

Artificial Intelligence Foundations and Applications

Introduction to Natural Language Processing
13 Nov 2022

Centre of Excellence in Artificial Intelligence
IIT Kharagpur

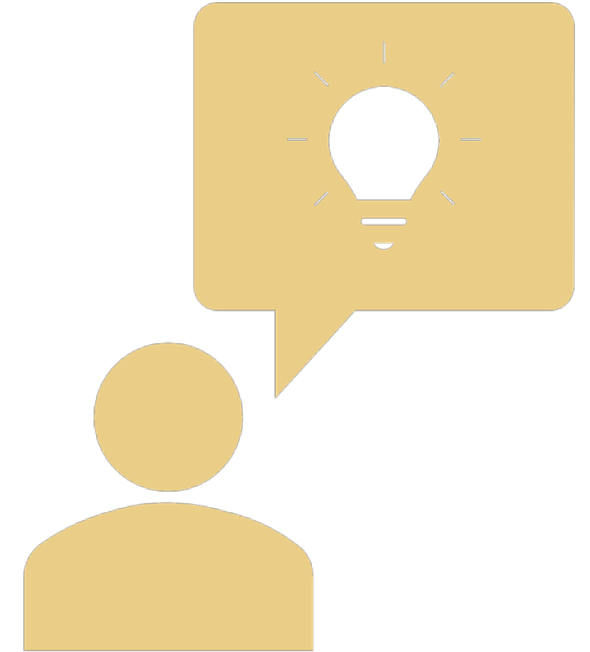
Language is the Tool for Communication

Language is the Vehicle for

- Learning knowledge
- Transmitting information
- Expressing thoughts, perceptions, feelings, information
- Making sense of complex and abstract thought

Communication is two-way

- Convey own ideas
- Receive thought of others

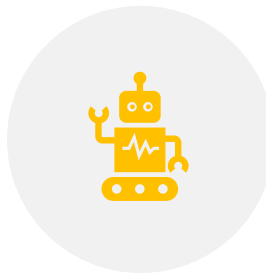


Natural Language Processing

Building computational systems for analyzing and understanding human language input and/or producing natural language output.



Allow computers to communicate with people using natural language.



Computational methods for understanding of human language.

- Automating **Language**

- **Analysis** Language → Representation
- **Generation** Representation → Language
- **Acquisition** Obtaining the representation and necessary algorithms, from knowledge and data



Important Skills

Interact with our world using natural language

- E.g., Conversational agents
- Have computers read all the text out there
 - Retrieve
 - Answer questions
 - Summarize
 - Find new insights, Intelligence

Some Applications



Search



Language Translation



Chatbots



Question Answering



Text Summarization



Sentiment Analysis



Topic Extraction



Named Entity Recognition

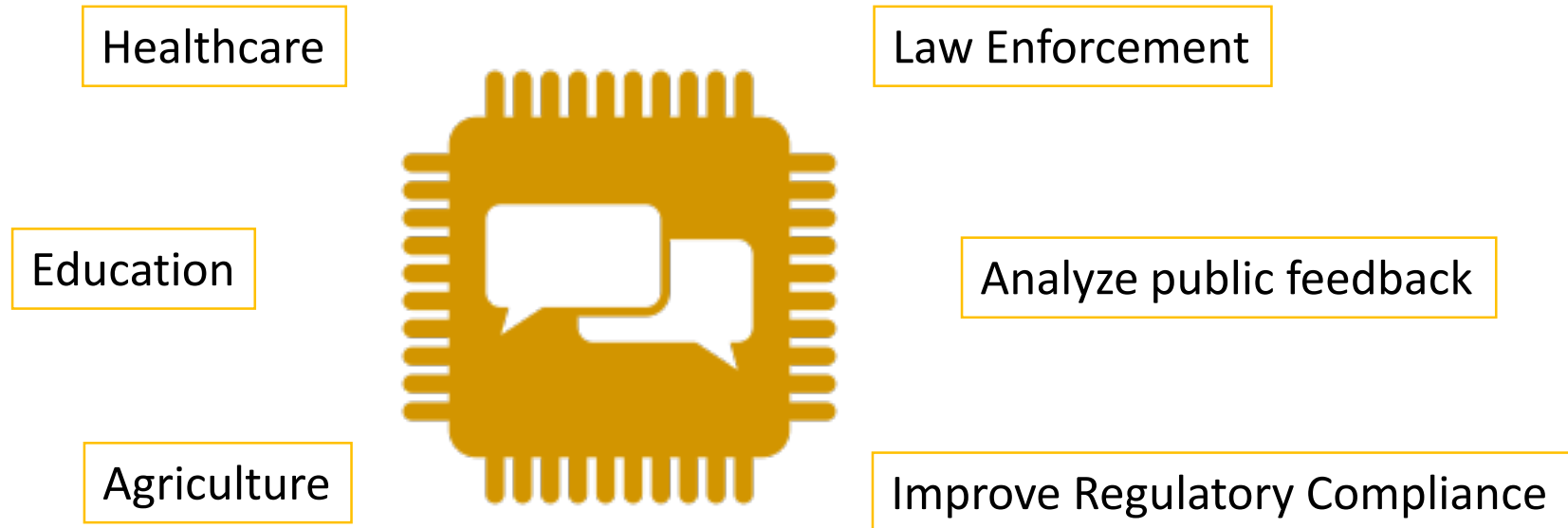


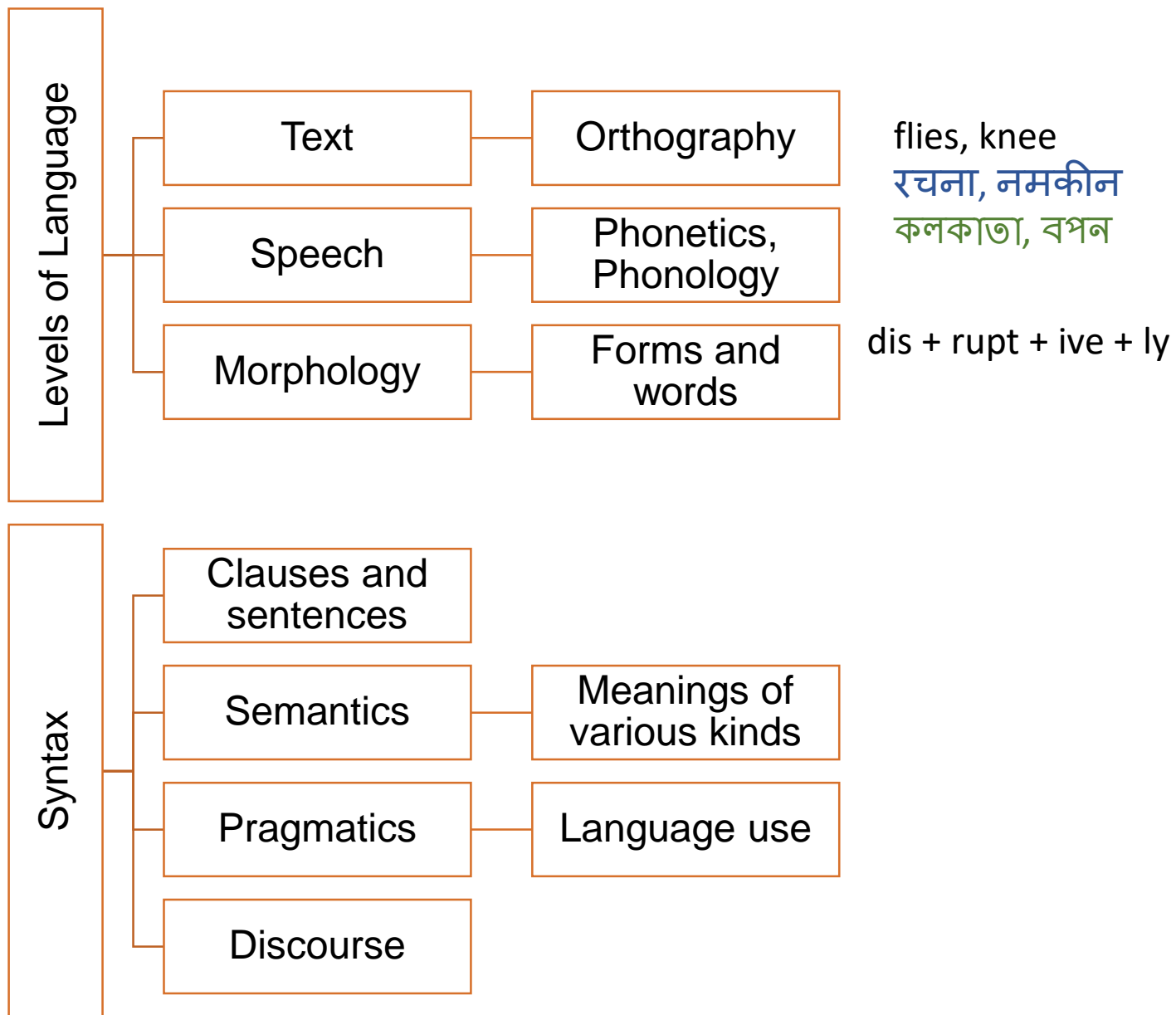
Relation Extraction

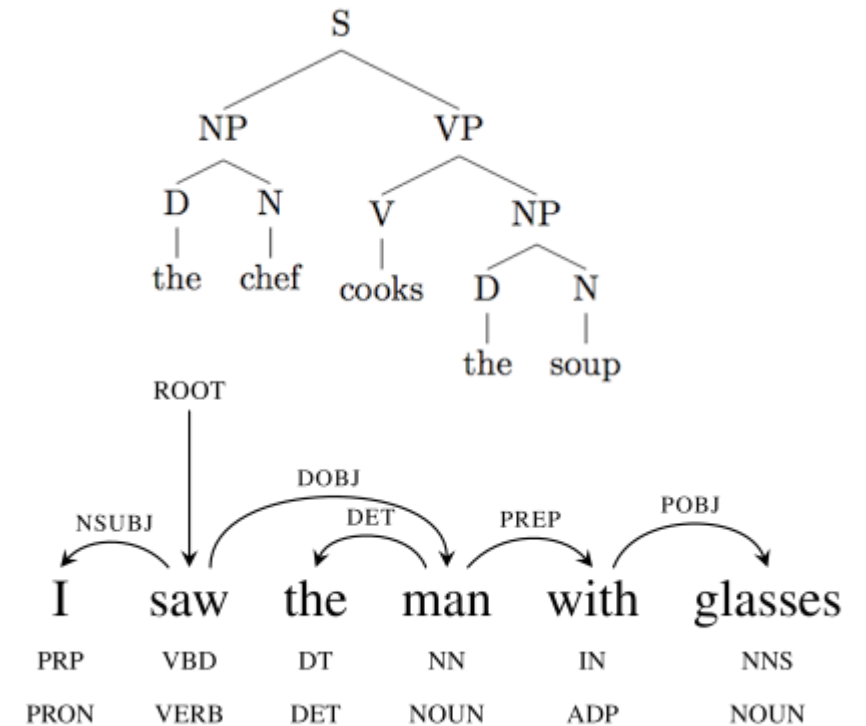
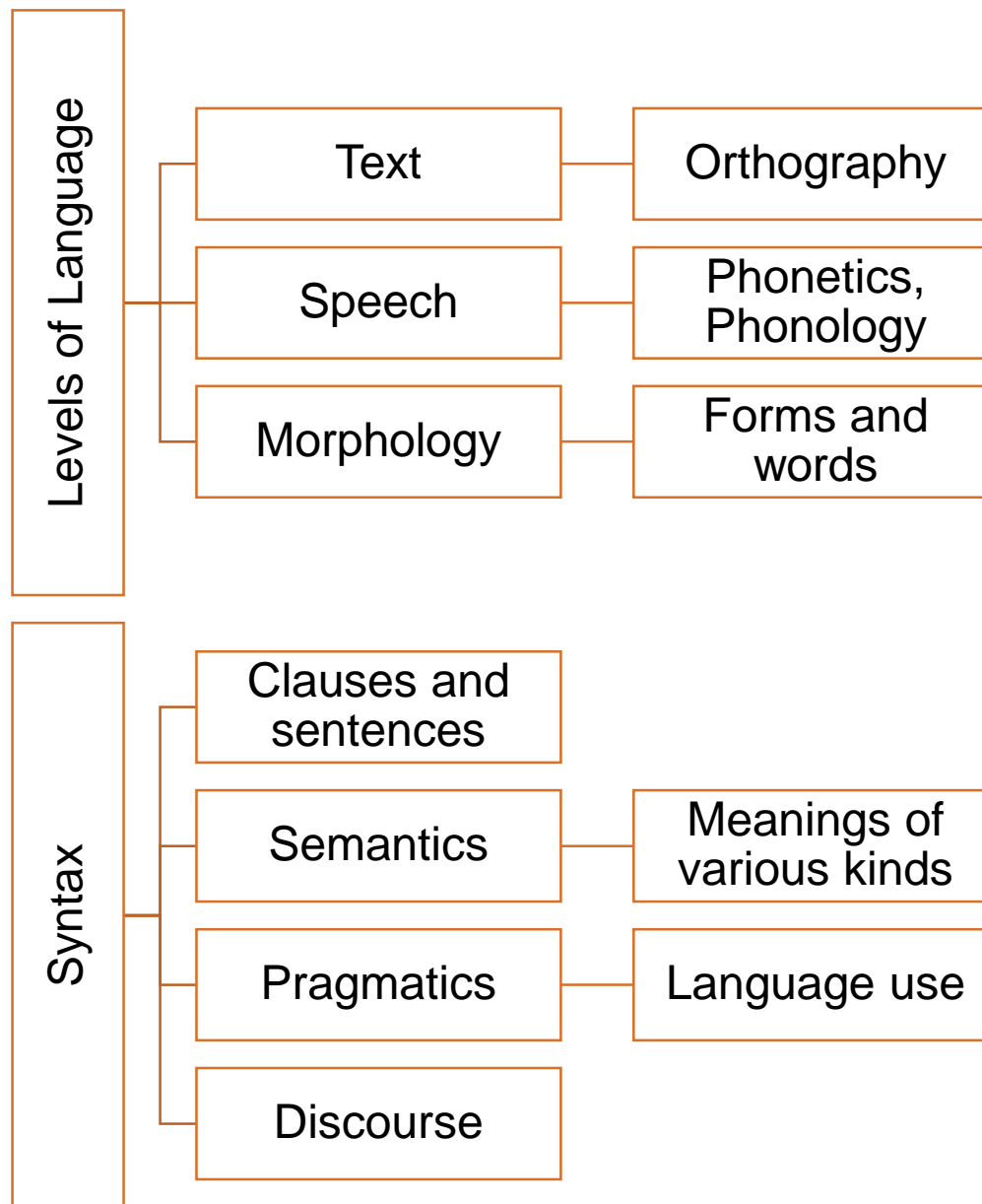


