

*An Introduction to Mathematical Linguistics for  
Historical Text Analysis*



Tokyo Tech



# Catalogue of Linguistics

*Hilofumi Yamamoto, Ph. D.  
Tokyo Institute of Technology*

言語学と日本語教育を研究するための教室

言語と文化

## 山元研究室カタログ

このカタログは、以下の URL もしくは QR コードで入手できます。

<https://cuckoo.js.ila.titech.ac.jp/yamagen/>



# 目次

第 1 章	科学研究費助成金	1
1.1	2010–13 年度	1
1.2	2014–17 年度	1
1.3	2018–21 年度	1
1.4	資料	1
第 2 章	受賞	9
2.1	2015 年度情報処理学会山下記念研究賞	9
2.2	2017 年度じんもんこん 2017 ベストポスター賞を受賞	9
2.3	2018 年度じんもんこん 2018 学生奨励賞を受賞	9
2.4	2017 年度東京工業大学教育賞最優秀賞を受賞	9
第 3 章	サイエンス・カフェ神戸「目で見てわかる歌ことばの姿」	17
3.1	サイエンス・カフェ神戸でのトーク	17
3.2	開催報告	17
第 4 章	ひらめき ときめきサイエンス	19
第 5 章	人文情報学月報	29
5.1	巻頭言	29
5.2	資料	29
第 6 章	大学研究室探検隊	35
6.1	取材	35
6.2	資料	35
第 7 章	JADH: 論文・ポスター	41
第 8 章	情報処理学会じんもんこん: 論文・ポスター・スライド	69

---

8.1	概要	69
8.2	発表年	69
8.3	論文・ポスター・スライド	69
第9章	Language Classes and Book list	75
9.1	辞書・専門書	75
9.2	日本語: 大学院	75
9.3	日本語: 学部	75
9.4	日本文化関連	76
9.5	言語学関連	76
9.6	MOOC TokyoTechX	76
9.7	その他	76

## 第1章

# 科学研究費助成金

### 1.1 2010–13 年度

研究題目「和歌形態素解析用辞書開発のための用語連接規則に関する基礎研究」

### 1.2 2014–17 年度

研究題目「和歌用語シソーラスの開発と用語空間記述に関する基礎研究」

### 1.3 2018–21 年度

研究題目「歌ことばの効果的可視化技術と通時的変化の記述に関する基礎研究」

### 1.4 資料

2018 年度申請書（抜粋）

## 1 作研究目的、研究方法など作

本研究計画調書は「小区分」の審査区分で審査されます。記述に当たっては、「科学研究費助成事業における審査及び評価に関する規程」(公募要領111頁参照)を参考にしてください。

本欄には、本研究の目的と方法などについて、3頁以内で記述してください。

冒頭にその概要を簡潔にまとめて記述し、本文には、(1)本研究の学術的背景、研究課題の核心をなす学術的「問い」、(2)本研究の目的および学術的独自性と創造性、(3)本研究で何をどのように、どこまで明らかにしようとするのか、について具体的かつ明確に記述してください。

本研究を研究分担者とともに行う場合は、研究代表者、研究分担者の具体的な役割を記述してください。

(概要) 10行程度で記述してください。

【背景と問題】本研究は可視化モデルを利用して、古代語の通時的語彙構造の変化を分析するものである。下記モデル(図1,2)は「吉野」と「桜」の関係を数理的手法により可視化し、300年間の比較を行った。2者間の通時的関係の変遷については明らかではあるが、すべての語についても同様に実現するには、1)単語の長さをすべて短い単位で分割したため、語句の比較が明確でない、2)多義語であるはずの語も一義的に分類されている、などの問題が残されている。

【目的】可視化技術を活かした古代語の通時研究はあまり多くはなく、本研究では古代語通時的変遷を効果的に可視化するシステムを構築し、基礎研究を行うことを目的とする。

【どこまで明らかに?】これまで(基盤研究C)のデータでは辞書という静的な方式で蓄積してきたが、本研究では和歌から直接的に1)単語の類似性情報の計算、2)語と語の関係データの生成を行い、これら動的データと静的データとの差分をとり、通時的変遷を可視化する要因を明らかにする。

## (本文) 歌ことばの効果的可視化技術と通時変化の記述に関する基礎研究 山元啓史(東京工業大学)

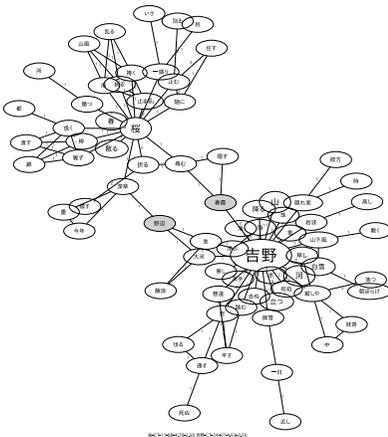


図1: 古今集(ca.905)の「吉野」と「桜」: 古今の時代では吉野は桜の関係よりもむしろ雪と吉野の関係の方が強いことがわかる[20]。

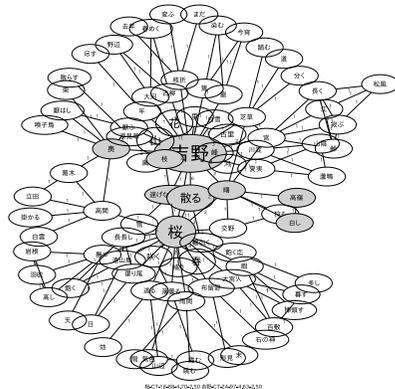


図2: 新古今集(1205)の「吉野」と「桜」: 専門家の間では桜と吉野の関係が一般的になるのは新古今になってからと言われている[20]。

## 【問題点】

現代語の理論的研究、自然言語処理の技術開発研究は、目覚ましく発展してきているが、古代語研究とりわけ通時研究は、コーパスのおかげで数こそ増え、古典資料の整備、人手による作業の多さ、評価の多様さなどの理由により、大規模かつ横断的な調査が実施された事例は多くない。

本研究の目的は、八代集(古今集905年頃から新古今集1205年)までの300年間の勅撰和歌集約9500首を対象に、古代語(和歌用語)を効果的に可視化するシステムの構築と通時的言語変遷の基礎研究を行うことである。これまでに基盤研究(C)により、二十一代集(古今集905年から新続古

## 【1 研究目的、研究方法など(つづき)】

今集1439年)までの20巻を対象に、単語を抽出するためのシステムと辞書、および意味分類辞書(シソーラス)の開発を、機械処理と人手による目視修正を行ってきた。

近年、自然言語処理技術の成果により、人工知能を利用したテキスト処理が盛んに行われるようになり、これらの技術を駆使した応用が多く見られるようになってきた。しかしながら、言語の分析、中でも古代語の研究については、研究者のコンピュータ技術、数理的思考が直接的に人文科学領域になじまず、まだ十分に生かされてきていない。

## 【可視化システム開発の必要性と問題】

上記のモデルは、吉野と桜について、古今集(図1)から新古今集(図2)の300年間の通時的変遷を示したものであるが、このように可視化技術は通時の変化を要約・分析するのに便利である。しかしながら、他の語についても同様に比較・実行するには、

- 1) 単語のサイズを一律的に決めため、当時の単語の成立が明確でない、
- 2) 多義語であるはずの語が一義的に分類されている、などの問題がある。

前者の解決法としては、Sentencepiece(ニューラル言語処理向けのトークナイザ; 教師なし、文脈依存、可逆式)単語分割アルゴリズムの利用である。後者の解決法としては、Word2VecなどのWord Embedding[1,2]、分散表現と呼ばれている文脈から単語の意味特性を計算し、その意味次元をニューラルネットで圧縮(次元数を減らす)し、ベクトルの近さを類義語もしくは同語と判定する方法である。

## 【解決法】

機械学習を利用した言語分析の研究において、従来からも指摘されている問題は、単語の類似性と関連性がうまく区別されていないことである。単語の類似性と関連性というのは、たとえば、(梅, 桜)は類似している一方、(梅花, 枝)は関連してはいるが、類似はしていないということである。これらを区別して処理すれば、人間が行うモデリングに近くなるという報告があることから、和歌の処理においても良い効果が期待できる。

また、分散表現には単語の曖昧性が考慮されていないという問題点がある。単語にはさまざまな意味がある。たとえば、英語の“spring”という語は「温泉」「スプリング」「春」という意味がある。単語の多義性を考慮せず、1つの“spring”という形態を1つのベクトルで表すのには限界がある。むしろ、表記は同じであったとしても、ある文脈に挟まれた語の表記を一旦伏せておき、仮に $x_i$ とし、文脈から得られたベクトルにしたがって、 $x_i$ の分散表現を与える方法を考えれば、異なるベクトルを同じ表記で示す必要はなくなる。これは、文脈の隔たりの大きいベクトルを $x_i$ の添字毎に分割し、多義性のある単語を用法・文脈ごとに記述する方法である。その結果、単語の用法の弁別性能が向上したことが報告されている[3,4]。

## 【目的】

古代語は現代語とは異なり、従来より可視化技術を利用した通時的言語体系の研究は今までに多くなく、限られた資料から、目視によって丹念に分析していくものが大半であった。

一方で、現代語の分析に大いに利用と期待が寄せられている自然言語処理技術は目覚ましい発展があり、人工知能技術、ニューラルネット、ベイズ統計学、時系列分析などの基礎技術と融合し、近年大きな成果を収めている。

古代語は言語資料に限られており、現代語のような新しいデータが次から次と出てくるものではないが、今までの成果を利用し、蓄積を取り込みつつ、総合することで、少なくとも考え方を取り入れることで成果を収められると考えている。学術的創造性として注目するのは、古代語通時研究のための効果的な可視化システムと語彙データベースの開発である。

## 【古代語へのチャレンジ】

上記で述べた自然言語処理技術の応用が解決策として有力ではあるが、和歌(古代語)というテキストの特性として、1)現代語のように大量のデータがあるわけではなく、データ量は限られた上

### 【1 研究目的、研究方法など(つづき)】

で研究を進めなければならないこと、2) 現存するテキストは何らかの理由(希少価値、読み継がれてきたほどの魅力、消失・散逸せずに残存している現状、長年にわたっても理解できる内容)で、テキストの内容、語の意味が限られている可能性はあること、から考えると、現代語でできることと、古代語にできることと隔たりがあることに注意すべきであり、簡単ではないことが予想される。ただし、上記の点が本研究のチャレンジであり、可能性が見えれば、通時的言語研究への貢献は大きいと考える。

### 【ゴール設定：何をどこまで?】

二十一代集のすべてについて行うのではなく、基本的な八代集についてのみを対象とし、これをこの4カ年のゴールとして設定し、着実に成果をあげる計画を実行する。ただし、単語の単位の切り出し推定実験には、できるだけ多くのデータを用いた方が有利なので、万葉集、二十一代集の和歌本文データを利用する。

古代語も現代語と同様にデータは辞書形式(静的)で蓄積されてきたが、本研究では、和歌データのみから動的に1) 単語の類似性情報の計算、2) 語と語の関係データの生成を行い、開発済みの静的データとの差分をとり、動的表示を可能にするための要因を明確にし、可視化を実現する。

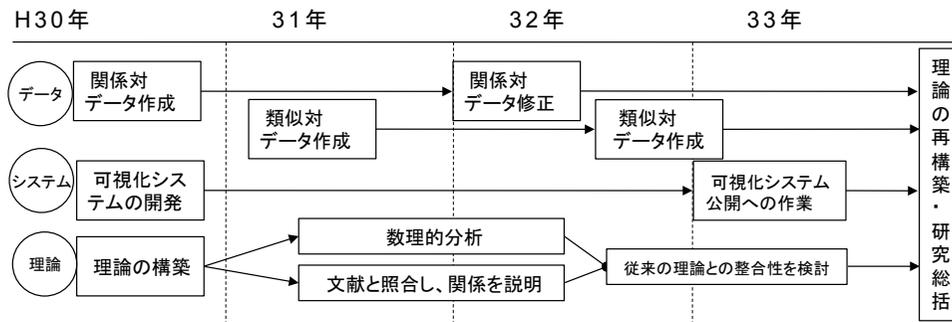


図 3: 研究計画・ロードマップ: データ、システム、理論の3要素で構成する。

### 参考文献

- [1] Le, Quoc V. and Tomas Mikolov (2014) "Distributed Representations of Sentences and Documents," CoRR, Vol. abs/1405.4053, URL: <http://arxiv.org/abs/1405.4053>.
- [2] Tomas Mikolov, Quoc V. Le and Ilya Sutskever (2013) Exploiting Similarities among Languages for Machine Translation, CoRR.
- [3] Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean (2013) "Efficient Estimation of Word Representations in Vector Space," CoRR, URL: <http://arxiv.org/abs/1301.3781>.
- [4] Řehurek, Radim and Petr Sojka (2010) "Software Framework for Topic Modelling with Large Corpora," in Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks, pp. 45-50, Valletta, Malta: ELRA, May, <http://is.muni.cz/publication/884893/en>.

## 2 作本研究の着想に至った経緯など作

本欄には、(1)本研究の着想に至った経緯、(2)関連する国内外の研究動向と本研究の位置づけ、(3)これまでの研究活動、(4)準備状況と実行可能性、について1頁以内で記述してください。  
 「(3)これまでの研究活動」の記述には、研究活動を中断していた期間がある場合にはその説明などを含めても構いません。

### 【1. 着想に至った経緯】

研究代表者がこれまでに15年以上の歳月(2001年より)をかけて作成してきた和歌用形態素解析辞書とシソーラス(語彙体系用語集)を使い、和歌の語彙体系を効果的に可視化するための技術を開発し、さらに通時的言語記述として適切であるかどうかを検証してきた(尚、形態素辞書とシソーラスは、25年度までの基盤研究(C)で二十一代集対応版が完成している)。

### 【2. 国内外の研究動向と位置づけ】

Word2vecを始めとするWord Embedding(分散表現)に関わる研究も、ツールも数多く発表されており、自然言語処理研究においては大いに理論化がされており、技術の応用も多々行われている。ただし、言語学、とりわけ古代語や通時の変化を分析するまでには至っておらず、今までの古代語研究の成果を言語処理に利用する方法論の検討が待たれている。海外の日本語・日本文化の研究者についても、技術の導入は徐々に行われてはいるものの、テキスト処理・機械学習などの基盤となる技術を使った研究成果はほとんど行われていない。

### 【3. これまでの研究活動】

漸近的語彙対応推定法 [4][6]: 単語対相互情報量により推定した語対応の技術を取り入れ、今までの人間によるシソーラス作りの弱点を改善し、シソーラス体系作りの自動化と理論化を試みた。

二十一代集シソーラスの開発 [8]: 表1に示すように思いも寄らぬ表記が多数出現するため、シソーラスを開発した。

表 1: シソーラスなしでは同じ語として計算できない例(一部)

語彙コミュニティの分析 [2]: R の Linkcomm を用い、単語をコミュニティとして、語群としての意味を検討した。

和歌用可視化システムの開発 [7]: D3.js を用いて、グラフ図形の生成とノードをクリックすることで、原典の和歌のリストが閲覧できる可視化システムを開発した。

自動タグ付けシステムの開発 [21]: 当時、和歌用の辞書

がなかったために、単語辞書、接続辞書を開発し、単語に切り分けるシステムを開発した。

かな表記	実際に和歌に出現する実例
たつた	立田, 竜田, 龍田, ...
たつらむ	立つらん, 立らん, 立覧, ...
ちぎりけむ	契りけん, 契けむ, 契けん, 契剣, ...
おもふへふ	思ふてふ, 思てふ, 思ふ蝶, 思蝶, ...
えてしがな	得てしかな, 得てし哉, ...

