

經濟部工業局

AI產業實戰應用人才淬煉計畫

人工智慧及資安應用實務訓練課程

遞歸神經網絡

**Recurrent Neural Network**

# aGENDA

- Sequence Data的分析(1)Time Series
- Sequence Data的分析(2)NLP
- 遞歸神經網絡  
Recurrent Neural Network
- NLP in security

Sequence Data

的分析(1)

Time Series

# Time Series

[https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series)

一組按照時間發生先後順序進行排列的數據點序列。

通常一組時間序列的時間間隔為一恆定值（如1秒，5分鐘，12小時，7天，1年），因此時間序列可以作為離散時間數據進行分析處理。

時間序列廣泛應用：

數理統計、

信號處理、

模式識別、

計量經濟學、

數學金融、

天氣預報、

地震預測、

腦電圖、

控制工程、航空學、通信工程以及絕大多數涉及到時間數據測量的應用科學與工程學。

## Time Series Analysis

時間序列模型基本概念：AR, MA, ARMA, ARIMA 模型

<https://mropengate.blogspot.com/2015/11/time-series-analysis-ar-ma-arma-arima.html>

<http://yongfeng.me/attach/time-series-analysis-zhang.pdf>

<https://ocw.mit.edu/courses/mathematics/18-s096-topics-in-mathematics-with-applications-in-finance-fall-2013/video-lectures/lecture-8-time-series-analysis-i/>

MIT OPEN COURSEWARE  
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

[Subscribe to the OCW Newsletter](#)

FIND COURSES ▾ For Educators ▾ Give Now ▾ About ▾ Custom ▾

Home » Courses » Mathematics » Topics in Mathematics with Applications in Finance » Video Lectures » Lecture 8: Time Series Analysis I

## Lecture 8: Time Series Analysis I

COURSE HOME

SYLLABUS

CALENDAR

INSTRUCTOR  
INSIGHTS

LECTURE NOTES &  
SLIDES

CASE STUDIES

VIDEO LECTURES

ASSIGNMENTS



Interactive Transcript



Previous

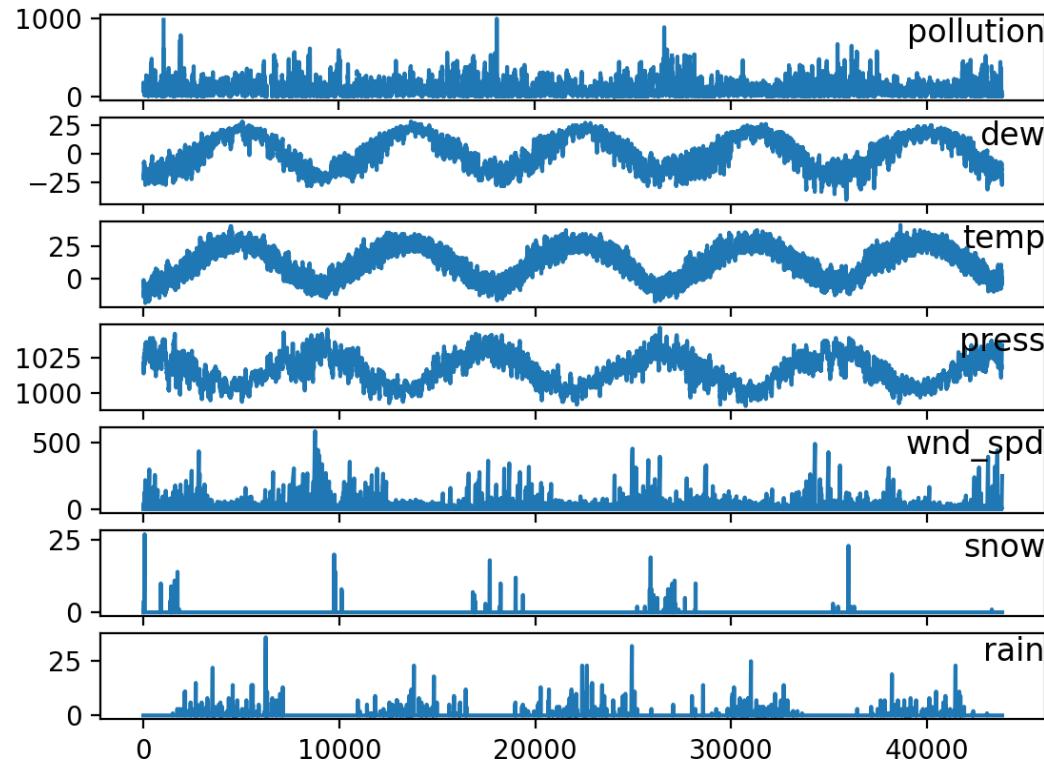
Next

# Deep Time Series Classification

- 預測股票每日最高價
- 空氣污染預測
- 地震預測
- PM2.5預測
- GDP預測
- 人口數預測

# 空氣污染預測

<https://raw.githubusercontent.com/jbrownlee/Datasets/master/pollution.csv>



Multivariate LSTM Forecast Model

<https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/>

<https://www.itread01.com/content/1544492764.html>

# <https://www.kaggle.com/>



## The 3rd YouTube-8M Video Understanding Challenge

Temporal localization of topics within video

[Research](#) · a month to go · 📹 video data, object detection

\$25,000  
222 teams



## Open Images 2019 - Object Detection

Detect objects in varied and complex images

[Research](#) · a month to go · 📸 image processing, image data

\$25,000  
453 teams



## Open Images 2019 - Visual Relationship

Detect pairs of objects in particular relationships

[Research](#) · a month to go · 📸 image processing, image data

\$25,000  
145 teams

# <https://faridrashidi.github.io/kaggle-solutions/>

## Kaggle Solutions

Fork me on GitHub

### The Most Comprehensive List of Kaggle Solutions and Ideas

This is a list of almost all available solutions and ideas shared by top performers in the past Kaggle competitions. This list will get updated as soon as a new competition finished. If you find other solutions beside the ones listed here I would suggest you to contribute to this repo by making a pull request. The symbols were used in this list is described [here](#).

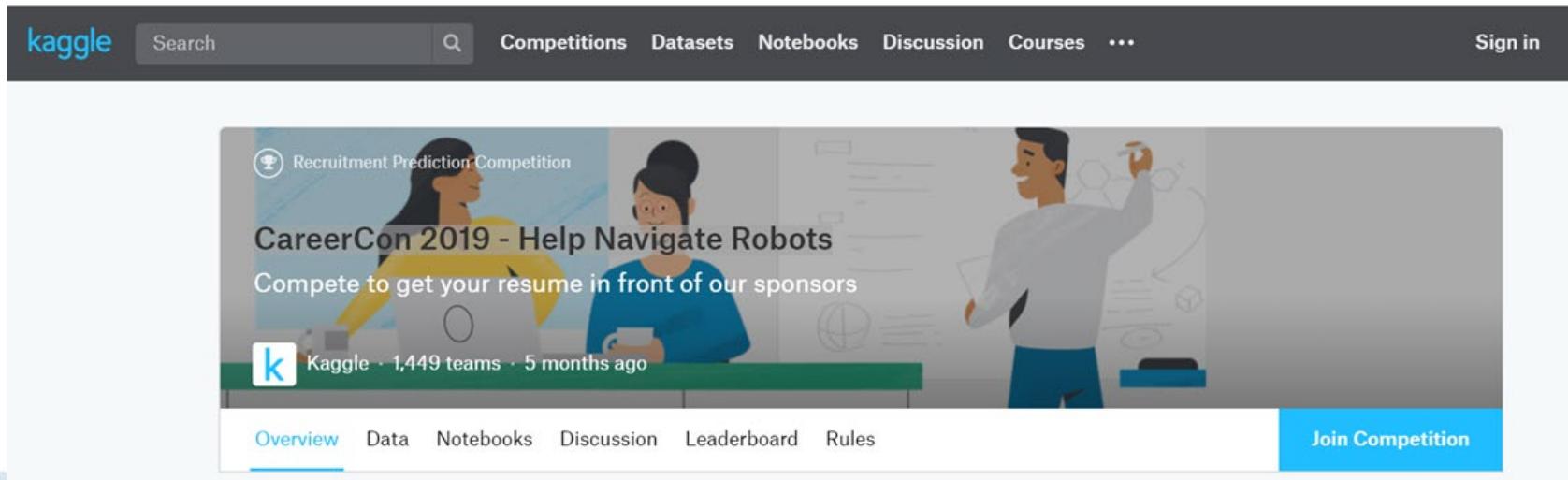
If you found it interesting you can give a star or make a fork

Star 20   Fork 6

Check out the following markdown pages about Top Kagglers Tips/Tricks and all Kernels of The week.

- [Top Kagglers Interviews and Lectures](#)
- [Kernels of The Week](#)

# <https://www.kaggle.com/c/career-con-2019>



<https://www.kaggle.com/purplejester/pytorch-deep-time-series-classification>

<https://github.com/ammar1y/My-Solution-to-Help-Navigate-Robots-Competition>

<https://www.youtube.com/watch?v=ageh45rxyXU>

# Kaggle Earthquake Prediction Challenge

<https://www.kaggle.com/c/LANL-Earthquake-Prediction/data>

The screenshot shows the main landing page for the Kaggle LANL Earthquake Prediction challenge. At the top left is a competition icon (a seismometer) and the text "Research Prediction Competition". In the center, the title "LANL Earthquake Prediction" is displayed above the question "Can you predict upcoming laboratory earthquakes?". To the right, a large yellow seismogram waveform is shown, with "\$50,000 Prize Money" written next to it. Below the title, the Los Alamos National Laboratory logo is present, along with the text "Los Alamos National Laboratory · 4,540 teams · 3 months ago". At the bottom, a navigation bar includes links for "Overview", "Data" (which is underlined in blue), "Notebooks", "Discussion", "Leaderboard", and "Rules". On the far right of the bar is a blue button labeled "Late Submission".

<https://github.com/Kaggle/kaggle-api>

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

[https://github.com/lISourcell/Kaggle\\_Earthquake\\_challenge](https://github.com/lISourcell/Kaggle_Earthquake_challenge)

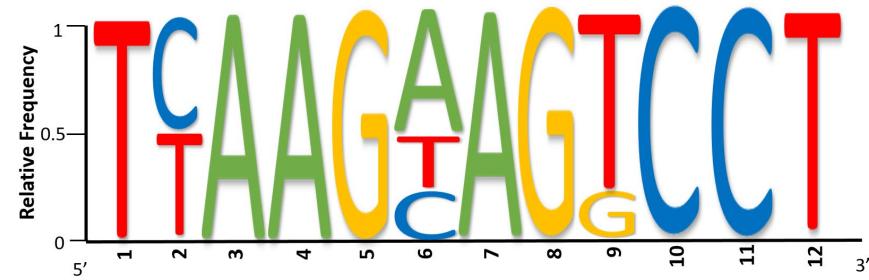
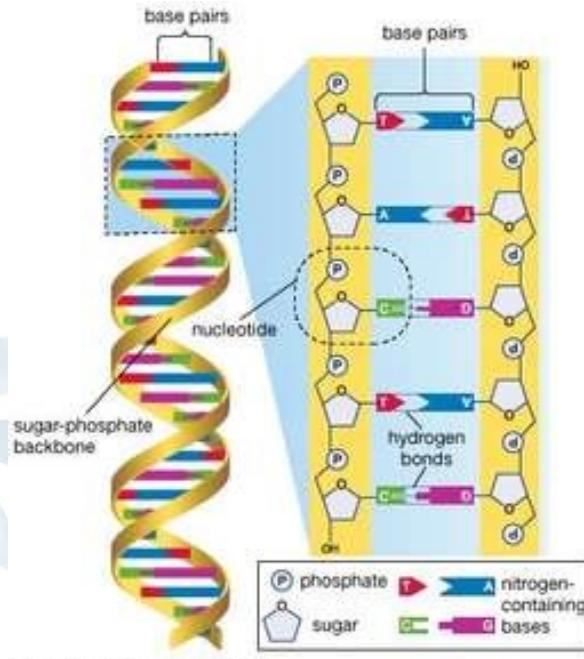
**DEMO**

Kaggle\_challenges/LANL\_Earthquake\_Challenge\_DNN\_1D\_CNN\_LSTM\_with\_TPU.ipynb

# DNA Sequence Data Analysis

Working with DNA sequence data for ML

<https://www.kaggle.com/thomasnelson/working-with-dna-sequence-data-for-ml>



# **Time series analysis in cybersecurity**

# Analyzing and Predicting Cyber Hacking with Time Series Models

C. Soundarya<sup>1\*</sup>, S. Usha<sup>2</sup>

<sup>1</sup>*PG Scholar, Department of Computer Science and Engineering, Anna University, BIT-Campus,  
Tiruchirappalli, India*

<sup>2</sup>*Assistant Professor, Department of Computer Science and Engineering, Anna University, BIT-Campus,  
Tiruchirappalli, India*

\*Corresponding author: csoundarya1995@gmail.com

# A detailed analysis of CICIDS2017 dataset for designing Intrusion Detection Systems

January 2018 · International Journal of Engineering & Technology 7(3):479-482

Authors:



**Ranjit Panigrahi**

Sikkim Manipal Institute of Technology



**Samarjeet Borah**

Sikkim Manipal Institute of Technology

[https://www.researchgate.net/publication/329045441\\_A\\_detailed\\_analysis\\_of\\_CICIDS2017\\_dataset\\_for\\_designing\\_Intrusion\\_Detection\\_Systems](https://www.researchgate.net/publication/329045441_A_detailed_analysis_of_CICIDS2017_dataset_for_designing_Intrusion_Detection_Systems)

## class imbalance problem

**Table 2:** Overall characteristics of CICIDS2017 dataset

Dataset Name	CICIDS2018
Dataset Type	Multi class
Year of release	2017
Total number of distinct instances	2830540
Number of features	83
Number of distinct classes	15

**Table 3:** Class wise instance occurrence of CICIDS2017 dataset

Class Labels	Number of instances
BENIGN	2359087
DoS Hulk	231072
PortScan	158930
DDoS	41835
DoS GoldenEye	10293
FTP-Patator	7938
SSH-Patator	5897
DoS slowloris	5796
DoS Slowhttptest	5499
Bot	1966
Web Attack – Brute Force	1507
Web Attack – XSS	652

**Table 4:** Class prevalence ratio of CICIDS2017 dataset

Sl No	Normal / Attack Labels	Number of instances	% of prevalence w.r.t. the majority class	% of prevalence w.r.t. the total instances
1	BENIGN	2359087	1	83.34406
2	Bot	1966	0.000833	0.06946
3	DDoS	41835	0.017734	1.47799
4	DoS GoldenEye	10293	0.004363	0.36364
5	DoS Hulk	231072	0.09795	8.16353
6	DoS Slow- httptest	5499	0.002331	0.19427
7	DoS slowloris	5796	0.002457	0.20477
8	FTP-Patator	7938	0.003365	0.28044
9	Heartbleed	11	0.000005	0.00039
10	Infiltration	36	0.000015	0.00127
11	PortScan	158930	0.067369	5.61483
12	SSH-Patator	5897	0.0025	0.20833
13	Web Attack – Brute Force	1507	0.000639	0.05324
14	Web Attack – Sql Injection	21	0.000009	0.00074
15	Web Attack – XSS	652	0.000276	0.02303

**Table 5:** Characteristics of new attack labels with their prevalence rate in CICIDS2017 dataset

<i>Sl No</i>	<i>New Labels</i>	<i>Old Labels</i>	<i>Number of instances</i>	<i>% of prevalence w.r.t. the majority class</i>	<i>% of prevalence w.r.t. the total instances</i>
1	Normal	Benign	2359087	100	83.34
2	Botnet ARES	Bot	1966	0.083	0.06
3	Brute Force	FTP-Patator, SSH-Patator	13835	0.59	0.48
4	Dos/DDos	DDoS, DoS GoldenEye, DoS Hulk, DoS Slowhttptest, DoS slowloris, Heartbleed	294506	12.49	10.4
5	Infiltration	Infiltration	36	0.001	0.001
6	PortScan	PortScan	158930	6.74	5.61
7	Web Attack	Web Attack – Brute Force, Web Attack – Sql Injection, Web Attack – XSS	2180	0.092	0.07

Sequence Data

的分析(2)

NLP

# NLP Natural Language Processing 自然語言處理

<https://zh.wikipedia.org/wiki/自然語言處理>

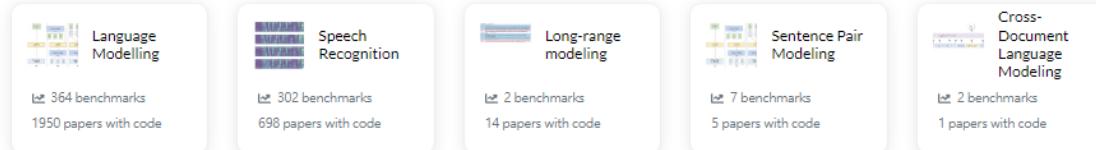
- 自然語言處理（Natural Language Processing，NLP）是人工智能和語言學領域的分支學科。
- 此領域探討如何處理及運用自然語言；自然語言處理包括多方面和步驟，基本有認知、理解、生成等部分。
- 自然語言認知和理解是讓電腦把輸入的語言變成有意思的符號和關係，然後根據目的再處理。自然語言生成系統則是把電腦資料轉化為自然語言。

# Natural Language Processing

1700 benchmarks • 517 tasks • 1498 datasets • 15621 papers with code

## Language Modelling

<https://paperswithcode.com/area/natural-language-processing>



## 2D Classification



## Question Answering



[▶ See all 14 tasks](#)

## Machine Translation



[▶ See all 9 tasks](#)

Information Gathering	Information Generation	Network Traffic Analysis	Malware Analysis	Event Detection and Prediction
Audio / Textual Impersonation	Identity / Private Information	Anomaly Detection	Code Analysis	Threat Intelligence
(Spear-) Phishing	Censorship and Disinformation	Domain Classification	Vulnerability Assessment	Risk Management

[linkedin.com/in/ursachi/](https://www.linkedin.com/in/ursachi/)



<https://www.masernet.com/project/role-and-applications-of-nlp-in-cybersecurity>

# RNN for.....

列到序列翻譯（sequence to sequence translation），包括將語音轉為文字或翻譯不同語言

## Sequential patterns

Text

Speech

Audio

Video

Physical processes

Anything embedded in time (almost everything)