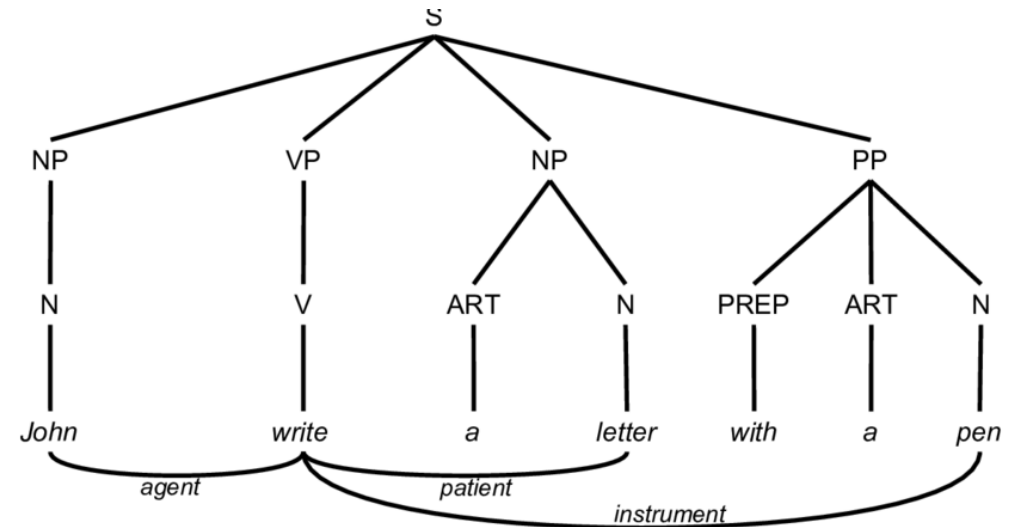
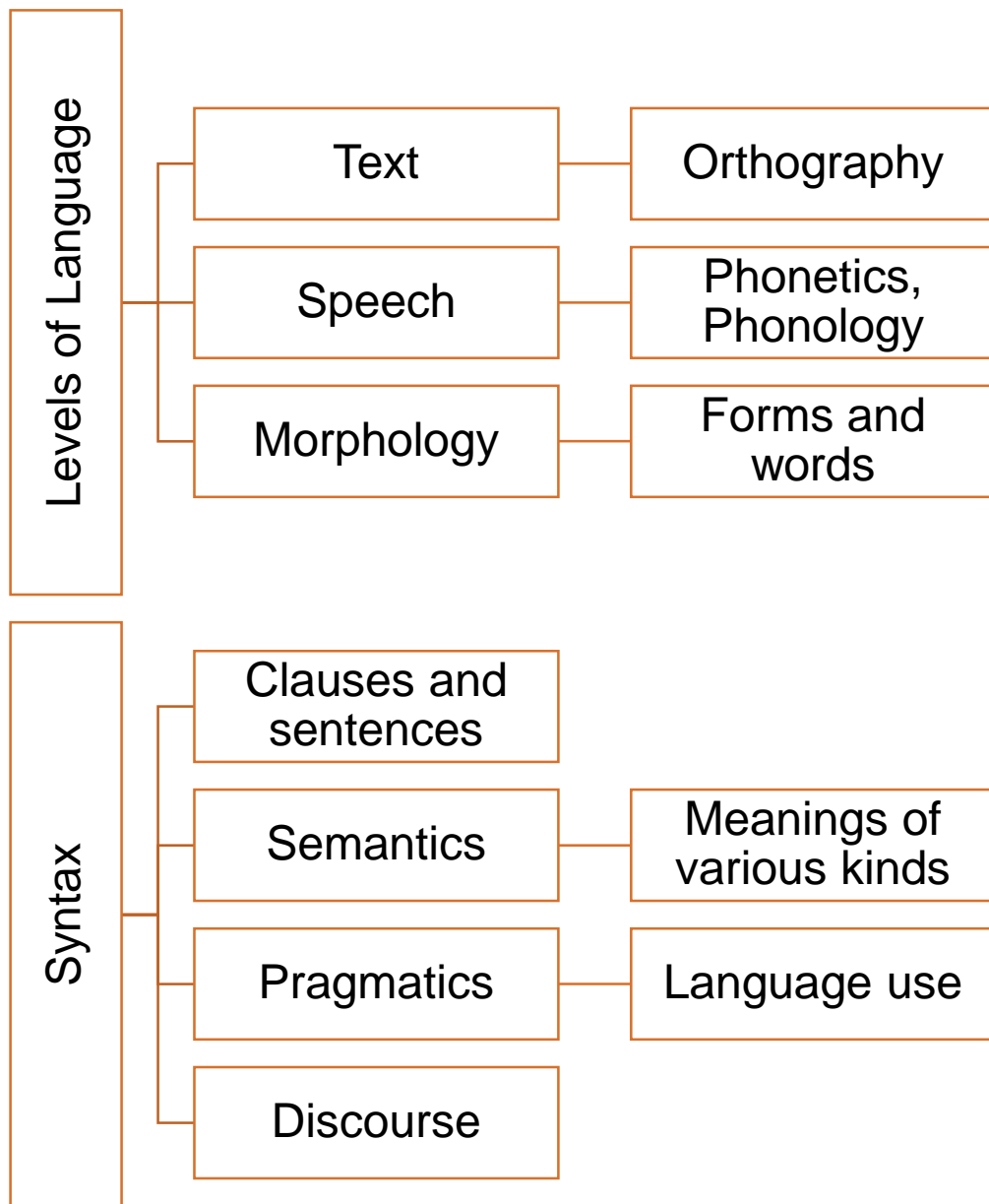
Social Media Monitoring

