TibaMe A.I. 應用工程師青年班第一期

10 things to know about

自然語言專題: Blog樂

組員名單:

袁學平、劉康彥、孫韶甫

林靖凱、陳豊壬、丁之琳

組長:吳秉軒







01 專題介紹

02 實作:分類

03 實作:推薦

04 實作:分析與摘要

05 應用部署

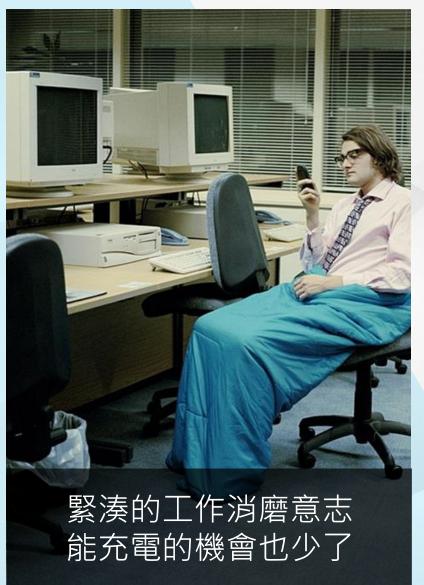
06 總結

# 專題介紹

報告人:**吳秉軒** 

# 日劇裡的小故事



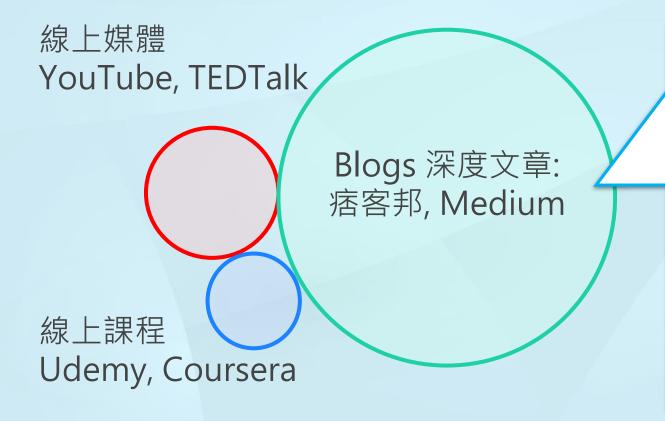




# 學習,職場只是開始

- □全球行動應用市場將從 2019 年, 21 億美金到 2026 年, 53 億美金。
- □終生學習、學習場域為重點。「微學習」趨勢:用少於10分鐘獲取知識:





快速、豐富 7,000 萬 +

每個月部落格文章新增數量

1億/250萬

全球/臺灣獨立造訪

# 使用部落格作為學習資源的挑戰

□讀者面:龐大的資源量和「微學習」 無法接軌。

調查 40 位歷屆 TibaMe 成員,最想要的功能:

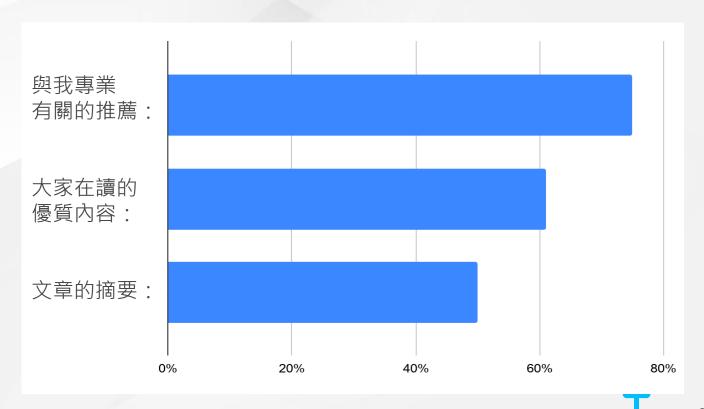
□ 作者面:全球最大的部落格 平台,不支援中文推薦。

#### 「我的文章沒人看」



再、再堅持一下Q\_Q

很現實的,如《讓自己的作品被看見,是作者的責任》這篇所說,寫文章就是想要有人看,不論是想靠這賺錢、賺知名度、或是單純宣揚理念都好,否則大可不用費心使用公開平台,只要用 Bear 或是 Evernote 等工具自己整理即可。



來源:medium.com

題秉

介軒

紹

# 當手機和 Chatbot 遇上部落格

- •••• FLOW ♀ 8:00 AM @ \* I CNN > Home (4) Typically replies instantly Manage Show me news stories from last December about Elon musk Here's something about december about elon musk A first look at Elon Musk's Elon N next grand idea save CNN Money: Elon Musk CNN N revealed Friday a video 'em, j depicting a network of makin underground ... Read this Story
- □ Medium 臺灣造訪量 42.24% 來自手機。
- □聊天機器人商業應用已成熟,臺灣LINE生態系完整。



來源:pro.similarweb.com, pixnet.net



專吳 題 新 新 紹

# Blog樂 功能設計

□以 LINE Bot 做部落格閱讀微學習:

