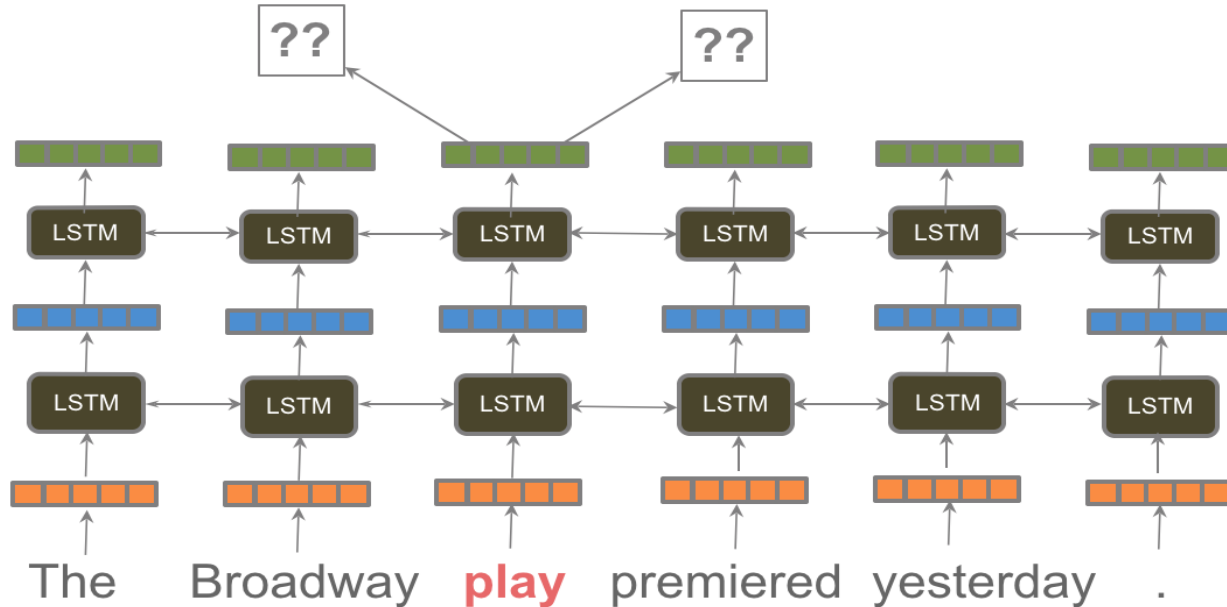


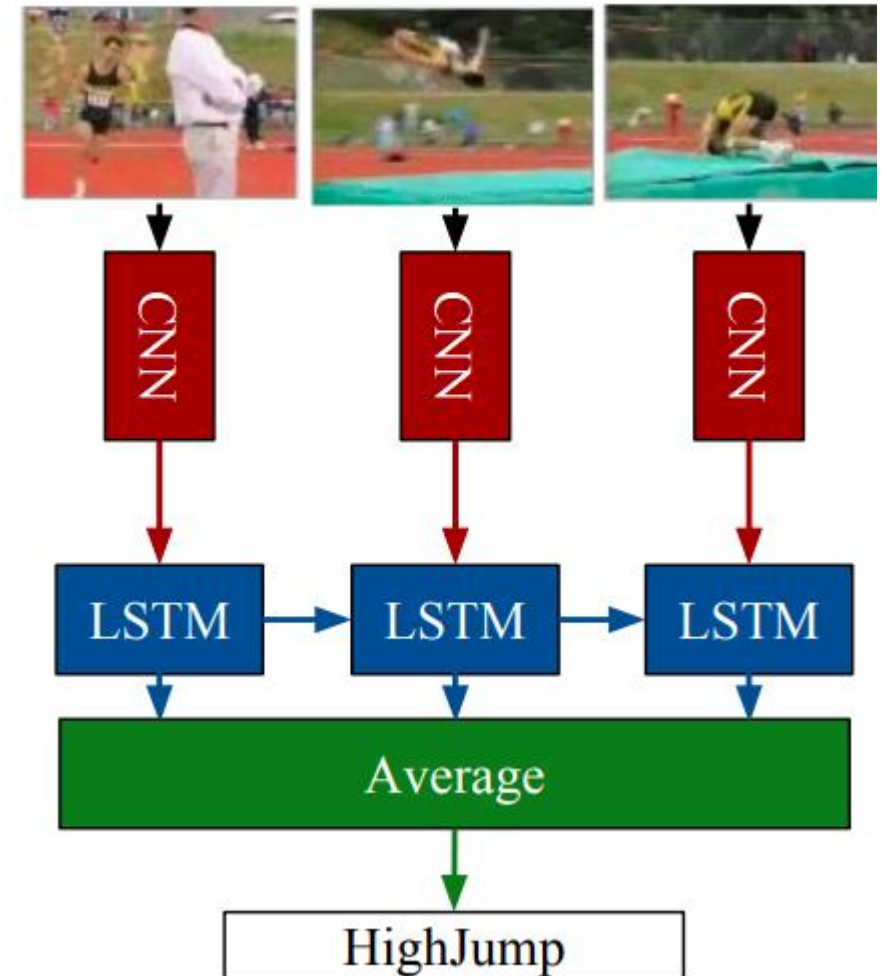
RNN

Recurrent neural network

- cho bài toán dữ liệu dạng chuỗi (sequence).
- video understanding
- Natural language processing



Activity Recognition Sequences in the Input

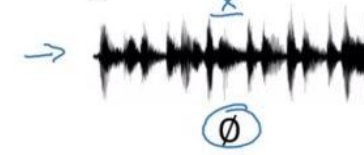


Sequence data

time-series data.

Examples of sequence data

Speech recognition



“The quick brown fox jumped
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis → AGCCCCTGTGAGGAACTAG

AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec
moi?

Do you want to sing with
me?

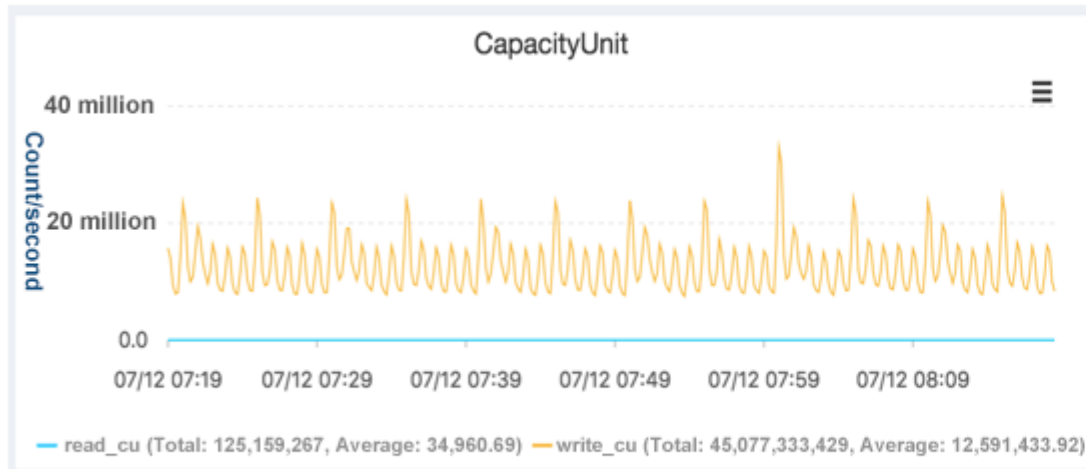
Video activity recognition



Running

Name entity recognition → Yesterday, Harry Potter
met Hermione Granger.

Yesterday, **Harry Potter**
met **Hermione Granger**.
Andrew Ng



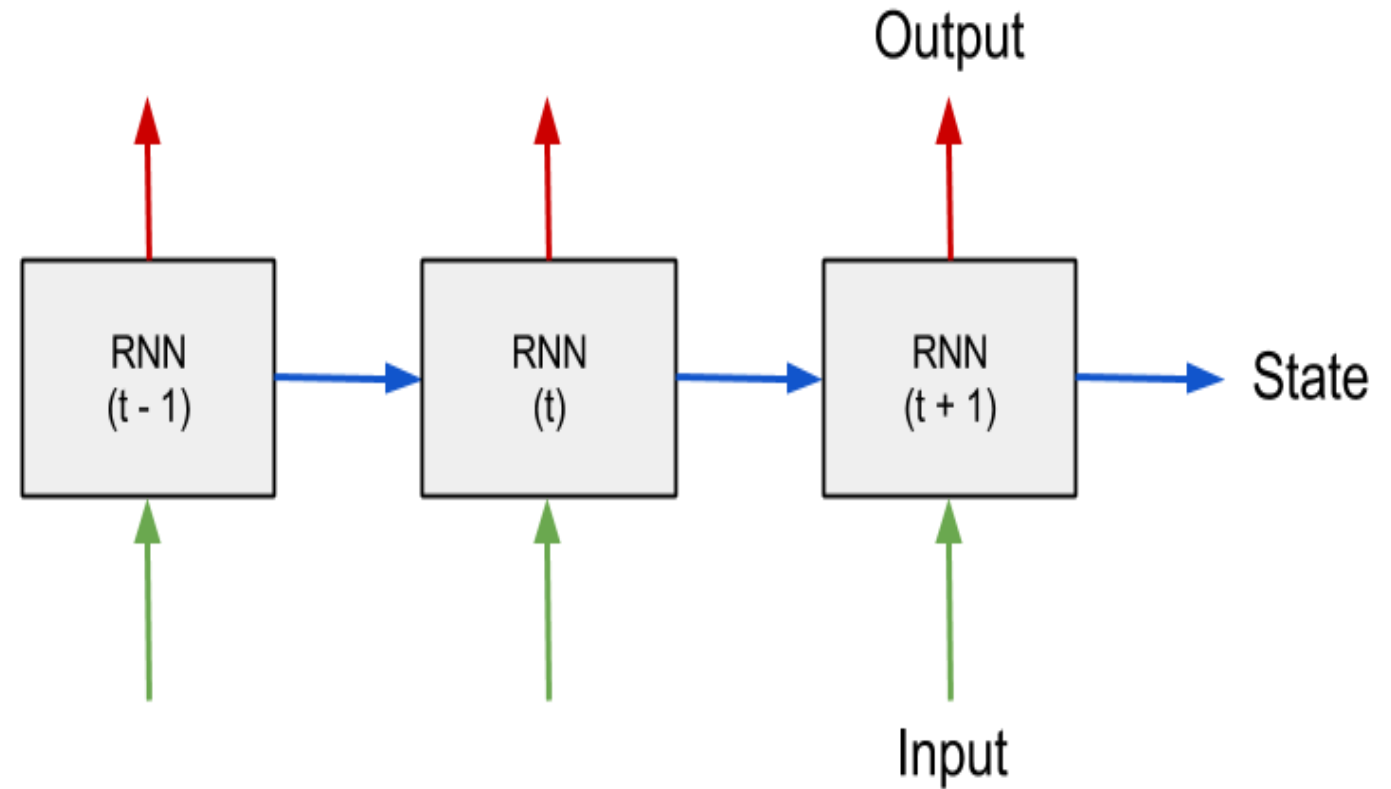
Monitoring time series data

Order No.	Time	Location	Event
77165205838	2018-07-10 10:00	Shanghai	Ship
77165205838	2018-07-11 12:00	Yuhang	Arrived in Hangzhou
77165205838	2018-07-11 14:00	Zhuantang	Start delivery and arrive at Apsara Park
77165205838	2018-07-11 16:00	Zhuantang	Accepted

Status time series data

Recurrent neural network

Recurrent Neural Networks (RNNs) are a kind of neural network that specialize in processing **sequences**.

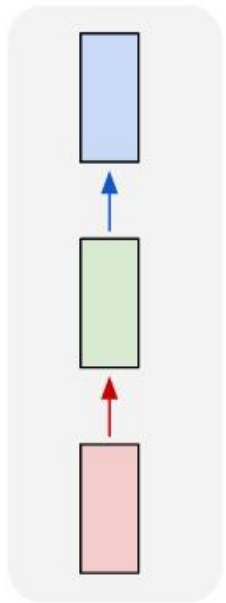


Phân loại bài toán RNN

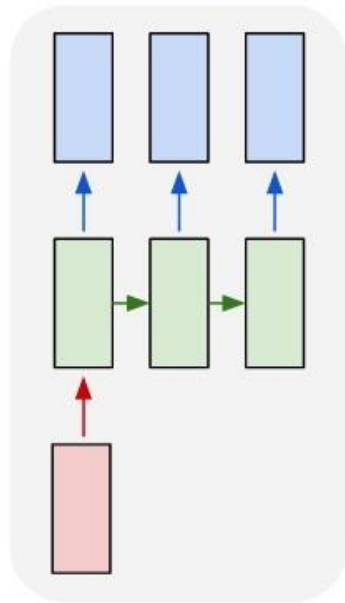
CNN: **fixed-size inputs and fixed-size outputs.**

RNNs are useful because they let us have **variable-length sequences** as both inputs and outputs.

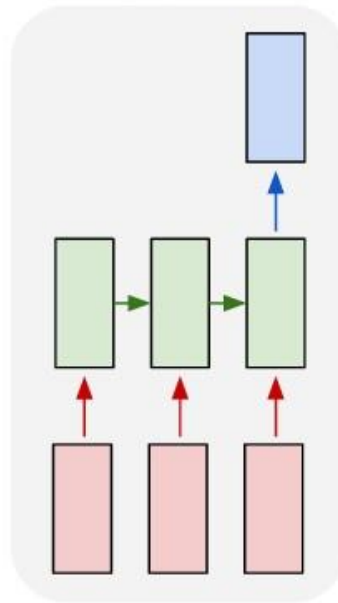
one to one



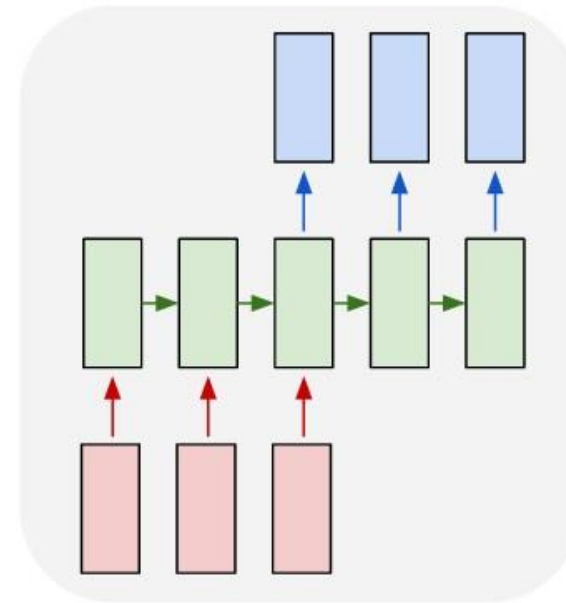
one to many



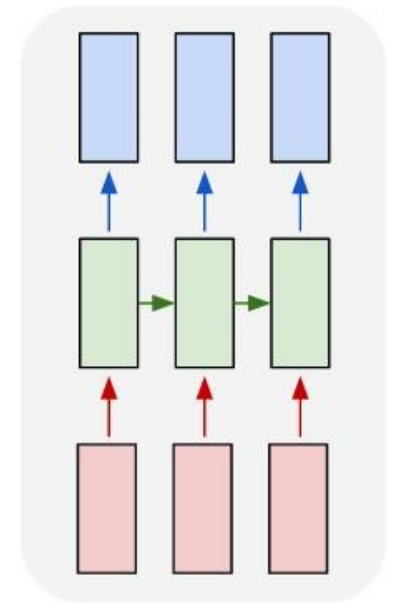
many to one



many to many



many to many



Process Sequences

- **One-to-one:** This is the classic feed forward neural network architecture, with one input and we expect one output.
- **One-to-many:** image captioning, one image as a fixed size input and the output can be words or sentences which are variable in length.



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."

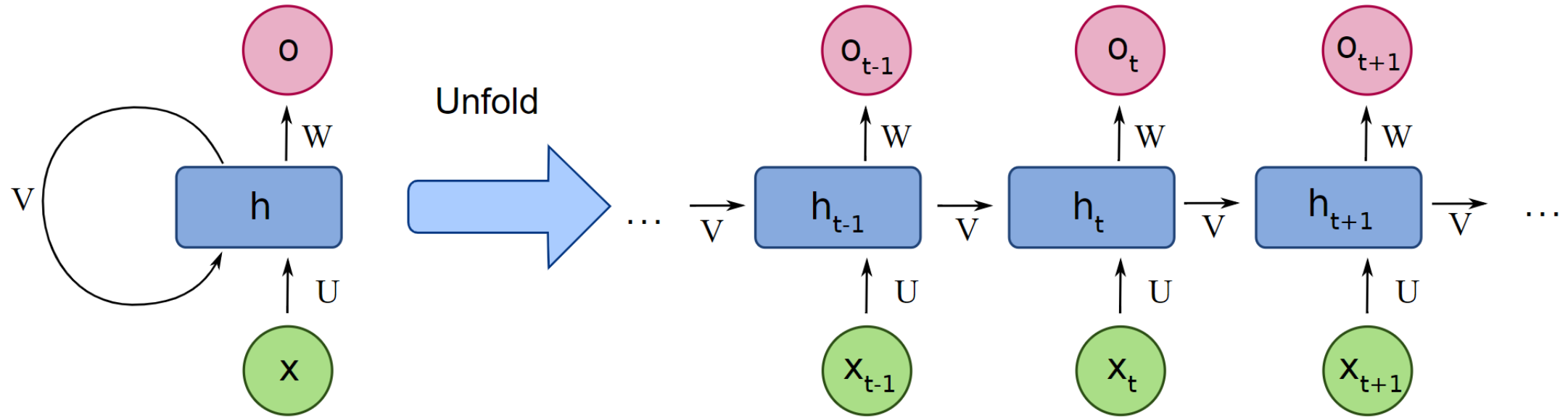


"two young girls are playing with lego toy."

