# Towards a sustainable digital infrastructure for historical Chinese texts

**Donald Sturgeon**

**Fairbank Center for Chinese Studies**

**Harvard University**

sturgeon@fas.harvard.edu

```
{
    "author" : "王真",
    "dynasty" : {
        "from" : {
            "id" : "16",
            "name" : "Tang"
        },
        "to" : {
            "id" : "16",
            "name" : "Tang"
        }
    },
    "edition" : {
        "collectionid" : "72",
        "title" : "指海",
        "url" : "http://ctext.org/library.pl?...
    },
    "lastmodified" : "2015-03-22 09:02:33",
    "tags" : "OCR_MATCH",
    "title" : "道德經論兵要義述",
    "toptitle" : "道德經論兵要義述",
```

# Chinese Text Project (ctext.org)

- Largest full-text digital library of pre-modern Chinese
  - 25 million pages of scanned primary source texts
  - 5 billion characters of transcription
  - Used daily by 25,000 people

- 2005-2013:       primarily *static* database
- 2013-present:    major *dynamic* components

- Originally: centrally edited and maintained
- Now: primarily maintained by crowdsourcing

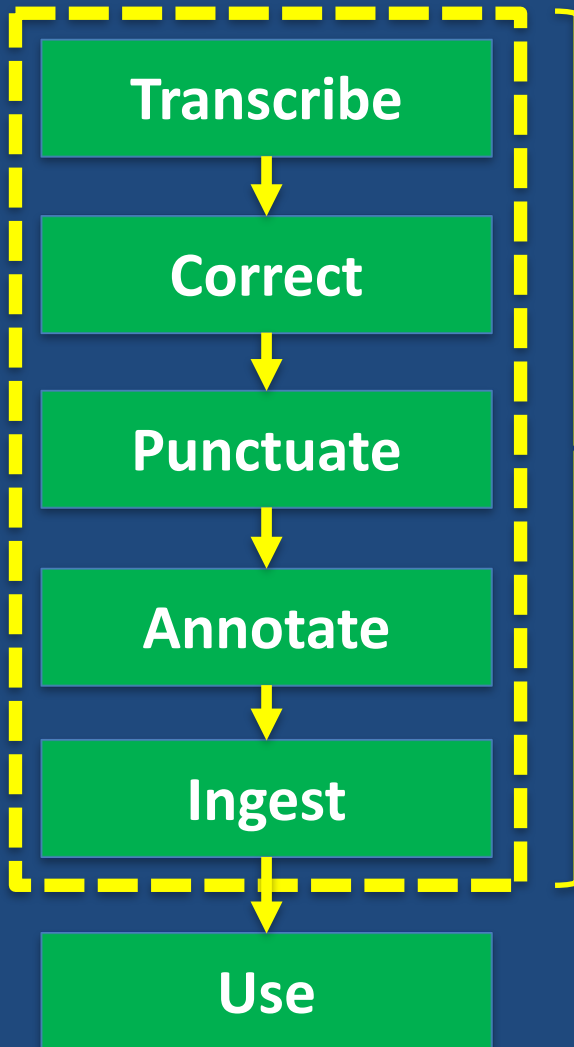# ctext.org as digital infrastructure

- Key components:
  - Optical Character Recognition
  - Crowdsourcing system
  - Application Programming Interface
- Example use cases:
  - Searchable database of texts
  - Scalable transcription tool
  - Repository for data mining

