

Interdisciplinary Study

—the story of a computer scientist turned linguist or visa versa

Jason S. Chang

Department of Computer Science
National Tsing Hua University, TAIWAN

2015-0306 Friday 11:00-12:30

Waseda University, Tokyo

Peeking into the Future of Writing

—Introducing WriteAway, an Interactive Writing Environment

Jason S. Chang

Department of Computer Science
National Tsing Hua University, TAIWAN

2015-0306 Friday 11:00-12:30

Waseda University, Tokyo

Entire contents © 2013 Jason S. Chang. All rights reserved.



I wrote this month in **Scientific American**, Taiwan

I wrote this month in Scientific American, Taiwan



I wrote this month in Scientific American, Taiwan

• 雲端上的寫作教練

從泥版、竹簡、紙張到互動式寫作環境，寫作的面貌一變再變。





Always Online
網路不打烊
撰文／張俊雄

雲端上的寫作教練

從泥版、竹簡、紙張到互動式寫作環境，寫作的面貌一變再變。

《語言之始》(The First Word)引述最新研究，指出語言源自口腔肌肉內的「運動文法基因」。生物為了生存，應當有程序性動作，以調整身體、覓食求偶，而這些動作與語言，也源自口腔肌肉——就像電腦老是用那些微處理器。互覺上，「演化之師」應當安排這樣的智慧設計，幫助人類在嚴酷的競爭下過著生存。在時間與空間的轉折點，運動文法基因突變重組，進而指揮咽喉發聲，傳情表意，語言於焉誕生。

居無定所的游牧先民，呼聲引伴傳遞訊息，凝聚家族聯繫社群，語言無形體也足堪勝任。然而群聚定居形成城市後，語言受限於時空，不足以應付更大規模、更長時期的互動。此時需要電報已能「閱讀」文章，而我們也讀了電腦合成的財經新聞而不自知。

其三，越來越多電腦情緒「閱讀」大量文章，例如Google的知識圖譜計畫；反之亦然，試圖有多少讀者已不自覺地瀏覽電腦合成的財經新聞，否則你很難踩上雲端。在互動式寫作環境中和電腦合作共筆，也不足為奇。

前兩者已成常態，只會越演越烈。第三項機器閱讀、電腦共筆，則力尚未及，值得深究。美國教育考試服務社(ETS)舉辦托福之外，也開發Criterion，透過自動評分與回饋，輔助學習寫作。英國學者喬登與史密德則開發Marking Mate網站，免費提供類似服務。

台灣清華大學不遠人後，也躍身其中開發出互動式寫作環境WriteAway，讓寫作不再靜謐，不時得到提示，WriteAway隨時觀察輸入文字，建議相關的文法與例句，供後續與編輯的參考。更有趣的是，它並非一成不變地呈現建議，而是自動把最相關者前移。如此，寫作者不至於分神而放緩寫作。WriteAway的雛形是由Citiusse學術資料庫與學術寫作文法，符合論文寫作之需，看來已有「極度可行產品」之姿。假設你在WriteAway輸入台式英文“This paper present method”，沒有提醒，也不知如何接續。沒關係，系統會顯示建議如method for doing something. This paper presents a method for generating solution”。不難想像，眼角一掃建議，你可立即改為“This paper presents a method for finding answers to...”，一舉避免三個動詞、冠、介詞的文法錯誤。

或許，未來會冒出更多更多互動式寫作環境。有一天，雲端教練會定時提醒你開始寫作；寫得精采時，為你擊掌鼓勵；寫到低潮時，幫你消除可能的作者障礙。今天，不妨先到writeaway.zipoweb.org體驗一下寫作的未來。

張俊雄是清華大學資訊系教授，研究興趣包括自然語言處理、電腦輔助學習、機器翻譯、研究之類，也翻譯了世界史小說《人與文明》、2004年出版《科學人》雜誌與科學史文集，也可以說「讀過半部史、讀過半部科學史、讀過半部文學史」。

2013.03

I wrote this month in Scientific American, Taiwan

- 雲端上的寫作教練 從泥版、竹簡、紙張到互動式寫作環境，寫作的面貌一變再變。
- Your Person Writing Trainer on the Cloud
Many faces of writing: from Clay Tablet to Interactive Writing Environment (IWE)





Always Online
網路不打烊
撰文／張俊雄

雲端上的寫作教練

從泥版、竹簡、紙張到互動式寫作環境，寫作的面貌一變再變。

《語言之始》(The First Word)引述最新研究，指出語言源自腦內「運動文法基因」。生物為了生存，需要具備程序性動作，以調整身體、覓食求偶，而這些動作與語言，也源自於相同的基因——就像電腦老是用那些微處理器。互覺上，「演化之師」應會安排這樣的智慧設計，幫助人類在嚴酷的競爭下過著生存。在時間與空間的轉折點，運動文法基因突變重組，進而指揮與調整，再往前進，語言於焉誕生。