Process Sequences

- **Many-to-one:** This is used for sentiment classification. The input is expected to be a sequence of words or even paragraphs of words. The output can be a regression output with continuous values which represent the likelihood of having a positive sentiment.
- **Many-to-many:** this model is ideal for machine translation like the one we see on Google translate. The input could be an English sentence which has variable length and the output will be the same sentence in a different language which also has variable length.

RNN model



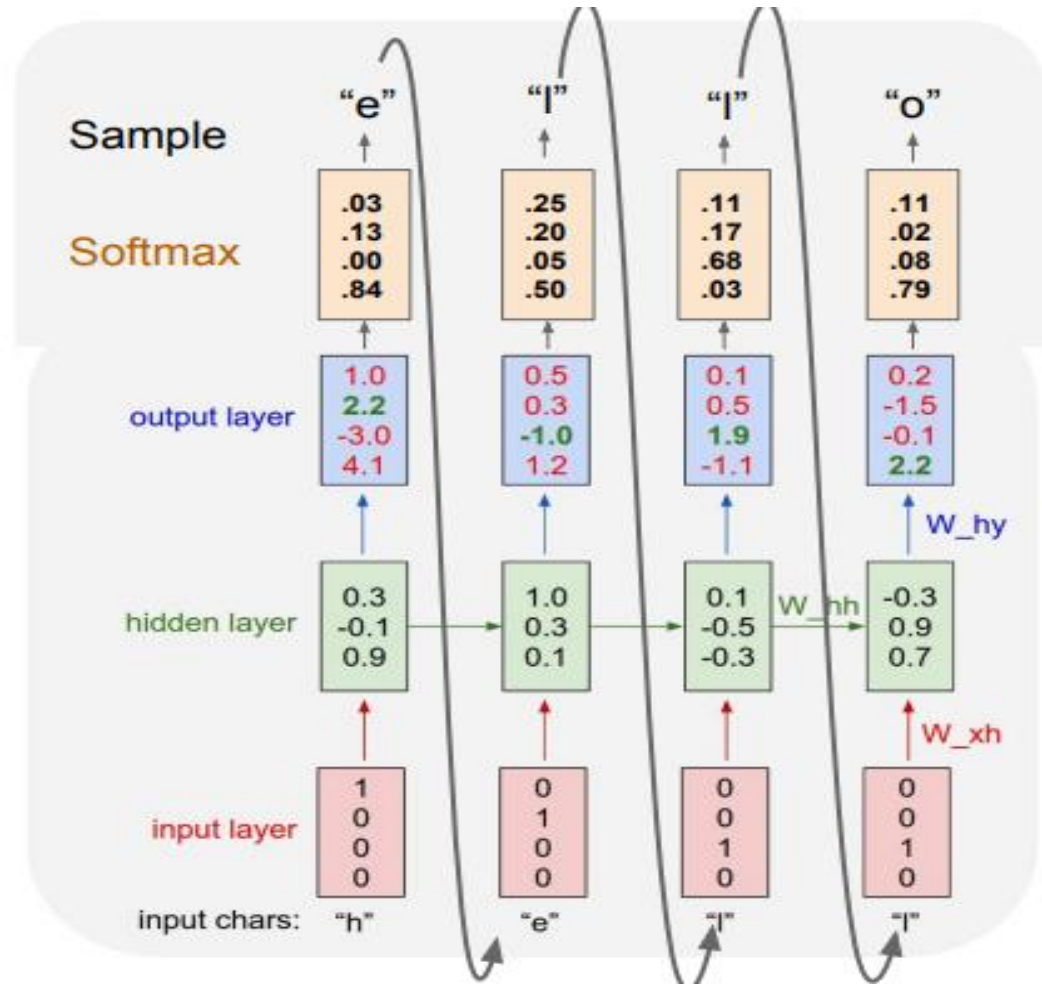
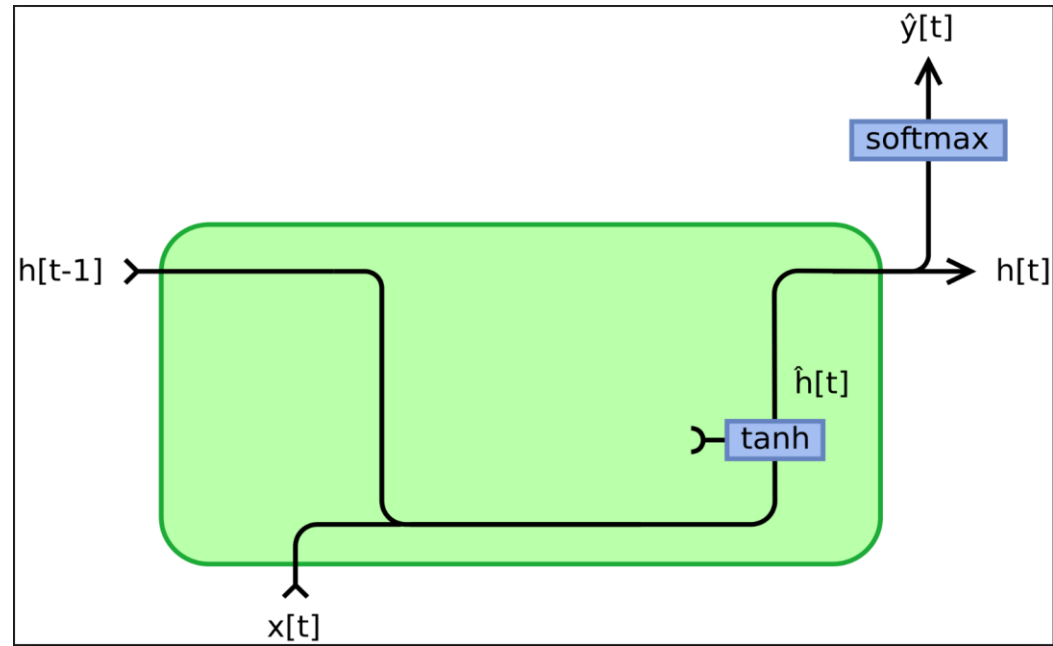
x : The input. It can be a word in a sentence or some other type of sequential data

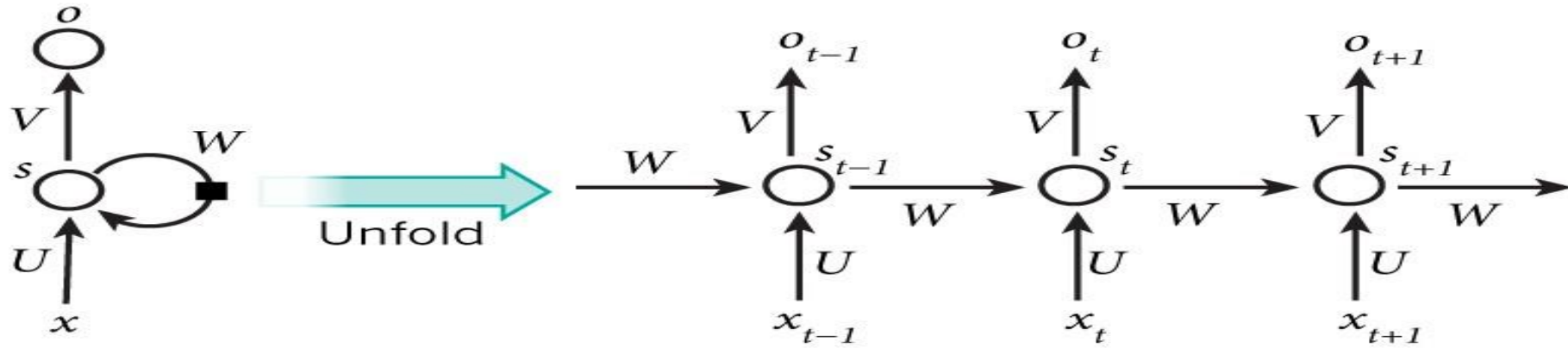
o : The output. For instance, what the network thinks the next word on a sentence should be given the previous words

h : The main block of the RNN. It contains the weights and the activation functions of the network

v : Represents the communication from one time-step to the other.

RNN model





$$s_t = \tanh(Ux_t + Ws_{t-1})$$

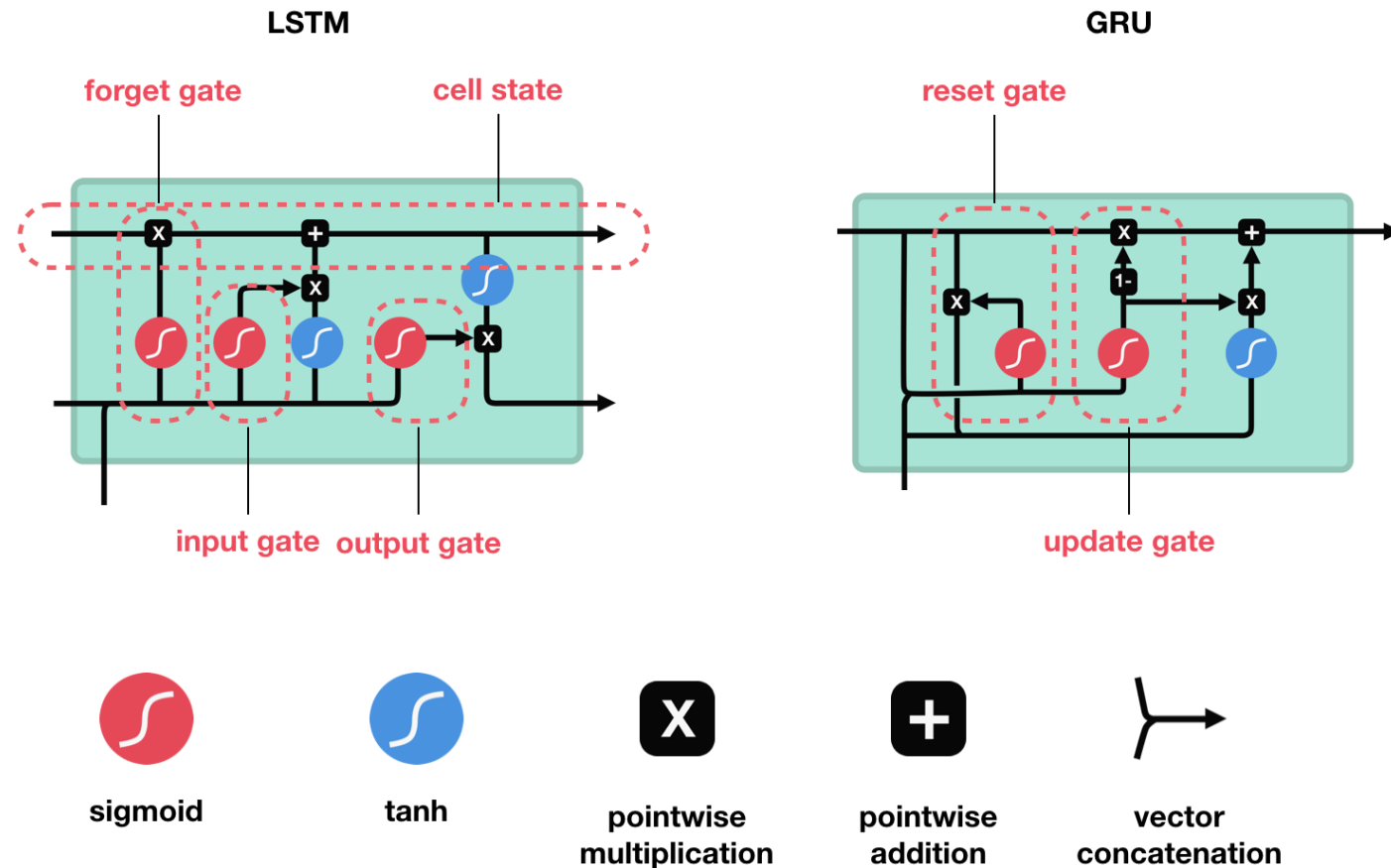
$$o_t = \text{softmax}(Vs_t)$$

Cross-Entropy Loss Function :
categorical cross entropy loss

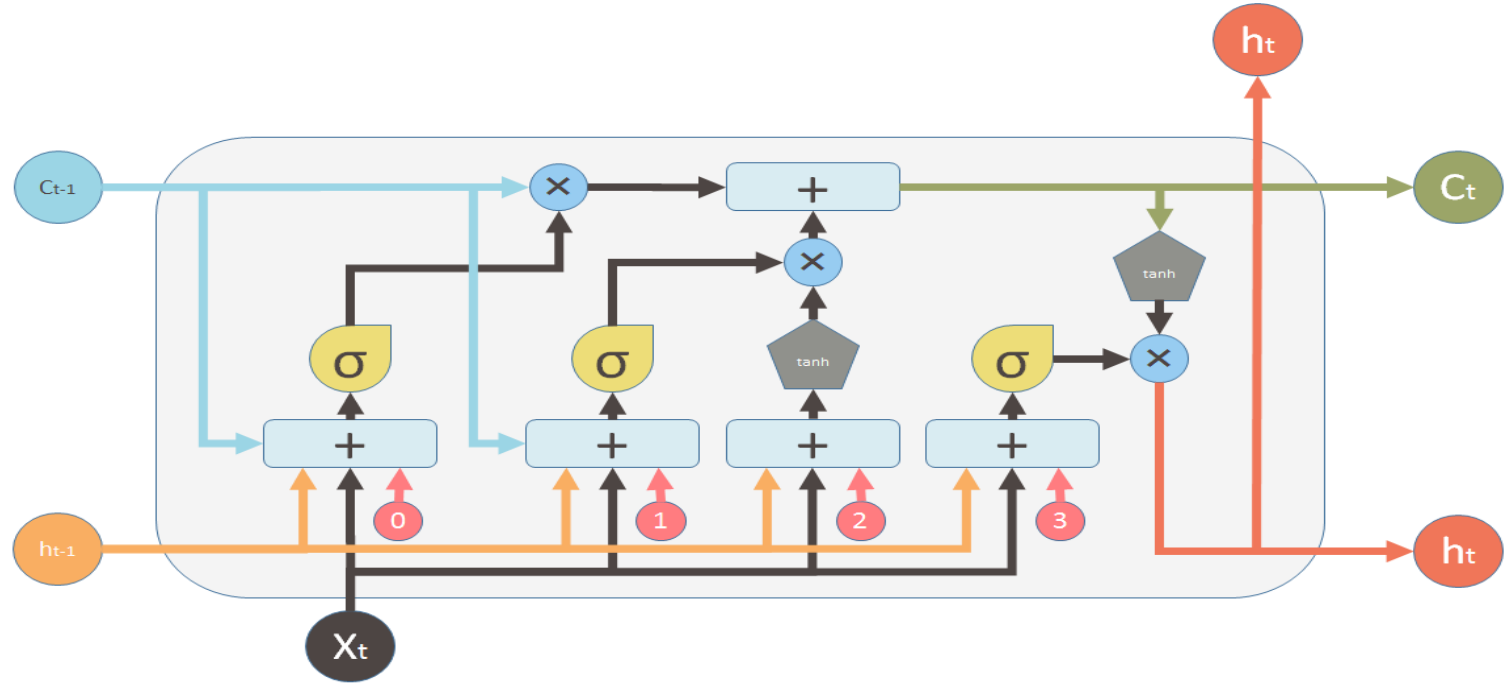
$$L(y, o) = -\frac{1}{N} \sum y_n \log o_n$$

RNN model

- Kiến trúc khá đơn giản nên khả năng liên kết các thành phần có khoảng cách xa trong câu không tốt do gradient bị thấp dần trong quá trình học (vanishing gradient)
- Không có cơ chế lọc những thông tin không cần thiết => Bộ nhớ của kiến trúc có hạn, nếu lưu tất cả những chi tiết không cần thiết thì sẽ dẫn đến quá tải
- Các kiến trúc để khắc phục các nhược điểm của RNN: LSTM và GRU



Mỗi module của 2 kiến trúc trên đều có trang bị các cổng (gate), giúp kiến trúc đánh giá được mức độ quan trọng của thông tin, từ đó đưa ra quyết định giữ lại hay bỏ đi.



Inputs:



Input vector



Memory from previous block



Output of previous block

outputs:



Memory from current block



Output of current block

Nonlinearities:



Sigmoid



Hyperbolic tangent

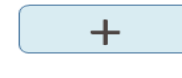
Bias:



Vector operations:



Element-wise multiplication



Element-wise Summation / Concatenation

LSTM

Output: c_t, h_t, c cell state, h là hidden state.