機器學習:文章分類、文章推薦。

深度學習:文章摘要、文章分析。

□建立產品黏著度正循環





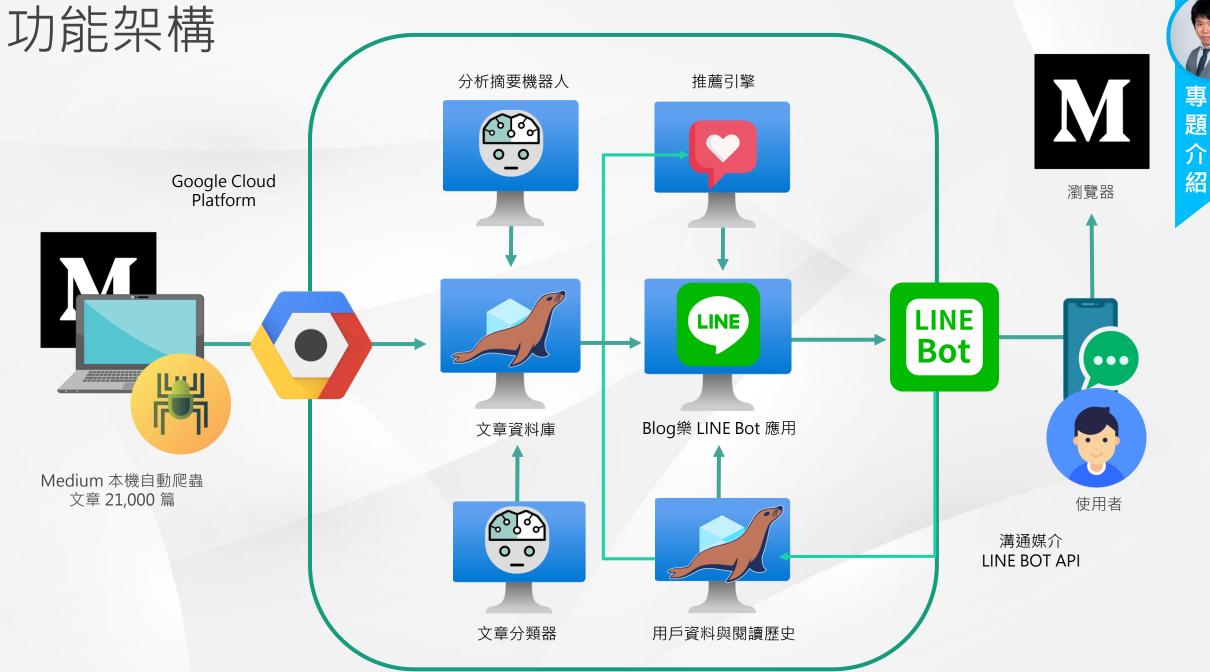


題 秉介 軒紹

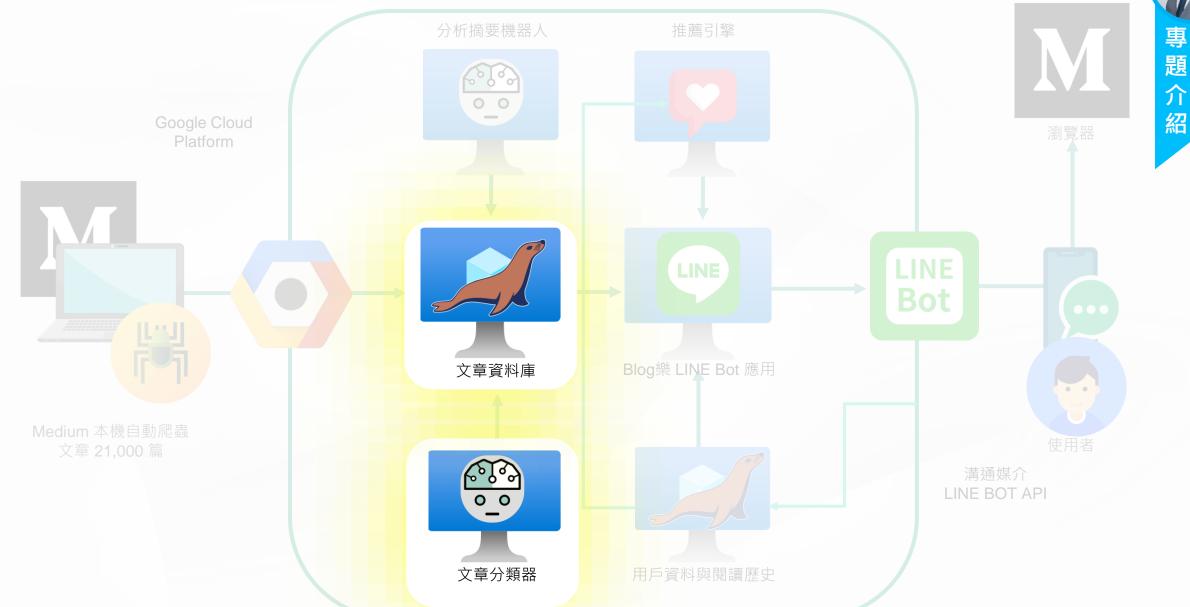
# Blog樂 功能展示







# 功能架構:分類



# 文章分類

報告人:**劉康彥** 

# 為什麼要做文章分類?

文章分

- → 不想看到廣告推銷文
- → 快速找到想看的文章





#### UX 設計師在敏捷團隊的因應之道

這篇文章適合:正跟著公司跑 Scrum 敏捷開發法,卻一直在流程裡載浮載沉的 UX 設計師

Read more...

(1.8K



#### 新冠疫情造成便利商店發展受阻 辦公室內的 無人便利店正夯

日本國內三大便利超商 7-11,全家便利以及羅森,從去年開始因為加盟店 主抗議超商強制 24 小時醫業不合理,經過了一年的反覆陳情協調,今年 9月日本公平交易委員會正式表示,三大超商強制加盟店 24...

Read more...

**(3)** 247



#### Email 的藝術

如何重建一隊警察

千奇百趣的日常

一直覺得公司電腦的 wallpaper 很沉悶,於是特意在 Unsplash找了幾張來用,換了好幾天才發現,好像很少見到新的 wallpaper,原來我的 wallpaper,就是長開的 Outlook...

希望你會同意,香港警察需要重建。過去兩星期的局勢發展,可見香港警察

的問題已經不是個別的警員或者指揮官的問題,而是整個警隊從一哥起計到

Read more...