# Digitization

- Scanning historical materials *relatively* cheap
  - Established technology
  - Storage & processing costs decrease exponentially
- Scanned images for conservation of material
  - Access to material (e.g. online image browsers)
  - No risk of damage to rare or unique objects
- Large-scale scanning projects undertaken
  - E.g. Harvard-Yenching: 5 million+ pages
  - Even larger scanning projects in Mainland China
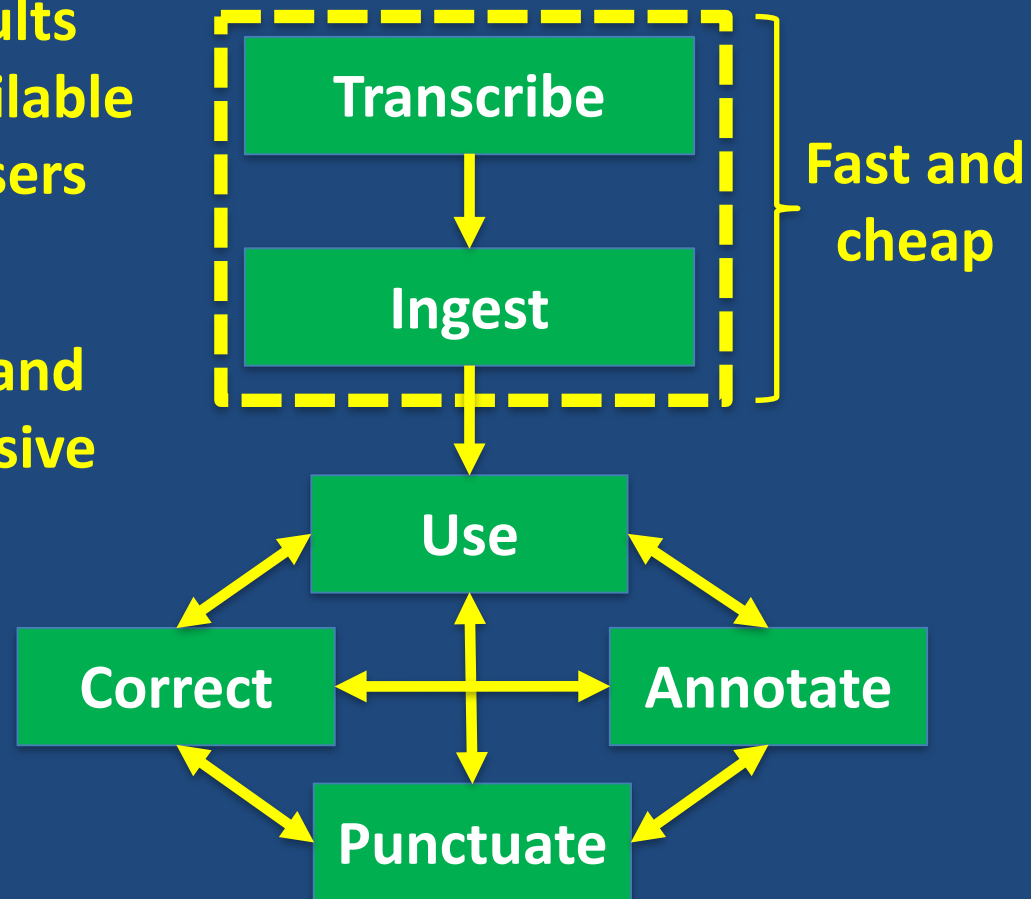
# Static versus dynamic