Input: c_{t-1}, h_{t-1}, x_t , trong đó: x_t - input ở state thứ t , c_{t-1}, h_{t-1} - output của layer trước.

f_t, i_t, o_t tương ứng với forget gate, input gate và output gate.

Forget gate: $f_t = \sigma(U_f * x_t + W_f * h_{t-1} + b_f)$

Input gate: $i_t = \sigma(U_i * x_t + W_i * h_{t-1} + b_i)$

Output gate: $o_t = \sigma(U_o * x_t + W_o * h_{t-1} + b_o)$

$0 < f_t, i_t, o_t < 1$; W, U - ma trận trọng số,

b_f, b_i, b_o là các hệ số bias.

$\tilde{c}_t = \tanh(U_c * x_t + W_c * h_{t-1} + b_c)$

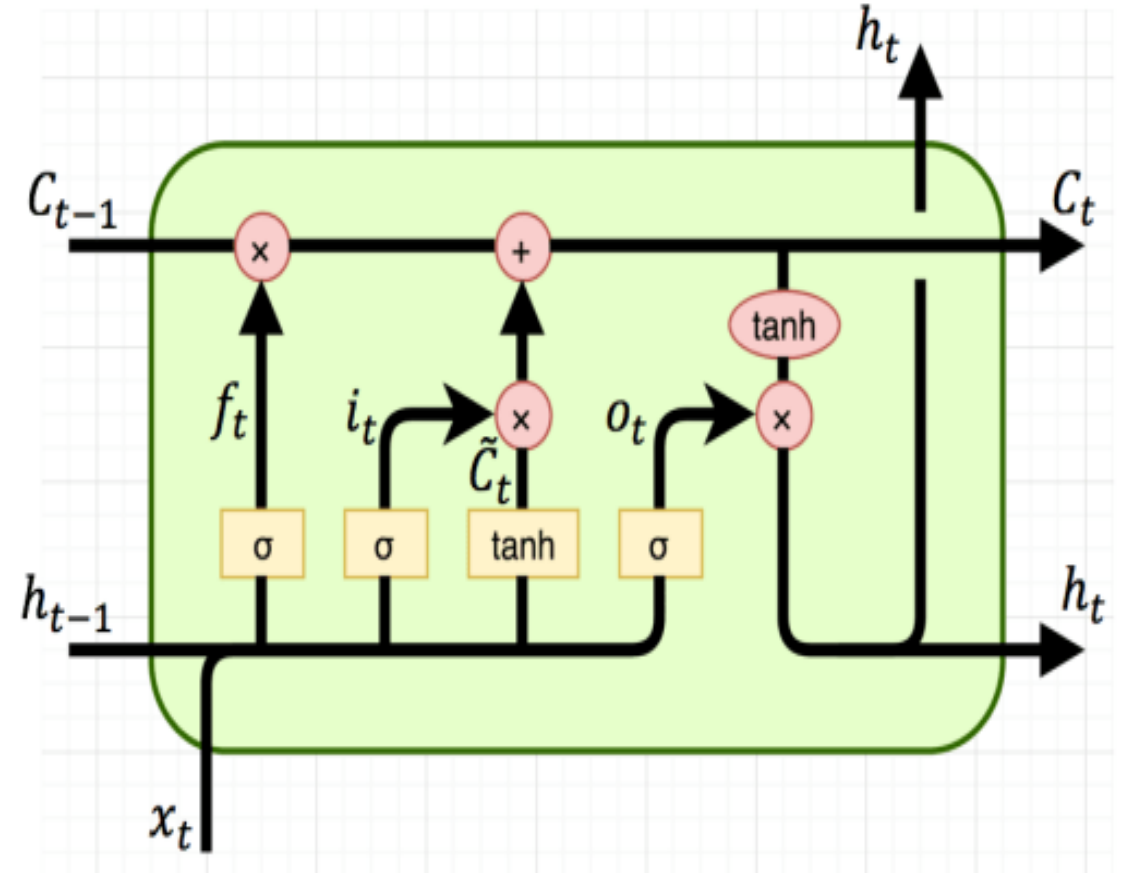
$c_t = (f_t * c_{t-1} + i_t * \tilde{c}_t)$

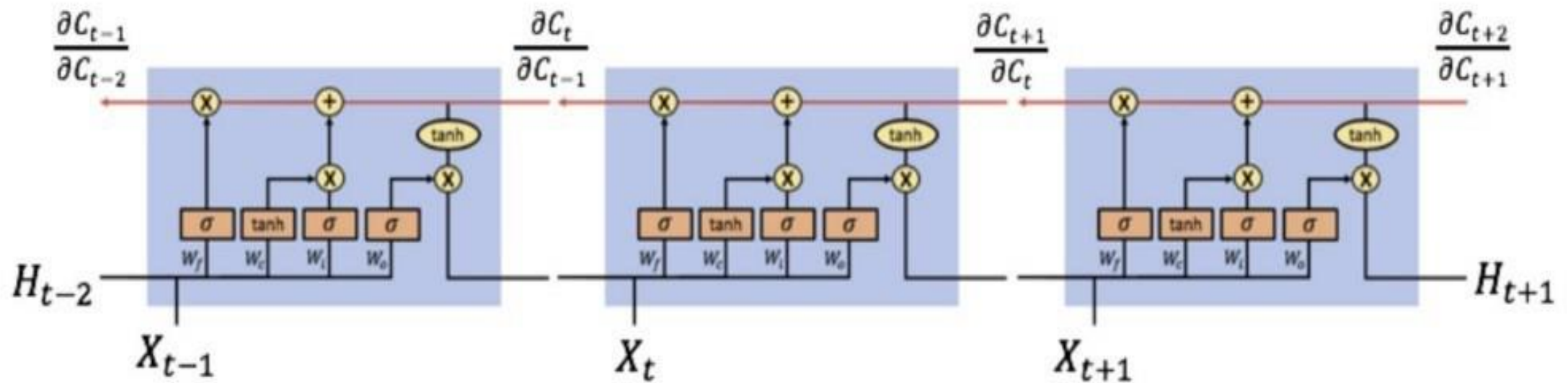
forget gate quyết định xem cần lấy bao nhiêu từ cell state trước và **input gate** sẽ quyết định lấy bao nhiêu từ input của state và hidden layer của layer trước.

$h_t = o_t * \tanh(c_t)$

h_t, \tilde{c}_t khá giống với RNN, nên model có **short term memory**.

thông tin nào cần quan trọng và dùng ở sau sẽ được gửi vào và dùng khi cần => có thể mang thông tin từ đi xa => **long term memory**.



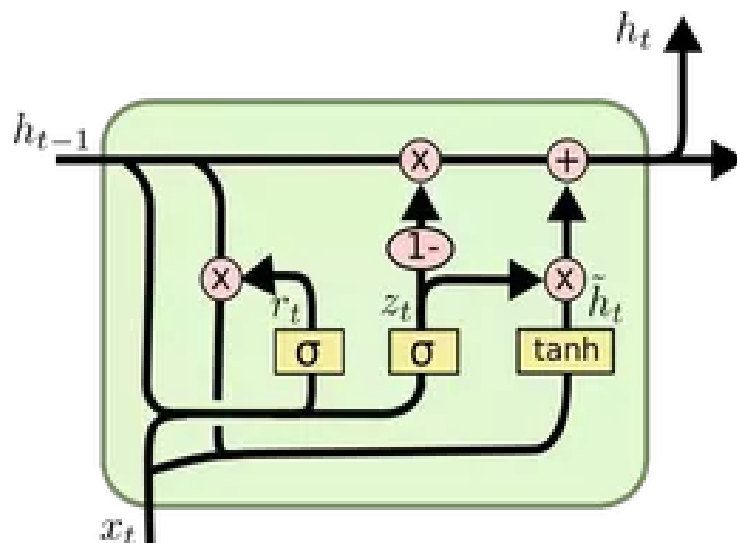


$\frac{\partial c_t}{\partial c_{t-1}} = f_t$, Do $0 < f_t < 1$, LSTM vẫn bị vanishing gradient nhưng bị ít hơn so với RNN.

khi mang thông tin trên cell state thì ít khi cần phải quên giá trị cell cũ, nên $f_t \approx 1 \Rightarrow$ **Tránh được vanishing gradient.**

GRU

- GRU là một phiên bản của LSTM với nguyên tắc hoạt động tương tự như LSTM nhưng có cấu tạo đơn giản hơn.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

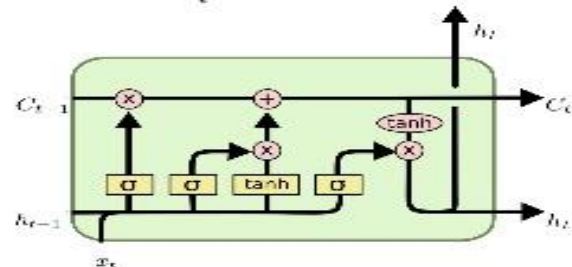
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

LSTM and GRU

- LSTM có thể lưu trữ thông tin với dữ liệu dài hơn so với GRU.
- Do cấu tạo đơn giản của mình, GRU thường xử lý nhanh hơn LSTM và có thể dễ dàng sử dụng để xây dựng các mạng có cấu trúc phức tạp.
- Hoạt động theo một chiều nhất định (forward direction) => chỉ mang thông tin tính tới thời điểm hiện tại.
- Trong bài toán NLP thì việc biết thông tin của các timesteps tiếp theo giúp cải thiện rất nhiều kết quả output (Translation, Speech recognition, Handwritten recognition,...)

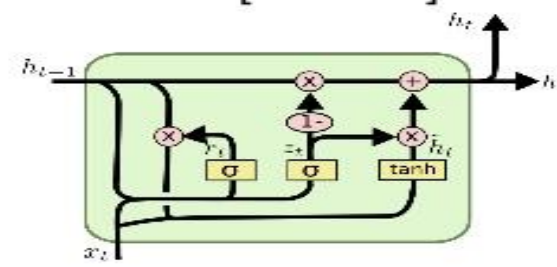
LSTM and GRU

• LSTM [Hochreiter&Schmidhuber97]



$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\
 \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \\
 C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t \\
 o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\
 h_t &= o_t * \tanh(C_t)
 \end{aligned}$$

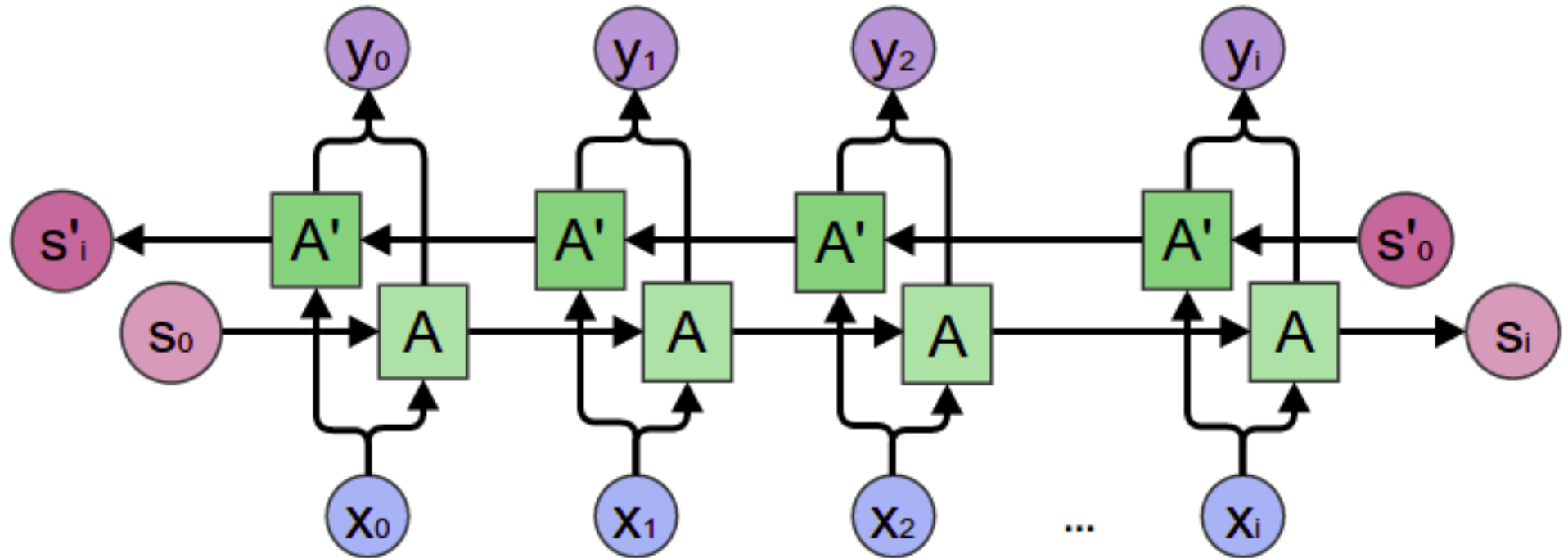
• GRU [Cho+14]



$$\begin{aligned}
 z_t &= \sigma(W_z \cdot [h_{t-1}, x_t]) \\
 r_t &= \sigma(W_r \cdot [h_{t-1}, x_t]) \\
 \tilde{h}_t &= \tanh(W \cdot [r_t * h_{t-1}, x_t]) \\
 h_t &= (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t
 \end{aligned}$$

Tohoku University, Inui and Okazaki Lab. (Biases are omitted.)
Sosuke Kobayashi

Bi-directional RNN



Một số bài toán Xử lý ngôn ngữ tự nhiên

- Answer questions using the Web
- Translate documents from one language to another
- Help make informed decisions
- Follow directions given by any user
- Fix your spelling or grammar
- Do library research; summarize
- Manage messages intelligently
- Grade exams
- Write poems or novels
- Listen and give advice
- Estimate public opinion
- Read everything and make predictions
- Interactively help people learn
- Help disabled people
- Help refugees/disaster victims
- Document or reinvigorate indigenous languages

Một số bài toán Xử lý ngôn ngữ tự nhiên

Rút trích thông tin văn bản (Information extraction):

- Web mining: rút trích tên người nổi tiếng, sản phẩm đang hot, so sánh giá sản phẩm, nghiên cứu đối thủ cạnh tranh, phân tích tâm lý khách hàng
- Biomedical, Business intelligent, Financial professional : đánh giá thị trường từ các nguồn khác nhau: giá xăng dầu tăng giảm, thông tin chiến tranh, chính trị giữa các nước,
- Terrorism event: sử dụng vũ khí gì, đối tượng tấn công là ai

Bài toán:

Rút trích tên thực thể (**Named entity recognition** – NER: people, organization, location)

Rút trích quan hệ giữa hai thực thể (**Relation extraction** – RE: founderOf, headquarteredIn).

The screenshot shows the Brat NER interface with a list of sentences and extracted entities and relations. The interface includes a search bar at the top with the text "/hello" and the Brat logo. The list of sentences is as follows:

- 1 Sự khác nhau giữa Ai, Machine learning và deep learning
- 2 Trí tuệ nhân tạo (AI) là tương lai.
- 3 Trí tuệ nhân không chỉ là khoa học viễn tưởng mà còn là một phần của cuộc sống hàng ngày của chúng ta.
- 4 Nó phụ thuộc vào mục tiêu phát triển AI của bạn
- 6 Ví dụ: Khi AlphaGo của Google đánh bại kì thủ cờ vây quốc tế người Hàn Quốc Lee Se-Dol vào năm 2016.
- 7 Thuật ngữ AI, machine learning và deep learning được giới truyền thông sử dụng để mô tả chiến thắng của DeepMind.
- 8 Cả AI, machine learning và deep learning đều góp phần tạo nên chiến thắng của AlphaGo trước kì thủ Se-Dol.
- 9 Nhưng chúng không giống nhau.
- 11 Cách đơn giản nhất để hình dung về mối quan hệ của 3 khái niệm trên là dùng sơ đồ Venn.
- 12 AI – ý tưởng đầu tiên – lớn nhất, sau đó là machine learning, và cuối cùng là deep learning – yếu tố thúc đẩy sự bùng nổ của AI ngày nay .

The entities and relations extracted are as follows:

- Entities: Product (AlphaGo), Location (Hàn Quốc), Person (Lee Se-Dol), Org (Google), Location (năm 2016).
- Relations: Produced (AlphaGo, Google).