### 【4. 準備状況と実行可能性】 技術・材料・資料関係

これまでの基盤研究(C)で培ってきた辞書・シソーラス、和歌本文データ、現代語訳データなどの基礎的な材料はすでに整備されている。一連のWord Embeddingの理論と技術は広く公開されており、入手済みであるので、データの的にも技術的には実行可能である。

#### 分担関係

ホドシチェック(阪大)は、機械学習技術を古代語の分析に応用、プログラミング、可視化技術を担当する。山元啓史(東工大)は、関連対・類似対データの開発、通時的分析手法の開発、研究総括を担当する。また、両者は、可視化システムをパイリンガルで表示するために、古代語の日英語対応関係を分析するシステムの開発を行う。

#### うまくいかない時の対応策

これまでの研究実績 [1][2]により機械学習で解析できる可能性はかなり高いが、和歌の根本的な限られたデータ量の都合により、うまくいかない時には、これまでの研究で利用したデータとの融合を考え、その上で、なぜうまく行かないのか、なぜシソーラスと連動させる必要があったのかを考察し、理論的な説明を構築し、研究の貢献とする。

以上を網羅した上で、データ処理による通時的な視点での古代語の空間記述研究が開始でき、この領域への貢献となるだろう。

### 3 作研究代表者および研究分担者の研究業績作

本欄には、研究代表者、研究分担者がこれまでに発表した論文、著書、産業財産権、招待講演のうち重要なものを選定し、現在もしくは過去から発表年次の順に、通し番号を付して2頁以内で記入してください。なお、学術誌へ投稿中の論文を記入する場合は、掲載が決定しているものに限ります。

学術誌論文の場合、論文名、著者名、掲載誌名、査読の有無、巻、最初と最後の頁、発表年(西暦)を記入してください。以上の項目が記入されていれば、各項目の順序の入れ替えや、著者名が多数の場合、主な著者名のみ記入しその他の著者を省略することは問題ありません。なお、省略する場合は、省略した員数と、研究代表者、研究分担者が記載されている順番を○番号と記入してください。

研究代表者には二重下線、研究分担者には一重下線を付してください。

1. H. Yamamoto, and B. Hodošček. Relationships between Flowers in a Word Embedding Space of Classic Japanese Poetry, Doshisha University, JADH2017 Proceedings of the 7th Conference of Japanese Association for Digital Humanities “ Creating Data through Collaboration ”, Faculty of Culture and Information Science, Doshisha University, Vol. 2017, pp. 70-72, (2017) (査読有).
2. H. Yamamoto and B. Hodošček. “Development of the dictionary of poetic Japanese description”, Digital Scholarship in History and the Humanities, the 6th conference of the Japanese Association for Digital Humanities, Japanese Association for Digital Humanities 2016 pp. 44-46, (2016) (査読有).
3. 山元啓史. “通時コーパスによる言語の研究”, コーパスと日本語史研究, ひつじ書房, pp. 17-35, (2015) (査読有).
4. 山元啓史, ホドシチェク・ボル, 村井源, “二十一代集シソーラスのための漸近的語彙対応システムの開発”, じんもんこんシンポジウム 2014, 人文科学とコンピュータシンポジウム論文集, Vol. 2014, No. 3, pp. 157-162, (2014) (査読有).
5. 山元啓史. “目で見てわかる歌ことば”, 日本語学, 明治書院, Vol. 33, no. 14, pp. 172-183, (2014) (査読無).
6. H. Yamamoto, B. Hodošček, and Hajime Murai. “Development of an Asymptotic Word Correspondence System between Classical Japanese Poems and their Modern Translations”, JADH Conference 2014, JADH Conference 2014 ABSTRACT, p.40, (2014) (査読有).
7. H. Yamamoto, B. Hodošček, and Makiro Tanaka, “A Visualization and Analysis System for Japanese Language Change: Quantifying Lexical Change and Variation using the Serial Comparison Model”, JADH Conference 2014, JADH Conference 2014 ABSTRACTS, p. 3, (2014) (査読有).
8. H. Yamamoto, and B. Hodošček. “Thesaurus of classical Japanese poetic vocabulary for the Nijuichidaishu (ca. 905-1439)”, 14th International Conference of European Association for Japanese Studies, 14th International Conference of European Association for Japanese Studies BOOK OF ABSTRACTS, p.86, (2014) (査読有).
9. B. Hodošček and H. Yamamoto, “A Diachronic and Synchronic Investigation into the Properties of Mid-Rank Words in Modern Japanese” The Japanese Association for Digital Humanities, the third annual conference at Ritsumeikan University, Kyoto, Japan, September 19-21, pp. 67-8. (2013) 査読有.
10. H. Yamamoto, “Lexical Modeling of Yamabuki (Japanese Kerria) in Classical Japanese Poetry”, The Japanese Association for Digital Humanities, the third annual conference at Ritsumeikan University, Kyoto, Japan, September 19-21, 62-3, (2013) 査読有.

## 【3 研究代表者および研究分担者の研究業績(つづき)】

11. H. Yamamoto, M. Tanaka, Y. Kondo, “Diachronic Corpus and Linguistic Space: New Methods for the Analysis of Language Change”, SNPD2012, Proceedings 13th ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing, IEEE, Vol. SNPD2012, No.101, 381–384, (2012) 査読有.
12. M. Tanaka, and H. Yamamoto, “Emotive Adjectives and Verbs of the Heian Japanese”, JADH 2012 conference abstracts, Vol. 2012, p. 52, (2012) 査読有.
13. H. Yamamoto, M. Tanaka, and Y. Kondo, “Design of Serial Comparison Model for the Diachronic Corpus Study of Japanese”, JADH 2012 conference abstracts, Vol. 2012, 51–2, (2012) 査読有.
14. 山元啓史. “グラフを用いた集合演算による和歌用語の解析”, 語彙研究, 語彙研究会, Vol. 9, 86–94, (2011) 査読有.
15. H. Yamamoto, and M. Tanaka, “Quantitative Analysis of Loanwords of Eight Literary Works in the Heian Period (794–1185)”, Osaka simposium on digital humanities 2011, Vol. 1, No. 1, 51–2, (2011) 査読有.
16. H. Yamamoto, “Graph Representation of the Connotations of Classical Japanese Poetic Vocabulary”, Osaka simposium on digital humanities 2011, Vol. 1, No. 1, p. 42, (2011) 査読有.
17. M. Tanaka, and H. Yamamoto, “An analysis of Sino-Japanese words of the Heian period for the development of the historical Japanese dictionary”, Asialex 2011, Lexicography: Theoretical and Practical Perspectives, 496–505, (2011) 査読有.
18. H. Yamamoto, and M. Tanaka, “Development of the thesaurus of classical Japanese poetic vocabulary”, Asialex 2011, Lexicography: Theoretical and Practical Perspectives, Vol. 2011, 576–585, (2011) 査読有.
19. 山元啓史, “「山吹」をめぐる和歌語彙の空間”, じんもんこんシンポジウム 2011, 人文科学とコンピュータシンポジウム論文集, 情報処理学会, Vol. 2011, No. 8, 141–146, (2011) 査読有.
20. 山元啓史, “八代集用語のモデリングシステム”, じんもんこんシンポジウム 2010, 人文科学とコンピュータシンポジウム論文集, 情報処理学会, Vol. 2010, No. 15, 247–254, (2010) 査読有.
21. 山元啓史, “分類コードつき八代集用語のソーラス”, 日本語の研究, 日本語学会, Vol. 5, No. 1, 46–52, (2009) 査読有.



## 第 2 章

# 受賞

### 2.1 2015 年度情報処理学会山下記念研究賞

2015 年度情報処理学会山下記念研究賞を受賞。山下記念研究賞は、情報処理学会が主催する研究会およびシンポジウムにおける研究発表のうち、特に優秀な論文の発表者に授与される賞。初代情報処理学会会長の故山下英男氏寄贈の資金にて運営されている。

### 2.2 2017 年度じんもんこん 2017 ベストポスター賞を受賞

12 月 9・10 日に大阪市立大学杉本キャンパスで開かれた人文科学とコンピュータシンポジウム「じんもんこん 2017」にて、リベラルアーツ研究教育院の山元啓史教授らの研究が、ベストポスター賞を受賞した。

### 2.3 2018 年度じんもんこん 2018 学生奨励賞を受賞

2018 年 12 月 1、2 日に東京大学で開催された、情報処理学会じんもんこん 2018 にて、社会理工学研究科価値システム専攻博士課程の平野充さんとリベラルアーツ研究教育院山元啓史教授による論文「低頻度ピッチクラスセットの 2-gram パターンを用いたモーツァルトの交響曲と弦楽四重奏曲の比較分析」が学生奨励賞を受賞した。

### 2.4 2017 年度東京工業大学教育賞最優秀賞を受賞

2019 年 2 月 1 日、「日本語・日本文化科目群運営のための統合的機能を有する教務システム (JCOS) の開発」にて、東工大教育賞最優秀賞を受賞した。





## リベラルアーツ研究教育院の山元啓史教授らの研究 が、じんもんこん2017ベストポスター賞を受賞

『歌ことば「橘」「梅」「桜」における関連対の抽出』

RSS

Like 2

ツイート

2017.12.28

12月9・10日に大阪市立大学杉本キャンパスで開かれた人文科学とコンピュータシンポジウム「じんもんこん2017」（主催・一般社団法人情報処理学会人文科学とコンピュータ研究会）において、リベラルアーツ研究教育院の山元啓史教授らの研究が、ベストポスター賞を受賞しました。

### 歌ことば「橘」「梅」「桜」における関連対の抽出





ホドシチエック ボル  
大阪大学  
boroslang.osaka-u.ac.jp

山元啓史  
東京工業大学  
yamagen@ila.titech.ac.jp

歌ことば辞典の開発にあたり、従来の「見出し語とその解説」の記述に加え、「見出し語—関連語」の形式による記述を提案し、「橘」「梅」「桜」の関連対の抽出を行った。

表1: 橘のサブクラスタ: 1番から10番までを抽出。average法、mcquitty法、single法のクラスタリングを用いた。括弧内はそれぞれのMaximum partition densityを示す。

No.	average (.43)		mcquitty (.43)		single (.38)	
	node	edge	node	edge	node	edge
1	昔	7	昔	7	昔	5
2	匂ふ	6	匂ふ	6	匂ふ	4
3	風	5	今年	4	夢	4
4	夢	5	辺	4	香る	3
5	今年	4	待つ	4	今年	3
6	辺	4	風	4	初む	3
7	待つ	4	夢	4	五月雨	3
8	香る	3	初む	3	折	3
9	五月雨	3	香	3	枕	3
10	初む	3	間	3	思ひ寝	3

ポストカードではさらに「梅」「桜」の結果...



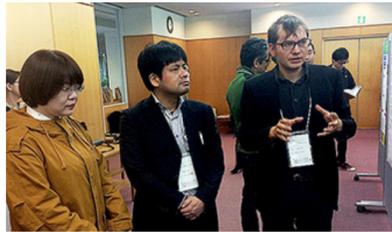

ライトニングトークスライド

受賞対象となったのは『歌ことば「橘」「梅」「桜」における関連対の抽出』と題する発表で、平安時代以来日本の伝統的な和語と和歌に用いられる歌ことばを演劇で言う主役と脇役をそれぞれ主役語と脇役語と位置づけ、従来の辞書では、語釈と和歌で構成されていた辞書記述に、脇役語を加え、当時の主役語「橘」「梅」「桜」の使用法や意味範囲を加えると用語の使われ方がわかりやすくなるというものです。

これを現代人の先入観からではなく、コンピュータによる機械的な技法によって客観的に和歌の中から注目すべき脇役語を選び出すことに成功しました。たとえば、主役語「橘」（一般的に夏の花と考えられている）は「昔」「香」「匂ふ」「枕」「夢」などの脇役語が現れ、それらを繋ぐと「昔の人の香りが匂い、夢に僕がしさを思い出される」が想像できます。この関係は、伊勢物語の第六十段でも有名な詠み人知らずの歌「五月待つ、花たち花の、香をかげば、昔の人の袖の、香ぞする」に近い関係であることがわかります。

この方法の原理は実は言語学によるものというよりも、コミュニティ分析（人と人がどのような関係においてなりたっているのかを分析する科学）と呼ばれる統計的方法論によります。たとえば、主役級の俳優ブラッド・ピット、トム・クルーズ、ジョニー・デップらがほぼ共演することがない一方で、脇役級の俳優はいろいろや主役級の俳優と共演し、さまざまな役割を演じています。単語も強い印象的な意味がある語もあれば、さほど印象的ではないが、さまざまな語と結びついて、一緒にメッセージを形成する語もあります。

このような語は一般的には多義語と言われていますが、多義性は他の語と結びついてその役割がはっきりすることがわかりました。



研究説明するホドシチェク講師

山元教授らの研究グループは、この方法の開発を通じて、言語における多義性の分析を歴コンピュータ記述による通時言語学という観点から構築しようとしています。

受賞対象となった研究は、JSPS科研費 基盤研究 (C) 「和歌用語シソーラスの開発と用語空間記述に関する基礎研究」によるもので、この内容はJSPS「ひらめき☆ときめきサイエンス: 目で見てわかる今の日本語、昔の日本語」でも中学生にもわかりやすく説明されています。



賞状授与式



受賞スピーチ

共同研究者のホドシチェク講師（大阪大学）は東京工業大学で工学博士を得て、現在大阪大学言語文化研究科で研究されています。今回の受賞対象のポスターはホドシチェク講師が、言語のもつ多義性をわかりやすく伝えるために創意工夫した点が評価されました。

📄 [じんもんこん2017サイト](#)

📄 [山元啓史研究室](#)



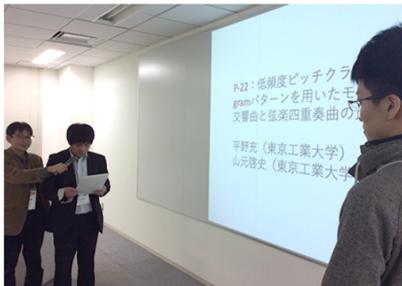
## 平野充さん（社会理工学研究科価値システム専攻・山元啓史研究室）が情報処理学会じんもんこん2018で学生奨励賞を受賞

RSS Like 0

ツイート

2018.12.21

2018年12月1、2日に東京大学で開催された、情報処理学会じんもんこん（人文科学とコンピュータシンポジウム）2018において、社会理工学研究科価値システム専攻博士課程の平野充さんとリベラルアーツ研究教育院山元啓史教授による論文「低頻度ピッチクラスセットの2-gramパターンを用いたモーツァルトの交響曲と弦楽四重奏曲の比較分析」が学生奨励賞を受賞しました。平野さんは山元教授の指導の下、修士課程・博士課程を通して5年の月日をかけ、モーツァルトの交響曲と弦楽四重奏曲をデータ化し、ピッチクラスセットを使って計量分析を行ったところ、同じ楽曲であっても、両者のピッチクラス・パターンに微妙な差があることを発見しました。



授賞式の模様

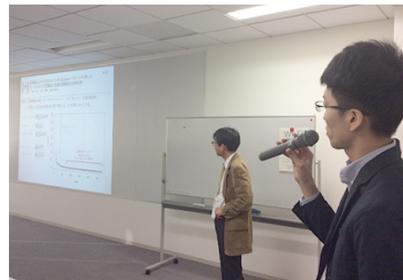


受賞後のインタビュー

モーツァルトは、すぐれた業績を残した偉大な作曲家です。その作曲家の同じ曲が、交響曲と弦楽四重奏曲の両方の形式で演奏されています。本当はモーツァルトは交響曲として作曲したのでしょうか、弦楽四重奏として作曲したのでしょうか。いろいろな解釈はあるかもしれませんが、実のところ、よくわかっていません。また、交響曲・弦楽四重奏の両者は質的にどう違うのでしょうか。プロの演奏家に問い合わせても、楽器編成が違うこと以上のことはよくわかっていません。