## Static

**Transcribe**

↓

**Correct**

↓

**Punctuate**

↓

**Annotate**

↓

**Ingest**

↓

**Use**

**Results unavailable to users**

**Slow and expensive**

## Dynamic

**Transcribe**

↓

**Ingest**

**Fast and cheap**

↓

**Use**

**Correct** ↔ **Annotate**

**Punctuate**

# The Long Tail

**Popularity: frequency of use**

**Mainstream texts and editions: already transcribed**

**Texts to be transcribed in next 10 years**

**Everything else: inaccessible and obscure**
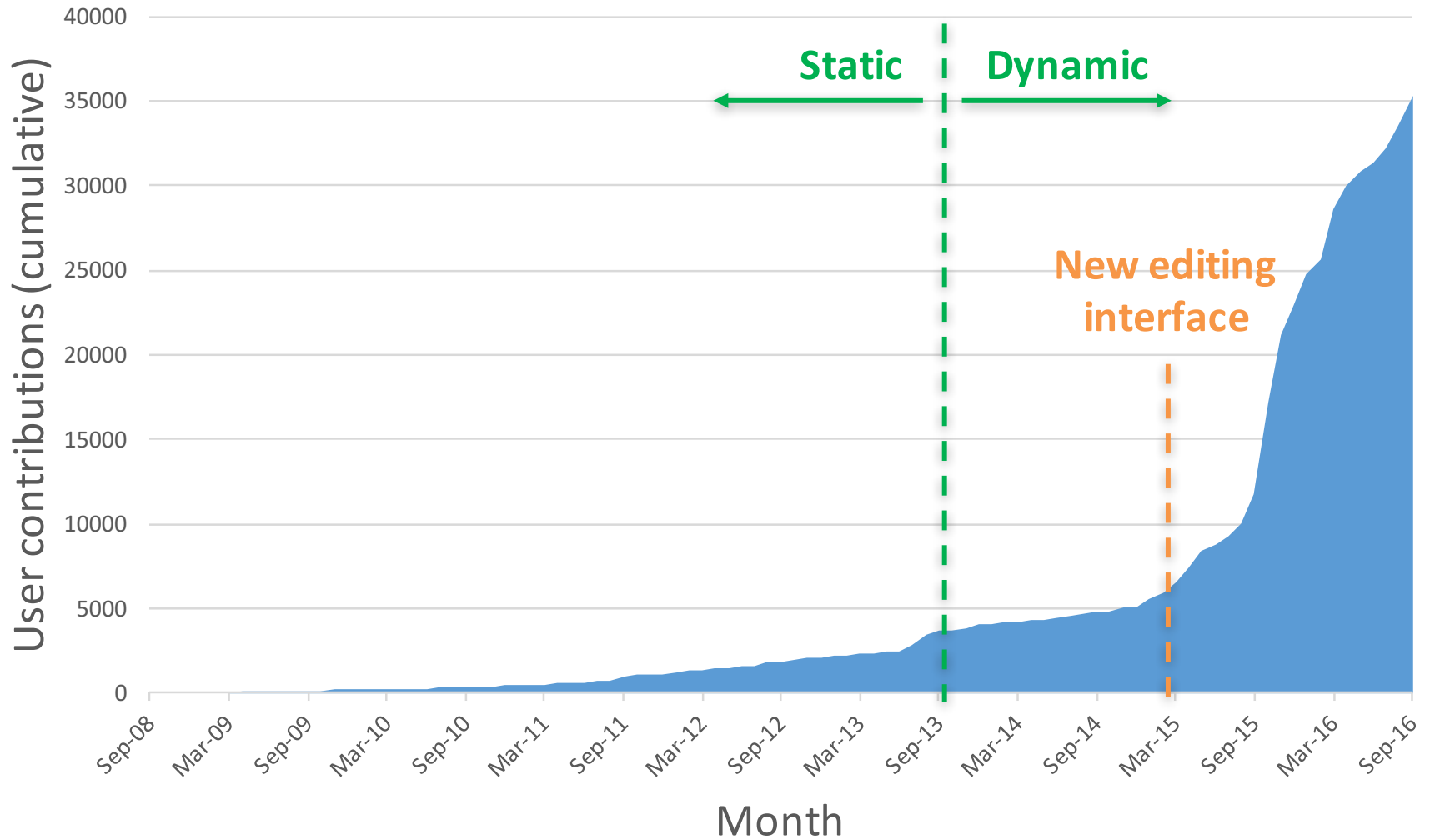
**Texts and editions**

# Database size

# Database usage

# User contributions

# Domain-specific OCR

- Leverage domain knowledge
  - Language, writing conventions
  - Existing transcriptions, text reuse, other editions
- Easier problem => more accurate results
  - Compare with e.g. Google Books OCR
- Enables image search & basis for transcription
  - Frequent use => more opportunities for correction

# ctext.org vs Google Books

**Google books scan**

**Google books OCR result**

一外佴進內他琇事關緊要該班之王大臣等每於日
1.

不成事體矣嗣後進班之王大臣等均＾午刻進內

進值班儻有午後＾行進賽者卽行＾溱棻此

旨從前因媼布明德年老是＾每逢內務府大臣

4
步
＾
13 力
, X ＾＾,
＾ ||
顷丙值变暴＾＾崔內.

# ctext.org vs Google Books

Google books scan

ctext.org OCR result (unedited)



壇齋戒及圓明園駐蹕之日除內務府大臣進班之鹽
均未年老理應親身進班嗣後凡股宿
之期令內務府郎中等補班矣現在內務府大臣體
百從前因編布明德年老是以每逢內務府大臣進四
　　奉
值班儻有午後始行進班者即行恭奏欽此下月
不成事體矣嗣後進班之王大臣等均於午刻進內
入始行進內次日天明開門即出此皆朕所深知童
外但進內值班事關緊要該班之王大臣等每於日

# Dynamic approach: advantages

- Imperfect data can still be useful
  - E.g. OCR-derived text enables full-text search

《二十二》 [View] [Edit] [History]

⚠ Transcribed automatically with OCR. Please help correct any errors.

