

Bài 6: Trích chọn thông tin với CRFs

Thông tin giảng viên

- Phd Nguyen Kiem Hieu
- Computer science department, School of Information and Communication Technology, HUST
- Email: hieunk@soict.hust.edu.vn

Nội dung buổi học

- Giới thiệu về trích chọn thông tin (Information Extraction)
- Sequence labeling
 - POS tagging
 - Named Entity Recognition
- Hidden Markov Models
- Conditional Random Fields (CRFs)
- Using CRFs tools for IE

Trích chọn thông tin

- Information Extraction
- Trích chọn thông tin là bài toán tự động trích chọn các thông tin có cấu trúc từ các văn bản không cấu trúc.
- Tổ chức lại thông tin một cách có hệ thống, các thông tin trích chọn được có thể đưa vào database làm đầu vào cho các thuật toán khác (data mining).

Firm XYZ is a full service advertising agency specializing in direct and interactive marketing. Located in Bigtown CA, Firm XYZ is looking for an Assistant Account Manager to help manage and coordinate interactive marketing initiatives for a marquee automotive account. Experience in online marketing, automotive and/or the advertising field is a plus. Assistant Account Manager Responsibilities Ensures smooth implementation of programs and initiatives Helps manage the delivery of projects and key client deliverables . . . Compensation: \$50,000-\$80,000 Hiring Organization: Firm XYZ

INDUSTRY	Advertising
POSITION	Assistant Account Manager
LOCATION	Bigtown, CA.
COMPANY	Firm XYZ
SALARY	\$50,000-\$80,000

Trích chọn thông tin

Cập nhật dữ liệu vào CSDL thông qua trích chọn thông tin từ các đoạn văn bản

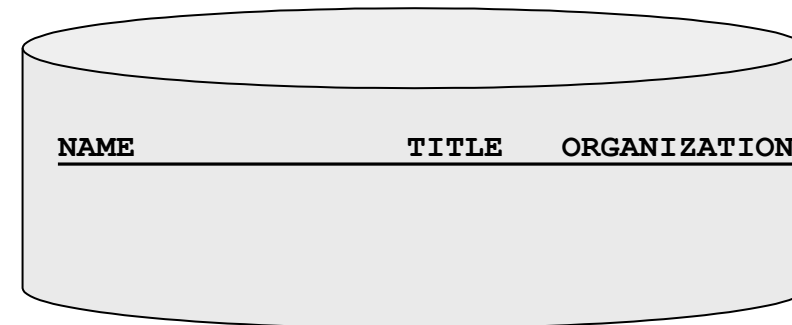
October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...



Trích chọn thông tin

Cập nhật dữ liệu vào CSDL thông qua trích chọn thông tin từ các đoạn văn bản

October 14, 2002, 4:00 a.m. PT

For years, [Microsoft Corporation](#) [CEO Bill Gates](#) railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...



NAME	TITLE	ORGANIZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft..

Độ phức tạp của bài toán

Closed set

U.S. states

He was born in Alabama...

The big Wyoming sky...

Regular set

U.S. phone numbers

Phone: (413) 545-1323

The CALD main office is 412-268-1299

Complex pattern

U.S. postal addresses

University of Arkansas

P.O. Box 140

Hope, AR

Headquarters:

1128 Main Street, 4th Floor

Cincinnati, Ohio 45210

Ambiguous patterns, needing context and many sources of evidence

Person names

...was among the six houses sold
by Hope Feldman that year.

Pawel Opalinski, Software
Engineer at WhizBang Labs.

Trích xuất các mối quan hệ

Jack Welch will retire as CEO of General Electric tomorrow. The top role at the Connecticut company will be filled by Jeffrey Immelt.

Single entity

Person: Jack Welch

Person: Jeffrey Immelt

Location: Connecticut

Binary relationship

Relation: Person-Title

Person: Jack Welch

Title: CEO

Relation: Company-Location

Company: General Electric

Location: Connecticut

N-ary record

Relation: Succession

Company: General Electric

Title: CEO

Out: Jack Welsh

In: Jeffrey Immelt

“Named entity” extraction

Các bài toán con

- Named Entity Recognition (NER)
- Coreference Resolution
- Entity Linking
- Relation Extraction
- Event Extraction

Relation Extraction: Disease Outbreaks

May 19 1995, Atlanta -- The Centers for Disease Control and Prevention, which is in the front line of the world's response to the deadly Ebola epidemic in Zaire, is finding itself hard pressed to cope with the crisis...

**Information
Extraction System**

<i>Date</i>	<i>Disease Name</i>	<i>Location</i>
Jan. 1995	Malaria	Ethiopia
July 1995	Mad Cow Disease	U.K.
Feb. 1995	Pneumonia	U.S.

Relation Extraction: Protein Interactions

“We show that CBF-A and CBF-C interact with each other to form a CBF-A-CBF-C complex and that CBF-B does not interact with CBF-A or CBF-C individually but that it associates with the CBF-A-CBF-C complex.”

CBF-A $\xleftrightarrow[\text{complex}]{\text{interact}}$ CBF-C

CBF-B $\xrightarrow{\text{associates}}$ CBF-A-CBF-C complex

Resolving coreference (both within and across documents)

John Fitzgerald Kennedy was born at 83 Beals Street in Brookline, Massachusetts on Tuesday, September 29, 1917, at 3:00 pm,[7] the second son of Joseph P. Kennedy, Sr., and Rose Fitzgerald; Rose Kennedy, in turn, was the eldest child of John "Honey Fitz" Fitzgerald, a prominent Boston political figure who was the city's mayor and a three-term member of Congress. Kennedy lived in Brookline for his first five years and attended Edward Devotion School, Noble and Greenough Lower School, and the Dexter School, through 4th grade. In 1927, the family moved to 5040 Independence Avenue in the Bronx, New York City; two years later, they moved to 294 Pondfield Road in Bronxville, New York, where Kennedy was a member of Scout Troop 2 (and was the first Boy Scout to become President).[8] Kennedy spent summers with his family at their home in Hyannisport, Massachusetts, and Christmas and Easter holidays with his family at their winter home in Palm Beach, Florida. For the 5th through 7th grade, Kennedy attended Riverdale Country School, a private school for boys. For 8th grade in September 1930, the 13-year old Kennedy attended Canterbury School in New Milford, Connecticut.



