



**words.hk**

國  
語  
典

**Building Cantonese Dictionaries  
Using Crowd-sourcing Strategies:  
The words.hk Project**

Chaak-ming Lau  
*words.hk / CUHK*

*2017.03.17*

# Cantonese



Is it a language?

- No → End of story
- Yes → Continue

國  
方  
典

# What is a dictionary?

- Length
  - ❖ 50 words?
  - ❖ 1000 pages?
- Content
  - ❖ complete list?



# Group Discussion

- in groups of four
- discuss (for 5 minutes)
  - ❖ What should be included in a Cantonese dictionary
    - ❖ length, what kinds of words, what should be excluded
  - ❖ Write a sample entry for your dictionary
  - ❖ Use 1 minute to present your ideas

# Group Presentation

國  
方  
典

Let's turn to English ...

國  
方  
典

# Which of these words...

- do you expect in an English dictionary?
  - ✦ the
  - ✦ man
  - ✦ attention
  - ✦ bling
  - ✦ dim-sum
  - ✦ pragmatics
  - ✦ Justin Bieber

# Found one...

## the

determiner • **UK**  STRONG /ði:/ WEAK /ðə/ **US**  STRONG /ði:/ WEAK /ðə/

**the** *determiner* (PARTICULAR)

- ★ **A1** used before nouns to refer to particular things or people that have already been talked about or are already known or that are in a situation where it is clear what is happening:

*I just bought a new shirt and some new shoes. The shirt was pretty expensive, but the shoes weren't.*

*Please would you pass the salt.*

*I'll pick you up at the airport.*

- ★ **A1** used before some nouns that refer to place when you want to mention that type of place, without showing exactly which example of the place you mean:

*We spent all day at the beach.*

*Let's go to the movies this evening.*

*I have to go to the bank and get some Euros.*

# We have ...

## Jesus Christ

*noun* • **UK**  /,dʒiː.zəs 'kraɪst/ **US**  /,dʒiː.zəs 'kraɪst/ (ALSO Christ, )(ALSO Jesus)

- ★ the man believed by his religious followers to be the son of God. Christianity is based on his life and teachings.

# but not ...

## Search suggestions for Justin Bieber

We have these words with similar spellings or pronunciations:

- 1 just in time
- 2 just in case
- 3 just-in-time
- 4 justifiable
- 5 justifiably
- 6 justified
- 7 justifies

# man

*noun* • **UK**  /mæn/ **US**  /mæn/ PLURAL **men** **UK**  /men/ **US** 

**man** *noun* (MALE)

★ **A1** [C] an adult male human being

# dim sum

*noun* [U] • /,dɪm 'sʌm/ /,dɪm 'sʌm/

★ a Chinese meal or snack of small dishes including different steamed or fried foods:

# pragmatics

*noun* [U] • /præg'mæt.ɪks/ /præg'mætɪ.ɪks/ SPECIALIZED

★ the study of how language is affected by the situation in which it is

- **Found**

- ✿ the (grammatical words)
- ✿ man (basic words)
- ✿ attention (harder words)
- ✿ bling (slangs)
- ✿ dim-sum (recent loanwords)
- ✿ pragmatics (specialised words)
- ✿ Jesus Christ (important proper names)

- **Not Found**

- ✿ Justin Bieber

# Look at another entry

*What language is this?*

**bling**

*noun* [U] • **UK**  /blɪŋ/ **US**  /blɪŋ/ INFORMAL

★ **jewellery or decoration that attracts attention because it is very noticeable and looks expensive:**

