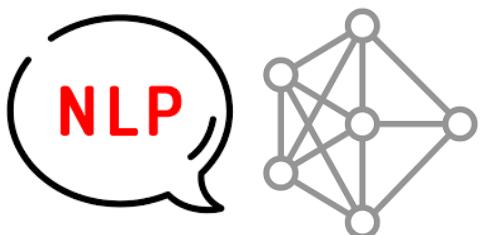




CHULA **ENGINEERING**
Foundation toward Innovation

COMPUTER



211594: Introduction to Natural Language Processing (NLP)

Peerapon Vateekul & Ekapol Chuangsawanich
Department of Computer Engineering,
Faculty of Engineering, Chulalongkorn University



Peerapon Vateekul, Ph.D.



Ekapol Chuangsuwanich, Ph.D.



Can Udomcharoenchaikit
P'Can (TA)



Nattachai Tretasayuth
P'Boss (TA)



Outlines

- What is NLP?
 - Definition
 - Levels of understanding in NLP
 - NLP today
- Why NLP is hard?
 - Ambiguity
 - Issues in Thai NLP
- NLP & Text mining
 - NLP pipeline (English & Thai)
 - History of NLP techniques
 - Case Study: Determiner placement
- Deep Learning
 - Definition
 - Reasons for exploring Deep Learning
- NLP Tools
- Course Logistics
- Google Cloud Demo

+

What is NLP?

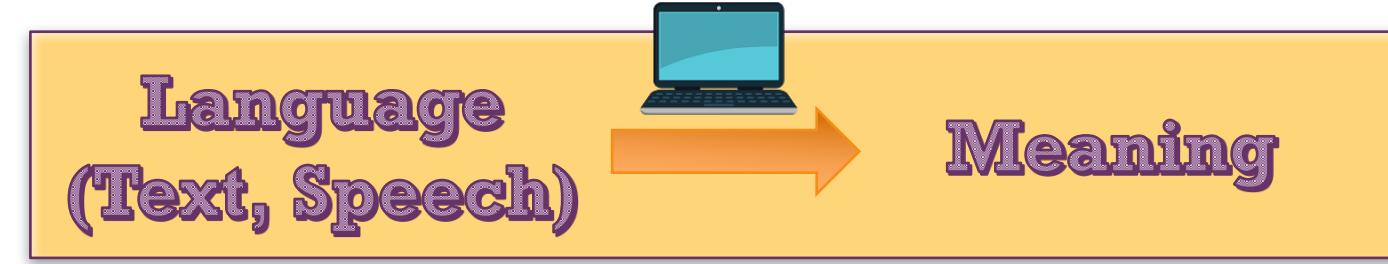


Natural Language Processing (NLP)

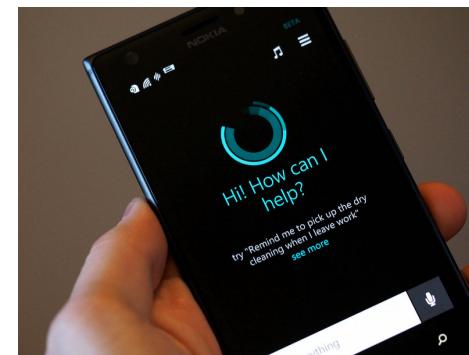
- Subfield of AI

- GOAL:

- Bridge the gap between **how people communicate** and **what machines understand** in order to perform useful tasks, e.g.
- Making appointments, buying things, question answering, etc.



NLP (Interpretation)
I want to eat Japanese food.



May I order Yayoi for you?

AI/ML (Generation)



NLP Goal

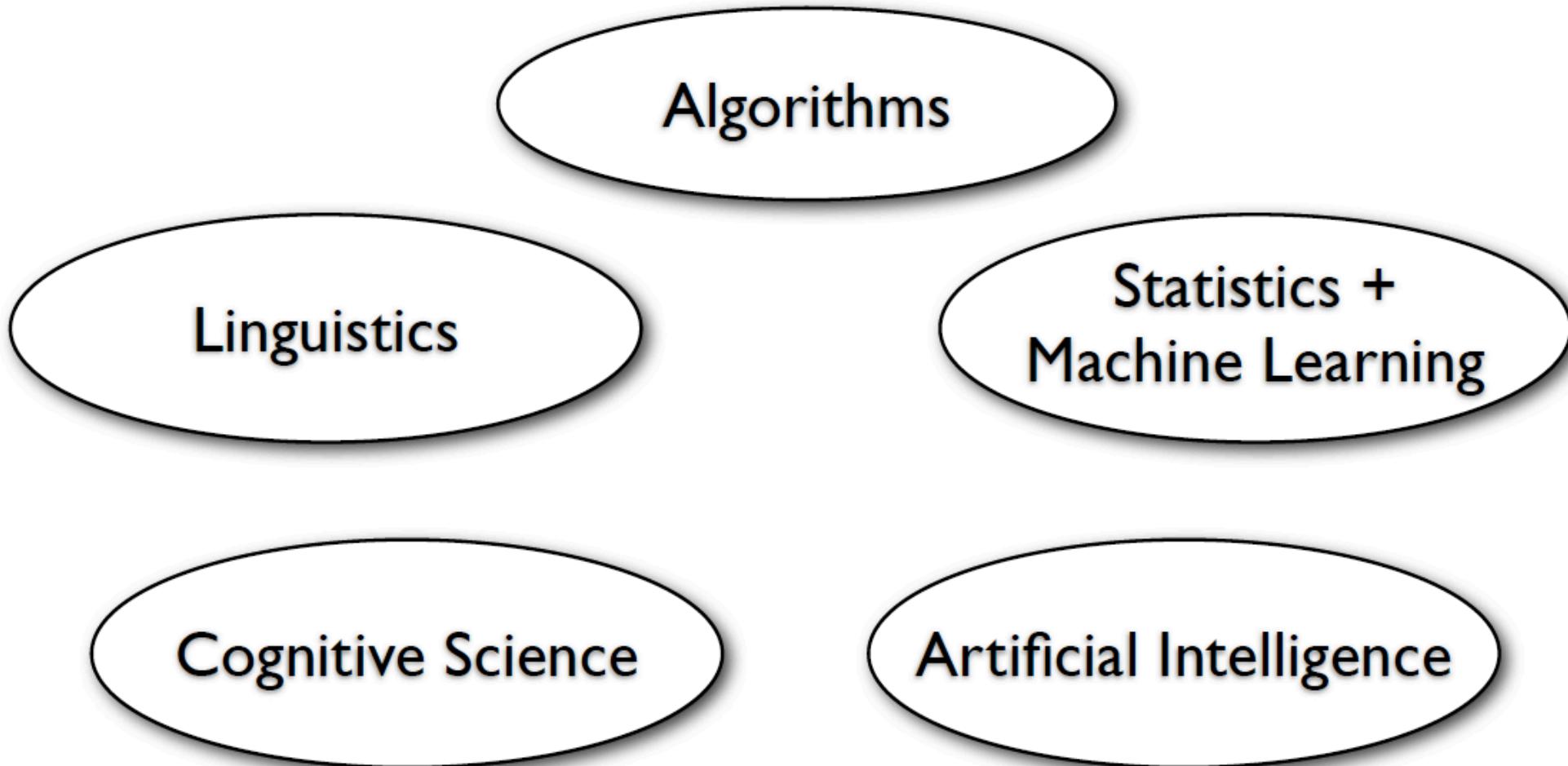
Goal: intelligent processing of human languages

- Not just string matching





NLP is interdisciplinary





Level of understanding in NLP

https://www.tutorialspoint.com/artificial_intelligence/artificial_intelligence_natural_language_processing.htm

Lexical Analysis:

Text → Paragraphs, Sentences, and Words

Syntactic Analysis (Parsing):

Grammar/Relationship between words

Semantic Analysis:

Exact meaning of the sentence

Discourse Integration:

Meaning of the sentence **based on the previous sentence
(pronouns)**

Pragmatic Analysis:

Actual Meaning based on **the context** and real-world knowledge

Discourse

Semantics

Syntax: Constituents

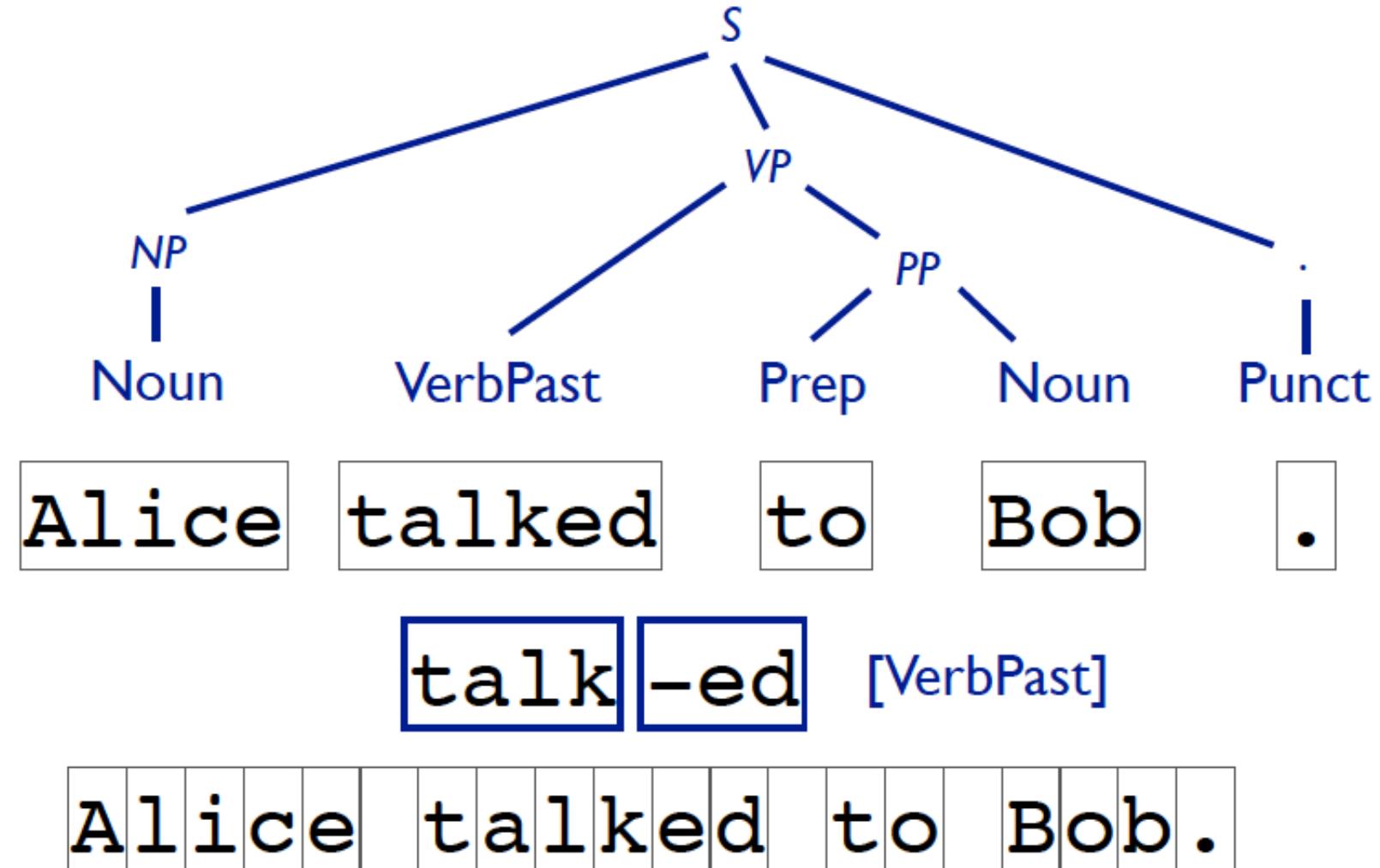
Syntax: Part of Speech

Words

Morphology

Characters

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)
SpeakerContext(s)
TemporalBefore(e, s)





NLP today: Technology



Dan Jurafsky

Language Technology

making good progress

mostly solved

Spam detection

Let's go to Agra!