Bài toán gán nhãn chuỗi

- Sequence labeling
- Nhiều bài toán NLP có thể đưa về bài toán gán nhãn chuỗi
- **Đầu vào:** một chuỗi các từ
- **Đầu ra:** chuỗi các từ đã được gán nhãn

VBG	NN	IN	DT	NN	IN	NN
Chasing	opportunity	in	an	age	of	upheaval

POS tagging

PERS	O	O	O	ORG	ORG
Murdoch	discusses	future	of	News	Corp.

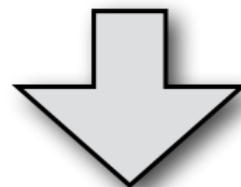
Named entity recognition

B	B	I	I	B	I	B	I	B	B
而	相	对	于	这	些	品	牌	的	价

Word segmentation

Sequence labeling

She promised to back the bill
 $\mathbf{w} =$ $w^{(1)}$ $w^{(2)}$ $w^{(3)}$ $w^{(4)}$ $w^{(5)}$ $w^{(6)}$



$\mathbf{t} =$ $t^{(1)}$ $t^{(2)}$ $t^{(3)}$ $t^{(4)}$ $t^{(5)}$ $t^{(6)}$
PRP VBD TO VB DT NN

- Cho một chuỗi các từ $\mathbf{w}=w^{(1)}...w^{(n)}$, tìm chuỗi các nhãn (tag) có khả năng xảy ra cao nhất $\mathbf{t}=t^{(1)}...t^{(n)}$

$$\mathbf{t}^* = \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t} \mid \mathbf{w})$$

Gán nhãn từ loại

- Part of Speech tagging – POS tagging
- Mỗi từ trong câu được gán nhãn thể từ loại tương ứng của nó
- **Đầu vào:** 1 đoạn văn bản đã tách từ + tập nhãn
- **Đầu ra:** cách gán nhãn chính xác nhất

Gán nhãn từ loại

- Các ứng dụng:
 - Tổng hợp tiếng nói: record - N: ['reko:d], V: [ri'ko:d];
 - Tiền xử lý cho phân tích cú pháp.
 - Nhận dạng tiếng nói, tìm kiếm, v.v...

Tập nhãn cho tiếng Anh

- Tập ngữ liệu Brown: 87 nhãn
- 3 tập thường được sử dụng:
 - Nhỏ: 45 nhãn - Penn treebank
 - Trung bình: 61 nhãn, British national corpus
 - Lớn: 146 nhãn, C7

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
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>	VCN	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	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	“	Left quote	<i>‘ or “</i>
POS	Possessive ending	<i>’s</i>	”	Right quote	<i>’ or ”</i>
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>[, (, {, <</i>
PRP\$	Possessive pronoun	<i>your, one’s</i>)	Right parenthesis	<i>],), }, ></i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>. ! ?</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>: ; ... - -</i>
RP	Particle	<i>up, off</i>			

Ví dụ

- **There/EX** are/VBP 70/CD children/NNS **there/RB**
- EX: từ chỉ sự tồn tại there
- RB: phó từ
- → Khó khăn trong gán nhãn từ loại: nhập nhằng.

Gán nhãn từ loại tiếng Việt

Câu tiếng Việt đã tách từ	Qua những lần từ Sài_Gòn về Quảng_Ngãi kiểm_tra công_việc , Sophie và Jane thường trò_chuyện với Mai , cảm_nhận ngọn_lửa_sống và niềm_tin mãnh_liệt từ người phụ_nữ VN này .
Câu tiếng Việt đã được gán nhãn từ loại	Qua những lần từ Sài_Gòn về Quảng_Ngãi kiểm_tra công_việc , Sophie và Jane thường trò_chuyện với Mai , cảm_nhận ngọn_lửa_sống và niềm_tin mãnh_liệt từ người phụ_nữ VN này .
Chú thích từ loại	<div> DANH TỪ ■ DANH TỪ ■ THÁN TỪ ■ ĐỘNG TỪ ■ PHỤ TỪ ■ TRỢ TỪ ■ TÍNH TỪ ■ GIỚI TỪ ■ TỪ ĐƠN LẺ ■ ĐẠI TỪ ■ CẢM TỪ ■ TỪ VIẾT TẮT ■ ĐỊNH TỪ ■ LIÊN TỪ ■ KHÔNG XÁC ĐỊNH ■ </div>

Nhận dạng thực thể có tên

- Named-entity recognition (NER)
- Bài toán con quan trọng của trích rút thông tin
- Thực thể (entity) là đối tượng hoặc tập hợp các đối tượng trong thế giới tự nhiên được mô tả bằng ngôn ngữ
- Phân loại:
 - Tên người
 - Tên địa điểm
 - Tên tổ chức
 - Giá trị số
 - Thời gian

Named-entity recognition

- nhận dạng trong văn bản các nhóm thực thể có tên đã được định trước như tên người, tổ chức, địa điểm, thời gian, ...
- Các nhãn (tag)
 - PERS
 - ORG
 - LOC
 - DATE



Named-entity recognition

Pierre Vinken , 61 years old , will join IBM 's board
as a nonexecutive director Nov. 29 .



[PERS Pierre Vinken] , 61 years old , will join
[ORG IBM] 's board as a nonexecutive director
[DATE Nov. 2] .

Nhãn BIO

- Định nghĩa các nhãn (tag) mới:
 - **B-PERS, B-DATE, ...**: Đánh dấu bắt đầu thực thể có tên (Begin)
 - **I-PERS, I-DATE, ...**: Đánh dấu các từ tiếp theo của thực thể có tên (Inside)
 - **O**: Đánh dấu các từ không thuộc thực thể có tên (Outside)

Nhãn BIO

[PERS Pierre Vinken] , 61 years old , will join
[ORG IBM] 's board as a nonexecutive director
[DATE Nov. 2] .



