

NLP Exercise

Part-of-Speech Tagging Medical Named Entity Recognition

AI VIET NAM
Nguyen Quoc Thai

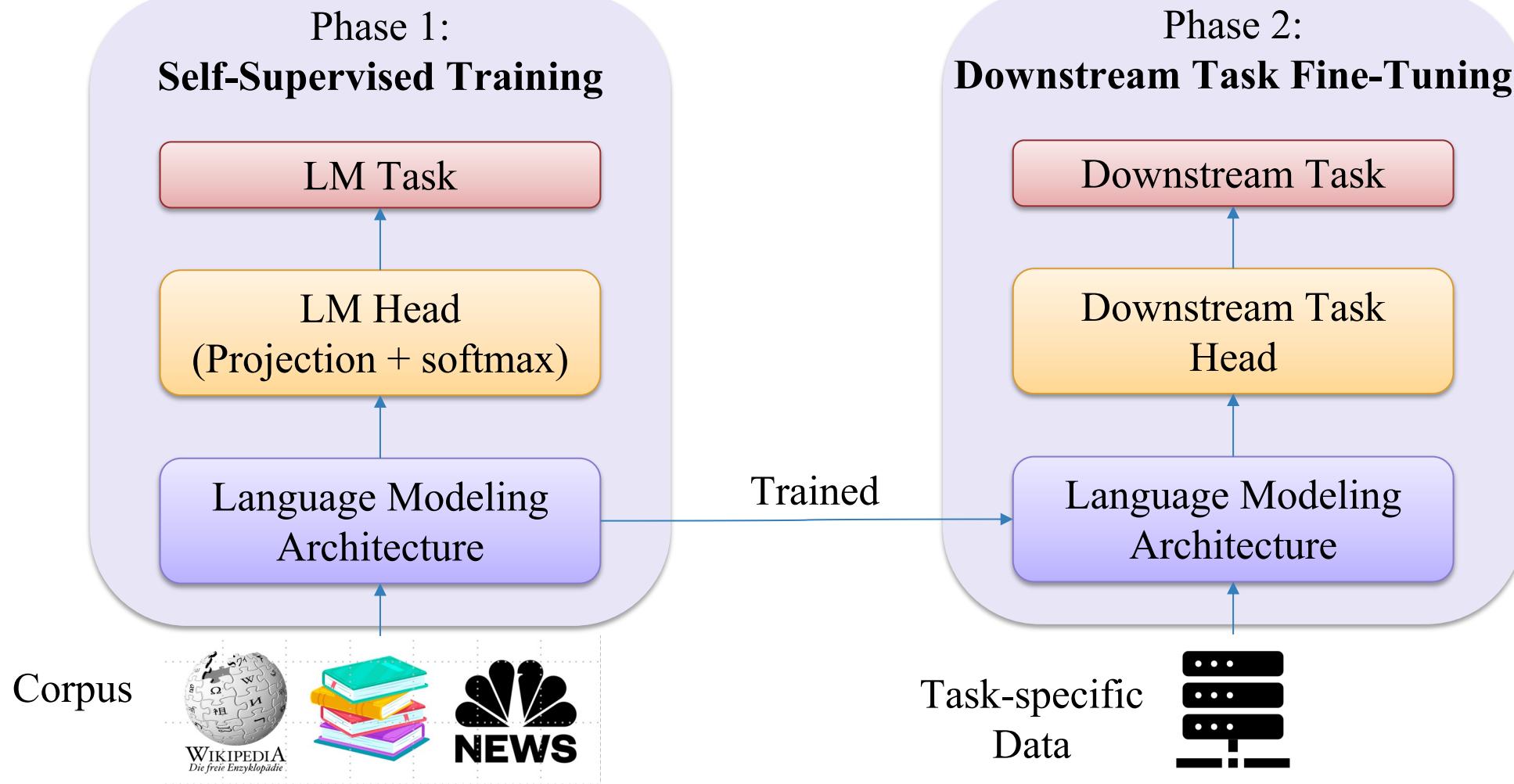
Outline

- **Introduction**
- **Part-of-Speech Tagging (POS)**
- **Named Entity Recognition (NER)**
- **Medical Named Entity Recognition**

Introduction



Pre-trained Models for Text

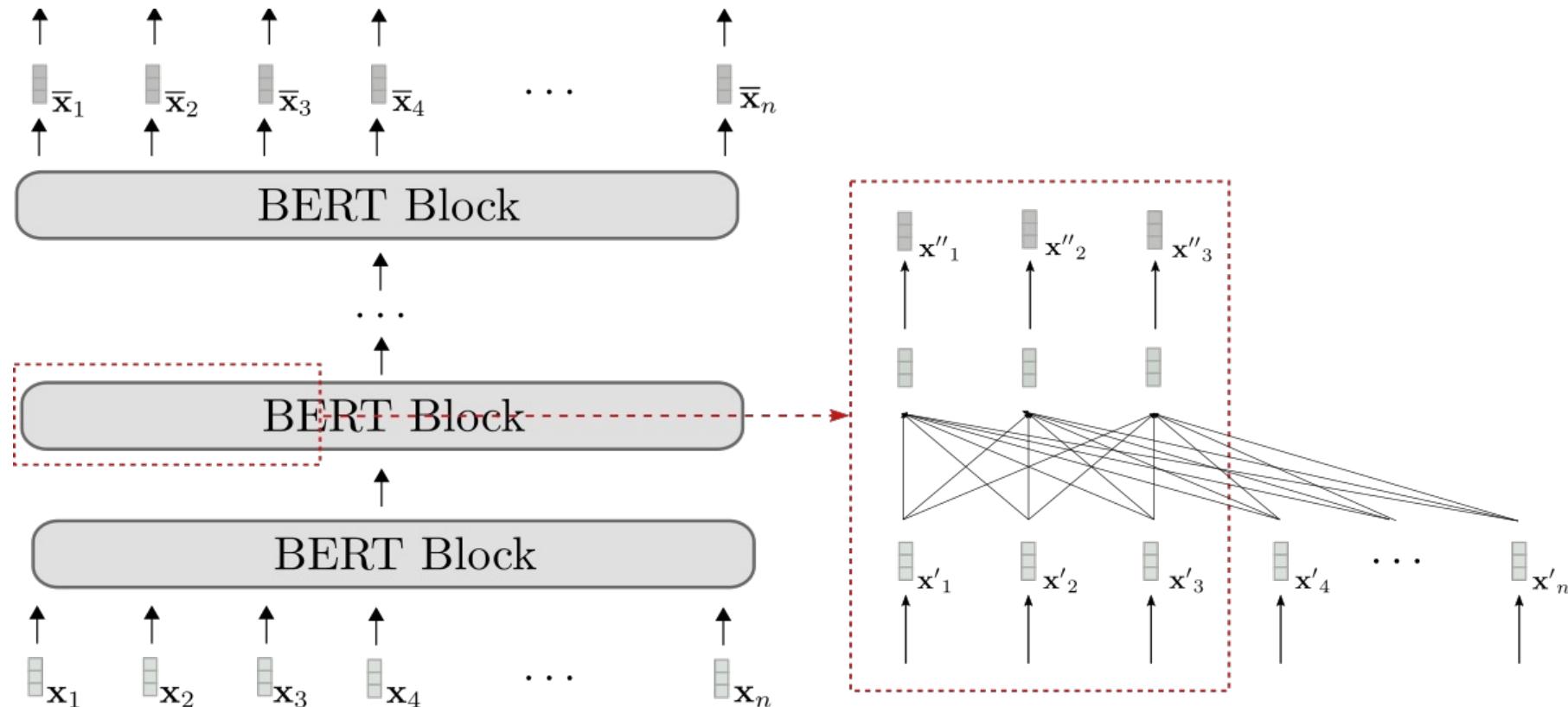


Introduction



BERT: Encoder Model

- ❖ Maps an input sequence to a contextualized sequence: $f_{\theta_{BERT}}: X_{1:n} \rightarrow \bar{X}_{1:n}$