Buy V1AGRA ...



Part-of-speech (POS) tagging

ADJ ADJ NOUN VERB ADV

Colorless green ideas sleep furiously.

Named entity recognition (NER)

PERSON ORG LOC

Einstein met with UN officials in Princeton

Sentiment analysis

Best roast chicken in San Francisco!



The waiter ignored us for 20 minutes.



Coreference resolution

Carter told Mubarak he shouldn't run again.

Word sense disambiguation (WSD)

I need new batteries for my **mouse**.



Parsing



I can see Alcatraz from the window!

Machine translation (MT)

第13届上海国际电影节开幕...



The 13th Shanghai International Film Festival...

Information extraction (IE)

You're invited to our dinner party, Friday May 27 at 8:30



May 27
add

still really hard

Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

Summarization

The Dow Jones is up



The S&P500 jumped

Housing prices rose

Economy is good

Dialog

Where is Citizen Kane playing in SF?



Castro Theatre at 7:30. Do you want a ticket?



NLP today: Machine Translation (MT)

Google google translate

All Images Maps News Videos More Settings Tools

About 1,180,000,000 results (0.39 seconds)

English ▾

As the new year gets underway, expert commentators give their view on what 2018 holds in store.

Here are three big themes to watch out for over the next 12 months.

Can the stock market rally go on? The new year has begun with stock markets in the UK and US hitting new record highs.

The Dow Jones Industrial Average rose above 25,000 points for the first time this week, while the broader S&P 500 is also at historic highs.

Thai ▾

เป็นปีใหม่ที่กำลังได้รับการแสดงความคิดเห็นของผู้เชี่ยวชาญให้มุ่งมองของพวกรเข้าเกี่ยวกับสิ่งที่ 2018 เก็บไว้ในร้าน

ต่อไปนี้เป็นหัวข้อใหญ่สามข้อที่ควรระวังในช่วง 12 เดือนข้างหน้า

การซัมมูมตลาดหุ้นสามารถดำเนินต่อไปได้หรือไม่?

ปีใหม่เริ่มมีตลาดหุ้นในสหรัฐอาณาจักรและสหราชอาณาจักรสูงเป็นประวัติการณ์

ดัชนีเฉลี่ยอุตสาหกรรมดาว 琼斯ปรับตัวสูงขึ้นกว่า 25,000 จุดเป็นครั้งแรกในสัปดาห์นี้ขณะที่ดัชนี S & P 500 ที่ใหญ่ขึ้นก็อยู่ในระดับสูงเป็นประวัติการณ์

Markets, Brexit and Bitcoin: 2018's themes

By Chris Johnston
Business reporter

5 January 2018

f t m Share



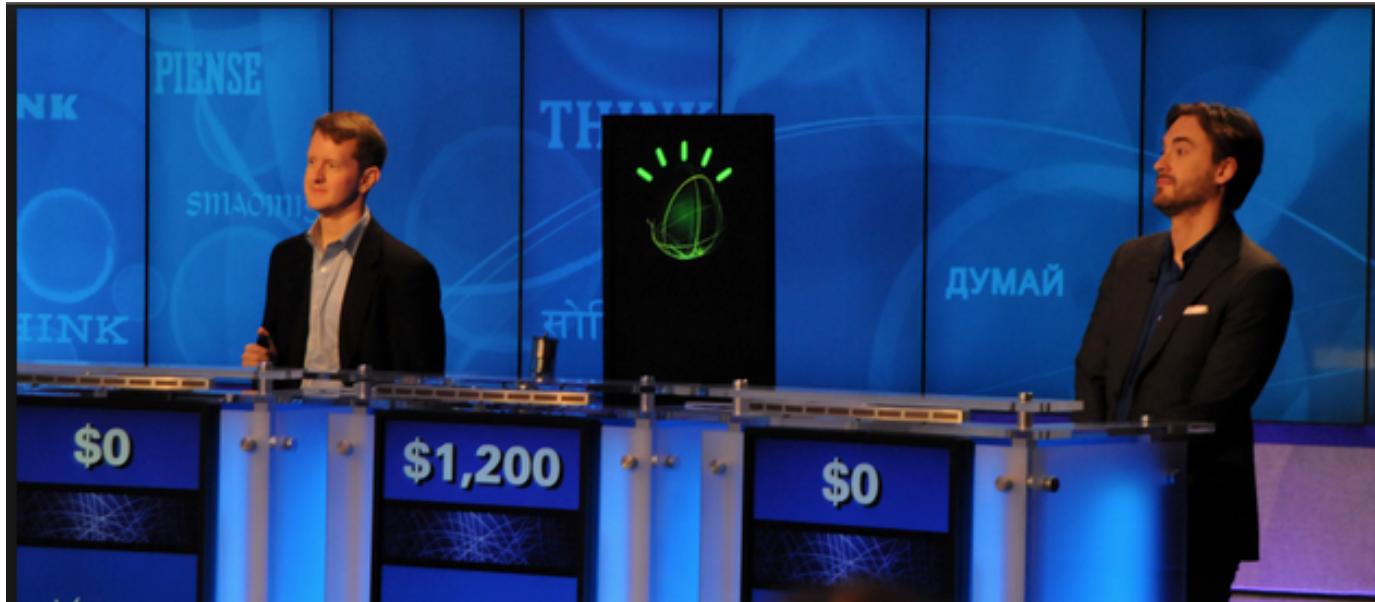
GETTY IMAGES

As the new year gets underway, expert commentators give their view on what 2018 holds in store.

<http://www.bbc.com/news/business-42581934>



NLP today: Question Answering (QA)



IBM Watson wowed the tech industry and a corner of U.S. pop culture with its 2011 win against two of Jeopardy's greatest champions. Here's how IBM pulled it off and a look at what Watson's real career is going to be.

<https://www.techrepublic.com/article/ibm-watson-the-inside-story-of-how-the-jeopardy-winning-supercomputer-was-born-and-what-it-wants-to-do-next/>



NLP today: Question Answering (QA) (cont.)

ปริศนาฟ้าแลน | ตั้ก, พายี่เก', โรเบิร์ต, เมล | 29 พ.ย. 60 Full HD



<https://www.youtube.com/watch?v=CAJAQUao7HU>



NLP today: Search/Summarization

Google wonder woman

All Images Videos News Maps More Settings Tools

About 91,900,000 results (0.78 seconds)

Top stories



Wonder Woman Scores a Huge Best Picture Nomination

IGN.com
1 day ago



Gal Gadot addresses James Cameron's controversial Wonder Woman comments

The Independent
23 hours ago



Gal Gadot Just Responded To James Cameron's "Wonder Woman" Comments I...

BuzzFeed
18 hours ago

→ More for wonder woman

[Wonder Woman \(2017 film\) - Wikipedia](#)
[https://en.wikipedia.org/wiki/Wonder_Woman_\(2017_film\)](https://en.wikipedia.org/wiki/Wonder_Woman_(2017_film)) ▾

Wonder Woman is a 2017 American superhero film based on the DC Comics character of the same name, distributed by Warner Bros. Pictures. It is the fourth installment in the DC Extended Universe (DCEU). The film is directed by Patty Jenkins, with a screenplay by Allan Heinberg, from a story by Heinberg, Zack Snyder, ...

Gal Gadot · Patty Jenkins · Elena Anaya · Doctor Poison

Wonder Woman 

2017 · Fantasy/Science fiction film · 2h 21m

[Play trailer on YouTube](#)

7.6/10 · IMDb

90% liked this film  

Before she was Wonder Woman (Gal Gadot), she was Diana, princess of the Amazons, trained to be an unconquerable warrior. Raised on a sheltered island paradise, Diana meets an American pilot (Chris Pine) who tells her about the massive conflict that's raging in the outside world. Convinced that she c... [MORE](#) ▾

Initial release: May 15, 2017 (Shanghai)
Director: Patty Jenkins



NLP today: Information Extraction (IE)

Data science perspective on clinical research



ID	AGE	RACE	STUDY	PROC	BIRTHS	MA_AGE	ASSESS	DENSITY	FINDING	FINDING T
9527	78	2	6/12/06	BIDXU-L	0	P		3	CALCS	N
32875	56	1	7/11/06	BIDXB-B	0	N		3		
2247	72	1	4/12/06	BIDXU-R	0	N		3		
45521	61	1	3/30/06	BIDXB-B	0	B		3	CALCS	S
48987	41	1	4/5/06	BIDXB-B	0	P		3	CALCS	N
4179	67	1	5/12/06	BIDXB-B	0	P		2	CALCS	N
26300	59	1	3/31/06	BIDXU-L	0	N		3		
67960	64	1	4/7/06	BIDXB-R	0	P		3	MASS	O
43283	61 W		7/21/06	BIDXB-B	0	B		3		
43319	51	1	4/7/06	BIDXB-B	0	N		3		

Abstract clinical records into a database

Pathology Report: REMOVED_ACCESSION_ID
ACCESSIONED ON: REMOVED_DATE
CLINICAL DATA: Carcinoma **right breast**.
*** FINAL DIAGNOSIS ***
LYMPH NODE (SENTINEL), EXCISION
(REMOVED_CASE_ID): METASTATIC
CARCINOMA IN 1 OF 1 LYMPH NODE.
NOTE: The metastatic deposit spans 0.19cm and
is identified on H&E and cytokeratin immunostains.
A second cytokeratin-positive but cauterized focus
likely also represents metastatic tumor (<0.1cm).
There is **no evidence of extranodal extension**.
BREAST (RIGHT), EXCISIONAL BIOPSY
(REMOVED_ACCESSION_ID :
REMOVED_CASE_ID -B): **INVASIVE DUCTAL
CARCINOMA (SEE TABLE #1).** DUCTAL
CARCINOMA IN-SITU, GRADE 1. ATYPICAL
DUCTAL HYPERPLASIA. LOBULAR NEOPLASIA
(ATYPICAL LOBULAR HYPERPLASIA).
TABLE OF PATHOLOGICAL FINDINGS #1