Pierre_B-PERS Vinken_I-PERS ,_O 61_O years_O old_O ,_O
will_O join_O IBM_B-ORG 's_O board_O as_O a_O
nonexecutive_O director_O Nov._B-DATE 29_I-DATE ._O

VLSP 2016

	POS tag	Chunking tag	NE	Nested NE
Anh	N	B-NP	O	O
Thanh	Np	I-NP	I-PER	O
là	V	B-VP	O	O
cán_bộ	N	B-NP	O	O
Ủy ban	N	B-NP	B-ORG	O
nhân_dân	N	I-NP	I-ORG	O
Thành_phố	N	I-NP	I-ORG	B-LOC
Hà_Nội	Np	I-NP	I-ORG	I-LOC
.	.	O	O	O

Các phương pháp

- Các phương pháp rule-based
 - Email, Thời gian, Số điện thoại, URL, Số lượng tiền
- Học máy thống kê
 - Mô hình Markov ẩn (Hidden Markov Model - HMM)
 - Maximum Entropy Markov Model (MEMM)
 - Conditional Random Field (CRF)
- Học sâu
 - RNN/LSTM
 - BERT

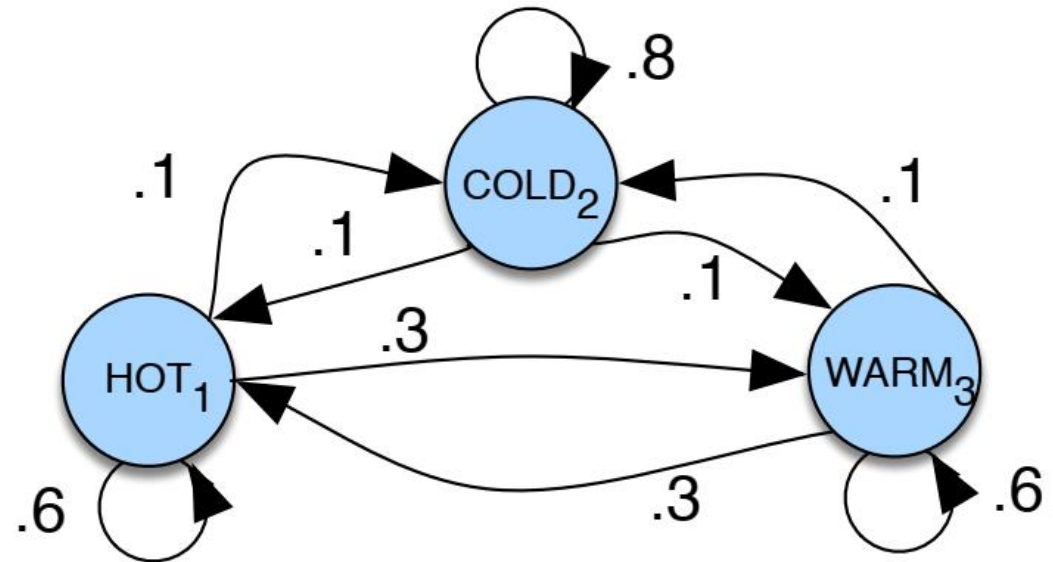
Các thư viện

- NLTK
- Spacy
- Stanford Core NLP
- Allen NLP
- Flair

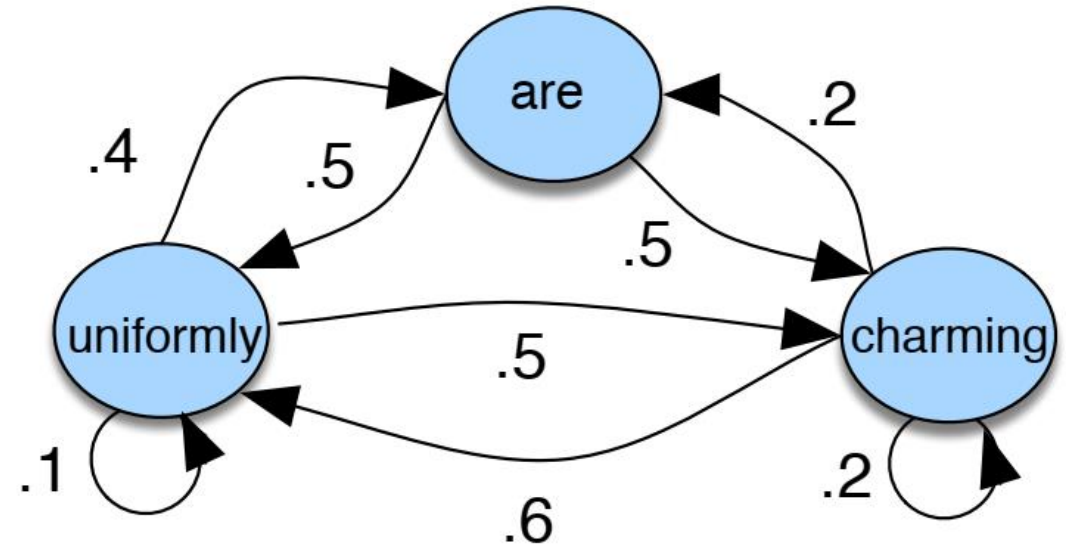
Các mô hình Markov

- Mô hình Markov ẩn là một trong những mô hình học máy quan trọng
- Các mô hình Markov cơ bản: mô hình chuỗi Markov và mô hình Markov ẩn
- Mô hình chuỗi Markov (Markov chain), hay còn được gọi là mô hình Markov có thể quan sát được, là mô hình Markov đơn giản nhất
- Mô hình chuỗi Markov và Markov ẩn đều được mở rộng từ Automat hữu hạn
- Một automat hữu hạn có trọng số có các cạnh được gắn với các xác suất, biểu diễn xác suất đi vào cạnh đó. Tổng tất cả các xác suất của các cạnh đi ra từ một đỉnh phải bằng 1.
- Chuỗi Markov là trường hợp đặc biệt của automat hữu hạn có trọng số, khi đó chuỗi đầu vào sẽ quyết định các trạng thái mà automat sẽ đi qua.