230

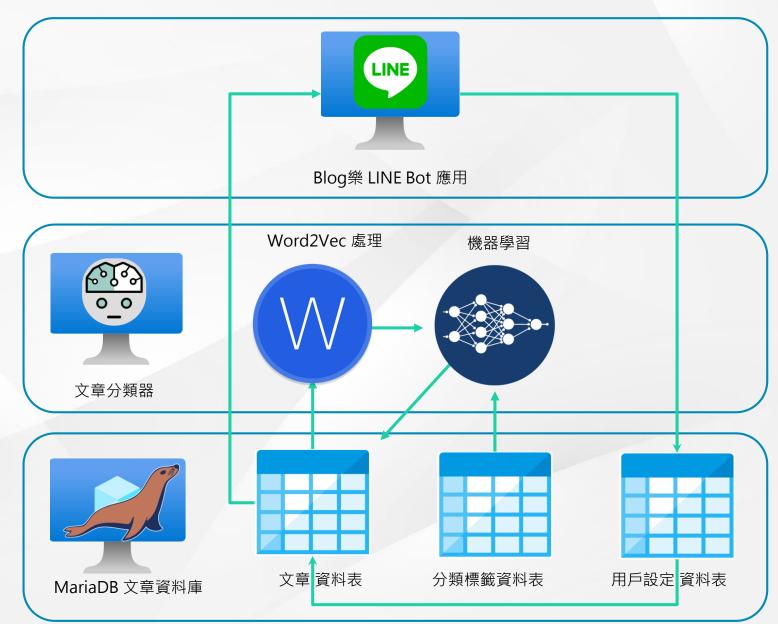


D v

13

# 功能構想

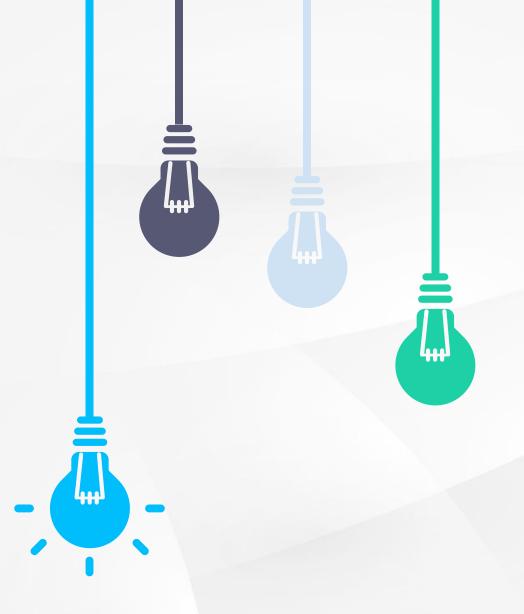






分彥

類



# 資料處理流程



瀏覽資料庫補足缺漏值



訂下各大分類標記訓練資料



文本預處理跟 TF-IDF 特徵提取



嘗試使用各種模型並做挑選



## 挑戰與解決方案



STORIES



未成省 lun 24, 2019 · 5 min read



#### 如何重建一隊警察

希望你會同意,香港警察需要重建。過去兩星期的局勢發展,可見香港警察的問題已經不是個別的警員或者指揮官的問題,而是整個警隊從一哥起計到 PC都出現了嚴重的文化問題。

Read more...



W

# 模型預測不佳 能不能有更多選擇?

========

第 24 篇 如何重建一隊警察 [3]

========



========

第 24 篇 如何重建一隊警察 [10]

=======



# 文本分類方式解說與結果

分彥

## 二分類

不符合分類要求

## 多分類

部分文章分類困難

## 多標籤

更為精準的分類



ID	標題	分 類	工程 開發	資料 科學	使用 設計	產品 管理	行銷 管理	商業 投資	職場 心得
580	UX設計師在敏捷團隊的因應之道	3	0	0	1	1	0	0	1