# NLP應用

文字朗讀 ( Text to speech ) / 語音合成 ( Speech synthesis )

語音識別 ( Speech recognition )

中文自動分詞 ( Chinese word segmentation )

詞性標註 ( Part-of-speech tagging )

句法分析 ( Parsing )

自然語言生成 ( Natural language generation )

文字分類 ( Text categorization )

資訊檢索 ( Information retrieval )

資訊抽取 ( Information extraction )

文字校對 ( Text-proofing )

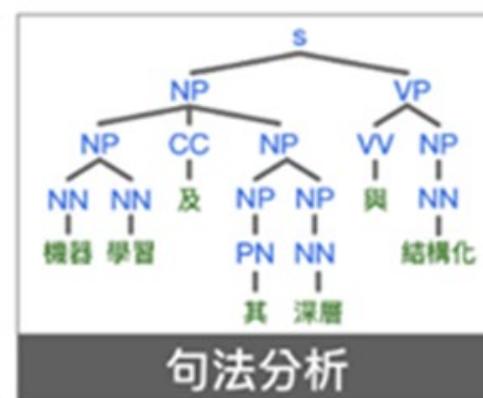
問答系統 ( Question answering )

機器翻譯 ( Machine translation )

自動摘要 ( Automatic summarization )

文字蘊涵 ( Textual entailment )

命名實體辨識 ( Named entity recognition )



# Chinese Word Segmentation 中文分詞

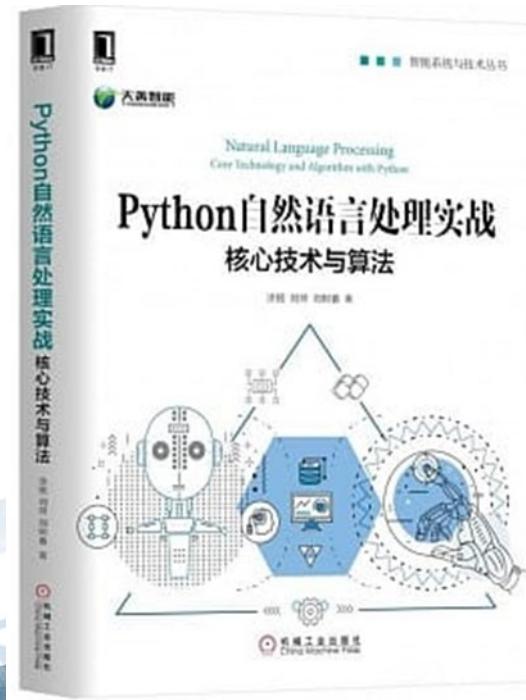
Chinese word segmentation is the task of splitting Chinese text (a sequence of Chinese characters) into words.

'上海浦東開發與建設同步' → ['上海', '浦東', '開發', '與', '建設', '同步']

♣ marks the system that uses character unigram as input. ♦ marks the system that uses character bigram as input.

- Huang et al. (2019): BERT + model compression + multi-criterial learing ♦
- Yang et al. (2018): Lattice LSTM-CRF + BPE subword embeddings ♦♦
- Ma et al. (2018): BiLSTM-CRF + hyper-params search♦♦
- Yang et al. (2017): Transition-based + Beam-search + Rich pretrain♦♦
- Zhou et al. (2017): Greedy Search + word context♦
- Chen et al. (2017): BiLSTM-CRF + adv. loss♦♦
- Cai et al. (2017): Greedy Search+Span representation♦
- Kurita et al. (2017): Transition-based + Joint model♦
- Liu et al. (2016): neural semi-CRF♦
- Cai and Zhao (2016): Greedy Search♦
- Chen et al. (2015a): Gated Recursive NN♦♦
- Chen et al. (2015b): BiLSTM-CRF♦♦

# Chinese Word Segmentation 中文分詞



chapter-3 中文分詞技術

chapter-4 詞性標注與命名實體識別

chapter-5 關鍵字提取

chapter-6 句法分析

chapter-7 文本向量化

chapter-8 情感分析

chapter-9 NLP中用到的機器學習演算法

chapter-10 基於深度學習的NLP演算法

<https://github.com/nlpinaction/learning-nlp>

# Chinese Word Segmentation 中文分詞

<https://github.com/fxsjy/jieba>

支持繁體分詞

支持自訂字典

MIT 授權協議

支持三種分詞模式：

全模式，把句子中所有的可以成詞的詞語都掃描出來，速度非常快，但是不能解決歧義；

精確模式，試圖將句子最精確地切開，適合文本分析；

jieba

“結巴” 中文分詞：做最好的 Python 中文分片語件

"Jieba" (Chinese for "to stutter") Chinese text segmentation: built to be the best Python Chinese word segmentation module.

```
import jieba
```

```
sent = '中文分詞是文本處理不可或缺的一步！'
```

```
seg_list = jieba.cut(sent, cut_all=True)
print('全模式：', '/ '.join(seg_list))
```

```
seg_list = jieba.cut(sent, cut_all=False)
print('精確模式：', '/ '.join(seg_list))
```

```
seg_list = jieba.cut(sent)
print('預設精確模式：', '/ '.join(seg_list))
```

```
seg_list = jieba.cut_for_search(sent)
print('搜尋引擎模式', '/ '.join(seg_list))
```

搜索引擎模式，在精確模式的基礎上，對長詞再一次切分，提高召回率，適合用於搜尋引擎分詞。

```
1 import jieba  
2  
3 sent = '中文分詞是文本處理不可或缺的一步！'  
4  
5 seg_list = jieba.cut(sent, cut_all=True)  
6  
7 print('全模式: ', '/'.join(seg_list))  
8  
9 seg_list = jieba.cut(sent, cut_all=False)  
10 print('精確模式: ', '/'.join(seg_list))  
11  
12 seg_list = jieba.cut(sent)  
13 print('預設精確模式: ', '/'.join(seg_list))  
14  
15 seg_list = jieba.cut_for_search(sent)  
16 print('搜尋引擎模式: ', '/'.join(seg_list))  
17
```

全模式： 中文/ 分/ 詞/ 是/ 文本/ 處理/ 不可/ 不可或缺/ 或缺/ 的/ 一步/ /  
精確模式： 中文/ 分詞/ 是/ 文本/ 處理/ 不可或缺/ 的/ 一步/ !  
預設精確模式： 中文/ 分詞/ 是/ 文本/ 處理/ 不可或缺/ 的/ 一步/ !  
搜尋引擎模式 中文/ 分詞/ 是/ 文本/ 處理/ 不可/ 或缺/ 不可或缺/ 的/ 一步/ !

搜索引擎模式，在精確模式的基礎上，對長詞再次切分，提高召回率，適合用於搜尋引擎分詞。

# 詞性標註Part-of-speech tagging

<http://terms.naer.edu.tw/detail/1678982/>

字詞（**word**）是語言系統中具有獨立語意或扮演特定語法功能，且可以自由使用的最小語言單位。

依據字詞在句法結構或語言形態上扮演的角色，經由詞性分類賦予語句中每個字詞適當之詞性符號或標記的過程，則稱為詞性標記（**part-of-speech tagging**，或稱POS tagging）。

隨著語言的不同，詞性分類的方式也有所差異。

基本上可分為動詞（**verb**）、名詞（**noun**）、形容詞（**adjective**）、副詞（**adverb**），以及其他類別，例如：代名詞（**pronoun**）、介係詞（**preposition**）、連接詞（**conjunction**）或感嘆詞（**interjection**）

# 詞性標註Part-of-speech tagging

<http://terms.naer.edu.tw/detail/1678982/>

詞性標記的範例

Tagging is the task of labeling (or tagging) each word in a sentence with its appropriate part of speech.



Tagging\_VBG is\_BEZ the\_AT task\_NN of\_IN labeling\_VBG (\_( or\_CC  
tagging\_VBG )\_) each\_DT word\_NN in\_IN a\_AT sentence\_NN with\_IN its\_PP\$  
appropriate\_JJ part\_NN of\_IN speech\_NN ..

VBG為動名詞或現在分詞，BEZ代表is，AT為冠詞，NN為名詞，IN為介係詞，CC為連接詞，DT為限定詞，PP\$為所有格，JJ為形容詞。

# 詞性標註 Part-of-speech tagging

```
!pip list | grep nltk
```

```
import nltk
from nltk import word_tokenize, pos_tag

nltk.download('punkt')
nltk.download('averaged_perceptron_tagger')

print(pos_tag(word_tokenize("I'm learning NLP")))
# [('I', 'PRP'), ('"m", 'VBP'), ('learning', 'VBG'), ('NLP', 'NNP')]
```

# pos\_tag load the Standard treebank POS tagger

1. CC Coordinating conjunction
  2. CD Cardinal number
  3. DT Determiner
  4. EX Existential there
  5. FW Foreign word
  6. IN Preposition or subordinating conjunction
  7. JJ Adjective
  8. JJR Adjective, comparative
  9. JJS Adjective, superlative
- .....

請參閱<https://blog.csdn.net/fxjtoday/article/details/5841453>

# 詞性標註 Part-of-speech tagging

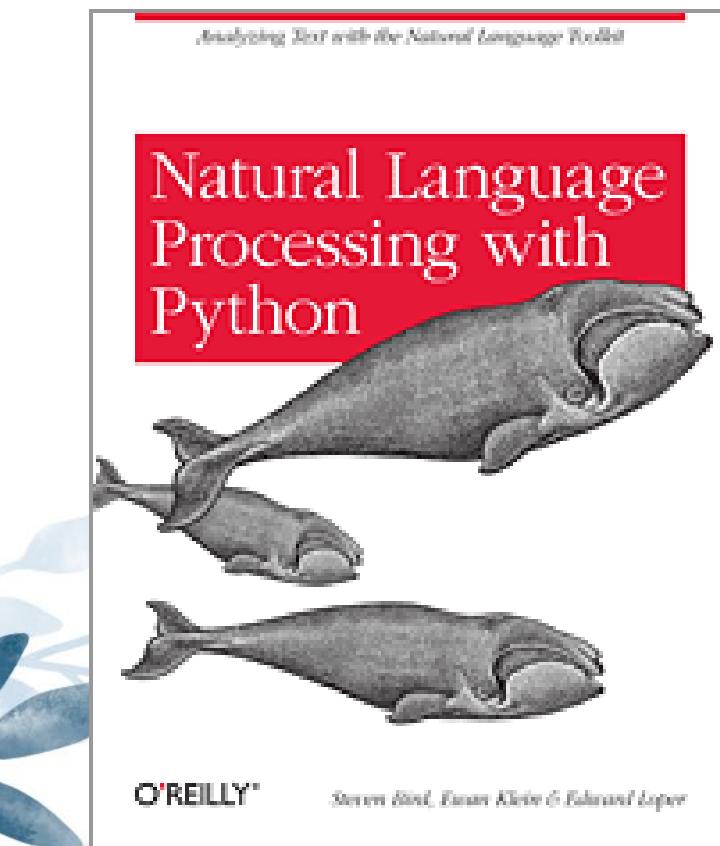
```
1 nltk.download()  
  
NLTK Downloader  
-----  
d) Download l) List u) Update c) Config h) Help q) Quit  
-----  
Downloader> l  
  
Packages:  
[ ] abc..... Australian Broadcasting Commission 2006  
[ ] alpino..... Alpino Dutch Treebank  
[*] averaged_perceptron_tagger Averaged Perceptron Tagger  
[ ] averaged_perceptron_tagger_ru Averaged Perceptron Tagger (Russian)  
[ ] basque_grammars.... Grammars for Basque  
[ ] biocreative_ppi..... BioCreAtIvE (Critical Assessment of Information  
Extraction Systems in Biology)  
[ ] bllip_wsj_no_aux.... BLLIP Parser: WSJ Model  
[ ] book_grammars..... Grammars from NLTK Book  
[ ] brown..... Brown Corpus  
[ ] brown_tei..... Brown Corpus (TEI XML Version)  
[ ] cess_cat..... CESS-CAT Treebank  
[ ] cess_esp..... CESS-ESP Treebank
```

# 詞性標註 Part-of-speech tagging

Complete guide for training your own Part-Of-Speech Tagger



# 詞性標註 Part-of-speech tagging



<http://www.nltk.org/book/>

## Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit

Steven Bird, Ewan Klein, and Edward Loper

This version of the NLTK book is updated for Python 3 and NLTK 3. The first edition of the book, published by O'Reilly, is available at [http://nltk.org/book\\_1ed/](http://nltk.org/book_1ed/). (There are currently no plans for a second edition of the book.)

0. [Preface](#)
1. [Language Processing and Python](#)
2. [Accessing Text Corpora and Lexical Resources](#)
3. [Processing Raw Text](#)
4. [Writing Structured Programs](#)
5. [Categorizing and Tagging Words](#) (minor fixes still required)
6. [Learning to Classify Text](#)
7. [Extracting Information from Text](#)
8. [Analyzing Sentence Structure](#)
9. [Building Feature Based Grammars](#)
10. [Analyzing the Meaning of Sentences](#) (minor fixes still required)
11. [Managing Linguistic Data](#) (minor fixes still required)
12. [Afterword: Facing the Language Challenge](#)

[Bibliography](#)  
[Term Index](#)

# 命名實體識別(Named Entity Recognition, NER)

在句子的序列中，定位並識別人名、地名、機構名等任務。

"There was nothing about this storm that was as expected," said **Jeff Masters**, a meteorologist and founder of **Weather Underground**. "**Irma** could have been so much worse. If it had traveled 20 miles north of the coast of **Cuba**, you'd have been looking at a (Category) 5 instead of a (Category) 3."

Person

Organization

Location

圖片來源：<https://blog.paralleldots.com/data-science/named-entity-recognition-milestone-models-papers-and-technologies/>

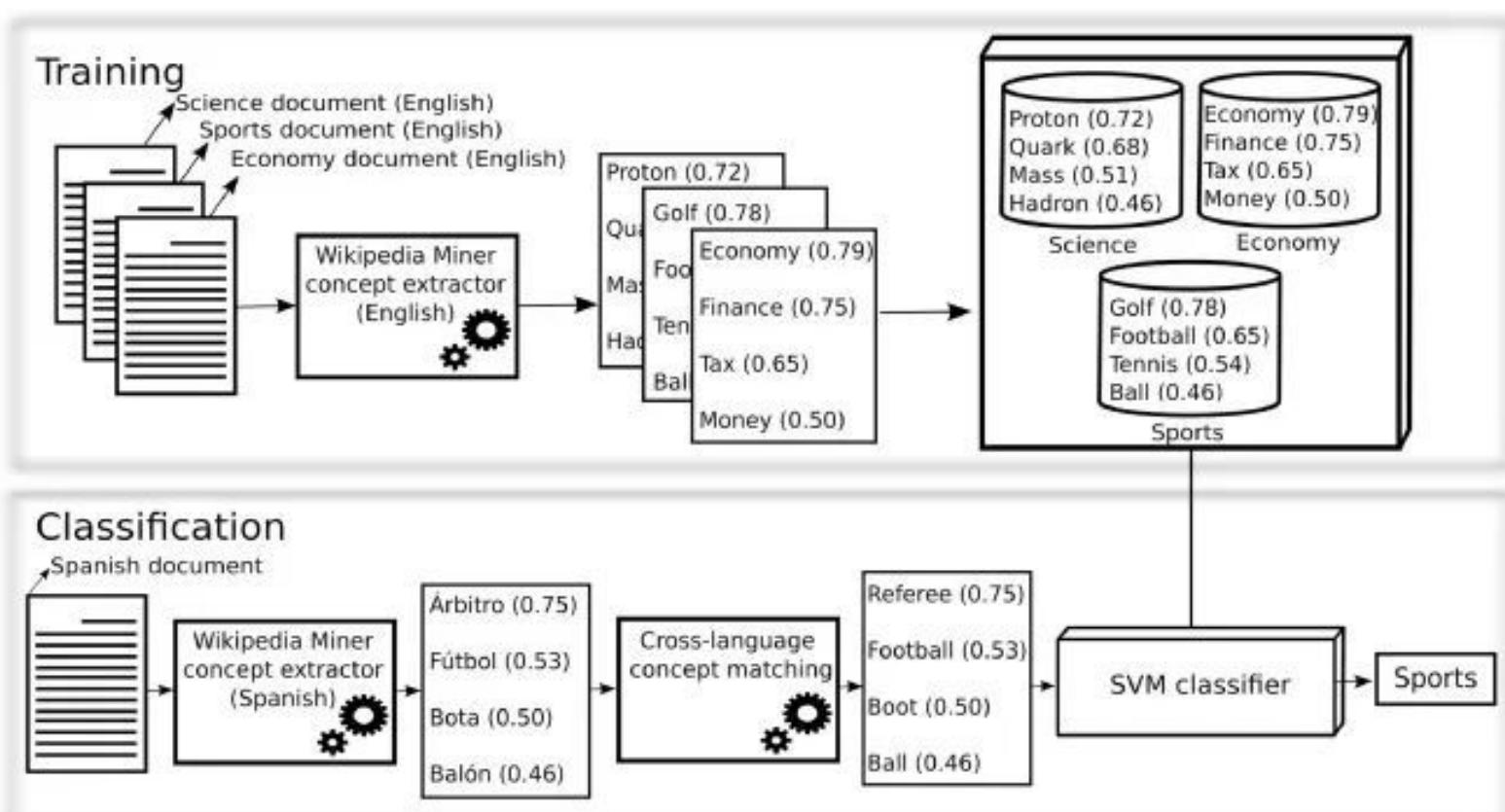
<https://ithelp.ithome.com.tw/articles/10209418>

# Text classification 文本分類

1. 詞嵌入向量化：word2vec, FastText等等
2. 卷積神經網路特徵提取：Text-CNN, Char-CNN等等
3. 上下文機制：Text-RNN， BiRNN， RCNN等等
4. 記憶存儲機制：EntNet， DMN等等
5. 注意力機制：HAN等等

# text classification 文本分類

傳統機器學習時代的主流分類器:Naive Bayes，Maximum Entropy，K-NN，和SVM。  
經典特徵模型：經典距離定義模型Vector Space Model (Rocchio)，N-grams等等。



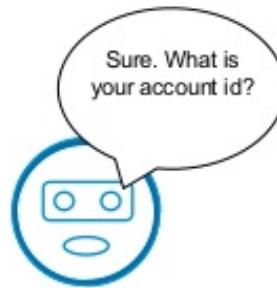
# Text classification 文本分類

1. 詞嵌入向量化：word2vec, FastText等等
2. 卷積神經網路特徵提取：Text-CNN, Char-CNN等等
3. 上下文機制：Text-RNN， BiRNN， RCNN等等
4. 記憶存儲機制：EntNet， DMN等等
5. 注意力機制：HAN等等

# Question-Answering system 智能問答

## Question-Answering Systems

- Systems that answer questions asked in natural language
- Goal-oriented (to complete one or more tasks):
  - Obtain information/facts
  - Seek recommendations
  - Make an appointment
  - Ask for help
  - ...