Chuỗi Markov



(a)



(b)

Phân bố ban đầu $\pi = [0.1, 0.7, 0.2]$

Chuỗi Markov

- Các thành phần:

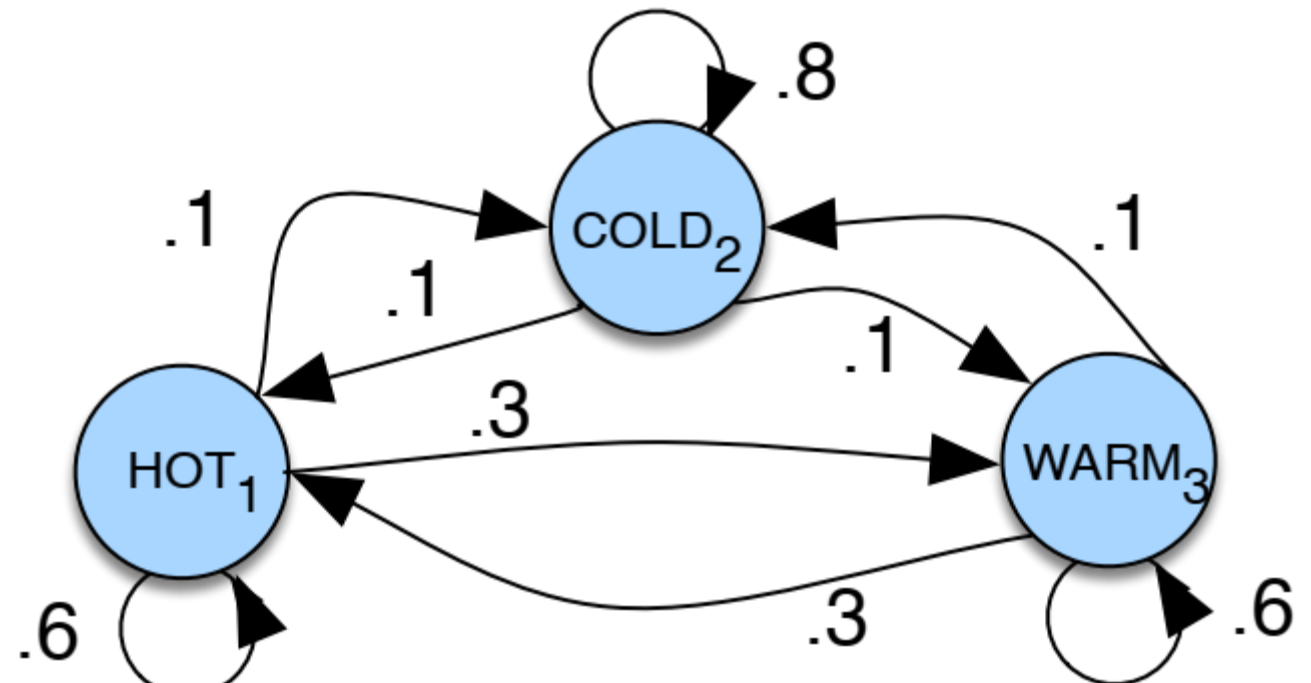
$Q = q_1 q_2 \dots q_N$	tập hợp N trạng thái
$A = a_{11} a_{12} \dots a_{N1} \dots a_{NN}$	ma trận xác suất chuyển trạng thái A, a_{ij} biểu diễn xác suất chuyển từ trạng thái i sang trạng thái j $\sum_{j=1}^N a_{ij} = 1 \forall i$
$\pi = \pi_1, \pi_2, \dots, \pi_N$	phân bố xác suất ban đầu của các trạng thái $\sum_{i=1}^N \pi_i = 1$

- Giả định Markov: xác suất của một trạng thái chỉ phụ thuộc vào trạng thái trước nó

$$P(q_i = a | q_1 \dots q_{i-1}) = P(q_i = a | q_{i-1})$$

Ví dụ

- Hãy tính xác suất của các chuỗi sau:
 - hot hot hot hot
 - cold hot cold hot



Phân bố ban đầu $\pi = [0.1, 0.7, 0.2]$

Mô hình Markov ẩn

- Mô hình chuỗi Markov dùng để tính xác suất của một chuỗi sự kiện mà chúng ta có thể quan sát được
- Tuy nhiên, trong nhiều trường hợp thì có những sự kiện chúng ta quan tâm có thể không quan sát trực tiếp được
- Mô hình Markov ẩn cho phép chúng ta xem xét cả các sự kiện quan sát được và các sự kiện ẩn.

VBG	NN	IN	DT	NN	IN	NN
Chasing	opportunity	in	an	age	of	upheaval

POS tagging

Mô hình Markov ẩn

- Các thành phần:

$Q = q_1 q_2 \dots q_N$	tập hợp N trạng thái
$A = a_{11} a_{12} \dots a_{N1} \dots a_{NN}$	ma trận xác suất chuyển trạng thái A, a_{ij} biểu diễn xác suất chuyển từ trạng thái i sang trạng thái j $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	chuỗi sự kiện quan sát được
$B = b_i(o_t)$	emission probabilities: xác suất sự kiện o_t được sinh ra từ trạng thái q_i
$\pi = \pi_1, \pi_2, \dots, \pi_N$	phân bố xác suất ban đầu của các trạng thái $\sum_{i=1}^N \pi_i = 1$

Xác suất chuyển trạng thái

- Xác suất chuyển trạng thái được tính bằng cách đếm số lần xuất hiện của các nhãn trên kho ngữ liệu