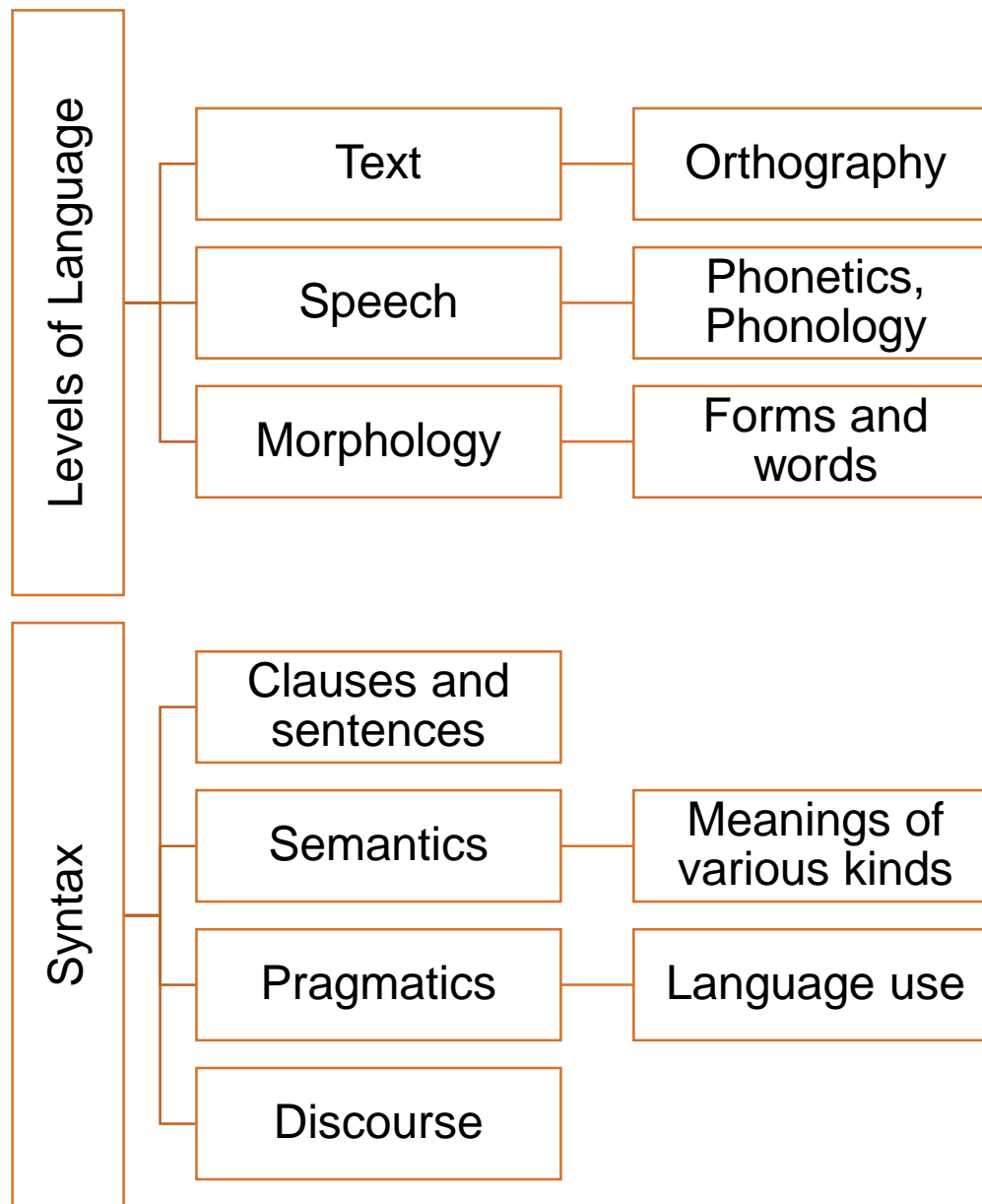
Some application domains

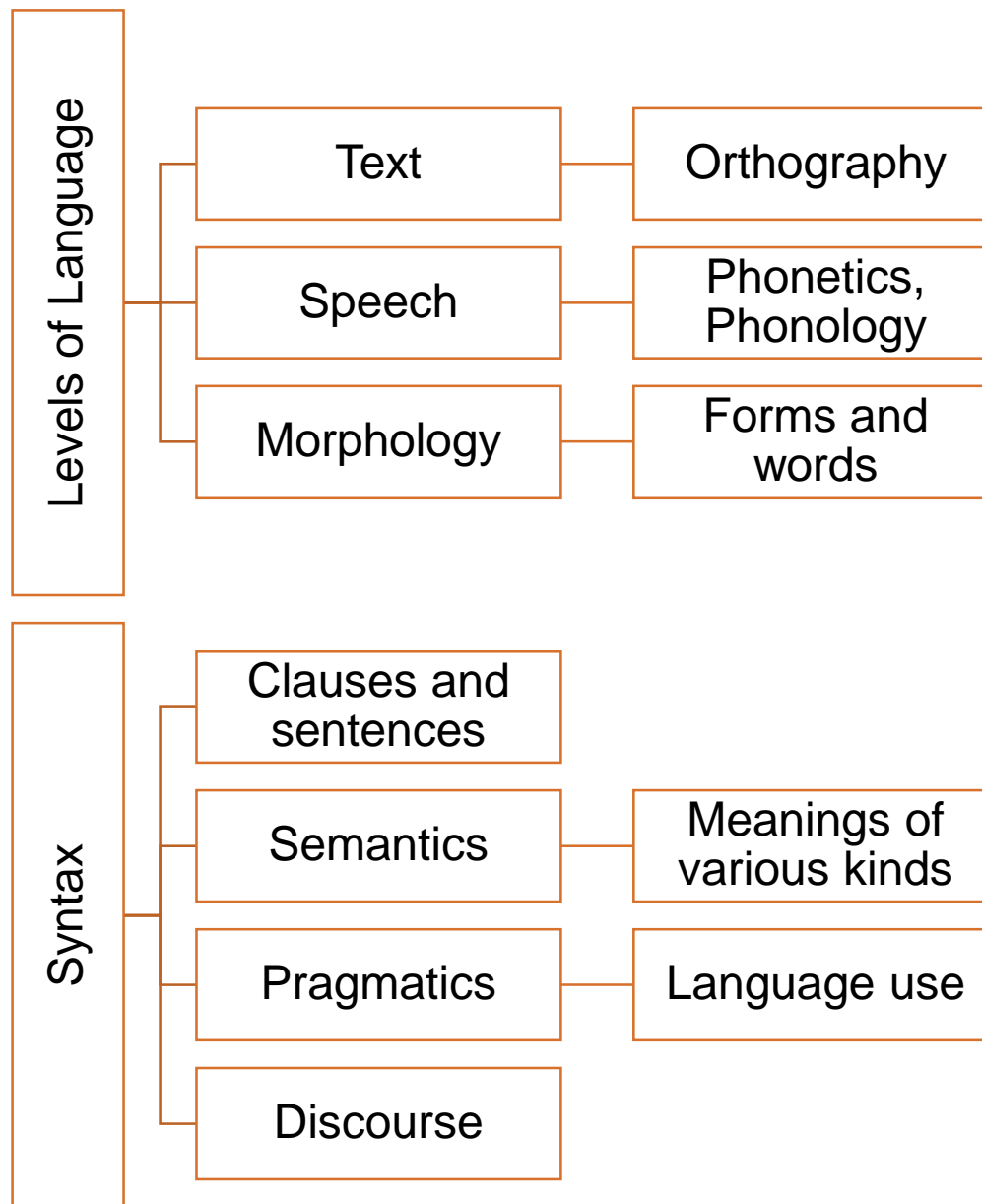






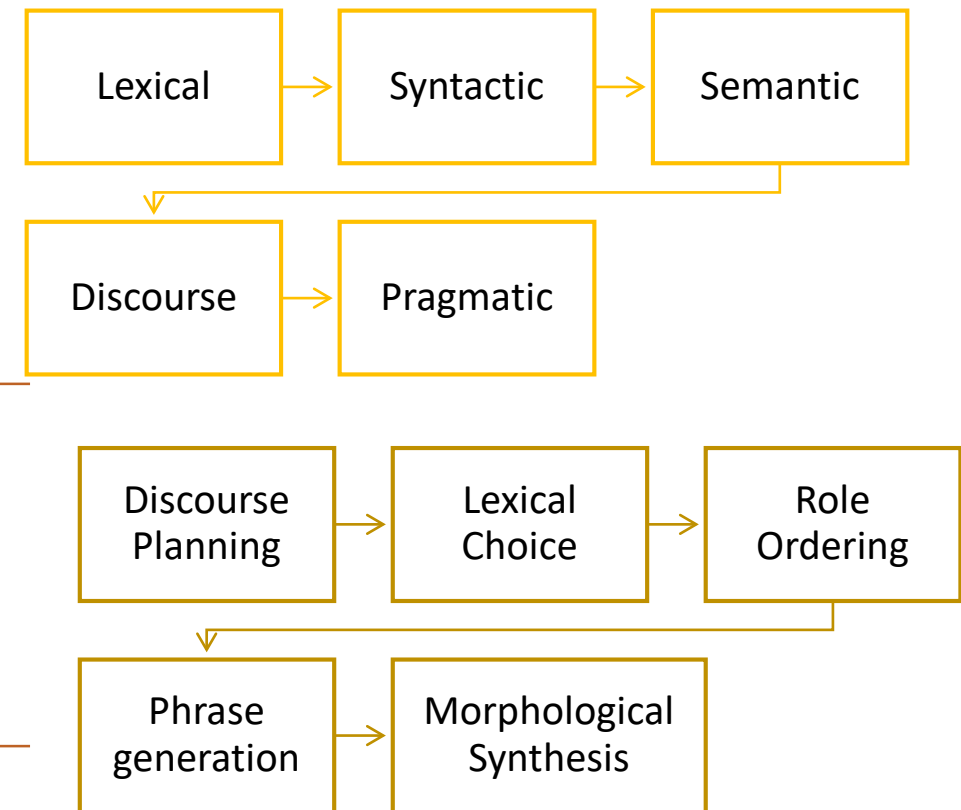
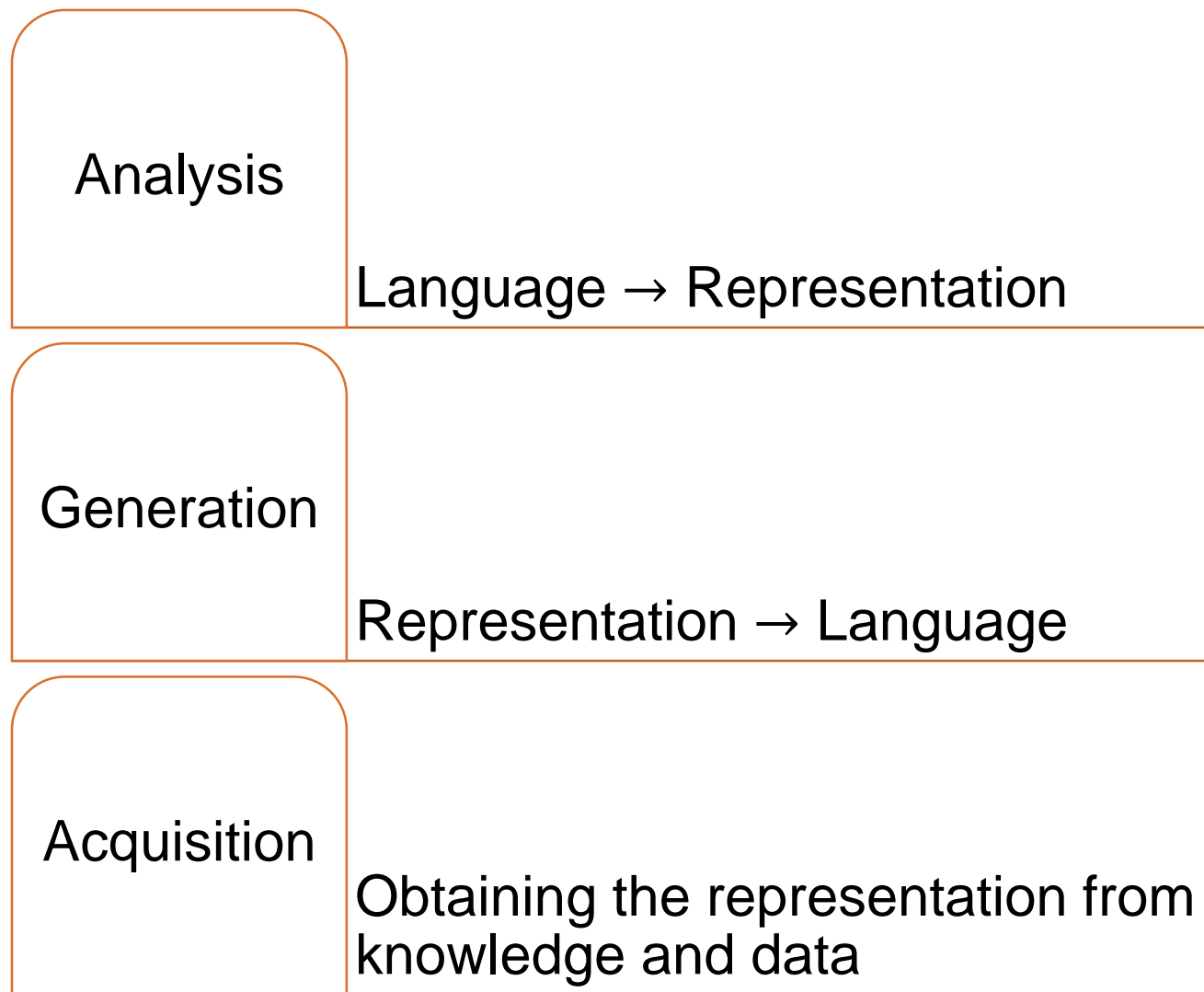






This lesson explains **syntactic analysis**.

It discusses the algorithms for **this task**.





Morphology

- The identification, analysis and description of the structure of words
- Challenges
 - Ambiguity (flies, bears)
 - Segmenting text into words (Thai)
 - Sandhi splitting (Sanskrit)
 - Morphological variations
 - Words with multiple meanings (based on context, domain)
 - Multiword expression

WORDS
MORPHOLOGY

This is a simple sentence

be
3sg
present

Part of Speech

Part of speech tagging

PART OF SPEECH

WORDS

MORPHOLOGY

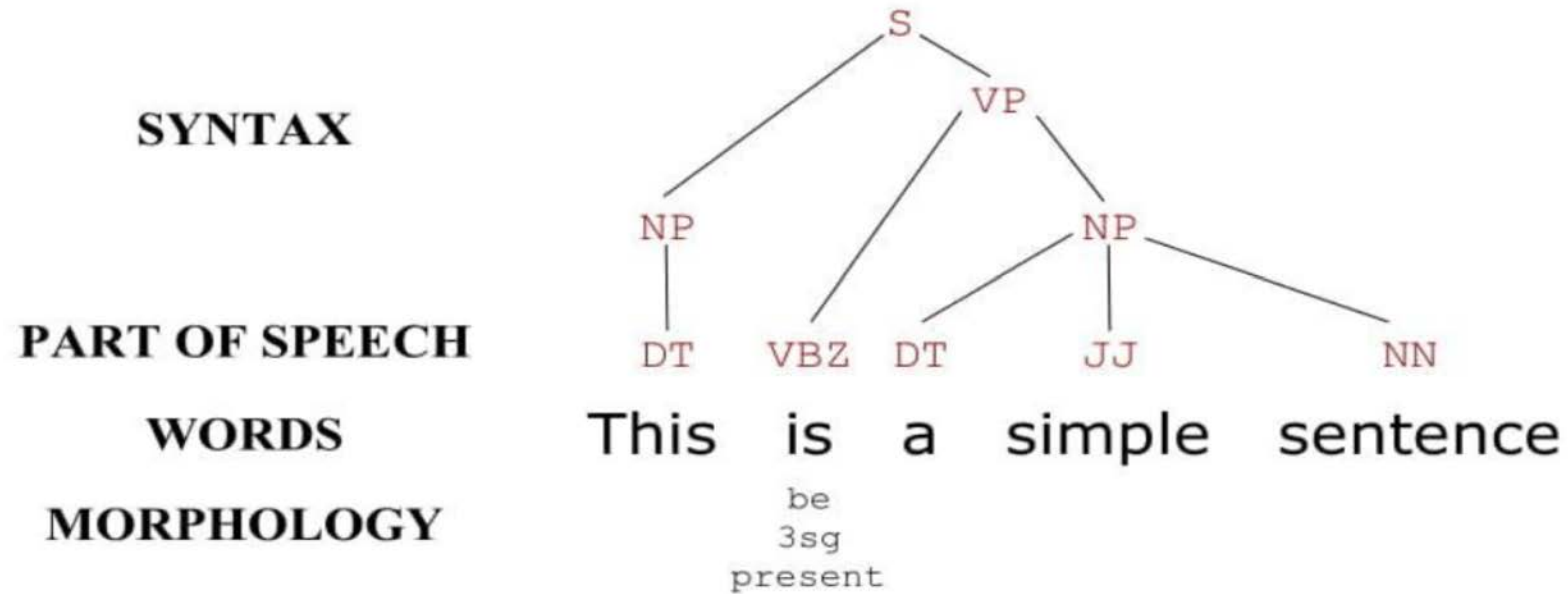
DT	VBZ	DT	JJ	NN
This	is	a	simple	sentence
	be			
	3sg			
	present			

Syntax

- Syntax concerns the way in which words can be combined together to form (grammatical) sentences
 1. revolutionary new ideas appear infrequently
 2. colourless green ideas sleep furiously
 3. *ideas green furiously colourless sleep
- Words combine syntactically in certain orders in a way which mirrors the meaning conveyed
- John gave her dog biscuits
 - (john (gave (her) (dog biscuits)))
 - (john (gave (her dog) (biscuits)))

Syntax

- Syntactic parsing



Semantics

- The manner in which lexical meaning is combined morphologically and syntactically to form the meaning of a sentence
 - Concerns the meaning of words, phrases and sentences
 - The meaning of a sentence is usually a productive combination of the meaning of its words
- Named entity recognition
- Word sense disambiguation
- Semantic role labeling

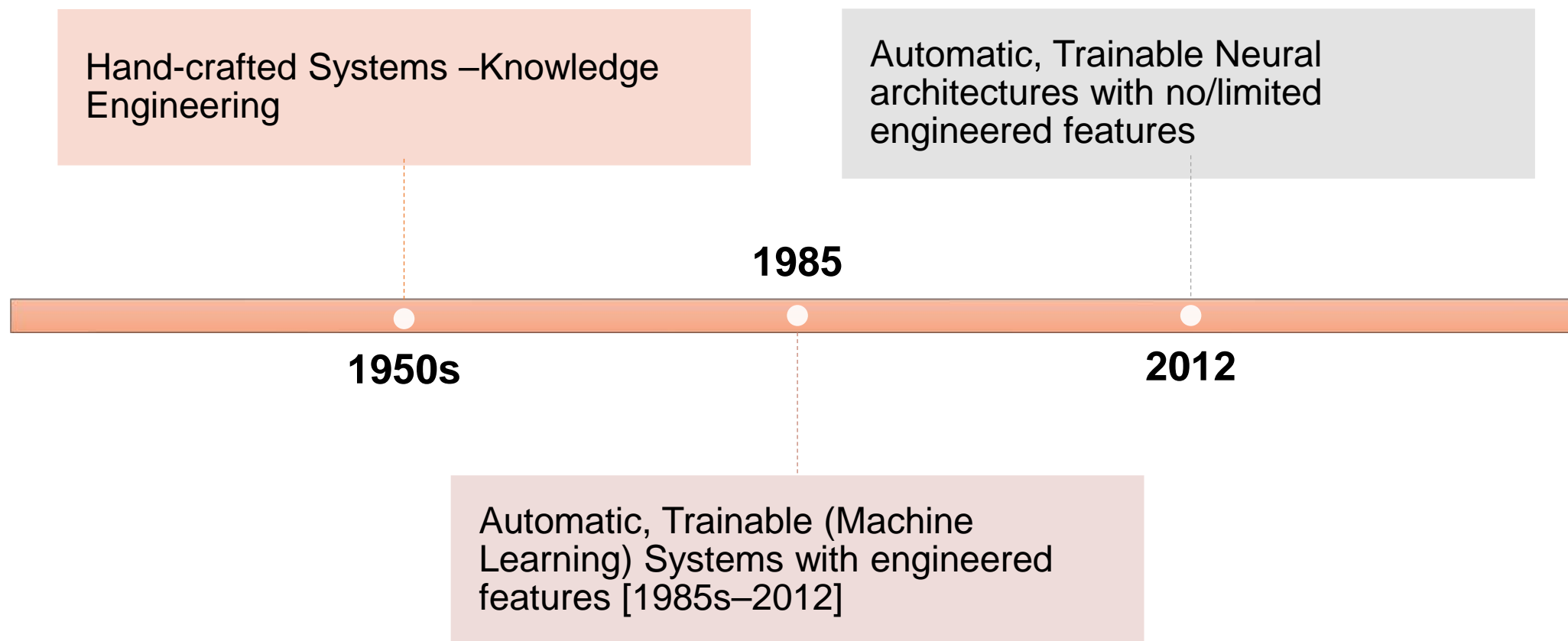
Discourse analysis

- The meaning of a sentence depends upon the sentences that preceded it and also invokes the meaning of the sentences that follow it.
- The discourse structure of connected text, i.e. the nature of the discourse relationships between sentences (e.g. elaboration, explanation, contrast)