[AAAI 2019 tutorial] End-to-end goal-oriented question answering systems  
<https://www.slideshare.net/QiHe2/aaai-2019-tutorial-endtoend-goaloriented-question-answering-systems>

# Question-Answering system 智能問答



<https://github.com/l11x0m7/book-of-qna-code>

智能問答與深度學習  
王海良, 等 電子工業出版社  
2018-11-01

<https://www.tenlong.com.tw/products/9787121349218>

## Question And Answer Demo Using BERT NLP

### Paragraph \*

California, nicknamed the Googleplex. In August 2015, Google announced plans to reorganize its various interests as a conglomerate called Alphabet Inc. Google is Alphabet's leading subsidiary and will continue to be the umbrella company for Alphabet's Internet interests. Sundar Pichai was appointed CEO of Google, replacing Larry Page who became the CEO of Alphabet.

\*Maximum 1000 characters

### Question 1 \*

Who is current CEO?

Sundar Pichai

<https://www.pragnakalp.com/demos/BERT-NLP-QnA-Demo/>

面向文本分類的通用語言模型微調

Universal Language Model Fine-tuning for Text Classification

Jeremy Howard & Sebastian Ruder

<https://arxiv.org/abs/1801.06146>

<https://kknews.cc/code/6qqggj3.html>

<https://medium.com/mlreview/understanding-building-blocks-of-ulmfit-818d3775325b>

Inductive transfer learning has greatly impacted computer vision, but existing approaches in NLP still require task-specific modifications and training from scratch.

We propose **Universal Language Model Fine-tuning (ULMFiT)**, an effective transfer learning method that can be applied to any task in NLP, and introduce techniques that are key for fine-tuning a language model.

Our method significantly outperforms the state-of-the-art on six text classification tasks, reducing the error by 18-24% on the majority of datasets. Furthermore, with only 100 labeled examples, it matches the performance of training from scratch on 100x more data. **We open-source our pretrained models and code.**

# 機器翻譯

Machine Translation

Transformer(2017)

BERT(2018)

XLNet(2019)

ELMO, BERT, GPT

<https://www.youtube.com/watch?v=UYPa347-DdE>

# Seq2Seq(2014)

---

## Sequence to Sequence Learning with Neural Networks

---

**Ilya Sutskever**

Google

[ilyasu@google.com](mailto:ilyasu@google.com)

**Oriol Vinyals**

Google

[vinyals@google.com](mailto:vinyals@google.com)

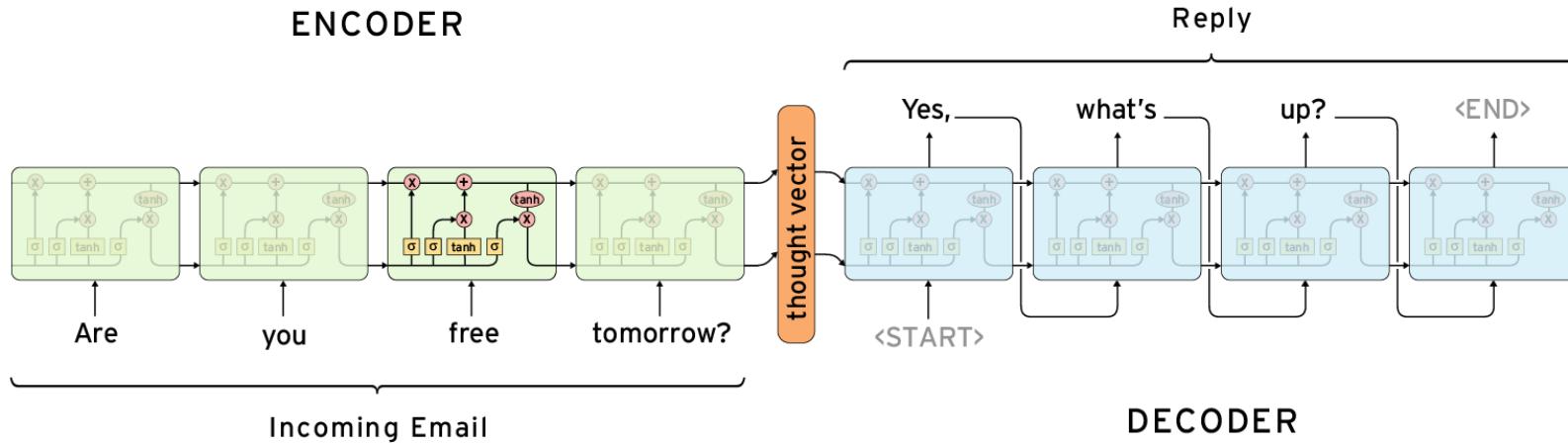
**Quoc V. Le**

Google

[qvl@google.com](mailto:qvl@google.com)



# Sequence to Sequence 是由 Encoder 與 Decoder 兩個 RNN 構成



運作原理:當看到一段話時，會先將這句話理解吸收，再根據我們理解的內容說出回覆，Sequence to Sequence 就是在模擬這個過程。

[https://zake7749.github.io/2017/09/28/Sequence-to-Sequence-tutorial/#Sequence\\_to\\_Sequence](https://zake7749.github.io/2017/09/28/Sequence-to-Sequence-tutorial/#Sequence_to_Sequence)

# Google Transformer(2017)

Attention Is All You Need

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
Lukasz Kaiser, Illia Polosukhin

(Submitted on 12 Jun 2017 (v1), last revised 6 Dec 2017 (this version, v5))

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an **encoder-decoder configuration**. The best performing models also connect the encoder and decoder through an **attention mechanism**.

We propose a new simple network architecture, **the Transformer**, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely.

Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

# pytorch-transformers

<https://github.com/huggingface/pytorch-transformers>



PASSED

PyTorch-Transformers (formerly known as `pytorch-pretrained-bert`) is a library of state-of-the-art pre-trained models for Natural Language Processing (NLP).

The library currently contains PyTorch implementations, pre-trained model weights, usage scripts and conversion utilities for the following models:

1. **BERT** (from Google) released with the paper [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#) by Jacob Devlin, Ming-Wei Chang, Kenton Lee and Kristina Toutanova.
2. **GPT** (from OpenAI) released with the paper [Improving Language Understanding by Generative Pre-Training](#) by Alec Radford, Karthik Narasimhan, Tim Salimans and Ilya Sutskever.
3. **GPT-2** (from OpenAI) released with the paper [Language Models are Unsupervised Multitask Learners](#) by Alec Radford\*, Jeffrey Wu\*, Rewon Child, David Luan, Dario Amodei\*\* and Ilya Sutskever\*\*.
4. **Transformer-XL** (from Google/CMU) released with the paper [Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context](#) by Zihang Dai\*, Zhilin Yang\*, Yiming Yang, Jaime Carbonell, Quoc V. Le, Ruslan Salakhutdinov.
5. **XLNet** (from Google/CMU) released with the paper [XLNet: Generalized Autoregressive Pretraining for Language Understanding](#) by Zhilin Yang\*, Zihang Dai\*, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, Quoc V. Le.
6. **XLM** (from Facebook) released together with the paper [Cross-lingual Language Model Pretraining](#) by Guillaume Lample and Alexis Conneau.
7. **RoBERTa** (from Facebook), a [Robustly Optimized BERT Pretraining Approach](#) by Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du et al.

These implementations have been tested on several datasets (see the example scripts) and should match the performances of the original implementations (e.g. ~93 F1 on SQuAD for BERT Whole-Word-Masking, ~88 F1 on RocStories for OpenAI GPT, ~18.3 perplexity on WikiText 103 for Transformer-XL, ~0.916 Pearson R coefficient on STS-B for XLNet). You can find more details on the performances in the Examples section of the documentation.

# pytorch-transformers

!pip install pytorch-transformers

```
import torch
from pytorch_transformers import *

MODELS = [ (BertModel,      BertTokenizer,      'bert-base-uncased'),
          (OpenAIGPTModel, OpenAIGPTTokenizer, 'openai-gpt'),
          (GPT2Model,       GPT2Tokenizer,       'gpt2'),
          (TransfoXLModel, TransfoXLTokenizer, 'transfo-xl-wt103'),
          (XLNetModel,     XLNetTokenizer,     'xlnet-base-cased'),
          (XLMModel,       XLMTokenizer,       'xlm-mlm-enfr-1024'),
          (RobertaModel,   RobertaTokenizer,   'roberta-base')]

for model_class, tokenizer_class, pretrained_weights in MODELS:
    print(model_class)
```

4  
5

```
1 import torch
2 from pytorch_transformers import *
3 MODELS = [ (BertModel, BertTokenizer, 'bert-base-uncased'),
4             (OpenAIGPTModel, OpenAIGPTTokenizer, 'openai-gpt'),
5             (GPT2Model, GPT2Tokenizer, 'gpt2'),
6             (TransfoXLModel, TransfoXLTokenizer, 'transfo-xl-wt103'),
7             (XLNetModel, XLNetTokenizer, 'xlnet-base-cased'),
8             (XLMModel, XLMTokenizer, 'xlm-mlm-enfr-1024'),
9             (RobertaModel, RobertaTokenizer, 'roberta-base') ]
10 for model_class, tokenizer_class, pretrained_weights in MODELS:
11     print(model_class)
12
<class 'pytorch_transformers.modeling_bert.BertModel'>
<class 'pytorch_transformers.modeling_openai.OpenAIGPTModel'>
<class 'pytorch_transformers.modeling_gpt2.GPT2Model'>
<class 'pytorch_transformers.modeling_transfo_xl.TransfoXLModel'>
<class 'pytorch_transformers.modeling_xlnet.XLNetModel'>
<class 'pytorch_transformers.modeling_xlm.XLMModel'>
<class 'pytorch_transformers.modeling_roberta.RobertaModel'>
```

# NLP學習資源

# <https://github.com/graykode/nlp-tutorial>

## 1. Basic Embedding Model

- 1-1. [NNLM\(Neural Network Language Model\)](#) - Predict Next Word
  - Paper - [A Neural Probabilistic Language Model\(2003\)](#)
  - Colab - [NNLM\\_Tensor.ipynb](#), [NNLM\\_Torch.ipynb](#)
- 1-2. [Word2Vec\(Skip-gram\)](#) - Embedding Words and Show Graph
  - Paper - [Distributed Representations of Words and Phrases and their Compositionality\(2013\)](#)
  - Colab - [Word2Vec\\_Tensor\(NCE\\_loss\).ipynb](#), [Word2Vec\\_Tensor\(Softmax\).ipynb](#), [Word2Vec\\_Torch\(Softmax\).ipynb](#)
- 1-3. [FastText\(Application Level\)](#) - Sentence Classification
  - Paper - [Bag of Tricks for Efficient Text Classification\(2016\)](#)
  - Colab - [FastText.ipynb](#)

## 2. CNN(Convolutional Neural Network)

- 2-1. [TextCNN](#) - Binary Sentiment Classification
  - Paper - [Convolutional Neural Networks for Sentence Classification\(2014\)](#)
  - Colab - [TextCNN\\_Tensor.ipynb](#), [TextCNN\\_Torch.ipynb](#)
- 2-2. DCNN(Dynamic Convolutional Neural Network)

## 3. RNN(Recurrent Neural Network)

- 3-1. [TextRNN](#) - Predict Next Step
  - Paper - [Finding Structure in Time\(1990\)](#)
  - Colab - [TextRNN\\_Tensor.ipynb](#), [TextRNN\\_Torch.ipynb](#)
- 3-2. [TextLSTM](#) - Autocomplete
  - Paper - [LONG SHORT-TERM MEMORY\(1997\)](#)
  - Colab - [TextLSTM\\_Tensor.ipynb](#), [TextLSTM\\_Torch.ipynb](#)
- 3-3. [Bi-LSTM](#) - Predict Next Word in Long Sentence
  - Colab - [Bi\\_LSTM\\_Tensor.ipynb](#), [Bi\\_LSTM\\_Torch.ipynb](#)

## 4. Attention Mechanism

<http://web.stanford.edu/class/cs224u/index.html>

<https://github.com/cgpotts/cs224u/>



The banner features a dark red background with white text. At the top, there are three navigation links: "Home", "Projects", and "Policies and requirements". Below these, the course title "CS224U: Natural Language Understanding" is displayed in a large, white, sans-serif font.

#### Course info

MW 4:30–5:50 pm in [Gates B1](#)

[Piazza site](#)

[Canvas site](#)

[Github](#)

All course correspondence:

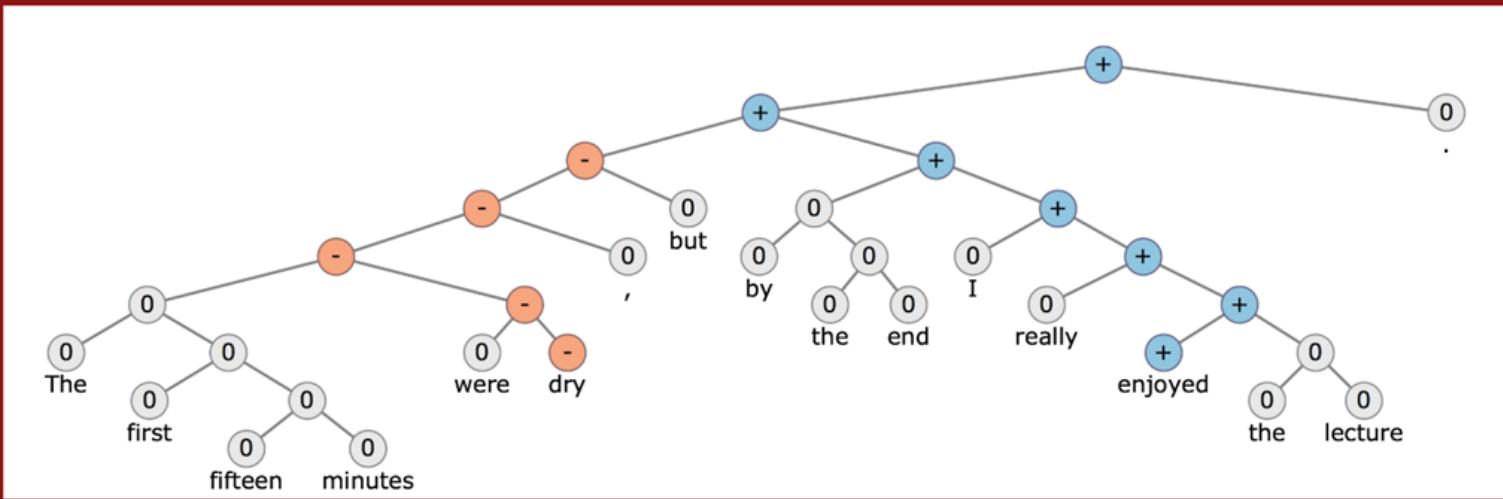
## Schedule

Plan	Resources	Assignments
<p><a href="#">Course overview</a> <a href="#">[video]</a> <a href="#">Course set-up</a></p> <p>Apr 1 <a href="#">Jupyter notebook tutorial</a></p>	<p><a href="#">Mitchell 2004</a> <a href="#">Levesque 2013</a> <a href="#">Manning 2015</a></p>	<p><a href="#">HW 1 due</a> Apr 15, 4:30 pm</p>