# 功能架構:推薦

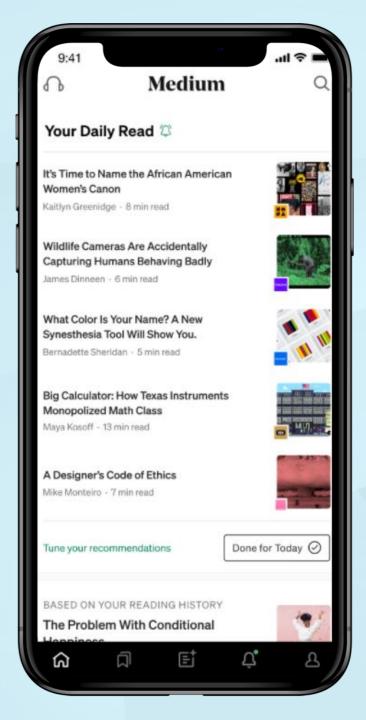




類

# 文章推薦

報告人:**林靖凱,丁之琳** 



# 為什麼要推薦系統?

Medium 官方 APP 的呈現



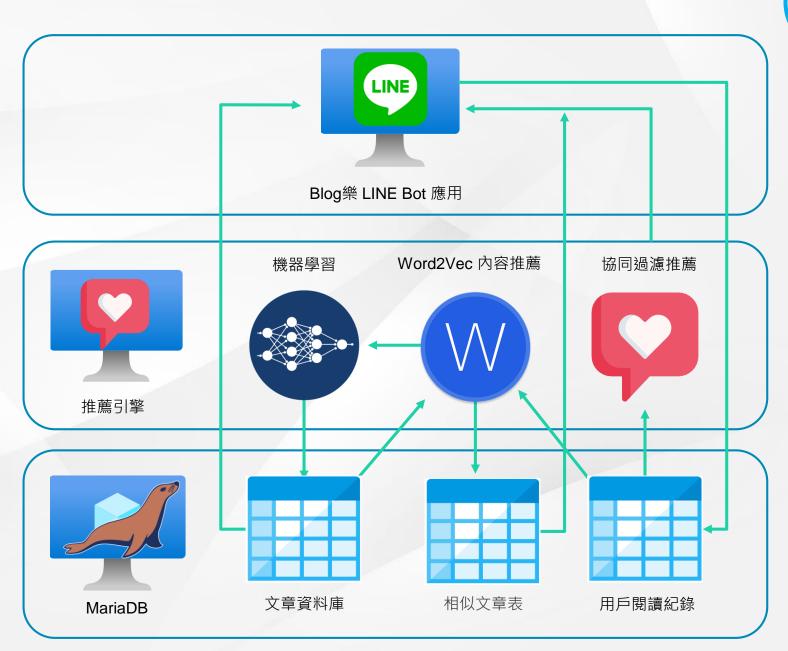


沒有依照讀者過往的閱讀喜好推薦文章



# 功能構想

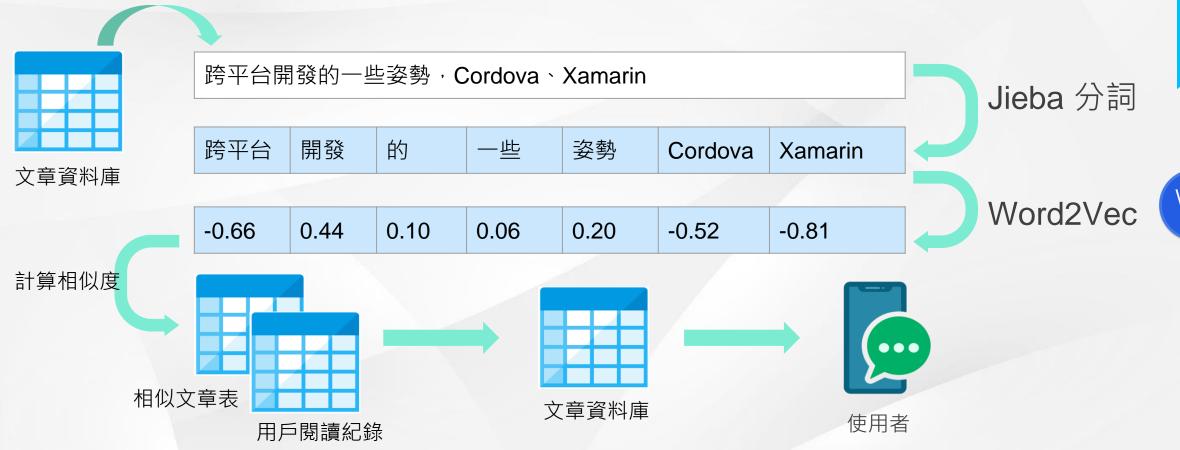






# 推薦文章:過去閱讀記錄推薦 (Content Based)





# 推薦文章:協同過濾 (Collaborative Filtering)



文林 章 推 薦

推薦系統套件:Surprise



用戶閱讀紀錄

	A article	B article	C article	D article	
A user	5	4	?	?	
B user	?	4	?	3	
C user	?	?	1	2	

	A article	B article	C articl	D article
A user	5	4	?	?
B user	5	4	2	3
C user	?	?	1	2

算法:SVD



# 挑戰與解決方案





內文相近度推薦, 計算時間較長





增加資料表,將計算的 結果提前存入



**冷啟動**狀況,使用者 喜好及人數還不夠多



熱門分類文章

閱讀記錄推薦

協同過濾推薦

依照使用者已閱讀 篇數逐步採用開啟 更多推薦方式



# 結果呈現



推凱

薦



### w2v 閱讀記錄推薦

輸入文章是

資料庫內第 130.0 篇

文章ID: 693

原文標題: [私心推薦] PM / UX / UI 學習資源總整理

相似文章排序:

資料庫內第 3596.0 篇

相似度: 0.962354302406311

文章ID: 4708

推薦標題: 如何推動設計文化與創新思維

資料庫內第 7706.0 篇

相似度: 0.9610608220100403

文章ID: 9708

推薦標題: 騰訊 UIUX 暑期實習分享全記錄



### surprise 協同過濾推薦

IntputFindCmd:

SELECT `title` FROM `main` WHERE `id` = 8045

文章ID: 8045

可能的喜好程度: 4.037292825543596

推薦標題: 適合練習照片牆 App 的 flickr API

IntputFindCmd:

SELECT `title` FROM `main` WHERE `id` = 4884

文章ID: 4884

可能的喜好程度: 3.798032858903403

推薦標題: 不只是抓資料,API 可以為我們提供各種服務

IntputFindCmd:

SELECT `title` FROM `main` WHERE `id` = 1625

文章ID: 1625

可能的喜好程度: 3.716170735378744 推薦標題: 線上支付體系的大戰略思考

# 推薦功能延伸:優質新文預測

優文預測

預測近期具有潛力的優質新文章,並在推薦功能內回推給使用者



# 以機器學習判斷優質新文推薦

- □以固定算式標記訓練集文章熱門程度。
- □ 不同熱門等級的內文做文本特徵處理。

們不乏面試過許多資深 PM , 但卻常常碰到下述狀況:

wallpaper·就是長開的 Outlook..

□使用機器學習模型預測近期文章熱門程度。





預 琳

# 文本預處理與模型選擇

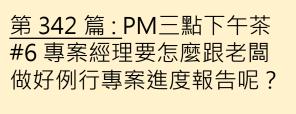
優丁文強琳

- □以TF-IDF 進行文本特徵提取
- □挑選合適模型



# 模型評估結果

□ 預測近六個月內的新文在不同文章評比等級的機率,並 將高品質文章回推給使用者



等級:1

高: 0.6724

中: 0.2345

低: 0.0931





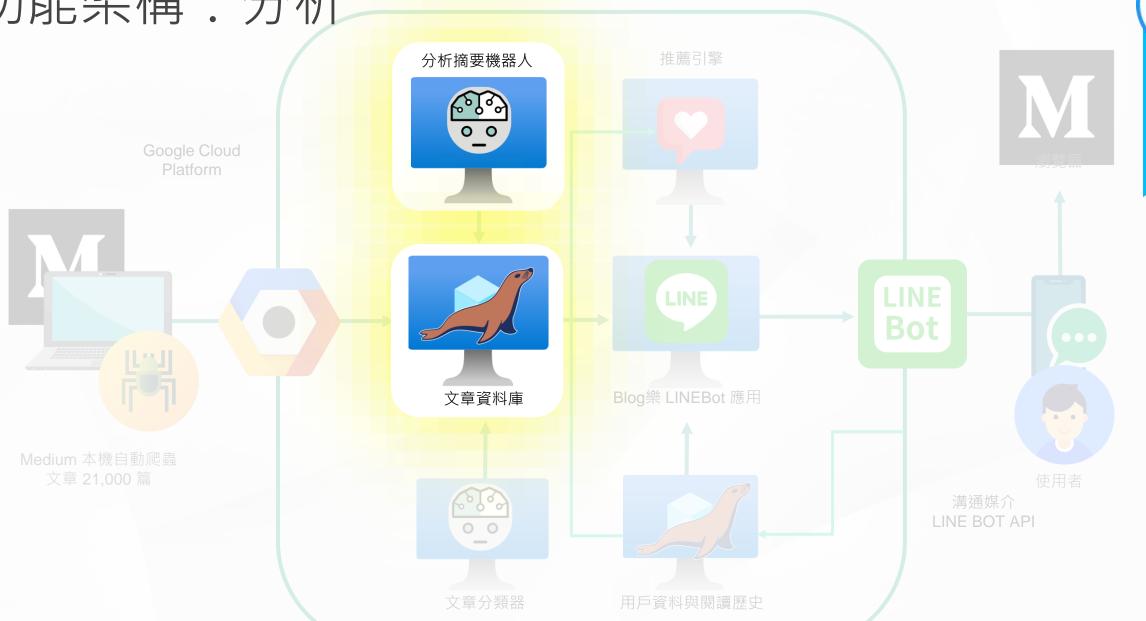


29

預 琳

測

# 功能架構:分析





優文 預 測

# 文章分析:趨勢分析



報告人:**袁學平** 

# 趨勢分析的需求







新的靈感



文章作者

# 功能構想:提供關鍵詞

□ 紀錄使用者搜尋關鍵詞:基本做法

□ 由 A.I. 分析文章內容抽取關鍵詞:隨時可用



趨勢分 析



# 功能架構

排序/選取關鍵詞



分析摘要機器人

趨勢關鍵詞提供



趨勢 分析

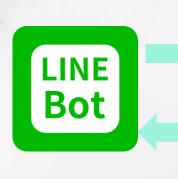
NLP處理 分詞工具/排序模型

新關鍵詞/文章

匯入文章



文章資料庫





使用者

匯入



用戶資料與閱讀歷史

搜尋和閱讀結果

# 分詞工具: Jieba / CKIP

分平

探尋顧客內心最真實的意圖-關 鍵字搜尋分析 (Keyword Search Analytics )





