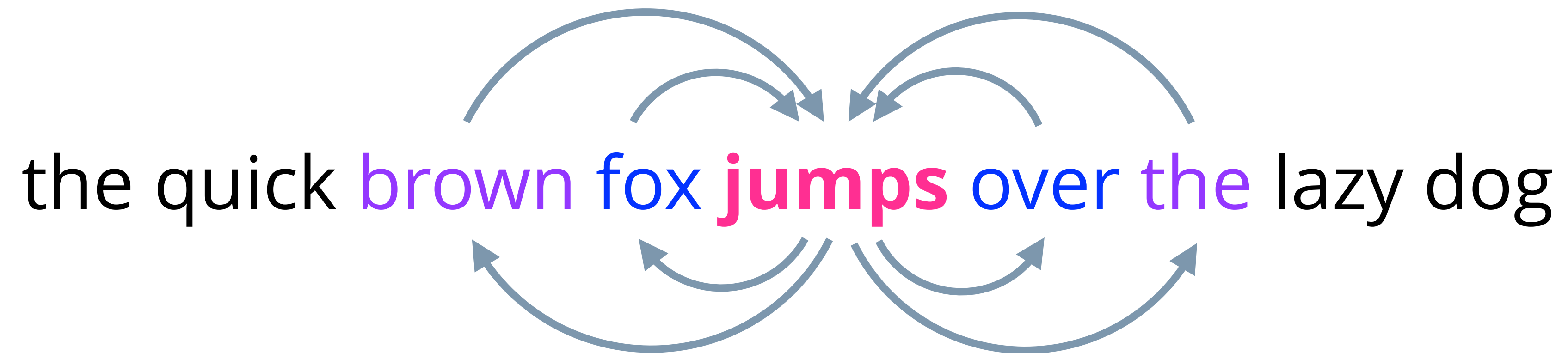


Word Embeddings



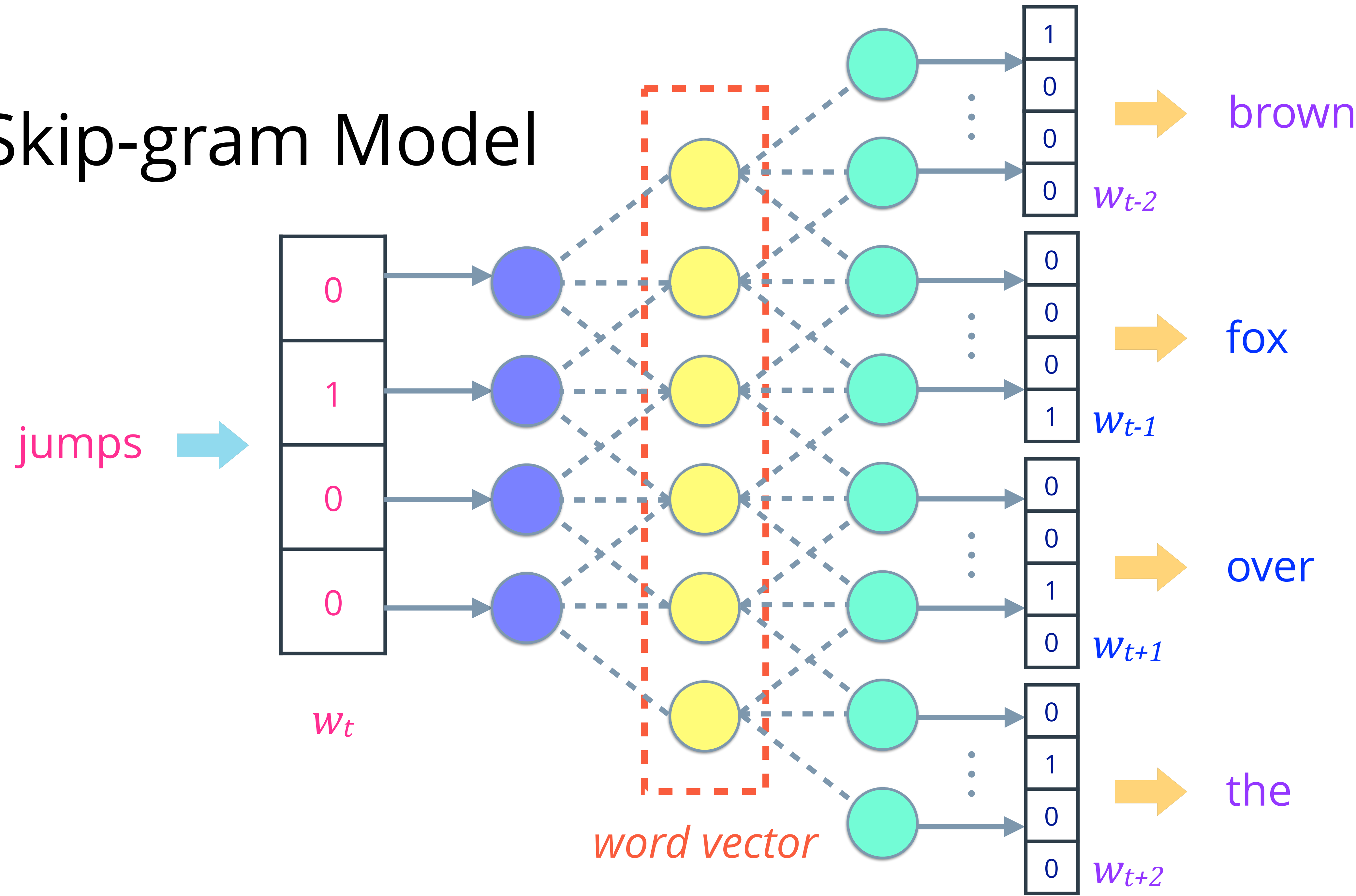
Word2Vec

Continuous Bag of Words (CBow)



Continuous Skip-gram

Skip-gram Model



Word2Vec: Recap

- Robust, distributed representation.
- Vector size independent of vocabulary.
- Deep learning ready!