話無定用的遺世先民，呼聲引伴傳達訊息，凝聚家族聯繫社群，語言無形中也承襲責任，然而群聚定居形成城市後，語言受限於時空，不足以應付更大規模、更長時間的互動。此時需要再往前一步，才得以維繫文明所需之記事、符號、契約，社會演化壓力加重，「錄音」技術突破——先民開始「錄製」語言，到龜甲、銅幣、泥版、竹簡，開啟了故事、文字、書籍的歷史新頁。語言的變遷如此神妙，引發「巴比倫塔興建崩壞」、「金銀字號鬼哭神號」的傳奇。

波多斯基(Henry Petroski)在《利器》倡言，凡工具必如生物般演化，語言演化已無庸贅言，寫作工具亦然。千年前之筆墨紙張儲存，時人卻多敲擊鍵盤、觸摸螢幕以為書寫；30年前打字機風光一時，如今只能在博物館偶遇，一代之間竟見證了書寫由紙筆變為打字，再跨入電腦列印、網路發傳。寫作媒體與工具的變遷演化，不可不察。

寫作的下一步為何？安得水晶球，一窺寫作的未來？那無「清晰、望遠」的水晶球，還是依稀可見未來。其一，越來越多人用非母語的英文溝通全球；隨之，學生也需為考試而學習外語寫作。其二，越來越多人透過電腦網路或寫作，有據或無據，在開放式線上課程聽講答題，在維基百科上編寫修改文章，或在StackOverflow為文解答問題。

其三，越來越多電腦情緒「閱讀」大量文章，例如Google的知識圖譜計畫；反之亦然，試著有多少讀者已不自覺地瀏覽電腦合成的財經新聞，若說你該對線上書報，在互動式寫作環境中和電腦合作共筆，也不足為奇。

前兩者已成常態，只會越演越烈。第三項機器閱讀、電腦共筆，則力尚未及，值得深究。美國教育考試服務社(ETS)舉辦托福之外，也開發Criterion，透過自動評分與回饋，輔助學習寫作。英國學者費登與史密德則開發Making Made網站，免費提供類似服務。

台灣清華大學不為人知，也隱身其中開發出互動式寫作環境WriteAway，讓寫作者不停擊鍵，不時得到提示，WriteAway隨時觀察輸入文字，建議相關的文法與例句，供後續與編輯的參考。更有趣的是，它並非一成不變地呈現建議，而是自動把最相關者前移。如此，寫作者不至於分神而放棄寫作。WriteAway的雛形是由Citiusc 學術資料庫與學術寫作文法，符合論文寫作之需，看來已有「極度可行產品」之姿，假設你在WriteAway輸入台式英文「This paper present method」，沒有提醒，也不知如何拼綴，沒關係，系統會顯示建議如method for doing something. This paper presents a method for generating solution'，不難想像，眼角一掃建議，你可立即改為「This paper presents a method for finding answers to...」，一舉避開了三個動詞、冠、介詞的文法錯誤。

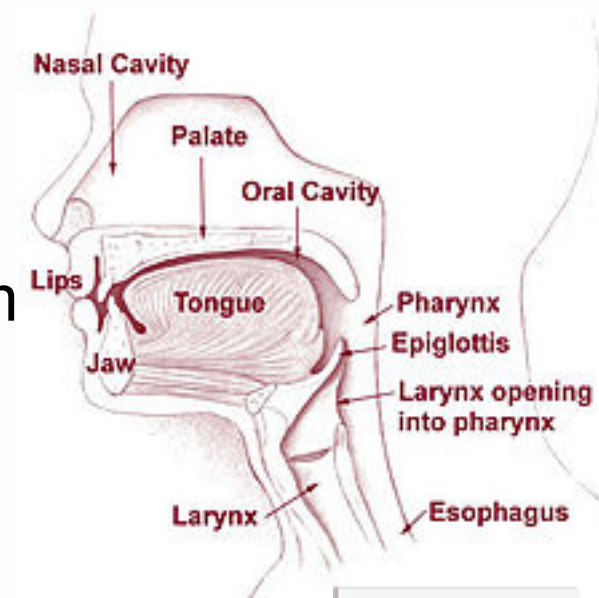
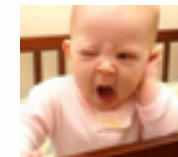
或許，未來會冒出更多更多互動式寫作環境。有一天，雲端教練會定時提醒你開始寫作；寫得精采時，為你擊掌鼓勵；寫到低潮時，幫你消除可能的作者障礙。今天，不妨先到 writeaway.zipoweb.org 體驗一下寫作的未來。

張俊雄是清華大學資訊系教授，研究興趣包括自然語言、電腦輔助學習、機器翻譯、研究之類，也翻譯了世界史小說《人與文明》、2004年出版《科學人》雜誌撰寫科技文章，也可以到「網路不打烊」專欄，「科技閱讀」欄刊他的文章。

2013 年 3 月

How did Language emerge? evolve?