### 以市場區隔與機器學習篩選有價值的顧客

己想要的資訊。企業則希望它的服務或產品,能夠出現在搜尋結果的第一 頁,讓消費者有更多的機會予以點擊,而「關鍵字」就是串連「搜尋引擎 天秤」兩端的媒介。

「關鍵字」指的是,網路使用者在搜尋引擎裡輸入所欲搜索查詢的字後, 經網站比對到的字,即稱為「關鍵字」(keyword);而「關鍵字搜尋分 析」則是網路行銷中很重要的一種分析工具,其有助於企業不斷優化網站 設計,以及達成網路行銷的目的。

#### 閱讀更多: 翻轉零售— 大數據帶來的零售業革命

從技術上來看,網路使用者在搜尋引擎裡所輸入的字稱為「搜索查詢」 (Search Query),而經網頁比對到的字被稱為「關鍵字」 (keyword) ·

□ 將中文文本分成詞以抽取關鍵詞



以 市場 區隔 與 機器 學習 篩 價值的 選有 顧客



除 Jieba 外再導入 CKIP 中英文分詞結果



# 排序模型:BERT

- □ 目標 NLP 模型 BERT 實作
- □ 先標籤關鍵詞以供訓練

市場 深度 機器 學習

區隔 Machine Learning



機器學習 市場區隔 簽約金 Machine Learning 深度學習



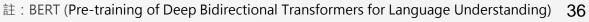
編碼後文章 作為標籤

0



所有訓練文章

以市場區隔與機器學習篩選...



## 分析結果: 關鍵詞重要度序列



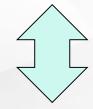
趨勢分 析 有



重要度編碼

... 0.8 0.8 0.8 0.8 0.1 0.5 0.5 0.5 .

對應文章篩選關鍵詞



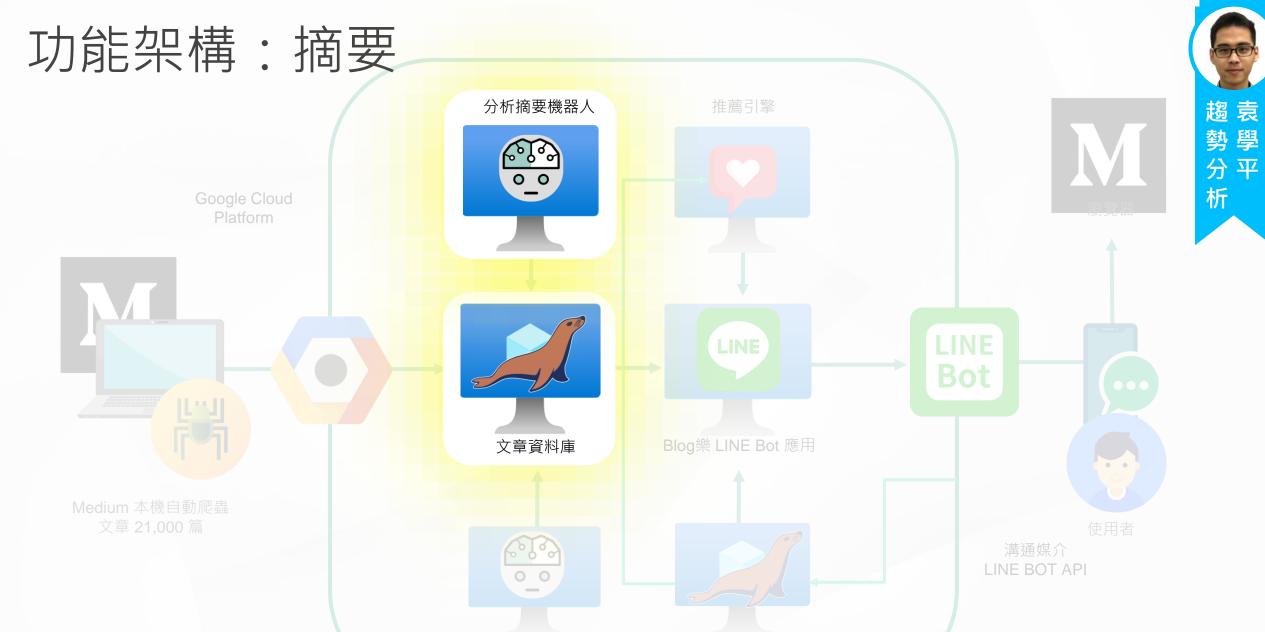
...機器學習或深度學習...

工程和資料分析類排序結果:

- 1. NativeScript
- 2. Flutter
- 3. Android
- 4. Electron
- 5. MacOS



統計所有近期文章後...



## 文章分析:摘要生成

報告人:**陳豊壬** 

## 摘要,有需要嗎?

敏捷式開發(Agile)、瀑布式開發 (Waterfall)、敏捷式UX、Lean UX。兜幾?



V 2 2 2 2 --

大家可能看掉標頭邮會覺得:到底是在工穀級。

今天小闆WLI要带大家破解这些名词、分享其中意则。

若人次有着過之前文章「<u>企品原型製作工具</u>(Prototyping Tools)使用分字』 小場有務物提刊Agile Methodology和Waterfall Methodology的差異・今天 要更深入採計一點・

. . .