<http://cs224d.stanford.edu/>



## CS224d: Deep Learning for Natural Language Processing

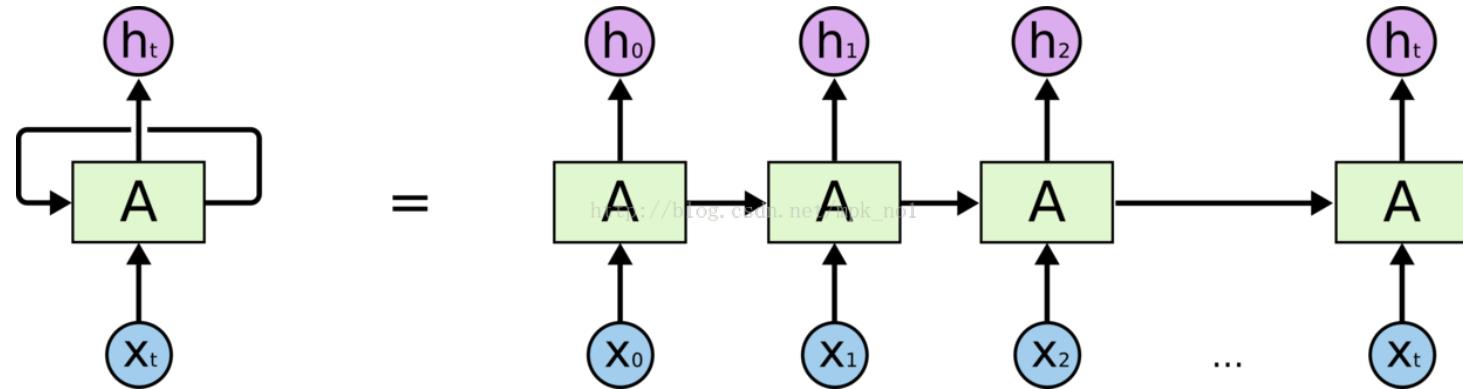


遞歸神經網絡

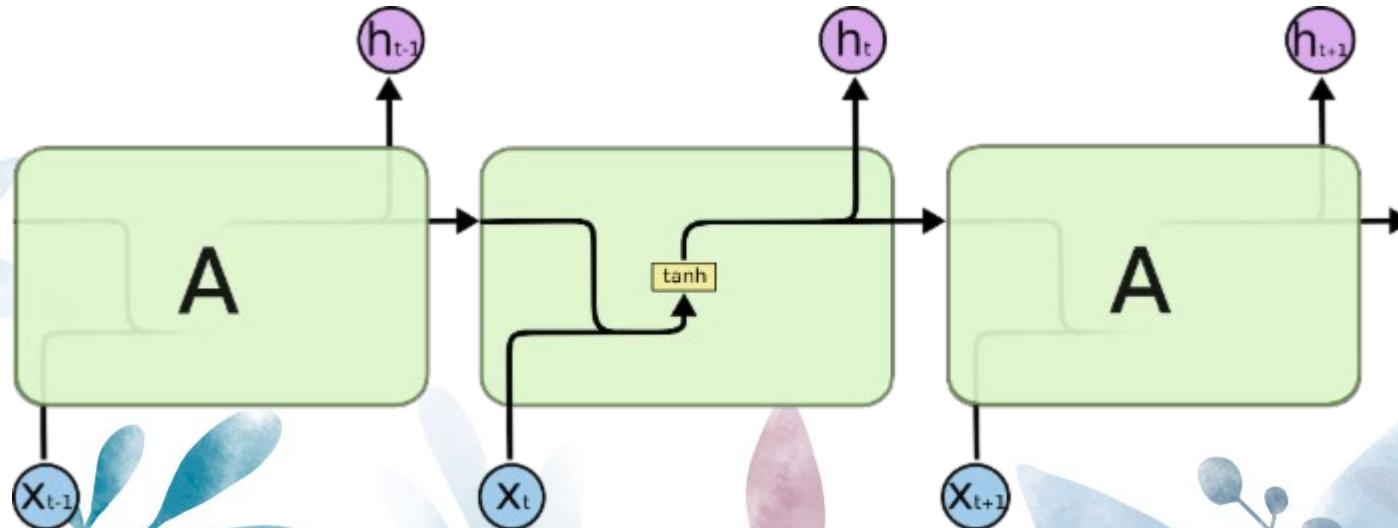
Recurrent  
Neural Network

模型

# SimpleRNN

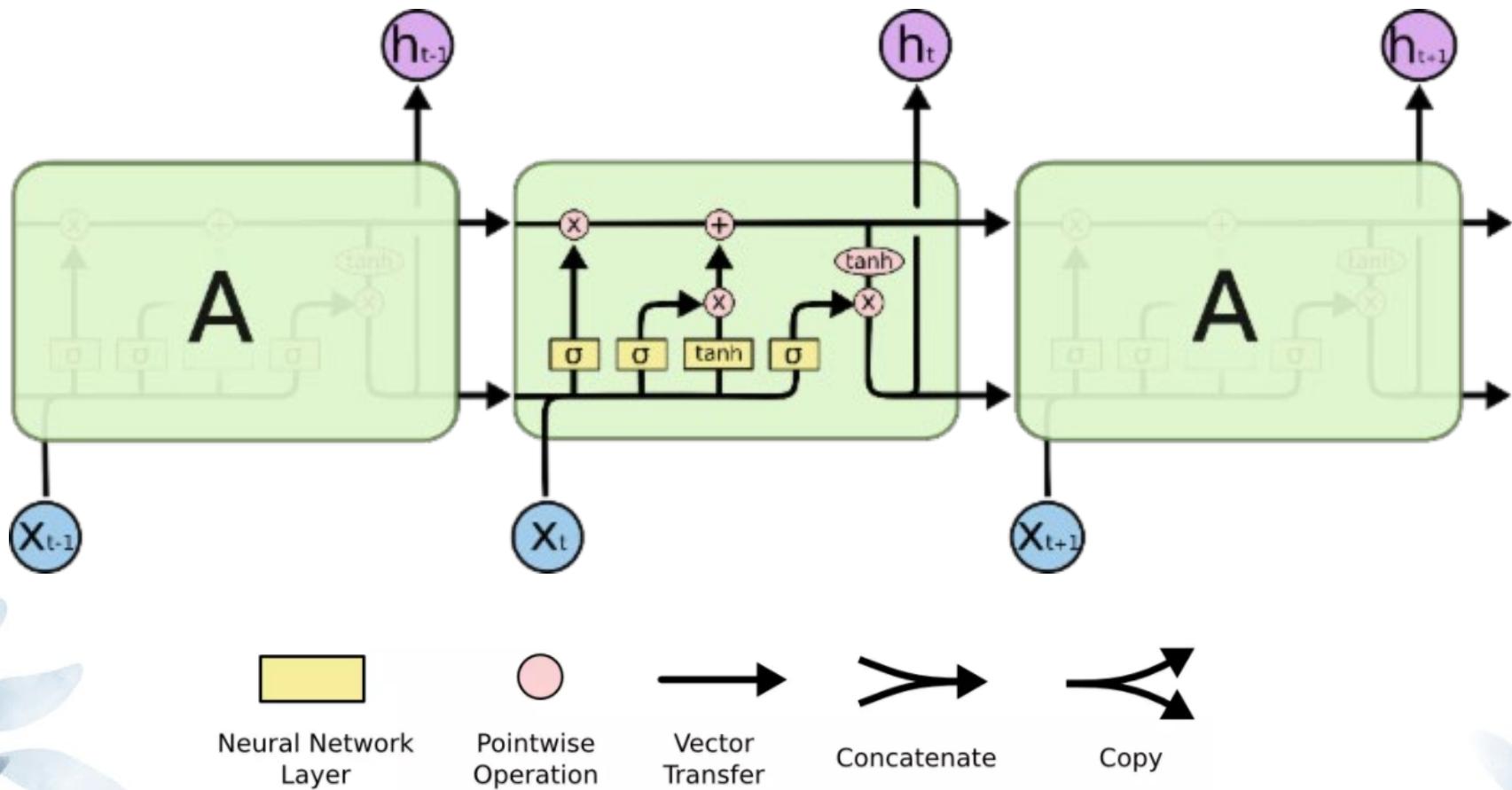


所有迴圈神經網路結構都是由完全相同結構的（神經網路）模組進行複製而成的。  
在普通的RNNs中，這個模組結構非常簡單，比如僅是一個單一的  $\tanh$  層



# LSTM

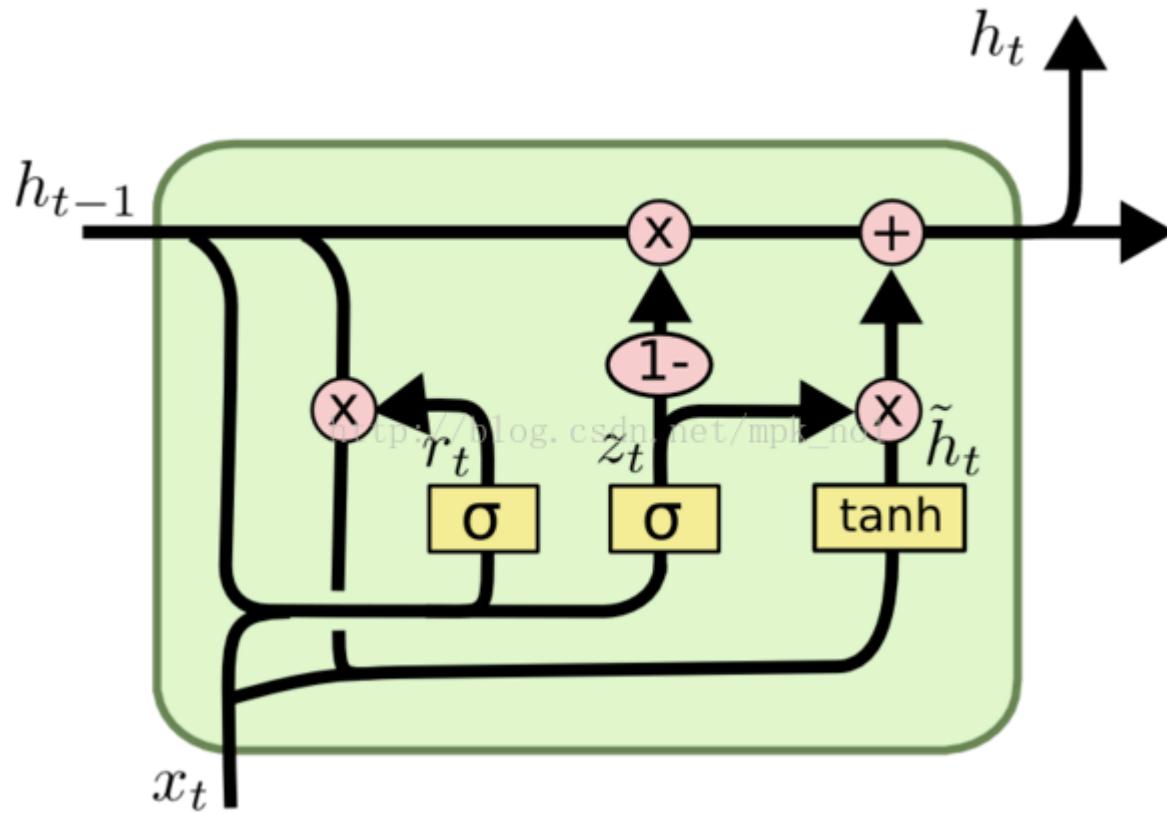
LSTMs 也有類似的結構,但是它們不再只是用一個單一的  $\tanh$  層, 而是用了四個相互作用的層



[https://brohrer.mcknote.com/zh-Hant/how\\_machine\\_learning\\_works/how\\_rnns\\_lstm\\_work.html](https://brohrer.mcknote.com/zh-Hant/how_machine_learning_works/how_rnns_lstm_work.html)

<https://blog.csdn.net/fendouaini/article/details/80198994>

# GRU



[https://blog.csdn.net/mpk\\_no1/article/details/72875185](https://blog.csdn.net/mpk_no1/article/details/72875185)

# Understanding LSTM Networks

*Posted on August 27, 2015*

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

<https://blog.csdn.net/menc15/article/details/71271566>

[https://blog.csdn.net/Jerr\\_y/article/details/58598296](https://blog.csdn.net/Jerr_y/article/details/58598296)



Andrej Karpathy blog

About

Hacker's guide to Neural Networks

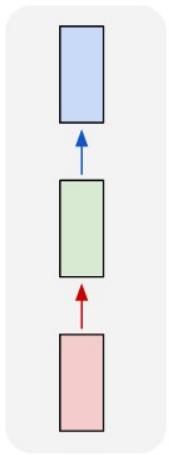
# The Unreasonable Effectiveness of Recurrent Neural Networks

*May 21, 2015*

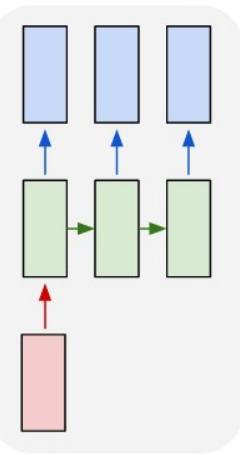
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

<https://github.com/karpathy/char-rnn>

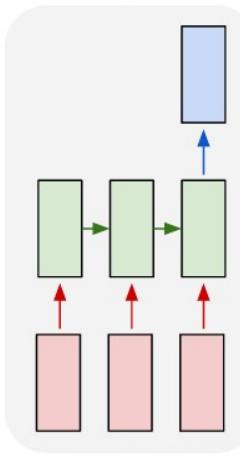
one to one



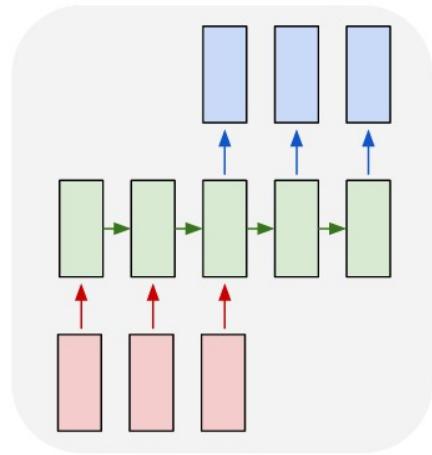
one to many



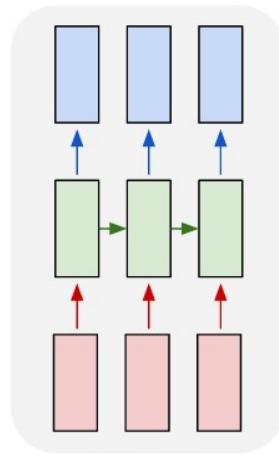
many to one



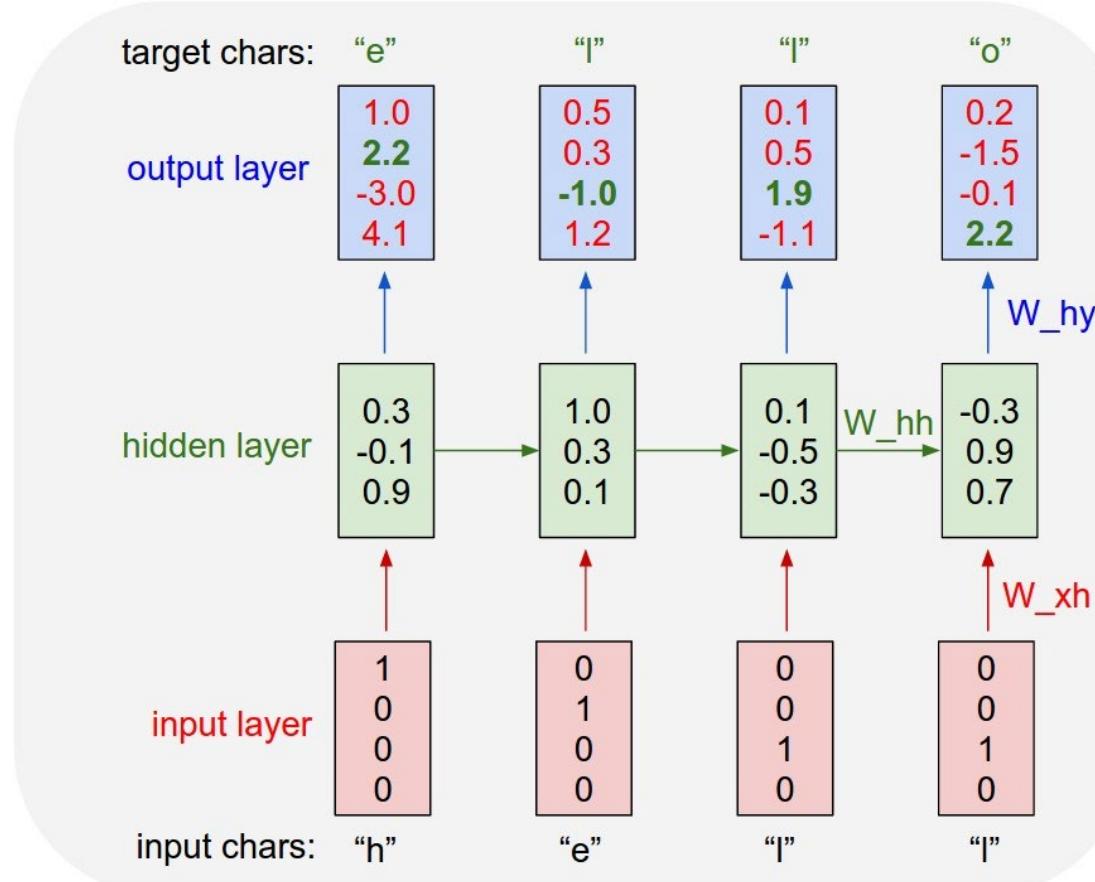
many to many

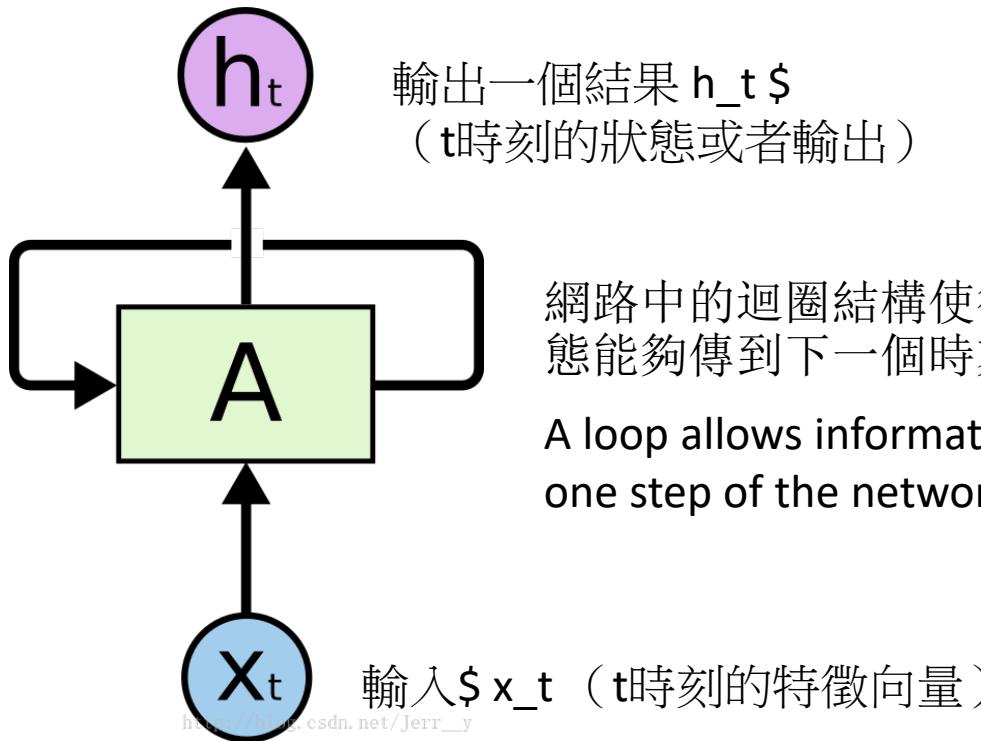


many to many

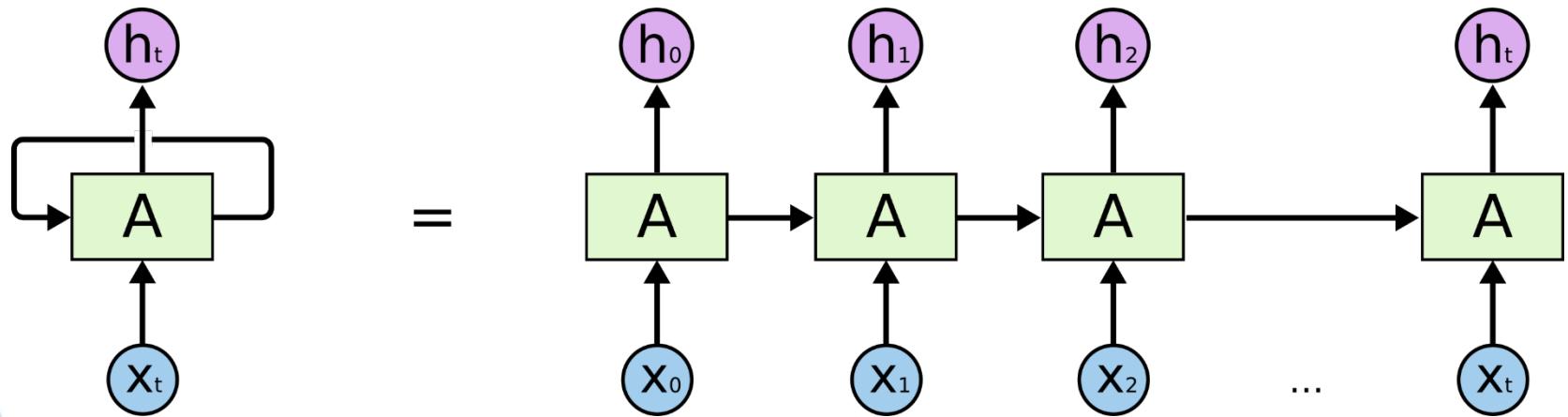


# An example RNN with 4-dimensional input and output layers, and a hidden layer of 3 units (neurons)





A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor.