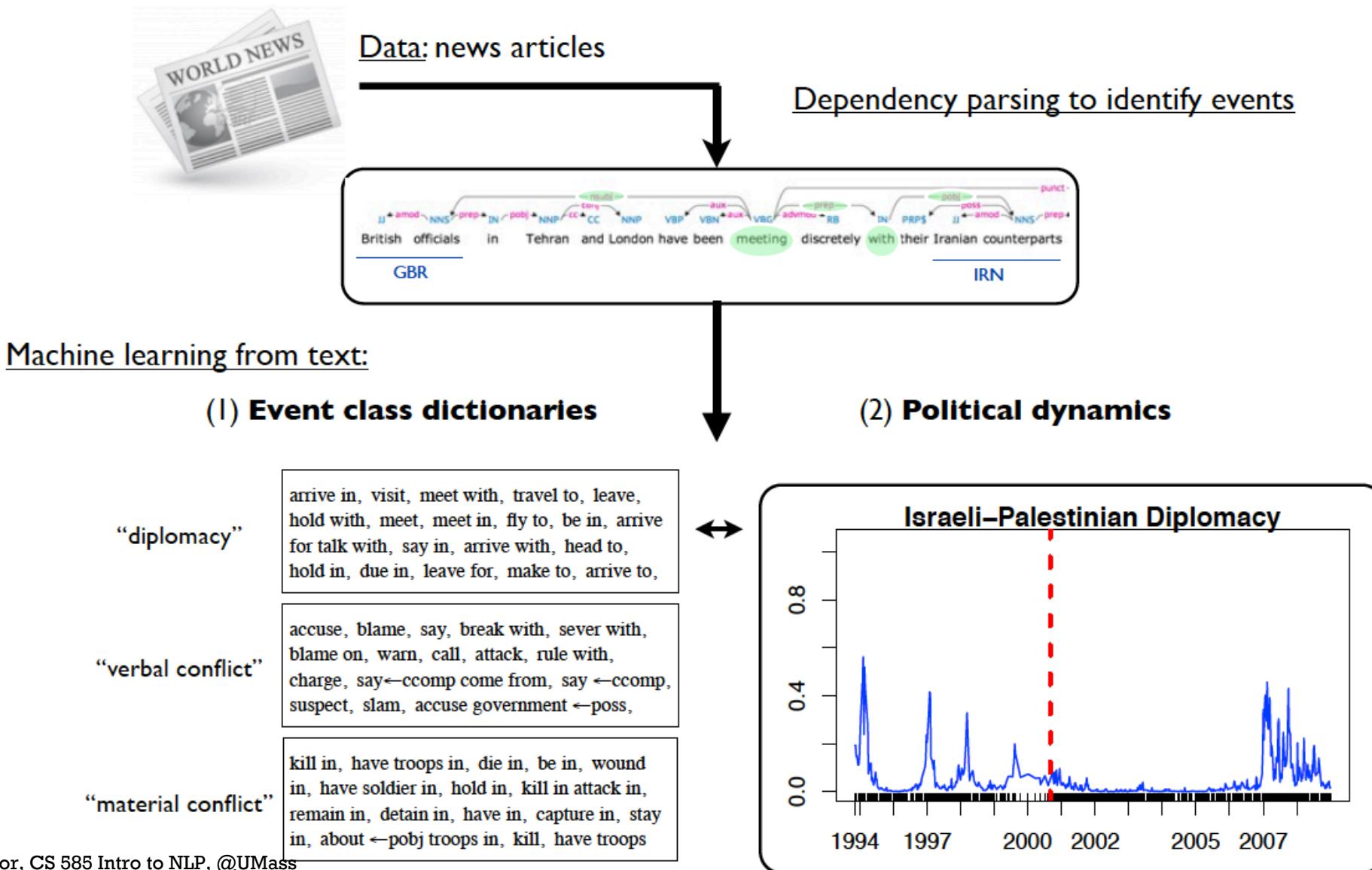


Name	Extraction
Breast Side	Right
Ductal Carcinoma in Situ	Present
Invasive Lobular Carcinoma	Absent
Invasive Ductal Carcinoma	Present
Cancer	Present
Lobular Carcinoma in Situ	Absent
Atypical Ductal Hyperplasia	Present
Atypical Lobular Hyperplasia	Present
Lobular Neoplasia	Present
Flat Epithelial Atypia	Absent
Blunt Adenosis	Absent
Atypia	Present
Positive Lymph Nodes	Present
Extracapsular Axillary Nodal Extension	Absent
Isolated Cancer Cells in Lymph Nodes	Absent
Lymphovascular Invasion	Absent
Blood Vessel Invasion	Absent
Estrogen Receptor Status	Positive
Progesterone Receptor Status	Positive
HER 2 (FISH) Status	Unknown

Parsing pathology reports into database



NLP today: Trend analysis



Hathaway Phenomenon



A couple weeks ago, Huffington Post blogger Dan Mirvish noted a funny trend: when Anne Hathaway was in the news, Warren Buffett's Berkshire Hathaway's shares went up. He pointed to [six dates going back to 2008](#) to show the correlation. Mirvish then suggested a mechanism to explain the trend: "automated, robotic trading programming are picking up the same chatter on the Internet about 'Hathaway' as the IMDb's StarMeter, and they're applying it to the stock market."

BERKSHIRE HATHAWAY INC.

3555 Farnam Street
Omaha, NE 68131
Official Home Page

- [A Message From Warren E. Buffett](#)
- [Annual & Interim Reports](#)
Updated November 3, 2017
- [Special Letters From Warren & Charlie RE:Past, Present and Future](#)
- [Link to SEC Filings](#)
- [Links to Berkshire Subsidiary Companies](#)
- [Corporate Governance](#)
- [Owner's Manual](#)
- [Letters from Warren E. Buffett Regarding Pledges to Make Gifts of Berkshire Stock](#)
- [News Releases](#)
Updated November 3, 2017
- [Warren Buffett's Letters to Berkshire Shareholders](#)
Updated February 25, 2017
- [Charlie Munger's Letters to Wesco Shareholders](#)
- [Annual Meeting Information](#)
- [Celebrating 50 Years of a Profitable Partnership](#)
(A commemorative book first sold at the 2015 Annual Meeting and now for sale on eBay.)
- [Comparative Rights and Relative Prices of Class A and B](#)
- [Berkshire Activewear](#)

GEICO
FOR A FREE CAR INSURANCE RATE QUOTE THAT COULD SAVE YOU SUBSTANTIAL MONEY
WWW.GEICO.COM OR CALL 1-800-395-6349, 24 HOURS A DAY



NLP is hard!
Word-level ambiguity

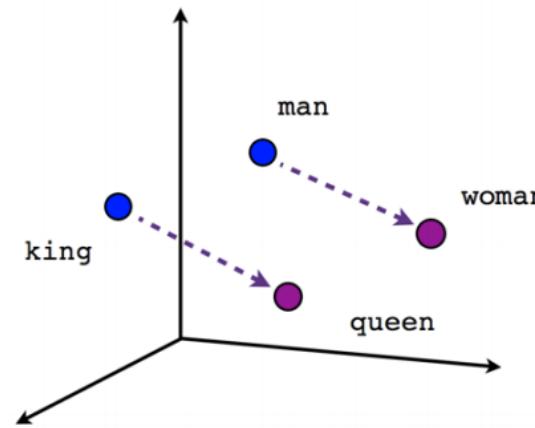
+

Why NLP is hard?

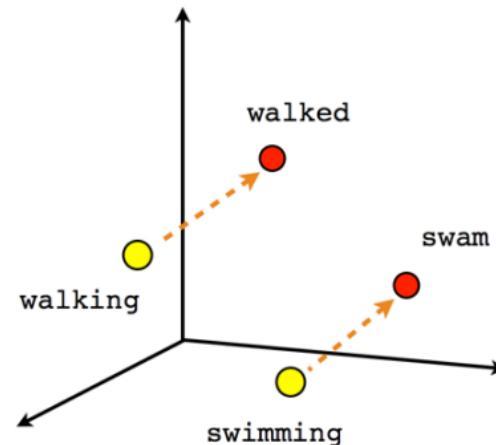


Why NLP is hard?

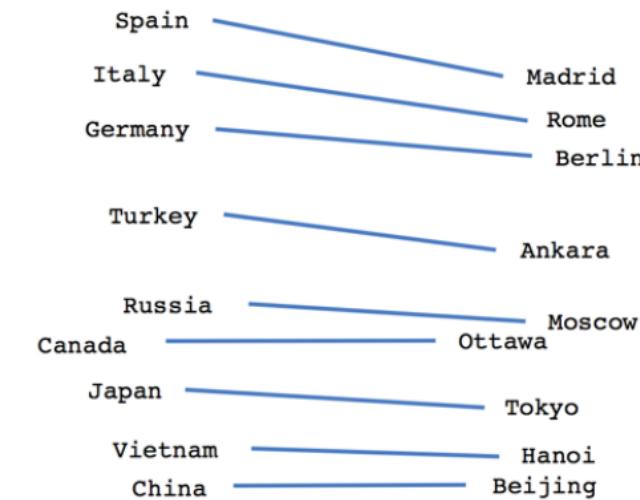
- Complexity in **representing**, learning and using linguistic/situational/world/visual knowledge



Male-Female



Verb tense



Country-Capital

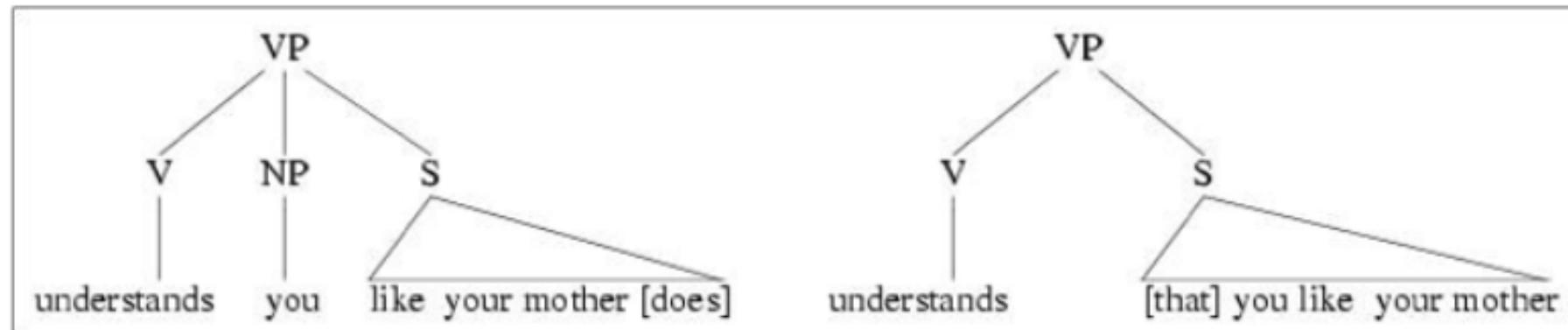


Why NLP is hard? (cont.)

- Human languages are **ambiguous** (unlike programming and other formal languages).
 - Some parts can be ignored
- Human languages interpretation depends on real world, common sense, and contextual knowledge (pragmatic analysis)

At last, a computer understands you like your mother”

Ambiguity at syntactic level: Different structures lead to different interpretations



The Pope's baby steps on gays. [Ref: Prof. Christopher Manning, CS224N/Ling284, 2017]



Issues in Thai NLP

■ Word segmentation

■ No word delimiters

- ฉัน|นำ|ดอก|ไม้|ไป|ให้ว|ศาล|พระ|ภูมิ|ที่|โรง|เรียน|ประจำ|
- ฉัน|นำ|ดอก|ไม้|ไป|ให้ว|ศาล|พระ|ภูมิ|ที่|โรง|เรียน|ประจำ|
- ฉัน|นำ|ดอก|ไม้|ไป|ให้ว|ศาล|พระ|ภูมิ|ที่|โรง|เรียน|ประจำ|
- ฉัน|นำ|ดอก|ไม้|ไป|ให้ว|ศาล|พระ|ภูมิ|ที่|โรง|เรียน|ประจำ|

■ Sentence segmentation

■ No sentence boundary markers

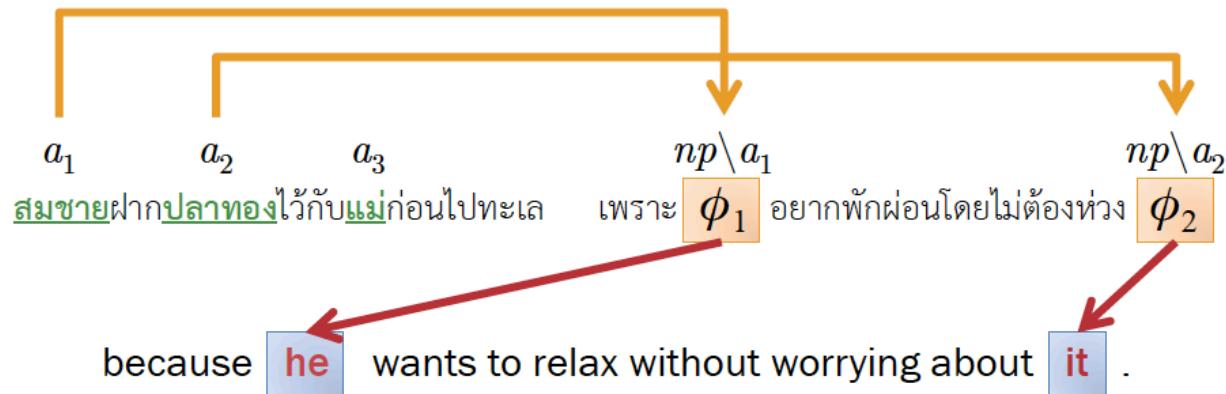
อย่างไรก็ตาม อดีตประธาน ทปอ. กล่าวว่า มีการหักหัวเรื่องนี้มาตลอดว่า มีช่วงเวลาว่างนานขนาดนี้ ทำไมถึงยังต้องมีการจัดสอบนอกเหนือจากนี้อีก เพราะการสอบล่วงล้าไปในเวลาระหว่างเรียนมั้ยมั่นนั้นกระทบกับเรื่องอื่นๆ โดยเฉพาะการเรียนในชั้นเป็นวงจรลูกโซ่ แนวโน้มที่เข้ามาแก้เรื่องนี้ เป็นความคิดที่ดี แต่ยังไม่เห็นเรื่องใช้ผลการเรียนในชั้นมาเป็นองค์ประกอบรับรอง ซึ่งอาจทำให้เด็กไม่สนใจห้องเรียน และมุ่งกวัดวิชา ทำให้การสอบเข้าอุดมศึกษา ตกเป็นจำเลยข้อหาทำลายระบบการศึกษาขั้นพื้นฐาน วนไปสู่ปัญหาเก่าๆ ได้



Issues in Thai NLP (cont.)