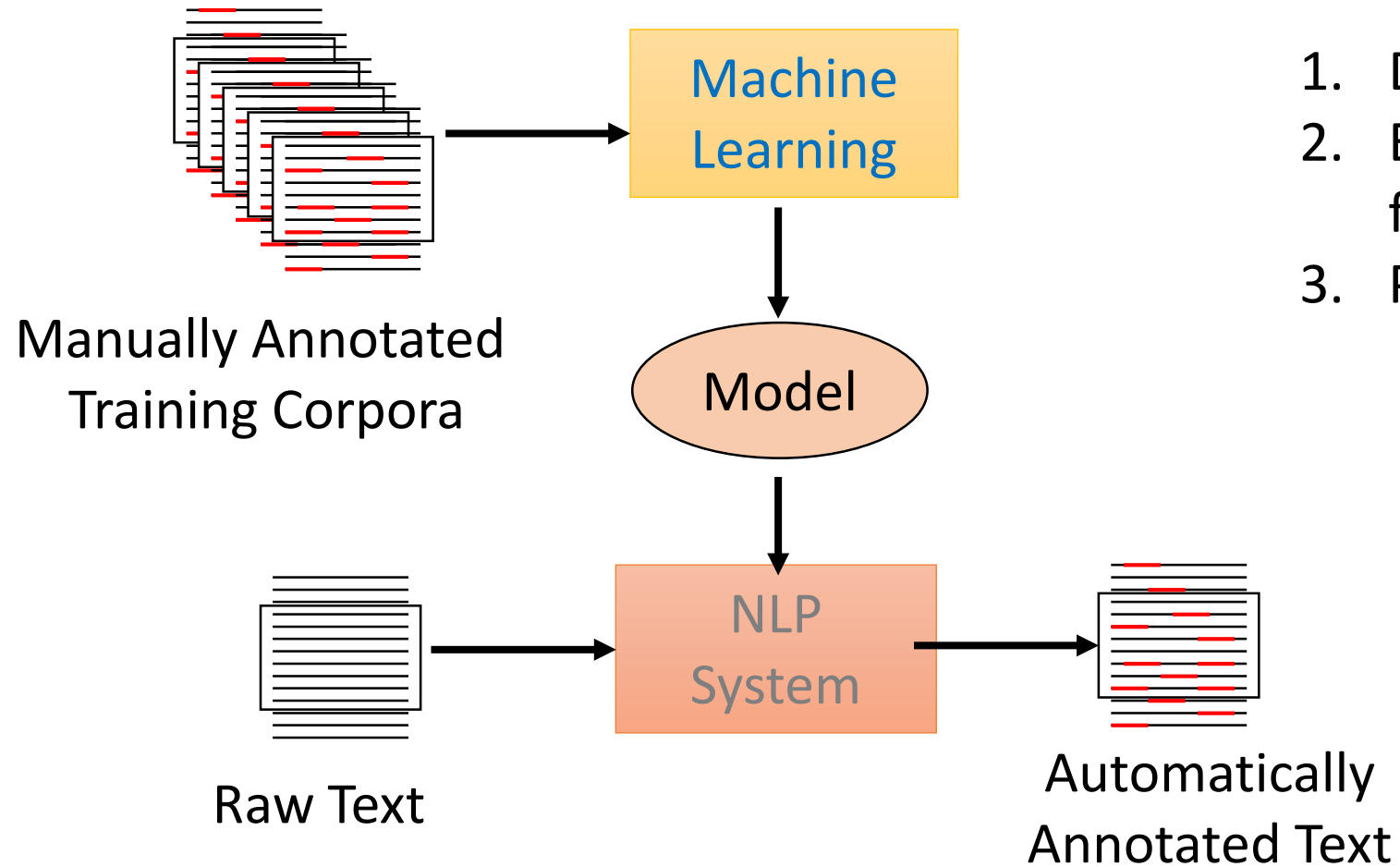
Hardness of NLP

- Ambiguity
- Richness (Variability)
 - Any meaning may be expressed many ways, and there are immeasurably many meanings.
- Linguistic diversity across languages, dialects, genres, styles

Three Generations of NLP



Machine Learning Approach to NLP



1. Data
2. Extraction of “features” from text data
3. Prediction of output

How do we represent words?

$$x_{1,1} = \begin{bmatrix} 0 \\ 0 \\ \cdot \\ 0 \\ 1 \\ 0 \\ \cdot \\ 0 \end{bmatrix}$$

dimensionality = number of possible words

index of this word

This means basically
nothing by itself

Vector Embedding Dimension

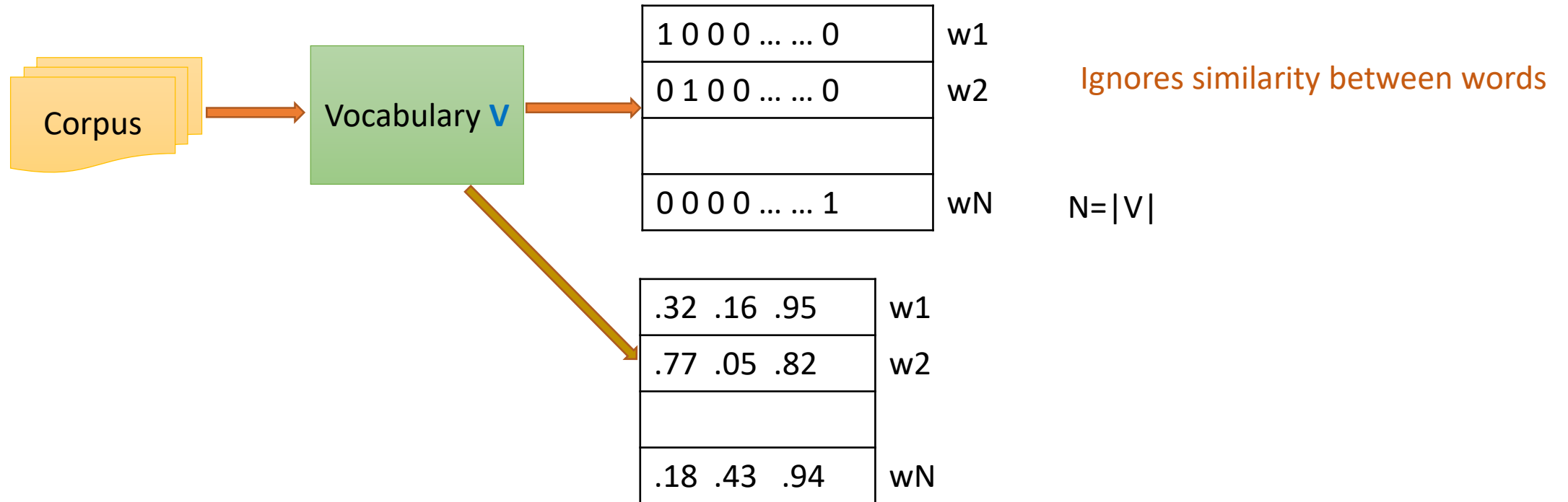
An embedding method should give different embeddings to different words, but a naïve (**one-hot** encoding) is very expensive:

- Dog [0, 0, 0, 0,..., 0, 0, 1, 0, 0, 0, 0,...,0]
 - Man [0, 0, 0, 0,..., 0, 0, 0, 0, 1, 0, 0, 0,...,0]
 - But words are too different (inner products always zero) in this representation.
- } Vector size = vocabulary size

Maybe if we had a more meaningful representation of words, then learning downstream tasks would be much easier!

Meaningful = vectors corresponding to similar words should be close together

Word Representations



Word2vec: Represent each word with a low-dimensional dense vector

Model more generalizable

Word Representation

- Continuous Representation: based on context
- Distributional hypothesis
You can get a lot of value by representing a word by means of its neighbors

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

One of the most successful ideas of modern NLP

government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge

↖ These words will represent *banking* ↗

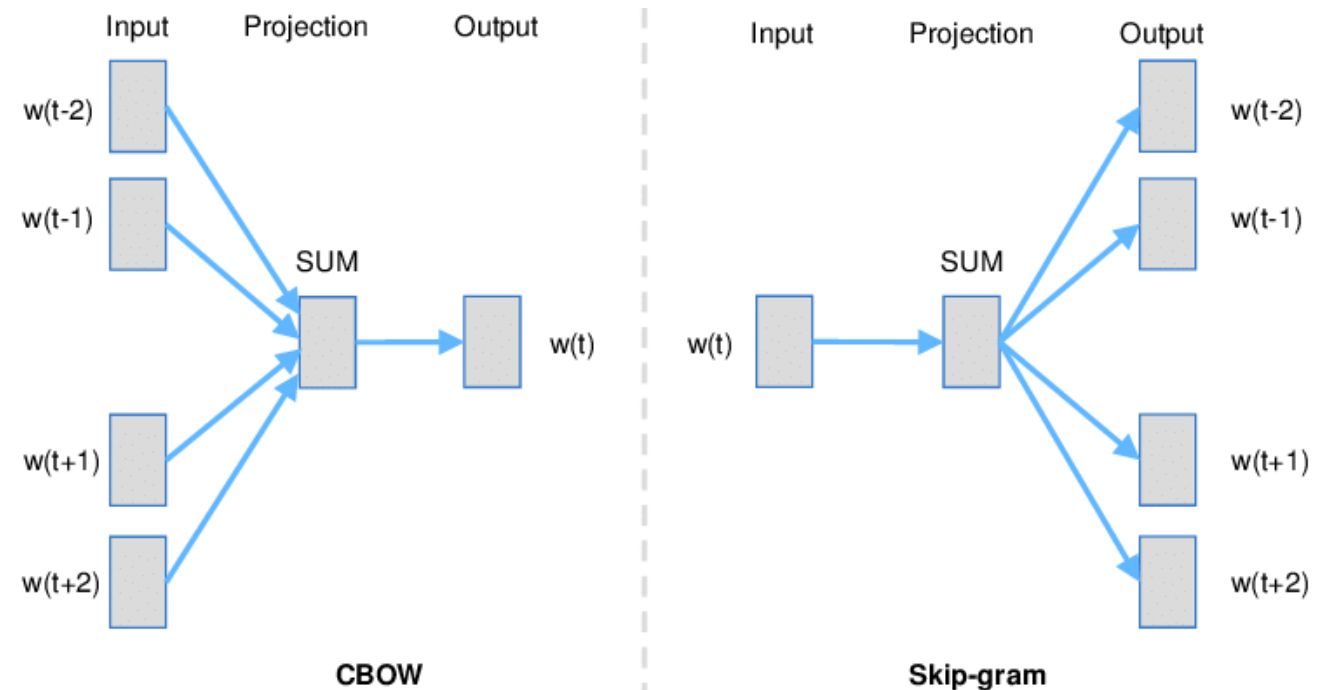
Word2vec Representations

“You shall know a word by the company it keeps”

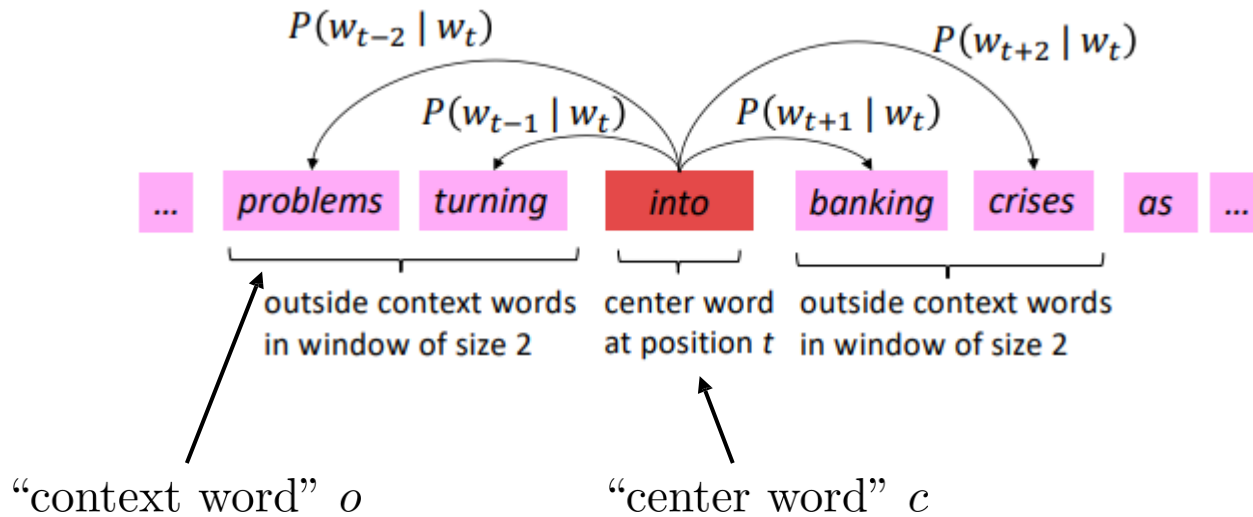
Key idea: Predict surrounding words of every word

Assign each word a vector such that similar words have similar vectors

1. CBOW: $P(\text{Word} | \text{Context})$
2. Skipgram: $P(\text{Context} | \text{Word})$



Can we **predict** the neighbors of a word from its **embedding**?



(learned) vector representation of o (learned) vector representation of c

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

all possible words (vocabulary)

looks a bit like a **logistic regression** model

how to train? $\arg \max_{u_1, \dots, u_n, v_1, \dots, v_n} \sum_{c, o} \log p(o|c)$

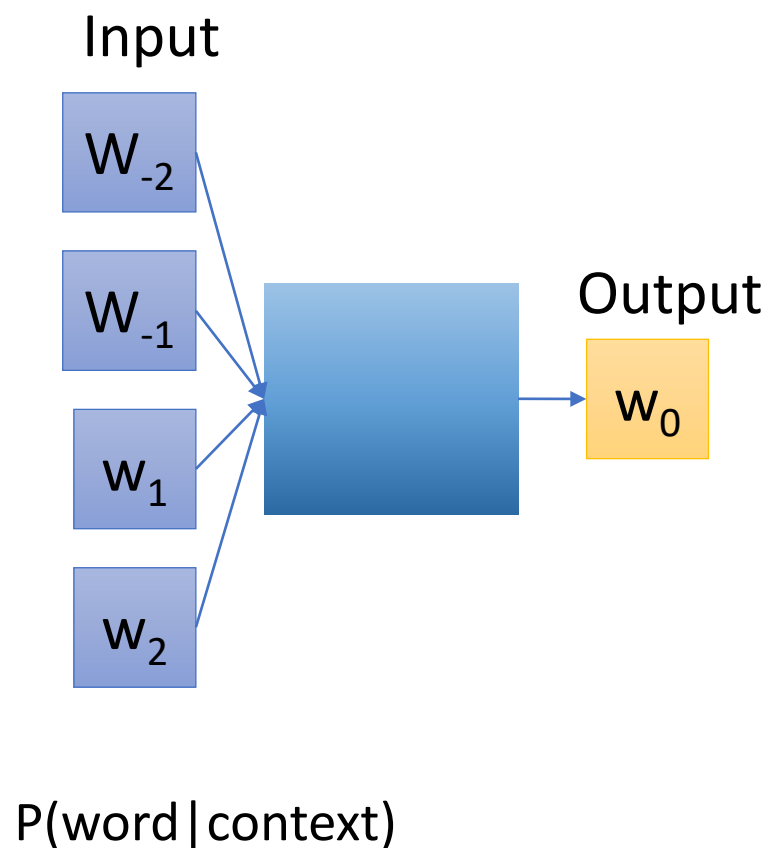
u and v vectors for all possible words

all possible c and o combinations
e.g., for each word c , pick all words o that are within 5 step