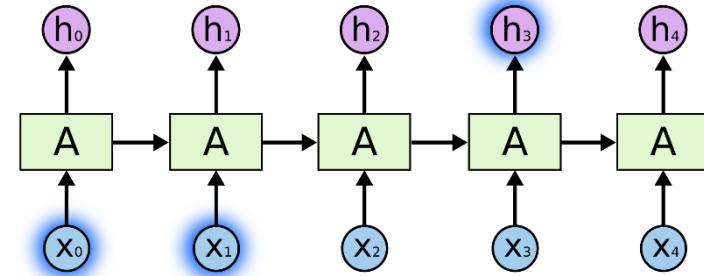
This chain-like nature reveals that recurrent neural networks are intimately related to sequences and lists.

They're the natural architecture of neural network to use for such data.

# The Problem of Long-Term Dependencies

長時期依賴的問題

consider a language model trying to predict the next word based on the previous ones.



the clouds are in the ?

the clouds are in the **sky**

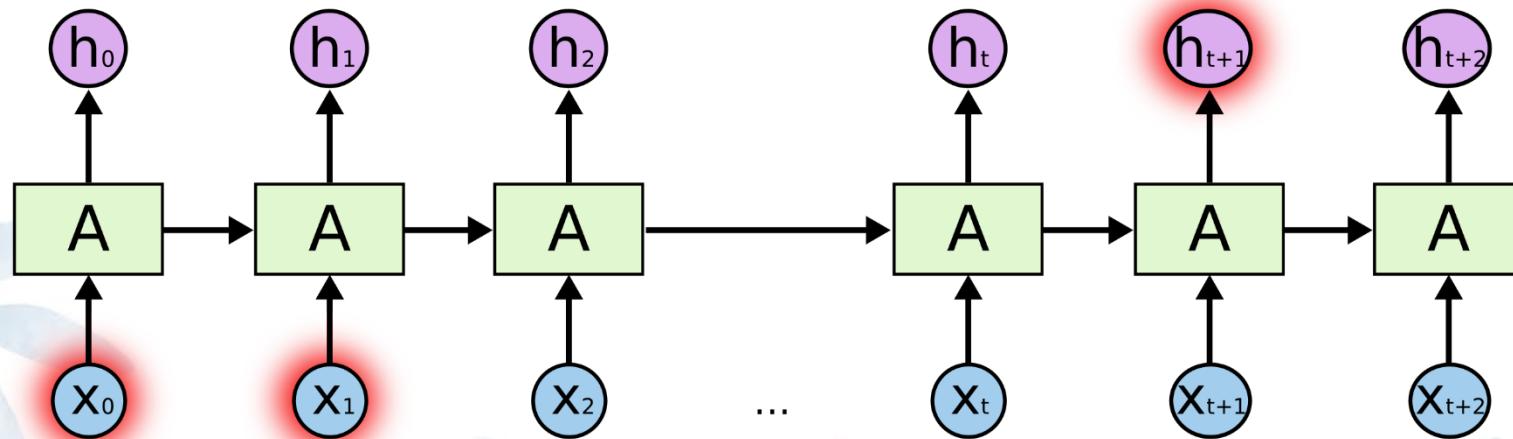
不需要更多的資訊，我們就能夠自然而然的  
想到下一個詞應該是“sky”。

所要預測的內容和相關資訊之間的間隔很小，這種情況下 RNNs 就能夠利用過去的資訊，很容易的實現

# 長時期依賴的問題

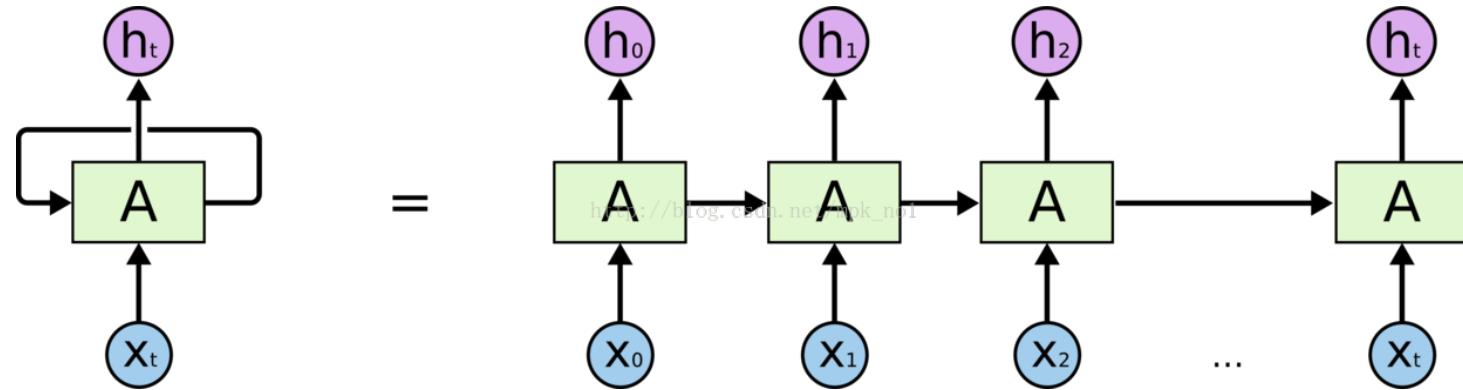
預測 “I grew up in France ... (此處省略1萬字)... I speak ?”

“I grew up in France... I speak fluent **French**. ”

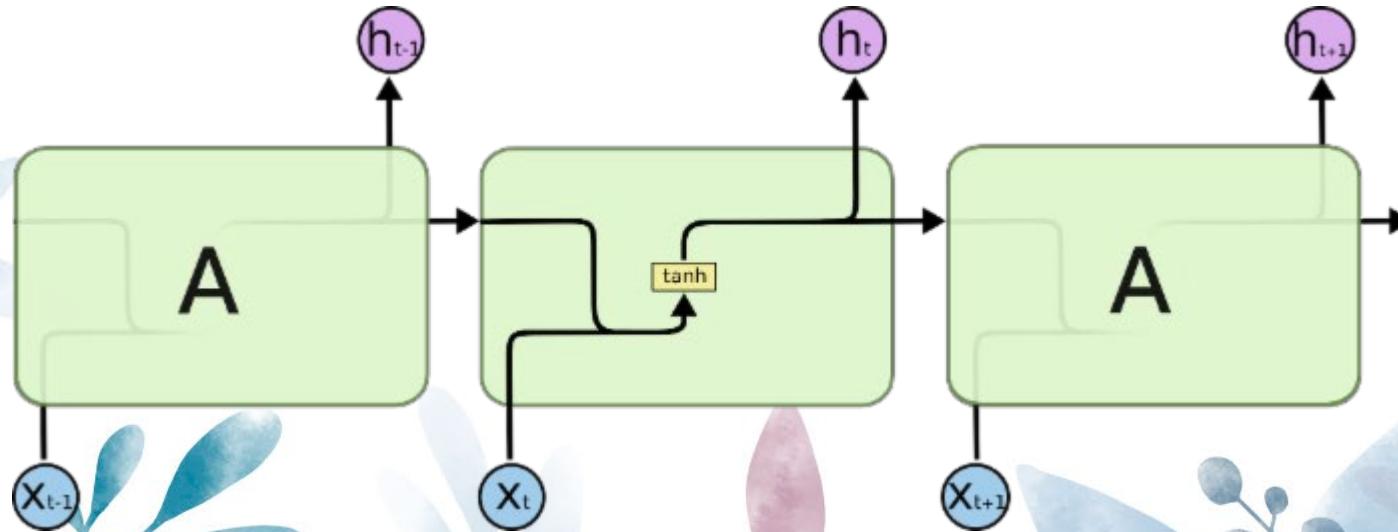


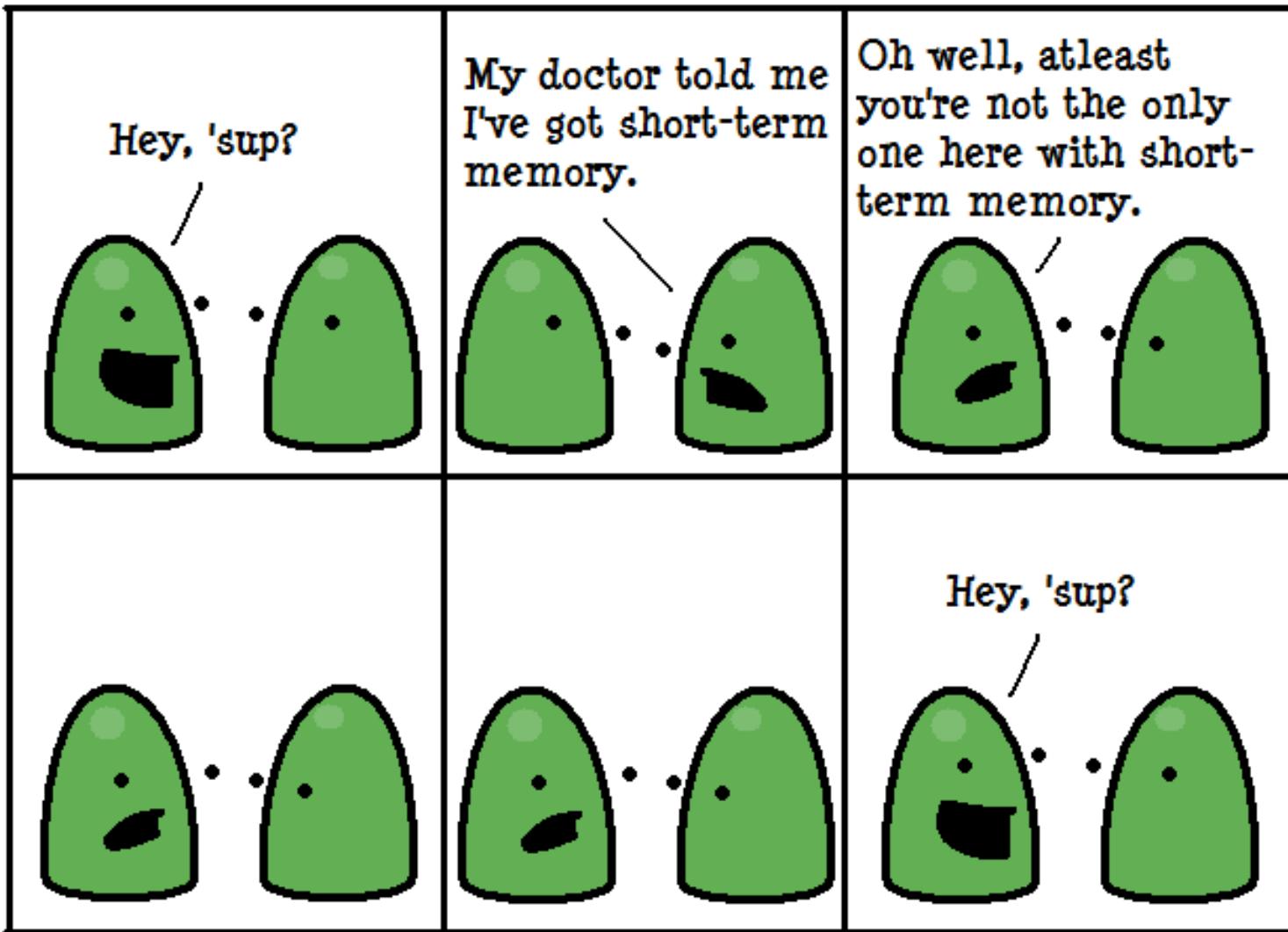
預測資訊和相關資訊間的間隔增大， RNNs 很難去把它們關聯起來

# SimpleRNN



所有迴圈神經網路結構都是由完全相同結構的（神經網路）模組進行複製而成的。  
在普通的RNNs中，這個模組結構非常簡單，比如僅是一個單一的  $\tanh$  層





<https://www.deviantart.com/toxicpaprika/art/Short-term-memory-88267010>

# LONG SHORT-TERM MEMORY

NEURAL COMPUTATION 9(8):1735–1780, 1997

Sepp Hochreiter

Fakultät für Informatik

Technische Universität München

80290 München, Germany

[hochreit@informatik.tu-muenchen.de](mailto:hochreit@informatik.tu-muenchen.de)

<http://www7.informatik.tu-muenchen.de/~hochreit>

Jürgen Schmidhuber

IDSIA

Corso Elvezia 36

6900 Lugano, Switzerland

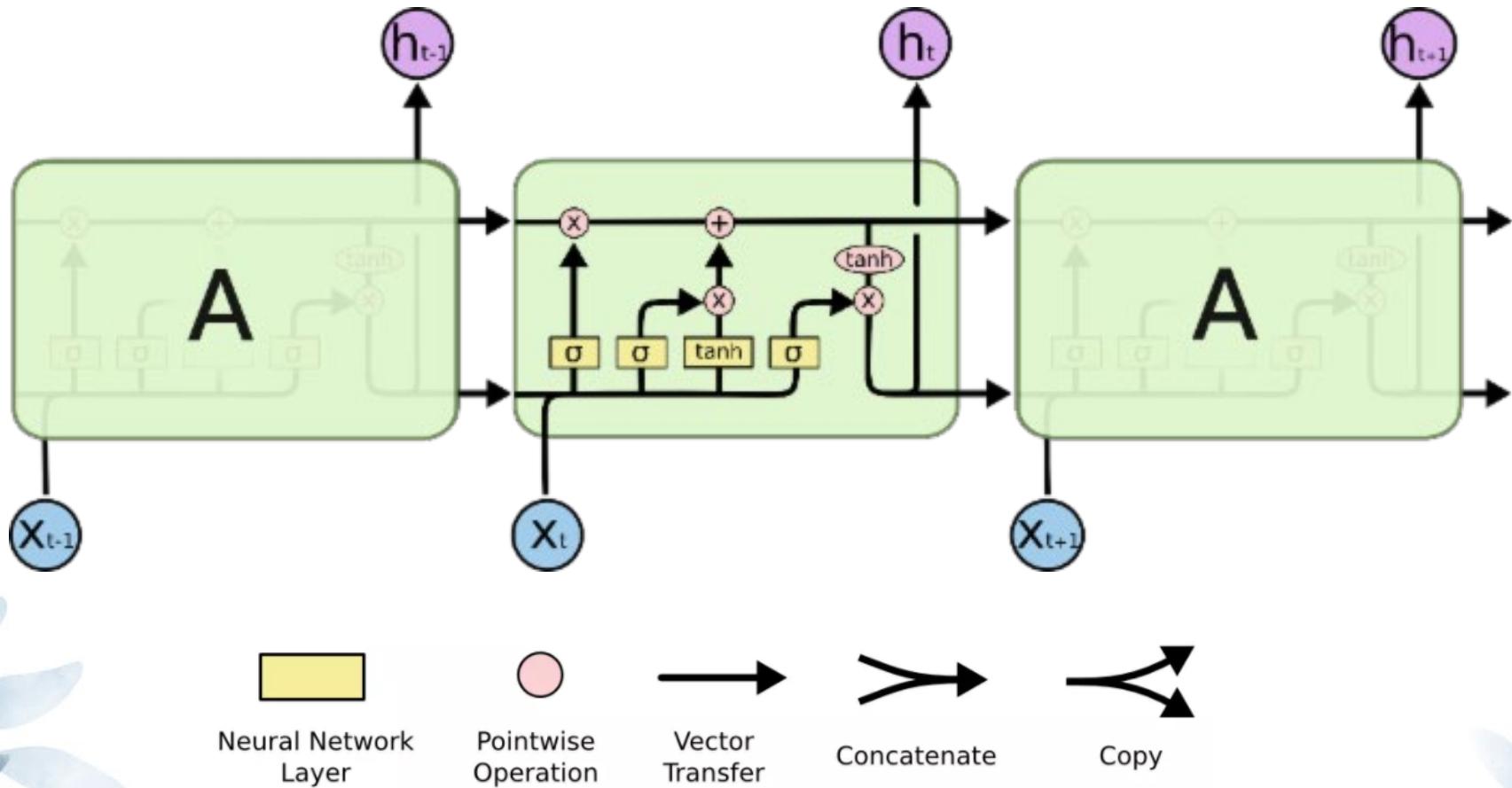
[juergen@idsia.ch](mailto:juergen@idsia.ch)