Syntax ambiguity

- Pronouns and some constituents can be **omitted** as long as they can be implied from the context



Nostalgic Thai slangs

เชัวร์ป้าดนี่มีซึมไปเลย เดี๋ดดวง งานาก้า
หะยะแหยง ซังกะบัวย โอคุซึ่ง เชือหัวไอเร่อง เสร็จโก^ะ.
บ่มไก ชาไปต่อย สะแಡວแห้ว อูฐไภก์ໄລກ์บอย เดดสະນອาร
ໂหลຍໂหຍ ทັງร้านราคาເກ່າໄຫວ້ ສະປະບະແວປ ຈີບຈອຍ
ຮັກສາຍັນທີ່ນອຍ່າ ແຕ່ຮັກໃຫນນາມະຄົນ ແອປເປົ້ວວອເຊີງຕັນ ຂອງແກ້ຕ້ອງນີ້ 5 ບຸ່ນ **ຈາບ**
ໄອຄຸກຮົກນິແບວັດືກ ນາຍຄົດເຫັນວັນໄຫມປີ1 ມີເຕີບນາກ.
ບອຍໄມ້ດືມຄະ ຄົກບຸ ອານຸແນະ ສູນມ.ຍທ. (ສາຍມາກ ອໍຍ່າຫົວງ) ເດືດສະຮະຕີ
ຫັນອົມແນ້ນ ເດີກຫາຣດ ຕັນ ຮດເປັນອາໄຮວະ ເຮົພວ້າ
ຕະຕັ້ງໂທນິ່ງ ໂນເວ ສະເຫັນ ຫັນແຕກຫນອໄມ່ຮັບເຢັບ ໃຫ້ຕາຍເກອະໂຮບັນ
ປະກັບໃຈຈອດ ຈ້ອຍແດກ ກີບເກົຍເຮັກ້າ ສຍື່ນກີ່ຍ
ຄຸນຫລອກດາວ



Princes Garden, Edinburgh

+

NLP & Text Mining



NLP & Text Mining

- NLP: Language → Meaning
- Text Mining
 - Text mining, which is sometimes referred to “text analytics” is one way to make qualitative or “**unstructured data**” **usable by a computer**.
 - Convert from unstructured to structured data
 - **NLP** techniques are the building blocks for text mining tasks

NBC Nightly News @nbcnightlynews
America's #1 evening news broadcast.
Tweets by @newsdel & @braddjaffy. Join us on Facebook <http://facebook.com/nbcnightlynews>

NBC News @NBCNews
A leading source of global news and information for more than 75 years. Have a news tip or question? Ask @rozzy, @lou_dubois, @jbaiata or @anthonyquintano.

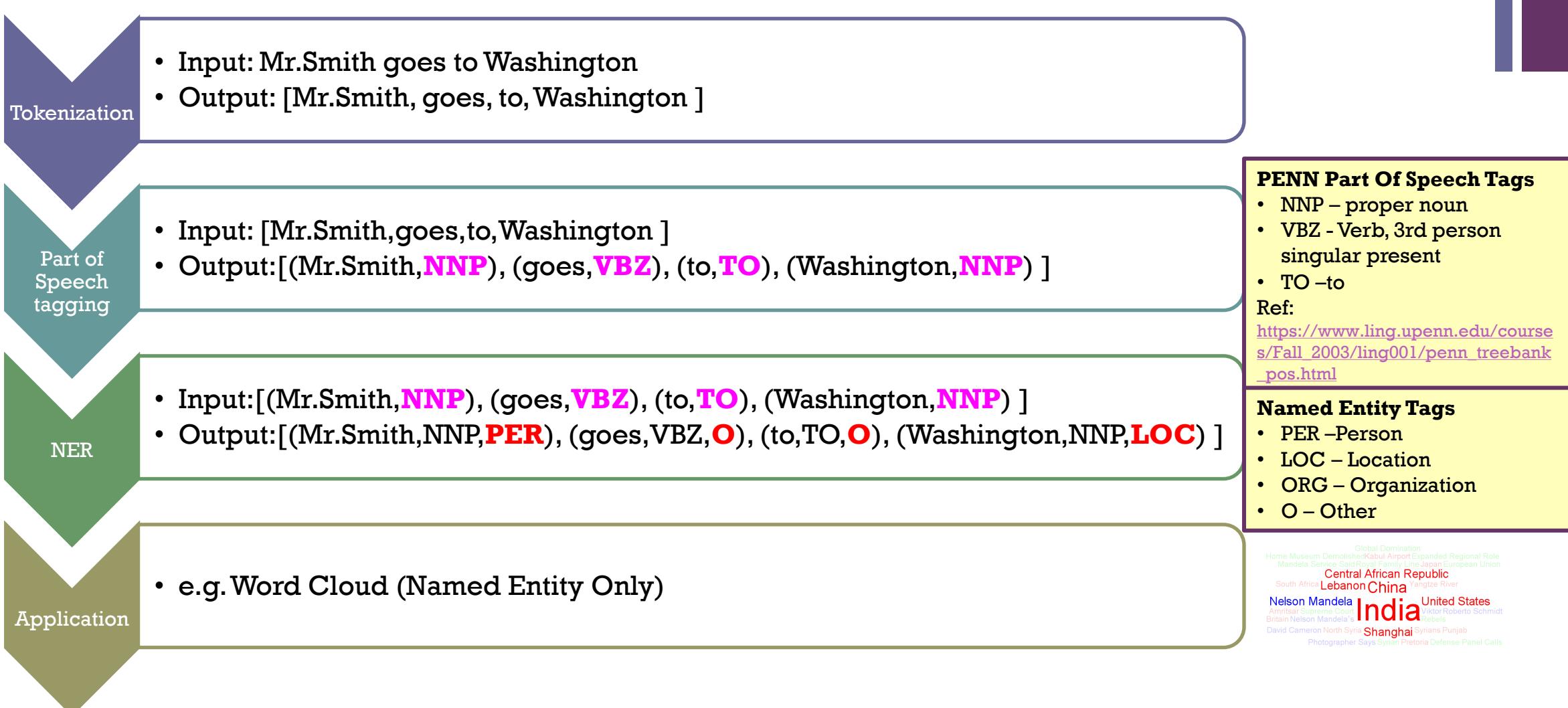
CNN Breaking News @cnnbrk
CNN.com is among the world's leaders in online news and information delivery.



Comment	Good	Like	Hate	#
Tweet1	7	8	0	😊
Tweet2	1	0	10	😢
Tweet3	2	9	1	😊



NLP pipeline example (English)





NLP pipeline example (Thai)



- Input: ขสมก. เลึง จัด หารณ
- Output: ข, สมก., เลึง, จัด, หา, ณ



- Input: [ข, สมก., เลึง, จัด, หา, ณ]
- Output: [(ข, NR), (เลึง, VV), (จัด, VV), (หา, NN), (ณ, NN)]

PENN Part Of Speech Tags

- NR – proper noun
 - VV - Main verbs in clauses, verb-form
 - NN – Non-proper noun
- Ref: BEST2010 dataset



- Input: [(ข, NR), (เลึง, VV), (จัด, VV), (หา, NN), (ณ, NN)]
- Output: [(ข, NR, ORG), (เลึง, VV, O), (จัด, VV, O), (หา, NN, O), (ณ, NN, O)]

Named Entity Tags

- PER –Person
- LOC – Location
- ORG – Organization
- O – Other



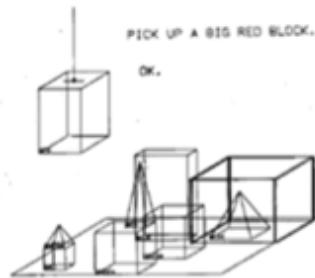
- e.g. Word Cloud (Named Entity Only)



History of NLP techniques

■ 1) Symbolic approach:

- encode all the required information into computer
- In 1960s and 1970s, Noam Chomsky (eminent linguist) believed that that statistical techniques would **never** be sufficient to gain a deep understanding on human language.
- This led to the dominance of **knowledge-based approaches**, requiring humans experts to encode knowledge into computers.
- **Disadvantage: It is required substantial human effort.**



The dialog that was used as a SHRDLU demo:

Person: PICK UP A BIG RED BLOCK.
Computer: OK.
Person: GRASP THE PYRAMID.
Computer: I DON'T UNDERSTAND WHICH PYRAMID YOU MEAN.
Person: FIND A BLOCK WHICH IS TALLER THAN THE ONE YOU ARE HOLDING AND PUT IT INTO THE BOX.
Computer: BY "IT", I ASSUME YOU MEAN THE BLOCK WHICH IS TALLER THAN THE ONE I AM HOLDING.
Computer: OK. (does it)
Person: WHAT DOES THE BOX CONTAIN?
Computer: THE BLUE PYRAMID AND THE BLUE BLOCK.
Person: WHAT IS THE PYRAMID SUPPORTED BY?
Computer: THE BOX.
Person: HOW MANY BLOCKS ARE NOT IN THE BOX?

```
(DEFTHEOREM TC-GRASP
  (THCONSE (X Y)  (#GRASP $?X))
  (THGOAL(#MANIP $?X))
  (THCOND ((THGOAL (#GRASPING $?X)))
    ((THGOAL (#GRASPING $_Y))
     (THGOAL (#GET-RID-OF $?Y)
      (THUSE TC-GET-RID-OF))))
  (T))
  (THGOAL (#CLEARTOP $?X) (THUSE TC-CLEARTOP))
  (THSETQ $_Y (TOPCENTER $?X))
  (THGOAL (#MOVEHAND $?Y)
    (THUSE TC-MOVEHAND))
  (THASSERT (#GRASPING $?X)))))

(DEFTHEOREM TC-PUT
  (THCONSE (X Y Z)  (#PUT $?X $?Y))
  (CLEAR $?Y (SIZE $?X) $?X)
  (SUPPORT $?Y (SIZE $?X) $?X)
  (THGOAL (#GRASP $?X) (THUSE TC-GRASP))
  (THSETQ $_Z (TCENT $?Y (SIZE $?X)))
  (THGOAL (#MOVEHAND $?Z) (THUSE TC-MOVEHAND))
  (THGOAL (#UNGRASP) (THUSE TC-UNGRASP))))
```

The Internals of SHRDLU



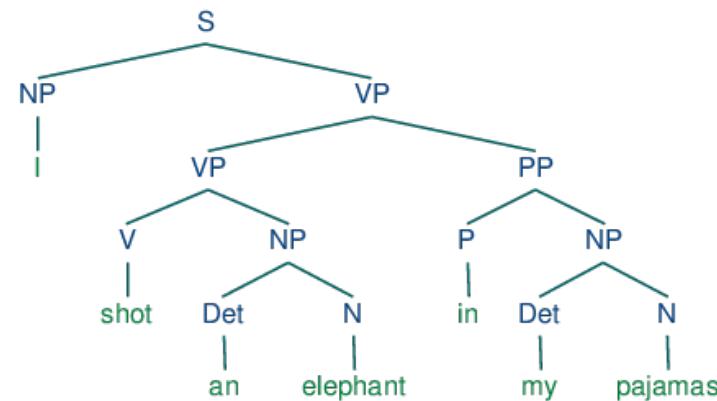
History of NLP techniques (cont.)