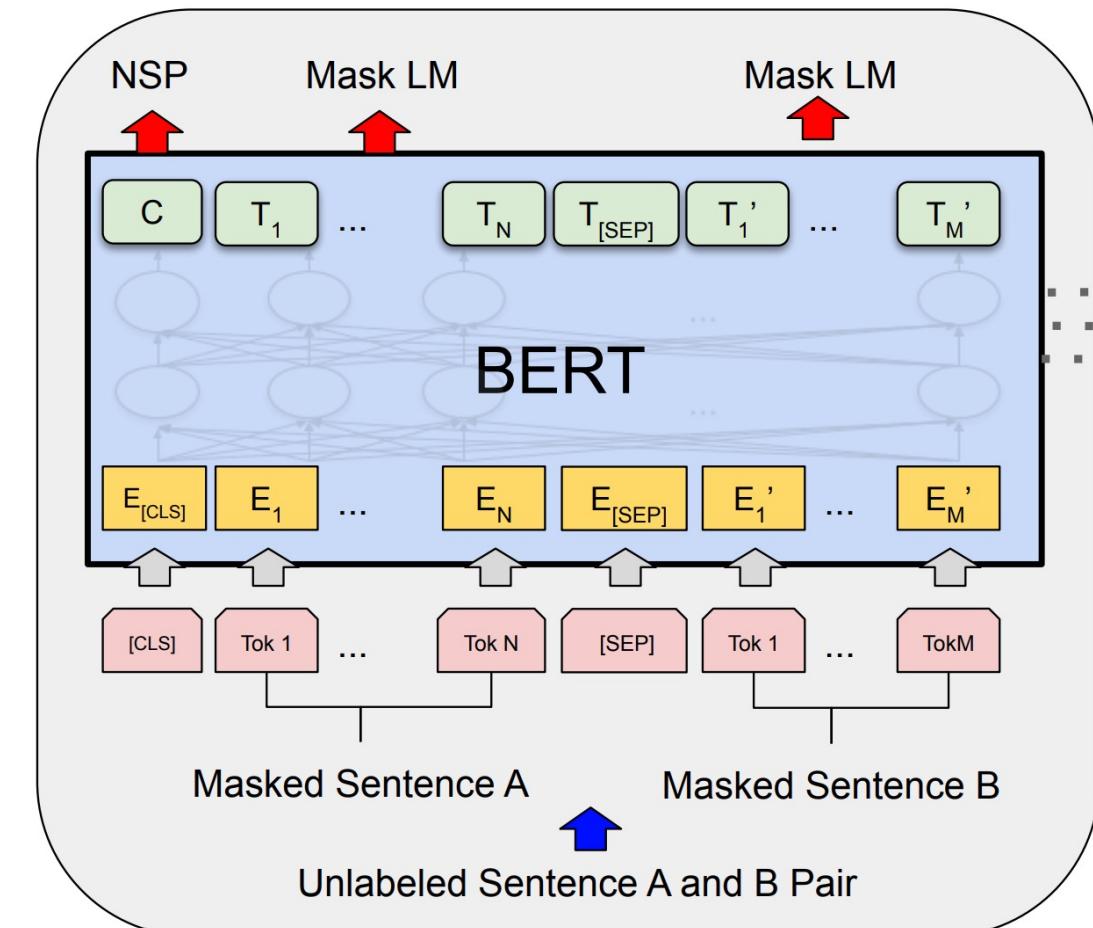


Introduction



BERT: Pre-Training

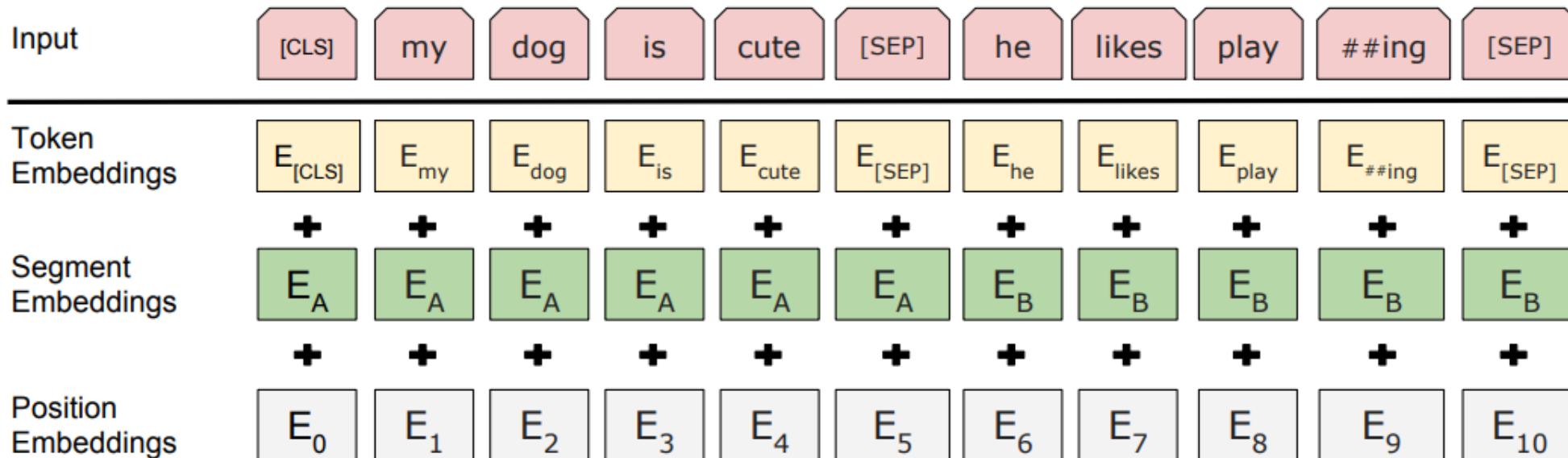
- ❖ Masked LM (15% token):
 - 80%: replace with [MASK]
 - 10%: replace with a random word
 - 10%: keep unchanged
- ❖ Next Sentence Prediction (NSP)
 - Classification Task
 - 2 Labels: IsNext and NotNext
 - Use [SEP] [CLS] token