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

- Ví dụ, MD xuất hiện 13124 lần, trong đó có 10471 lần VB xuất hiện ngay sau nó

$$P(VB|MD) = \frac{C(MD, VB)}{C(MD)} = \frac{10471}{13124} = .80$$

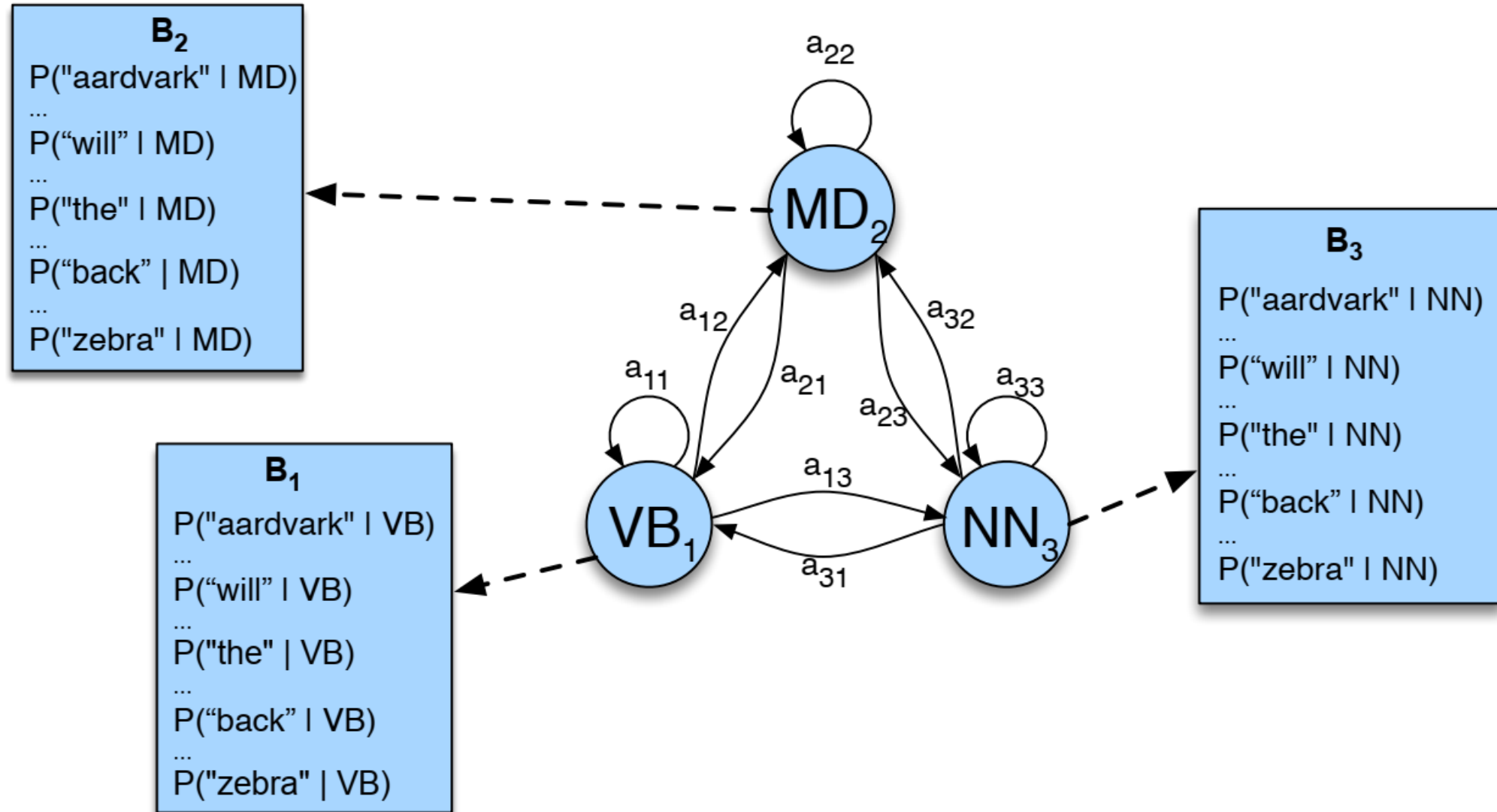
Emission probabilities

- Xác suất sự kiện o_t được sinh ra từ trạng thái q_i
- Xác suất một nhãn đi với một từ

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

$$P(will|MD) = \frac{C(MD, will)}{C(MD)} = \frac{4046}{13124} = .31$$

Mô hình Markov ẩn



HMM cho sequence labeling

- Xác định chuỗi trạng thái ẩn tương ứng với chuỗi quan sát được, được gọi là decoding
- Cho đầu vào là một HMM $\lambda = (A, B)$ và một chuỗi quan sát được $O = o_1 o_2 \dots o_T$, tìm chuỗi các trạng thái $Q = q_1 q_2 \dots q_N$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Định lý Bayes