<http://www.idsia.ch/~juergen>

設計初衷是希望能夠解決神經網路中的長期依賴問題，讓記住長期資訊成為神經網路的預設行為，而不是需要很大力氣才能學會。

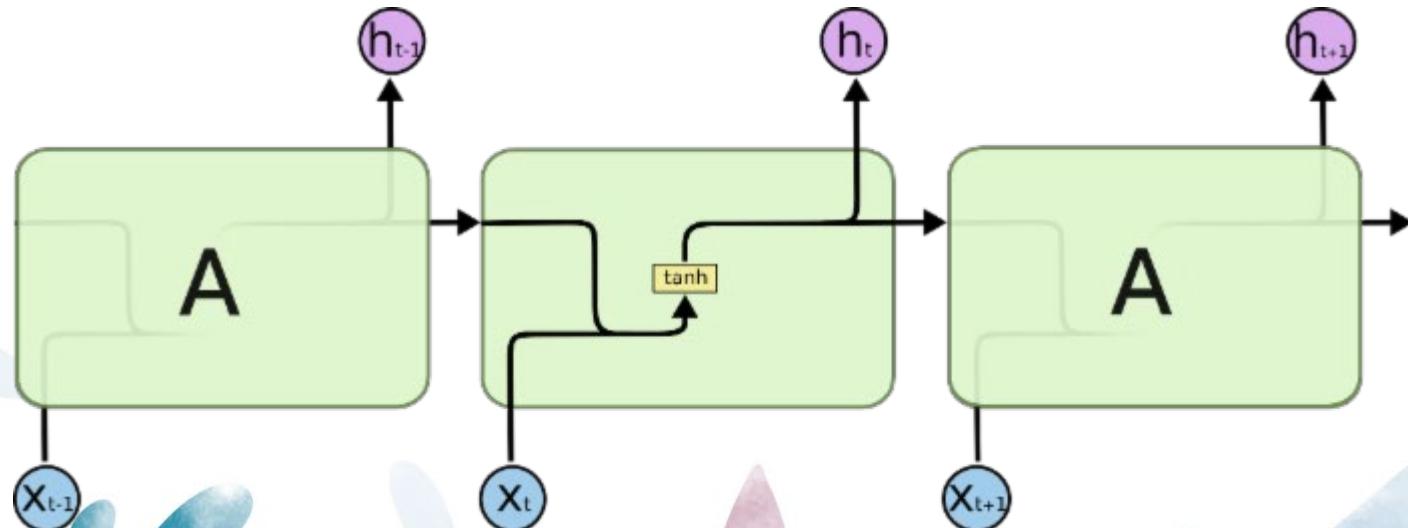
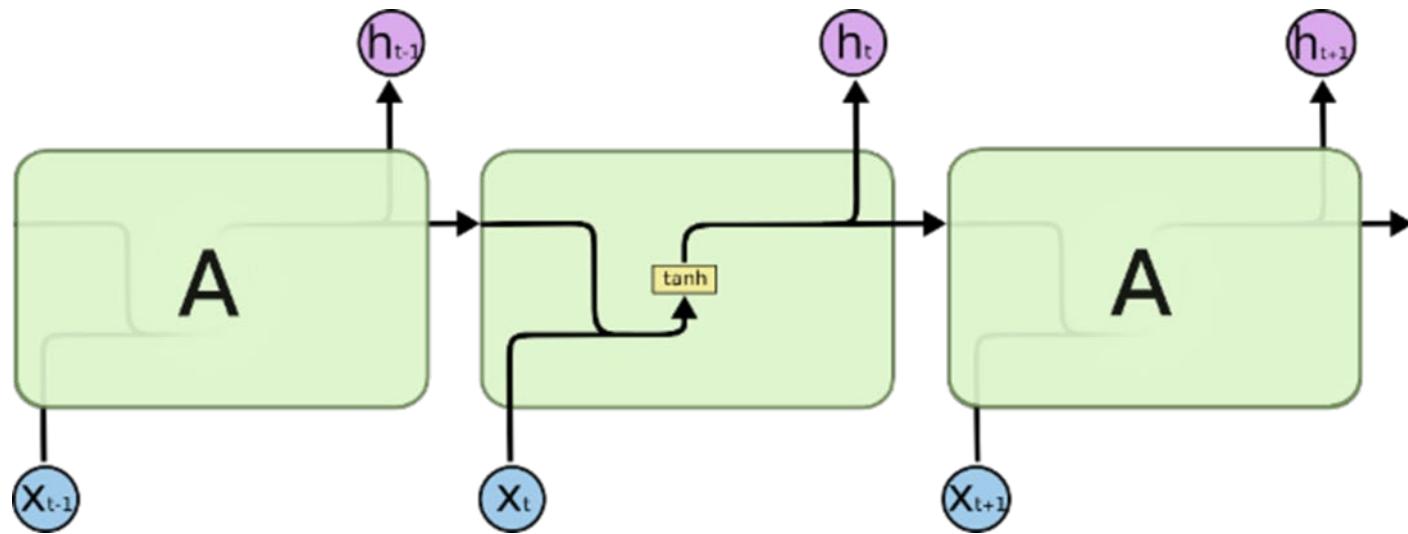
# LSTM

LSTMs 也有類似的結構,但是它們不再只是用一個單一的  $\tanh$  層, 而是用了四個相互作用的層



[https://brohrer.mcknote.com/zh-Hant/how\\_machine\\_learning\\_works/how\\_rnns\\_lstm\\_work.html](https://brohrer.mcknote.com/zh-Hant/how_machine_learning_works/how_rnns_lstm_work.html)

<https://blog.csdn.net/fendouaini/article/details/80198994>

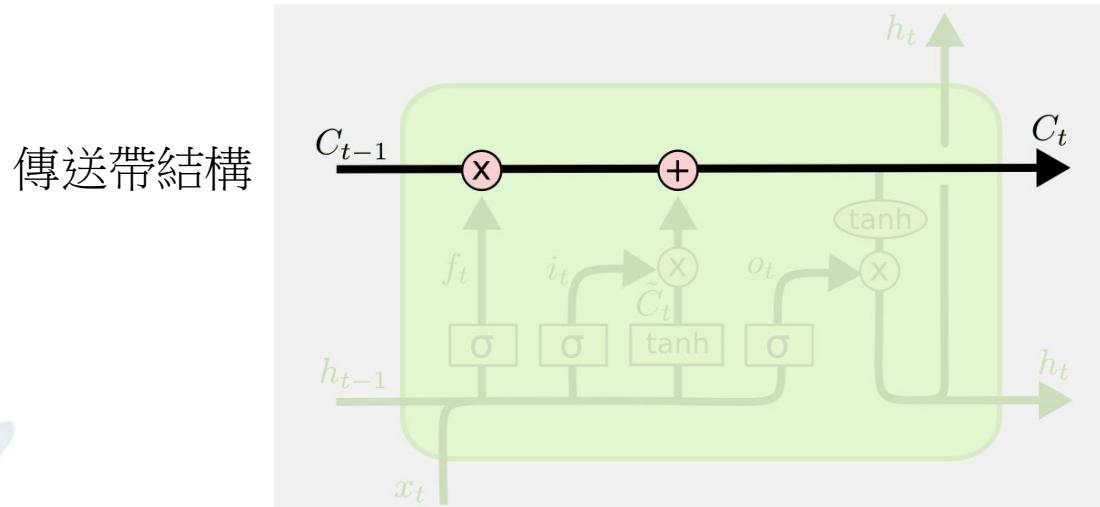


# cell state

長時期的記憶保留

LSTMs 最關鍵的地方在於 cell 狀態 和 結構圖上面的那條橫穿的水平線。

cell 狀態的傳輸就像一條傳送帶，向量從整個 cell 中穿過，只是做了少量的線性操作。這種結構能夠很輕鬆地實現資訊從整個 cell 中穿過而不做改變。  
(這樣就可以實現長時期的記憶保留)



## 使用Gate來添加或者刪除資訊

每個 LSTM 有三個這樣的門結構，來實現保護和控制資訊。  
(分別是 “forget gate layer”, 遺忘門； “input gate layer”, 傳入門；  
“output gate layer”, 輸出門)

## 遺忘門 (forget gates)

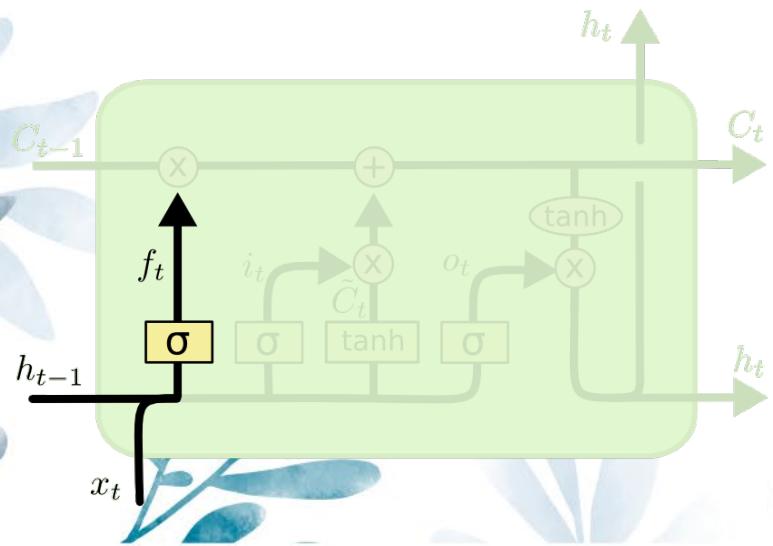
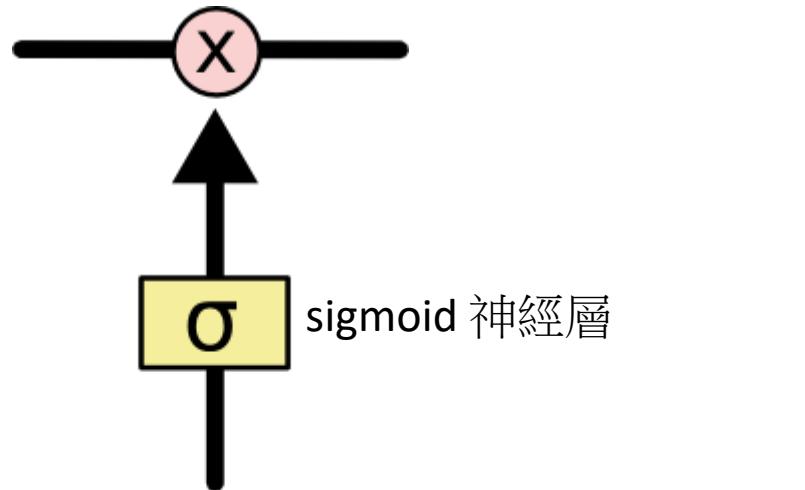
決定要讓那些資訊繼續通過這個 cell?

sigmoid 神經層

輸出是一個數值都在 0, 1 之間的向量（向量長度和 cell 的狀態  $C_{t-1}$  一樣），表示讓  $C_{t-1}$  的各部分資訊通過的比重。

0 表示“不讓任何資訊通過”，1 表示“讓所有資訊通過”。

pointwise multiplication operation



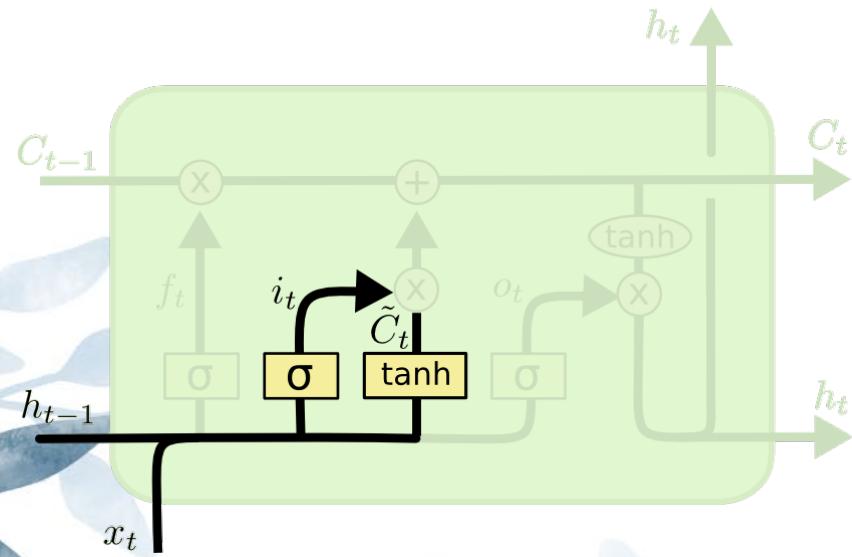
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

# 傳入門 (input gates)

決定讓多少新的資訊加入到 cell 狀態

[1] a sigmoid layer called the “input gate layer” decides which values we’ll update.

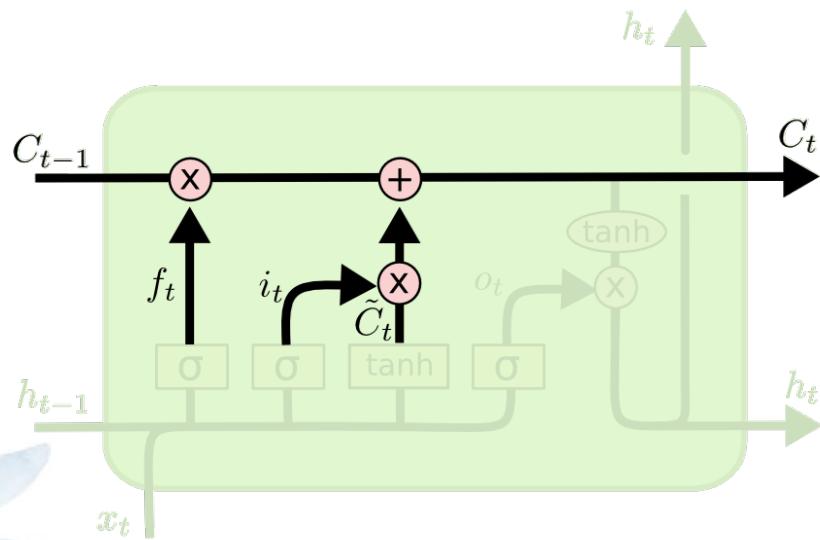
[2] a tanh layer creates a vector of new candidate values,  $C_t$ , that could be added to the state



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

combine these two to create an update to the state



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

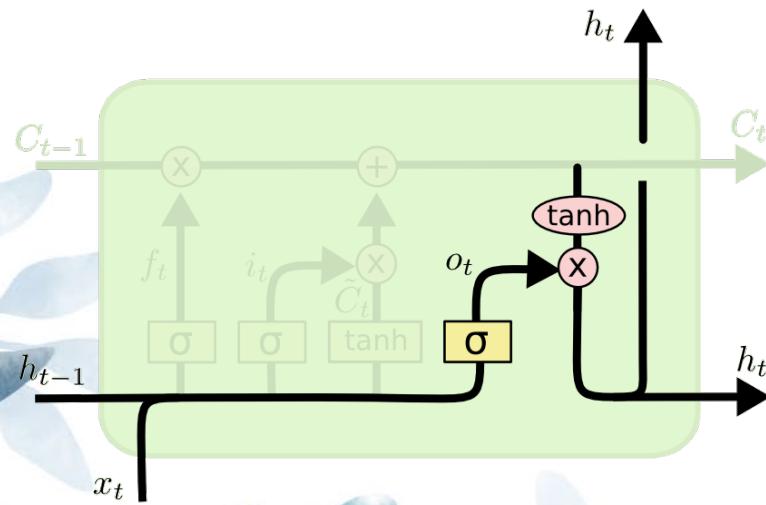
# 輸出門

決定輸出什麼值

This output will be based on our cell state, but will be a filtered version

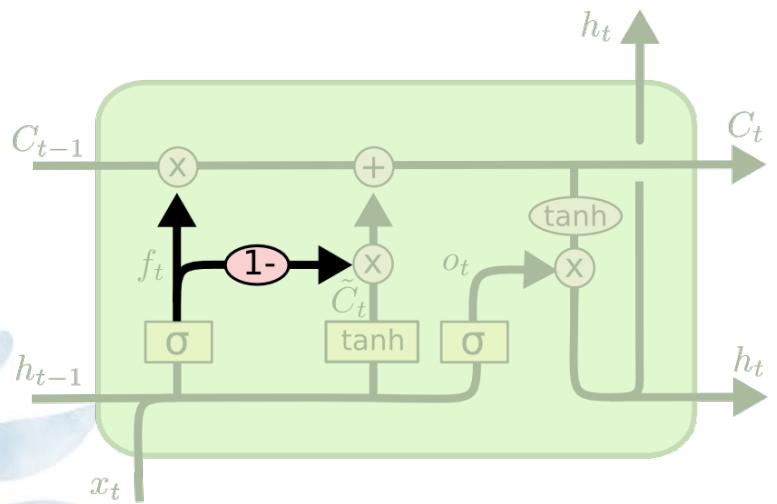
First, we run a sigmoid layer which decides what parts of the cell state we're going to output.

Then, we put the cell state through tanh (to push the values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.