Introduction



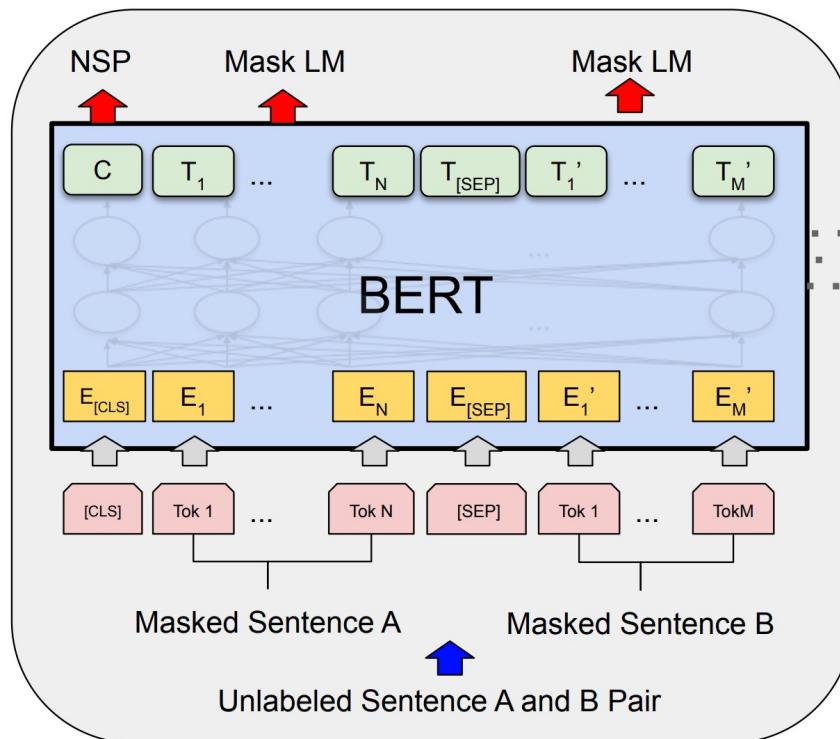
BERT: Input Representation



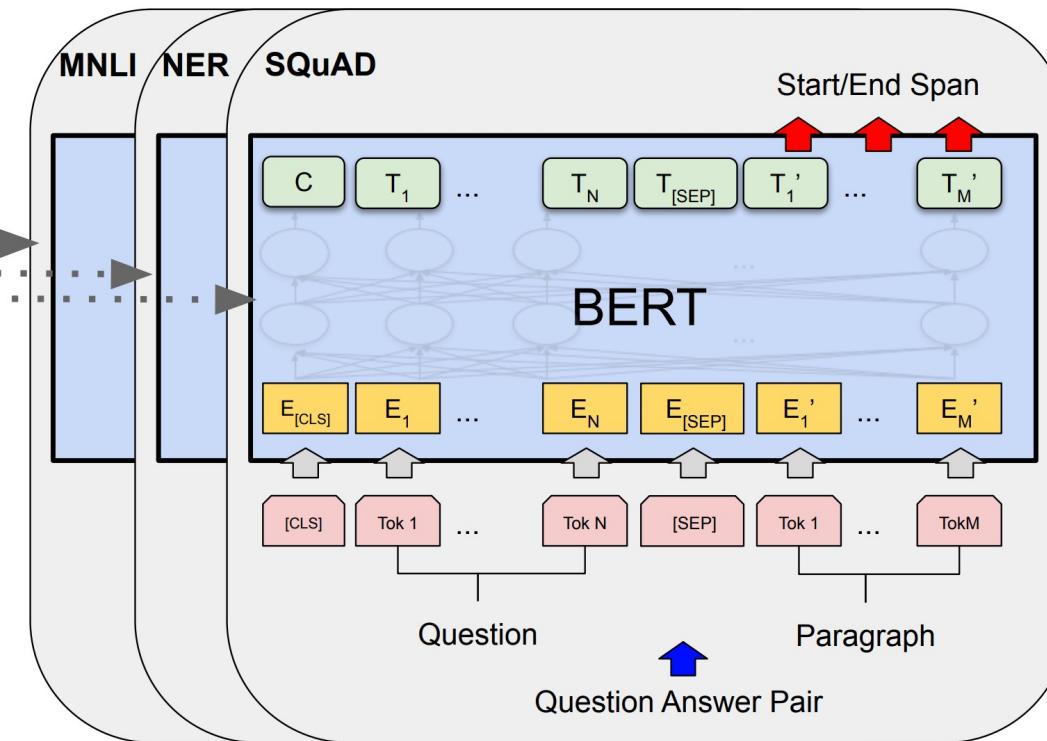
Introduction



BERT: Fine-Tuning



Pre-training



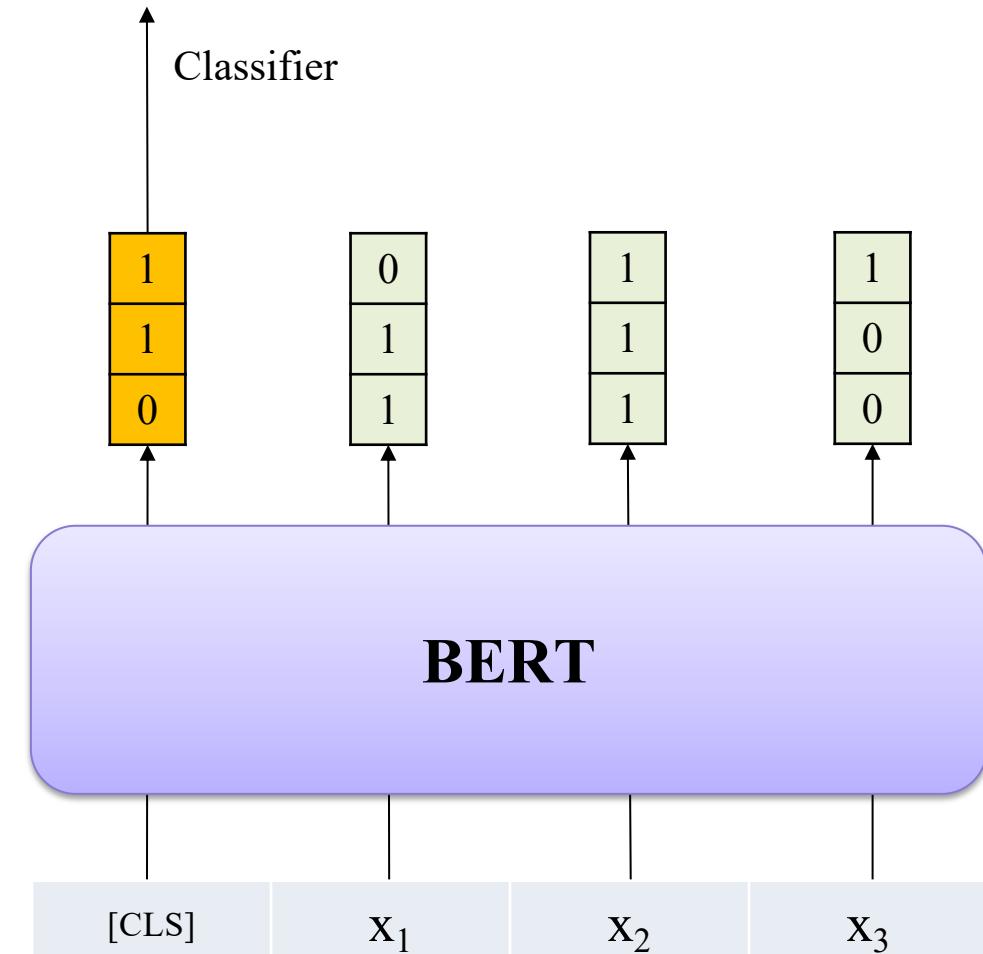
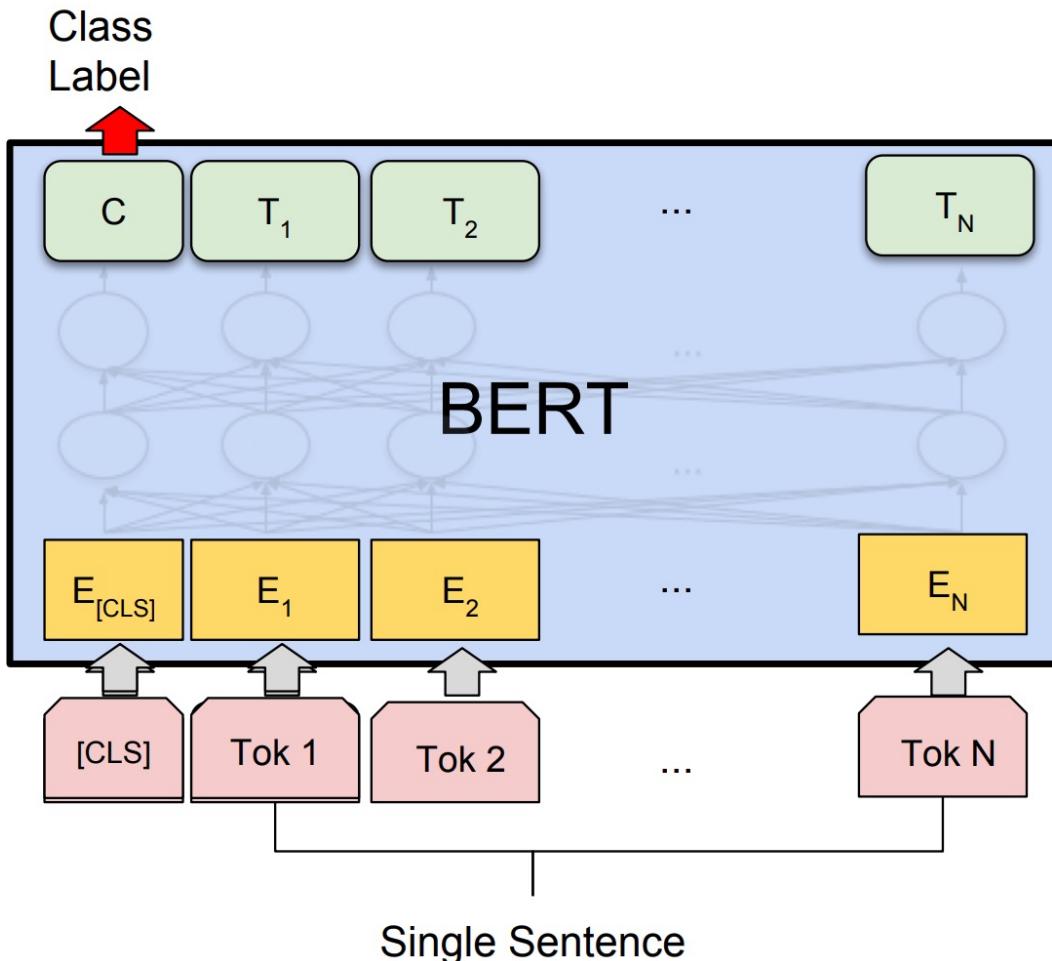
Fine-Tuning

Introduction



BERT: Text Classification

- ❖ Document-Level Text Classification

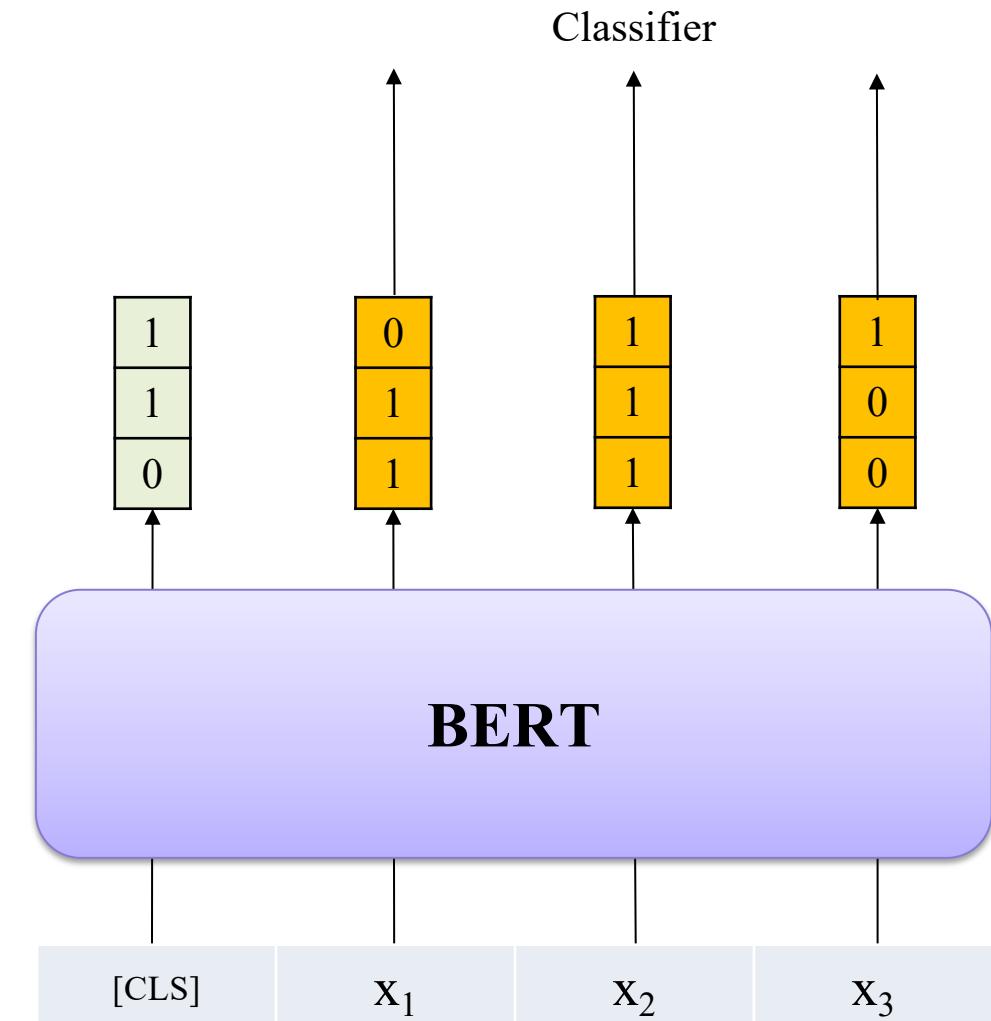
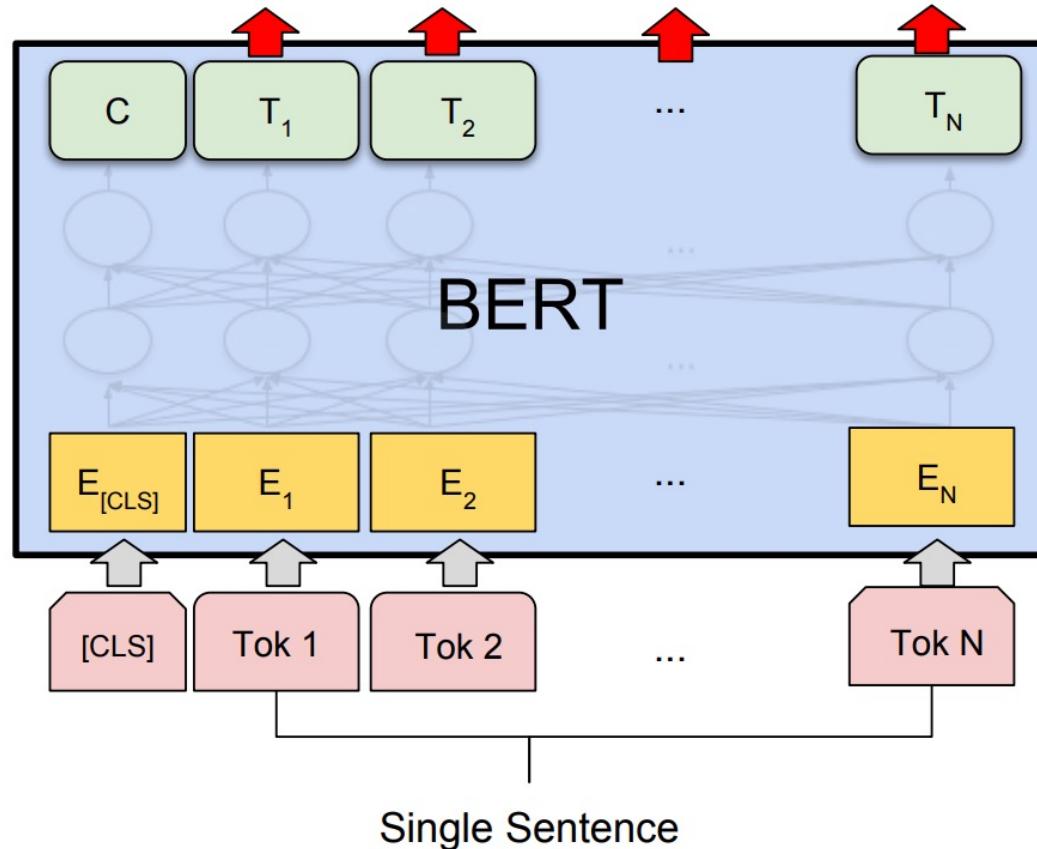