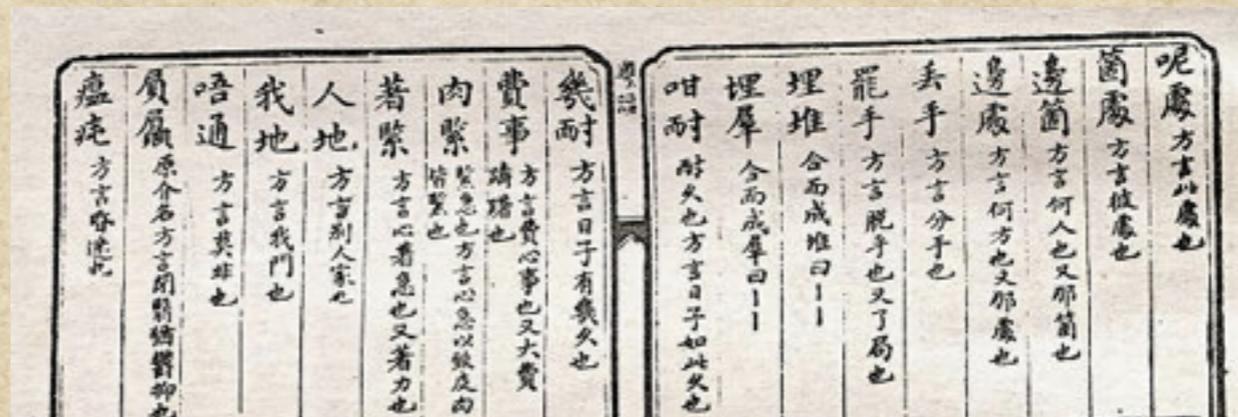
*She wore a fake-fur coat, big sunglasses and lots of bling.*

*and what language is this?*

# Expectation

- English dictionaries
  - ❖ Comprehensive
  - ❖ Written in English
- Cantonese dictionaries
  - ❖ Same principles?
- Question: What should a Cantonese dictionary look like?

# Existing Cantonese Dictionaries



【蔽翳】pei<sup>4</sup>-ai<sup>4</sup> ①憂愁，發愁：使乜咁～啊？ ②使人發愁的事：冇咁多～啊！

【冇得】mou<sup>5</sup>dak<sup>1</sup> [動][口](動詞の前に用いて)…する事・物がない、…のしようがない、…できない、無理である。[語法]動詞の前に用いて動作・行為の及ぶ対象がないことを表す。|| 冇好多怪事發生係冇得解嘅。Yau<sup>5</sup>hou<sup>2</sup>do<sup>1</sup>gwaai<sup>3</sup>si<sup>6</sup>faat<sup>3</sup>sang<sup>1</sup>hai<sup>6</sup>mou<sup>5</sup>dak<sup>1</sup>gaai<sup>2</sup>

國  
 方  
 典

# Monolingual dictionaries?

- Issues

- ❖ not comprehensive, only contain:

- ◆ words that do not overlap with "mainstream Chinese", i.e. Mandarin
- ◆ or harder words

***objection***

- ◆ Sorry we don't have that word, have you tried a French dictionary

- ❖ Descriptions written for Chinese / Mandarin users

# Online resources

- There is a nice one, don't waste your time!
- Character-based, not word-based

粵語審音配詞字庫

漢語多功能字庫  
古文字彙、多義詞解  
英語索引、粵語審音配詞

設定語

使  
用凡  
例

意見  
箱

香

部首: 香 [186] 筆畫: 9 字音分類: 單讀音字

大五碼: ADBB 倉頡碼: 竹木日 頻率 / 頻次: 853 / 2363

中文字源 國語辭典

CEDICT 林語堂

音節 (香港語言學學會)	粵音	根據	同音字	相關音節	詞例(解釋) / 備註
hoeng1		黃(p.43) 周(p.200) 李(p.68) 何(p.313)	腳, 鄉, 蕓	--選擇--	香水, 香味, 香 火[2.]