平野さんは、主観的な分析を一切交えず、楽譜から得られたデータにより、よく使われるピッチクラスセットのパターン、稀にしか使われないパターンに分け、よく使われるピッチクラスパターンを浮き立たせ、遷移パターンが交響曲と弦楽四重奏の間で有意に異なることを発見しました。この研究成果が同シンポジウムで認められ、今回の受賞となりました。

また、平野さんは、この研究をきっかけにもっと多くの音楽計量分析を行う科学者が今後増えることを期待し、学位取得後も音楽の計量分析を通して、楽曲の研究だけでなく、聞き方、楽しみ方を伝えていきたいと考えています。



平野さんによる受賞論文のライトニングトーク



## 「より優れた教育の推進に」平成29年度東工大教育賞授与式を実施

Like 105

Share

ツイート

受賞・表彰

教育

RSS

2019.04.03

2月1日、大岡山キャンパス本館で平成29年度東工大教育賞授与式が行われました。

この賞は、教員の教育方法及び教育技術等の向上を図り、より優れた教育を推進することを目的として制定されたもので、今回で16回目となります。

授与式では、最優秀賞に選ばれたリベラルアーツ研究教育院 山元啓史教授ら受賞者に対して、益一哉学長から賞状及び報奨金(目録)が授与されました。



あいさつする益学長



受賞者を代表してあいさつする山元教授

### 平成29年度東工大教育賞受賞者一覧

教育に関して優れた業績を挙げたとして、次の58名(10件)が選ばれました。

(所属は受賞当時、所属順・敬称略)

(代表者の所属順)

#### 最優秀賞

受賞者(所属)	対象業績
山元啓史教授(代表者・リベラルアーツ研究教育院)、佐藤礼子准教授(リベラルアーツ研究教育院)、平川八尋准教授(リベラルアーツ研究教育院)、森田淳子准教授(リベラルアーツ研究教育院)	「日本語・日本文化科目群運営のための統合的機能を有する教務システム(JCOS)の開発」



受賞者の記念撮影

## 第3章

# サイエンス・カフェ神戸 「目で見てわかる歌ことばの姿」

### 3.1 サイエンス・カフェ神戸でのトーク

サイエンスカフェ (Science Café) は、1997年から1998年にかけて、イギリスとフランスで同時発生的に行われたのが起源とされる、カフェのような雰囲気の中で科学を語り合う場、もしくはその場を提供する団体の名前である。英国での呼称に倣ってカフェ・シアンティフィック (Café Scientifique) と呼ぶこともある。

[サイエンスカフェ—Wikipedia](#)

2008年3月24日、神戸大学主催の[ようこそサイエンスカフェ神戸へ](#)で発表した。

### 3.2 開催報告

次ページ pdf。



## 第4章

# ひらめき　ときめきサイエンス

略称「ひらとき」と呼ばれるセミナーは、大学や研究機関で「科研費」(KAKENHI)により行われている最先端の研究成果に、小学5・6年生、中学生、高校生が、直に見る、聞く、触れることで、科学のおもしろさを提供する科学日本学術振興会からの委託事業である。研究費助成金と同じ取り扱いである。

整理番号	HT30068	分野	その他・人文	キーワード	言語学、情報科学
------	---------	----	--------	-------	----------

研究機関名	東京工業大学				
プログラム名	目で見てわかる昔の日本語と今の日本語：タイムマシンに乗らずに行ける昔の世界				
先生（代表者）	山元 啓史（やまもと ひろふみ） リベラルアーツ研究教育院・教授				
自己紹介	大学卒業から今までずっと外国人に対して日本語を教えてきました。教えているうちに「ことばはどんな形をしているのだろう」と思い、言語学を勉強しはじめました。いろいろなことばを知りたくなって、アメリカとオーストラリアに留学しました。世界のことばと日本のことばを比較したり、昔のことばがどうして今の形になったかを調べたりするようになりまして。ことばは誰もが使える楽しい宝物のように思います。さあ、みんなでことばについてお話ししましょう。				
開催日時・募集対象	①平成30年8月1日（水） ②平成30年12月26日（水）	受講対象者	①中学生 ②小学6年生	募集人数	①20名 ②10名
集合場所・時間	東京工業大学西1号館1階ラウンジ		（集合時間）	9:50	
開催会場	東京工業大学大岡山キャンパス 住所：〒152-8550 東京都目黒区大岡山 2-12-1 W1-8 アクセスマップURL： <a href="http://www.titech.ac.jp/maps/index.html">http://www.titech.ac.jp/maps/index.html</a>				
<b>内 容</b>					
ことばは時代につれて変化します。私たちの知っていることばの意味は今の意味で、昔の意味とはまったく同じではありません。もしタイムマシンに乗って昔の日本語が聞けたなら「あれえ～何か変だ！違うぞ？」と思うことでしょう。大昔の録音は残っていませんから、実際に聞くことはできません。しかし、昔の文章からことばの使われ方を図に描いて見ることはできます。そんな目で見てわかる昔のことばの世界についてお話しします。					
<b>スケジュール</b>				<b>持ち物</b>	
09:50～10:00 受付（大岡山キャンパス西1号館1階ラウンジ集合） 10:00～10:20 開講式：あいさつ、科研費の説明 10:20～11:00 自己紹介：参加者、ご父兄の皆様、山元研究室学生（終了後15分休憩） 11:15～12:00 講義：ことばの意味を図で見る仕組み 12:00～13:00 ランチタイム：サンドイッチ、みんなでおしゃべりしながら、楽しく食べましょう。 13:00～14:00 実習：コンピュータで自分のネットワークを描こう。 14:00～15:00 休憩：クッキータイム 15:00～16:00 お散歩：鳥人間コンテストのマイスターを訪問しよう！ 16:00～16:40 発表会：みんなで意見と感想を述べよう！ 16:40～17:00 修了式：アンケート記入、未来博士号授与、写真撮影 17:00 終了・解散：お疲れさま。 (8月、12月とも同じスケジュールで行います)				筆記用具	
				<b>特記事項</b>	
				中学生：夏ですので、飲料水、汗を吸いやすいハンカチなど忘れずに。 小学生：冬ですので、暖かい服装で。	

## 《お問合せ・お申込先》

所属・氏名：	リベラルアーツ研究教育院・山元啓史
住所：	東京都目黒区大岡山 2-12-1 東京工業大学 W1-8
TEL 番号：	03-5734-2324
FAX 番号：	03-5734-2324
E-mail：	yamagen@ila.titech.ac.jp
申込締切日：	中学生：平成 30 年 7 月 18 日（水） 小学生：平成 30 年 12 月 5 日（水）

※セミナーは参加者が話し合いをしながら、進めていきます。**おしゃべり**が大好きな皆さん、お待ちしております。ぜひご応募ください。

※当プログラムは先着順にて受付を行います。

※毎年かなり多数のご応募がございます。締切日前に締め切りとさせていただくこともございます。お早めにお申込みください。お申込みの際には、ぜひ日本学術振興会申し込みフォームのコメント欄に「**参加の動機**」をお書きください。応募者多数の場合は、**その文面で選考**させていただきます。楽しいコメントをお待ちしております。選考結果は、中学校（8月）・小学校（12月）のそれぞれ開催 2 週間前に電子メールにてご連絡いたします。あらかじめご了承ください。

## 《プログラムと関係する先生（代表者）の科研費》

研究代表者	研究期間	研究種目	課題番号	研究課題名
山元啓史	H30-H33	基盤研究(C)	18K00528	歌ことばの効果的可視化技術と通時的言語変化記述に関する基礎研究
山元啓史	H26-H29	基盤研究(C)	26370530	和歌用語シソーラスの開発と用語空間記述に関する基礎研究
山元啓史	H22-H24	基盤研究(C)	22520458	和歌形態素解析用辞書開発のための用語接続規則に関する基礎研究



★この科研費について、さらに詳しく知りたい方は、下記をクリック！

<http://kaken.nii.ac.jp/>

※国立情報学研究所の科研費データベースへリンクします。

平成29年度  
ひらめき☆ときめきサイエンス～ようこそ大学の研究室へ～KAKENHI  
(研究成果の社会還元・普及事業)  
実施報告書

HT29084

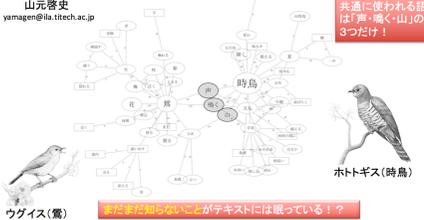
目で見えてわかる昔の日本語と今の日本語：タイムマシンに乗らずに行ける昔の世界

## 言語学とコンピュータ 山元啓史研究室



テキスト処理で言語のありさま、調べましょう！

ともに春の二鳥(「ウグイス」「ホトギス」)だけど...  
古今和歌集(905年頃)では、言葉の使われ方がグラフで描いてみると、  
こんなに違う！

山元啓史  
yamagen@ila.titech.ac.jp

ウグイス(鶯)

ホトギス(時鳥)

昔と今では使われ方が、がテキストには眠っている！

開催日：平成29年8月2日(水)

実施機関：東京工業大学  
(実施場所) (大岡山キャンパス)実施代表者：山元 啓史  
(所属・職名) (リベラルアーツ研究教育院・准教授)

受講生：中学生21名

関連URL：<https://cuckoo.js.ila.titech.ac.jp/~yamagen/hirameki2017.html>

## 【実施内容】

■受講生に分かりやすく研究成果を伝えるために以下のような工夫をしました。

- 全員(参加者、協力者、保護者、事務局)で自己紹介をし、互いに話しやすい雰囲気を作りました。
- 聴講するのではなく、保護者も含め、3～4名のグループに別れ、ディスカッションを進めました。
- ワークブックを作成し、参加者全員が自ら考え、自分で研究の要所が書き込めるようにしました。



- ワークブックをすべて書き込んだら自由研究レポートができあがるようにしました。
- 大学の研究も中学校の勉強と関係づけながら、ディスカッションを進めました。
- 和歌については中学の国語便覧を用い、具体的にページ数を示し、後日復習できるようにしました。
- 言語学でも数学を使うことを示し、関数電卓を用いて、単語の重み計算を実習しました。



- 研究内容だけでなく、言語学の基礎(世界に言語はいくつあるか)、数学の基礎(対数とは、心理尺度とは)、研究の基礎(「特徴とは何か」、「似ている」と「同じ」「違う」とは)など皆で考えました。
- 散歩の時間をとり、学内の建物、ものづくりセンターを見学し、鳥人間コンテストで有名なサークルの協力を得て、人力飛行機、ものづくりの実物に触れることができました。



- 保護者の皆様のお席もご用意し、参加者と同じワークブックを使って討論に参加いただきました。
- 参加者が考えている間、保護者の皆様には研究内容や大学で行われている教育の紹介を行いました。
- 復習できるように、研究室のウェブの特設ページに、当日の記録と写真を掲載しました。



#### ■当日のスケジュール

- 09:50～10:00 受付（大岡山キャンパス西1号館1階ラウンジ）
- 10:00～10:15 開講式：あいさつ、科研費の説明
- 10:20～11:00 自己紹介：参加者、保護者の皆様、研究室学生、研究企画課職員
- 11:00～12:00 講義：ことばの意味を目で見る仕組みとは何か。
- 12:00～13:20 ランチタイム（サンドイッチを食べました）
- 13:20～14:30 実習：コンピュータで自分のネットワークを描こう。
- 14:30～15:00 休憩：クッキータイム
- 15:00～16:00 お散歩：鳥人間コンテストのマイスターを訪問しよう
- 16:00～16:30 発表会：みんなで意見と感想を述べよう！
- 16:30～17:00 修了式：アンケート記入、未来博士号授与、写真撮影
- 17:00 終了・解散

#### ■実施の様子

実施はスケジュールのとおりですが、学術的には以下の内容を盛り込みました。

- ①言語学概論：世界の言語、日本語と外国語、昔の日本語と今の日本語の違い。
- ②計量言語学：頻度とは何か、文書頻度という考え方、重み付けとは何か。
- ③数学と言語：言語を数理的に捉える、数学の成果を言語学に利用する。
- ④研究方法論：仮説をもとに方法を考え、考察をまとめ、結論を導き出す。
- ⑤可視化技法：グラフ理論とグラフ記述言語を学び、モデルを作って目に見える状態を作り出す。



### ■事務局との協力体制

研究推進部研究企画課と事前に打ち合わせを行い、プログラム実施にあたって必要となる準備を確認して下さったほか、当日は事務担当者として研究企画課の2名が参加し、配布物の袋詰作業等の事前準備および受付・写真撮影等を担当していただきました。



### ■広報活動

東京工業大学のウェブサイトにて「東工大の夏休みイベント 2017」カレンダーに実施プログラムの情報を掲載したほか、リベラルアーツ研究教育院のWebサイトにも実施の告知を行いました。

<http://www.titech.ac.jp/outreach/community/summer2017.html>

[http://educ.titech.ac.jp/ila/event\\_information/2017/054042.html](http://educ.titech.ac.jp/ila/event_information/2017/054042.html)



### ■安全配慮

保険に加入し、それを参加者に周知しました。昼食は夏場であることを考慮し、温度による賞味変化の少ないもの、中学生の分量として適切なものを選び、食物アレルギーが起こらぬよう、成分表示を行いました。水分補給には注意を促し、自分で飲み物を持参するようお願いしました。

### ■今後の発展性、課題

ネットワークモデルを作る実習の他に、簡単なプログラムを書いて動かしてみる、小さいレポートを書いてみるなど、参加者同士のディスカッションを今回よりも多く活発にできればと考えています。

1回あたりの受講人数は限度がありますが、毎年、抽選に漏れる方が多いので、チャンス拡大のために回数を夏冬2回にすることを考えています。また、受講者の幅を広げるために、中学生だけでなく、小学生の部も検討に入れています。



【実施分担者】 該当なし

【実施協力者】 6名

【事務担当者】 田中 愛彩美・齋藤 順子 研究推進部研究企画課・事務職員

## ひらめきときめきサイエンス「目で見てわかる昔の日本語、今の日本語」開催

いいね! 24

ツイート 6

[社会連携](#)
[開催報告](#)
[RSS](#)

2015.09.17

8月5日朝10時から夕方5時まで、中学生向けの言語学セミナー「目で見てわかる昔の日本語と今の日本語：タイムマシンに乗らずに行ける昔の世界」が開催されました。このプログラムは、独立行政法人日本学術振興会による「ひらめき☆ときめきサイエンス」事業の支援を受け、実施されました。

### はじめに

会場となった東京工業大学大岡山キャンパスの西1号館に、17名の中学生とその保護者の方々が集まりました。

「ひらめき☆ときめきサイエンス」は、科学研究費による研究成果を、社会に還元・普及するための事業です。小・中・高校生に研究成果をわかりやすく伝える体験プログラムを募集・支援しています。

「ひらめき☆ときめきサイエンス」に採択されたプログラムの多くは、理系の研究テーマを取り扱っていますが、本セミナーは歴史言語学がテーマです。理系の学問は、難しさはあるものの、実験を通してその変化を見ることがごく普通ですが、人文系の学問はとかく概念的で、どこで何が行われたのかわかりにくいものです。参加者の多くが「いったい何が行われるのかわからないが、とにかく参加してみた」と述べていました。言語の移り変わりを参加者がいかに伝えるかがポイントとなりました。



計量言語学っていったい何？「うぐいす」と「ほととぎす」どう違う？

### セミナーを楽しくわかりやすくするために

セミナーに先立ち全員（参加者、協力者、保護者）で自己紹介をし、互いに話しやすい雰囲気を作られました。また、受け身で聴講するだけにならないよう、3～4名のグループでディスカッションをしながら、セミナーは進みました。