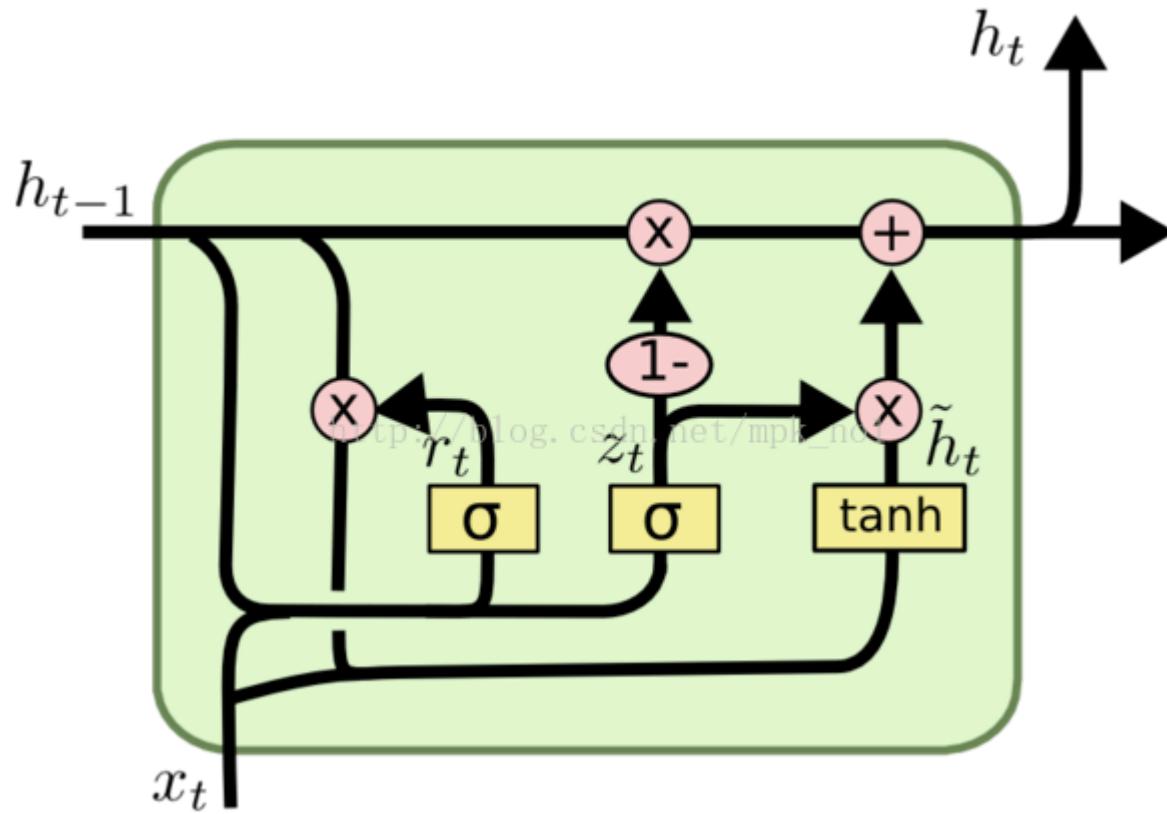
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh (C_t)$$

# Depth-Gated Recurrent Neural Networks



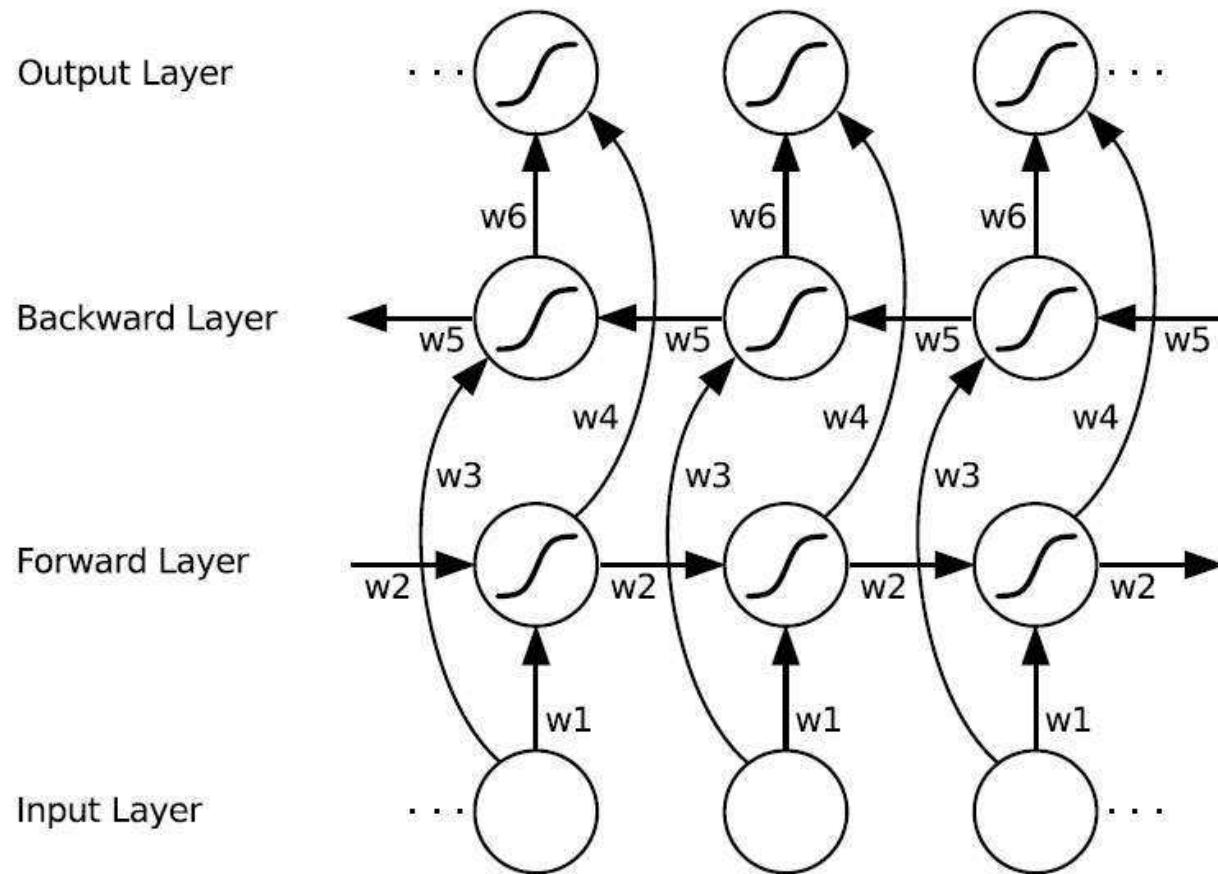
$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

# GRU



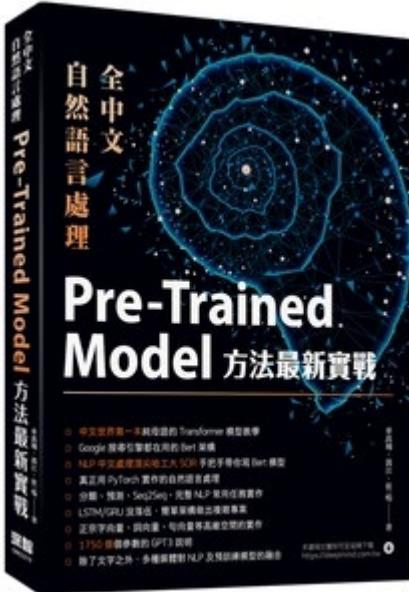
[https://blog.csdn.net/mpk\\_no1/article/details/72875185](https://blog.csdn.net/mpk_no1/article/details/72875185)

# 雙向LSTM(Bi-directional LSTM)



# 預訓練自然語言模型 (Pre-trained Models for NLP)

# Google的關鍵字排名、Google翻譯、購物網站推薦系統、Siri/OK Google



- ★ Word2vec 詞向量、GloVe 詞向量
- ★ 靜態詞向量預訓練模型
- ★ 動態詞向量預訓練模型  
以語言模型為基礎的動態詞向量預訓練
- ★ 預訓練語言模型、GPT、**BERT**
- ★ 模型蒸餾與壓縮、DistilBERT、TinyBERT、MobileBERT、TextBrewer
- ★ 生成模型、BART、UniLM、T5、**GPT-3**、可控文字生成
- ★ 多語言融合、多媒體融合、異質知識融合
- ★ VideoBERT、VL-BERT、DALL·E、ALIGN

[https://www.tenlong.com.tw/products/9789860776942?list\\_name=srh](https://www.tenlong.com.tw/products/9789860776942?list_name=srh)

<https://www.ithome.com.tw/tags/gpt-3>

ithome 新聞 產品&技術 專題 AI Cloud DevOps 資安 研討會 社群 IT EXP

# GPT-3

## New GPT-3 Capabilities: Edit & Insert

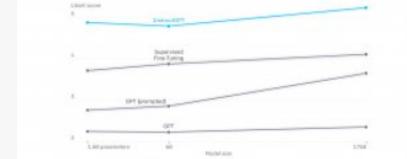
We've released new versions of GPT-3 and Codex which can edit or insert content into existing text, rather than just completing existing text. These new capabilities make it practical to use the OpenAI API to revise existing content, such as rewriting a paragraph of text or refactoring code. This unlocks new use cases and improves existing ones; for example, insertion is already being piloted in GitHub Copilot with promising early results.

| OpenAI | 語言模型 | GPT-3 | Codex

## 語言模型GPT-3現新增文字編輯和插入功能

語言模型GPT-3現在可以參考文字前後段落，插入連結兩段落的文字，或是修改現有內容的語氣和結構

2022-03-16



| OpenAI | GPT-3 | NLP

## OpenAI改進GPT-3使其更能聽懂人類指示，並減少輸出有毒內容

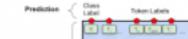
OpenAI API服務現在使用經微調過的GPT-3模型InstructGPT，與原始GPT-3相比，錯誤率更低並且產生幻覺

2022-01-28

Google AI Blog

A Fast WordPiece Tokenization System

Using InstructGPT, Multi-Modal Inference, and Text-to-Image, Google Research has developed a fast wordpiece tokenizer that is based on the latest research in neural network architecture and training techniques.



| AI | google | 斷詞 | 英特爾 | oneAPI | IBM | 對抗攻擊 | 商湯科技 | OpenAI | GPT-3 | IT周報

## AI趨勢周報第183期：斷詞等好久？Google開源新工具加速8倍工作時間

Google開源一套英文斷詞新工具，比現有方法快8倍；英特爾免費釋出跨架構工具包，一口氣整合900多項技術優化...-01-06

## Customizing GPT-3 for Your Application

Fine-tune with a single command

| OpenAI | GPT-3 | API

## OpenAI API服務供用戶自定義GPT-3模型

OpenAI API用戶現在只要準備額外的資料集，就可以使用OpenAI工具自定義GPT-3，並透過API取用自定義後的模型

2021-12-15

## Build next-gen apps with OpenAI's powerful models.

OpenAI's API provides access to GPT-3, which performs a wide variety of natural language tasks, and Codex, which translates natural language to code.

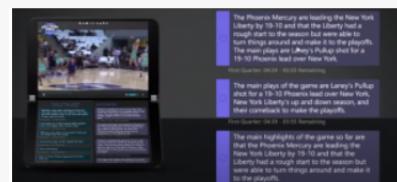
GET STARTED READ DOCUMENTATION

| OpenAI | GPT-3 | 自然語言

## OpenAI擴大開放以GPT-3為基礎的語言生成API

用戶現在不需要申請排隊，就可以使用OpenAI所開發，以GPT-3為基礎的語言生成API，官方還釋出內容過濾器，助用

2021-11-19

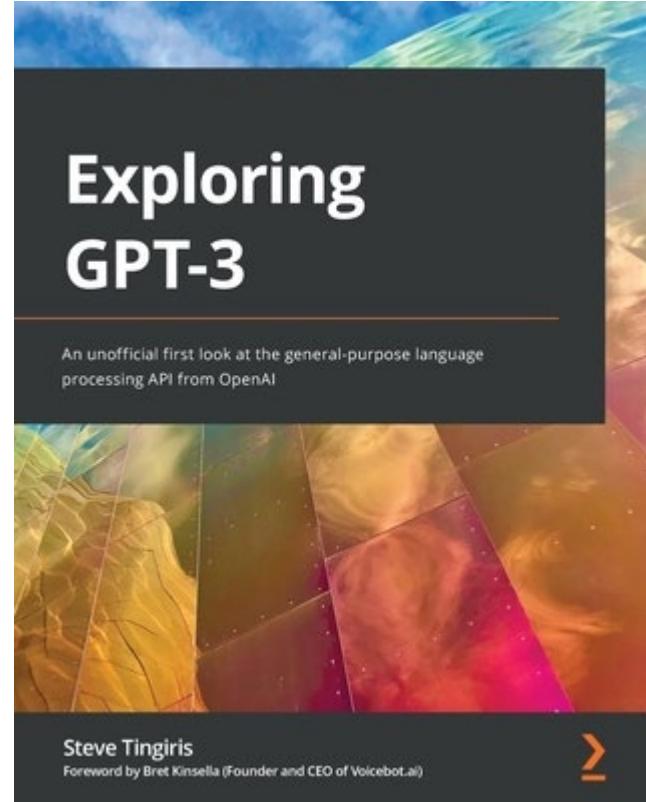
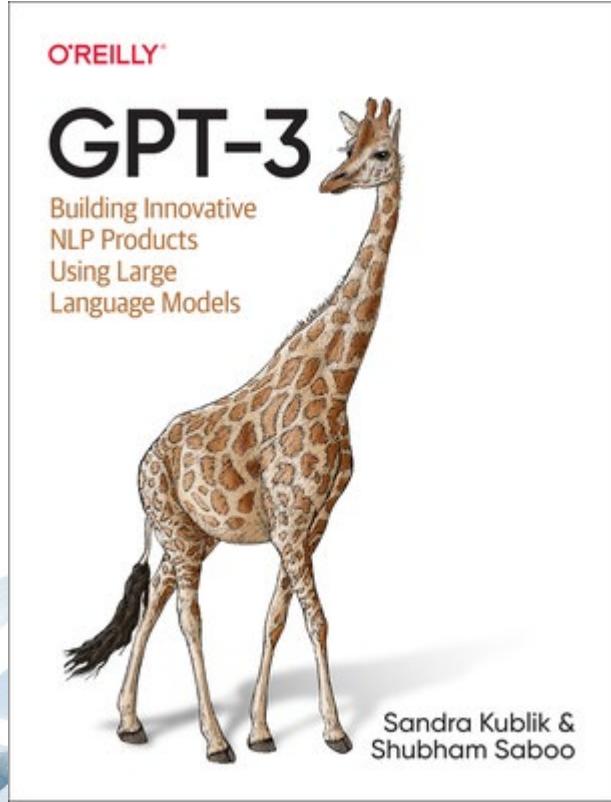


| 微軟 | Ignite | GPT-3 | OpenAI | RPA | 玉山AI挑戰賽

## AI趨勢周報第179期：微軟新推Azure OpenAI服務，企業使用GPT-3的門檻更低了！

微軟在近日Ignite大會上宣布推出Azure OpenAI服務，讓企業更容易以API使用GPT-3系列模型，有別於以往的邀請制。-01-03

<https://en.wikipedia.org/wiki/GPT-3>



Gpt-3: Building Innovative Nlp Products Using Large Language Models

<https://www.scientificamerican.com/article/we-asked-gpt-3-to-write-an-academic-paper-about-itself-then-we-tried-to-get-it-published/>

The screenshot shows the header of the Scientific American website. On the left, there are buttons for "Subscribe" and "Latest Issues". In the center, the "SCIENTIFIC AMERICAN" logo is displayed. On the right, there are links for "Cart 0", "Sign In", and "Newsletters". Below the header, a navigation bar includes links for "Coronavirus", "Health", "Mind & Brain", "Environment", "Technology" (which is underlined), "Space & Physics", "Video", "Podcasts", "Opinion", "Store", and a search icon. A dark banner across the middle features the text "Digital & 4-Year Archive" and a red "Subscribe" button.

ARTIFICIAL INTELLIGENCE | OPINION

# We Asked GPT-3 to Write an Academic Paper about Itself—Then We Tried to Get It Published

An artificially intelligent first author presents many ethical questions—and could upend the publishing process

By Almira Osmanovic Thunström on June 30, 2022

<https://techcrunch.com/2021/05/25/microsoft-uses-gpt-3-to-let-you-code-in-natural-language/>

## **Microsoft uses GPT-3 to let you code in natural language**

Frederic Lardinois @fredericl / 11:00 PM GMT+8 • May 25, 2021

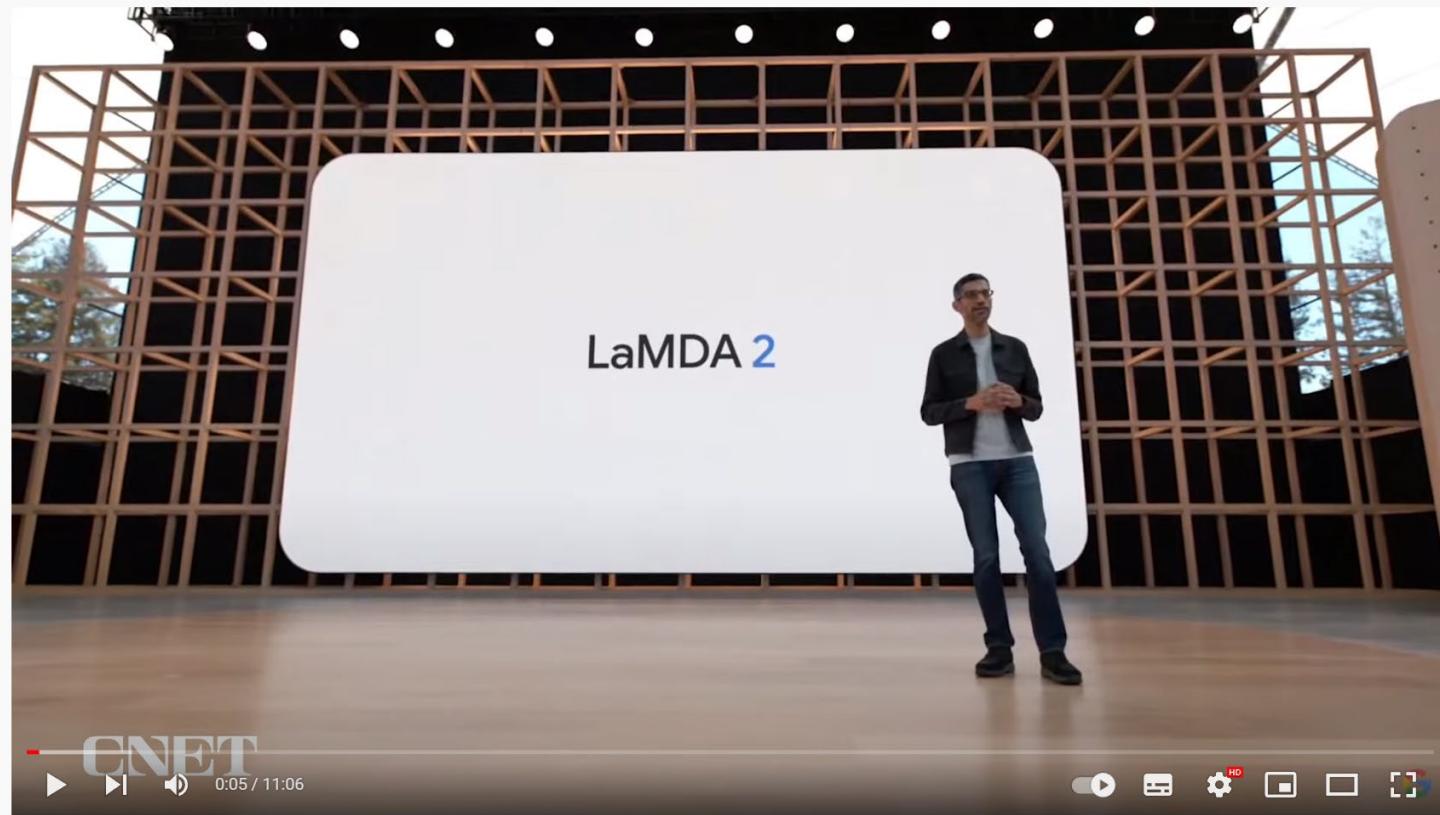
 Comment



 Image Credits: Getty Images

說AI有感知能力的Google工程師被解僱了  
<https://www.ithome.com.tw/news/152100>

<https://www.youtube.com/watch?v=u8t4-eqEYds>



Google LaMDA 2 vs. OpenAI GPT 3: (Watch the AI Demos)

觀看次數 : 32,955 次 2022年6月22日 Google's LaMDA and OpenAI's GPT 3 are both making big leaps in the development of conversational AI. See how the tech...更多內容

529 不喜歡 分享 下載 剪輯片段 儲存 ...