Một số bài toán Xử lý ngôn ngữ tự nhiên

Gán nhãn từ loại (Part-of-Speech tagging POS):

Con ruồi đậu mâm xôi đậu

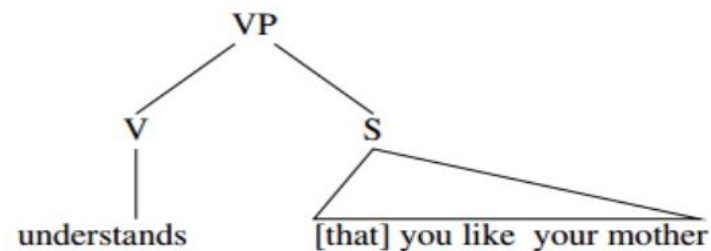
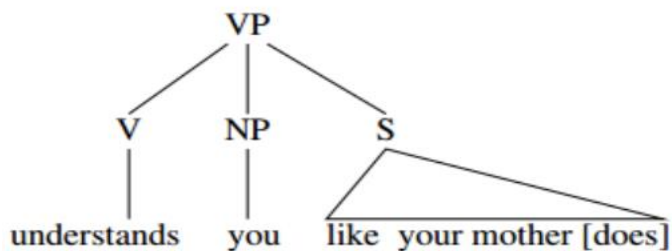
Với : D: determinator (định từ), N: noun (danh từ), V: verb (động từ). Ta có các cặp tương ứng: *Con/D ruồi/N đậu/V mâm/N xôi/N đậu/N*

Bài toán:

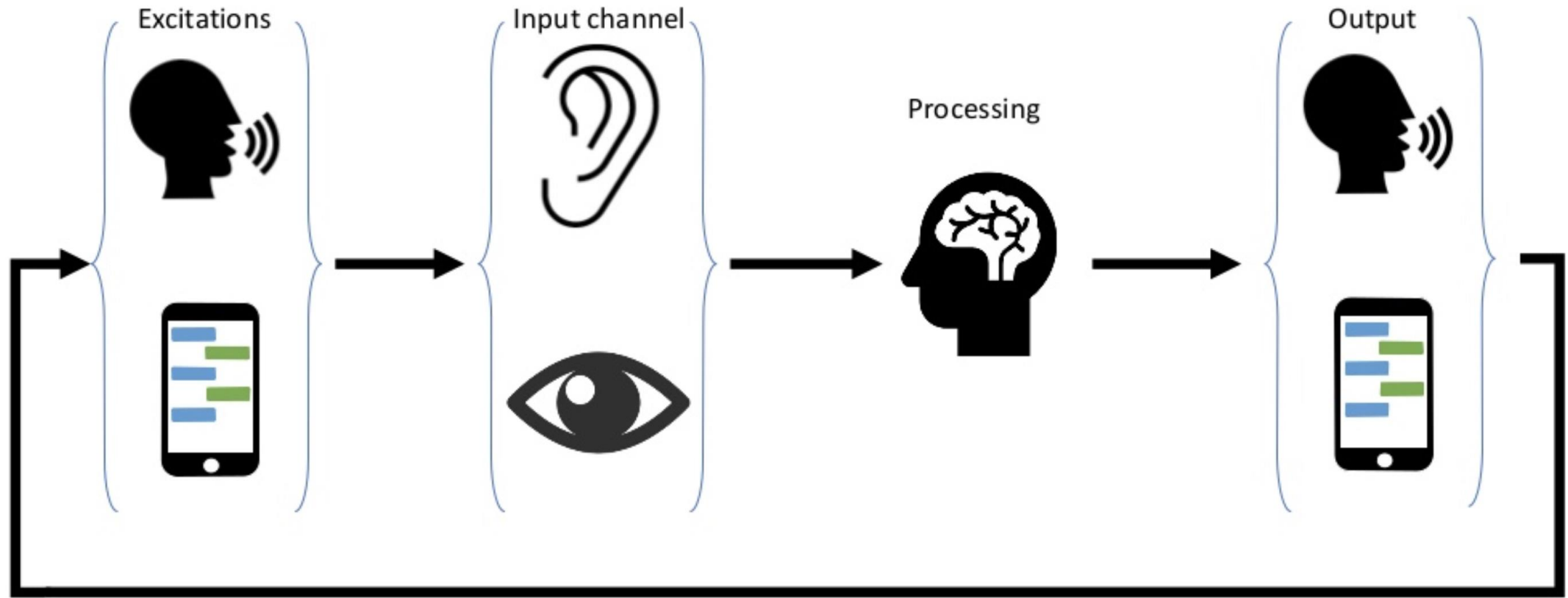
- Named-Entity recognition (gán nhãn tên thực thể).
 - *bà ba [CON NGUOI] bán bánh mì [THUC PHAM] ở phường mười ba [DIA DIEM]*.
- Machine translation (dịch máy)
- Speech recognition (nhận diện tiếng nói).

Những vấn đề trong xử lý ngôn ngữ tự nhiên

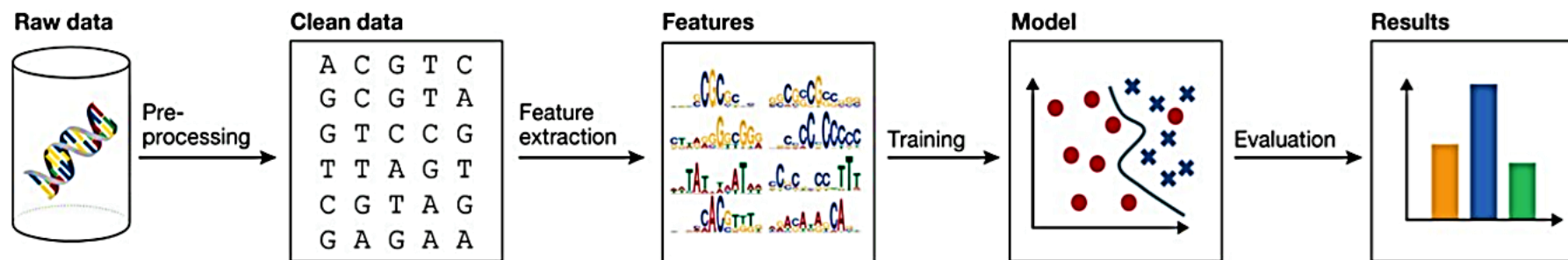
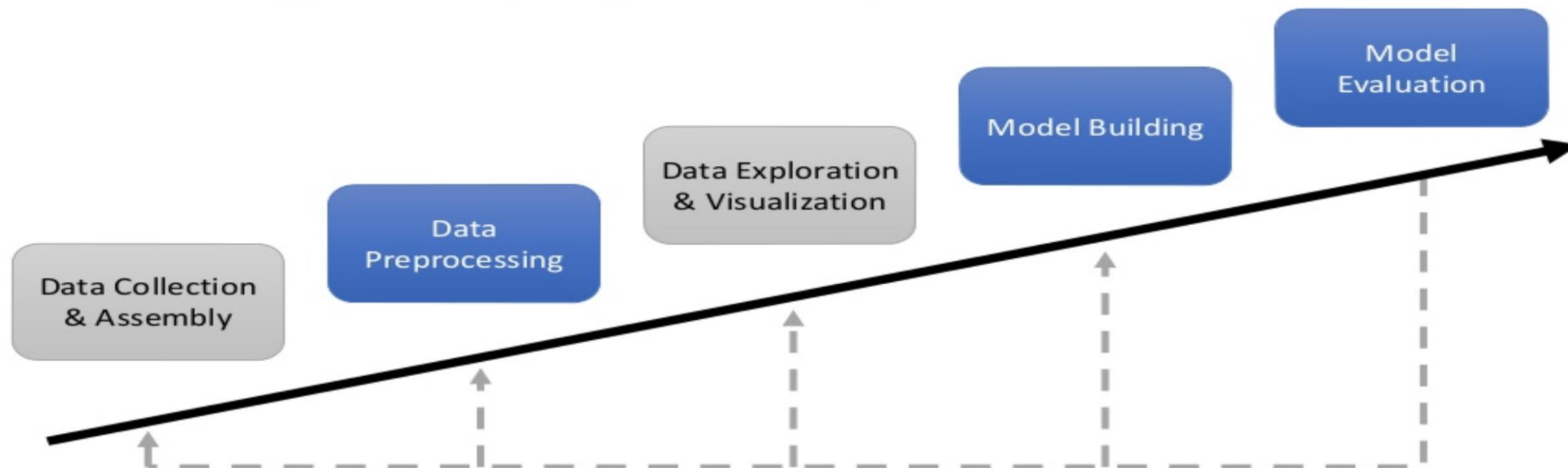
- Sự nhập nhằng trong ngôn ngữ
- Ngôn ngữ tự nhiên sử dụng ngữ cảnh một cách phức tạp và tinh tế để truyền đạt ý nghĩa.
- Ngôn ngữ tự nhiên thường gây nhầm lẫn.
- Ngôn ngữ tự nhiên liên quan tới suy luận về thế giới.
- Ngôn ngữ tự nhiên là một phần quan trọng trong việc tương tác giữa con người với nhau (một hệ thống mang tính xã hội).
- Ví dụ:
 - Ông già đi nhanh quá
 - They book that hotel. They read that book.
 - A computer understands you like your mother



Bài toán Xử Lý Ngôn Ngữ Tự Nhiên

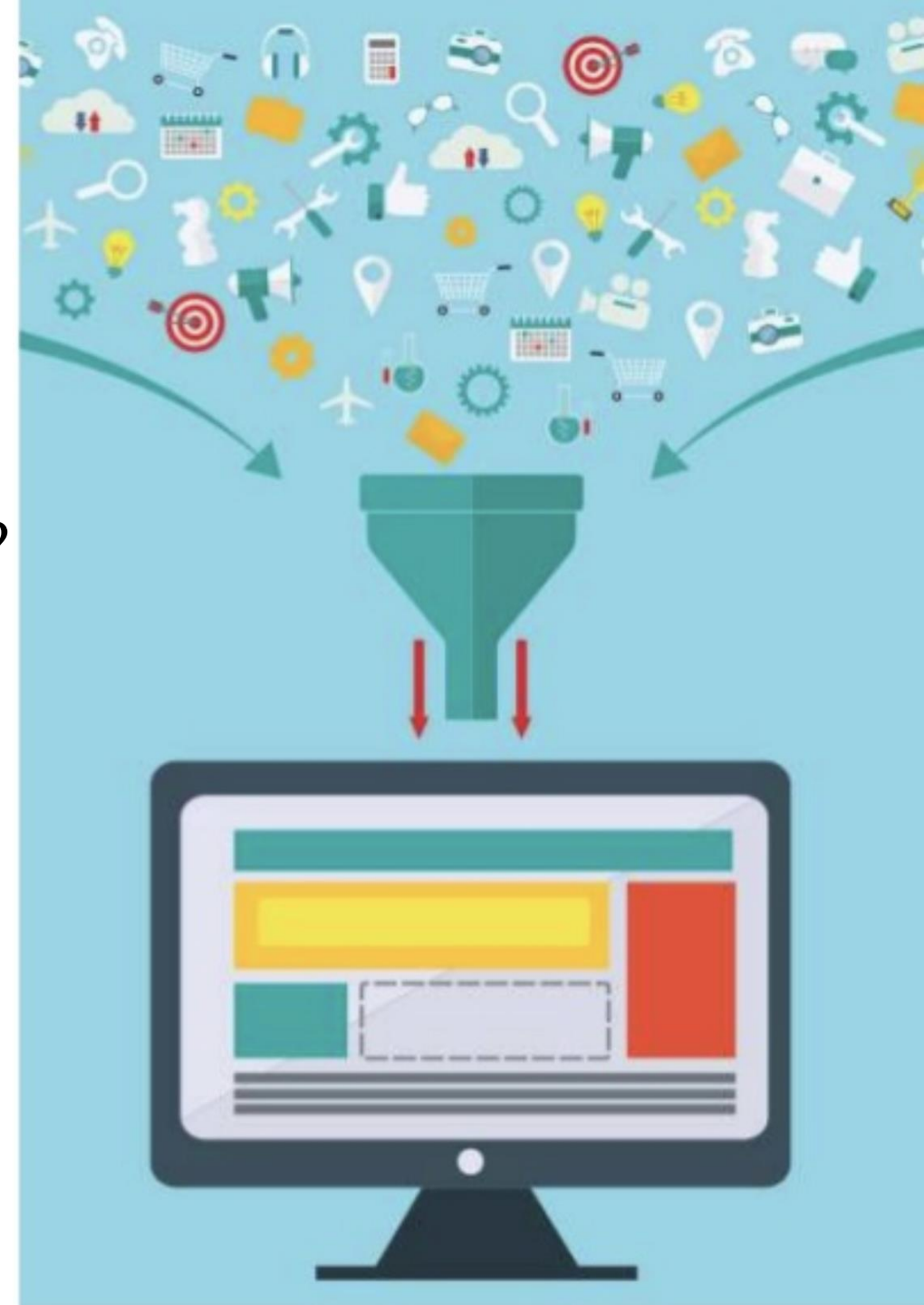


Bài toán Xử Lý Ngôn Ngữ Tự Nhiên

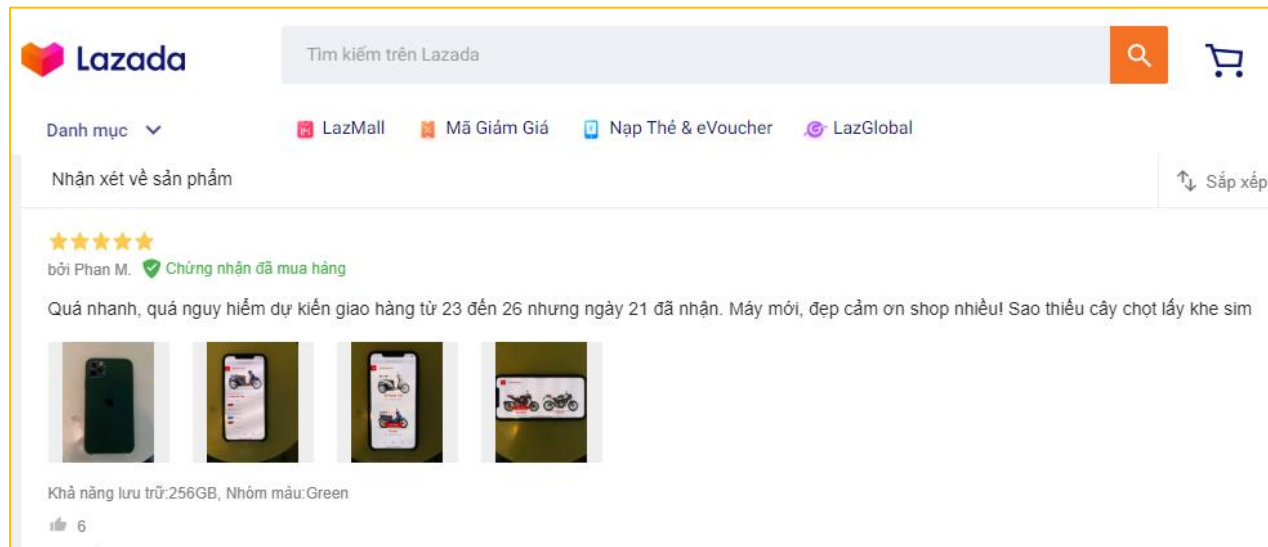


Thu thập dữ liệu

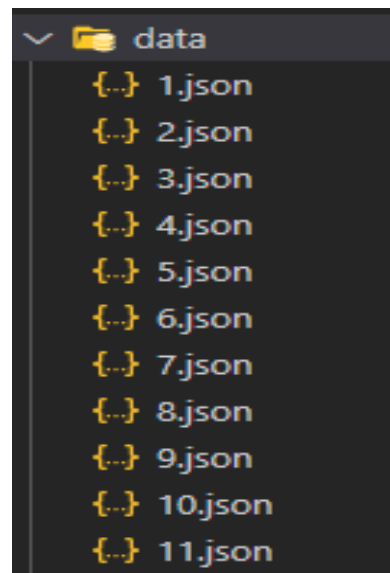
- Thu thập từ nhiều nguồn khác nhau
 - Nguồn dữ liệu tin cậy
 - Dữ liệu nào thực sự quan trọng cho mô hình??
- => Quá trình tốn nhiều thời gian, công sức



Thu thập dữ liệu



Name	Sta...	Ty...	Initiator	Size	Time
Lazada_PDP....	200	gif	VM23:6	97 B	41 ms
getReviewLis...	200	xhr	index.j...	5.1...	84 ms
getReviewLis...	200		Other	0 B	55 ms



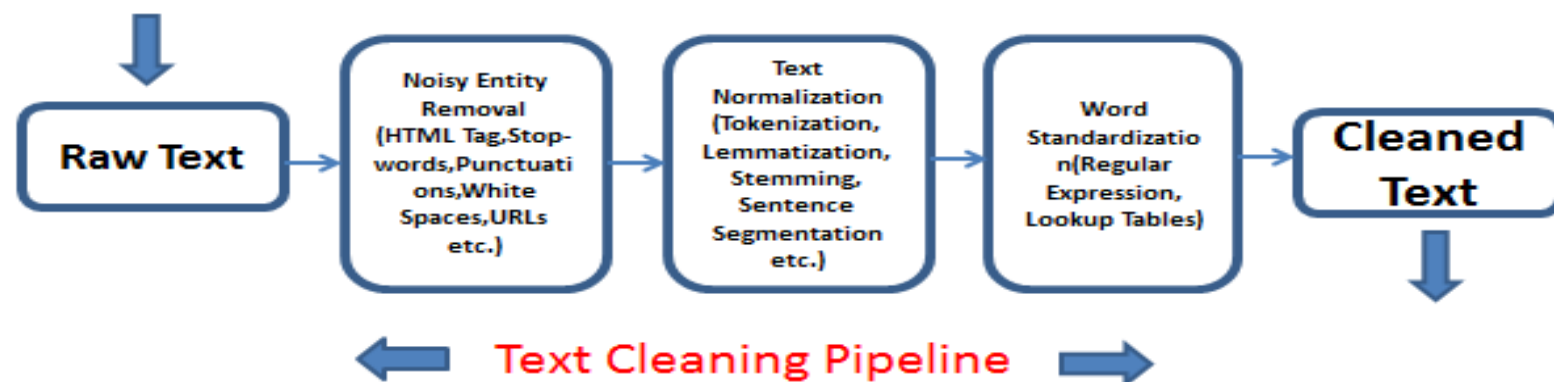
1	bt tai nghe một bên nghe được một bên không bình thường nghe tạm ổn	
2	xịn sò ghê	
1	hàng mẫu ma đẹp tam on	
0	Ko ưng ý lắm, nhiều luk sd hay bị đơ	
1	bphone sản xuất tại VN thì chất lượng như z là ổn	
0	hàng tệ thẻ nhớ chập chon thừa các cửa hàng khác	