搜索次數: 65764 (管理人員專用區)

國  
方  
典

國  
語  
典

# The Project

Goal

Compilation

The Role of Social Networks

# words.hk

- **Build a true Cantonese dictionary**
  - ❖ **C1: Written in Cantonese ;**
  - ❖ **C2: Comprehensive, i.e. all common Cantonese words, spoken or written / unique to Cantonese or shared with Mandarin, should be included; and**
  - ❖ **C3: Words (rather than characters) will be the default unit**

# Compilation

- NOT corpus-based compilation
- Existing bilingual dictionaries
  - ❖ HKUDict (Luke et al., unpublished dictionary data)
  - ❖ A Dictionary of Cantonese Slang (Hutton & Bolton, 2005)
  - ❖ Lexical Lists for Chinese Learning in Hong Kong (Education Bureau, 2007, 2009)
- Existing word-lists
  - ❖ 現代標準漢語與粵語對照資料庫 (CUHK, 2001)
  - ❖ Word lists from various Input Method Engines

# Compilation (con't)

- Imported entries are
  - ✿ used as reference only
  - ✿ rewritten or restructured manually by editors

# The role of social networks

- Facebook

- ❖ Page (words.hk)  **words.hk**

- ◆ 10,000+ likes as of today

- ❖ Member recruitment

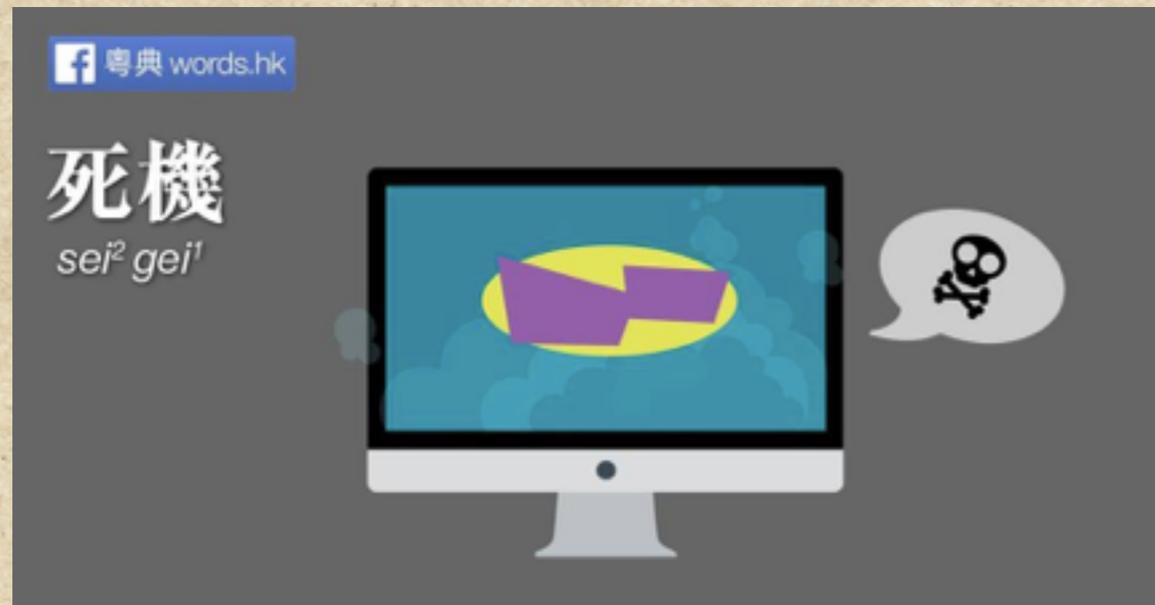
- ◆ Almost all members are recruited on Facebook

- ❖ A word a day

- Wordathon meetings

- ❖ Next: 18 Mar 2017 (Sat.) 3pm @Mong Kok

# A word a day - 每日一字



Illustrator: Can's Illustrations

words.hk | Cantonese Dictionary | Cantonese Corpus  
brought to you by Hong Kong Lexicography Limited

# A word a day - 每日一字



Illustrator: P of KatLau

APR 22 金魚佬  
gam<sup>1</sup> jyu<sup>4</sup> lou<sup>2</sup>

名 以「睇金魚」為名，誘拐、非禮細路女嘅男人；亦都指有變童傾向嘅人（量詞：個）  
men who abduct or harass young girls under the pretext of inviting them to 'see goldfish'; also refers to pedophiles

妹妹仔，想唔想跟叔叔上樓睇下我啲金魚？  
Hey little girl, do u want to come upstairs with uncle and see some goldfish?