Introduction



BERT: Text Classification

- ❖ Token-Level Text Classification



Introduction



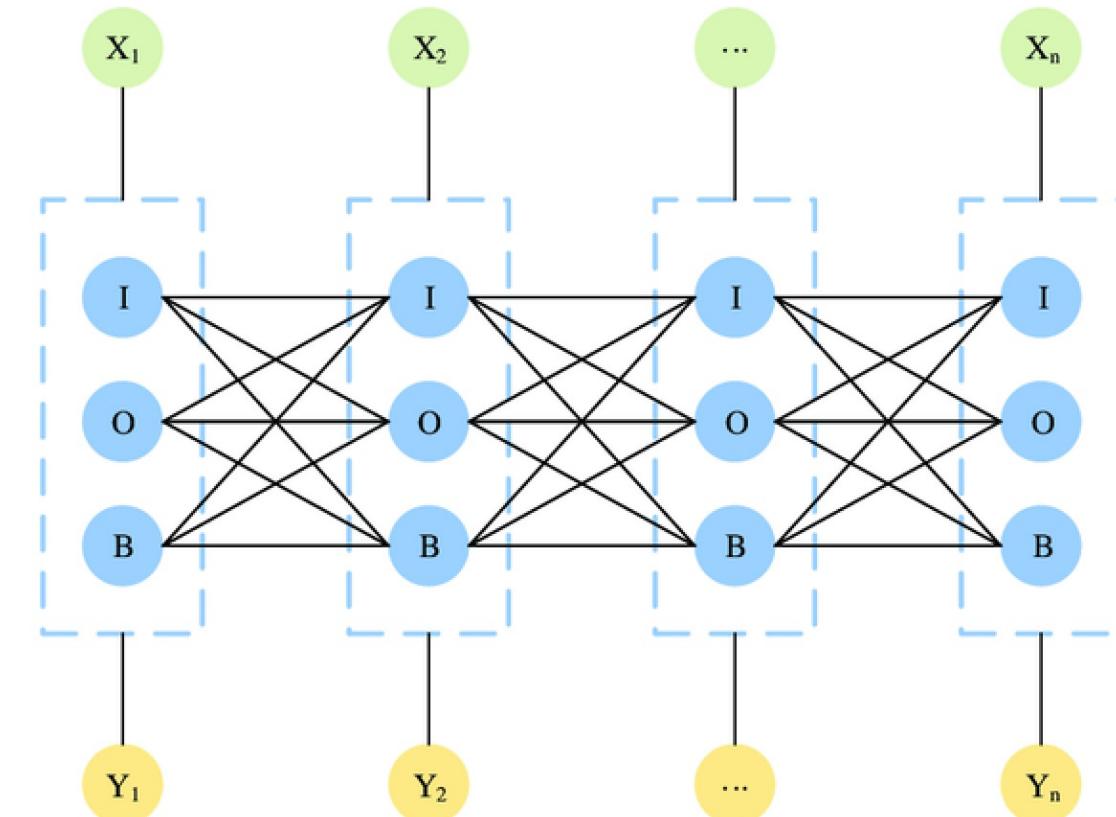
Token-Level Text Classification

❖ Input:

Sequence n tokens: $\{w_1, w_2, \dots, w_n\}$

❖ Output:

Sequence n tokens: $\{y_1, y_2, \dots, y_n\}$



Introduction



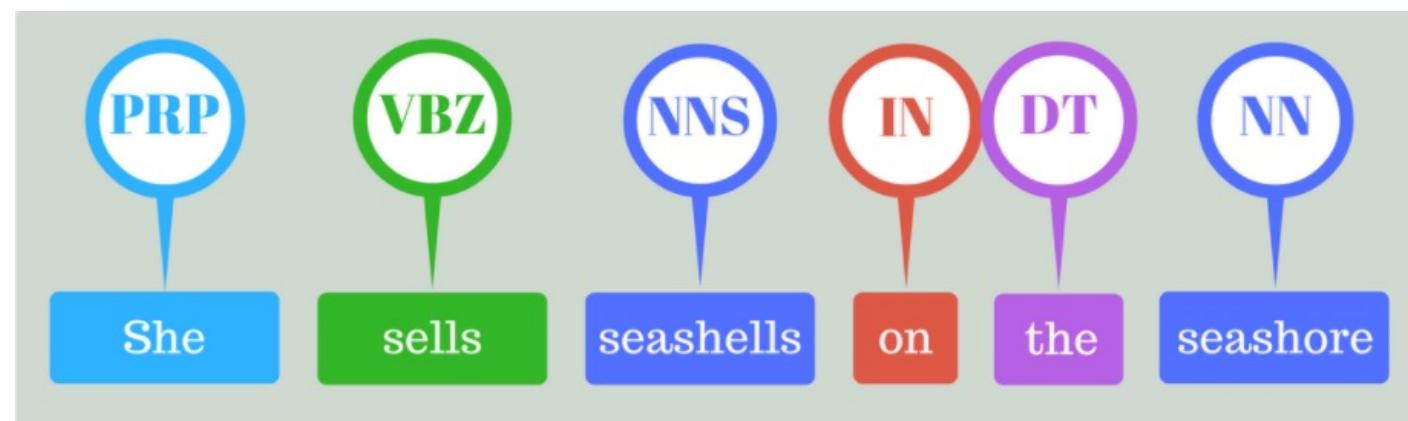
Token-Level Text Classification

❖ Input:

Sequence n tokens: $\{w_1, w_2, \dots, w_n\}$

❖ Output:

Sequence n tokens: $\{y_1, y_2, \dots, y_n\}$



Part-of-Speech (POS) Tagging

Introduction



Token-Level Text Classification

- ❖ Input:

Sequence n tokens: $\{w_1, w_2, \dots, w_n\}$

- ❖ Output:

Sequence n tokens: $\{y_1, y_2, \dots, y_n\}$

When Sebastian Thrun PERSON started at Google ORG in 2007 DATE, few people outside of the company took him seriously. "I can tell you very senior CEOs of major American NORP car companies would shake my hand and turn away because I wasn't worth talking to," said Thrun PERSON, now the co-founder and CEO of online higher education startup Udacity, in an interview with Recode ORG earlier this week DATE.

A little less than a decade later DATE, dozens of self-driving startups have cropped up while automakers around the world clamor, wallet in hand, to secure their place in the fast-moving world of fully automated transportation.