さらに、参加者自身が考え、一つ一つ自分で研究の要所を書き込めるように、専用のワークブックが配布されました。このワークブックに全て書き込むと、自由研究のレポートができあがっているようになっています。ま

た、大学の研究は決して中学校の勉強とかけ離れているわけではないので、中学校の勉強と関係づけながら説明が進んでいきました。

### 大学の勉強って何？

はじめに「今日は1日、大学生になったつもりで勉強しよう」と呼びかけられた中学生たち。「大学の勉強には答えがあるかどうかわからない、いや、むしろ答えよりも問題を作るのが大学の勉強だ」との説明に、一足早く大学生の気分を味わいました。研究内容だけでなく、言語学の基礎（世界に言語はいくつあるか）、数学の基礎（対数とは、感覚尺度とは）、研究の基礎（「特徴とは何か」、「似ている」と「同じ」、「違う」とは）など、さまざまな問いをグループで話し合い、ワークブックに鉛筆を走らせました。

### 本セミナーの題材と方法

セミナーのテーマは、平安時代の言語の意味が現在の意味とどう違っているのか、それを可視化を通して見てみようというものです。分析対象は平安時代の古今和歌集です。和歌については中学の国語資料集が用いられました。具体的にページ数を示し、学校や自宅で復習できるよう、工夫されました。また、言語学でも数学を使うことを示し、関数電卓を用いて、単語の重み計算を実習しました。参加した中学生たちは、国語の本を開いたり、数学の対数を勉強したりと、学校の勉強ではあまり経験したことのない文理融合型の学習を経験しました。



はじめて使う関数電卓。なぜ和歌の研究で電卓を？

### もっともっと大学のキャンパスを見て歩こう！

セミナーの合間をぬって、散歩の時間を設け、学内を見学しました。鳥人間コンテストや、ロボットコンテストで有名なサークルの協力を得て、人力飛行機、ロボットに触れることができました。人力飛行機の翼とプロペラを抱えて「わぁ！こんなに軽いついてびっくり」との感想がたくさん寄せられました。



人力飛行機の実物に触れる

### 保護者もディスカッション

保護者の方々にも座席を用意し、参加者と同じワークブックを配布し、ご見学いただきました。ワークブックに沿って、自主的に保護者同士でディスカッションをしてくださっていました。また中学生の参加者が考えている間、保護者の方々には研究内容や大学で行われている活動の紹介が行なわれました。

## おわりに

1日の終わりには未来博士号の授与とアンケートを実施し、終了しました。アンケートでは「学校の授業とは異なり、考える作業で頭をフル回転できた」「東工大で、なぜ和歌を？と思って参加したが、その意味がわかり、目からウロコ」などのご意見をいただきました。



未来博士号（言語学）の授与式



みんなで記念写真

なお、本セミナーの内容を復習できるよう、山元研究室のwebサイトに、当日の記録とワークブックのpdfが掲載されています。

ひらめき☆ときめきサイエンス「目で見てわかる昔の日本語・今の日本語」第1回（2015.08.05）情報

お問い合わせ先  
 山元啓史  
 Email : [yamagen@ryu.titech.ac.jp](mailto:yamagen@ryu.titech.ac.jp)  
 Tel : 03-5734-2324



## 第5章

# 人文情報学月報

### 5.1 巻頭言

巻頭言なるものをはじめて依頼され、執筆した。

### 5.2 資料

## 目次

## 【前編】

《巻頭言》「言語学とコンピュータ」(山元啓史：東京工業大学)

【人文情報学 / Digital Humanities に関する様々な話題をお届けします。】

《巻頭言》「言語学とコンピュータ」  
(山元啓史：東京工業大学)

特別なことがない限り、論文以外の文章は書かないことにしています。もちろん、巻頭言を書くのははじめてです。論文ではないことをいいことに、今までボツになった本の内容について書くことにしましょう。

今まで本を出版しようとして、ボツになった企画が 2 つあります。1 つはコーパス言語学の入門書シリーズの 1 冊で、これは依頼された原稿でしたが、ボツになりました。もう 1 つは東工大の学生のために書いた言語学の教科書でした。

コーパス言語学の本は概論的なものを依頼されました。それを私はコーパス言語学概論と勝手に勘違いして、書き進めていきました。編集の方からは読者は Windows を使っていることを前提に、との注文はありましたが、私自身 Windows を使わないこともあって、Linux のコマンドやパイプ、簡単なスクリプトを中心に説明したテキスト処理の原稿を書きました。Windows であっても、cygwin を使えば大差ないと思ったからです。しかし、Windows 前提でないと本は売れないとのことで NG でした。「ディレクトリとは」「ファイルとは」「OS とは」などのコンピュータの基本用語を説明するように、と書き直しを告げられました。それらを説明した本はたくさんあるので、私自らがボツにしました。GUI のメニュー表示や用語が変わることはあっても、UNIX 由来のコマンドはずっと変わらないし、何をしているのかが、わかりやすいので、その方が息の長い記事になると思ったのですが、編集者さんはそうは思わなかったようです。

とにかく、テキスト処理は、手を動かさないことには、何も始まらないので、その本には次のような例題と練習を載せました。

1. 例題：文の長さのデータの平均値を求めよ。
2. 例題：任意の用語の文脈がわかるようにリストを作れ。
3. 例題：前後の文脈がわかるように文字順に並べ替えよ。
4. 例題：形態素解析器をインストールして、使ってみよ。
5. 例題：形態素解析器を使って名詞だけを選び出せ。
6. 例題：単語の頻度を計算せよ。

コマンドの基本的な原理を説明した上で、どのコマンドを使い、どのプログラム

を組み合わせれば、自分の意図する出力が得られるか、考えてもらう演習です。これの行き着く先は、いわゆるシェル芸というものです。シェル、キーボード・ショートカット、コマンドの組み合わせで、縦横無尽にテキストを料理するってやつです。誰もが同じことを考えるもので、この本を書いた後に、「言語処理 100 ノック」( <http://www.cl.ecei.tohoku.ac.jp/nlp100/> ) というものがあるのを知りました。私の方向性は間違っていないことはわかりましたが、同時にいまさら私が書く必要もないなあとも思いました。

さて、もう 1 つのボツになった本は、「みんなで考える言語学」と題する教科書です。どうせ出版されないのだから「言語学の素」という調味料に似た題名をつけたこともありました。この本は東工大の大学院生に向けた授業が元になっています。ある出版社の担当者さんが「下書きでも良いので内容を見せてほしい」というので、お見せしたところ「オーソドックスな言語学でない」との返答でやんわり断られました。日本語教育能力試験などの検定試験対策になりそうなものを期待したのかもかもしれません。

- ・ 1 章「ワイトゲンシュタインと言語ゲーム」
- ・ 2 章「チューリングとチューリングマシン」
- ・ 3 章「ジップとジップの法則」
- ・ 4 章「ダニエル・ジョーンズの 18 の基本母音」
- ・ 5 章「ソシユールと記号論」
- ・ 6 章「フィルモアと格文法」
- ・ 7 章「チョムスキーと生成文法」

確かに「オーソドックス」ではありません。ワイトゲンシュタインからはじまる言語学の教科書なんてありません。ワイトゲンシュタインは哲学者。チューリングは数学者。ジップでやっとトークンを取り扱うので言語学かな？とも。ダニエル・ジョーンズ(マイフェアレディのヒギンズ博士のモデル)が出てきたあたりから、言語学のように見えます。音韻論を教えるのにダニエル・ジョーンズを出す教科書はほとんどないでしょう。たとえば、かの有名な George Yule の *The Study of Language* の索引でも、“ Jones, Daniel ” の索引項目は見られません。おおむねアメリカの大学の教科書は版を重ねて、演習問題をどんどん新しくしていきます。演習問題はさまざまな観点から入れ替えられます。もっと勉強したい人のための *Further Readings* のリスト差し替えも頻繁です。どんどん版を重ねるので、古い版は面白いくらい安く入手できます。この本の第 3 版は 380 円(新品)でアマゾンから購入できます。

なぜオーソドックスでない構成になったのか？これにはいろいろな理由がありますが、一番の理由は、対象が東工大の大学院生だったということです。数理、計算、物理、化学などの専門家ではあっても、言語学は決して彼らの専門ではありませんし、彼らも言語学を自分の専門として勉強しようとは思っていません。こういう学生に「そもそも言語学とは」などと紋切り型で授業をはじめても眠くなるばかりです。言語学の知識はなくても、授業初日から、ディスカッションがしたくなるような授業を考えました。自分が話すことばと比べながら、「言語学の歩み」を教師が語るのではなく、ディスカッションによって学生さん自身に考えてもらう授業にしました。

どの章にも簡単な紹介・導入を記載しましたが、それ以外は「演習問題」です。これを 3、4 名のグループで「ああだ」「こうだ」とディスカッションしては、それをグループごとに発表していきます。

たとえば、1章の演習問題（言語ゲーム）は、

「私が通りかかったとき、すでにゲームは進行中だった」の「私が通りかかったとき」を「私が生まれたとき」に、「ゲーム」を「言語」に言い換えたら、言語とはどんなものと言えるだろうか？

チェスや将棋、ポーカーのルールを知らなくても、見ているうちにそのルールがわかり、なんとなくゲームに参加できるのは、なぜだろう？

などです。人間が生まれたとき、すでに言語は存在し、いつのまにか、人間はそのルールを身につけ、それに参加し、それを発展させ、死んでいく。そして、つぎの世代の人間がその言語を使い、少しずつじわじわ形を変えていく。確かに言語は人間の口から出たものですが、人間が作ろうと思って作ったものではありません。何らかの力学によって、自然な仕組みで言語ができてきます。それは常に一定なものではなく、むしろ動的なものです。混沌としているようですが、その形には法則性があります。どういう例がわかりやすいでしょうか、あまりいい例ではありませんが、たとえば、人間の肘の関節は、内側には曲がるが、外側には曲がりませんよ！ というような「なあって」というような法則性です。その「なあって」というものが本当は何であるのかがよくわからないので、それを見つける研究をしているのですね。

2章の「チューリングマシン」では、

日常に見られるテープとヘッドに似たものを見つけて、そのどの部分がテープ、ヘッドに当たるかを述べよ。

というものです。ここでは、得体のしれない言語というものを、記述するには具体的に何をすればいいのか、そもそも記述するとはどんなことかなどを話し合います。言語もリアなものであり、その抽象的な姿を整理するには計算機モデルが役立つそうだというお話です。

3章ではジップの第二法則を紹介し、

人名の出現頻度、新聞記事に見られる単語の頻度がそれに従うのはなぜでしょう。また言語だけでなく、他の自然界にも見られるのはなぜでしょう。

と問いかけます。たとえば、人口の多い都市の数は少なく、人口の少ない町や村はめちゃくちゃ多い。ガラスの割れた大きい破片の数は少ないが、だんだん小さくなっていて、粉々になった破片の数はもう数えられないほどたくさんである。ジップ則を通して、単語の分布と自然の摂理にはどういう関係があるのかを議論してもらいます。実際に今も、なぜそれらがジップ則にしたがうのかはよくわかっていないものですから、この議論はそう簡単には終わりません。おそらく、とことんその理由を説明しなければ気がすまない理系の学生にはうってつけのトピックであったのでしょ。

東工大は伝統的に自然言語処理の研究者を多く輩出していることで有名です。その意味では東工大には、言語を扱う素地はあったと言えます。2016年4月、東京工業大学は日本で初めての学部と大学院を一緒にした学院を設置しました。そして、この4月より東工大では、正式に学士課程の科目名として「言語学」を設け、理系の学生のための言語学の授業がはじまります。理学・工学を学ぶ新入生の目には、東

工大の言語学はどうつるのでしょうか。まだ始まったばかりです。非常に楽しみです。

33

### 執筆者プロフィール

山元啓史（やまもと・ひろふみ）専門は言語学、言語変化、外国語としての日本語教育。オーストラリア国立学大学院博士課程修了。Ph. D in Linguistics。1993年筑波大学文芸・言語学系留学生センター助手、1995年同講師、1997年カリフォルニア大学サンディエゴ校客員研究員、2006年オーストラリア国立大学客員研究員、2009年東京工業大学留学生センター准教授、2017年東京工業大学リベラルアーツ研究教育院教授。著書は、“Japanese A Comprehensive Grammar” Routledge, 「コーパスと日本語史研究」ひつじ書房、などがある。

Copyright (C) YAMAMOTO, Hilofumi 2016- All Rights Reserved.

### 編集後記（編集室：ふじたまさえ）

第57号前編、後編ともにいつも以上に読み応えのある内容となりました。巻頭言をはじめ、ご寄稿いただいた皆さま、ありがとうございます！

どの内容も素晴らしかったのですが、特に個人的な興味としては、特別寄稿をいただいた OMNIA のことが気に入っています。また、巻頭言として掲載している山元先生の文章も、大変興味深い内容でした。

後編のイベントレポートの中では、国立国会図書館関西館の菊池さんがおっしゃっていた「DHの現状や課題などを体系的にまとめた日本語の解説書」が気になります。個人的な感想ですが、本メルマガが扱っている話題も含め、DHについて体系的にまとめたおすにはどういった媒体が良いのか考えてみると、印刷物よりは Wikipedia のようなデジタルのもののほうが合っているようにも思いました。

人文情報学月報編集室では、国内外を問わず各分野からの情報提供をお待ちしています。

情報提供は人文情報学編集グループまで...

DigitalHumanitiesMonthly[&]googlegroups.com

[&]を@に置き換えてください。

人文情報学月報 [DHM057]【後編】 2016年04月29日(月刊)

【発行者】"人文情報学月報"編集室

【編集者】人文情報学研究所 & ACADEMIC RESOURCE GUIDE (ARG)

【ISSN】2189-1621

【E-mail】DigitalHumanitiesMonthly[&]googlegroups.com

[&]を@に置き換えてください。

【サイト】<http://www.dhii.jp/>

Copyright (C) "人文情報学月報" 編集室 2011- All Rights Reserved.