# Sample Entries

食  
sik<sup>6</sup>

【動】1. 經食道吞下 **to eat** 我未食過臭豆腐，好想試下。 *I've never eaten stinky tofu. I really want to give it a try.* 你食咗晏未？ *Have you had lunch yet?* 你食咗藥未？ *Have you taken your medicine yet?* 2. 經氣管或其他方法吸入或注入身體 **to inhale, inject, or otherwise absorb into the body** 食煙 *to smoke (cigarettes)* 食白粉 *to shoot heroin* 3. 承受痛苦 **to suffer; to be on the receiving end of bad things** 我食咗佢好多拳。 *I took a lot of punches from him* 業主又要加租，小商戶冇議價能力，唯有硬食。 *The landlord is increasing the rent again. With no bargaining power, the tenants have no choice but to accept it.* 4. 消耗、耗用 **to consume, to expend** 呢隻app好食電 *This app consumes a lot of power* 食晒我啲時間 *taking up all my time* 5. 征服、吞併、控制、抽取 **to conquer; to embrace, extend and extinguish; to control; to extract something from** 呢個明星恃住自己有錢又靚仔，食咗好多條女 *Using his*

# 「食」

出fb post版本

收返埋

修改

刪除

解釋 #1

讀音： sik<sup>6</sup> 

詞性： 動詞

解釋： 1. (yue) 將一啲嘢擺入口到咬，然後經食道吞入肚裏面  
(eng) to eat

配詞 / 用法：

(yue) 食飯 (sik<sup>6</sup> faan<sup>6</sup>)

(eng) to have a meal

(yue) 食晏 (sik<sup>6</sup> aan<sup>3</sup>)

(eng) to have lunch

例句：

(yue) 我未食過臭豆腐，好想試下。 (ngo<sup>5</sup> mei<sup>6</sup> sik<sup>6</sup> gwo<sup>3</sup> cau<sup>3</sup> dau<sup>6</sup> fu<sup>6</sup>, hou<sup>2</sup> soeng<sup>2</sup> si<sup>3</sup> haa<sup>5</sup>.)

(eng) I've never eaten stinky tofu. I really want to give it a try.

2. (yue) 服用藥物  
(eng) to take (medicine)

例句：

(yue) 你食咗藥未？ (nei<sup>5</sup> sik<sup>6</sup> zo<sup>2</sup> joek<sup>6</sup> mei<sup>6</sup>?)

(eng) Have you taken your medicine yet?

3. (yue) 經氣管或者其他方法吸入或注入身體  
(eng) to inhale, inject, or otherwise absorb into the body

配詞 / 用法：

(yue) 食煙 (sik<sup>6</sup> yan<sup>1</sup>)

## Online Interface

- <http://beta.words.hk/食>

# Status of the project

- On-going project
  - ❖ Active discussion every day
  - ❖ Bi-weekly wordathon meetings
- ~500 volunteers
- 42000+ entries
  - ❖ 34000+ edited w/ Cantonese definitions

詞  
方  
典

# Methodology

Agile vs. waterfall

Core Principles

*Usage-over-etymology*

*Decision Problem Avoidance*

# Methodology

- Not executed by step-by-step plan
- Agile vs. waterfall
  - ✿ Concept borrowed from software engineering and lean start-ups
  - ✿ Ensure maximum flexibility
- Adhere to Core Principles