#### 般排式開發 ( Agile Methodology )

Agile是现在故籍走品開發的維勢(紫然根據走品性質會有許不同) · 顧名 思致 · 似是秘訣 · 秘訣 · 秘訣 ·

根據定義 - Apile的原則是減少沒費 - 產出迅速 - 不斷循環 - 快速學習 -

實際使用上海例、Scrum是Agile Methodology·但蒸著名的香浆、一般來說 一個Scrum Team與成角產品經環、Scrum Mester (主導Scrum的人、可能 是project manager)、工程師和開發者。有可能以產品功能來信或主要任 務、譬如說交及散體的「治國散體更過衰可以經典看到所有蘇門的剩余並 位數程,(我們在兩面我說他想要的我不能)、以治為例一般來說以兩獨為單 位(報圖是他所為Speint)來的新列,每天早上開除可能有每日會過 (Stand up meeting)來更新說太的狀況。看到完設所有固定品有關的人 使報告,不斷述代於的循環、學習和調整。

並不身級稅技式階發就是好的,有時候會因為衝到到後面跟一開始的目標 價儲,讓產品缺乏一致性,硬體產品也較不適用於此方法,畢竟硬體做出 來紙不能一百更改。

#### 港布式開發(Waterfall)

瀑布式似比較不达代式(Iteratvic) · 像瀑布般從上往下:

產品需求→設計→開發→驗證→保修更新

開發時一般不會四五上一步驅,若有需要更改成是修正時較為困難,所以 每一步驅動會非常小心仔細的紊出。

<u>这篇文章</u>非常清楚指出何時較適合使用Agilc或Waterfall

### 精準選讀!!







## 功能構想



文 章 題 語 要

#### 內文

地球表層均由堅 硬的岩石構成,是人類最常 使用的自然資源之一。用肉 眼觀察,可以發現各種岩石 的外觀、顏色與質地各自不 同,依成因可以分為火成岩 、沉積岩及變質岩三大類。 由岩漿冷卻凝固形成的岩石 是火成岩,若岩漿噴發至地 表快速冷卻,其物 較小,甚至肉眼無 如玄武岩;若岩漿 處緩慢冷卻,則形成的 顆粒較大,如花岡岩...

文章資料庫



#### 摘要

#### (條列式:)

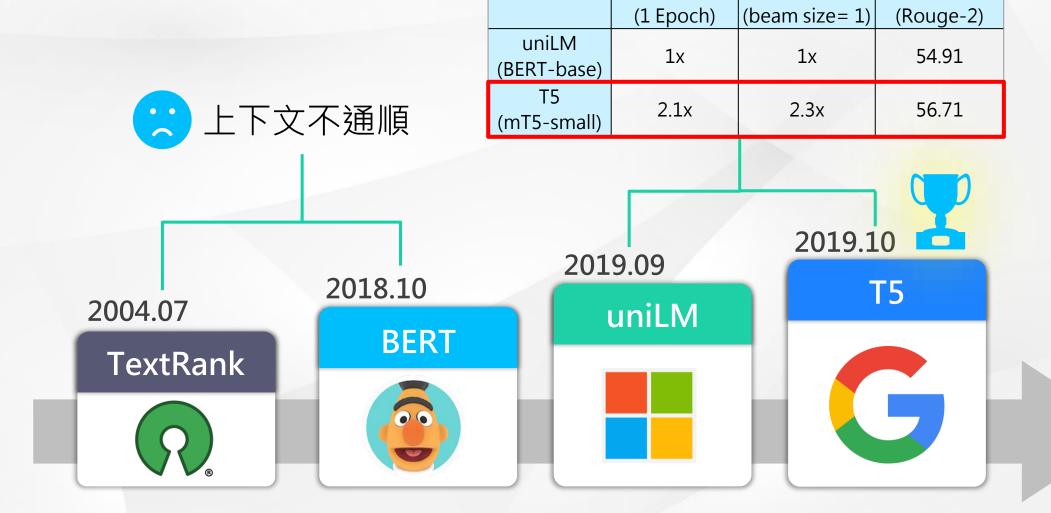
- 1.地球表層由岩石所構成。
- 2.依成因可以分為火成岩、沉 積岩及變質岩三大類。
- 3.岩漿冷卻凝固形成的岩石是 火成岩,如玄武岩、花岡 岩。

#### (敘述式:)

地球表層係由岩石 所構成,依其成因可以分為火 成岩、沉積岩及變質岩三大類 。火成岩係由岩漿冷卻凝固形 成,如玄武岩、花岡岩。

## 演算法/模型選用





訓練速度

生成速度

摘要效果

註: Rouge( Recall-Oriented Understudy for Gisting Evaluation )
uniLM( Unified Language Model Pre-training for Natural Language Understanding and Generation )
T5( Text-To-Text Transfer Transformer )