## 第 6 章

# 大学研究室探検隊

大学研究室探検隊 Vol. 6: 東京工業大学 山元啓史研究室

### 6.1 取材

2018年2月号「サクセス15」pp.16-19. グローバル教育出版より、インタビュー記事が発行された。日本学術振興会に取材があったことを報告。

中学生のみなさんにはあまりなじみがないかもしれませんが、多くの人が進むであろう大学の研究室では、文系・理系を問わず、日々さまざまな研究が行われています。このコーナーでは、そうした研究室や研究内容を紹介していきます。ここで見つけた研究がみなさんの視野を広げ、将来の目標への道標となるかもしれません。第6回は、言語の可視化に関する研究を行う東京工業大学の山元教授の研究室を紹介します。

### 6.2 資料

(4ページ先から戻る順でご覧ください)







山元啓史(やまもと・ひろふみ) / 民間の日本語学校で教鞭をとった後、筑波大学文芸言語系助手、講師。カリフォルニア大学サンディエゴ校、オーストラリア国立大学で客員研究員を経て東京工業大学リベラルアーツ研究教育院、環境・社会理工学院 社会・人間科学系教授。夏休み中学生向けセミナーは毎年すぐに満員。2016年情報処理学会山下記念賞受賞。

ゆきのうちに はるはきにけり うくひすの...

snow of inside at spring (topic/ome (past) (perfect) warbler of

図1: 共出現パターンの作り方

「雪のうちに 春はきにけり 鶯の こほれる涙 いまやとくらん」という和歌のなかに出てくる単語でペアを作ります。このように同じ文に現れたペアを「共出現パターン」と呼びます。



約20名の小中学生が参加したワークショップの様子。和気あいあいとした雰囲気なかで言語学のおもしろさに触れられるプログラムです。

筆を入れる人はほとんどいませんが、ことばの形はそのまま、意味だけが変化して現在も使われることばになっていきます。

「昔の文章に出てくる単語はいまとは異なる意味を持っている可能性があります。昔の文章を読むときはどうしても現代の常識にあてはめて考えてしまいがちです。例えば『食べる』という単語の前にある単語は食べものをさすと思うでしょう。でも本当はまったく関係がないかもしれません。また、現存していないものの名前も私たちにはわかりません。人間はあくまで推測しか導き出せないのです」(山元先生)

そこで役立つのがコンピュータです。「はな(花)」という単語が花び

らだけをさすのか、つぼみや茎も含めた全体をさすのかわからなくても、コンピュータで多くの和歌のなかから「はな」という単語がどんな単語といっしょに使われているかを探し出し、その結果を図に書き出せば、そこから花の意味はもちろん、香りや感触が確認できるというのです。コンピュータによる分析は客観的で信ぴょう性のある結果としてとらえられます。

「私たちは昔のことばを直接聴いたり、昔の食べものを直接味わったりすることはできません。でも、昔のことばを分析すれば、タイムマシンに乗らずに昔の世界を感じることができるのです」(山元先生)

ここでの「昔」とは、平安時代をさします。現存する平安時代の書物から、日本語は千年以上も前から使われていた言語だということがわかっていきます。ここまで古い歴史を持つ言語は世界でも日本語とアイスランド語だけだそう。日本語はそれだけ過去をさかのぼって調査できる貴重な言語なのです。

### 和歌を科学的に分析して ことばの形をとらえる

では実際にどう図を作っていくか

というと、まず和歌から任意で2つずつ単語を取り出し、ペアを作っていきます。図1の「雪のうちに 春はきにけり 鶯の こほれる涙 いまやとくらん」という和歌からは、「雪、うち」「雪、春」「雪、き」「雪、けり」「雪、うくひす」...という形でどんなペアができていきます。

次に言葉の重みを調べていきます。「どこにも出てくる単語は『探す 価値のない単語』です。例えば千首の和歌すべてに出てくる単語があったら、このような単語は検索しても、なにかを特定するのに役立ちません。つまり情報量0です。一方、たまにしか出てこない単語は当時の人がなにかを伝えるために使った『探す価値のある単語』です。こういう単語の情報量は多くなります。

このように単語の重みを計算していきます。そして2単語の重みの平均値が大きいペアから図に出力していきます。これらの作業はすべてコンピュータで行います。

図2は古今和歌集から『春』に関することばのペアを集めたものです。図にすることで、目に見えないはずの『春』の形が見えるようになってきます」(山元先生)

図3は鶯(ウグイス)が出てくる

視野が広がる!?

# 大学研究室<sup>39</sup>

# 探

# 検

# 隊

Vol.6

東京工業大学  
山元啓史  
研究室

研究内容

コンピュータを  
利用して言語を  
可視化する研究

中学生のみなさんにはあまりなじみがないかもしれませんが、  
多くの人が進むであろう大学の研究室では、文系・理系を問わず、日々さまざまな研究が  
行われています。このコーナーでは、そうした研究室や研究内容を紹介していきます。  
ここで見つけた研究がみなさんの視野を広げ、将来の目標への道標となるかもしれません。  
第6回は、言語の可視化に関する研究を行う東京工業大学の山元教授の研究室を紹介します。

(画像・資料提供：東京工業大学 山元啓史研究室)

今

回取り上げるのは言語学に関する研究です。といってもみなさんがイメージする文系の学問としての研究とはひと味違います。お話を伺った山元啓史教授が在籍するのは理系トップレベルの大学として名高い東京工業大学。研究も「言語を可視化する」、つまり言語をコンピュータで分析して図として表現することで、言語の形を考えると、なんと面白いものなのです。

「言語を可視化する」と聞くと、なんだか難しそうですが、山元先生が「ひらめき☆ときめきサイエンス(※)」の一環として中学生向けに開いているワークショップ「目で見えてわかる昔の日本語と今の日本語×タイムマシンに乗らずに行ける昔の世界」で扱う内容はみなさんにもわかりやすいものとなっています。まずはその内容を見ていきましょう。

単語の意味は  
時代によって変わる

日本語に限らず、言語は一般的に一度形や音、つづりが定着するとそれらは変わりにくいのですが、その意味は時代によって変化していくといわれています。「下駄箱」や「筆箱」は、一般的にはもうそれらに下駄や

(※) 日本学術振興会が科学研究費助成事業(科研費)の一環として主催するプログラム。大学や研究機関で科研費により行われている最先端の研究を、小5・小6、中学生、高校生が体験できる。全国の大学・研究機関で行われている。



## 第7章

# JADH: 論文・ポスター

1. OSDH2011: Osaka Symposium on Digital Humanities 2011 “Graph Representation of the Connotations of Classical Japanese Poetic Vocabulary”
2. JADH2012: Inheriting Humanities; Program PDF; Design of Serial Comparison Model for the Diachronic Corpus Study of Japanese; Emotive Adjectives and Verbs of the Heian Japanese .
3. JADH2013: Bridging GLAM and Humanities through Digital Humanities; Lexical Modeling of Yamabuki (Japanese Kerria) in Classical Japanese Poetry; A Diachronic and Synchronic Investigation into the Properties of Mid-Rank Words in Modern Japanese;
4. JADH2014: Bridging GLAM and Humanities through Digital Humanities; A Visualization and Analysis System for Japanese Language Change: Quantifying Lexical Change and Variation using the Serial Comparison Model; Development of an Asymptotic Word Correspondence System between Classical Japanese Poems and their Modern Translations;
5. JADH2015: Encoding Cultural Resources;
6. JADH2016: Digital Scholarship in History and the Humanities;
7. JADH2017: Creating Data through Collaboration;
8. JADH2018: Leveraging Open Data;
9. UCLA workshop in 2016–2018: UCLA Department of Asian Culture and Language invited me for the UCLA workshop and lectures of Computer and Humanities 2016–2018; In 2016, I present two lectures for faculty members, Graduate students, and one class for undergraduate students; In 2018, I present two lectures for faculty members as well.
10. Posters and flyers 2011–2018:



Hilo Yamamoto  
Ph. D. Linguistics

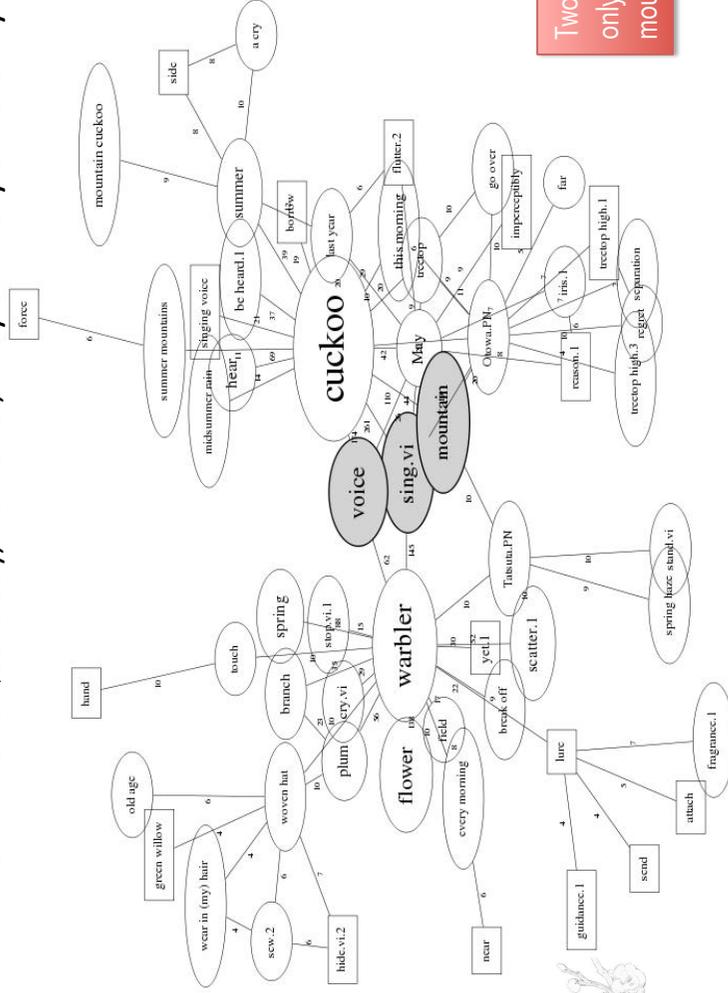
# Graph Representation of the Connotations of Classical Japanese Poetic Vocabulary

## Computer modeling using co-occurrence patterns

Two birds, warbler and cuckoo, come to Japan in spring. In the *Kokinshu* (ca. 905), however, they are very differently used.



Tokyo Institute of Technology



Cuckoo

Two birds' network shares only three words: i.e., sing, mountain, and voice.



Warbler

Poem texts must contain more information we have not known yet!



# Design of Serial Comparison Model for the Diachronic Corpus Study of Japanese

**Hilofumi Yamamoto**  
Tokyo Institute of Technology

**Makiro Tanaka**  
National Institute for Japanese Language and Linguistics, Japan

**Yasuhiro Kondo**  
Aoyama Gakuin University

## Development of Diachronic Corpus

Project by the National Institute for Japanese Language and Linguistics, Japan, NINJAL: 2009–13, 4 year project.

Main purpose: Study of Japanese language  
(sub) purpose: Study of Japanese (classic) literature

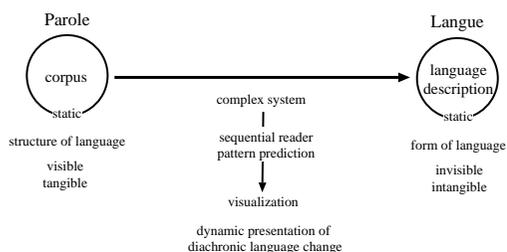


Figure 1: Corpus and Description, Langue and Parole: The nature of language is dynamic and always changing while the phenomena of language might be static. We should consider the dynamic change of language as a component comprised of various elements. The feature of language we usually observe is a complex system and tangled with wide-ranging elements.

## Contents of Diachronic Corpus

1. The Tale of the Bamboo-Cutter (ca. 890; Taketori monogatari; 12,583 tokens)
2. Tales of Ise (ca. 901; Ise monogatari; 15,900 tokens)
3. Tales of Yamato (ca. 950; Yamato monogatari; 26,733 tokens)
4. The Tosa Diary (ca. 935; Tosa nikki; 8,113 tokens)
5. The Pillow Book (ca. 996; Makura no sōshi; 79,861 tokens)
6. Tale of Genji (ca. 1100; Genji monogatari; 510,711 tokens)

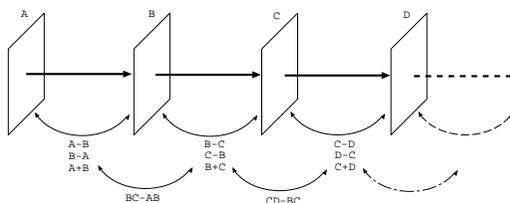
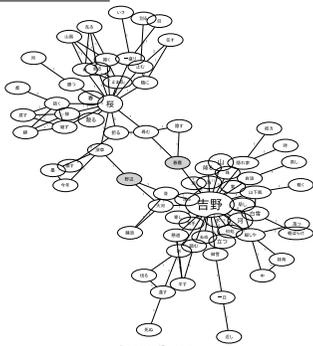


Figure 2: Extraction of delta from each synchronic layer: A, B, C and D are arbitrarily-assigned synchronic layers on the time axis. Examination of linguistic transitions is achieved through the comparison of lexical items in each layer with those in other layers, and the discovery of common principles appearing in the delta of data extracted from both systems as well.

**A case study:** use of **SAKURA** (cherry blossoms) in **Mt. Yoshino** → Kokinshū (ca. 905) vs Shinkokinshū (1205)



Sakura (桜) and Yoshino (吉野), a place name in Nara prefecture  
← Kokinshū (ca. 905)  
Shinkokinshū (1205) →  
during 300 years differences.

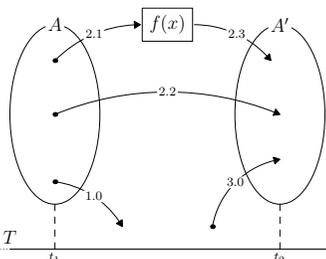
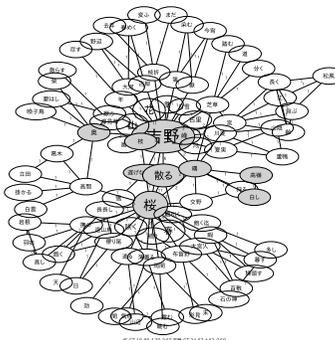


Figure 3: Serial comparison model; differential model of transitional linguistic elements of target texts; A is a set of elements that occurred at Time  $t_1$ ; A' is a set of elements that occurred at Time  $t_2$ ; T is the time axis;  $f(x)$  is a function for converting an element  $x$  of A into that of A'.

## Future Task

- To define linguistic units suitable for each era
- To develop a dictionary for machine analysis  
→ it allows us syntagmatic and paradigmatic anal

## Conclusion

- Addressed basic concepts and framework of diachronic corpus
- Illustrated the serial comparison model for historical analysis  
→ Lexical differences between any two groups of t





# Lexical Modeling of *Yamabuki*, Japanese Kerria in Classical Japanese Poetry

Hilofumi Yamamoto / Tokyo Institute of Technology  
yamagen@ryu.titech.ac.jp



## Introduction

- We conduct a lexical study of classical Japanese poetry using network modeling.
- The terms *yamabuki* (kerria), *kahazu* (frog), and *Ide* (placename) are contained in some poetic dictionaries as entry items or collocations, and we have confirmed that they have strong relationships with each other.
- We have discovered the hub node term *yahe* in network models. The term *yahe* is, however, not recorded in any poetic dictionaries even as a single term.

### Material: *Hachidaishū*

the eight anthologies compiled by the order of Emperors (ca. 905–1205), which contains about 9,500 poems.

### Calculation methods:

$$w(t, d) = (1 + \log tf(t, d)) \cdot idf(t)$$

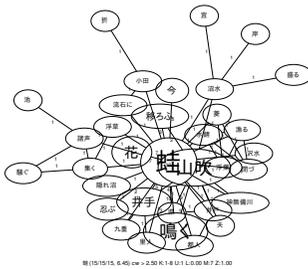
$$cw(t_1, t_2, d) = (1 + \log ctf(t_1, t_2, d)) \cdot cidf(t_1, t_2)$$

$$cidf(t_1, t_2) = \sqrt{idf(t_1) \cdot idf(t_2)}$$

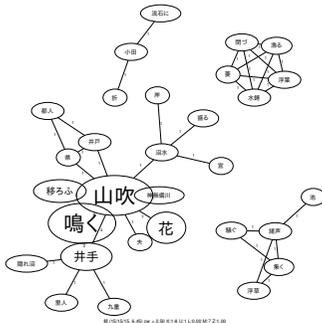
$$idf(t) = \log \frac{N}{df(t)}$$

**Figure 1:** The picture of “*Yamabuki To Kahazu*” (kerria and frog) by Hiroshige Utagawa (<http://www.gekkanbijutsu.co.jp/shop/goods/030761011.htm>).