■ 2) Statistical approach:

- infer language properties from language samples
- In 1980s, an empirical revolution took place. Inspired by information theory, it began using **probabilistic approaches** in NLP.
- **Disadvantage: It is required hand-crafted features.**

■ 2.5) Deep Learning approach:

- It is a **feature-engineering embedded** neural approach.
- Since 2010s, it has been gaining a lot of attentions and showing many successes.



PennTree Bank (1993): one million words from WSJ, manually annotated with syntactic structure



Case Study: Determiner placement

Symbolic vs. statistical approaches

- Goal: It aims to place “the” (determiner).

Scientists in United States have found way of turning lazy monkeys into workaholics using gene therapy. Usually monkeys work hard only when they know reward is coming, but animals given this treatment did their best all time. Researchers at National Institute of Mental Health near Washington DC, led by Dr Barry Richmond, have now developed genetic treatment which changes their work ethic markedly. "Monkeys under influence of treatment don't procrastinate," Dr Richmond says. Treatment consists of anti-sense DNA - mirror image of piece of one of our genes - and basically prevents that gene from working. But for rest of us, day when such treatments fall into hands of our bosses may be one we would prefer to put off.

Types of Determiner		
Articles	Demonstrative	Possessive Adjectives
the an A	this that these those	my, your his, her its, our your, their
Quantifiers	Numbers	Ordinals
some, any few, little more, much any, every	one, two three, four twenty, hundred	First, Second Third, Last next

www.links2learn.co.uk



Case Study: Determiner placement (cont.)

Symbolic vs. statistical approaches

Symbolic approach

- Determiner placement is largely determined by:
 - Type of noun (countable, uncountable)
 - Uniqueness of reference
 - Information value (given, new)
 - Number (singular, plural)
- However, **many exceptions** and special cases play a role:
 - The definite article is used with newspaper titles (The Times), but zero article in names of magazines and journals (Time)
- **Hard to manually encode this information!**

Statistical approach

- Consider it as **classification**
- Predictions: $\{-1, +1\}$
- Features:
 - plural?
 - first appearance in text?
 - head token
 - ...

Minnen et al.	83.58%
Turner&Charniak	86.74%
Knight&Chander	78%

“lazy monkeys”

$$\begin{matrix} \downarrow \\ [1 \ 1 \ 0 \ 0 \ 0 \ \dots \ 1]^T \end{matrix}$$

-1

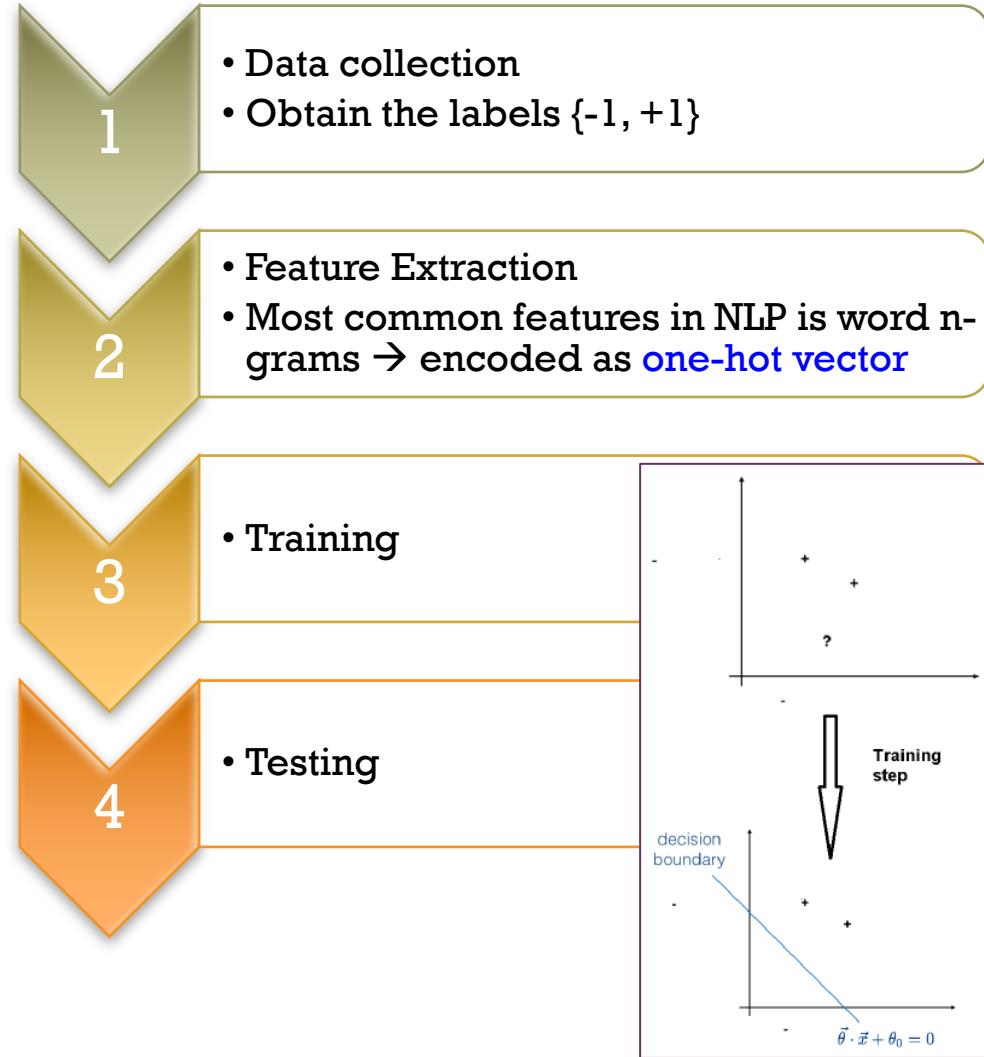
“the United States”

$$\begin{matrix} \downarrow \\ [1 \ 1 \ 0 \ 0 \ 0 \ \dots \ 0]^T \end{matrix}$$

+1



Limitation of traditional statistical approach



- Sparsity:
 - feature vectors are typically high-dimensional and sparse (i.e. most elements are 0).
- Feature engineering:
 - Need experts to manually design features



Map discrete, one-hot vectors into low-dimensional continuous representations.
Self learned features → Deep Learning

pear

[1 0 0 0 ... 0]



[0.4 0.1 0.1]

apple

[0 0 1 0 ... 0]



[0.6 0.2 0.3]

+

Deep Learning



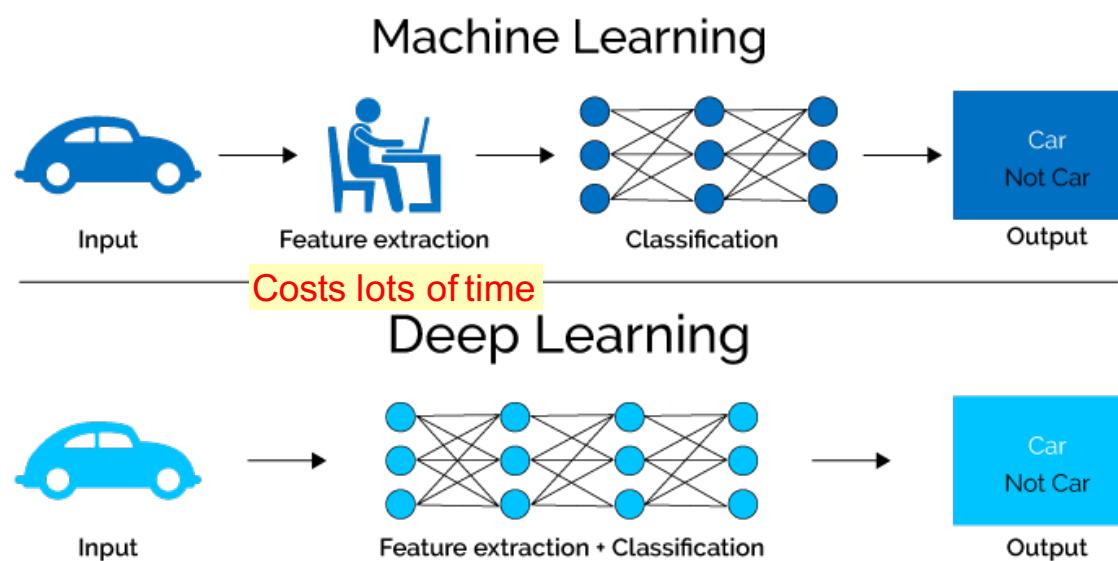
What is Deep Learning (DL)?



Part of the machine learning field of learning representations of data. Exceptional effective at learning patterns.



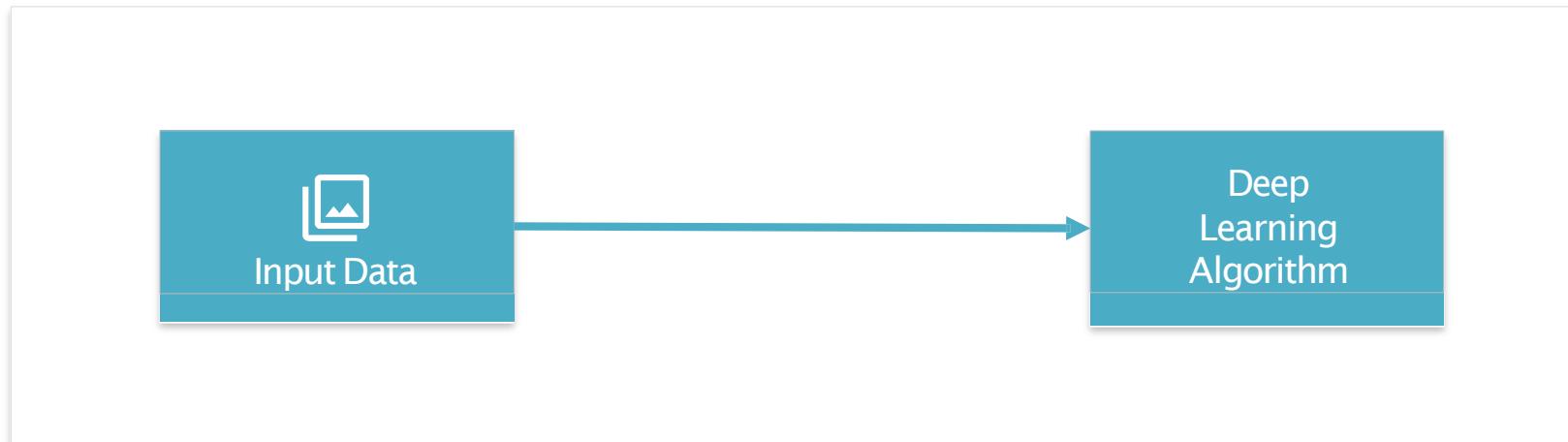
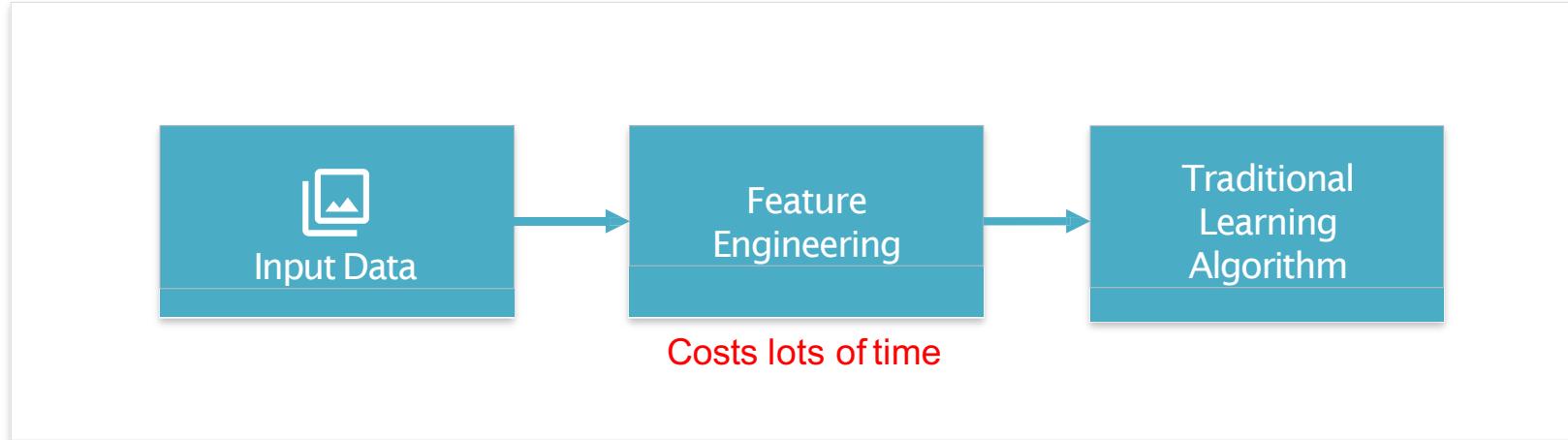
Utilizes learning algorithms that derive meaning out of data by using a hierarchy of multiple layers that mimic the neural networks of our brain.