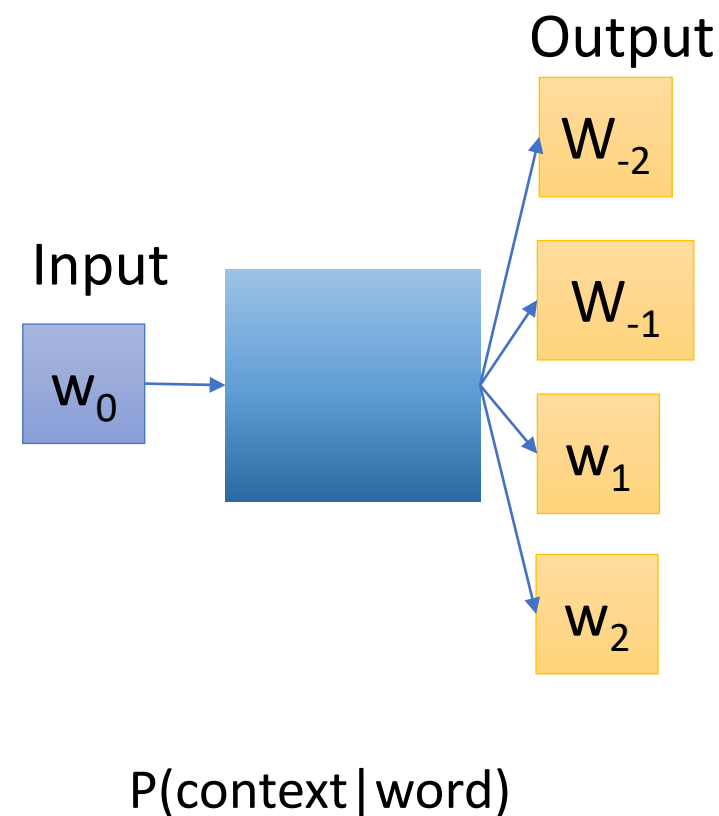


Word2Vec: main context representation models

Continuous Bag of Words (CBOW)

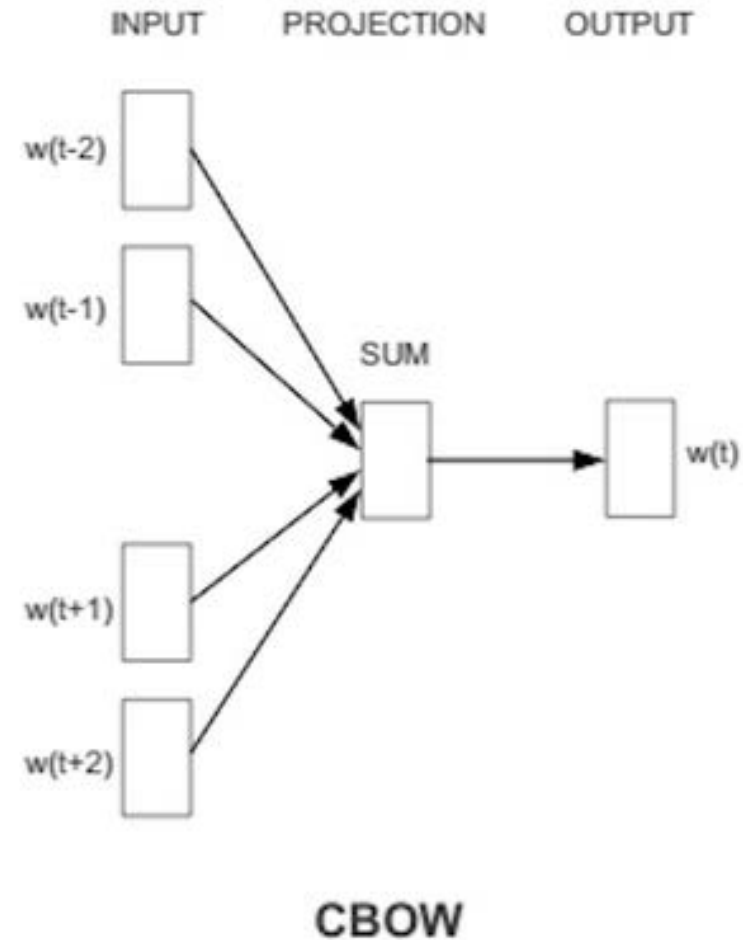


Skip-gram



CBOW

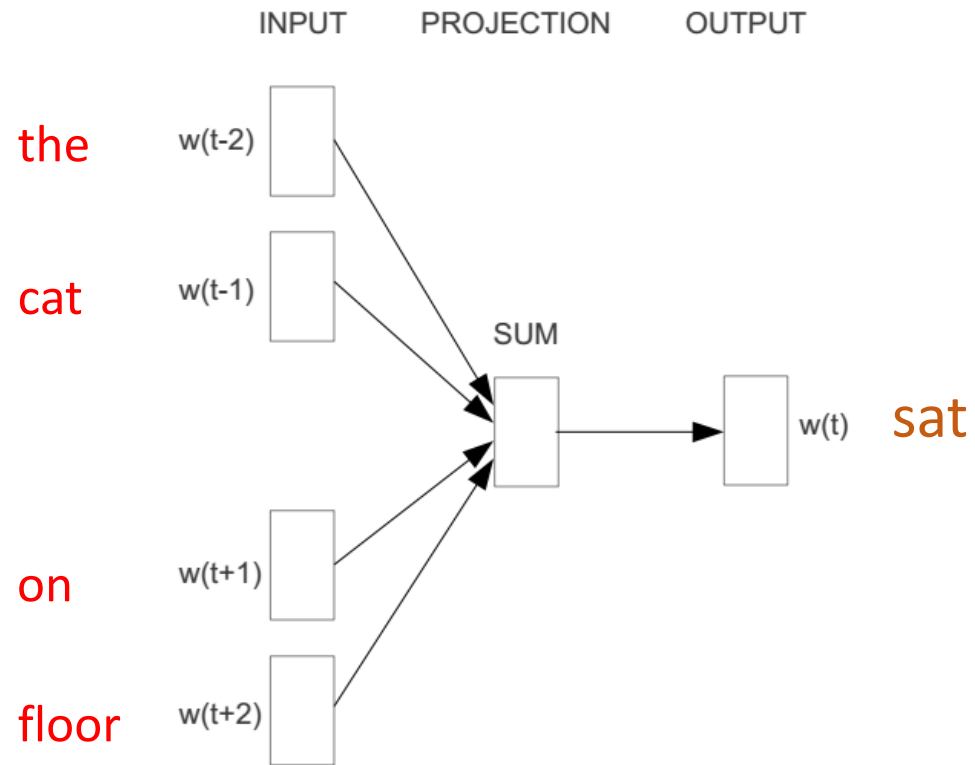
- Bag of words
- Gets rid of word order.
- Takes vector embeddings of n words before target and n words after and adds them (as vectors).

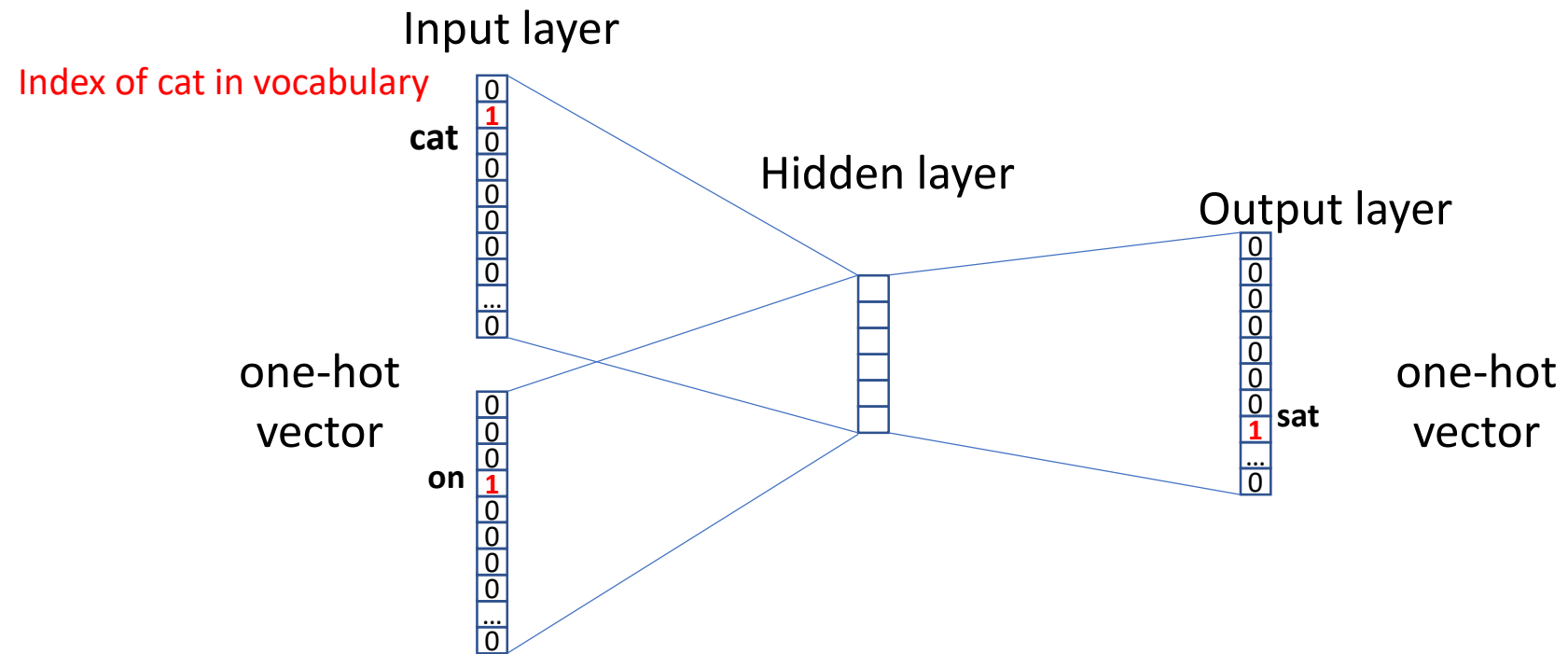


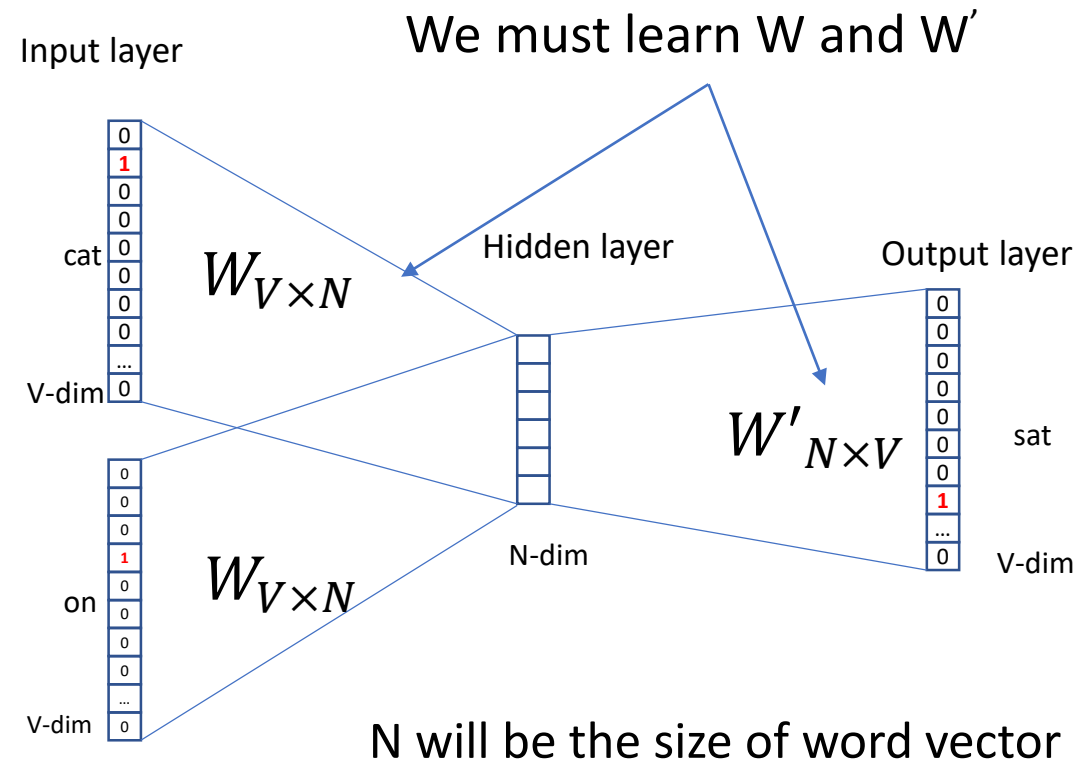
Word2vec – Continuous Bag of Word

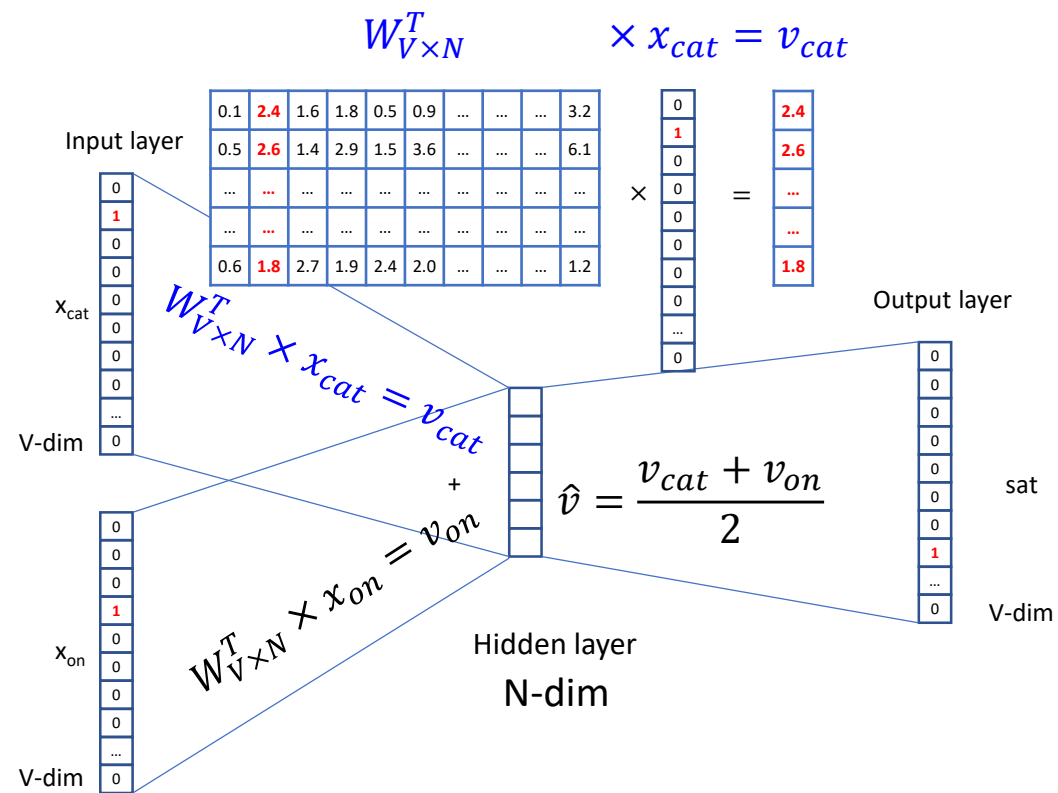
E.g. “The cat sat on floor”