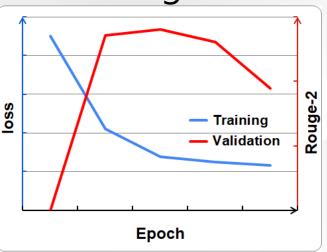
## 遭遇挑戰與解決方式

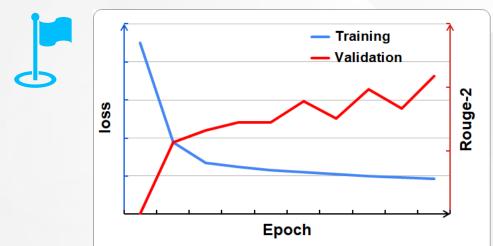






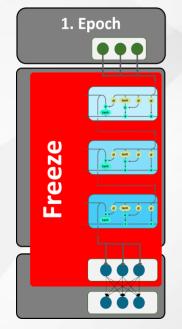
### Overfitting

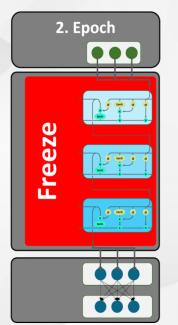


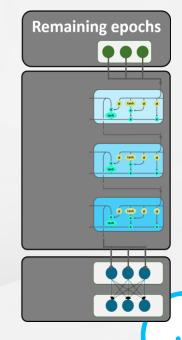


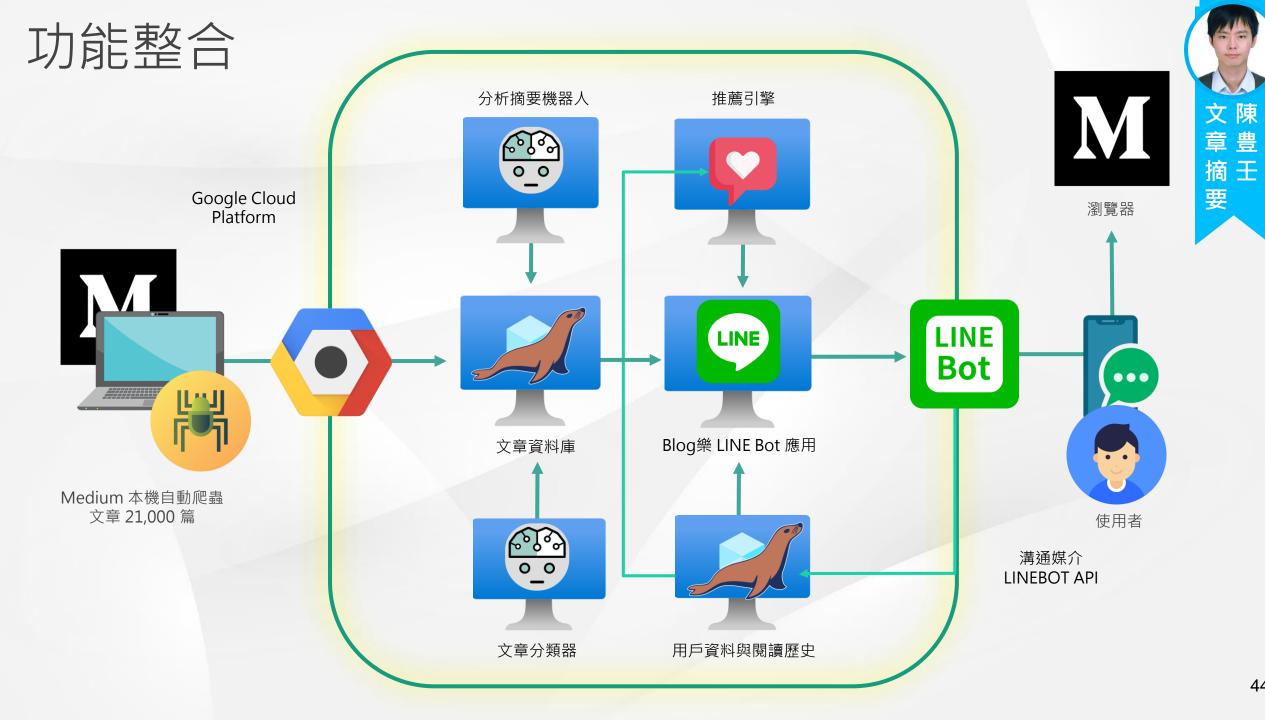


### Gradual Unfreezing 逐步解凍







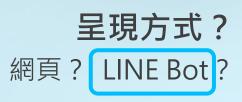


# 應用部署

報告人:**孫韶甫** 

## 上述功能怎麼統合完成?





### 資料庫型式?

MySQL? MariaDB? MongoDB?

團隊程式開發?

如何整合?Git? Code review?



Azure ? GCP ? AWS ?



Google Cloud Platform

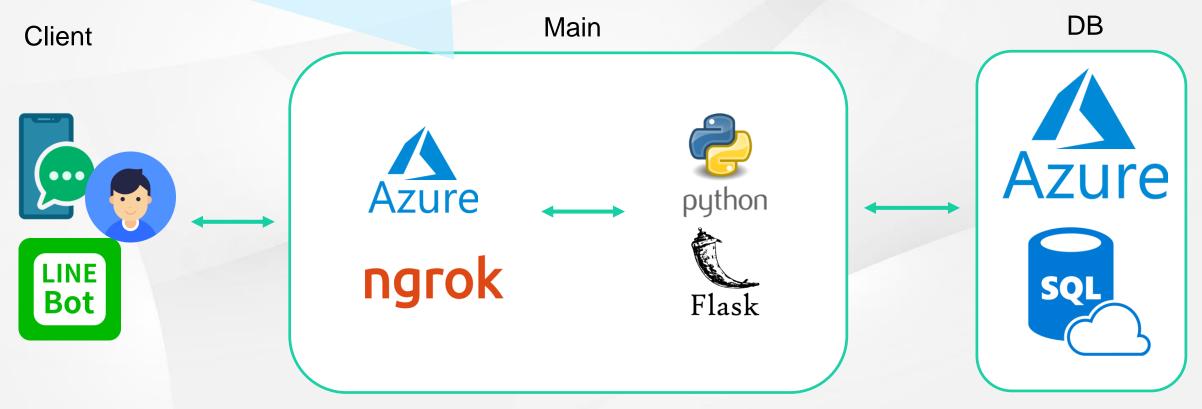
#### 網路架構?

NGINX? Apache? uWSGI? ngrok? Flask? Django?...

## 測試階段架構圖

- □ Azure App Service:長時間沒使用需喚醒。雲端機器收費較昂貴。
- □ngrok:每分鐘連線次數有限制。8小時需重新開啟服務(免費版)。
- □ Azure SQL Server收費:雲端機器費用+資料庫容量費用。





## 部署階段架構圖

□ Certbot:免費產生 SSL 憑證。

□ NGINX: 高效的 Web Server。

□ uWSGI:接收 NGINX 動態請求並處理後,發給 Web Application。

☐ WSGI: Python Web Server GateWay Interface ∘

Main Client Google Compute **Engine** Google Cloud Platform \*\*Certbot LINE Web Server **Bot** Flask



DB

## 團隊軟體開發

# 庭 2%

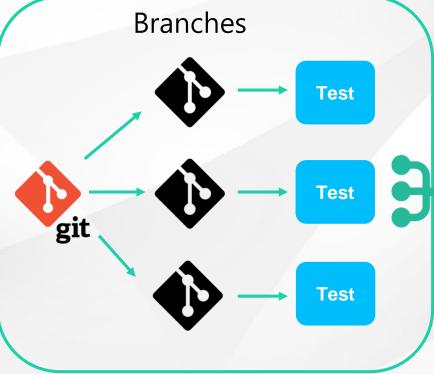
應用架構 構

#### 前期準備:

- □建立功能架構圖。
- □建立基礎軟體架構。
- □訂定團隊合作規則。
- □運行環境測試。



### 團隊開發 (GitHub flow)



#### C.R. (Code Review):

- →檢閱開發程式碼。
- 」程式碼討論優化。
- □ 知識共享,相互提升。

C.R. & Merge

All Test



#### All Test:

- □長時間運行測試。
- □單一功能測試。
- **〕**同時間多人使用測試。

# 總結

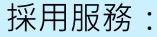
## 技術總覽



專 題 類 無 結

### 使用技術:

- □ 分詞: Jieba、CKIP
- □ 詞轉向量: Gensim、TF-IDF
- → 機器學習:Naive Bayes \Logistic Regression
- □ 分類系統: Label Power Set
- □ 推薦系統: Scikitsurprise
- □ 自然語言: BERT、UniLM、T5

