Named Entity Recognition (NER)

Outline

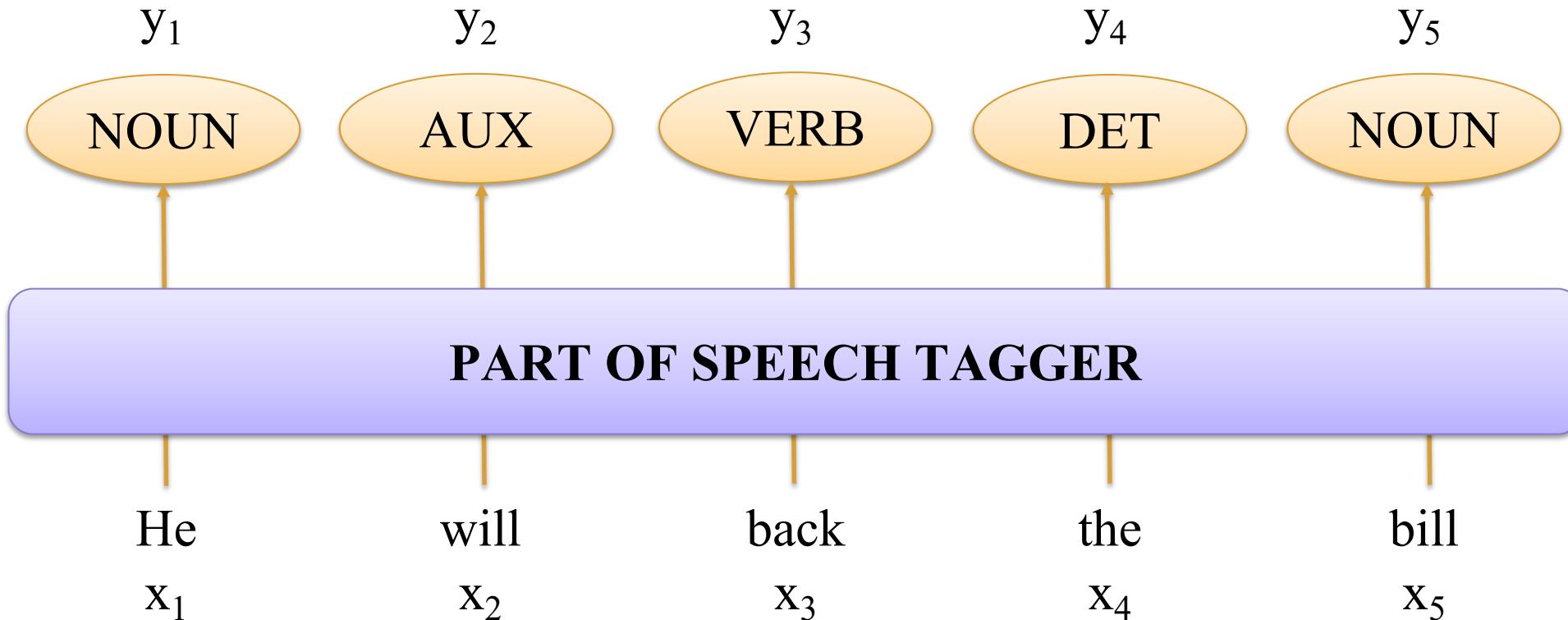
- **Introduction**
- **Part-of-Speech Tagging (POS)**
- **Named Entity Recognition (NER)**
- **Medical Named Entity Recognition**

POS Tagging



Token-Level Text Classification

- ❖ Assign a POS tag to each token in text



POS Tagging



“Universal Dependencies” Tagset

	Tag	Description	Example
Open Class	ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, first, second</i>
	PART	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
	SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
Other	PUNCT	Punctuation	<i>, , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

POS Tagging



“Penn Treebank” Tagset

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

POS Tagging



Penn Tree Bank Dataset

❖ Samples: 3914

A	lorillard	spokewoman	said	,	this	is	an	old	story	.
---	-----------	------------	------	---	------	----	----	-----	-------	---



DT	NNP	NN	VBD	,	DT	VBZ	DT	JJ	NN	.
----	-----	----	-----	---	----	-----	----	----	----	---

There	is	no	asbestos	in	our	products	now	.
-------	----	----	----------	----	-----	----------	-----	---