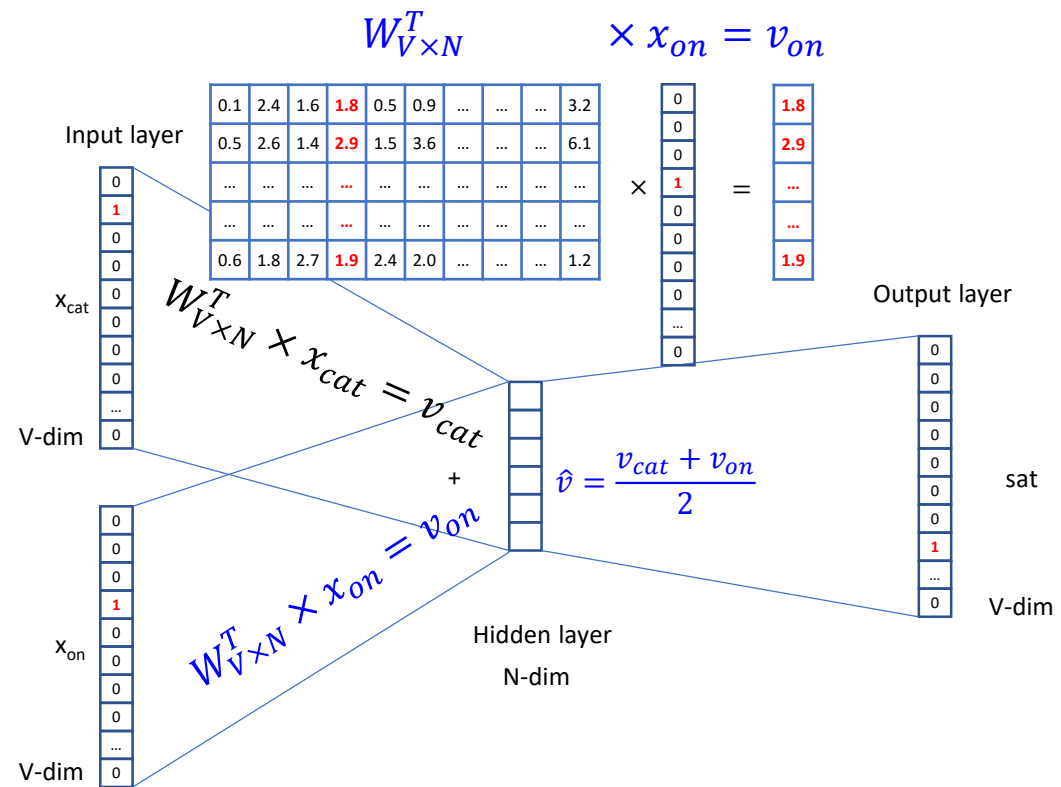
- Window size = 2

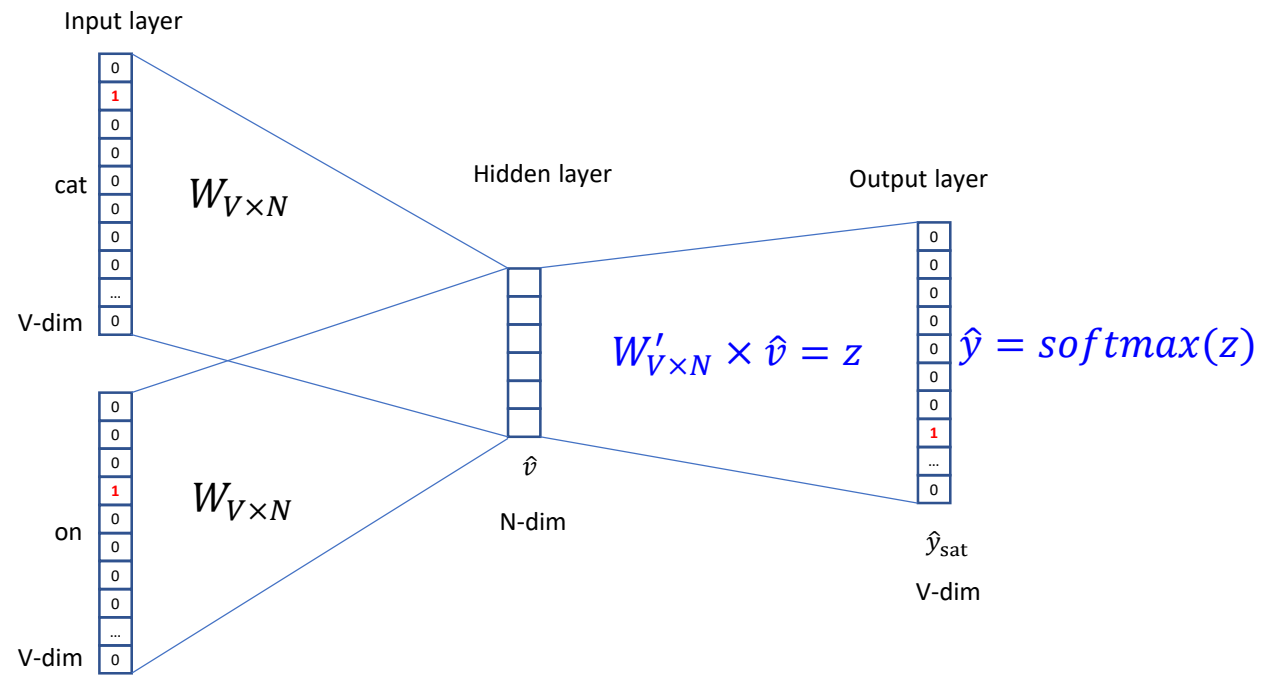


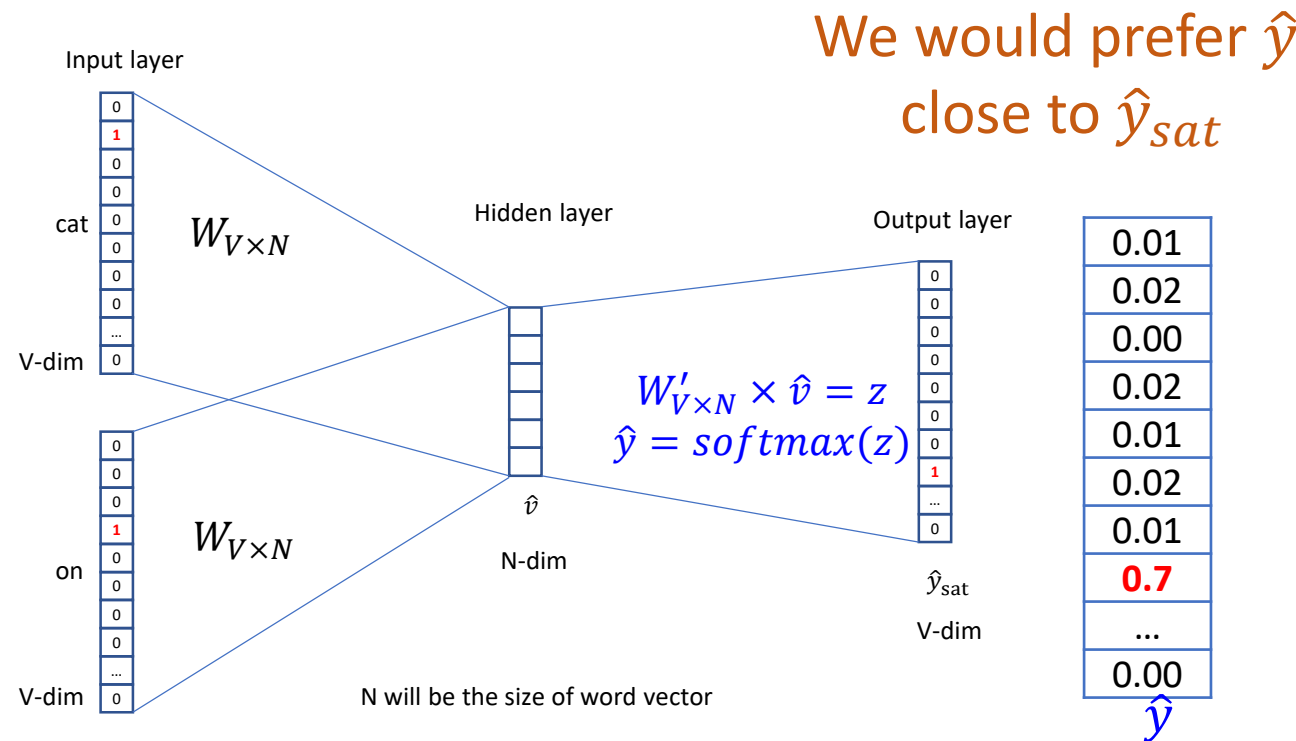


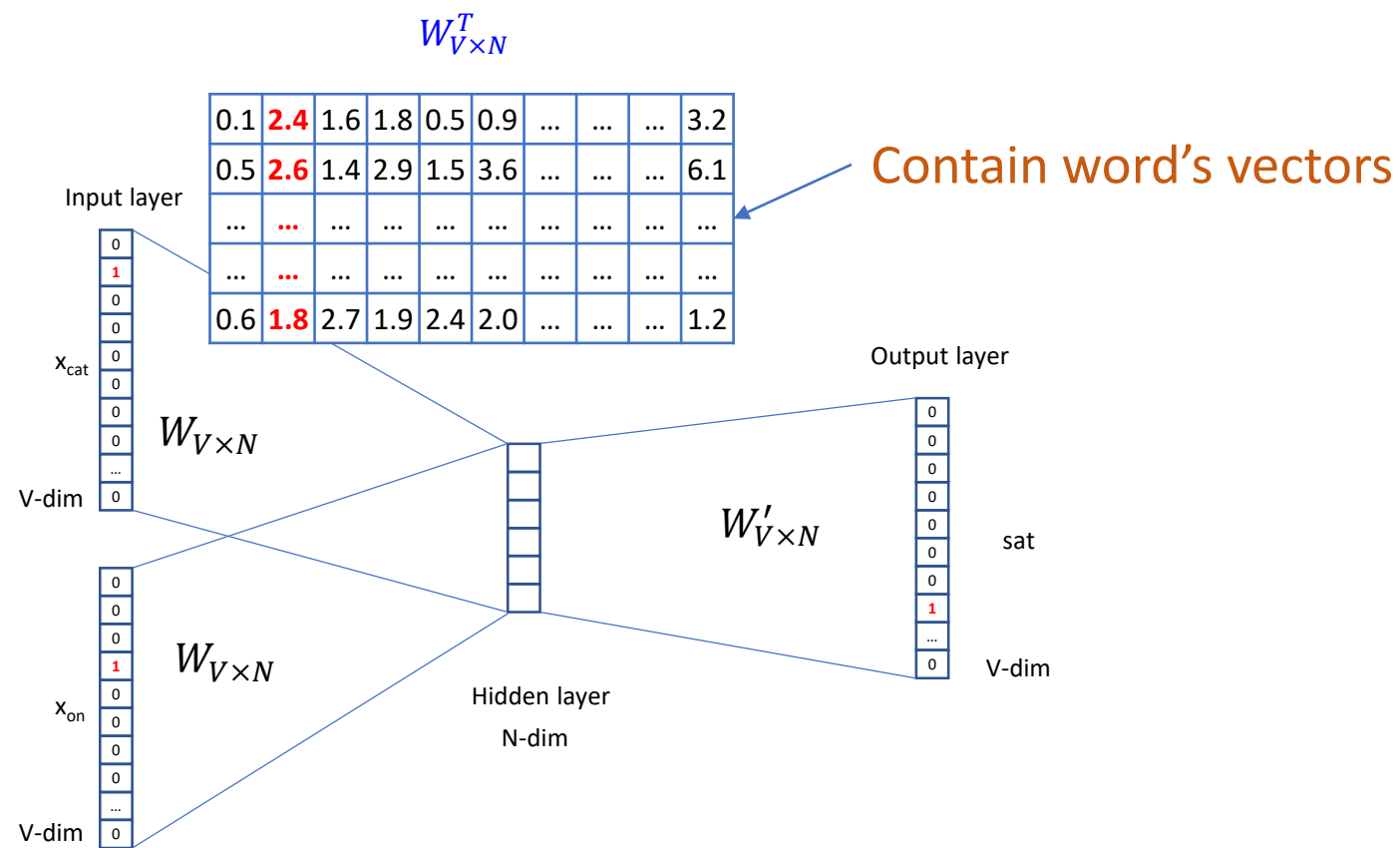








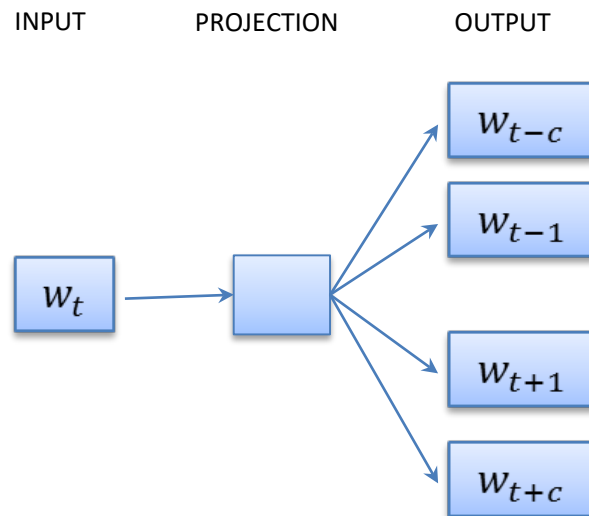




We can consider either W or W' as the word's representation. Or even take the average.