- Experts suggest that **Language** emerged (*The First Word* by Christine Kenneally)
 - Genes that control **MOTION GRAMMAR** (e.g., baby yawning)
- Motion of muscles --> motion of the **vocal chord**
- Language started with shapeless, volatile **Speech**
- **Civilization** needed something more concrete: **Writing**
- **Language** has been evolving much like plants and anim
- **Writing** changed more drastically than speech
 - Finger, stylus, quill, pen, ball pen
 - Clay, tablet, paper
- **Today** we use fingers
 - Keyboards, soft keyboards, touch screens
 - Writing appears in personal disk files, emails, blogs, Facebook (paperless)



New faces of Writing

- We have new ways of using **Writing**
 - Collaborative writing: [Wikipedia](#)
 - Community-based QA: [Stackoverflow.com](#)
 - Presenting yourself: [Blogs](#)
- Academic writing:
 - [Marking Mate](#)
 - [Criterion](#)
 - *[WriteAway](#)*

Writing programs in the 70s

- Writing in a more non-native language (e.g., Fortran, Pascal)
- Grammar is simple and well-defined
- When I learned how to program in college
- I wrote program using a card punching machine
- Computer read the cards
- One day later, an operator handed out the results
 - Syntax errors 401
- I punch a few new cards. And start over again

One more thing

- When people started to use this *new language*
- People were told to write Structural Programs
 - Use strictly 3 structures: *sequence*, *condition*, *repetition*
- Using a good programming language
 - Not *Fortran* (I would put *C*, *C++* here)
 - *Algol*, *Algol 68*, *Pascal*, *Python*
- Go to's are considered harmful (*Dijkstra 1968*)
- Academic analog of programming
 - *Rhetoric structure*
 - *Linking phrases*

Programing in the past is like Writing today

- Students write in a non-native language (e.g., **English**)
- Grammar is complex and ill-defined
- Students write their essays in a word processor
 - Limited interactions (**spelling**, **S-V agreement**)
- Students hand in their assignment to the teacher
- One week later,
 - The teacher hands back graded papers, with
 - **Grade**
 - **Corrective feedback** (there is debate about this)

Writing programs TODAY

- People write programs in **Interactive Development Environment**
- In **IDE**
 - You type a few words (tokens)
 - The system will show you **suggestions**, what lies ahead
 - In other words, IDE help you with **prompts** and **autocomplete**
 - You **test drive** your program and get instant corrective feedback
 - **Syntax errors**
 - **Semantic errors**

My main point

- I will show you
 - **Future of Writing = Present of Programming**
 - **Future = Interactive Writing Environment**
- Why?
 - History repeats itself
 - We have better theory of the language
 - Lexical grammar
 - **Pattern Grammar**
 - We have much more **Data** to derive Pattern Grammar
 - **Statistical methods** are mature

Better Theory of Language

- **Old theory**: Two sides of a coin: **Vocabulary** + **Rules**
 - Words have parts of speech (POS)
 - Noun: object; Verb: action
 - Preposition/order: relation and event
 - Rules are related to POS, and almost independent of words
- **New theory**
 - Lexical grammar (e.g., Pattern Grammar)
 - Rules are intimately linked to words

Pattern Grammar

- PG is a model for describing the syntactic environments of individual words
- Each word has a set of patterns describing word usage in typical contexts
- One sense per pattern (often patterns are different for different word senses)

Sources:

- http://en.wikipedia.org/wiki/Pattern_grammar
- Hunston and Francis (2000): A corpus-driven approach to the lexical grammar of English

Pattern Grammar

Skim (v.) includes the following patterns in the COBUILD dictionary

- **V n off/from n:** *Skim the fat off the soup.* (limited prep. allow)
- **V n:** *Skim the wall surface smooth and ready for painting*
- **V over/across:** *Water skiers skimmed across the bay.*
- **V through n:** *Skim through the report and check for spelling mistakes?*

Sources:

- http://en.wikipedia.org/wiki/Pattern_grammar
- Hunston and Francis (2000): A corpus-driven approach to the lexical grammar of English

Dictionaries Embrace Pattern Grammar

MACMILLAN
DICTIONARY

- **have difficulty with something:** *She's having difficulty with her schoolwork this year.*
- **have difficulty (in) doing something:** *Six months after the accident, he still has difficulty walking.*
- **great/considerable difficulty:** *We had considerable difficulty finding anywhere to park.*
- **do something with/without difficulty:** *Seb was speaking with great difficulty.*

More Data Make It Easier to Derive PG

- **BNC**, *British National Corpus* 100 million words
- **COCA**, *Corpus of Contemporary American English* 450 million words
- **CiteseerX**, *Scholarly Big Data* 460 million words (2.8G)

Table 3: Collection and Usage Statistics

Statistic	Value
#Documents	3.5 million
#Unique documents	2.5 million
#Citations	80 million
#Authors	3-6 million
#docs added monthly	300,000
#docs downloaded monthly	300,000-2.5 million
Individual Users	800,000
Hits per day	2-4 million

Source: Wu, Zhaohui, et al. "Towards building a scholarly big data platform: Challenges, lessons and opportunities." in *Digital Libraries 2014*.

Take a closer look at CiteSeerX (for WriteAway)

```
$ head -3 citeseerx | tail | cut -f1
```

Because of their deployment in critical applications , the dependability modeling and analysis of Multiple-Phased Systems is a task of primary relevance .

```
$ wc citeseerx
```

20,000,000 sentences; 460,000,000 words; 2.8 G bytes

```
$ head -3 citeseerx | tail -1 | cut -f3,4
```

IN IN PRP\$ NN IN JJ NNS , DT NN NN CC NN IN JJ NNS VBZ DT NN IN JJ NN .

(Part of Speech, read: prep, prep, possessive pronoun,)

I-PP H-PP I-NP H-NP H-PP I-NP H-NP O I-NP I-NP H-NP O H-NP H-PP I-NP H-NP H-VP I-NP H-NP H-PP I-NP H-NP O

(Base phrase, read: PP-start, PP -start, NP-start, NP-end, PP, NP-start, NP-end, ...)

PP = Prepositional Phrase, NP = Noun Phrase

How Big is the Scholar Data

- Three scholarly big data sets
 - Google Scholar
 - Microsoft Academic Search
 - CiteSeer (mainly computer science)

Table 1: The Estimated Number of Scholarly Documents on the Web in Different Fields.

Discipline	Size in MAS	Estimated Size	public
Agriculture Science	447,134	1,088,711	12%
Arts & Humanities	1,373,959	5,286,355	24%
Biology	4,135,959	8,019,640	25%
Chemistry	4,428,253	10,704,454	22%
Computer Science	3,555,837	6,912,148	50%
Economics & Business	1,019,038	2,733,855	42%
Engineering	3,683,363	7,947,425	12%
Environmental Sciences	461,653	975,211	29%
Geosciences	1,306,307	2,302,957	35%
Material Science	913,853	3,062,641	12%
Mathematics	1,207,412	2,634,321	27%
Medicine	12,056,840	24,652,433	26%
Physics	5,012,733	13,033,269	35%
Social Science	1,928,477	6,072,285	19%
Multidisciplinary	9,648,534	25,798,026	43%
Total Sum		121,223,731	36,703,036

Increasing Effective Statistical Tools

- Parsing and shallow parsing (based on the old theory)
 - *Genia tagger* (Tsuruoka et al. 2005)
- Collocation Extraction
 - *XTract* (Smadja 1993)
- Extracting Good Dictionary Examples
 - *GDEX* (Kilgariff et al. 2008)
- Extracting Pattern Grammar from big data
 - *WriteAway* (Chang et al. 2015)

Source: Tsuruoka, Yoshimasa, et al. "Developing a robust part-of-speech tagger for biomedical text." *Advances in informatics*. Springer Berlin Heidelberg, 2005. 382-392.