# Core Principles

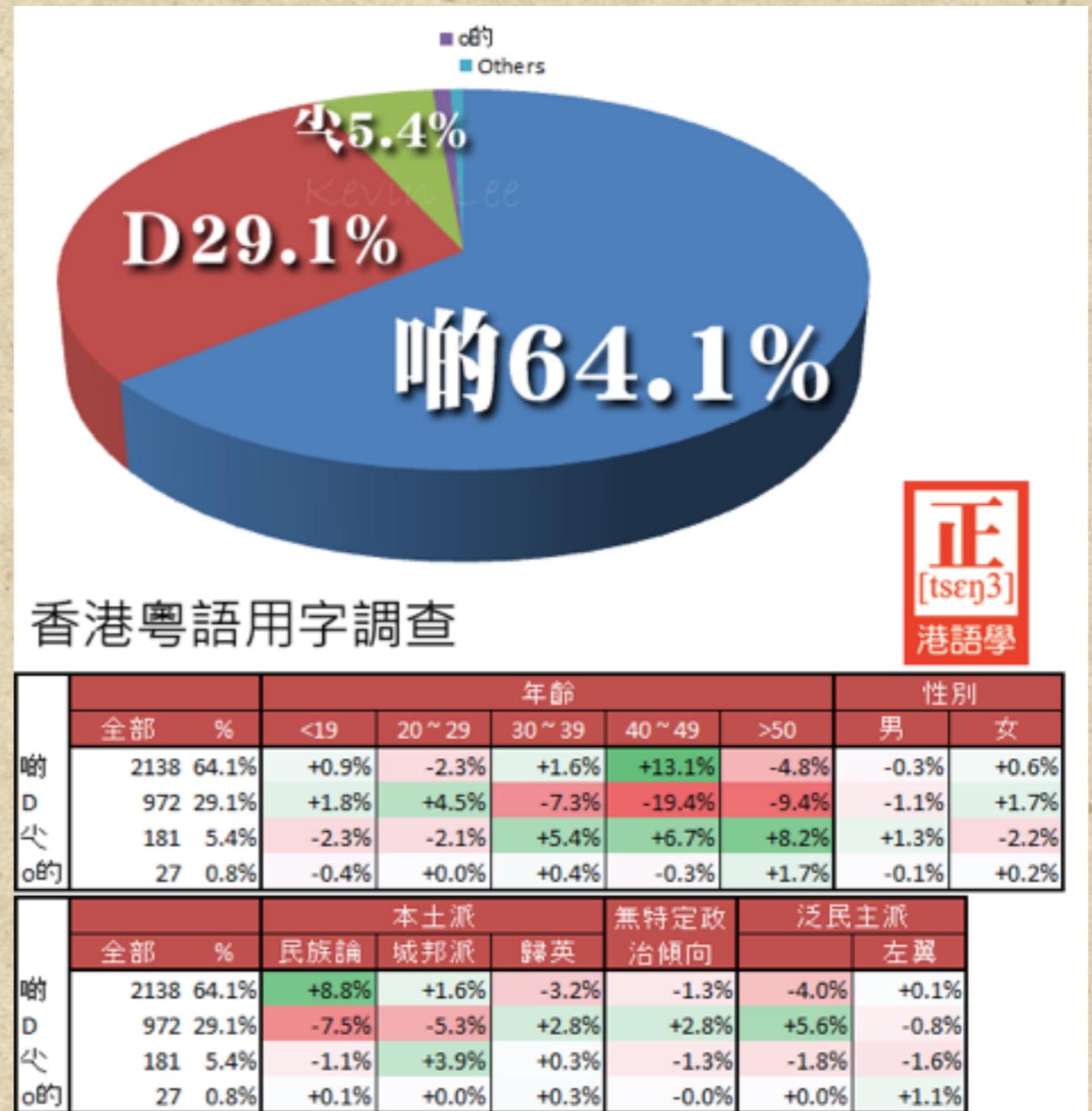
- Usage-over-etymology
  - ✦ ignore etymology
  - ✦ have everything based on usage
  - ✦ explaining actual usage is preferred over existing dictionaries or etymologies
- Decision Problem Avoidance
  - ✦ Decision problems:
    - ◆ something that requires a judgement
    - ◆ e.g. classification
  - ✦ Avoid these problems