Tiền xử lý văn bản

Dữ liệu lớn, đa dạng, nhiều noise=> xử lý dữ liệu thô ban đầu

Ví dụ :

- Xử lý các ký tự lặp lại : "*gooooood*" -> "*good*"
- Xử lý các ký tự cho đồng nhất: "*\$stupid*" -> "*stupid*"
- Xử lý các đầu vào đặc biệt: "*http://www.foo.com/bar*" -> "*[URL]*"
- Chuẩn hoá Unicode
- Chuyển chữ hoa thành chữ thường
- => **Bước làm sạch dữ liệu**



Tiền xử lý dữ liệu

- Tiền xử lý
- Gán nhãn dữ liệu

0, Cua toi ko co đây sặc cuc pin du phong 0, Sao cục pin của mình sặc 3 tiếng mà còn chưa đầy 1 chấm xanh

1, Giao hàng rất chậm. Kh có cáp đi kèm

2, Wá tuyệt zờiiiiiiii

2, Rất hai long

1, Sp lỗi mỗ nguồn khôg lên Xac vào đt khôg vào Yc đỏi cho mình
sp khác

2, Không có dây sạc chỉ có sạc cục với bao đựng pin dự phòng



chung

Tệp Chính sửa Xem Chèn Định dạng Dữ liệu Công cụ Tiện ích bổ sung Trợ giúp [Chỉnh sửa lần cuối cách đây 1 giờ](#)

 100% đ % .0 .00 123 Arial 10 B I U A

	A	B	C	D
1		1 bt tai nghe một bên nghe được một bên không bình thường nghe tạm ổn		
2		2 xịn số ghê		
3		1 hàng mẫu ma đẹp tam on		
4		0 Ko ưng ý lắm, nhiều luk sd hay bị đơ		
5		1 bphone san xuất tại VN thì chất lượng như z là ổn		
6		0 hàng tệ thể nhỏ chấp chon thừa các cửa hàng khác		
7		0 tệ sao bếp không có tnh có gì biết gọi đâu		
8		0 Đã mua hàng Nhưng không kết nối được với điện thoại Đề nghị được hỗ trợ		
9		0 táo cạn kiệt ý tưởng không qua nổi 5s thần thánh thất vọng		
10		0 Ngon trong tầm giá nhưng có đáng mua Điện thoại về ngày 6 12 nguyên seal tất cả đều ổn nhưng không thể chấp nhận được khi điện thoại lúc nhận sim 67 12 và không nhận sim 8 12 và máy ảnh v		
11		1 bình thường đồng gói cần thận giao đúng sản phẩm nhưng mau hết pin quá chất lượng với giá tiền như vậy cũng tạm ổn tạm ổn		
12		2 Bàn phím và chuột nhạy sử dụng tốt về độ bền thì phải dùng thêm thời gian với mức giá này thì combo này khá rẻ		
13		0 hàng tệ thể nhớ chấp chơn hàng thừa các cửa hàng khác		
14		0 đã đặt mua..nhưng không may 2 tai kết nối thông được..minh đã đổi trả hàng		
15		1 khá bình thường chất lượng âm thanh khá bình thường sennheiser còn có mẫu mã hộp rất đẹp mắt tai nghe sử dụng khá bền nghe rất rõ bình thường lẫn tạp âm mỗi khi nghe sennheiser có nhiều lo		
16		0 máy để giải trí mà giá này thì chắc là doanh số không cao rồi d		
17		1 Đồng Nai Về phía bản thân mình sử dụng thứ nhất mình cảm thấy sóng gió yếu mình sử dụng sim data vinaphone pin hực ngay khi không sử dụng 1 ngày là hết sạch pin Không kết nối wifi với TP Lir		
18		2 mình cảm thù nokia vì mình mua con 1200 từ 2007 đến giờ vẫn tốt chịu để thay cái khác		
19		0 nói chung là hên xui tcs hiện tại theo cảm nhận của mình chỉ là thời gian hỗ trợ hơi lâu nhưng một khi đã liên lạc được thì hỗ trợ tận tình thái độ tốt còn thời gian giao hàng hàng tcs gửi hàng tương đ		
20		2 sản phẩm chất lượng về chất lượng sản phẩm tai nghe ổn âm thanh to rõ nghe đều cả hai bên và chất lượng xứng đáng với tầm giá nhạc nền hơi lớn hơn giọng ca sĩ nhưng cũng không đáng kể bù		
21		2 điểm 10 cho chất lượng		
22		1 cũ nhả vùng		

Thay đổi **lượng** thành: ✕

Thay đổi ▼

Bỏ qua ▼

Thêm vào từ điển ▼

Một số dataset

- Tiếng Anh:
- <https://datasets.quantumstat.com/>
- <https://www.kaggle.com/datasets?search=text+classification>

Hebrew Parallel Movie Subtitles	06.25.20	Hebrew	Dataset comes from subtitles of movies and television shows for the purpose of semantic role labeling in Hebrew. It includes both FrameNet and PropBank annotations.	30,789	n/a	Semantic Role Labeling	2020	Eyal et al.	LINK PAPER
MEDIQA-Answer Summarization	06.25.20	English	Dataset containing question-driven summaries of answers to consumer health questions.	156	JSON	Summarization	2020	Savery et al.	LINK PAPER
NEJM-enzh	06.25.20	Chinese, English	Dataset is an English-Chinese parallel corpus, consisting of about 100,000 sentence pairs and 3,000,000 tokens on each side, from the New England Journal of Medicine (NEJM).	100,000	n/a	Machine Translation	2020	Liu et al.	LINK PAPER
Wikipedia Current Events Portal (WCEP) Dataset	06.25.20	English	Dataset is used for multi-document summarization (MDS) and consists of short, human-written summaries about news events, obtained from the Wikipedia Current Events Portal (WCEP), each paired with a cluster of news articles associated with an event.	10,200	JSON	Summarization	2020	Ghalandari et al.	LINK PAPER
Worldtree Corpus	06.25.20	English	Dataset contains multi-hop question answering/explanations where questions require combining between 1 and 16 facts (average 6) to generate detailed explanations for question answering inference. Each explanation is represented as a lexically-connected "explanation graph" that combines an average of 6 facts drawn from a semi-structured knowledge base of 9,216 facts across 66 tables.	5,114	Text, TSV	Question Answering, Knowledge Base	2020	Xie et al.	LINK PAPER
ScienceExamCER	06.25.20	English	Dataset contains 133k mentions in the science exam domain where nearly all (96%) of content words have been annotated with one or more fine-grained semantic class labels including taxonomic groups, meronym groups, verb/action groups, properties and values, and synonyms.	133,000	Text, TSV	Named Entity Recognition (NER)	2019	Smith et al.	LINK PAPER
RuBQ	06.25.20	Russian	Dataset consists of 1,500 Russian questions of varying complexity, their English machine translations, SPARQL queries to Wikidata, reference answers, as well as a Wikidata sample of triples containing entities with Russian labels.	1,500	JSON	Question Answering, Knowledge Base	2020	Korablinov et al.	LINK PAPER

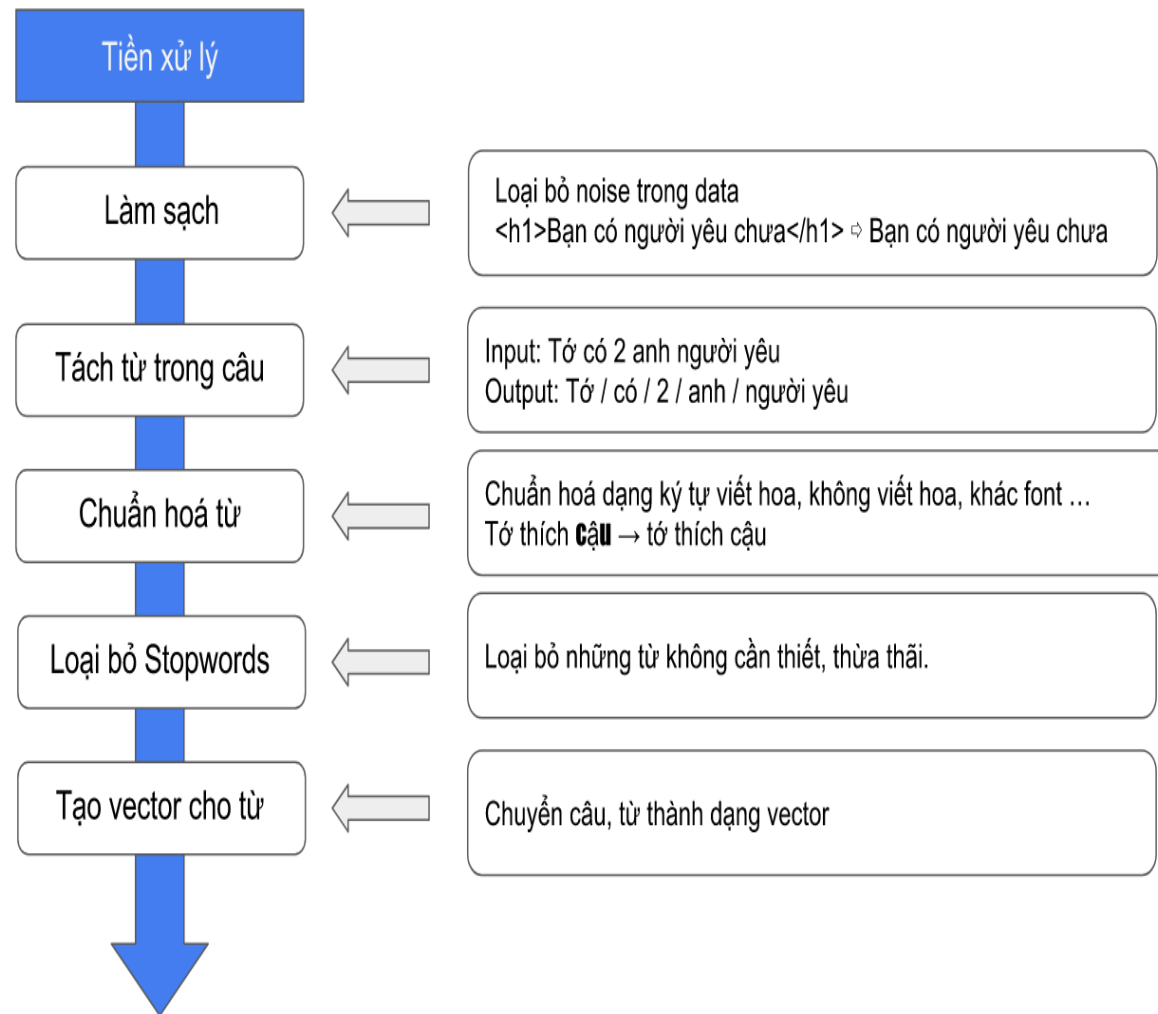
Tiếng Việt:

- Câu lạc bộ xử lý ngôn ngữ và tiếng nói Việt
- <https://www.vlsp.org.vn/resources>
- Các thư viện: underthesea,...
- Tự crawl...

Tiền xử lý văn bản

Tiếng Việt:

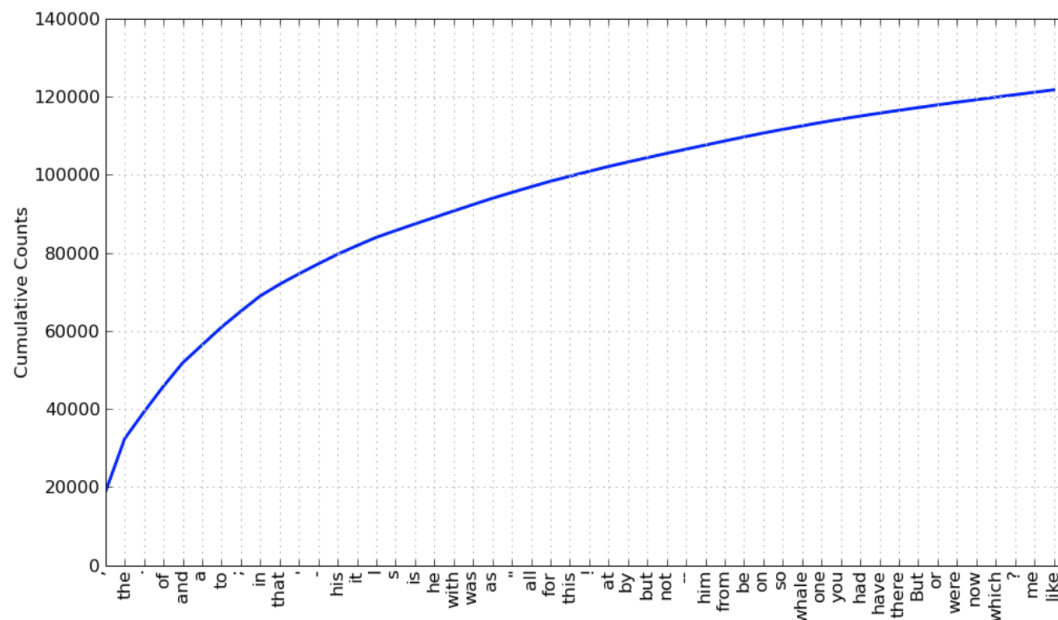
- Xóa HTML code (nếu có)
- Chuẩn hóa bảng mã Unicode (đưa về Unicode tổ hợp dựng sẵn)
- Chuẩn hóa kiểu gõ dấu tiếng Việt (dùng *òà úý* thay cho *oà uý*)
- Thực hiện tách từ tiếng Việt (sử dụng thư viện tách từ như pyvi, underthesea, vncorenlp,...)
- đưa về văn bản lower (viết thường)
- Xóa các ký tự đặc biệt: “.”, “,”, “;”, “)”, ...



StopWords

- **StopWords** là những từ xuất hiện nhiều trong ngôn ngữ tự nhiên, tuy nhiên lại không mang nhiều ý nghĩa.
 - Tiếng Việt: để, này, kia... .
 - Tiếng Anh: is, that, this...
- Có rất nhiều cách để loại bỏ StopWords nhưng có 2 cách chính là:
 - Dùng từ điển
 - Dựa theo tần suất xuất hiện của từ

cậu
của
cứ
dù
nọ
phóc
này
kia
để
...



Stopwords

Tiếng Anh:

```
> stopwords("english")
[1] "i"      "me"      "my"      "myself"  "we"
[6] "our"    "ours"    "ourselves" "you"     "your"
[11] "yours"  "yourself" "yourselves" "he"      "him"
[16] "his"    "himself" "she"      "her"     "hers"
[21] "herself" "it"      "its"      "itself"  "they"
[26] "them"   "their"   "theirs"   "themselves" "what"
[31] "which"  "who"     "whom"    "this"    "that"
[36] "these"  "those"   "am"      "is"      "are"
[41] "was"    "were"    "be"      "been"    "being"
[46] "have"   "has"     "had"     "having"  "do"
```

Tiếng Việt:

```
['bị', 'bởi', 'cà', 'các', 'cái', 'cần', 'càng', 'chi', 'chiếc', 'cho', 'chứ', 'chưa', 'chuyện', 'có', 'có_thể', 'cứ', 'của',  
'cùng', 'cũng', 'đã', 'đang', 'đây', 'để', 'đến_nỗi', 'đều', 'điều', 'do', 'đó', 'được', 'dưới', 'gì', 'khi', 'không', 'là', 'l  
ại', 'lên', 'lúc', 'mà', 'mỗi', 'một_cách', 'này', 'nên', 'nếu', 'ngay', 'nhiều', 'như', 'nhưng', 'những', 'nơi', 'nữa', 'phà  
i', 'qua', 'ra', 'rằng', 'rằng', 'rất', 'rất', 'rồi', 'sau', 'sẽ', 'so', 'sự', 'tại', 'theo', 'thì', 'trên', 'trước', 'từ', 'từ  
ng', 'và', 'vẫn', 'vào', 'vậy', 'vì', 'việc', 'với', 'vừa', '!', '"', '#', '$', '%', '&', "'", '(', ')', '*', '+', ',', '-',  
'.', '/', ':', ';', '<', '=', '>', '?', '@', '[', '\\', ']', '^', '_', '`', '{', '|', '}', '~']
```


Word representation

One hot vector

- Xây dựng một bộ từ vựng.
- Mỗi vector đại diện cho một từ có số chiều bằng số từ trong bộ từ vựng.
- Trong đó, mỗi vector chỉ có một phần tử duy nhất khác 0 (bằng 1) tại vị trí tương ứng với vị trí từ đó trong bộ từ vựng.