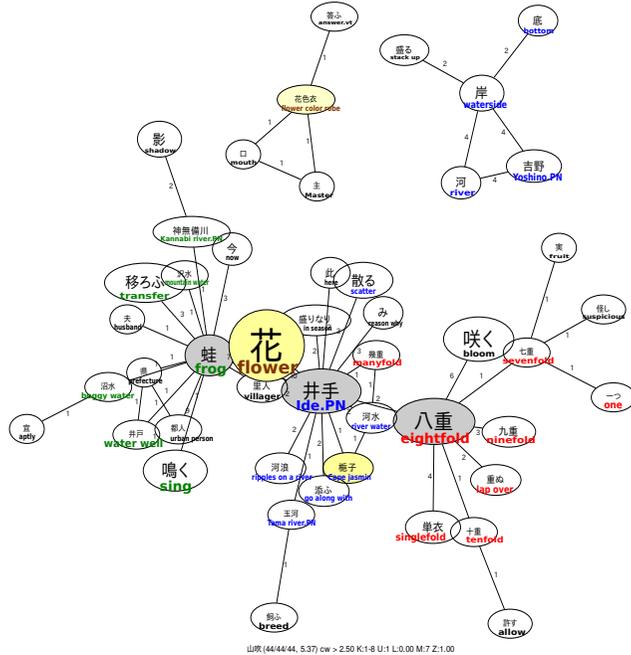
## Result



**Figure 2:** Graph model of *kahazu* (蛙, frog) before pruning node 蛙.



**Figure 3:** Graph model of *kahazu* (蛙, frog) after pruning node 蛙.



**Figure 4:** Graph model of *Yamabuki*: a core node, 山吹 *yamabuki*, is pruned. *kahazu* (蛙, frog), *Ide* (井手, place name, proper name), and *yahe* (八重, eightfold or double flower) are observed as hub nodes.

A minor term *yahe* (eightfold) can be shown as a hub node which plays a major role in connecting a topic word with other peripheral words which support/demonstrate poem stories. These minor words are not seen in poetic term dictionaries.

## Conclusion

1. Discern not only patterns described by experts but also patterns yet undescribed, and
2. Identify not only specific or tangible words but also abstract or conceptual words which have a tendency to be left out of dictionaries.



**Figure 5:** Single petal (left), white petal (center), and plena petal (right) of *yamabuki*. (<http://mkfarm.blog118.fc2.com/blog-entry-27.html>)





# Development of an Asymptotic Word Correspondence System between Classical Japanese Poems and their Modern Translations

Hilofumi Yamamoto\*†

\* University of California, San Diego †Tokyo Institute of Technology

Hajime Murai†

Bor Hodošček‡

‡Meiji University

## Introduction

- This project will develop an automatic word alignment system for parallel texts comprising of Classical Japanese poems and their associated modern translations.
- By using these parallel texts, we will clarify the details of language change within Japanese in an objective procedural manner that is not influenced by human observations.
- Our aim is to develop a thesaurus of classical Japanese poetic vocabulary using the system.

## Problem

What is Waka:



Tatsuta-Hime.. (5 syllables)  
tamakuru KAMI no (7)  
arebakoso (5)  
aki no konoha no (7)  
nusa to chirurame (7)

because Princess Tatsuta  
has a god to whom she offers brocades,  
the leaves of trees  
in autumn will scatter as an offering.

### 1. Orthography Problem

龍田, 立田, 竜田, たつた all indicate same placename: 'Tatsuta' in Nara pref.

### 2. Unit size Problem

Does 卯の花 consist of one word or 卯/の/花 three words?

### 3. Attribution Problem

Is 卯の花 the name of a flower or bean curd refuse?

### 4. Polysemy/PUN Problem

海松藻 'mirume' a kind of sea weed; also means 見る目 (human eyes).



## Methods

**Material:** *Kokinshū* a.k.a. *Kokinwakashū* is: the first anthology compiled by the order of Emperor Daigo (ca. 905), which contains about 1,100 poems. And 10 sets of their **Contemporary Japanese Translations (CT)**

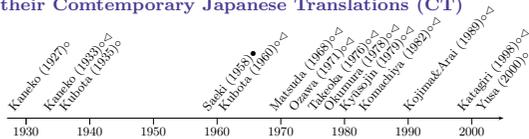


Figure 1: Dates of publication of annotations of the *Kokinshū*: ○ indicates that it has CT; ● indicates that it does not include CT; ▷ indicates that it is used in this project.

**Mutual Co-occurrence Rate:** Murai (2010)

$$mcr(o, t) = p(o|t) p(t|o)$$

where,  $o$  indicates a token in original texts;  $t$ , a token in translation texts;  $mcr(o, t)$ , the mutual co-occurrence rate;  $p(o|t)$ , the rate when a token  $o$  and  $t$  occur at the same time in corresponding texts which are original texts and translation texts.

→ when  $mcr$  is large enough, it will be estimated that token  $o$  and  $t$  are **contextually equivalent**.

## Result

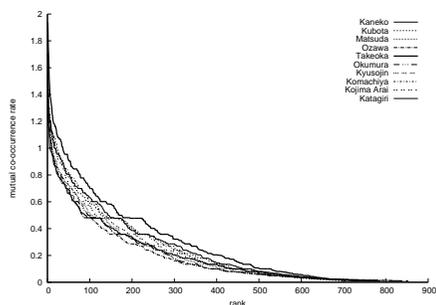


Figure 2: Distribution of Mutual Co-occurrence Rate: original text *Kokinshū* and ten sets of its translation texts.

### Good or poor estimated pairs

Table 1: Good estimated pairs and poor estimated pairs; the values of good pairs are the first ten items (over 1.3); and the values of poor pair items are the last ten items (lower 0.01).

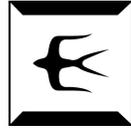
no.	good	pairs	poor	pairs
1	鳴く	鳴く cry	異なり	あの
2	風	風 wind	雫	どうして
3	世の中	世の中	此の	この
4	人	人 human	随に	まま
5	春	春 spring	匂ふ	美しい
6	秋	秋 autumn	見る	せい
7	時鳥	時鳥 cuckoo	連れ	つく
8	時鳥	ほととぎす	立ち返る	言う
9	散る	散る fall	有り	つく
10	見る	見る see	有り	まさしく

## Conclusion

1. This project has already begun: the parallel corpus of the *Kokinshū* has been constructed.
2. We are now working on the development of computer software and the optimization of the calculation methods.

## Reference

- Murai, Hajime. 2010 Extracting the interpretive characteristics of translations based on the asymptotic correspondence vocabulary presumption method: Quantitative comparisons of Japanese translations of the Bible. *Journal of Japan Society of Information and Knowledge* Vol. 20, No. 3, 293-310.



# Development of an Asymptotic Word Correspondence System between Classical Japanese Poems and their Modern Translations

Hilofumi Yamamoto\*†  
\* University of California, San Diego

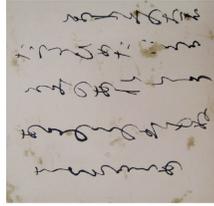
Hajime Murait†  
†Tokyo Institute of Technology

Bor Hodošček‡  
‡Meiji University

## Objectives:

The development of the **THESAURUS** of Classical Japanese Poetry. (10–13th century; Heian period )

### Waka: Japanese poetry



*Tatsuta-Hime...*

*tamukuru kami no / arebakoso*

*aki no konoha no / nusa to chirurame*

because princess Tatsuta has a god whom she offers  
brocades, the leaves of autumn will scatter as an offering.



## Problems:

- Unit: 卵の花 or 卵 / の / 花
- Orthography: sad!  
さびしい / さみしい / 寂しい / 淋しい
- Attribution:  
卵の花 ∈ plant or 卵の花 ∈ food



## Solution:

- Using Pararell Corpora:  
Original Poems and Contemporary Translation Texts
- Using Mutual Co-occurrence Ratio:  $mc_r(o, t) = p(o|t)p(t|o)$   
To estimate the rate of contextual similarity between classical words and contemporary words.



# The differences of connotations between two flowers, plum and cherry, in classical Japanese poetry, 10th century.

Hilofumi Yamamoto Tokyo Institute of Technology

## Introduction

- This project addresses an analysis of connotations of flowers in classical poetry: i.e., ‘ume’ (plum) and ‘sakura’ (cherry) .
- We will identify the characteristics of two flowers by computer modeling.
- Using parallel texts of original texts and contemporary translations of classical Japanese poetry, *the Kokinshū*, we will clarify the details of connotations in an objective procedural manner that is not influenced by human observations.
- The aim is to examine whether or not the residual of *CT* – *OP* gives information on the non-literal elements of *OP*.

## Problem

1. What is the difference between *ume* (plum) and *sakura* (cherry)?
2. What kind of connotations does each flower contain?
3. Which picture is that of cherry flowers?



## Methods

**Material: *Kokinshū*** a.k.a. *Kokinwakashū* is: the first anthology compiled by the order of Emperor Daigo (ca. 905), which contains about 1,111 poems. And 10 sets of their **Contemporary Japanese Translations (CT)**

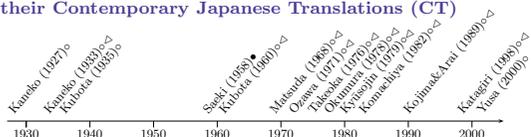


Fig. 1: Dates of publication of annotations of the *Kokinshū*: ◦ indicates that it has CT; • indicates that it does not include CT; ▷ indicates that it is used in this project.

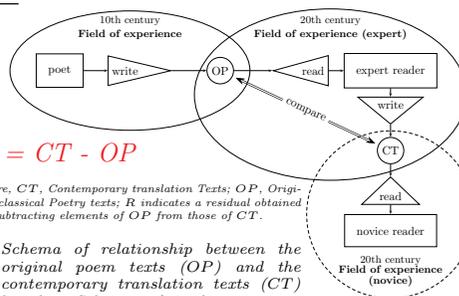


Fig. 2: Schema of relationship between the original poem texts (*OP*) and the contemporary translation texts (*CT*) based on Schramm (1954).

## Result

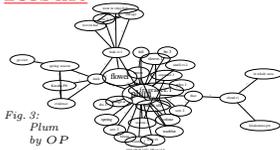


Fig. 3: Plum by OP

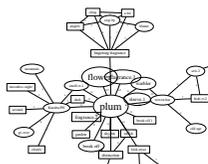


Fig. 4: Plum by CT

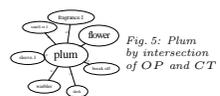


Fig. 5: Plum by intersection of OP and CT

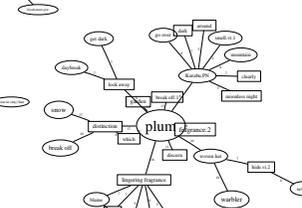


Fig. 6: Plum by subtracting OP from CT

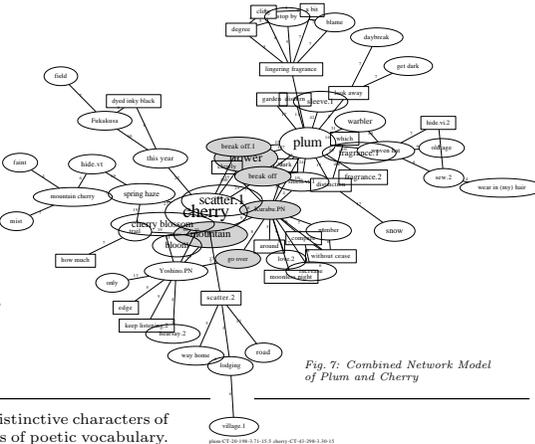


Fig. 7: Combined Network Model of Plum and Cherry

## Conclusion

- It will be necessary to examine not only common nouns but also the distinctive characters of proper nouns in order to further examine the connotative associations of poetic vocabulary.
- We observed proper nouns such as place names, *Kurabu*, *Tatsuta*, *Otowa*, *Yoshino* in the network models of common nouns, and concluded that they seem to strongly influence the associations of poetic vocabulary.
- The relative salience clearly indicates that both *ume* (plum) and *sakura* (cherry) share *Kurabu yama* (Mt. Kurabu), which comprises a cluster of nodes in the sub-network.

## Reference

- Schramm, W. L. 1954. How communication works. *The process and effects of mass communication*. 3–26. University of Illinois Press.
- Yamamoto, H. 2006. Extraction and Visualisation of the Connotation of Classical Japanese Poetic Vocabulary. Symposium for Computer and Humanities, 2006. The information processing society of Japan. Vo. 2006, No. 17, 21–8.



The differences of connotations between two flowers,  
plum and cherry, in classical Japanese poetry, 10th century

Hilofumi Yamamoto Tokyo Institute of Technology / University of California, San Diego

## Objectives:

To reveal the connotative differences between **Plum** and **Cherry** in Classical Japanese Poetry. (10–13th century; Heian period )

## Problem:

What kind of connotative meaning does each flower contain?



## Solution:

- Using Pararell Corpora:
  - i. e., Original Poems (OP) and Contemporary Translation (CT) Texts, and using Network Modeling, ....
  - we will analyze the residuals when subtracting OP from CT co-occurrences.

## Development of the Dictionary of Poetic Japanese Description



Hilofumi Yamamoto  
Tokyo Institute of Technology



Bor Hodošček  
Osaka University

### Introduction

- This paper proposes to further the development of a dictionary of classical Japanese poetry using pairwise term information (Yamamoto et al., 2014).
- Information on pairwise terms between an index and related term such as “flower–spring” is not included within traditional modern and classical Japanese dictionaries, even though this information connects terms with their contexts in a transparent way and thus offers an unbiased method for inferring the meaning of old Japanese terms.
- An R package for the analysis of linked communities in networks, linkcomm (Kalinka and Tomancak, 2011), is used to extract subordinate terms. Average, McQuitty, and single linkage methods are evaluated for the quality of their extraction of subordinate clauses of terms representing the ‘cherry’, ‘plum’, and ‘orange’ flowers. All methods extracted similar subordinate terms, which were quite natural in the context of classical Japanese poetry.

### Problem

1. Many scholars of Japanese poetry have tried to explain poetic vocabulary based on their **intuition** and **experience**.
2. As scholars can only describe constructions that they can **consciously** point out, those that they are **unconscious** of will **NEVER** be uncovered. ⇒ In order to conduct more exact and unbiased descriptions:
  - 1) using computer-assisted descriptions;
  - 2) using co-occurrence weighting methods on corpora of Japanese poetry; and
  - 3) using linkcomm R package, extract the lists of words grouping sub communities.
 ⇒ allows one to **BETTER GRASP** the construction of poetic words.

### Methods

**Calculation:** *Linkcomm* (Kalinka and Tomancak, 2011) for sub communities of three flowers: *ume* (plum), *sakura* (cherry), and *tachibana* (mandarin orange).

**Material:** *Hachidaishū* (ca. 905–1205) from *Kokkatakaikan* (Shin-pen Kokkatakaikan Henshū Committee, 1996), *Nijūichidaishū* database published by NIJIL (Nakamura et al., 1999), *Shin-Nihon Koten Bungaku Taikei* (Kojima and Arai, 1989), and *Shin-kokinshū* (Kubota, 1979).

### Result

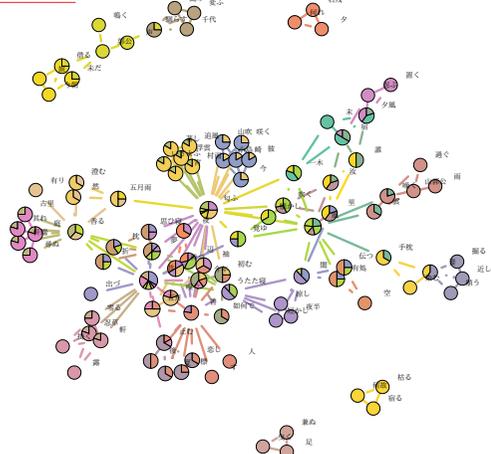


Fig. 1: Network of Words; mandarin orange.

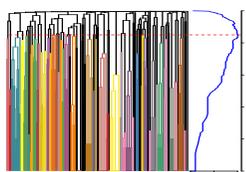


Fig. 2: Link Community Dendrogram.

Table 1: Sub-clusters of orange.