# Usage-over-etymology

- e.g. Which is the correct character?
  - ❖ The case of 啲 di1
    - ◆ Possible characters
      - ❖ 啲、D are commonly used
      - ❖ 𠵼 is the proposed character by amateur etymologists
      - ❖ 的 is used in older texts
      - ❖ 滴 is probably the true etymology, but no longer used / recognised by native speakers

# Usage-over-etymology (con't)

- ◆ 啲 / D / 㗎 are listed, ordered by frequency
- ❖ based on editors' subjective judgement
- ❖ agrees with Cantonese character usage survey conducted by 港語學



國語典

# Decision Problem Avoidance

- Cantonese vs. Non-Cantonese
  - ❖ Major problem faced by previous Cantonese dictionaries
    - ◆ Three-way distinction (鄧思穎 2015, Ch. 3)
      - ❖ 通用詞 shared vocabulary
      - ❖ 社區詞 community vocabulary
      - ❖ 方言詞 topolectal vocabulary
- Issues
  - ❖ implies a decision problem for every entry created
  - ❖ no easy test: prone to dispute
- Solution
  - ❖ Only label items not used at all in spoken Cantonese

# Decision Problem Avoidance (con't)

- Other problems

- ❖ word-class classification

- ◆ noun or verb? verb or adjective?

- ❖ headword problem

- ◆ should variants, e.g. 擦膠 (caat3 gaau1, eraser), 膠擦 (gaau1 caat3), 擦紙膠 (caat3 zi2 gaau1), listed separately?

- ◆ should K (kei1, karaoke) listed separately?

- ❖ or just under K房 (kei1 fong2, ~room),  
唱K (coeng3 kei1, to sing~),  
劈K (pek3 kei1, to get drunk at ~)?

# Decision Problem Avoidance (con't)

## ✿ proper nouns

◆ 周公 zau1 gung1

◆ 青山 cing1 saan1

## ✿ jargons

◆ e.g. jargons within political circles?

❖ e.g. 左膠 zo2 gaau1

- (Look it up on [words.hk](http://words.hk))

國  
語  
典

# Conclusion

Issues

Future of Cantonese Lexicography  
On Open Data

# Issues

- **Lexicography**
  - ❖ **Over-emphasis on Etymology / Characters**
  - ❖ **Prescriptivism**
  - ❖ **Fail to catch up with latest language developments**
- **Lack of Available resources**
  - ❖ **Lock-ups**
  - ❖ **Paper-based**

# Future of Cantonese Lexicography?

- **Word-based dictionaries**
- **Use both word-list compilation and corpus-based approach**
- **Crowd-sourcing to generate content**
- **Let the users decide the preferred written form**

# Open Data Policy

- **words.hk License Agreement**
  - ❖ <http://beta.words.hk/base/hoifong/>
  - ❖ dictionary data are
    - ◆ freely accessible to the public
    - ◆ can be freely modified, adopted, and published for non-commercial and certain small-business purposes
- **avoid lock-ups**

# Appeal to the community

- **Locked-up resources**
  - ❖ never been published, protected by copyright
  - ❖ unusable
- **Say goodbye to the old “proprietary” age!**

# Finally

- Get the priorities right
  - ❖ w.r.t. Cantonese lexicography
  - ❖ it is language description
    - ◆ i.e. keep a record of:
      - ❖ usage
      - ❖ diversity
  - ❖ NOT decision problems
    - ◆ folk etymology (like where do the ?)
    - ◆ existing 本字考?



# Reference

- Education Bureau (2007, 2009). *Lexical Lists for Chinese Learning in Hong Kong*. Available on <http://www.edbchinese.hk/lexlist/>
- Ferguson, C. (1959). *Diglossia*.
- Hutton, C., & Bolton, K. (2005). *A dictionary of Cantonese slang: the language of Hong Kong movies, street gangs and city life*.
- Snow, D. B. (2004). *Cantonese as written language the growth of a written Chinese vernacular*.
- 招子庸. (1821). 粵謳.
- 鄧思穎. (2015). 粵語語法講義. 香港: 商務印書館.