1 余常以道德扣諸老宿乃曰道何物耶依之而心修從之而理順德何物
耶布之而利傳積之而行圓反是二者則聖賢不取焉伽蓋者著明道德
之大宅也酉竺聖賢自圓覺而為之充三際遍十虛實一心成禹德故在
處僧伽藍與天地相為無始者無他蓋道德之自任也嘉定州在吳郡之
東南百里形勢平尋旱湖暮汝風帆浪船宜接城治大報國圓通寺際州
治之東北相距咫尺開山沙門明了族高氏壯年極厭塵氛禮枕州般若
寺住持愚叟賢公薙髮至元丙戊手鋤福翳開招經營廣堂以宇宛若化
成大德已束春欽奉璽書賜圓通總額越七年丙午入覲明年丁未冬武
宗皇帝加賜今額錫妙明圓悟佛心之號及欽受今上潛邸腸旨讀持至
皇慶壬子造物欲大其規制一夕祝融卷人加何明年癸丑奮志異修首
定向方竭匠氏之智取東山之材持漫陶冶百爾成臻延祐丙申復舉奉
旨如腸妙明圓悟普濟佛心大禪師之號感恩優異恩報無所一閱暑寒

# Dynamic approach: advantages

- Scans can be matched to existing transcription
  - Errors limited to minor alignment errors
  - Punctuation and structure preserved

# Dynamic approach: challenges

- User-submitted content may be incorrect
  - Need simple way of visualizing & verifying any edit

# Dynamic approach: challenges

- Typically need "human readable" serialization
  - Can be simplified using task-specific visual editors

# From database to platform

- Enable use as transcription & annotation tool
  - Image and textual data can be uploaded by users
  - OCR results can be corrected collaboratively
  - Future: semantic markup, translation, etc.
  - Full-text data can be exported via API
- Time and effort invested is preserved
  - Corrections *also* become part of public database
  - "Long tail" of non-mainstream material available

# Psychological factors

- Creating *new* content vs correcting *incorrect* content
- Unpredictable response vs instant gratification
- Incremental improvement vs all-or-nothing approach



https://xkcd.com/386/

# OCR + Crowdsourcing

# Task-specific visual editing tools

# Collation of crowdsourced data



ctext:2494

| | |
|---|---|
| **Composition:** | ⿲冂衤奠 |
| **Radical:** | 衣 |
| **Stroke count:** | 16 |

**Example glyphs**

《虞東學詩》 禖禖禖

**Example usage**

| | | |
|---|---|---|
| 《周禮述註·卷十六》 | : | ...以蒲為蔽天子喪服之車漢儀亦然犬禖以犬皮為覆笒玄謂蔽車旁禦風塵者...⊡ |
| 《周禮述註·卷十六》 | : | 註曰素車以白土堊車也�machine讀為蘋蘋麻以為蔽其禖服以素繒為緣此卒哭所乘為君之道 |
| 《周禮述註·卷十六》 | : | ○藻車藻蔽鹿淺禖革飾註曰藻水草蒼色以蒼土堊車以蒼繒為蔽也鹿淺禖以鹿夏皮為 |
| 《周禮述註·卷十六》 | : | ○駹車雚蔽然禖髤飾⊡ |
| 《周禮述註·卷十六》 | : | ○漆車藩蔽犴禖雀飾⊡ |
| 《周禮述註·卷十六》 | : | ...禫即乘漆車與吉同者禮窮則同也　禖莫歷反㬚音羔藻音藻雚音丸髤香求...⊡ |
| 《毛詩注疏·卷二十五》 | : | ...持車使牢固也幭字禮記作幦周禮作禖字異而義同玉藻言羔幭鹿幭春官巾...⊡ |
| 《御定駢字類編》 | : | 犬禖周禮巾車木車蒲蔽｜｜尾�013疏飾小服皆疏注｜｜以犬皮為覆笒 又素車�013蔽｜｜素飾小服 |
| 《讀詩略記·卷五》 | : | ...祀韓侯有錫此異數也幭孔疏云與幦禖同義玉藻羔幦鹿幦周禮巾車犬禖犴...⊡ |