No. node	average (.43)	mcquitty (.43)		single (.38)	
	edge	node	edge	node	edge
1 mukashi (old days)	7 mukashi	7 mukashi	7 mukashi	5	
2 nihofu (smell)	6 nihofu	6 nihofu	6 nihofu	4	
3 kaze (wind)	5 kotoshi	4 yume	4 yume	4	
4 yume (dream)	5 atari	4 kaoru	4 kaoru	3	
5 kotoshi (this year)	4 matsu	4 kotoshi	4 kotoshi	3	
6 atari (aroud)	4 kaze	4 somu	4 somu	3	
7 matsu (to wait)	4 yume	4 samidare	4 samidare	3	
8 kaoru (fragrance)	3 somu	3 ori	3 ori	3	
9 samidare (summer rain)	3 kaori	3 makura	3 makura	3	
10 somu (to dye)	3 yami	3 omohine	3 omohine	3	

### Conclusion

- Pairwise term information generated by the community centrality procedure works well.
- R package “linked communities” could extract proper sub cluster terms which contribute to the description of classical Japanese poetry.

### Reference

- Kalinka, A. T. and Tomancak, P. 2011. linkcomm: an R package for the generation, visualization, and analysis of link communities in networks of arbitrary size and type. *Bioinformatics*. 2011–2. **27** (14).
- Yamamoto, H., Hajime Murai, Bor Hodošček. 2014. Development of an Asymptotic Word Correspondence System between Classical Japanese Poems and their Modern Translations. Symposium for Computer and Humanities, 2014. The information processing society of Japan. Vol. 2014, No.3, 157–62.



## Development of the Dictionary of Poetic Japanese Description



Hilofumi Yamamoto  
Tokyo Institute of Technology



Bor Hodošček  
Osaka University

### Objectives:

To develop the dictionary of **YAMATO Japanese** description.

**Problem:** (10–13th century; Heian period )

- Missing an unbiased method for **inferring** the meaning of old languages.
- **Pairwise terms** such as ‘flower-spring’ are **NOT** included in dictionaries.

**QUIZ:** which picture was taken in winter?



A.



B.



C.

### Solution:

- Using linkcomm, the calculation of sub communities, we extract pairwise terms (index-relating term).

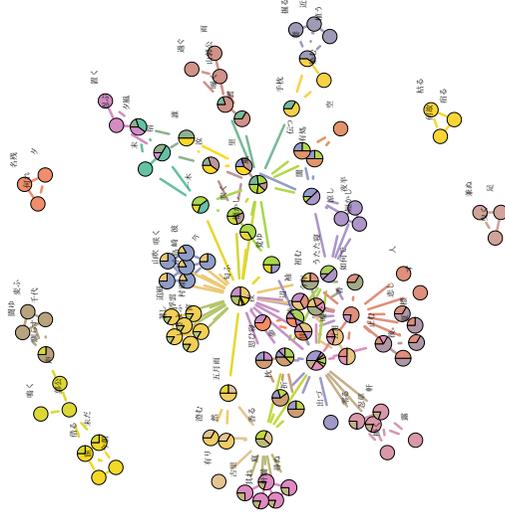


Fig. 1: Network of terms; mandarin orange.

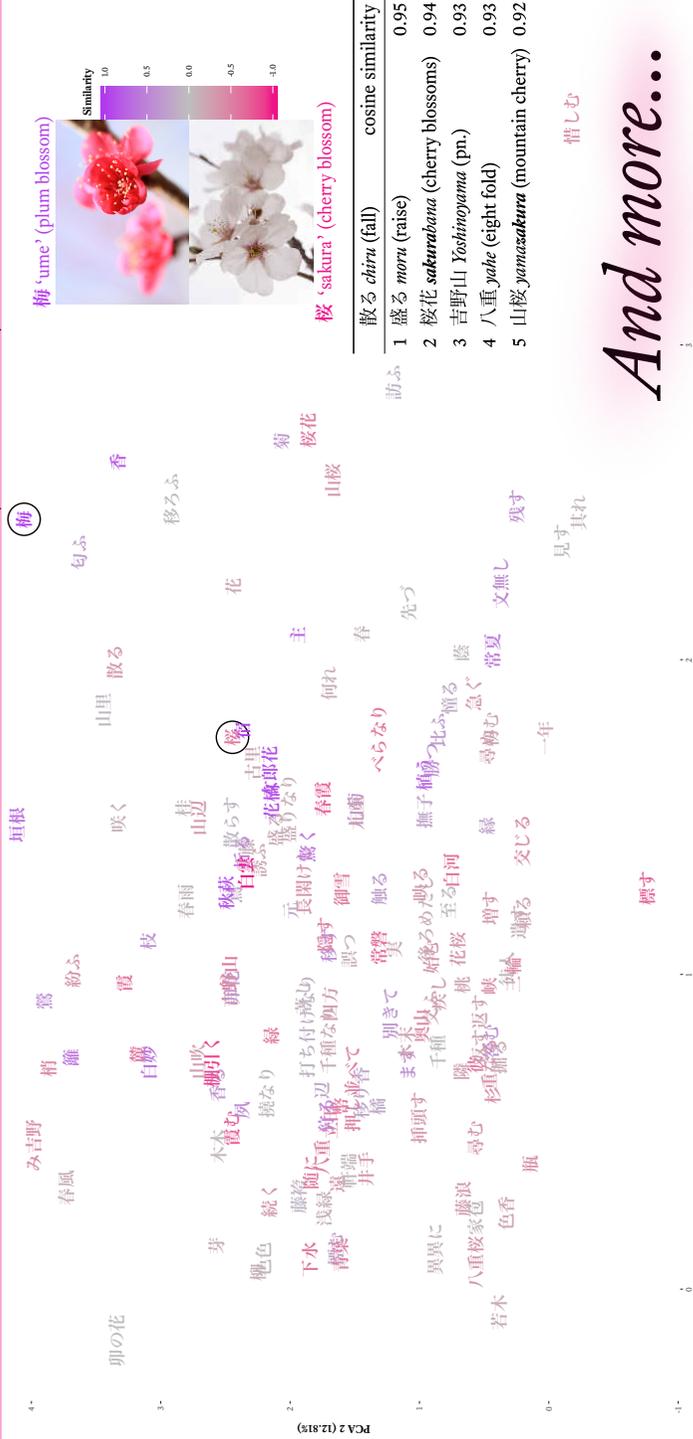
# RELATIONSHIPS BETWEEN FLOWERS IN A WORD EMBEDDING SPACE OF CLASSIC JAPANESE POETRY


 HiloFumi Yamamoto, Tokyo Institute of Technology  
 yamagen@ila.titech.ac.jp


 Bor Hodošček, Osaka University  
 bor@lang.osaka-u.ac.jp

@JADH2017  
 September 11

Examine the possibility of word embedding spaces (Word2Vec) to explain the semantic relationships between classical Japanese poetic terms within the *Hachidaishū* poem anthology. (ca. 905–1205)



*And more...*

Figure 1: PCA of word embedding space (4157 words × 50 dimensions) filtered to include only top 100 similar words for each of ume and sakura (150 total). Similarity is represented by the difference in similarity scores between ume and sakura, scaled to [-1, 1].

# RELATIONSHIPS BETWEEN FLOWERS IN A WORD EMBEDDING SPACE OF CLASSIC JAPANESE POETRY

Hilofumi Yamamoto, Tokyo Institute of Technology  
yamagen@ila.titech.ac.jp



Bor Hodošček, Osaka University  
bor@lang.osaka-u.ac.jp

## INTRODUCTION

- Word embedding methods such as Word2Vec (Mikolov et al., 2013; Le and Mikolov, 2014) have been shown effective in extracting semantic knowledge from large corpora.
- Quantify the relationship between the content of a word and its word embedding vector.
- Examine the possibility of word embedding spaces to explain the semantic relationships between classical Japanese poetic terms.

## PROBLEM

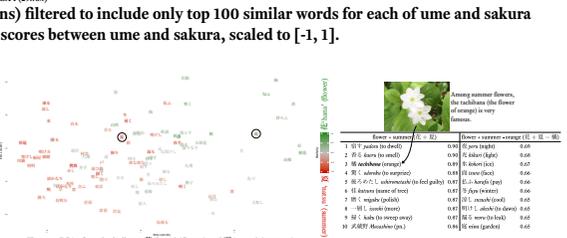
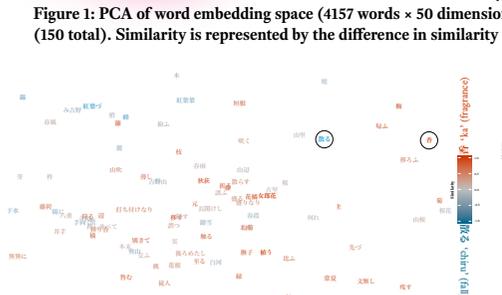
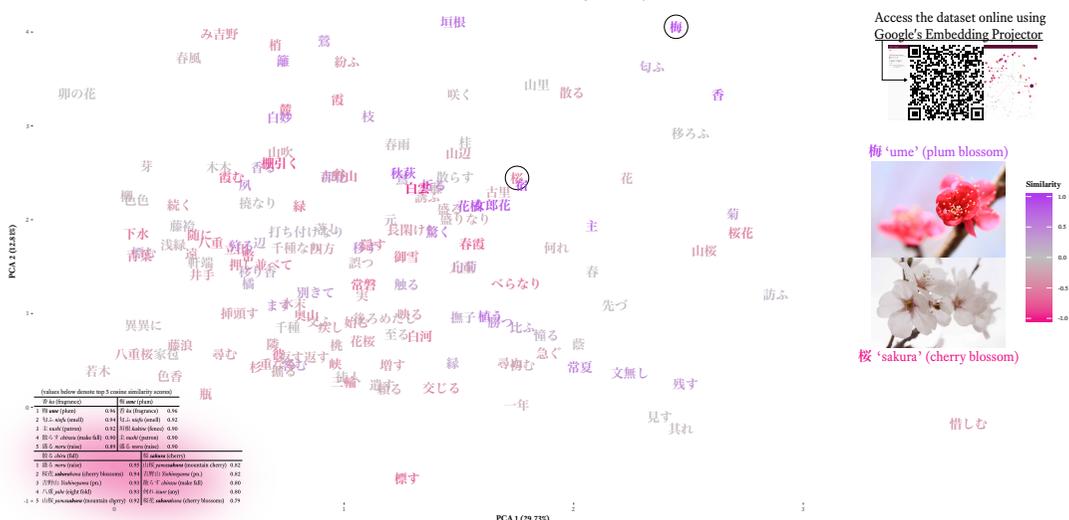
- Can word embeddings trained on the Hachidaiishu encode enough semantic information to find subordinate words via their superordinate concept?

## MATERIALS

- Hachidaiishu*: classical Japanese poem anthologies compiled under decree by Emperors (ca., 905–1205), comprising approximately 9,500 poems and 159,183 tokens (Source: *Kokkakaitan/Nijūichidaiishū* database published by NIJIL).
- Each poem is tokenized into lemma forms by kh (Yamamoto, 2007) which divides poem texts into tokens using a classical Japanese dictionary.

## METHODS

- 50-dimensional skip-gram model with negative sampling, context window covering the whole poem using Gensim 2.3.0 (Rehurek & Sojka, 2010).
- In order to examine the notable relationships between 'ka' (fragrance), 'chiru' (fall), we look at the cosine similarity scores between terms in the word embedding space generated by Word2Vec.



## RESULTS

- 'ka' (fragrance) is related to 'ume' (plum) (replicating Mizutani, 1983).
- Falling flowers denote 'sakura' (cherry) and not 'ume' (plum); 'sakura' (cherry) relates to chiru (fall), which indicates that people at the time lamented falling sakura (falling cherry blossom petals) (replicating p. 84 in Katagiri, 1983).
- Subtracting tachibana out from the summer vectors reveals a vector space devoid of relationships between natsu (summer) and hana (flower). These relational expressions (summer + flower; summer - flower - tachibana) reproduce our current understanding of the relationships between flowers and seasons as well as some emotions associated with them in the word embedding space.

## CONCLUSION

- Word embeddings allowed us to extract specific subordinate words based on the superordinate concept of classical terms → when the distance between two terms such as 'tachibana' (orange) and 'natsu' (summer) is close enough, the superordinate concept A indicates the subordinate concept a.
- We could therefore verify that it allows us to extract the concrete name from its superordinate concept.

## REFERENCES

Keighly, Vicki (1983) *Etymology of words*. Vol. 31 of *Kokusho* series, Tokyo: Kokusho Shoin.

Le, Quoc V. and Tomas Mikolov (2014) "Distributed Representations of Sentences and Discourse", *CoRR*, Vol. abs/1405.4053, URL: <http://arxiv.org/abs/1405.4053>.

Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean (2013) "Efficient Estimation of Word Representations in Vector Spaces", *CoRR*, URL: <http://arxiv.org/abs/1301.3781>.

Mizutani, Satoru (1983) *Gin (Fragrance)*, Vol. 2 of *Ankoku Shūkyō Shū-Kan*, Tokyo, Japan: Ankoku Shoin.

Rakata, Radosław and Tomáš Štěpánek (2015) "Subspace Framework for Topic Modeling with Large Corpora", in *Proceedings of the LREC 2015 Workshop on New Challenges for NLP Frameworks*, pp. 45–50, Valletta, Malta: ELRA, May, <http://www.lrec-conf.org/proceedings/lrec2015/>.

Rehurek, Jan and Radim Sojka (2010) *Katmatok - A Collection of Poem Indices and Metadata*, Boston MA USA: Cheng and Tsai Company.

Yamamoto, Hilofumi (2007) "Waka no seme no Hangeki no shūki shōmei / JIS-Tagger for Classical Japanese Poems", *Nihongo no Kenkyū / Studies in the Japanese Language*, Vol. 3, No. 3, pp. 31–39.



A

DIGI

WITH  
YAMAMOTO  
HILOFUMI

TAL

HUM

ANITIES

11/  
29  
9a

CLINIC

WEST ELECTRONIC CLASSROOM,  
CHARLES E. YOUNG  
RESEARCH LIBRARY

RICHARD C. RUDOLPH  
EAST ASIAN LIBRARY  
THE YANAI INITIATIVE  
UCLA LIBRARY

TOKYO INSTITUTE  
OF TECHNOLOGY  
Image courtesy of the  
Museum of Fine Arts, Boston  
www.mfa.org

3p



### WORKSHOP on Japanese Text Mining

**11/29, 11am-1pm @Presentation Room (YRL Rm 11348A)**

- Japanese corpora for mining: NINJAL, etc.
- Digital text-mining & analysis tools: MeCab, etc.; word2vec, text reuse analysis, and topic modeling.

### Consultation on JDH Projects:

**11/28-30, 1-5pm @West Classroom (YRL Rm 23167)**

- 11/28: Extensive reading
- 11/29: Ryukyuan dialects, Speech corpora, *Genji Monogatari*
- 11/30: Other topics

*Cosponsored by*

*Tadashi Yanai UCLA-Waseda Initiative for Globalizing Japanese Humanities,  
UCLA Terasaki Center for Japanese Studies,  
Tokyo Institute of Technology, and UCLA Library*

**11/29/2017**

**Wednesday**

**11 am – 1 pm**

**FREE  
WORKSHOP  
Japanese Text  
Mining**

**Hilofumi Yamamoto**  
Linguist and Professor  
Tokyo Institute of Technology

**Peter Broadwell**  
Academic Project Developer  
UCLA Digital Library

RESEARCH LIBRARY (CHALRES E. YONG)

**@ Presentation Room**  
(YRL Rm 11348A)  
Open to UCLA faculty,  
students and staff

Light Lunch Provided

**RSVP required**

Contact: East Asian Library  
Tomoko Bialock



# A Study on the Distribution of Cooccurrence Weight Patterns of Classical Japanese Poetry

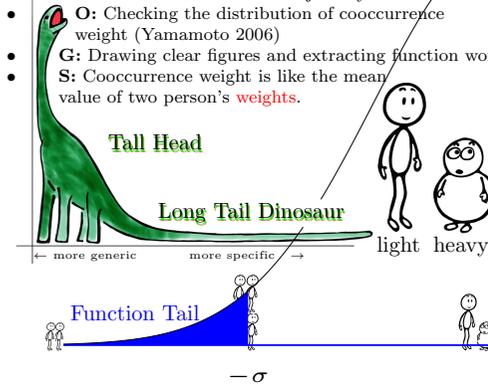
Hilofumi Yamamoto  
Tokyo Institute of Technology

Bor Hodošček  
Osaka University



## Introduction

- **P:** Hairball effect or spoke effect (Yamamoto 2005)
- **P:** Difficult to observe all word adjacency features.
- **O:** Checking the distribution of cooccurrence weight (Yamamoto 2006)
- **G:** Drawing clear figures and extracting function words.
- **S:** Cooccurrence weight is like the mean value of two person's weights.