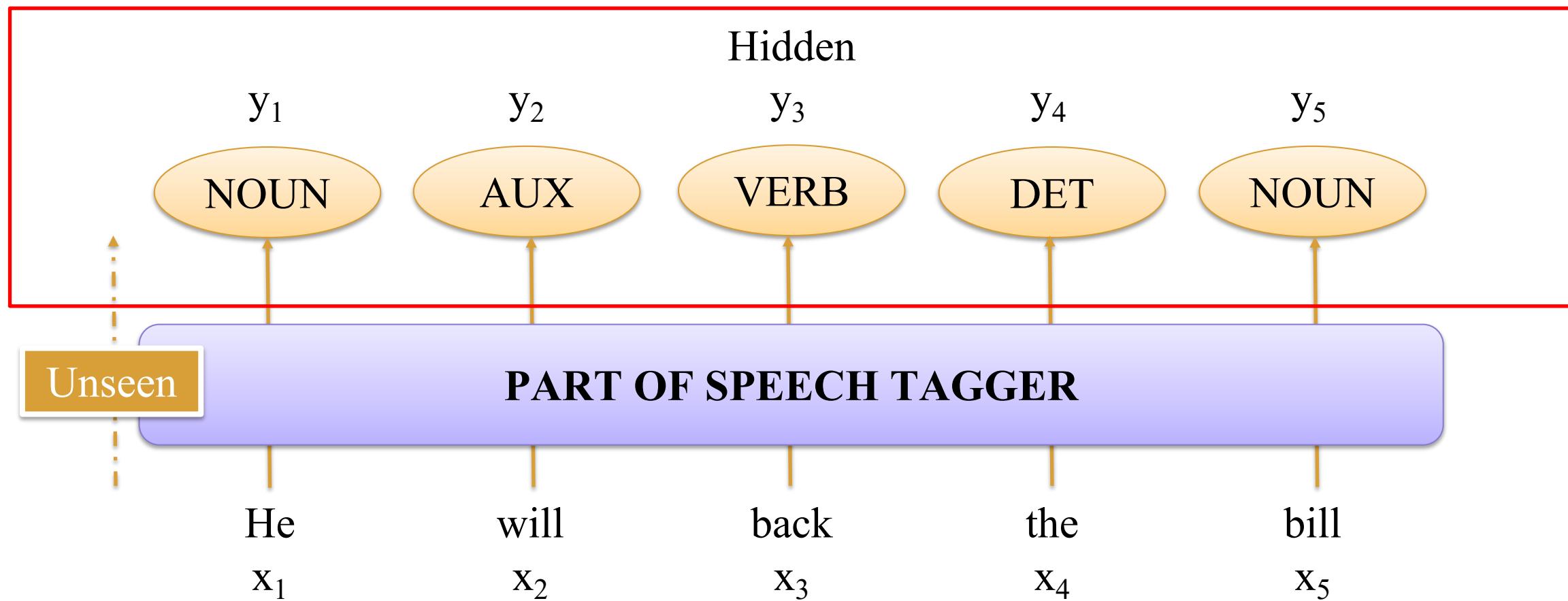


EX	VBZ	DT	NN	IN	PRP\$	NNS	RB	.
----	-----	----	----	----	-------	-----	----	---

POS Tagging



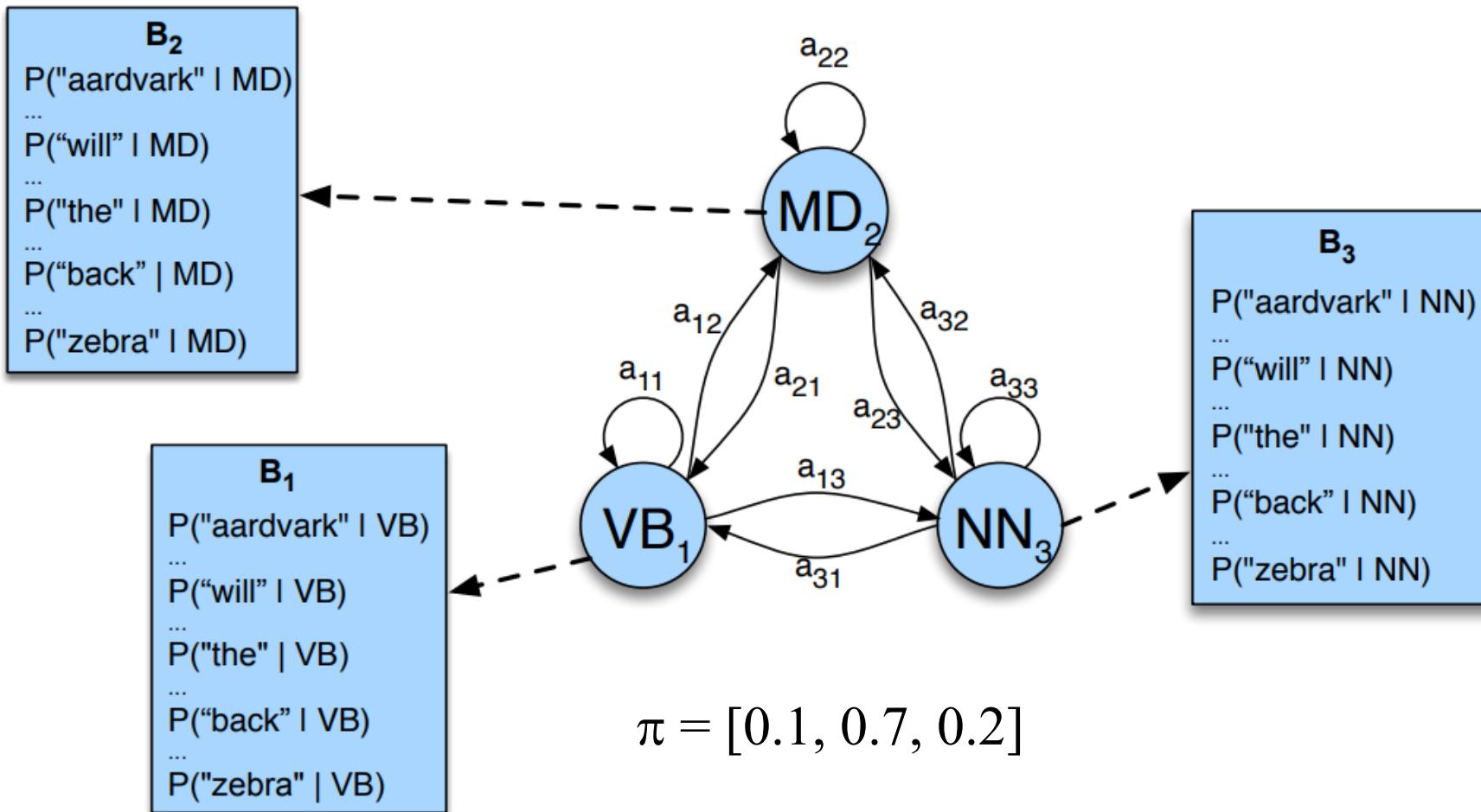
POS Tagging using Hidden Markov Model



POS Tagging



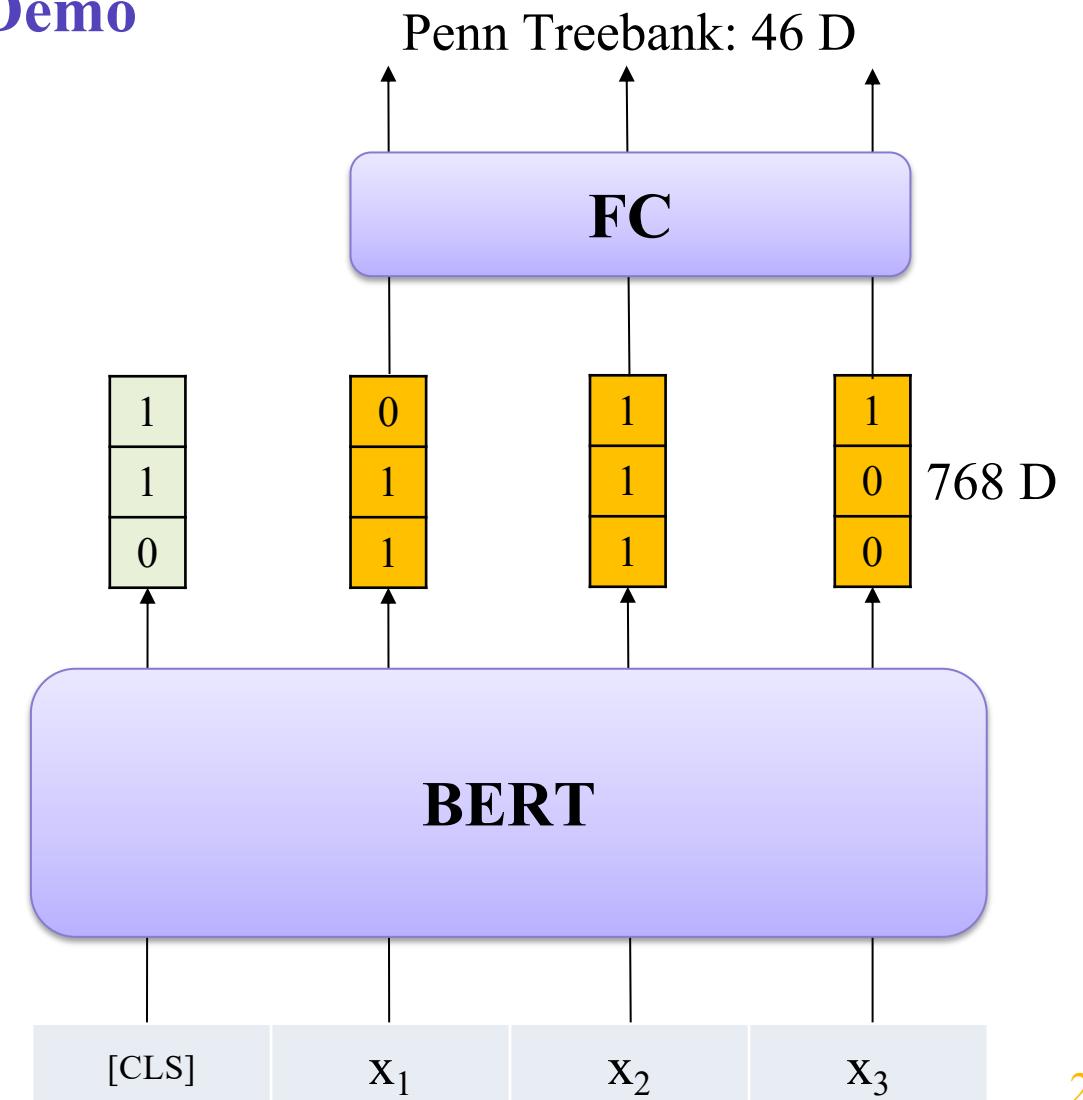
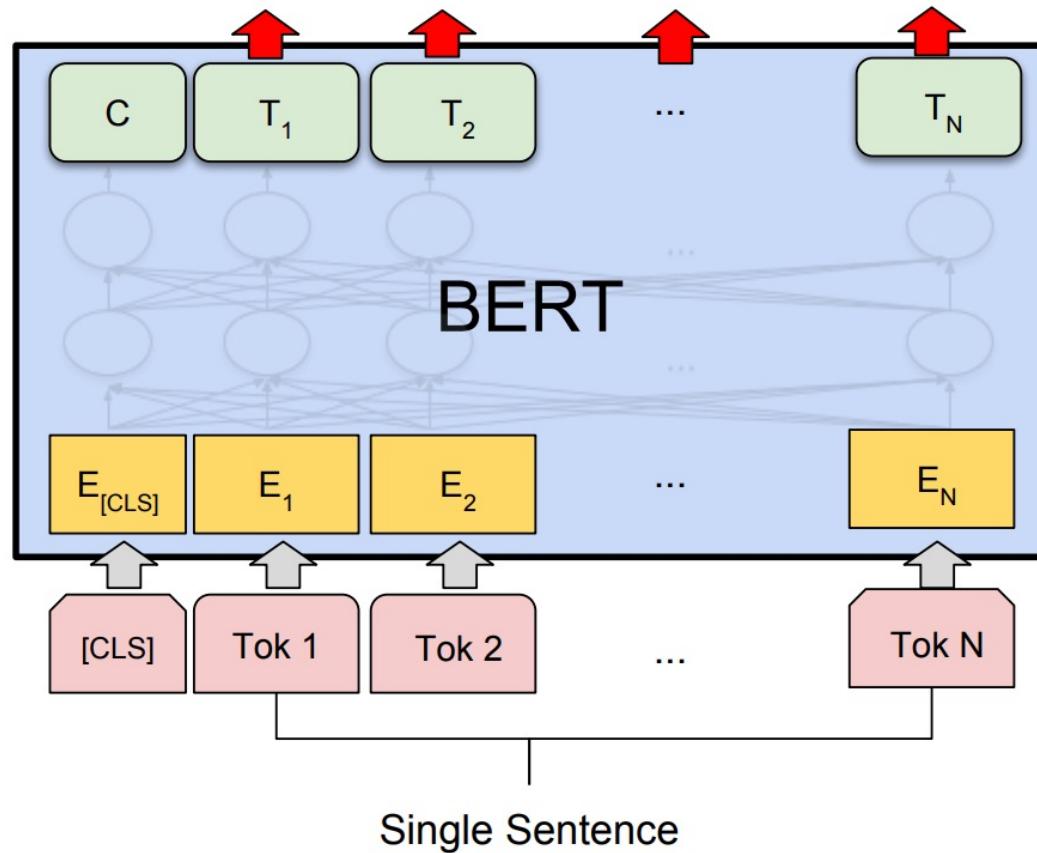
POS Tagging using Hidden Markov Model



POS Tagging



POS Tagging using Pre-trained Model - Demo



POS Tagging



POS Tagging using Spacy

I go to school.

text	POS	TAG	POS explained
= = =	= = =	= = =	= = =
I	PRON	PRP	pronoun
go	VERB	VBP	verb
to	ADP	IN	adposition
school	NOUN	NN	noun
.	PUNCT	.	punctuation

He will back the bill.

text	POS	TAG	POS explained
= = =	= = =	= = =	= = =
He	PRON	PRP	pronoun
will	AUX	MD	auxiliary
back	VERB	VB	verb
the	DET	DT	determiner
bill	NOUN	NN	noun
.	PUNCT	.	punctuation

Outline

- **Introduction**
- **Part-of-Speech Tagging (POS)**
- **Named Entity Recognition (NER)**
- **Medical Named Entity Recognition**

Named Entity Recognition



Example

In fact, the Chinese NORP market has the three CARDINAL most influential names of the retail and tech space – Alibaba GPE , Baidu ORG , and Tencent PERSON (collectively touted as BAT ORG), and is betting big in the global AI GPE in retail industry space . The three CARDINAL giants which are claimed to have a cut-throat competition with the U.S. GPE (in terms of resources and capital) are positioning themselves to become the ‘future AI PERSON platforms’. The trio is also expanding in other Asian NORP countries and investing heavily in the U.S. GPE based AI GPE startups to leverage the power of AI GPE . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one CARDINAL , with an anticipated CAGR PERSON of 45% PERCENT over 2018 - 2024 DATE .

To further elaborate on the geographical trends, North America LOC has procured more than 50% PERCENT of the global share in 2017 DATE and has been leading the regional landscape of AI GPE in the retail market. The U.S. GPE has a significant credit in the regional trends with over 65% PERCENT of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google ORG , IBM ORG , and Microsoft ORG .