Vocabulary

index:	Word:
0	aardvark
1	able
...	...
2409	black
2410	bling
...	...
3202	candid
3203	cast
3204	cat
...	...
5281	is
5282	island
...	...
8676	the
8677	thing
...	...
9999	zombie

the

cat

is

black

Vocabulary

index:	Word:
0	aardvark
1	able
...	...
2409	black
2410	bling
...	...
3202	candid
3203	cast
3204	cat
...	...
5281	is
5282	island
...	...
8676	the
8677	thing
...	...
9999	zombie

the

cat

is

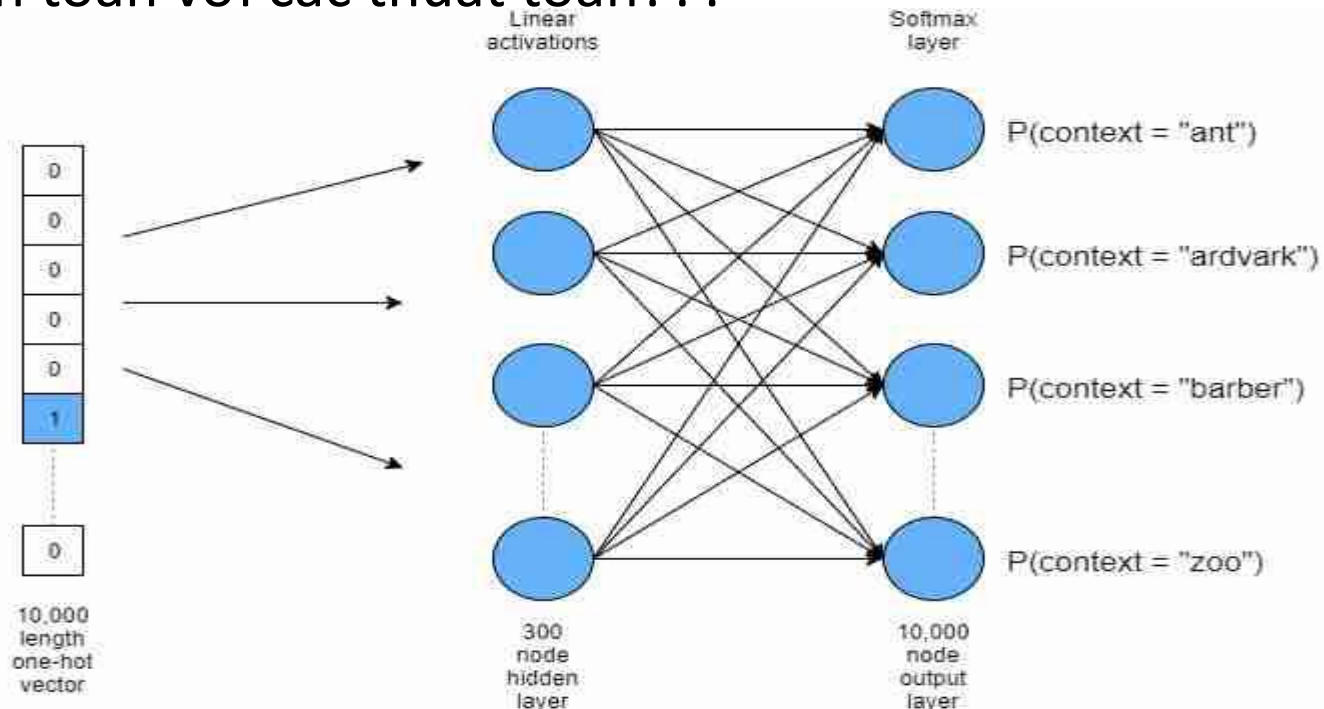
black

1 is at
index:
8676

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 0 \\ 1 \\ 0 \\ \dots \end{pmatrix} \begin{matrix} 0 \\ 1 \\ 2 \\ \dots \\ 8675 \\ 8676 \\ 8677 \\ \dots \end{matrix}$$

One hot vector

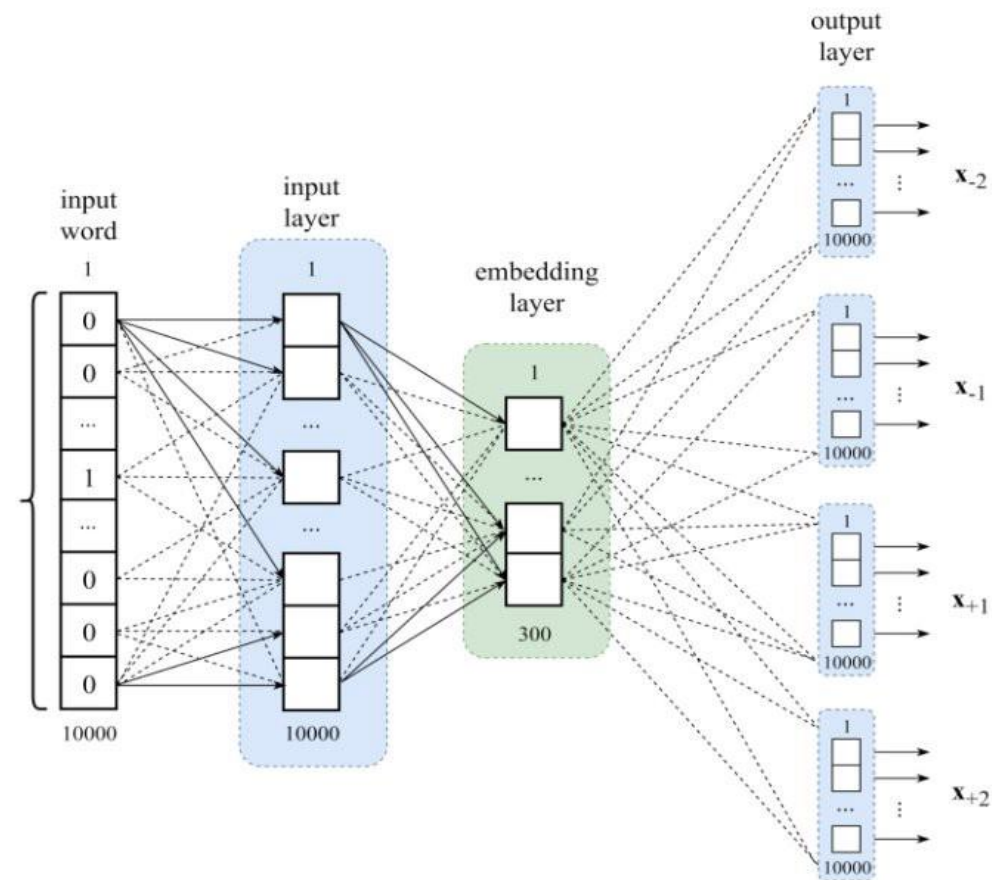
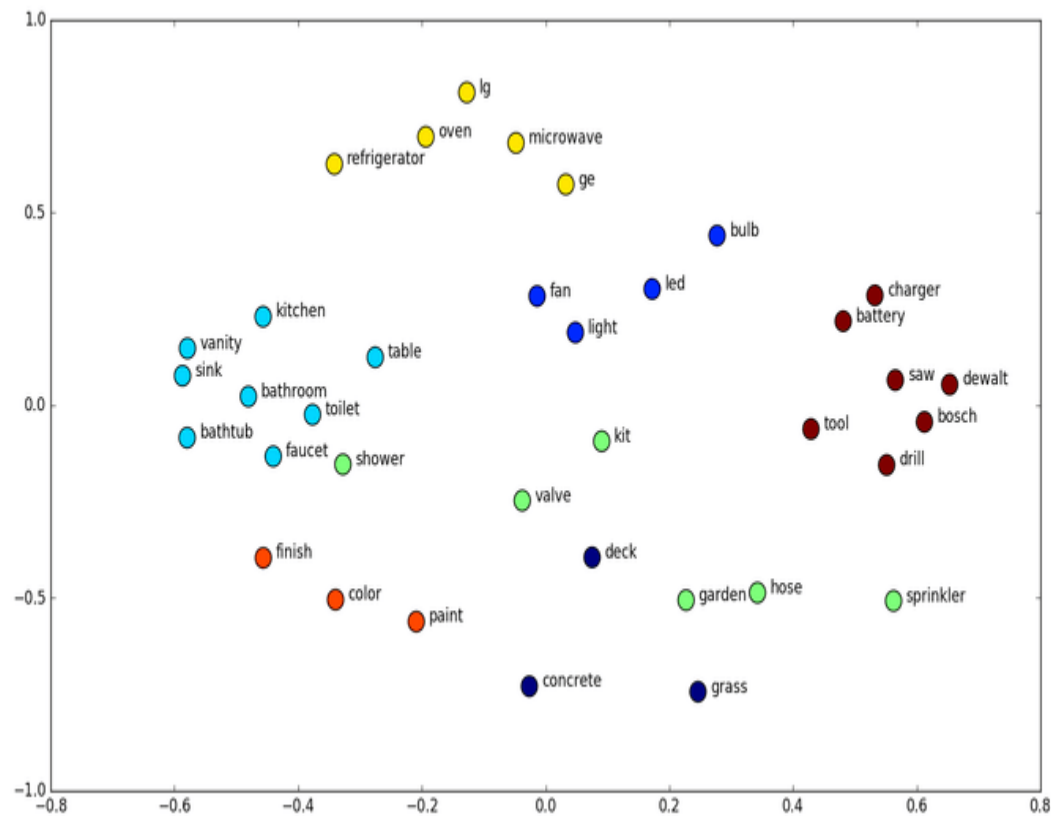
- Hạn chế
 - độ dài của vector là quá lớn
 - Không xác định được sự tương quan ý nghĩa giữa các từ
 - Tính toán với các thuật toán???



Phương pháp khắc phục

- Phương pháp truyền thống:
 - Bag of Words, IF-IDF
 - Matrix-Factorization: Giảm chiều vector SVD, PCA
- Word Embedding:
 - Word2Vec: Skip Gram, Continuous Bag of Words (CBOW).
 - GloVe
 - FastText
- ..

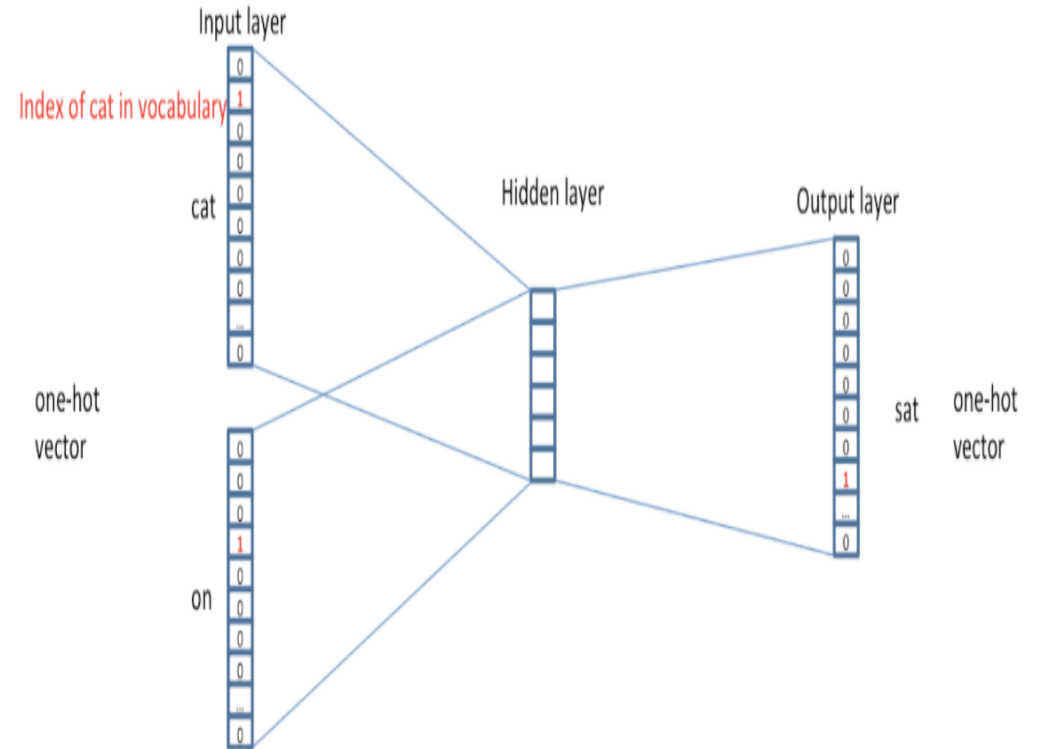
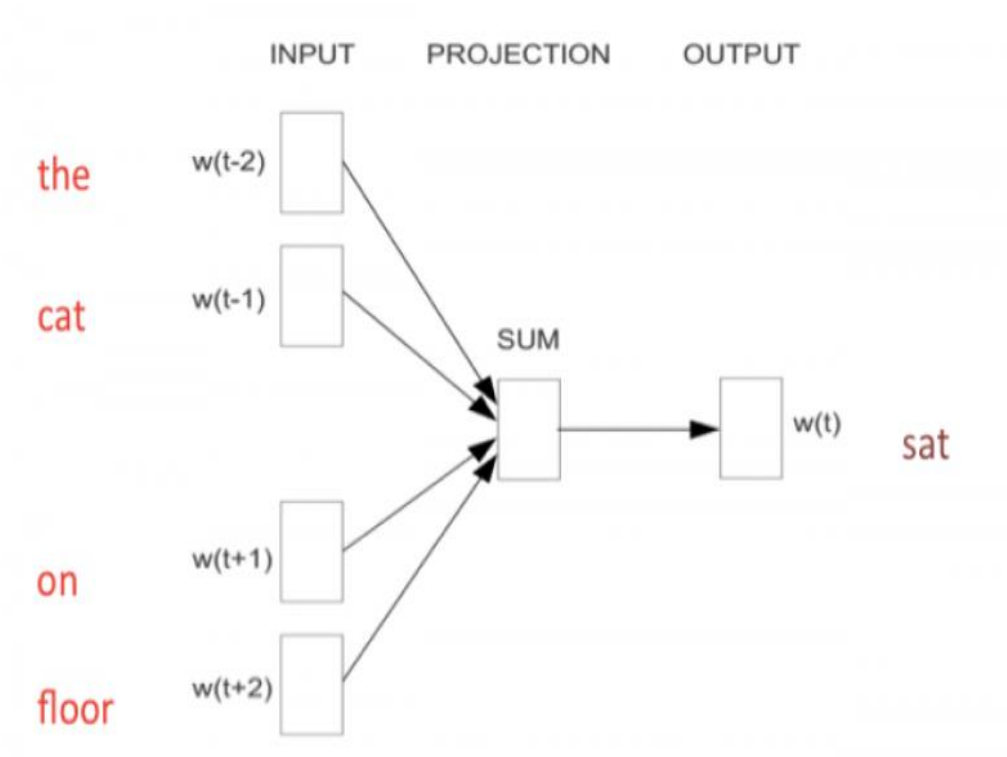
Word Embedding



Word Embedding: Word2Vec

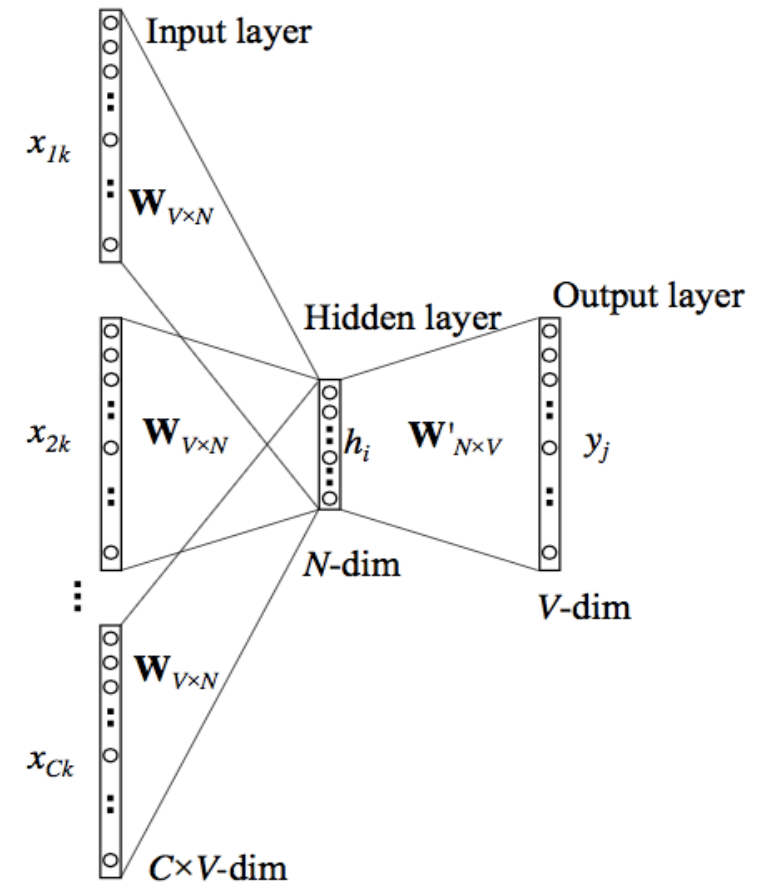
Continuous Bag of Words(CBOW)

Ví dụ: *The cat sat on floor*



Continuous Bag of Words(CBOW)

- Với một corpus thật lớn, mạng CBOW học lần lượt từng từ trong corpus với context tương ứng
- 4 vector one-hot của bốn từ trong context sẽ được cộng lại thành một vector input duy nhất (có 4 giá trị 1) và đưa vào hệ thống như một input duy nhất.
- Với ma trận encoder $W_{V \times N}$, ma trận decoder $W'_{N \times V}$, x là one-hot vector của một từ, embedding vector tương ứng $w = x.W$.