# NLP in security

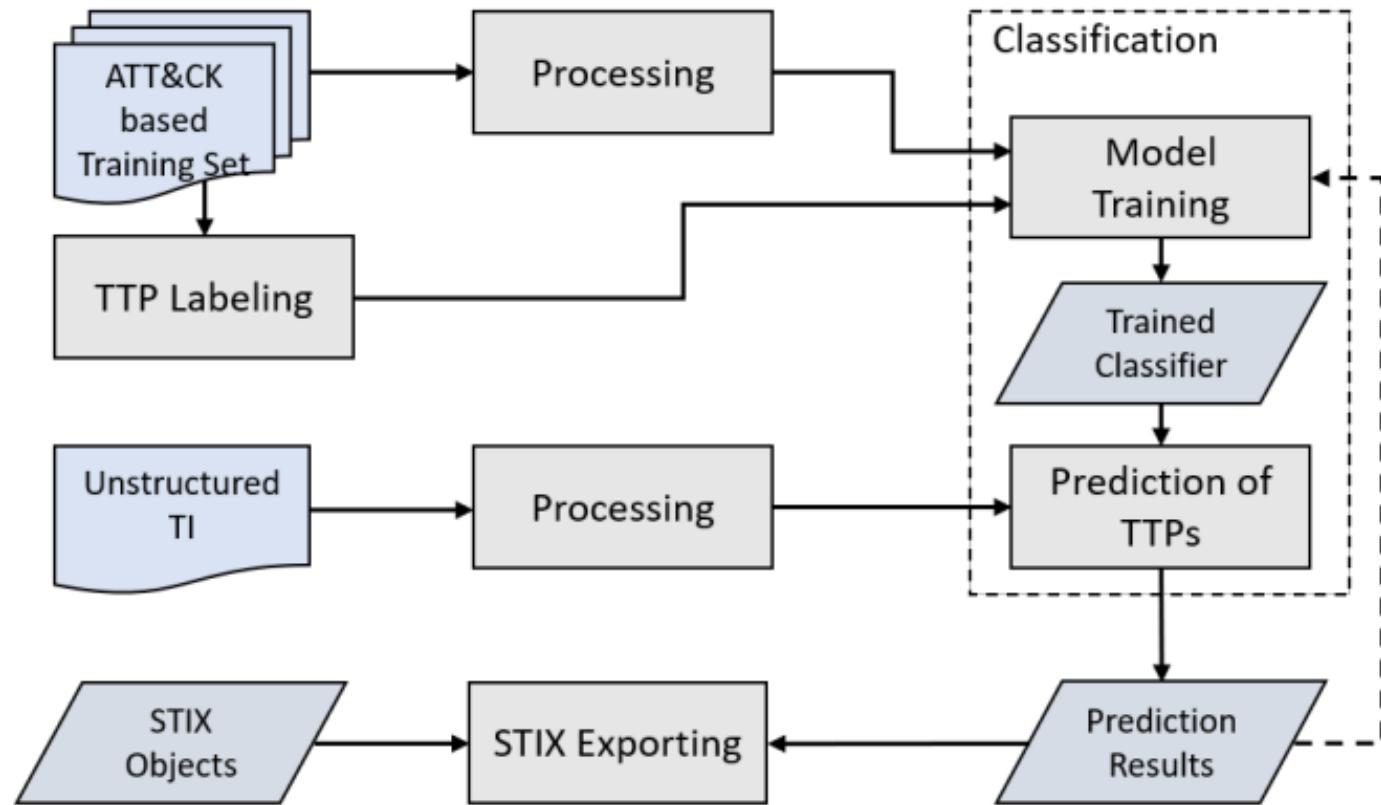
<https://arxiv.org/abs/2207.08478>



The image shows a screenshot of the arXiv search interface. At the top, there is a red header bar with the arXiv logo on the left and a search bar with the placeholder "Search..." and a "Help | Advanced" link on the right. Below the header, the URL "arXiv > cs > arXiv:2207.08478" is displayed. The main content area has a light gray background and features the title "Computer Science > Cryptography and Security" in bold black text. Below the title, the submission date "[Submitted on 18 Jul 2022]" is shown in a smaller gray font. The main title of the paper, "Towards Automated Classification of Attackers' TTPs by combining NLP with ML Techniques", is displayed in large, bold, black text. The authors' names, "Clemens Sauerwein, Alexander Pfohl", are listed below the title in blue text.

The increasingly sophisticated and growing number of threat actors along with the sheer speed at which cyber attacks unfold, make timely identification of attacks imperative to an organisations' security.

Consequently, persons responsible for security employ a large variety of information sources concerning emerging attacks, attackers' course of actions or indicators of compromise. However, a vast amount of the needed security information is available in unstructured textual form, which complicates the automated and timely extraction of attackers' Tactics, Techniques and Procedures (TTPs). In order to address this problem we systematically evaluate and compare different Natural Language Processing (NLP) and machine learning techniques used for security information extraction in research. Based on our investigations we propose a data processing pipeline that automatically classifies unstructured text according to attackers' tactics and techniques derived from a knowledge base of adversary tactics, techniques and procedures.



<https://arxiv.org/abs/2207.08478>

## NLP AND ML TECHNIQUES USED IN RELEVANT WORKS TO EXTRACT SECURITY-RELATED INFORMATION FORM UNSTRUCTURED TEXTS

	[44]	[49]	[50]	[30]	[51]	[17]	[18]	[20]	[21]	[26]	[31]	[23]	[22]	[32]	[52]	[53]	[34]	[37]	[38]	[24]	[27]	[29]	[54]	[41]	[42]
NLP	<b>Tokenization</b>	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	<b>POS Tagging</b>	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	<b>NER</b>				1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	<b>Stemming</b>	1		1				1				1		1										1	
	<b>Lemmatization</b>	1													1										
	<b>Bag-of-Words</b>						1				1		1				1	1	1						1
	<b>TF-IDF</b>	1		1				1	1	1	1		1	1	1	1	1	1	1	1					1
	<b>Word Embedding</b>		1			1		1		1	1	1	1	1	1	1	1	1	1	1	1				
ML	<b>SVM</b>	1				1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	<b>NB</b>						1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	<b>Other</b>	1	1	1	1		1		1	1	1	1				1	1			1	1	1	1	1	1

<https://arxiv.org/abs/2207.08478>

```
081111 083419 24621 INFO dfs.DataNode$DataXceiver: Receiving block blk_5214640714119373081 src:  
/10.251.121.224:47915 dest: /10.251.121.224:50010  
081111 083419 35 INFO dfs.FSNamesystem: BLOCK* NameSystem.allocateBlock:  
/user/root/rand7/_temporary/_task_200811101024_0014_m_001575_0/part-01575. blk_5214640714119373081  
081111 083420 24633 INFO dfs.DataNode$DataXceiver: Receiving block blk_5214640714119373081 src:  
/10.251.121.224:57800 dest: /10.251.121.224:50010  
081111 083422 24621 INFO dfs.DataNode$DataXceiver: writeBlock blk_5214640714119373081 received  
exception java.io.IOException: Could not read from stream  
081111 104136 26436 INFO dfs.DataNode$DataXceiver: Receiving block blk_-3208483482800741142 src:  
/10.251.111.209:34510 dest: /10.251.111.209:50010  
081111 104136 26954 INFO dfs.DataNode$DataXceiver: Receiving block blk_-3208483482800741142 src:  
/10.251.203.80:40000 dest: /10.251.111.209:50010  
081111 104136 27196 INFO dfs.DataNode$DataXceiver: Receiving block blk_-3208483482800741142 src:  
/10.251.111.209:40712 dest: /10.251.111.209:50010
```

Automatically detected anomaly

<https://www.youtube.com/watch?v=At19CBGpbMI>

# DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning

Authors:  [Min Du](#),  [Feifei Li](#),  [Guineng Zheng](#),  [Vivek Srikumar](#) [Authors Info & Claims](#)

CCS '17: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security • October 2017 • Pages 1285–1298 • <https://doi.org/10.1145/3133956.3134015>

Online: 30 October 2017 [Publication History](#)

 442  9,274



<https://www.youtube.com/watch?v=At19CBGpbMI>

DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning - Min Du - November 1st, 2017

DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning

Min Du, Feifei Li, Guineng Zheng, Vivek Srikumar  
University of Utah

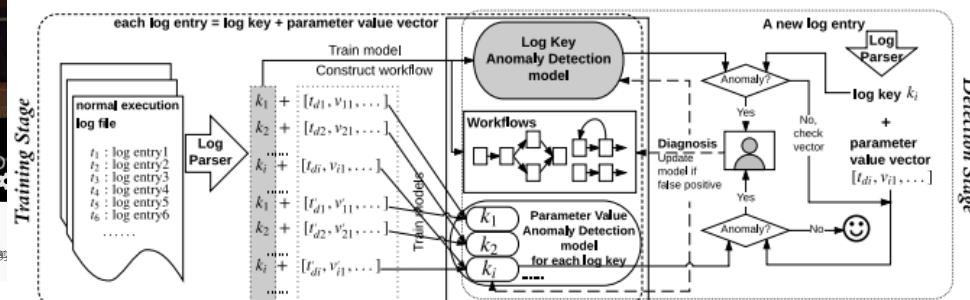
ACM CCS 2017. ©ACM, Inc. All Rights Reserved. www.acm.org

0:09 / 28:40 • 章節 >

ACM CCS 2017 - DeepLog: Anomaly Detection and Diagnosis from System Logs [...] - Min Du

觀看次數：8,888次 2018年2月7日 Presented by Min Du. November 1st, 2017.

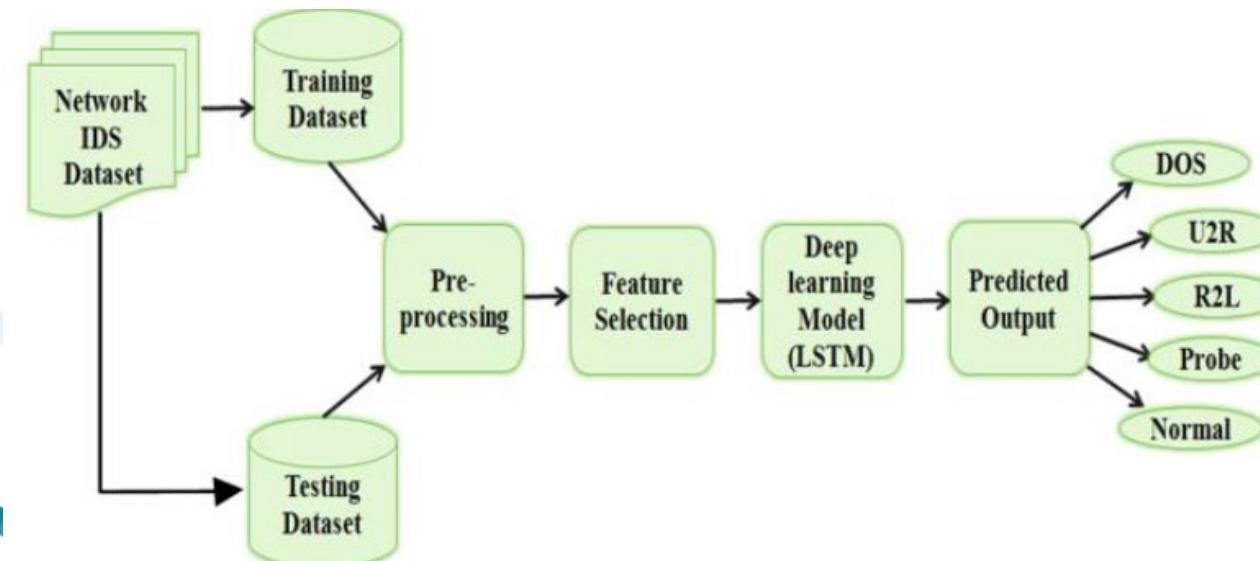
© 2017 ACM, Inc. All Rights Reserved. www.acm.org ...更多內容



# Long Short-Term Memory (LSTM) Deep Learning Method for Intrusion Detection in Network Security

Supriya Shende <sup>1st</sup>  
Department of Electronics Engineering  
Government College of Engineering  
Amravati, India

Samrat Thorat <sup>2nd</sup>  
Department of Electronics Engineering  
Government College of Engineering  
Amravati, India



## **DEEPCASE: Semi-Supervised Contextual Analysis of Security Events**

Thijs van Ede<sup>\*†</sup>, Hojjat Aghakhani<sup>†</sup>, Noah Spahn<sup>†</sup>, Riccardo Bortolameotti<sup>‡</sup>, Marco Cova<sup>§</sup>, Andrea Continella<sup>\*</sup>,  
Maarten van Steen<sup>\*</sup>, Andreas Peter<sup>\*</sup>, Christopher Kruegel<sup>†</sup>, Giovanni Vigna<sup>†</sup>

<sup>\*</sup>University of Twente, <sup>†</sup>University of California, Santa Barbara, <sup>‡</sup>ReaQta, <sup>§</sup>VMware, Inc.

{t.s.vanede, a.continella, m.r.vansteen, a.peter}@utwente.nl, {hojjat, ncs, chris, vigna}@cs.ucsb.edu,  
r.bortolameotti@reaqta.com, covam@vmware.com

<https://vm-thijs.ewi.utwente.nl/static/homepage/papers/deepcase.pdf>

## Computer Science &gt; Cryptography and Security

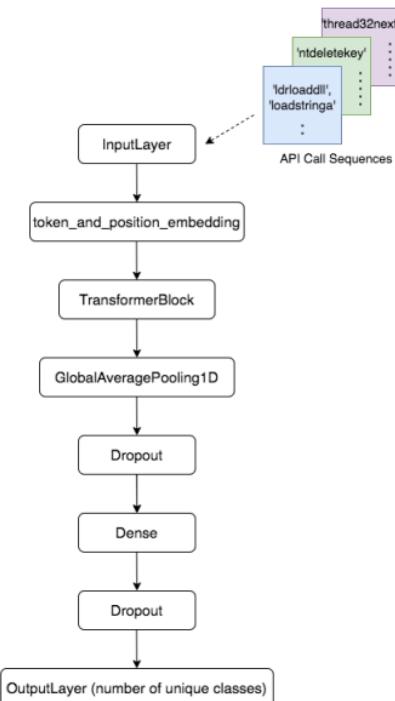
[Submitted on 25 Dec 2021 (v1), last revised 22 Jun 2022 (this version, v4)]

# An Ensemble of Pre-trained Transformer Models For Imbalanced Multiclass Malware Classification

Ferhat Demirkiran, Aykut Çayır, Uğur Ünal, Hasan Dağ

<https://arxiv.org/abs/2112.13236>

Classification of malware families is crucial for a comprehensive understanding of how they can infect devices, computers, or systems. Thus, malware identification enables security researchers and incident responders to take precautions against malware and accelerate mitigation. API call sequences made by malware are widely utilized features by machine and deep learning models for malware classification as these sequences represent the behavior of malware. However, traditional machine and deep learning models remain incapable of capturing sequence relationships between API calls. On the other hand, the transformer-based models process sequences as a whole and learn relationships between API calls due to multi-head attention mechanisms and positional embeddings. Our experiments demonstrate that the transformer model with one transformer block layer surpassed the widely used base architecture, LSTM. Moreover, BERT or CANINE, pre-trained transformer models, outperformed in classifying highly imbalanced malware families according to evaluation metrics, F1-score, and AUC score. Furthermore, the proposed bagging-based random transformer forest (RTF), an ensemble of BERT or CANINE, has reached the state-of-the-art evaluation scores on three out of four datasets, particularly state-of-the-art F1-score of 0.6149 on one of the commonly used benchmark dataset.



## Computer Science &gt; Machine Learning

[Submitted on 17 Mar 2022]

# The Analysis of Online Event Streams: Predicting the Next Activity for Anomaly Detection

Suhwan Lee, Xixi Lu, Hajo A. Reijers

<https://arxiv.org/abs/2203.09619>

Anomaly detection in process mining focuses on identifying anomalous cases or events in process executions. The resulting diagnostics are used to provide measures to prevent fraudulent behavior, as well as to derive recommendations for improving process compliance and security. Most existing techniques focus on detecting anomalous cases in an offline setting. However, to identify potential anomalies in a timely manner and take immediate countermeasures, it is necessary to detect event-level anomalies online, in real-time. In this paper, we propose to tackle the online event anomaly detection problem using next-activity prediction methods. More specifically, we investigate the use of both ML models (such as RF and XGBoost) and deep models (such as LSTM) to predict the probabilities of next-activities and consider the events predicted unlikely as anomalies. We compare these predictive anomaly detection methods to four classical unsupervised anomaly detection approaches (such as Isolation forest and LOF) in the online setting. Our evaluation shows that the proposed method using ML models tends to outperform the one using a deep model, while both methods outperform the classical unsupervised approaches in detecting anomalous events.