21

# CTP API Overview

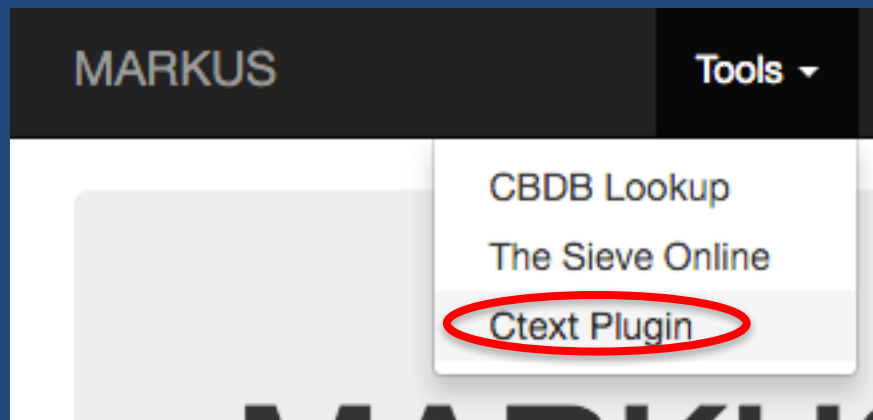- CTP URNs
  - Unique identifier for each edition & chapter of a text
- JSON API
  - Extract machine-readable data from ctext.org
  - Textual data, metadata, etc.
- Plugin system
  - XML description of how to connect to another site
  - Allows users to connect ctext.org to other projects
  - User-definable, sharable, automatically updated

# Plugins (User Perspective)

| Plugin | Description | Type | Example | Install |
|---|---|---|---|---|
| Plain text | Export a chapter as plain text. | chapter | [Plain text] | [Installed] |
| MHDB | MHDB character lookup. | character, word | [MHDB] | [Install] |
| MARKUS | Export a chapter to MARKUS (Requires Google Chrome) | book, chapter | [MARKUS] | [Install] |
| Frequencies | Compare character frequencies to corpus averages. | chapter | [Frequencies] | [Installed] |
| Glyphwiki | Glyphwiki character lookup. | character | [Glyphwiki] | [Installed] |
| CHISE | CHISE character lookup. | character | [CHISE] | [Installed] |

- Point-and-click installation within ctext.org
- Point-and-click installation from 3rd party site

筆 篔 笔 笔

U+7B46  Seal script  Semantic variant Simplified character

[MoE] [zdic] [Google Search] [Glyphwiki] [CHISE] [TLS] [Lin Yutang] [Unihan] [WWW-JDic] [Zhongwen.com] [MDBG] [CJKV-E] [MHDB] [Ricci] [Manage plugins]

**Radical:** 竹 + 6 strokes = 12 strokes total. 70#09

**References**

筆  1360C55  ㄅ丨 [bi3]. [Var. 笔]
92A.10-2/118

中央研究院近代史研究所 | 英華字典
INSTITUTE OF MODERN HISTORY, ACADEMIA SINICA

Re TL

An His
Concep
Genera
Jiang Sh

首頁  搜尋  瀏覽  說明
搜尋結果

1865馬禮遜五車韻府  1筆
1874司登得中英袖珍字典  1筆
1912翟理斯華英字典  3筆
全部列出

Chara

Unico

Radic

漢語

Look

詞目條列

搜尋「筆」共找到3筆詞目：

## Definitions

1. **Pinceau** (pour écrire); (*p. ext.*) tout instrument servant à écrire : **porte-plume, crayon**, stylographe, morceau de craie, *etc.*

2. a. **Écrire**; noter; mettre par écrit; composer; rédiger; décrire. Composition; style littéraire; (*fig.*) plume; écriture (de *qn*); genre de calligraphie.
   b. (*Litt.*) Ds 筆名 bǐ míng Nom de plume.

3. *Spécif. des coups de pinceau, des traits d'un car.*

4. *Spécif. des sommes d'argent.*

5. (*Litt. chin.*) Prose :
   a. *Ds* le 文心雕龍 *Wen Xin Diao Long*, 劉勰 Liu Xie classe sous le terme 筆 bǐ les genres littéraires en prose en les distinguant des genres poétiques qu'il classe sous la rubrique 文 wén.
   b. *Ds* 筆記 bǐ jì Essai en prose sous forme de notes.