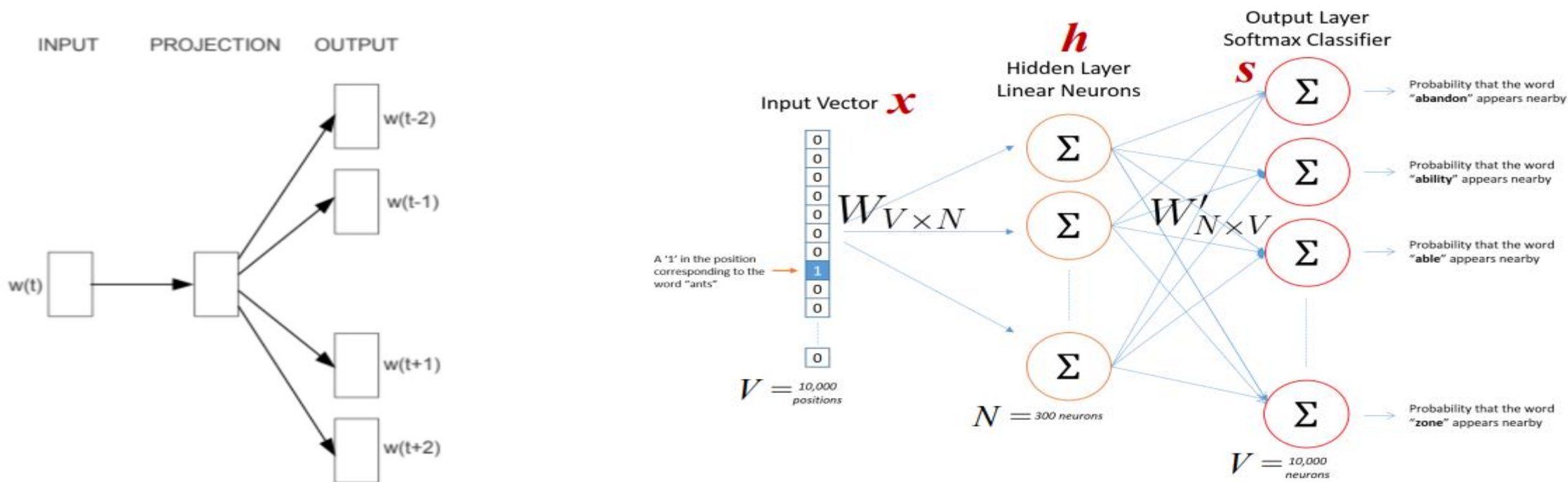
Continuous Bag of Words(CBOW)

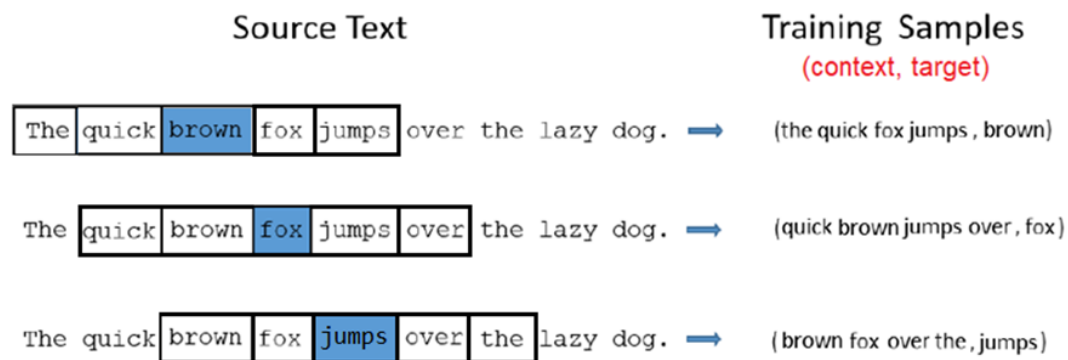
- embedding vector của một từ sẽ có các tính chất:
- Số chiều của w sẽ nhỏ hơn nhiều so với số chiều gốc của one-hot vector x .
- Một one-hot vector sẽ được encode không phải dựa vào chính nó như AutoEncoder mà dựa vào các từ *thường hay xuất hiện quanh nó* trong các văn bản.

=> từ thường hay xuất hiện cạnh nhau trong các văn bản sẽ được encode thành các vector tương tự nhau.

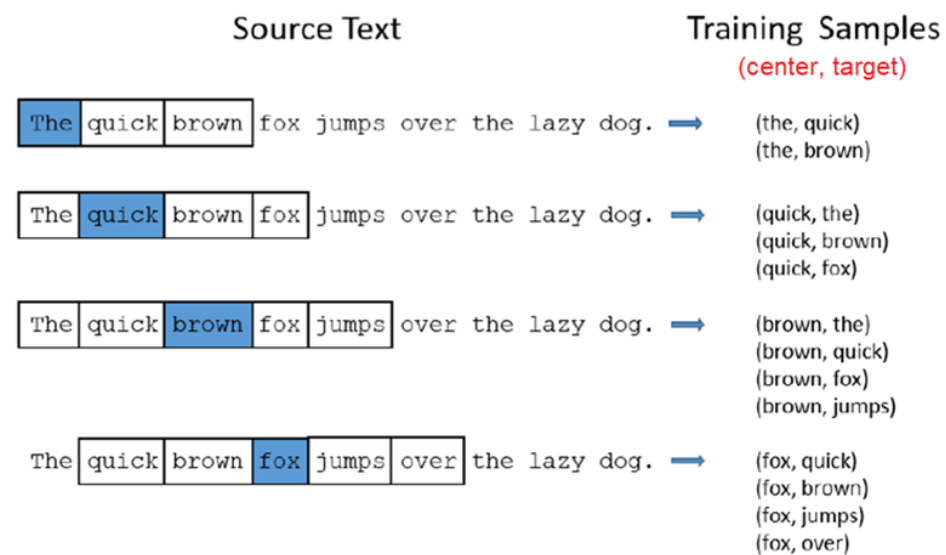
Mô hình Skipgram

- Mô hình Skipgram gần giống như mô hình CBOW, chỉ thay đổi vai trò:
- context được dùng làm output và từ trung tâm sẽ dùng làm input => 1 ma trận encoder W và 4 ma trận decoder W' .

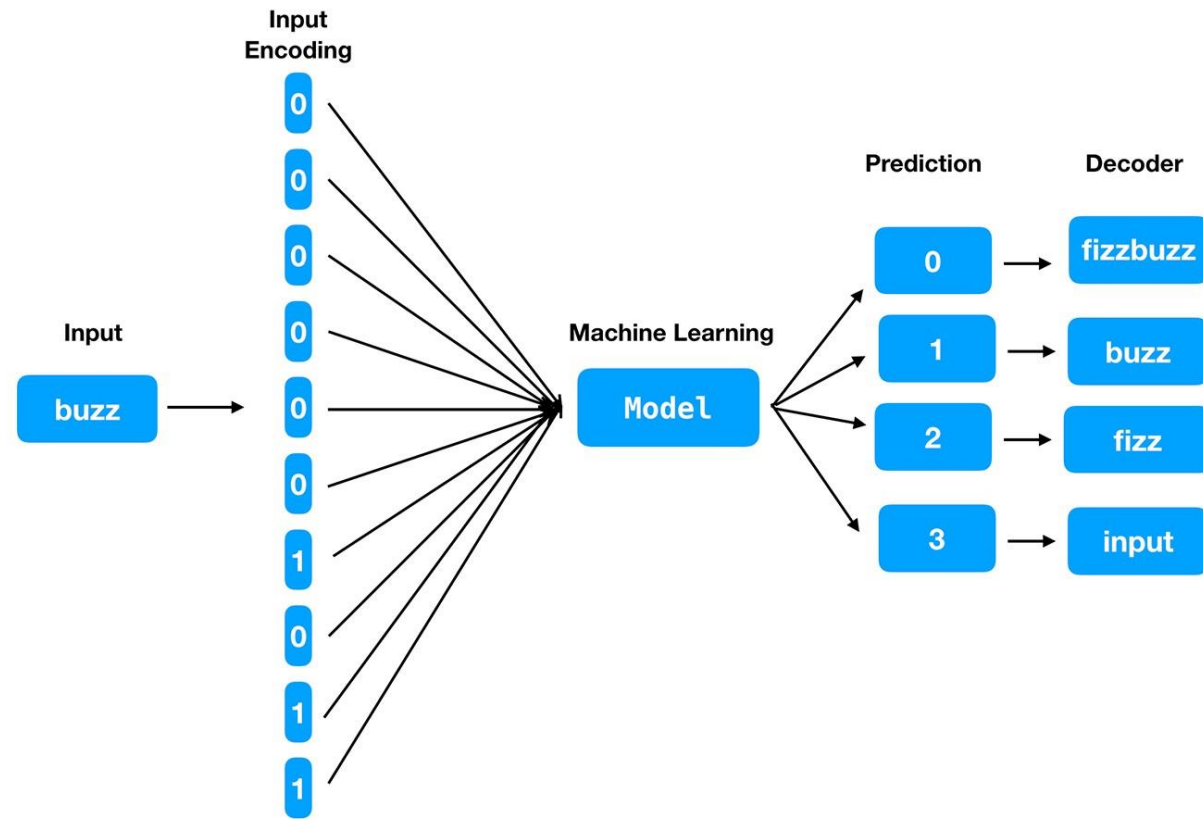




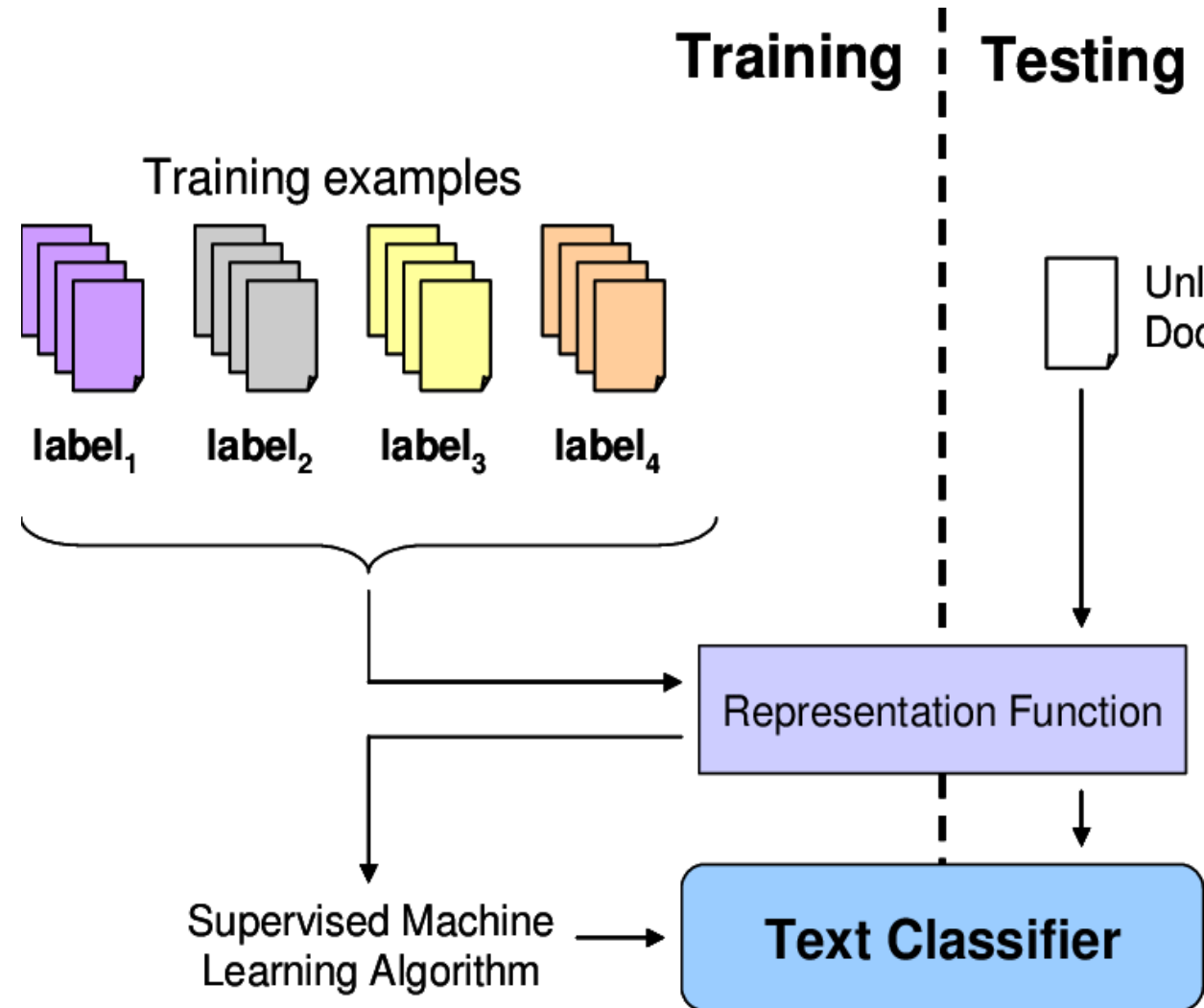
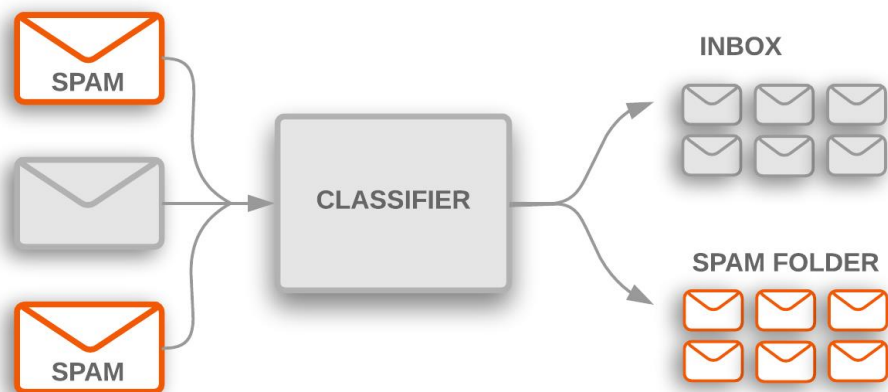
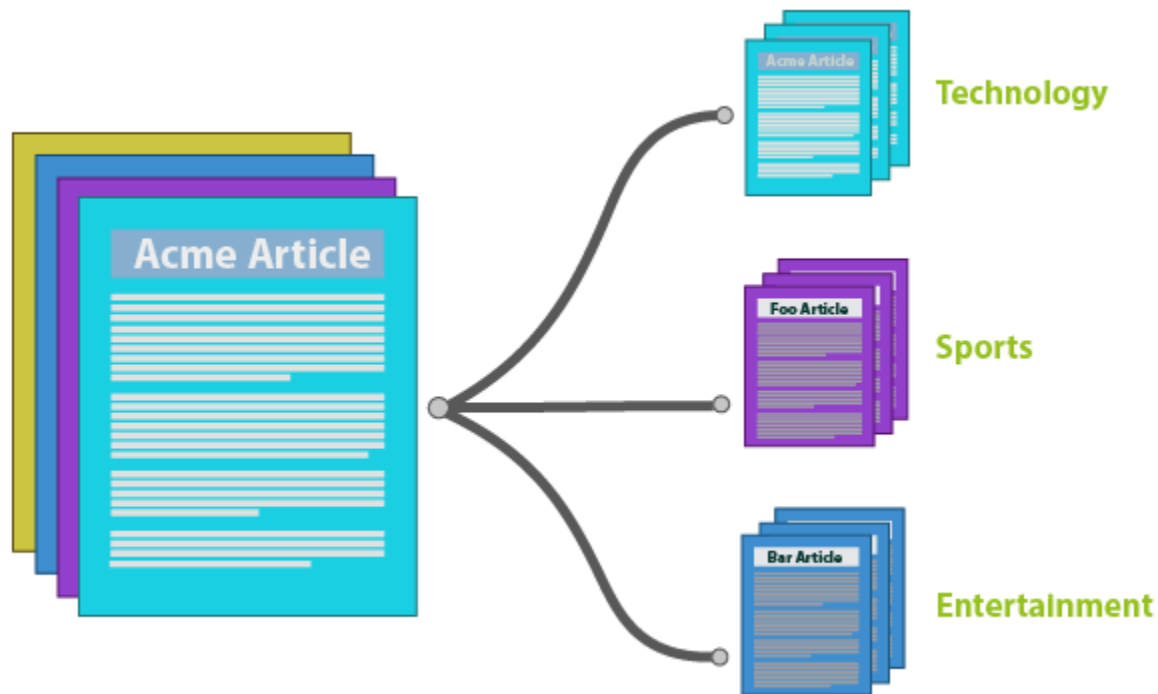
CBOW