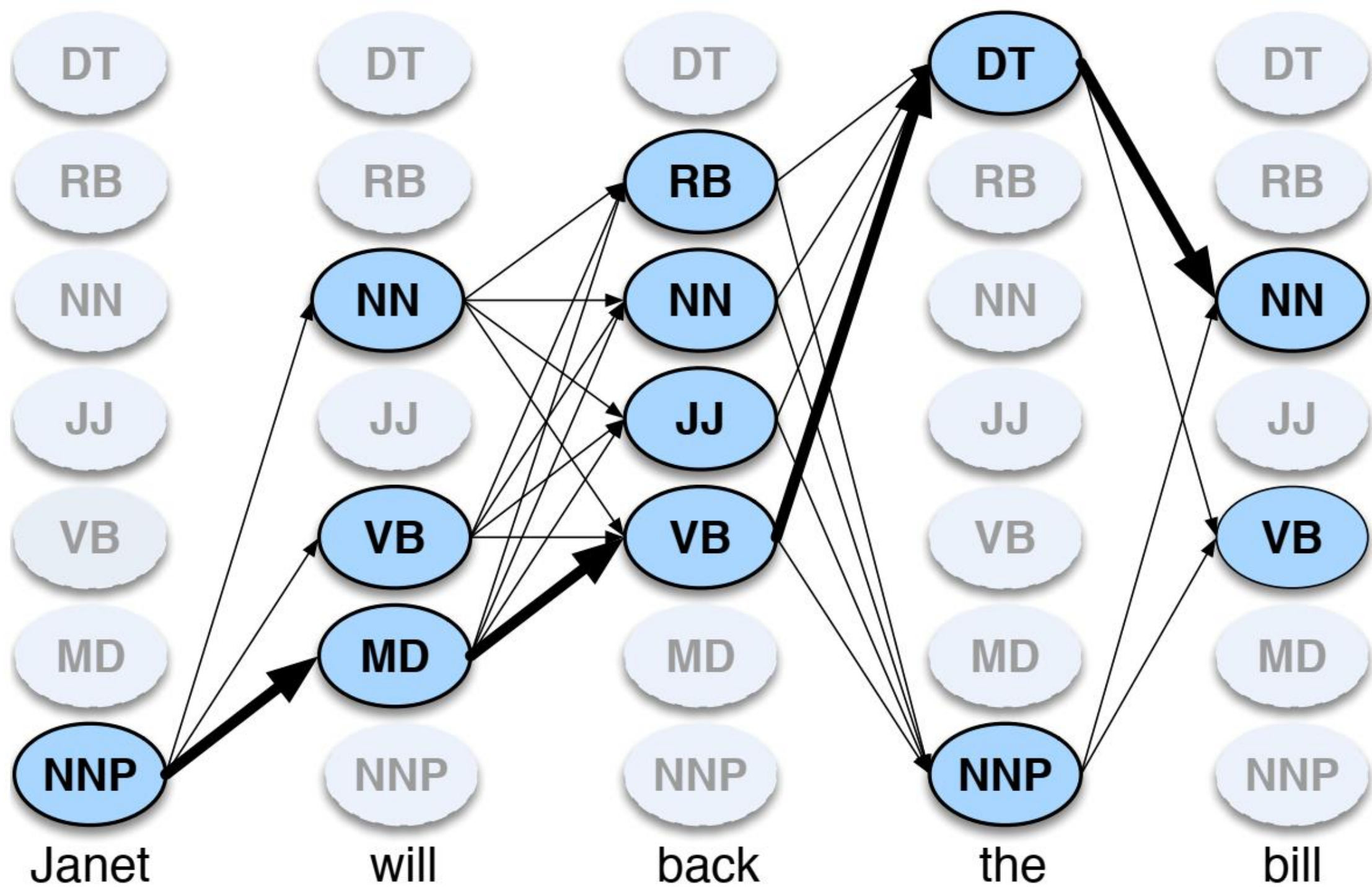
$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$$

HMM cho sequence labeling

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) P(t_1^n)$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$$

- Giải quyết bằng thuật toán Viterbi (quy hoạch động)



Bài tập 1

- Tính các xác suất $P(AT|PERIOD)$, $P(NN|AT)$, $P(BEZ|NN)$, $P(IN|BEZ)$, $P(AT|IN)$, và $P(PERIOD|NN)$. Giả sử tổng số lần xuất hiện của một thẻ là tổng theo hàng.

First tag	Second tag					
	AT	BEZ	IN	NN	VB	PERIOD
AT	0	0	0	48636	0	19
BEZ	1973	0	426	187	0	38
IN	43322	0	1325	17314	0	185
NN	1067	3720	42470	11773	614	21392
VB	6072	42	4758	1476	129	1522
PERIOD	8016	75	4656	1329	954	0

Bài tập 2

- Tính các xác suất $P(\text{bear}|\text{VB})$, $P(\text{bear}|\text{AT})$ lấy số lần xuất hiện của thẻ từ bài tập 1

	AT	BEZ	IN	NN	VB	PERIOD
<i>bear</i>	0	0	10	0	43	0
<i>is</i>	0	10065	0	0	0	0
<i>move</i>	0	0	0	36	133	0
<i>on</i>	0	0	5484	0	0	0
<i>president</i>	0	0	0	382	0	0
<i>progress</i>	0	0	0	108	4	0
<i>the</i>	69016	0	0	0	0	0
<i>.</i>	0	0	0	0	0	48809

Bài tập 3

- Dựa trên đó tính:
- $P(AT\ NN\ BEZ\ IN\ AT\ NN | \text{The bear is on the move})$
- $P(AT\ NN\ BEZ\ IN\ AT\ VB | \text{The bear is on the move})$

Conditional Random Fields

- Lafferty et al. 2001
- Được áp dụng rộng rãi trong nhiều lĩnh vực từ XLNNTN đến thị giác máy, phân tích chuỗi trong sinh học
- CRF là phương pháp thống kê nhưng thường vẫn được sử dụng kết hợp với các mô hình học sâu

Mô hình generative và discriminative

- Luôn cố gắng mô hình hóa phân bố xác suất trên (y, x)
- HMM: mô hình generative của chuỗi đầu vào x , mô tả phân bố “sinh” ra x khi đã biết nhãn y (áp dụng định lý Bayes)

$$\hat{y} = \operatorname{argmax}_y P(y|x) = \operatorname{argmax}_y P(x|y)P(y)$$

- Các mô hình discriminative (CRF) trực tiếp mô hình hóa $P(y|x)$ bằng các hàm đặc trưng

Conditional Random Fields

- Phân bố $P(\mathbf{y}|\mathbf{x})$ trong CRF được định nghĩa như sau

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{t=1}^T \exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}$$

- Vector tham số $\theta = \{\theta_k\} \in \Re^K$
- Hàm chuẩn hóa

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \prod_{t=1}^T \exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}$$

NER dựa trên CRF

- [1]: Sử dụng PoS và phân cụm chuẩn
- [2, 3]: PoS và phân cụm tự động bởi công cụ *NNVLP* và *Underthesea*
- [4]: Không sử dụng PoS và phân cụm

Table 4. Accuracy of our NER system with default and generated PoS, chunking tags; and without PoS and chunking tags

Setting	Precision	Recall	F_1
Default PoS and chunking tags	93.87	93.99	93.93
PoS and chunking tags generated by NNVLP 7	90.21	86.72	88.43
PoS and chunking tags generated by Underthesea	90.28	88.35	89.3
Without PoS, chunking tags	89.91	90.15	90.03

Đánh giá kết quả

- [1]: Sử dụng PoS chuẩn
- [2-6]: PoS tự động từ các công cụ
- [7]: Không sử dụng PoS và phân cụm

Table 5. Proposed NER systems without chunking tag-based features. We compare default PoS with PoS generated by other tools.