My main point

- I will show you
 - **Future of Writing = Present of Programming**
- Why?
 - History repeats itself
 - We have better theory of the language
 - Lexical grammar
 - Pattern Grammar
 - We have much more data to derive Pattern Grammar

Steve Jobs famously said

iPod + phone + Internet communicator = iPhone



I'll say

Typewriter + Dictionary + Concordance = WriteAway

WriteAway

We have difficulty|

less patterns

more patterns

less examples

more examples

[N] difficulty with something **755**

have difficulty with problems involving **25 5**

difficulty with this approach/method is **82 13**

[N] difficulty doing something **718**

have difficulty using it correctly **3 5**



I'll say

Typewriter + Dictionary + Concordance = WriteAway

WriteAway

Typewriter



We have difficulty|

less patterns

more patterns

less examples

more examples

[N] difficulty with something **755**

have difficulty with problems involving **25 5**

difficulty with this approach/method is **82 13**

[N] difficulty doing something **718**

have difficulty using it correctly **3 5**



I'll say

Typewriter + Dictionary + Concordance = WriteAway

WriteAway

Typewriter



We have difficulty|

Dictionary



less patterns

more patterns

more examples

[N] difficulty with something **755**

have difficulty with problems involving **25 5**

difficulty with this approach/method is **82 13**

[N] difficulty doing something **718**

have difficulty using it correctly **3 5**



I'll say

Typewriter + Dictionary + Concordance = WriteAway

WriteAway

Typewriter



We have difficulty|

Dictionary



less patterns

more patterns

more examples

[N] difficulty with something 755

have difficulty with problems involving 25 5

difficulty with this approach/method is

[N] difficulty doing something 718

have difficulty using it correctly 3 5

Concordance





Concordance for Walden / Thoreau, Henry David

Use the features on this page to analyze and evaluate the text.

[View full text](#)

Words beginning with: [Go](#)

Number of words: [Go](#)

Number of phrases: [Go](#)

Search: [Go](#)

Radius: Sort: ☐ None ☐ Left ☐ Right ☐ Match

Specialized searches: colors; adverbs; gerunds; "big names"; "great ideas"

1. ctualness of the **civilized man**? According to Liebig, man's b
2. ittle, while the **civilized man** hires his commonly because he
3. is tax, the poor **civilized man** secures an abode which is a p
4. commonly a poor **civilized man**, while the savage, who has th
5. tion between the **civilized man** and the savage; and, no doubt
6. ings. And if the **civilized man**'s pursuits are no worthier th
7. contact with the **civilized man**. Yet I have no doubt that tha
8. a blessing. The **civilized man** is a more experienced and wis



WriteAway

This paper

less patterns

more patterns

less examples

more examples

[N] paper does 123527

paper presents/proposes a 35626 11894

paper describes/discusses the 24462 7026

[N] paper does something 84532

paper describes/presents a system 1960 72

paper presents an approach to 1219 113

WriteAway

This paper presents|

less patterns

more patterns

less examples

more examples

[V] present something of something 52075

paper presents the results of a study conducted 630 7

paper presents the design and implementation of a mobile storage system called 259 6

[V] present something for doing something 27089

paper presents a parallel algorithm for solving the region growing problem based 235 7

we present a method for solving the following problem 142 11

This paper presents a method

[less patterns](#) [more patterns](#) [less examples](#) [more examples](#)

[N] method for something 25928

is a multistart type stochastic method for bound constrained global optimization problems 1330 6

by standard methods for one-dimensional systems 765 7

[N] method be 24538

method is applicable 17015 442

[N] method of something 19839

method of analysis/proof is 1325 115

on the method of moments 433 20

[N] method to do something 19659

of heuristic methods to solve hard computational search problems 528 8

presents a single-pass , view-dependent method to solve the general rendering equation 160 8

[N] method for doing something 18495

present a method for solving the following problem 960 11

of numerical methods for solving ordinary differential equations 326 5

Highlights of WriteAway

- Writing suggestions in the form of
 - grammar patterns
 - typical examples
 - Show less/more patterns/examples on demand
- Editing suggestions for the selected word
 - the same adjustable grammar patterns with examples
- Ranking patterns/examples according to context
 - move 'difficulty doing something' to the top when preceded by 'have'
 - Attempt to resolve part of speech ambiguity (report n. vs report v.)
- Adjustable text box
 - You can type the whole paragraph

Advantages of Using WriteAway

- Countability
 - Noticing the singular/plural forms in examples
- Verb tense/form
 - Noticing the verb forms in examples
- Article
 - Noticing the use or lack of articles
- Preposition
 - Noticing n.+prep. or prep.+n. patterns
 - E.g., evaluation on something
- Prep.+Verb Form
 - Noticing grammar patterns
 - E.g., method for doing something v.s. method to do something

What happens after this

- Improve WriteAhead
- Enter a competition to go to Silicon Valley Emersion Program
- Get more government funding (Taiwn MITI)
- Start a KickStarter campaign
- Collaborate Waseda University
- We are going after YOU

If you use Word, not getting
enough help, why not use

WriteAway

to get a glimpse of the
future of writing

If you use Word, not getting
enough help, why not use

<http://writeaway.nlpweb.org>

to get a glimpse of the
future of writing