Skipgram Model

- Input: Central word w_t
- Output: Words in its context: w_{con}
 $\{w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}\}$
- Each input word represented by a 1-hot encoding of size V



Source Text:

Deep Learning attempts to learn multiple levels of representation from data.

Input output pairs :

Positive samples:

(representation, levels)

(representation, of)

(representation, from)

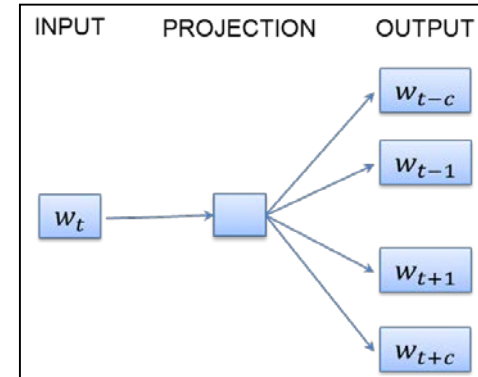
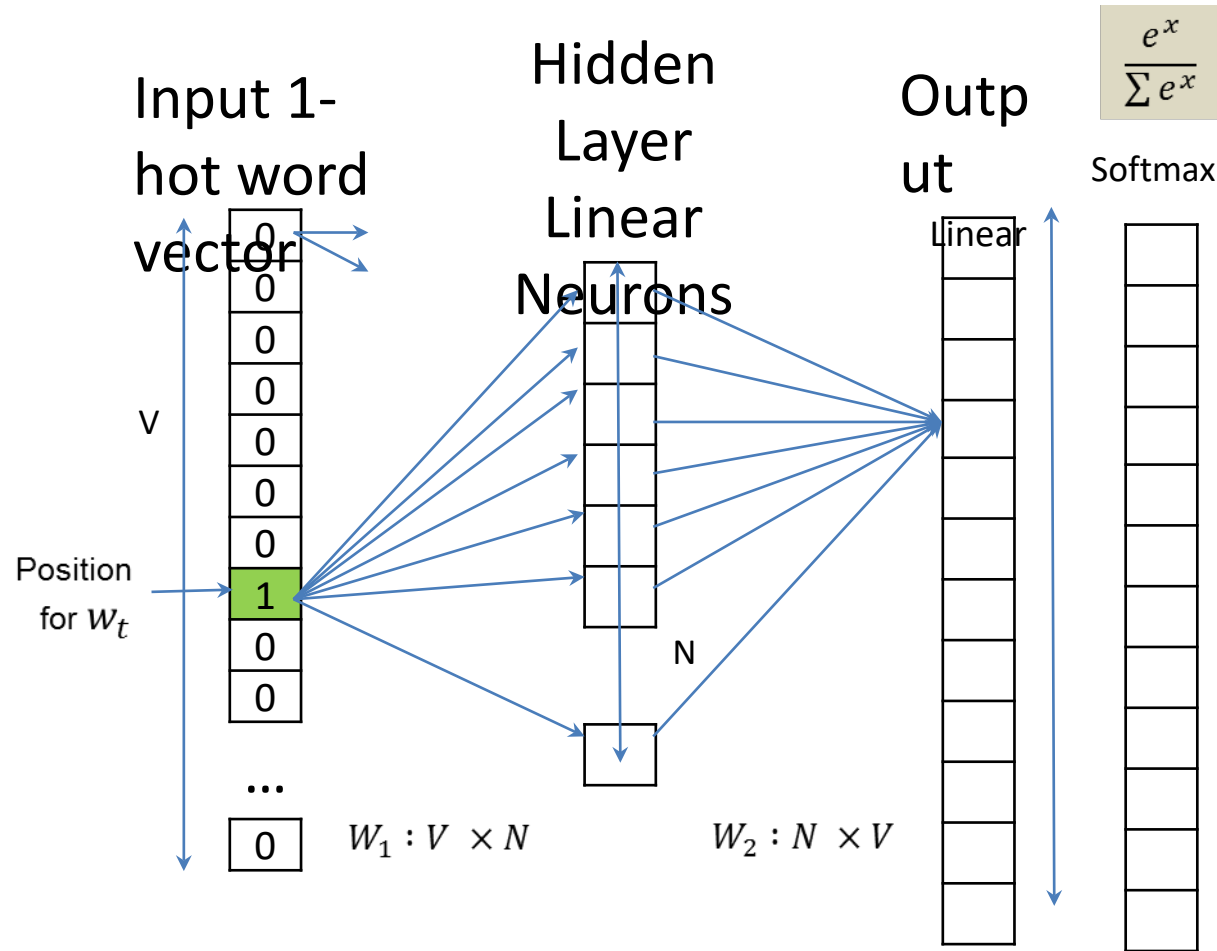
(representation, data)

Negative samples:

(representation, x)

[x: all other words except the 4 positive]

Skipgram Model



Probability that the word in a context position is w_i

Positive sampling
Negative sampling

Skipgram: Loss function

$$\text{Maximize} \quad \frac{1}{T} \sum_{t=1}^T \sum_{\text{context}} \log p(w_{\text{context}} | w_t)$$

$p(w_{\text{con}} | w_t)$ is the output of softmax classifier

$$p(w_{\text{con}} | w_t) = \frac{\exp(v'_{w_{\text{con}}} \cdot v_{w_t})}{\sum_{w=1}^W \exp(v'_w \cdot v_{w_t})}$$

Let the model parameters be θ . The solution is given by

$$\begin{aligned} & \underset{\theta}{\operatorname{argmax}} \sum_{(w_t, w_c) \in D} \log p(w_{\text{con}} | w_t; \theta) \\ &= \sum_{(w_t, w_c) \in D} \left(\log e^{(v'_{w_{\text{con}}} \cdot v_{w_t})} - \log \sum_x e^{(v'_x \cdot v_{w_t})} \right) \end{aligned}$$



Skip-Grams with Negative Sampling (SGNS)

Marco saw a furry little **wampimuk** hiding in the tree.



Skip-Grams with Negative Sampling (SGNS)

Marco saw a furry little wampimuk hiding in the tree.

words

wampimuk

wampimuk

wampimuk

wampimuk

...

contexts

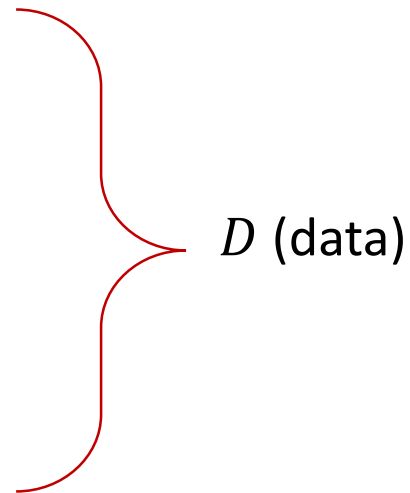
furry

little

hiding

in

...



“word2vec Explained...”
Goldberg & Levy, arXiv 2014

Skip-Grams with Negative Sampling (SGNS)

Maximize: $\sigma(\vec{w} \cdot \vec{c})$

- c was **observed** with w

words

wampimuk

wampimuk

wampimuk

wampimuk

contexts

furry

little

hiding

in

Minimize: $\sigma(\vec{w} \cdot \vec{c}')$

- c' was **hallucinated** with w

words

wampimuk

wampimuk

wampimuk

wampimuk

contexts

Australia

cyber

the

1985

Some interesting results

Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{||w_b - w_a + w_c||}$$

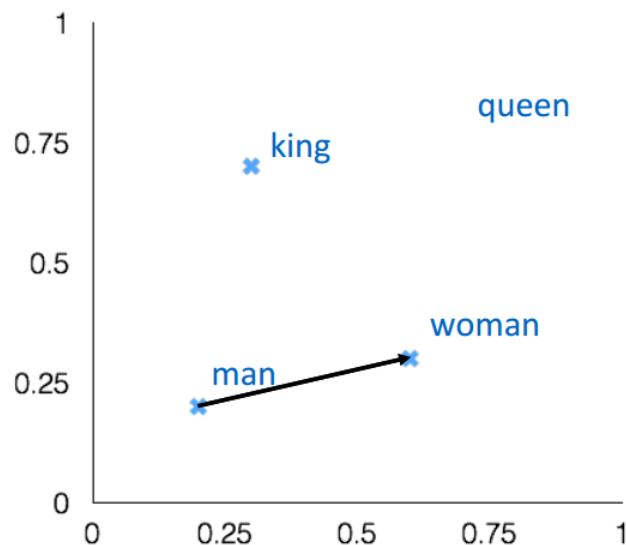
man:woman :: king:?

+ king [0.30 0.70]

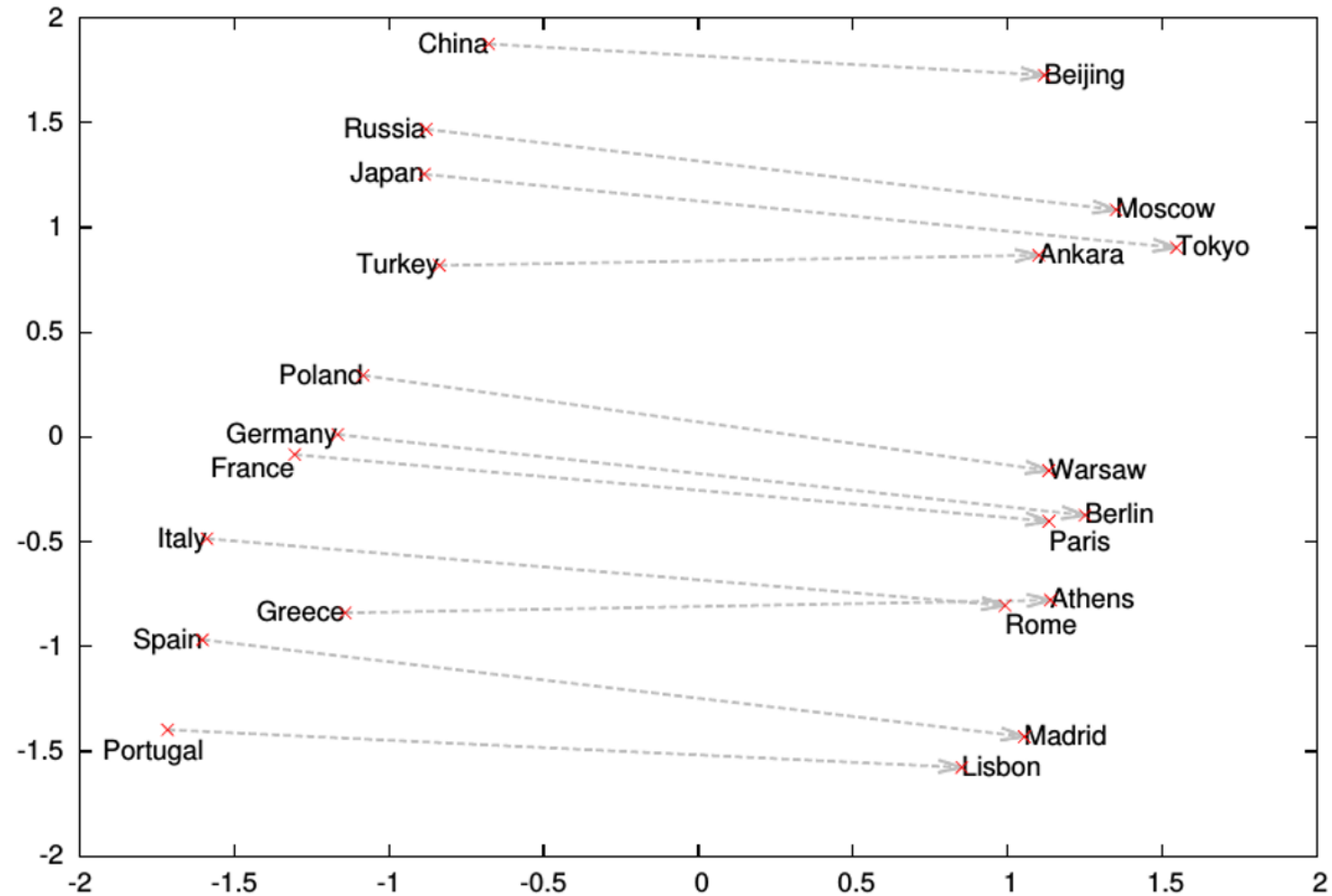
- man [0.20 0.20]