Deep Learning - Basics

No more feature engineering

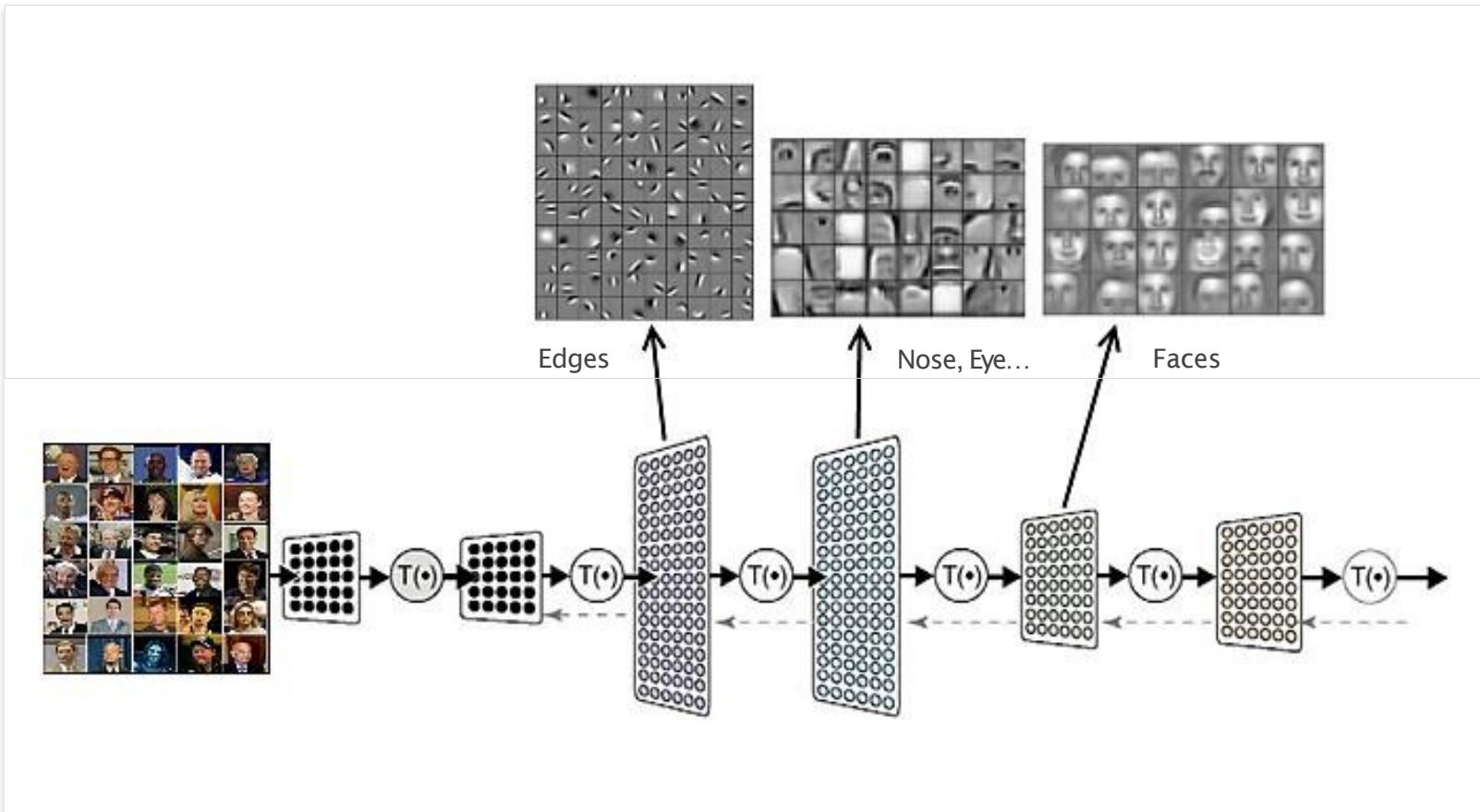




Deep Learning – Basics (cont.)

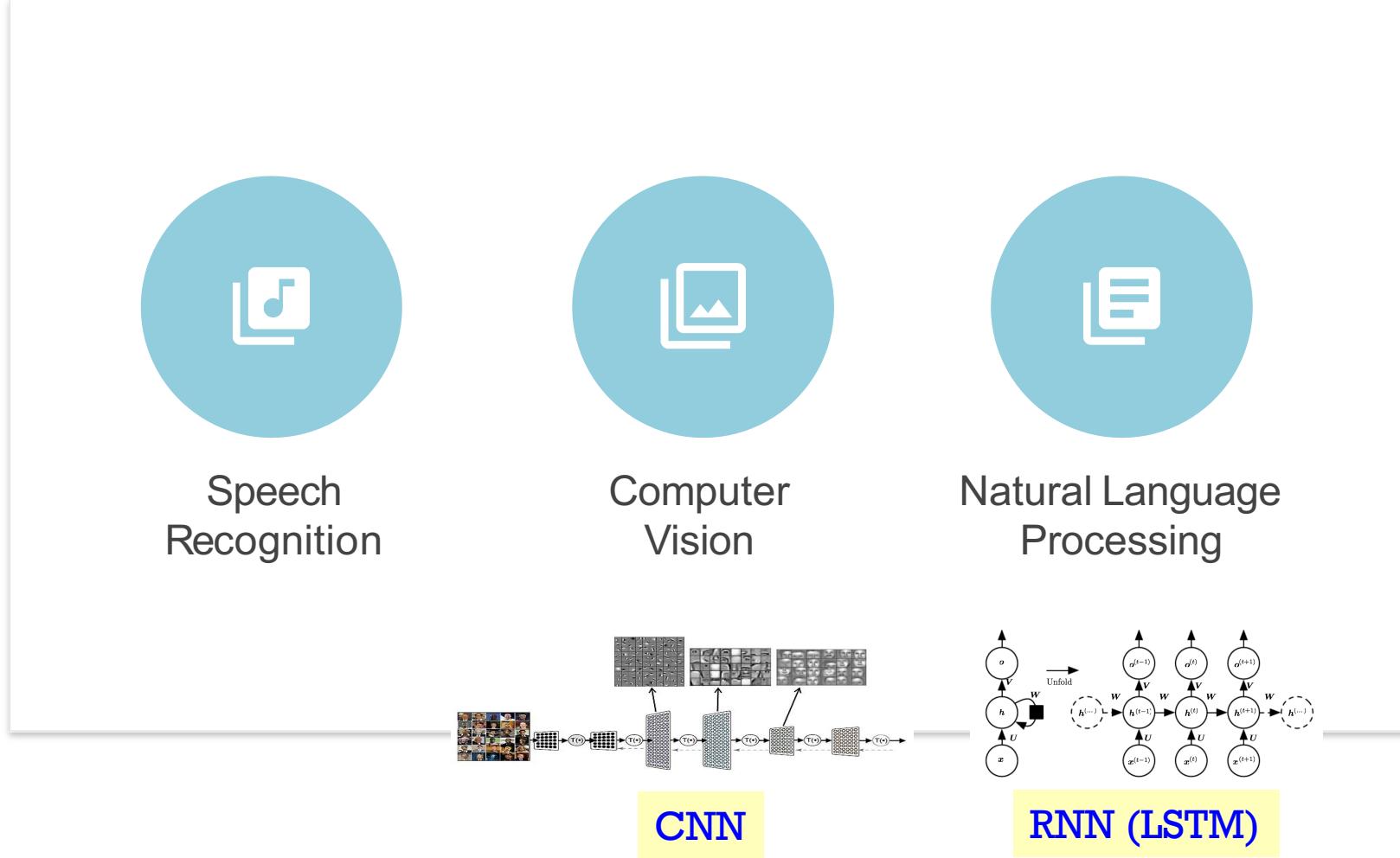
What did it learn?

A deep neural network consists of a **hierarchy of layers**, whereby each layer **transforms the input data** into more abstract representations (e.g. edge \rightarrow nose \rightarrow face). The output layer combines those features to make predictions.





Deep Learning Application

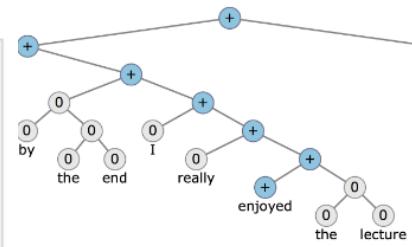




NLP + Deep Learning = Deep NLP

- Modern NLP techniques are based on deep learning models.
- These models have obtained very high performance across various NLP tasks.
- They often **do not** require traditional linguistic feature engineering to perform well.

 CS224d: Deep Learning for Natural Language Processing
 วิชา NLP with Deep Learning ของ Stanford ของ Winter 2017 ล่าสุดครับ
 Lecture Collection | Natural Language Processing with Deep Learning (Winter 2017) - YouTube
 Natural language processing (NLP) deals with the key artificial intelligence technology of understanding...
[YOUTUBE.COM](#)



pucktada/cutkum

[cutkum - Thai Word-Segmentation with Deep Learning in Tensorflow](#)

Pucktada Treeratpituk
RNN, F1 = 0.93 on BEST2010

[GITHUB.COM](#)