#### 適用領域:

- □ 大數據資料處理
- □ 行銷資料分析
- □ 電商推薦系統
- □ 數位教學系統
- □ 數位客服系統





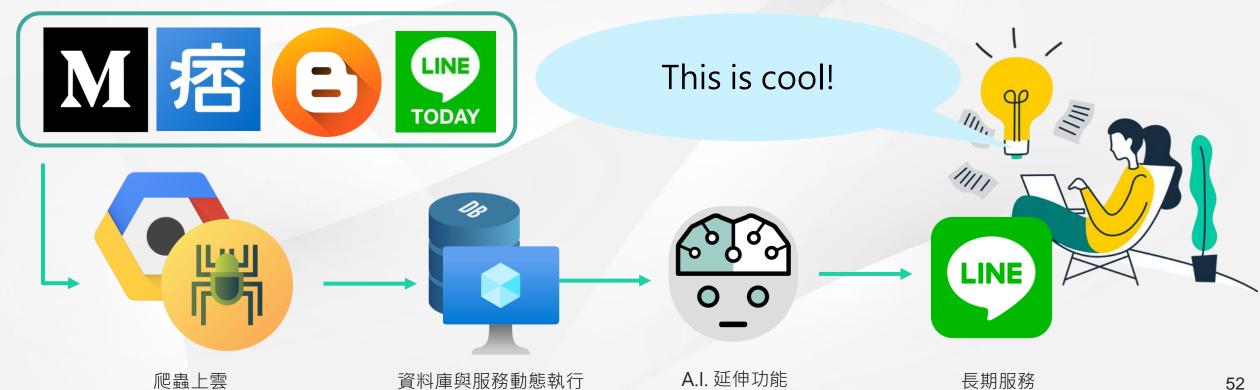




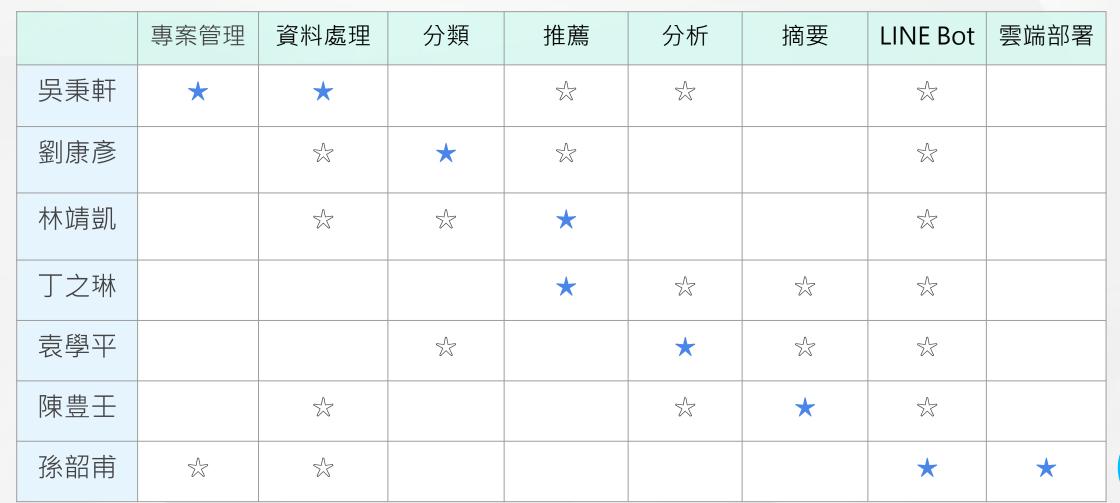
## 專案總結與未來展望

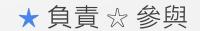
- □來源擴展與動態執行。
- □ 探索更多 A.I. 自然語言的功能延伸。





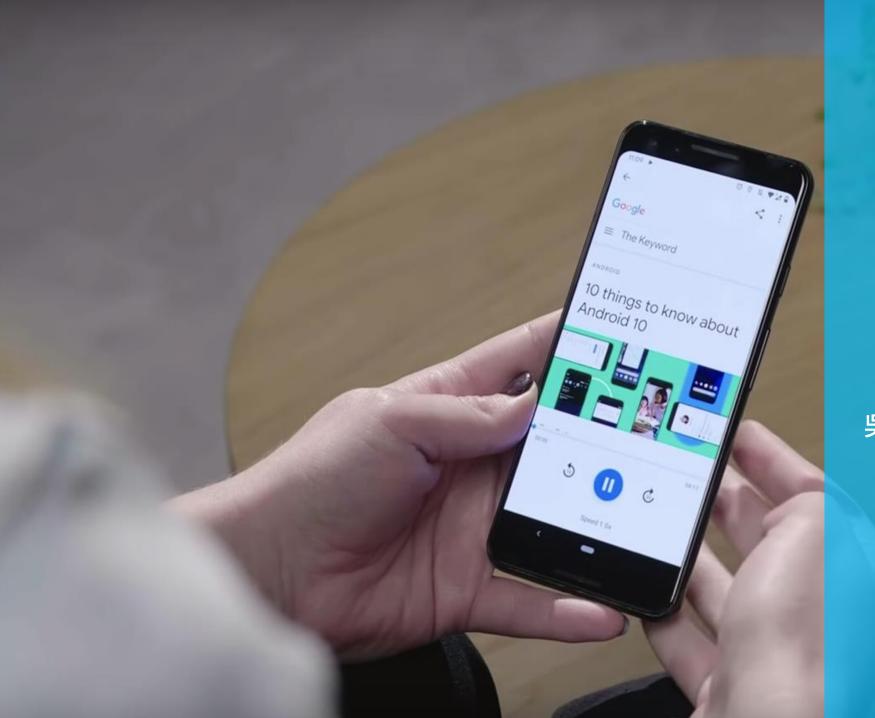
## 團隊分工











# 報告結束

吳秉軒、袁學平、劉康彥、孫韶甫 林靖凱、陳豊王、丁之琳

感謝您的聆聽