Setting	Precision	Recall	F_1
Default PoS tags	90.13	90.55	90.34
PoS by NNVL [7]	90.05	85.65	88.31
PoS by Underthesea	90.27	88.58	89.42
PoS by Pyvi	90.16	88.72	89.43
PoS by Vtik	89.62	86.42	87.99
PoS by VnMarMoT [19]	90.51	89.15	89.83
Without PoS, chunking tags	89.91	90.15	90.03

Đánh giá kết quả (tiếp)

- [1]: Sử dụng tách từ chuẩn
- [2,3]: Tách từ tự động sử dụng *UETSegmenter* và *RDRSegmenter*

Table 6. Accuracy of NER system with default and generated word segmentation. We did not use features based on PoS, chunking tags here.

Setting	Precision	Recall	F_1
Default Word segmentation	89.91	90.15	90.03
Word segmentation generated by UETSegmenter	87.67	84.95	86.29
Word segmentation generated by RDRsegmenter	89.05	84.98	86.97

Đánh giá kết quả (tiếp)

- [1]: Mô hình dựa trên tiếng (không tách từ)
- [2]: Sử dụng tách từ chuẩn
- [3]: Tách từ tự động bằng công cụ *RDRSegmenter*

Table 7. Accuracy of NER system with syllable-based and word-based model. We do not use features based on PoS and chunking tags. “ws” stands for word segmentation

Setting	Precision	Recall	F_1
Syllable-based model	88.78	82.94	85.76
Word-based model with gold ws	89.91	90.15	90.03
Word-based model with ws generated by RDRsegmenter	89.05	84.98	86.97

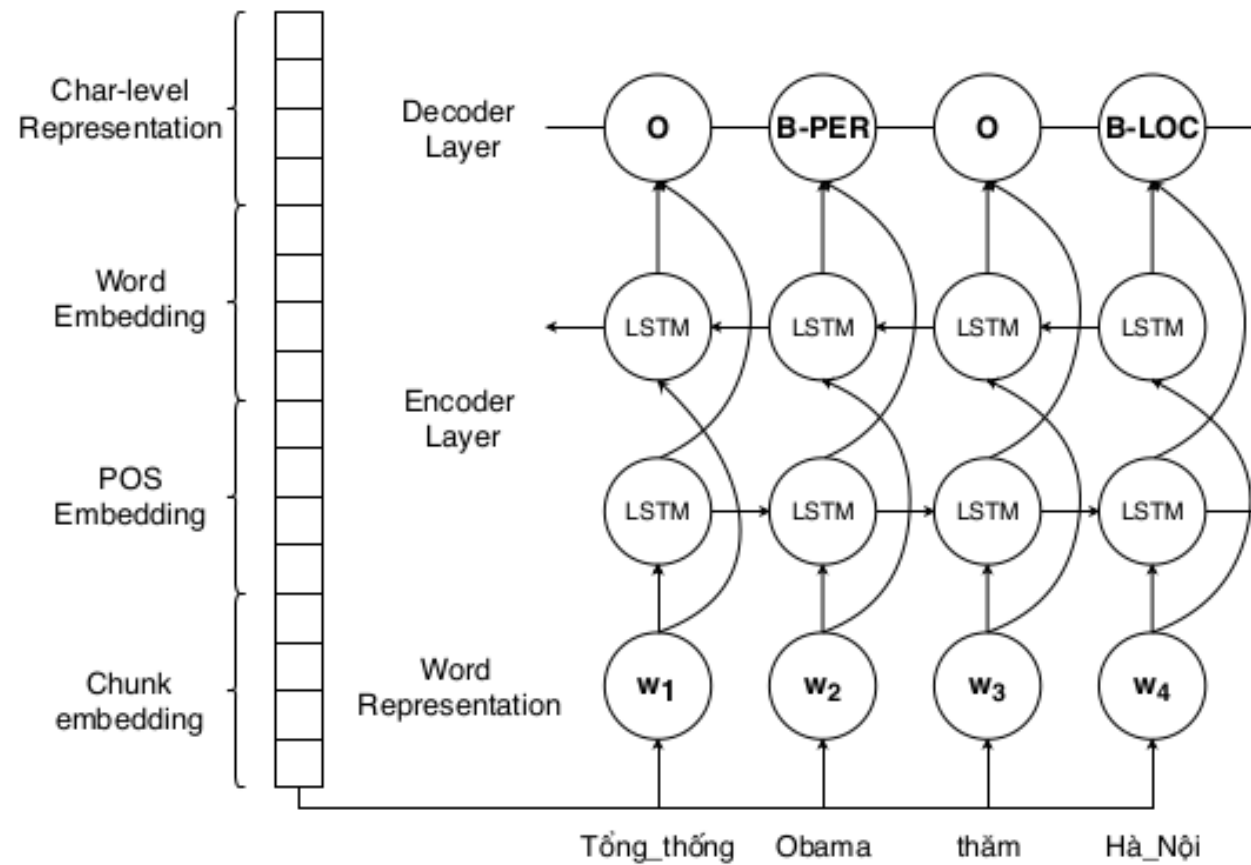
Đánh giá kết quả (tiếp)

- Word: Các từ trong cửa sổ
- Word shapes: Hình thái từ
- w2v: Biểu diễn từ nhúng
- Cluster: Biểu diễn phân cụm Brown

Table 8. Impact of word representation-based features. w2v denotes features based on word embeddings. “cluster” denotes cluster-based features.