Named Entity Recognition



Named Entity Recognition (NER) – Information Extraction

- ❖ Mapping span of text into entity tag
- ❖ Four entity tags are most common

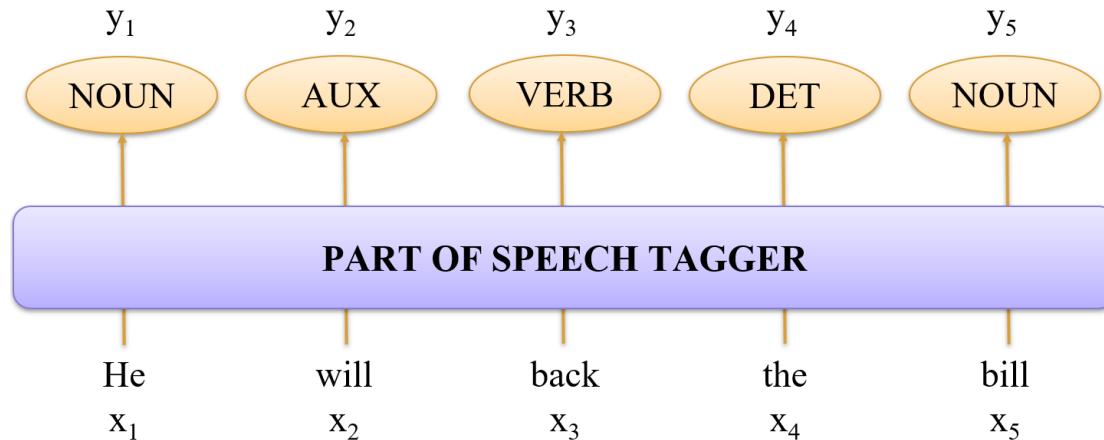
Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	Mt. Sanitas is in Sunshine Canyon.
Geo-Political Entity	GPE	countries, states	Palo Alto is raising the fees for parking.

- ❖ Other entity: Date, Time, Price,...

Named Entity Recognition



Span Recognition Problem



How to segment ?

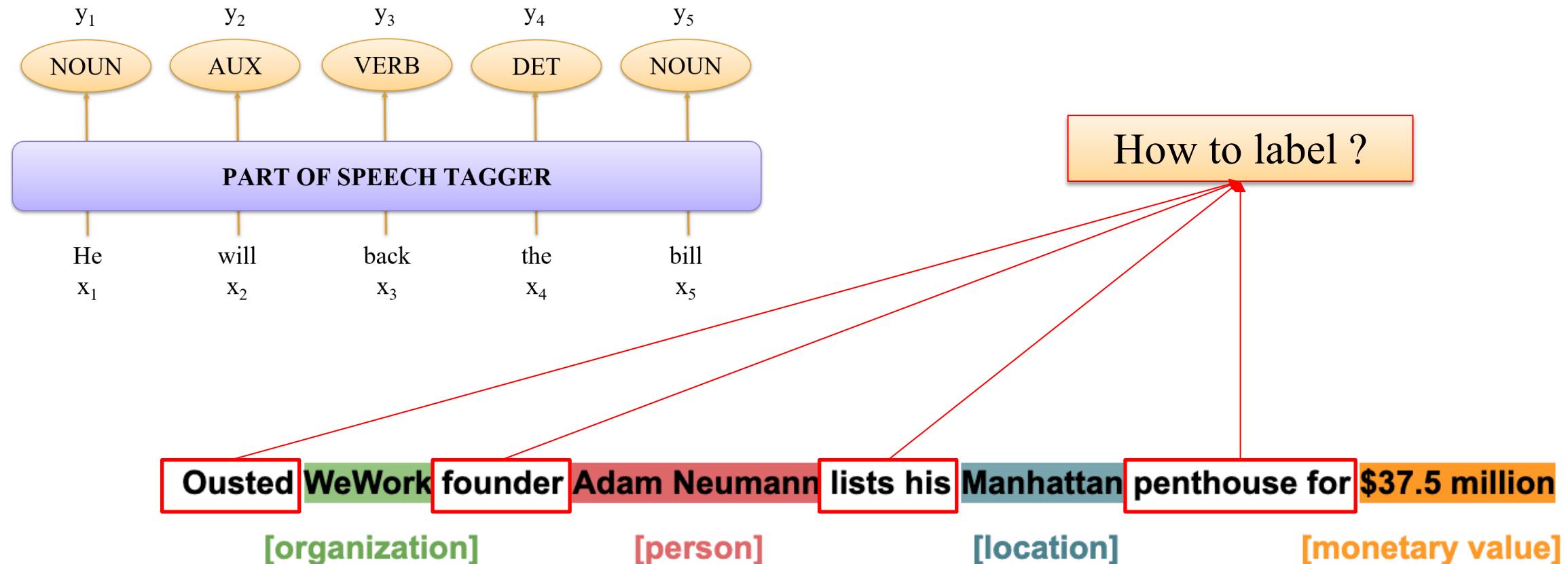
Ousted WeWork founder Adam Neumann lists his Manhattan penthouse for \$37.5 million

[organization] [person] [location] [monetary value]

Named Entity Recognition



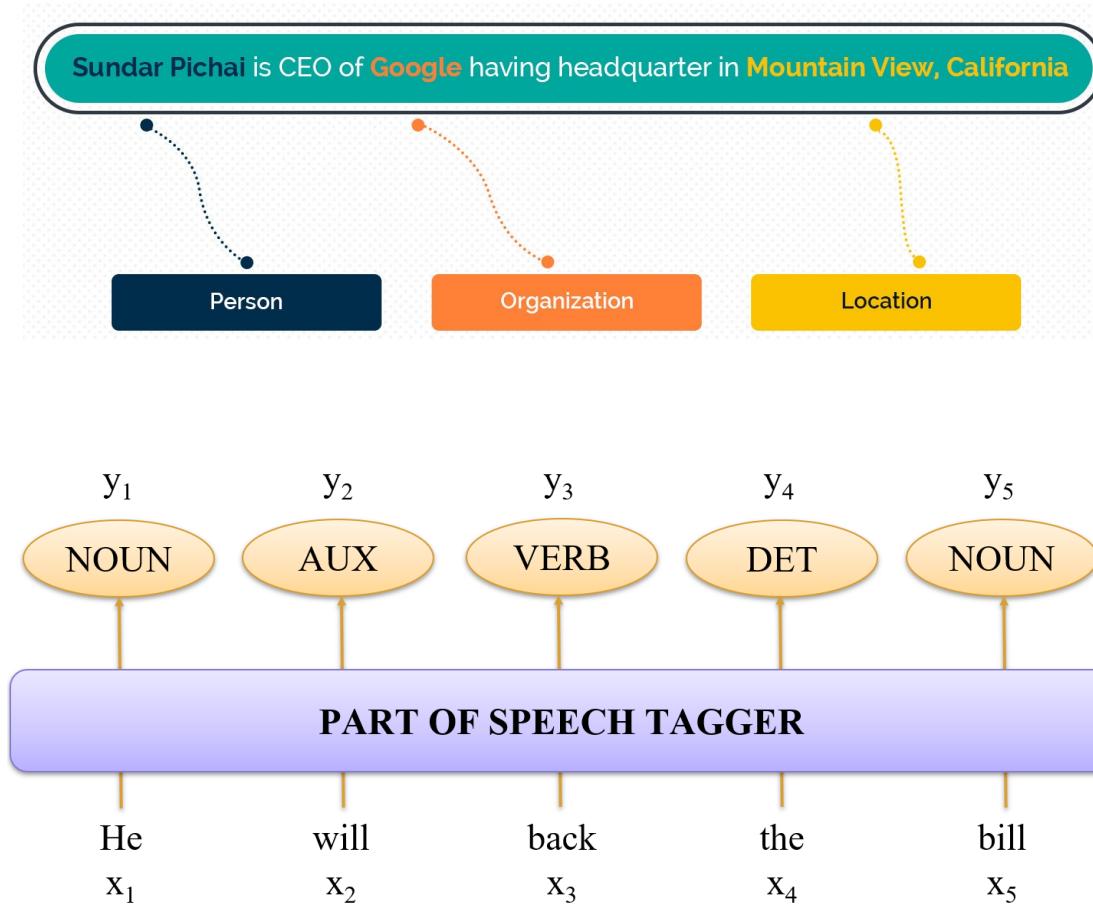
Span Recognition Problem



Named Entity Recognition



BIO Tagger



Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

Named Entity Recognition



BIO Tagger

B: token that *begins* a span

I: tokens *inside* a span

O: tokens outside of any span

of tags (where n is #entity types):

1 O tag,

n B tags,

n I tags

=> $2n+1$ tags

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

Named Entity Recognition



BIO Tagger Variants: IO, BIOES

Words	IO Label	BIO Label	BIOES Label
Jane	I-PER	B-PER	B-PER
Villanueva	I-PER	I-PER	E-PER
of	O	O	O
United	I-ORG	B-ORG	B-ORG
Airlines	I-ORG	I-ORG	I-ORG
Holding	I-ORG	I-ORG	E-ORG
discussed	O	O	O
the	O	O	O
Chicago	I-LOC	B-LOC	S-LOC
route	O	O	O
.	O	O	O

Named Entity Recognition



BIO Tagger Variants: IO, BIOES

I-Pr	I-Pr	O	O	O	I-Or	O	O	O	I-Lo	I-Lo	I-Lo
------	------	---	---	---	------	---	---	---	------	------	------

Sundar Pichai is CEO of Google having headquarter in Mountain View, California

B-Pr	I-Pr	O	O	O	B-Or	O	O	O	B-Lo	I-Lo	I-Lo
------	------	---	---	---	------	---	---	---	------	------	------

Sundar Pichai is CEO of Google having headquarter in Mountain View, California

B-Pr	E-Pr	O	O	O	B-Or	O	O	O	B-Lo	I-Lo	E-Lo
------	------	---	---	---	------	---	---	---	------	------	------