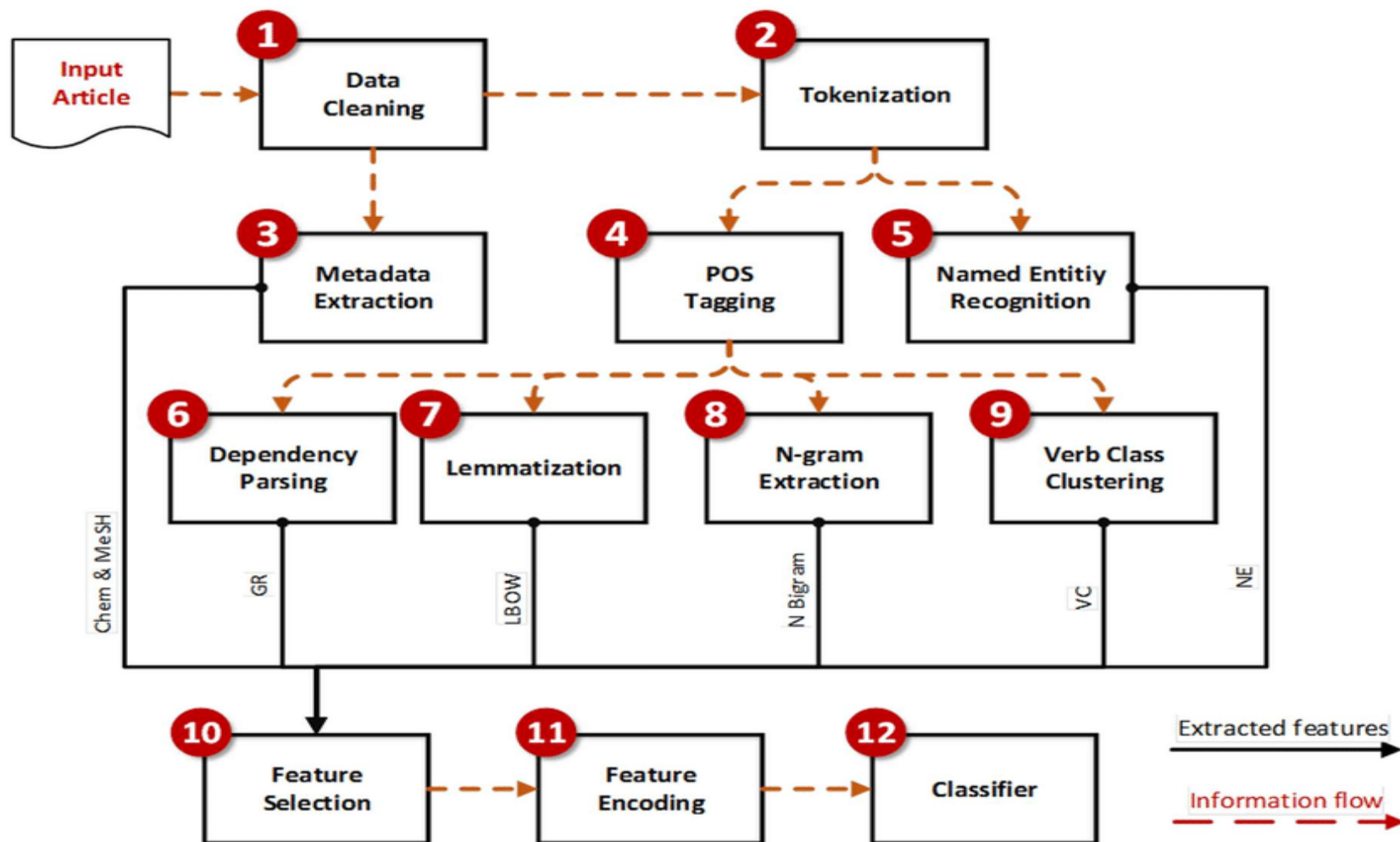
Skipgram



Bài toán phân loại văn bản



Phương pháp cổ điển



Deep Learning: CNN

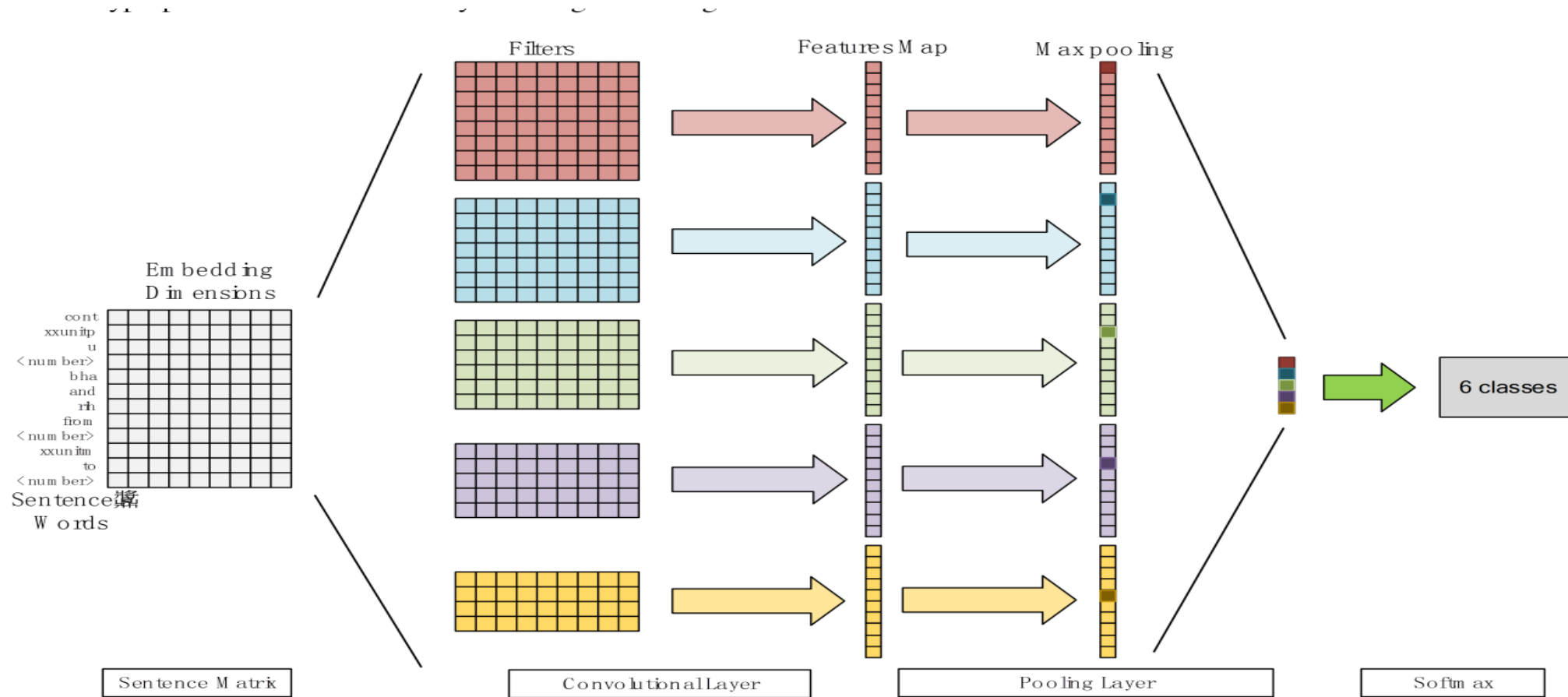
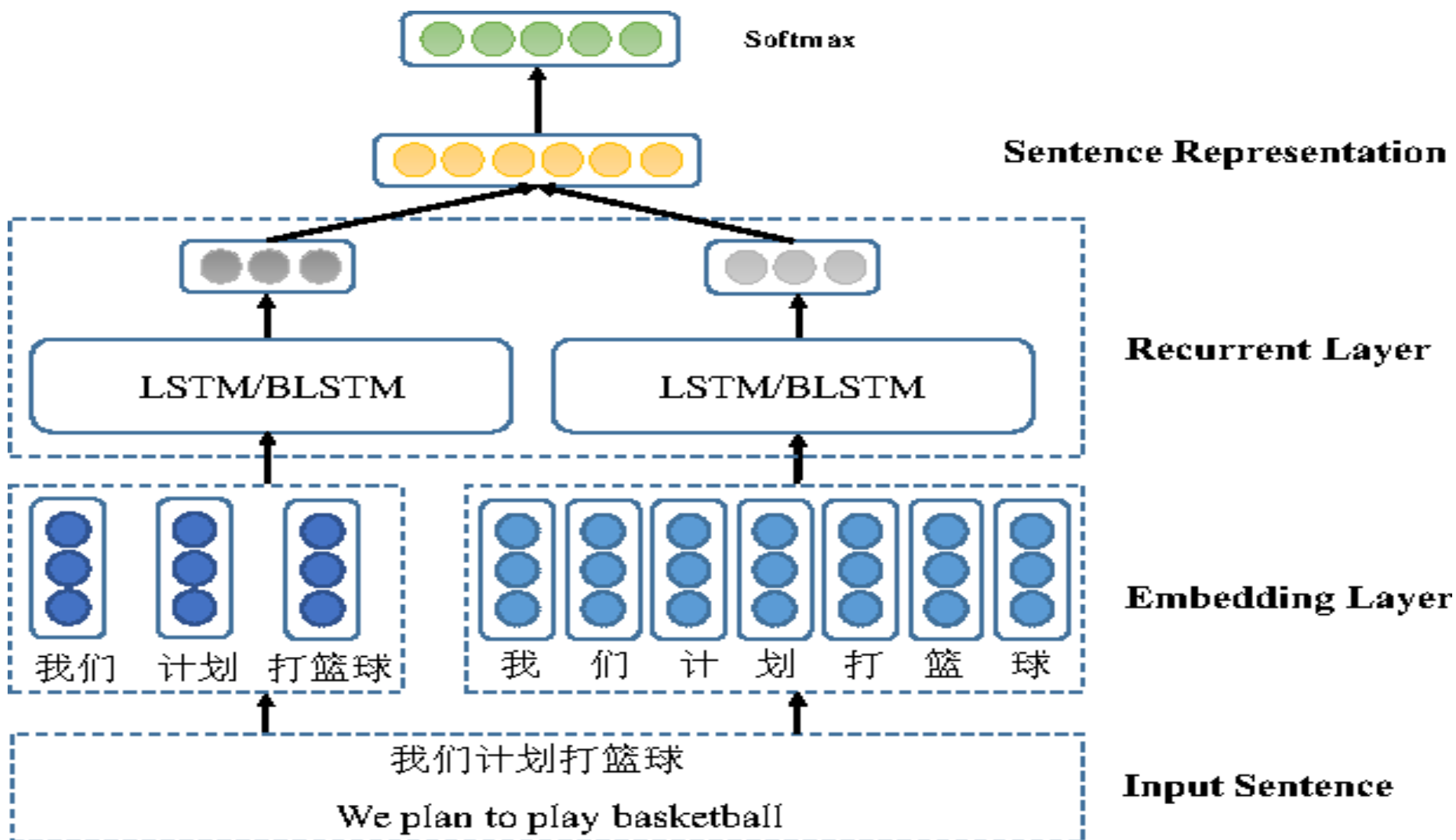
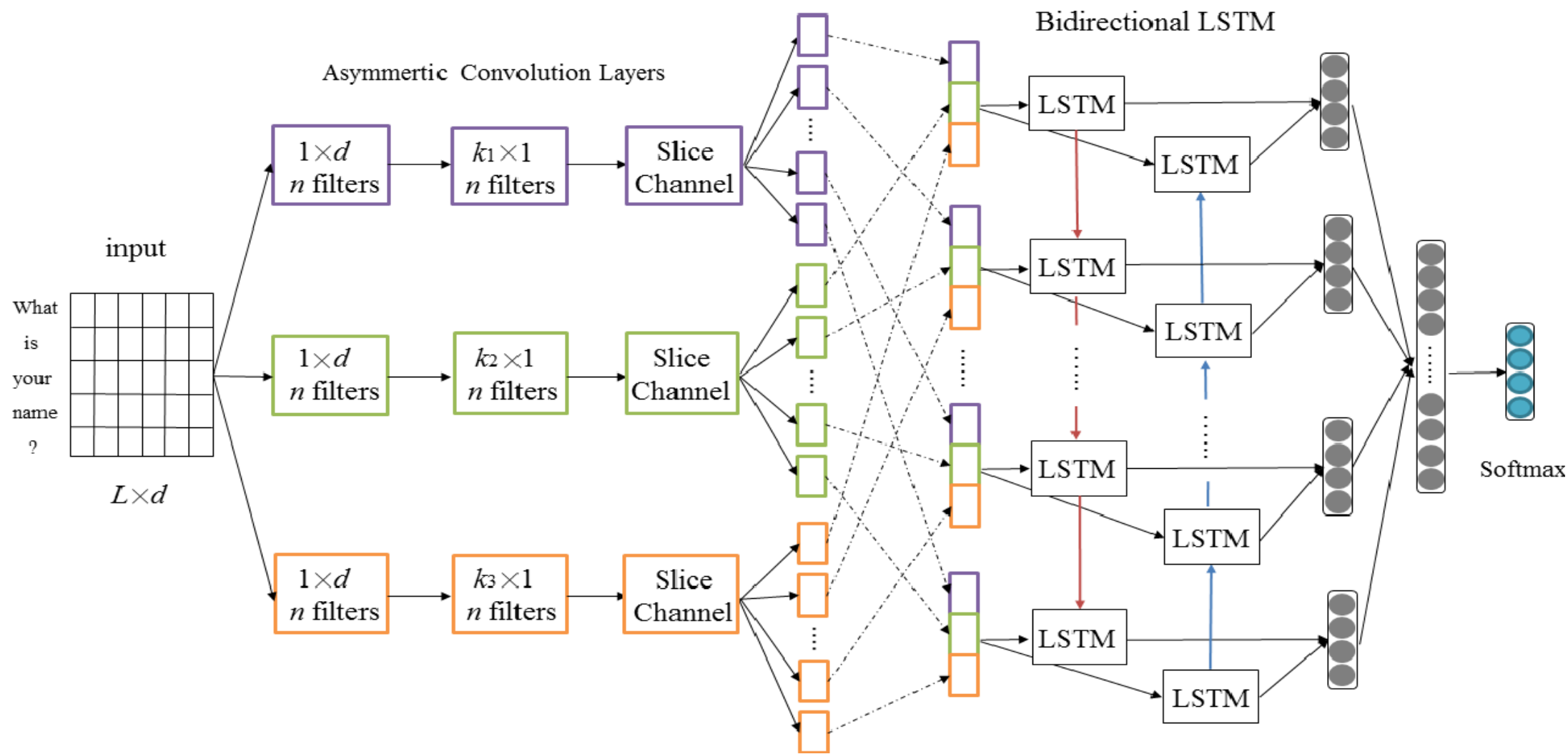


Figure 2. CNN-based text classification network.

Deep Learning: RNN



Deep Learning: CNN+RNN



Implementation of LSTM for Sentimental Analysis

Dữ liệu

	text	stars	sentiment
1411425	i went here and ordered the barbiq burger and i got the meat in medium well i had to admit it was really good although the fries got in the way it was salty and i had to dip the fries with mayo and ketchup all and all i would give the fries 35 out of 5 the service of the burger was right on time i like the look of the restaurant with the fun and relaxed look of the burger joint worth the visit especially if you just want to try a good burger	3	neg
971153	after calling 3 difference companies number one were the only one that had the least wait time for our broken ac in the middle of summer we called on monday morning and they were able to send a technician out that night at 8 pm when the technician came he determined it was our blower motor that was the problem and told us they will have to order the parts for it it will take couple days the repair will probably get done on wednesday and they will call us on tuesday to let us know the time they told us the repair will be between 11 am 2 pm on wednesday the repair guy pulled up in uhaul and after going up to the attic to check out the blower told us he needed to go get more parts for it and will be back in the afternoon\n\ncome 5 pm still no sign of the guy so we called them back and they had to call us back and track down the repair guy the receptionist called back and said he will be back soon come 7 pm still no sign so i called again and they told me they are on their way the guys did not show up until 845 pm took them about 30 minutes for the repair \n\nin the end yes my ac was repaired in a timely manner i understand summer is the busiest time for them my issue is the lack of communication we were waiting and waiting and not knowing whats going on and had to keep calling them back also we never got a receipt for the repair which they said they would email us	3	neg
1270461	went here for dinner and left very full and satisfied the family that runs the restaurant is originally from globe so the mexican food here is similar to what youd find in the globemiami area stepping into this restaurant it feels similar to how other restaurants such as serranos and rositas feels comfy atmosphere serving homestyle no frills cuisine i had the special tonight 2 chicken enchiladas with green sauce with rice and beans 850 while the enchiladas looked a bit mangled since chicken pieces were sticking out of the tortilla and the tortilla itself looked prodded and broken the enchiladas themselves were quite tasty with the green sauce and the rice and beans im usually a lightweight when it comes to finishing meals but in this case i cleaned my plate the salsa tastes good but is very watery which makes it hard to eat with chips we ordered iced teas and they were refilled promptly as needed id definitely be interested in going here again for some tasty and filling meals	4	pos

Embedding Layer

```
graph TD; A[Embedding Layer] --> B[LSTM units]; B --> C[Softmax];
```

LSTM units

Softmax