Setting	Precision	Recall	F_1
(1) = all features with default PoS, Chunk	93.87	93.99	93.93
(2) = (1) - cluster - w2v	91.66	92.02	91.84
(4) = word + word shapes + default PoS	88.01	87.95	87.98
(5) = word + word shapes + cluster + w2v	89.91	90.15	90.03
(6) = word + word-shapes	88.17	88.08	88.13
(7) = word + word-shapes + w2v	88.69	88.72	88.70
(8) = word + word-shapes + cluster	88.96	89.99	89.97

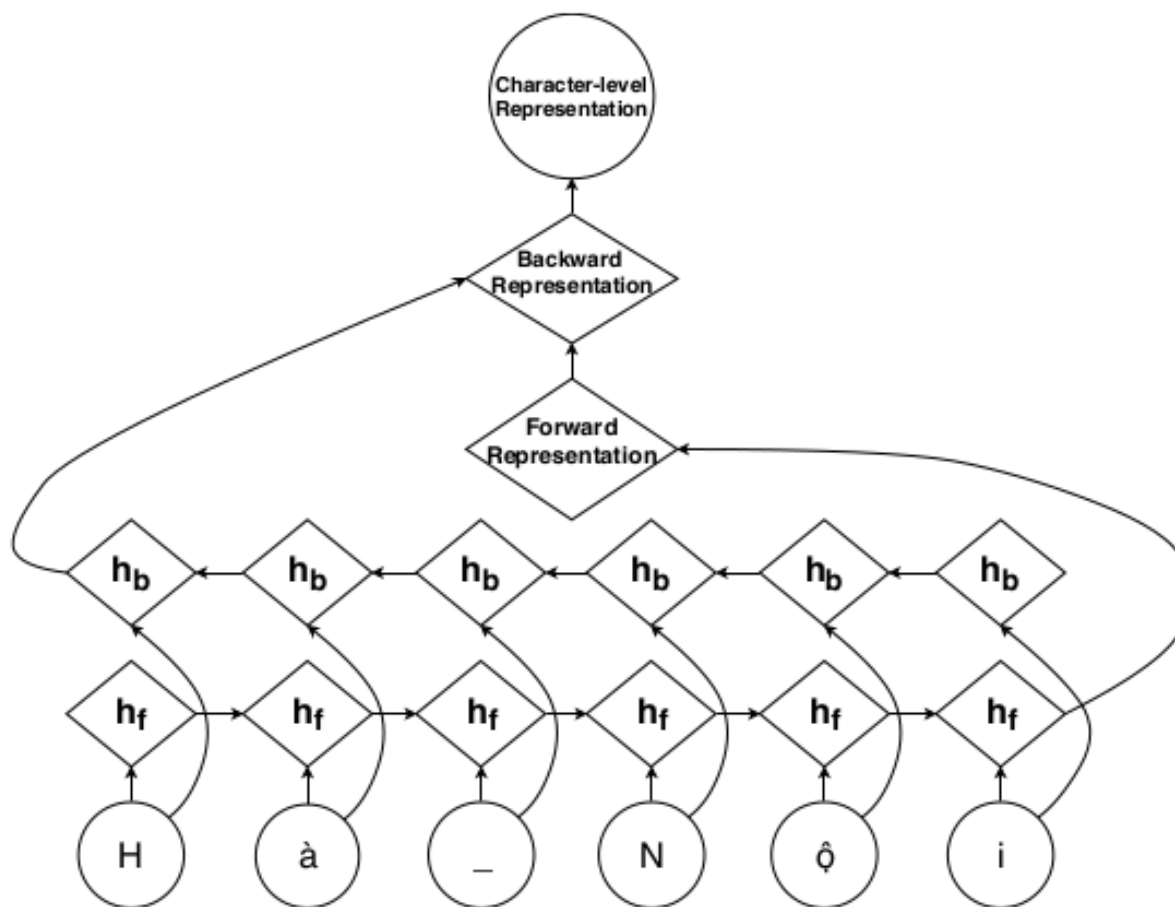
NER dựa trên RNN



Tầng đầu vào

- Biểu diễn nhúng kết hợp:
 - Biểu diễn từ: Sử dụng từ nhúng huấn luyện trước bởi word2vec trên 2 triệu văn bản
 - Biểu diễn ký tự: Sử dụng mạng LSTM hai chiều để học biểu diễn ký tự với khởi tạo ngẫu nhiên
 - Biểu diễn từ loại: Biểu diễn one-hot
 - Biểu diễn cụm: Biểu diễn one-hot

Học biểu diễn ký tự



LSTM hai chiều

- Sử dụng hai mạng LSTM theo chiều tiến và chiều lùi
 - Mục đích: Các từ ở đầu câu có thể sử dụng cả thông tin ở cuối câu để dự đoán và ngược lại
- Đầu ra được ghép nối để đưa vào tầng đầu ra

Tầng đầu ra

- Dự đoán các nhãn BIO ứng với các loại thực thể
 - VD: Với 3 loại thực thể ORG, PER, LOC, tập nhãn có 7 nhãn (B-ORG, I-ORG, B-PER, I-PER, B-LOC, I-LOC, O)
- Tầng đầu ra có thể được đưa vào một mô hình CRFs để thể hiện quan hệ với nhãn ở thời điểm trước thông qua xác suất chuyển đổi

Đánh giá kết quả

Method	P	R	F1	F1 (w.o char)
Feature-rich CRFs [25]	93.87	93.99	93.93	-
NNVLP [7]	92.76	93.07	92.91	-
BiLSTM-CRFs	90.97	87.52	89.21	76.43
BiLSTM-CRFs + POS	90.90	90.39	90.64	86.06
BiLSTM-CRFs + Chunk	95.24	92.16	93.67	87.13
BiLSTM-CRFs + POS + Chunk	95.44	94.33	94.88	91.36

BiLSTM-CRFs sử dụng thêm các thông tin PoS và phân cụm

BiLSTM-CRFs không kết hợp biểu diễn mức kí tự

Bài thực hành

- Sử dụng thư viện `sklearn_crfsuite` để huấn luyện mô hình CRF cho bài toán NER
- Sử dụng mô hình CRF để trích chọn thông tin từ văn bản

Thank you!