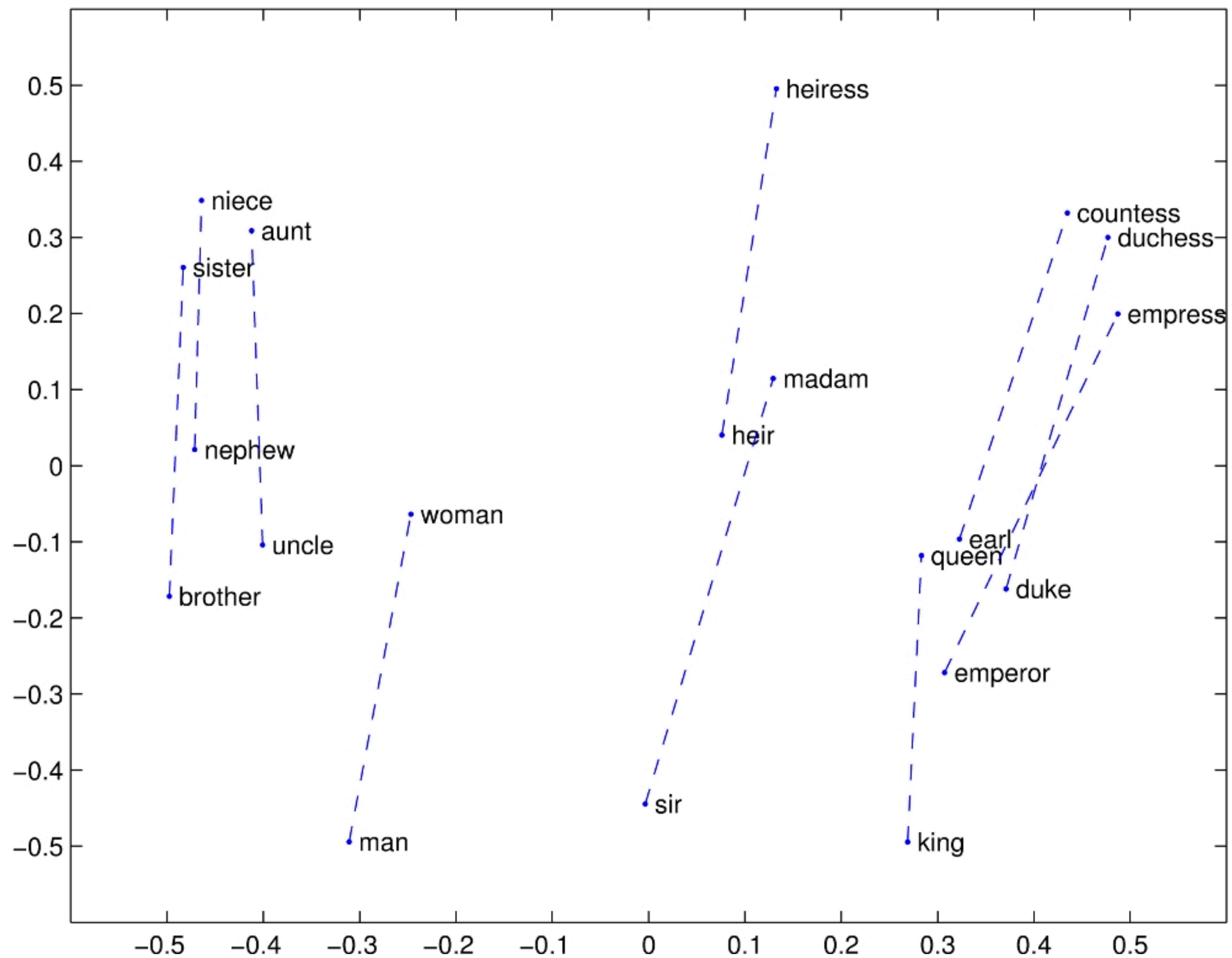
Word2Vec: Further Reading

Tomas Mikolov, et al., 2013. [Distributed Representation of Words and Phrases and their Compositionality](#), In *Advances of Neural Information Processing Systems (NIPS)*, pp. 3111-3119.

Adrian Colyer, 2016. [The amazing power of word vectors](#).

More Word Embeddings

- [GloVe](#): Global Vectors for Word Representation
- [t-SNE](#): t-Distributed Stochastic Neighbor Embedding



Workshop Summary

Text Processing

- Stop word removal, stemming, lemmatization

Basic NLP

- Bag-of-Words, TF-IDF, document classification

Advanced NLP

- LDA, topic modeling, word embeddings

What's Next?

- [Recurrent Neural Networks](#) (RNNs)
- [Long Short-Term Memory Networks](#) (LSTMs)
- [Visual Question-Answering](#)



Adarsh Nair



Luis Serrano



Arpan Chakraborty



udacity.com/ai
udacity.com/ml