Sundar Pichai is CEO of Google having headquarter in Mountain View, California

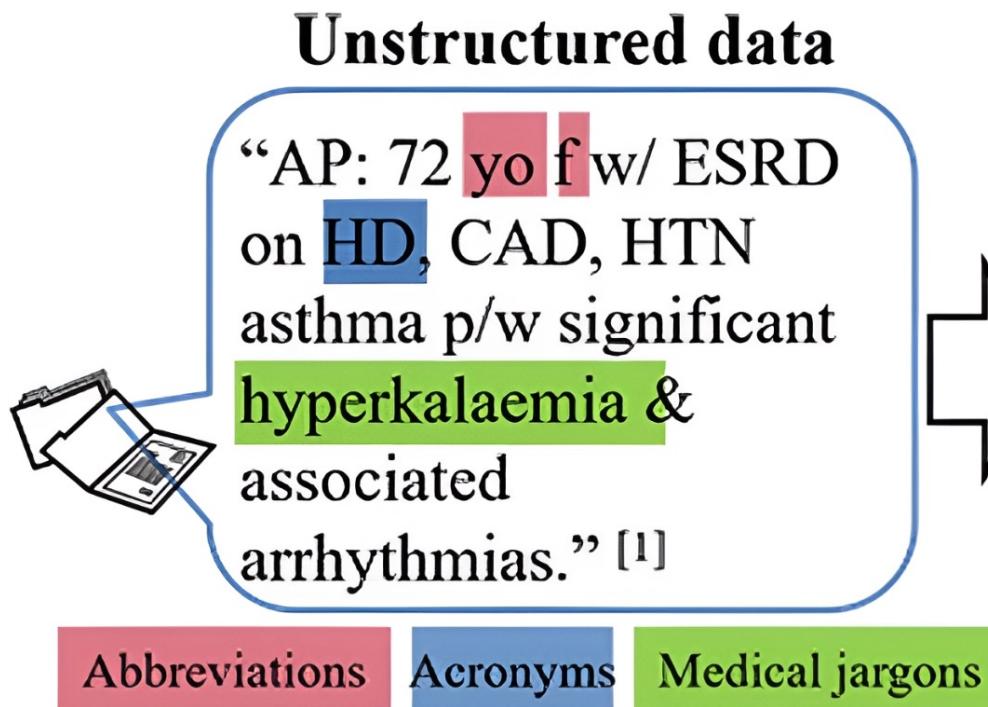
Outline

- **Introduction**
- **Part-of-Speech Tagging (POS)**
- **Named Entity Recognition (NER)**
- **Medical Named Entity Recognition**

Medical NER



NER Application for Medical Information Extraction



Meaningful information

Attributes	Medical text
Age	72 years old (yo)
Gender	female (f)
Condition	hypertensive disease (HD)
Symptom	hyperkalemia

Medical NER

!

MACCROBAT2018

15939911.ann

15939911.txt

16778410.ann

16778410.txt

17803823.ann

17803823.txt

18236639.ann

18236639.txt

Sony formed a joint venture with Ericsson, a mobile phone company based in Sweden.
Sony announced today that ...

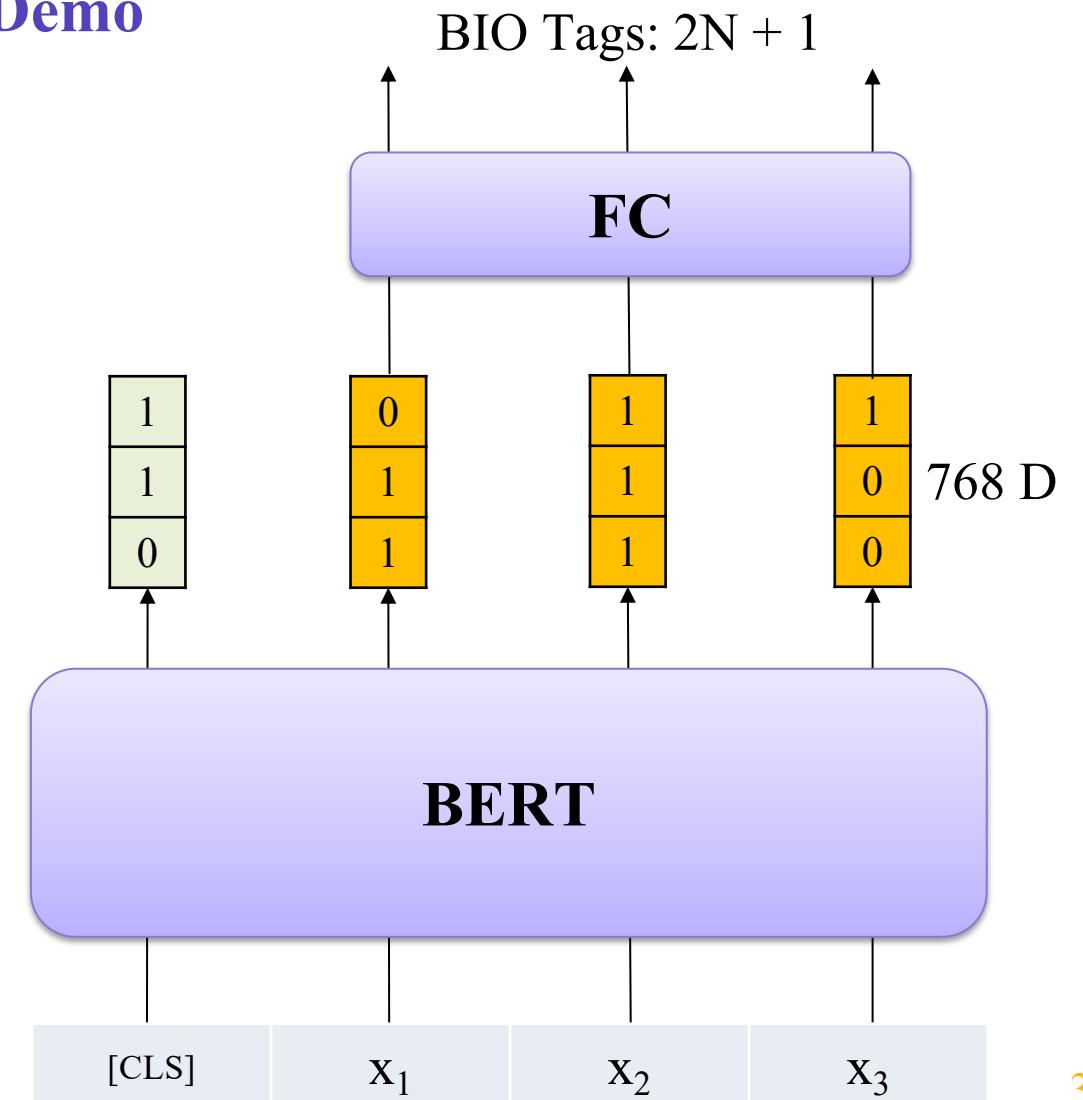
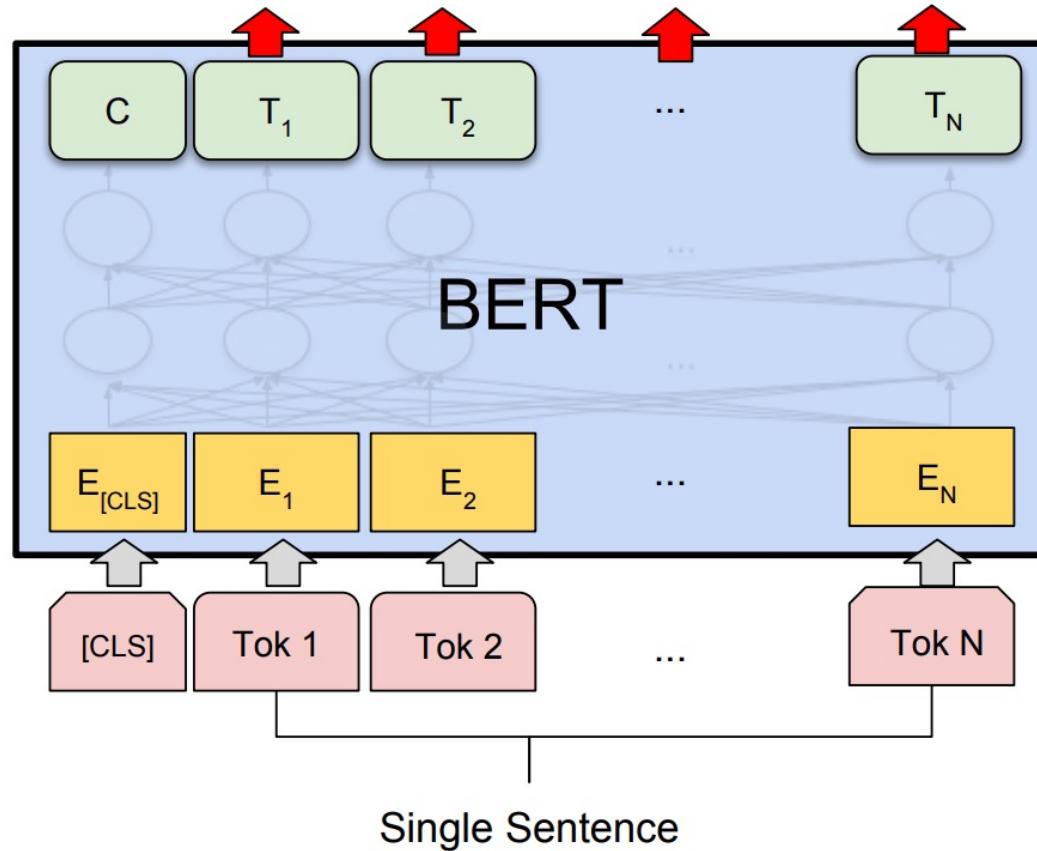
T1	Organization	0 4	Sony
T2	MERGE-ORG	14 27	joint venture
T3	Organization	33 41	Ericsson
E1	MERGE-ORG:T2	Org1:T1 Org2:T3	
T4	Country	75 81	Sweden
R1	Origin	Arg1:T3 Arg2:T4	

Standoff Format

Medical NER



Medical NER using Pre-trained Model - Demo





Thanks!

Any questions?