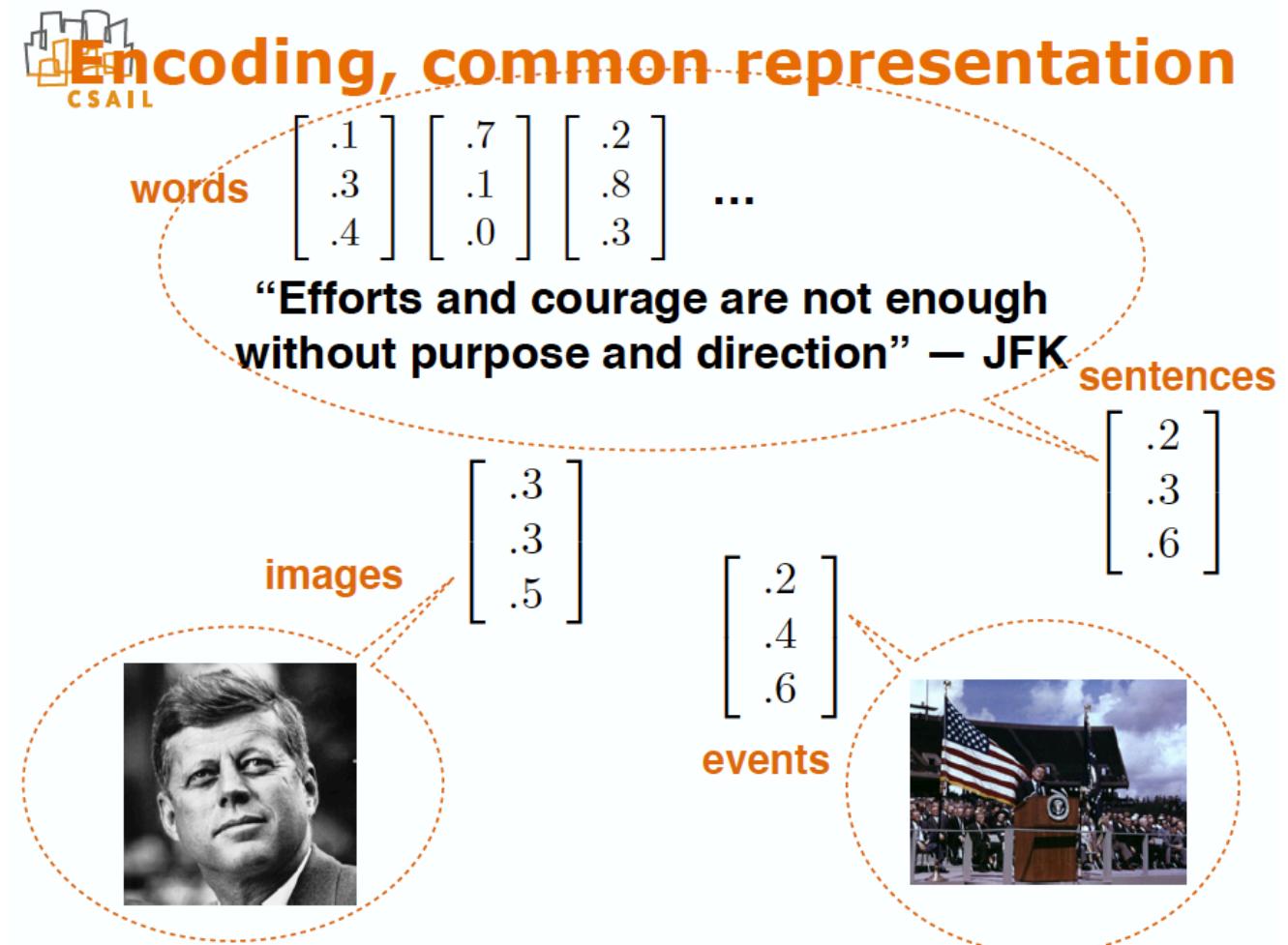
Thai word segmentation with bi-directional RNN F1=99.18%

This is code for preprocessing data, training model and inferring word segment boundaries of Thai text with bi-directional recurrent neural network. The model provides precision of 99.04%, recall of 99.31% and F1 score of 99.18%. Please see the [blog post](#) for the detailed description of the model.



Reasons for exploring Deep Learning

- Learned features are easy to adapt, fast to learn
- Deep learning provides a very flexible? Universal, learnable framework for **representing** world, visual, and linguistic information

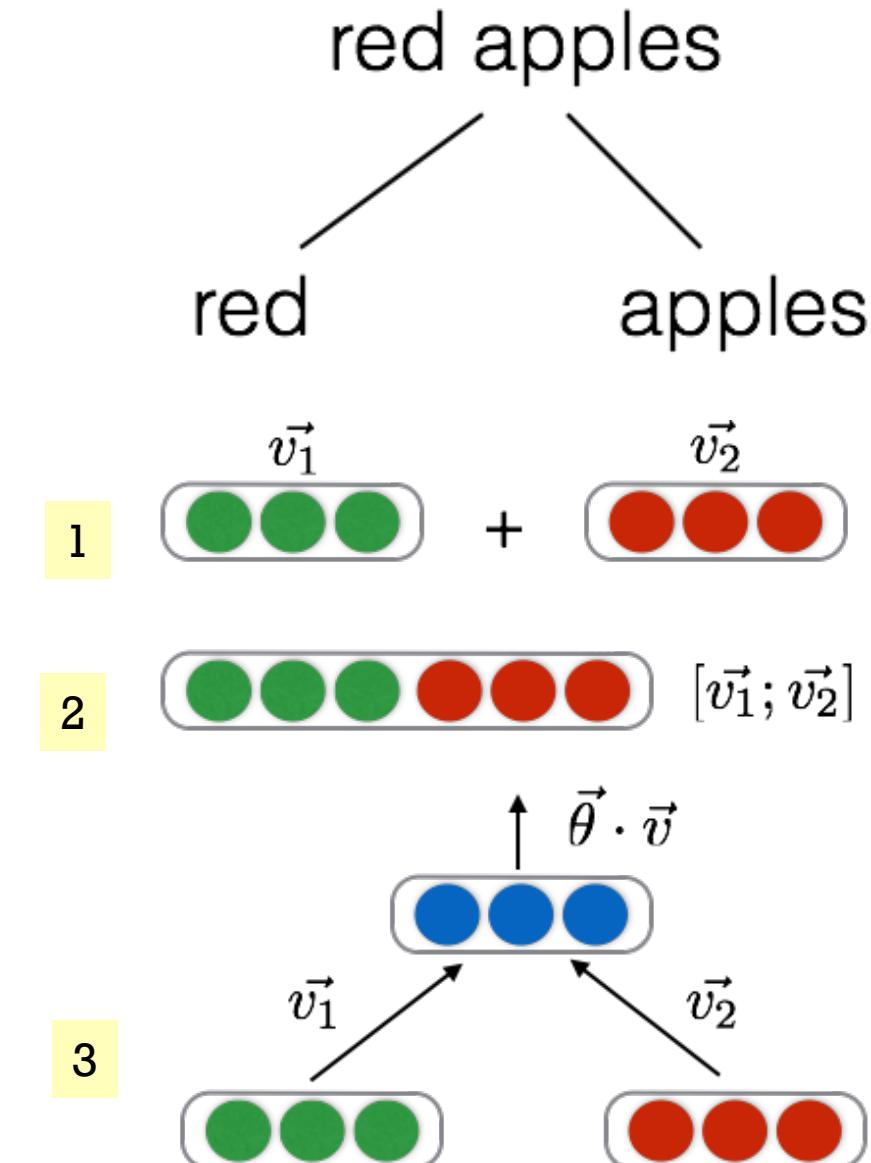




Reasons for exploring Deep Learning (cont.)

- Flexible neural “Lego pieces”
 - Common representation, diversity of architectural choices

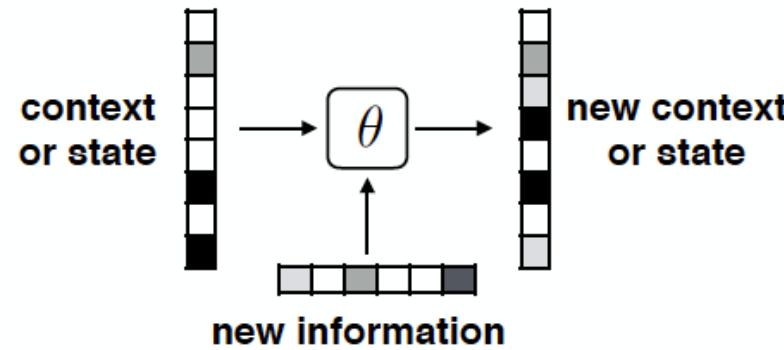
- Can represent any levels of NLP
 - Word
 - Phrase
 - Sentence
 - Paragraph (document)



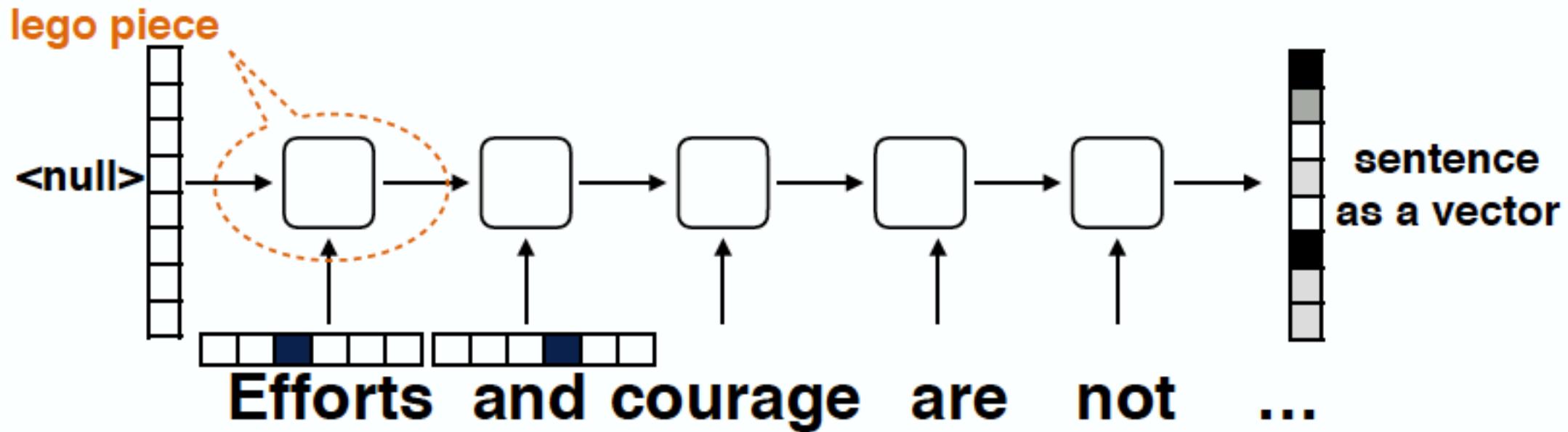


Reasons for exploring Deep Learning (cont.)

Example of encoding sentences

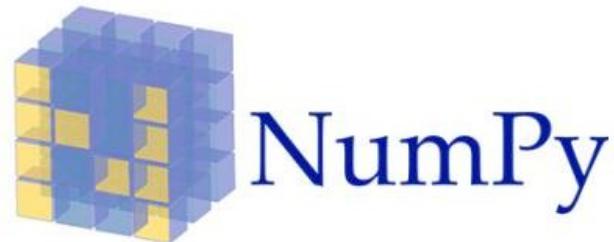


**RNN
(LSTM, GRU)**



+

NLP Tools



- <http://nltk.org>
- Tokenization (Parser)
- Stop words removal
- Stemming
- N-gram
- Part-of-Speech Tagging
- Named Entity Recognition
- Etc.