1  to begin to write……（完整顯示）  投筆而書……（完整顯示）

# API and Plugins

# APIs in Teaching

```
In [3]: from ctext import *

        laozi = gettextasstring("ctp:dao-de-jing")

        for match in re.finditer(r"足[^。，、；！？]", laozi):
            matched_text = match.group(0)
            print(matched_text)
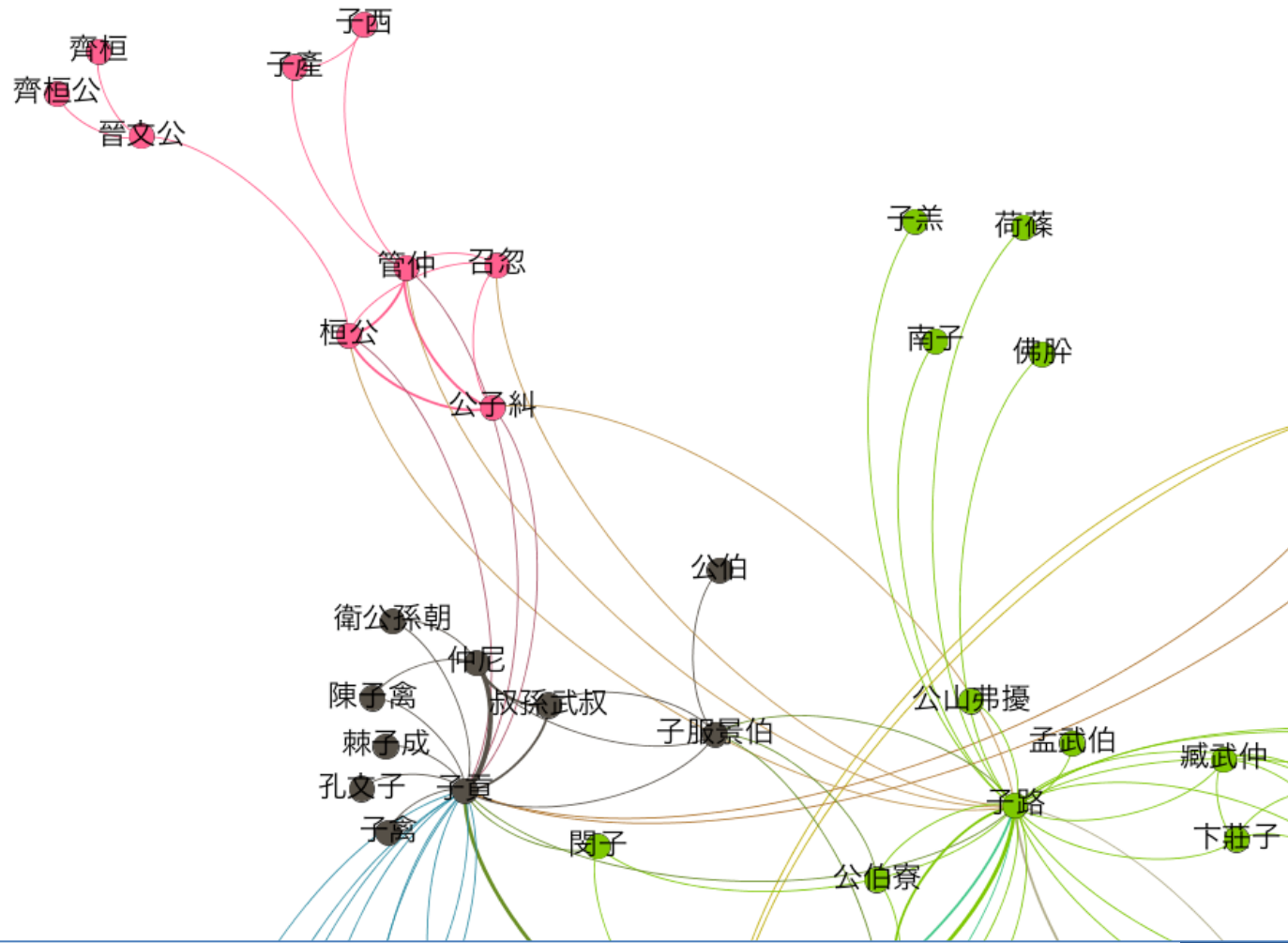```

足者
足見
足聞
足既
足以
足不
足之
足矣
足以
足下
足者
足以

# APIs in Teaching

# APIs in Teaching

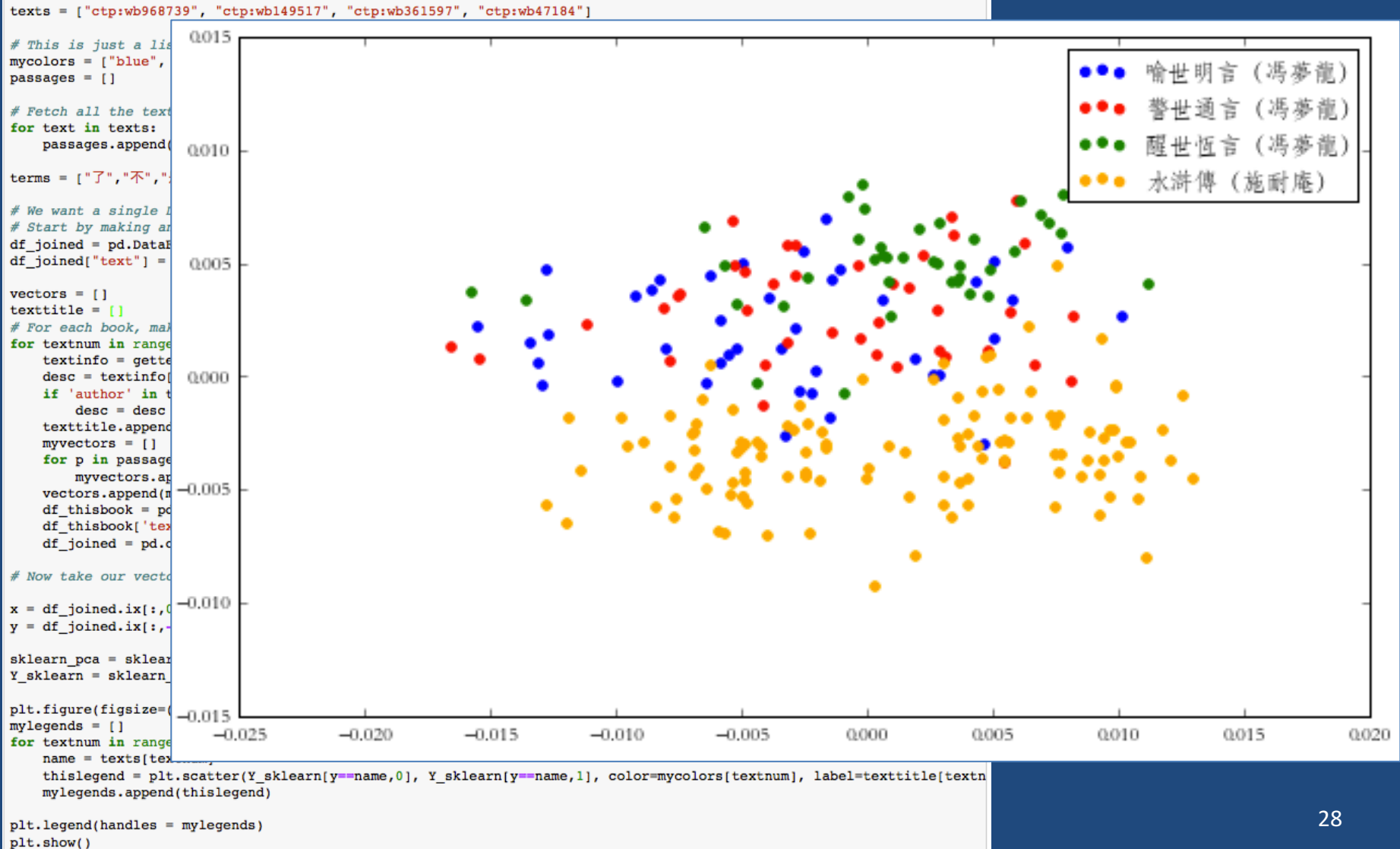# Conclusion

- OCR & crowdsourcing allow access to long tail
  - Obscure texts many of which are untranscribed
  - Infrastructure helps users help themselves
    - Side effect: their efforts benefit everyone else too
- API access
  - Machine-friendly interface enables new use cases
  - Allows use in many workflows
  - Streamline DH teaching and research
    - Avoids dilemma of toy examples vs data wrangling

More information:

**ctext.org**

Practical introduction:

**dsturgeon.net/ctext**