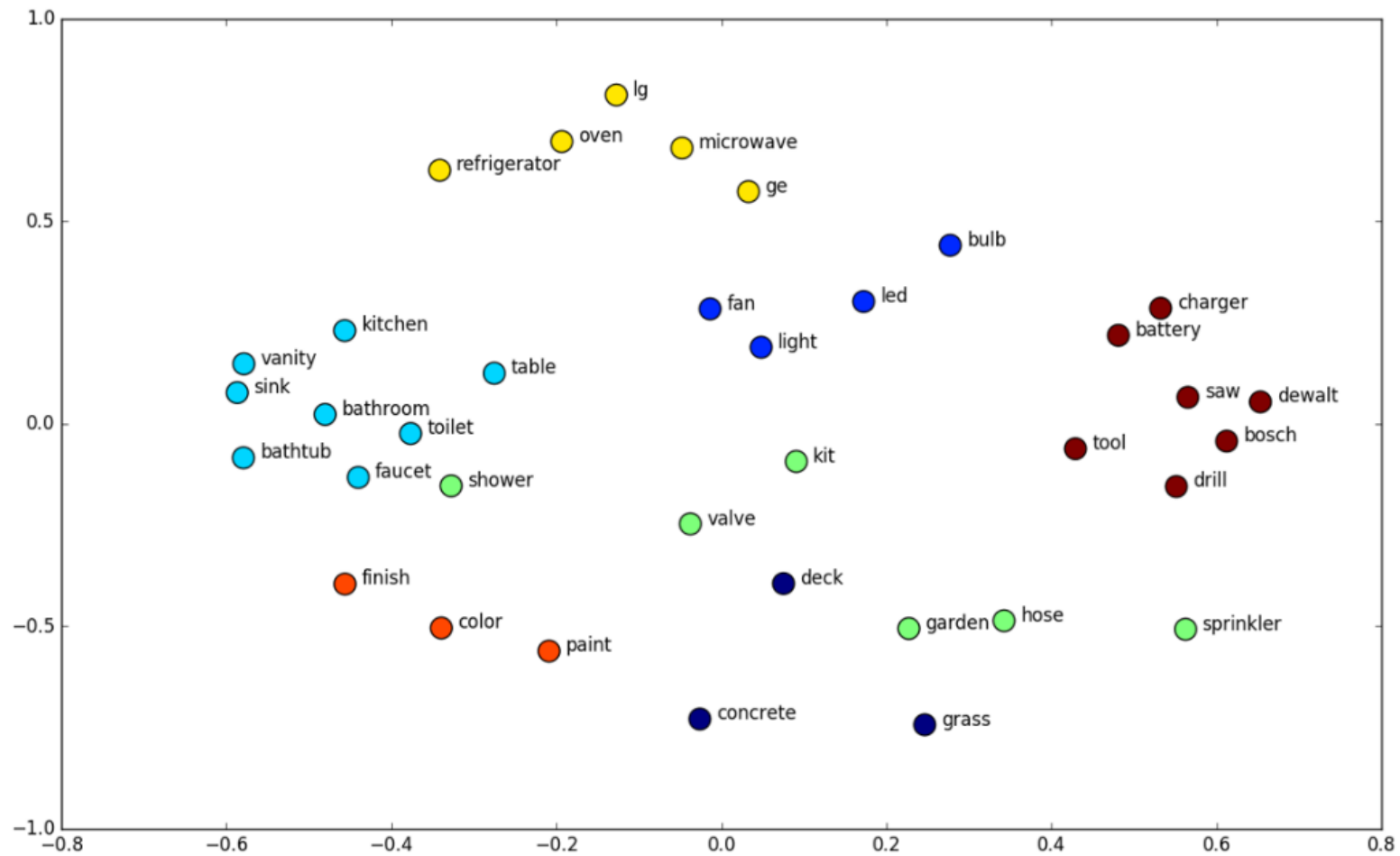
+ woman [0.60 0.30]

queen [0.70 0.80]

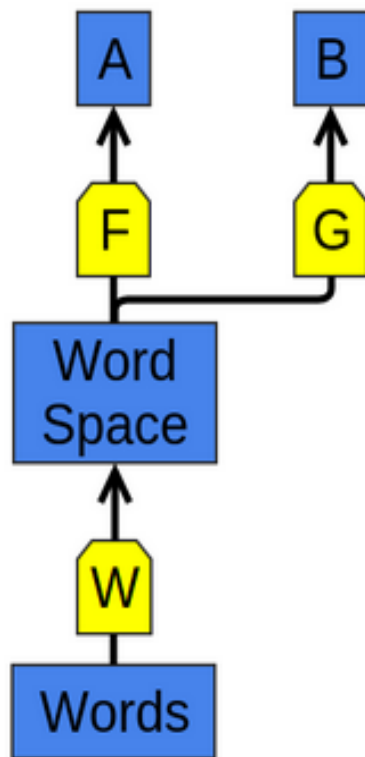


Word analogies





Word embedding applications

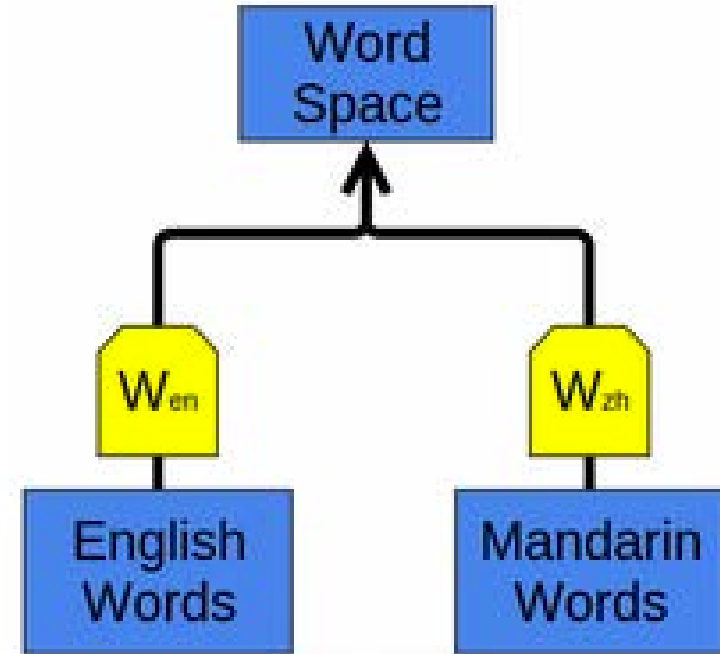


W and F learn to perform task A. Later, G can learn to perform B based on W .

- The use of word representations... has become a key “secret sauce” for the success of many NLP systems
 - across tasks including named entity recognition, part-of-speech tagging, parsing, and semantic role labeling
- Learning a good representation on a task A and then using it on a task B is one of the major tricks in the Deep Learning toolbox.
 - Pretraining, transfer learning, and multi-task learning.
 - Can allow the representation to learn from more than one kind of data.

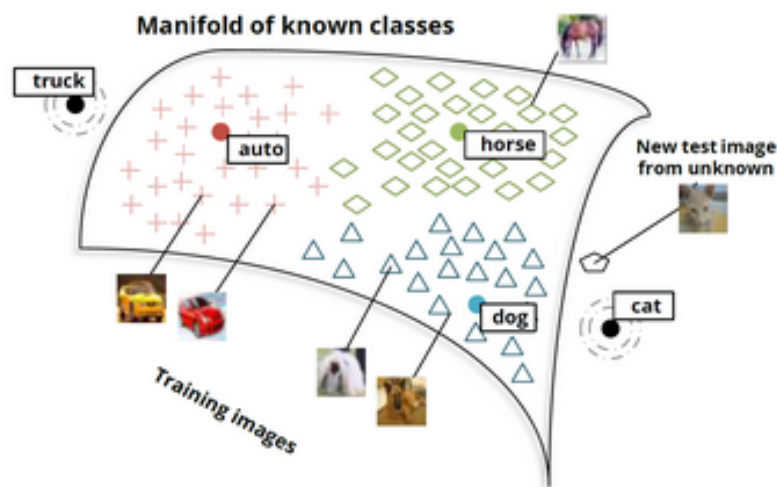
Word embedding applications

- Can learn to map multiple kinds of data into a single representation.
 - E.g., bilingual word-embedding
- Embed as above, but words that are known as close translations should be close together.
- Words we didn't know were translations end up close together!
- Structures of two languages get pulled into alignment.

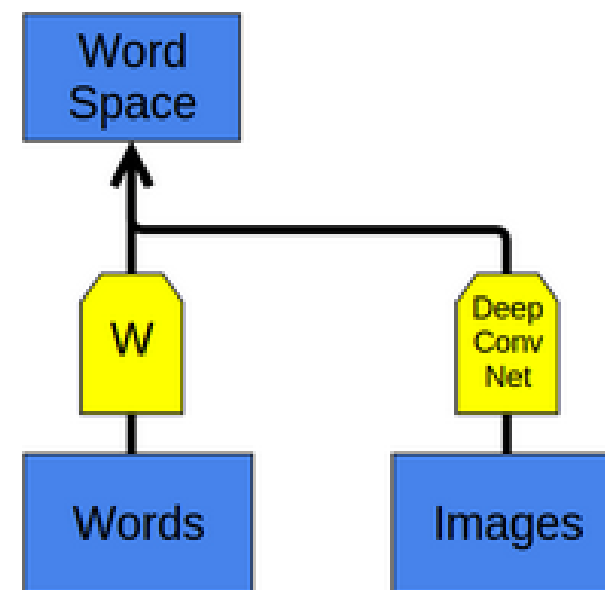


Word embedding applications

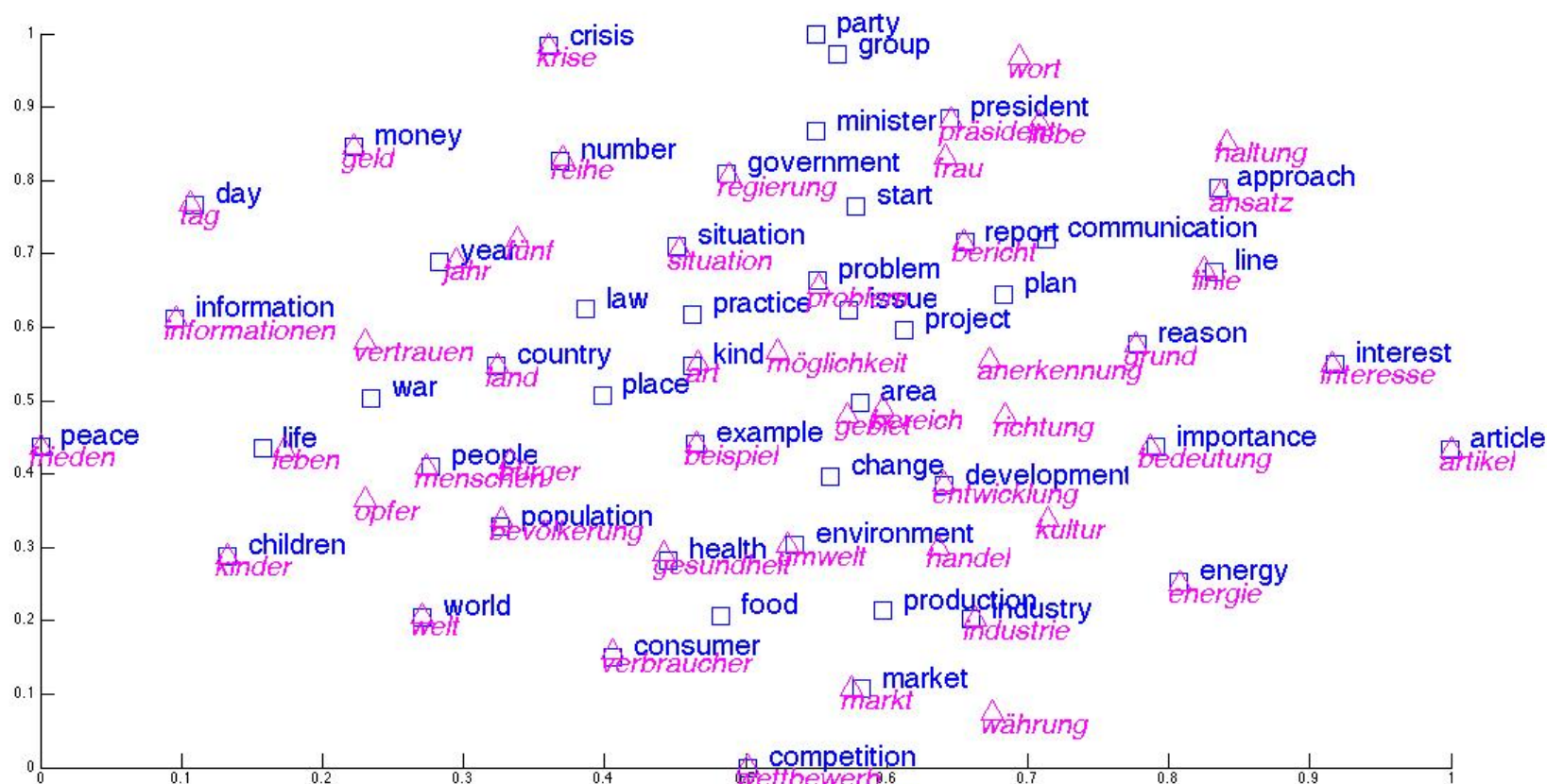
- Can apply to get a joint embedding of words and images or other multi-modal data sets.
- New classes map near similar existing classes: e.g., if 'cat' is unknown, cat images map near dog.



(Socher *et al.* (2013b))



Multilingual Embeddings



Luong et al 2015

An interesting application

- Lawrence Berkeley lab material scientists applied word embedding to 3.3 million scientific abstracts published between 1922-2018.
 - 500k words. Vector size: 200 dimension, skip-gram model
- Captured things like periodic table and structure-property relationship in materials:
 $\text{ferromagnetic} - \text{NiFe} + \text{IrMn} \approx \text{antiferromagnetic}$
- Discovered new thermoelectric materials

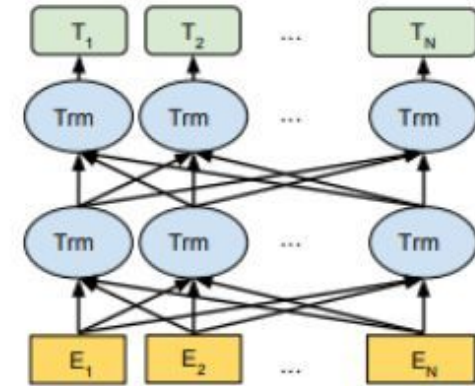
Nature, July 2019 V. Tshitonya et al, “Unsupervised word embeddings capture latent knowledge from materials science literature”.

Contextualized Word Vectors

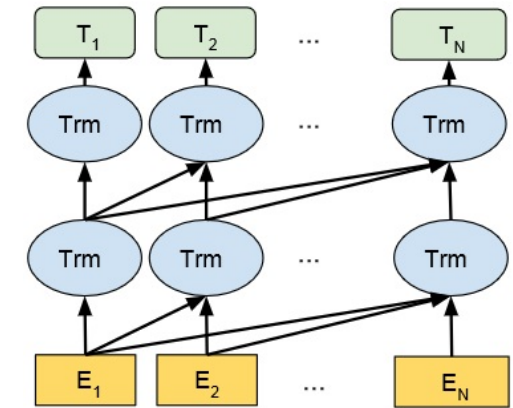
Incorporating context into word embeddings
a watershed idea in NLP

- BERT: Bidirectional Encoder Representations from Transformers (BERT, 2018)
- GPT-2/3

Led to significant improvements on virtually every NLP task.



BERT Architecture



GPT

Language Models

How likely is a sentence (w_1, w_2, \dots, w_n) ?

- Predict the next word
- Complete the sentence

$P(\text{I saw a bus}) \gg P(\text{eyes awe a boss})$

Pre-Trained Language Models

- Instead of training the model from scratch, you can use another pre-trained model as the basis and only fine-tune it to solve the specific NLP task.

Multilinguality

mBERT

Google's multilingual BERT model
generates language-independent cross-
language sentence embeddings for 109
languages