Tools for Thai NLP

Open-source

- Word Segmentation
 - CutKum
 - LexTo (NECTEC)
 - SWATH (CMU) – 80's
 - ICU (IBM) – 90's
- POS tagging
 - SWATH (CMU)

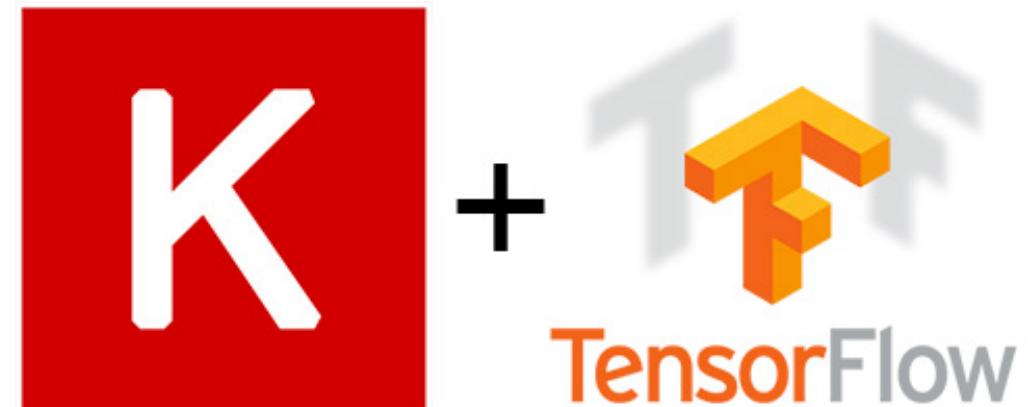
Commercial

- Word Segmentation
 - SegIt (NECTEC)
- POS tagging
 - PosIt (NECTEC)

+

Deep Learning tools

- Tensorflow (Google): Python, etc.
- Torch (The Idiap Research Institute):
 - PyTorch: Python
- Theano (University of Montreal): Python
- Keras (François Chollet): Python
- Lasagne (Sander Dieleman): Python
- MXNet (Microsoft)





Course Logistics



Class schedule

Word Model

Sequence

Tree

Application

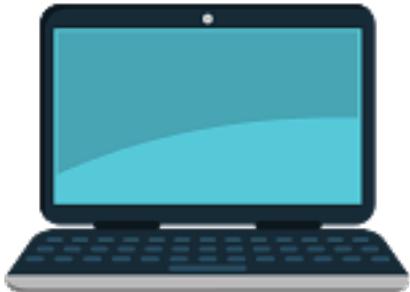
Implementation

Week	Topics (E=Ekapol, P=Peerapon)
1 - 8/1	Intro to NLP (P)
2 - 15/1	Tokenization LexTo, DNN, CNN, LSTM with Keras. (E)
3 - 29/1	Language modeling (P) N-grams, smoothing, Neural LM
4 - 5/2	Representation (P) TF-IDF, word and sentence embeddings, adaptation
5 - 12/2	PoS tagging and information extraction (E) CRF and beamsearch
6 - 19/2	Parsing (E) PCFG, Recursive neural networks
	Midterm week
7 - 26/2	In class midterm
8 - 12/3	Document/sentiment classification (E) LDA, Naive Bayes, EM
9 - 19/3	Question Answering (P&E) Reading comprehension, text generation
10 - 26/3	Chatbots (P)
11 - 2/4	Industry Talk (TBA)
12 - 9/4	Other NLP applications (E) ASR, MT, Summarization
13 - 23/4	Project presentation

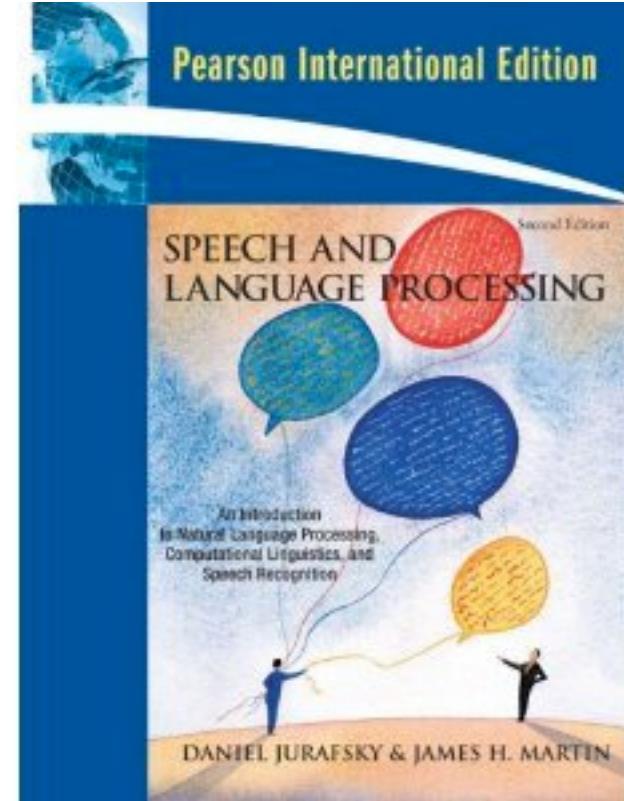


Course Grading

- Assignments 30% (5% each with 5% extra)
- Midterm 35%
- Project 35%



Google Cloud Platform

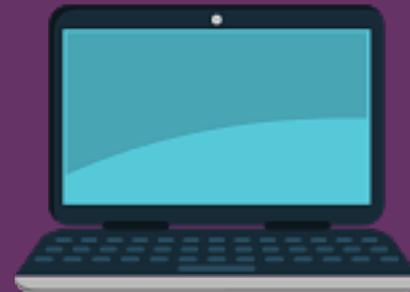


Speech and Language Processing, 2nd Edition 2nd Edition
by [Daniel Jurafsky](#) (Author), [James H. Martin](#) ▾ (Author)

<https://nlp.stanford.edu/~manning/xyzzy/JurafskyMartinEd2book.pdf>

+

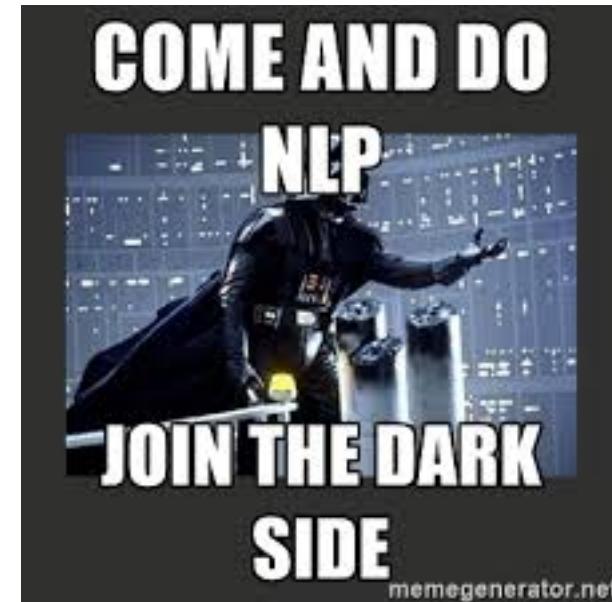
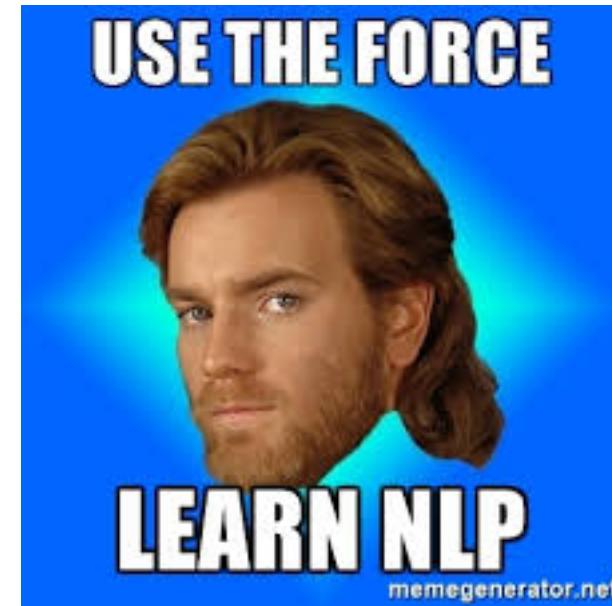
Google Cloud Demo



+

NLP

- Neuro-Linguistic Programming
 - Top hit if you google NLP
 - Pseudoscience
- Natural Language Processing



นวัตกรรมเปลี่ยนโลก ‘เอไอ-ดิจิทัลทิวน-บล็อกเชน’

• **บุษกร ภู่แสง**
กรุงเทพธุรกิจ

สำนักงานคณะกรรมการนโยบายวิทยาศาสตร์ฯ เปิดรายงาน ๖ อันดับนวัตกรรมเด่นปี 2561 เมย์เทคโนโลยีปัญญาประดิษฐ์ ติดอันดับที่ห้ามาระมารังต่อเนื่อง ตามด้วยบล็อกเชน จีโนมิกส์ โอลิมิกส์ เทคโนโลยีการพิมพ์สามมิติ และชาร์ริงอิโคโนมี ที่ต้องจับตามอง ขณะที่ ยุปกรณ์ความจำเพื่อเก็บข้อมูลเชิง ดิจิทัล แพลทฟอร์ม อาร์คิดิสก์แบบพกพา กลายเป็น เทคโนโลยีดาวรุ่ง เหตุผู้ใช้เปลี่ยนไปเก็บ ข้อมูลในระบบปฏิบัติการคลาวด์แทน

กิติพงศ์ พ้อมองค์ เผยวิธีการสำนักงานคณะกรรมการนโยบายวิทยาศาสตร์ฯ เทคโนโลยี และนวัตกรรมแห่งชาติ (สวทน.) กระทรวงวิทยาศาสตร์และเทคโนโลยี กล่าวว่า แนวโน้ม เทคโนโลยีและนวัตกรรมที่มีโอกาสพัฒนาและ ก้าวหน้าในปี 2561 จะมุ่งเน้นความเป็นอัจฉริยะ และความปลอดภัยในทุกด้าน ที่เพิ่มมากขึ้น

โลกดิจิทัลจะก้าวเข้าไปเกี่ยวข้องกับ

วิถีชีวิตและการดำเนินธุรกิจอย่างป্রาย坪หากเพิ่มสูงขึ้น ดังนั้นนวัตกรรมและเทคโนโลยี ที่เป็นดาวรุ่ง ที่มีการนำมาประยุกต์ใช้ใน ชีวิตประจำวันและการดำเนินธุรกิจต่างๆ ในปี 2561 ได้ดังนี้

‘เอไอ’มีบทบาทในทุกวงการ

- ก. กลุ่มอัจฉริยะเทคโนโลยีและนวัตกรรม ที่มีแรงและเริ่มมีบทบาทสำคัญในโลกปัจจุบัน จนถึงอนาคต คือ เทคโนโลยีปัญญาประดิษฐ์ (Artificial Intelligence หรือ AI) คือ การใช้ เทคโนโลยีคอมพิวเตอร์ร่วมกับเครื่องจักรกล เพื่อการเรียนแบบสติปัญญาและการกระทำ ของมนุษย์ เช่น การเข้าใจภาษา การวิเคราะห์ ตัวเหตุและผล และการรับรู้สภาพแวดล้อม ต่างๆ บริษัทเทคโนโลยีและนวัตกรรมที่ขึ้นนำ ระดับโลกมีการนำเทคโนโลยีเอไอเข้ามาเพื่อ พัฒนานวัตกรรมและการบริการเป็นอย่าง มาก เช่น ภูมิสืบ เพชบุค และแอปเปิล ที่นำ เอไอ มาใช้งานด้านเพื่อเป็นเครื่องมือเพื่อ การตัดสินใจ

มีการคาดการณ์ว่า หากมีการนำ เทคโนโลยีเอไอมาใช้งานจำนวนมากจะ ทำให้การใช้แรงงานมนุษย์ลดลง ตัวอย่าง การใช้เทคโนโลยีเอไอแล้วมีผลกระทบต่อ การลดจำนวนแรงงานคน เช่น ยานพาหนะ ที่ขับเคลื่อนได้ด้วยตัวเอง

หนึ่งในตัวอย่างของอุปกรณ์และสิ่งของ อัจฉริยะ ที่หลายบริษัทบัญชีใหญ่ประกาศเปิด ตัวและทดสอบชั้บ ทำให้อาชีพคนขับรถที่ใช้ แรงงานมนุษย์ลดลง และจากเทคโนโลยีเอไอจะ นำไปสู่การพัฒนาเทคโนโลยีที่เกี่ยวข้อง ได้แก่ Machine Learning, Neuro-Linguistic Programming (NLP) ในหลายประเทศ แยกบอเมริกาเหนือและลاتินอเมริกามีการนำ เทคโนโลยีนี้เข้าไปใช้งานและแก้ไขปัญหาใน ธุรกิจและอุตสาหกรรมต่างๆ เพิ่มขึ้น

เช่น ในบริษัท Progressive Environmental & Agricultural Technologies จำกัด (PEAT) ได้พัฒนา ซอฟต์แวร์ชื่อ Plantix ที่อ่านพืชที่ดินฐานเทคโนโลยี เอไอ เพื่อการนับตรวจน้ำพืชอย่างแม่นยำ



Computational Linguistics vs NLP

- If you use the term **Computational Linguistics (CL)**
 - You are probably a linguist
 - You use computers to study and analyze languages
- If you use the term **NLP**
 - You are probably an engineer/computer scientist
 - You work on applications that involve languages
- Most of the time they're very related (or the same thing)