## Methods

Material: *the Hachidaishū (ca. 905–1205)*

Calculation of Cooccurrence Weight:  $cw$

$$w(t, d) = (1 + \log tf(t, d)) \cdot idf(t)$$

$$cw(t_1, t_2, d) = (1 + \log ctf(t_1, t_2, d)) \cdot cidf(t_1, t_2)$$

$$cidf(t_1, t_2) = \sqrt{idf(t_1) \cdot idf(t_2)}$$

$$idf(t) = \log \frac{N}{df(t)}$$

Distribution of  $cw$  becomes **Bell curve**.

- Over  $\sigma \Rightarrow$  Content Tail.
- Under  $-\sigma \Rightarrow$  Function Tail.



C. F. Gauss (1777–1855)

## Result

Table 1: Upper cutoff patterns of *ame* (sakura):  $cw =$  co-occurrence weight;  $z =$  z-value (normalized value of frequency). word annotations: ari(be), ba(cond.), ha(topic), hana(flower), hito(human), keri(post.), ki(past.), koso(emphatic.), miru(see), mo (absol.), nani(no exist), mi(erg.), o(obj.), onom(think), rami(aux.will), so(do), te(p.), to(adv), ware(we), zo(emphatic.), zu(erg.)

$cw$	$z$	pattern	$cw$	$z$	pattern	$cw$	$z$	pattern			
1	0.62	-0.91	mo-keri	11	0.59	-0.96	nasi-ha	21	0.52	-1.05	mu- o- zo
2	0.62	-0.92	hana-o	12	0.57	-0.98	o-ramu	22	0.52	-1.05	o- zo
3	0.62	-0.92	o-kono	13	0.57	-0.98	mo-ramu	23	0.52	-1.05	miru- o
4	0.60	-0.94	zu-keri	14	0.57	-0.98	ha-ki	24	0.48	-1.09	ba- mo
5	0.60	-0.94	su-ha	15	0.56	-1.00	zu-mo	25	0.48	-1.09	o-keri
6	0.60	-0.94	to-ha	16	0.56	-1.00	o-te	26	0.43	-1.16	zu- ha
7	0.59	-0.96	ari-ha	17	0.55	-1.01	hito-mo	27	0.43	-1.16	te- o
8	0.59	-0.96	ari-mo	18	0.54	-1.02	zu-te	28	0.43	-1.16	te- ha
9	0.59	-0.96	ware-mo	19	0.52	-1.05	zo-ha	29	0.34	-1.27	o- ha
10	0.59	-0.96	nasi-o	20	0.52	-1.05	omou-o	30	0.34	-1.27	o- mo

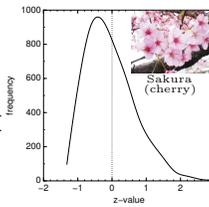


Table 2: Lower cutoff patterns of *ame* (sakura) in Kokinshū: 30 out of 164 patterns extracted;  $cw =$  co-occurrence weight;  $z =$  z-value (normalized value of frequency) word annotations: ba(cond.), bakari(only), besi(should be), chiru(fall), fukakoso(deep green), hana(flower), isa(already), kakusu(hide), koi(win), koto(pull), komoru(go deep inside), magiru(mix), makasu(entrust), maku(wind up), manimani(as it is), masi(as), masu(mix), me(eye), minami(south), miyako(city), mono(thing), nagara(even if), sakura(cherry), si(emphatic), sumi(black ink), tate(start/stand), tazumu(being around), tu(past), uturu(change), watausu(give), yamakaze(mountain wind), yamu(stop), yanagi(willow), yononaka(world)

$cw$	$z$	pattern	$cw$	$z$	pattern		
1	3.86	3.18	yamu-manimani	106	2.38	1.31	si-fukakusa
2	3.75	3.04	minami-magiru	107	2.38	1.31	sakura-hana
3	3.67	2.93	minami-maku	108	2.38	1.31	sakura-isa
4	3.61	2.86	maku-magiru	109	2.38	1.31	sakura-ha
5	3.42	2.62	yanagi-koku	110	2.38	1.30	sakura-me
6	3.38	2.57	yamu-makasu	—	—	—	—
7	3.38	2.56	mau-koku	155	2.17	1.04	chiru-katu
8	3.27	2.43	yanagi-mazu	156	2.17	1.04	bakari-sumi
9	3.26	2.42	sakura-yamu	157	2.16	1.03	maku-besi
10	3.25	2.40	minami-yamakaze	158	2.16	1.03	tatu-miku
101	2.40	1.33	uturu-komoru	159	2.16	1.03	tatu-tazumu
102	2.40	1.33	sakura-watausu	161	2.16	1.03	miyako-sakura
103	2.40	1.33	katu-nagara	162	2.16	1.02	kakusu-si
104	2.39	1.32	sakura-masi	163	2.14	1.00	yononaka-sakura
105	2.39	1.31	sakura-makasu	164	2.14	1.00	monou-sakura

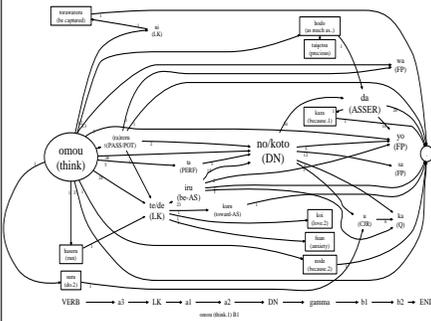


Figure 2: Construction of the predicate of *omou* (think) with Function Tail

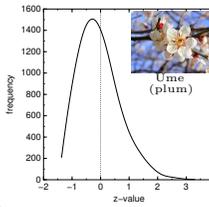


Figure 1: Bell curves

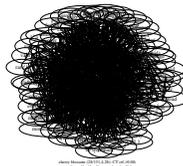


Figure 3: Hairball effect

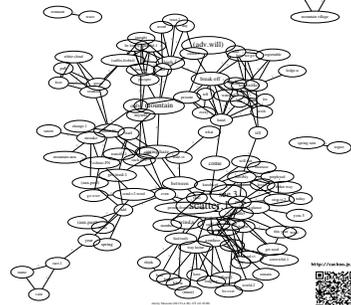


Figure 4: Only with Content Tail

## Conclusion

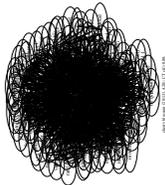
- 1) the distribution of classical texts fits a **Gaussian (Bell) curve** as well as in modern texts (Hodošček and Yamamoto 2013);
- 2) the  $cw$  value can separate patterns into three layers (low-, mid-, and high-range) using **inflection points** ( $-\sigma$  and  $\sigma$ );
- 3) of the three layers, the high-range could be extracted **without a list of stop words**;
- 4) the mid-range lexical layer might include mathematical traits not yet revealed in the present study.

## Reference

- Yamamoto, H. (2005), Visualisation of the construction of poetic vocabulary using the database of the *Kokinshū*, Jinbun kagaku to detabētōsu (Humanities and Database) the 11th symposium, 81–8, The council of humanities and database.
- Yamamoto, H. (2006), Extraction and Visualisation of the Connotation of Classical Japanese Poetic Vocabulary, Symposium for Computer and Humanities, vol. 2006, 21–28, The information processing society of Japan.
- Hodošček, B. and H. Yamamoto (2013) “Analysis and Application of Midrange Terms of Modern Japanese”, in *Computer and Humanities 2013 Symposium Proceedings*, No. 4, pp. 21–26.

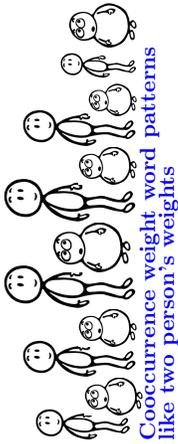
## Problem

Hairball Effect!!!



## Tall Head

Word distribution by frequency-rank



## Solution

Cooccurrence Weight

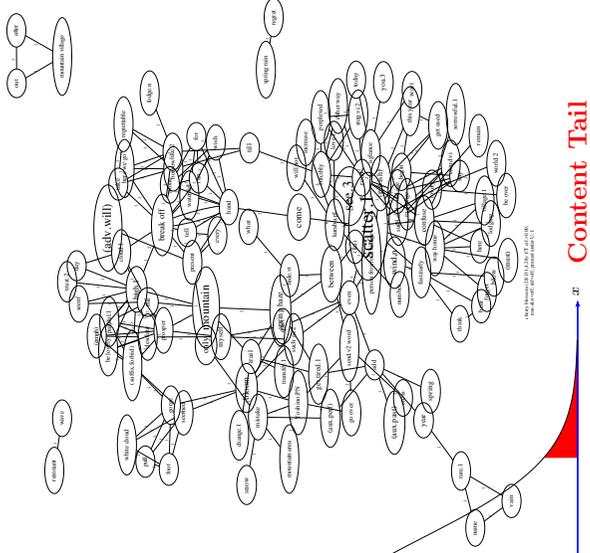
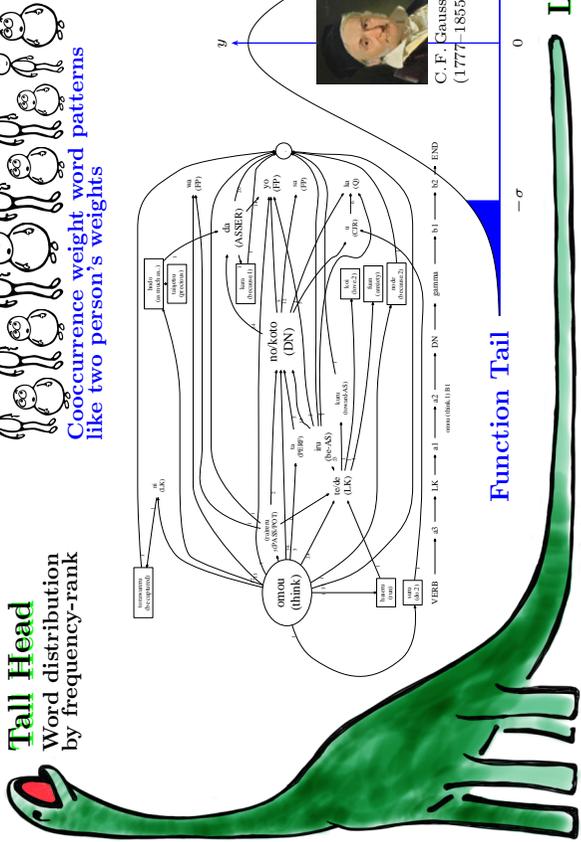
⇒ Bell curve!

It allows us to divide words into 3 kinds: ⇒ high, middle, and low.

## See what's happen!

Over  $1\sigma$  ⇒ **Content Tail**

Under  $-1\sigma$  ⇒ **Function Tail**



← more generic words

more specific words →

UCLA EAST ASIAN LIBRARY PRESENTS 2 JAPANESE DH LECTURES—

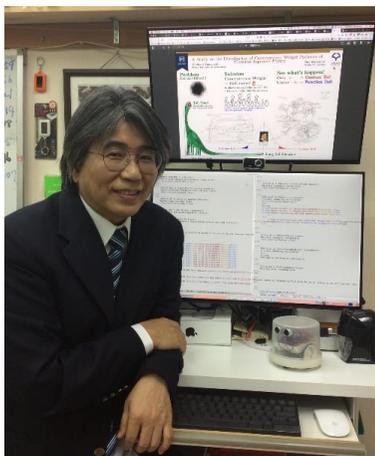
Free &amp; Open to the Public

**NOV. 27, 2018 (Tues)****11:00 am – 12:00 pm****3:30 pm – 4:30 pm**

Presentation Room (YRL 11348)  
in the UCLA Research Library  
(Charles E. Young)

Register online: <http://bit.ly/2OSTDMY>

Any questions, contact [tbialock {at} library.ucla.edu](mailto:tbialock@library.ucla.edu)



**Hilofumi YAMAMOTO** is

a professor of Linguistics at Tokyo Institute of Technology. He applies mathematical approach to historical changes of language. Please find more information about his research from his website:

<https://cuckoo.is.ila.titech.ac.jp/~yamagen/index-e.html>.

**11:00 am – 12:00 pm**

### Implications of Digital Humanities Research on Japanese Language and Literature Studies: Balancing between Technology and Research Inquiry

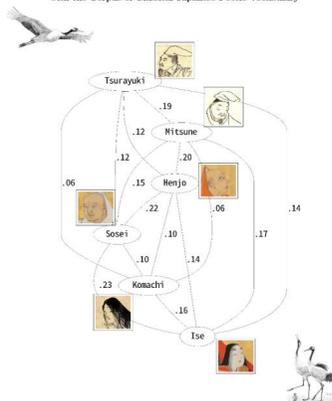
Critical issues in conducting digital humanities research are addressed. The presenter will emphasize the importance of balancing between the use and choice of suitable technology and the formulation of research objectives appropriate to technology in digital humanities research. Based on his corpus study (on the *Hachidaishū*) he will suggest a model of digital humanities approach for Japanese language and literature studies.

**3:30 pm – 4:30 pm**

### The Future Research Framework of Corpus Research Contributing to Language Learning/Education

In the second lecture the presenter will introduce a framework of corpus research for future language acquisitions and language education of Japanese exploring a variety of digital data, tools, and preservation infrastructures available on the Web.

UCLA EAST ASIAN LIBRARY      UCLA 2018 Midsumo Symposium      Nov. 27, 2018  
Mathematics, Literature, and Historical Linguistics  
with the Corpus of Classical Japanese Poetic Vocabulary



SPONSORED BY: EAST ASIAN LIBRARY (RICHARD C. RUDOLPH), TOKYO INSTITUTE OF TECHNOLOGY, AND UCLA LIBRARY

# An Analysis of the Differences Between Classical and Contemporary Poetic Vocabulary of the Kokinshū

Hilofumi Yamamoto  
Tokyo Institute of Technology

Bor Hodošček  
Osaka University

2019.5.7

## 1 Introduction

The purpose of the current project is to clarify the relationship between literal (or written) elements and non-literal elements (or connotation) of an ancient language. We will first clarify the differences between the original ancient language and modern language translations of poems in the same literary work, the Kokinshū. In particular, we will examine whether the translations of the Kokinshū use the same words as in a poem (or words corresponding to the modern language) or whether they use words not corresponding to words in a poem. To specify elements written only in the translations, we subtract the elements of original poems (OP: the Kokinshū) from the elements in their contemporary translations (CT), and analyze the residual elements. The differences, therefore, may include two kinds of elements: 1) elements resulting from chronological differences in language; 2) elements added for interpretation. We will subtract the elements of OP from those of CT to account for these differences. While similar attempts have done the subtraction processing manually for the analysis of modern language (Miyazima 1979, 1980, Suzuki 1988, Hasumi 1991), this is the first attempt to subtract a set of linguistic elements from another set by the computer.

## 2 Methods

We will use the corpus of the Kokinshū by Nakamura et al. (1999). As shown in Figure 1, the poems and the translations are stored as corpora databases

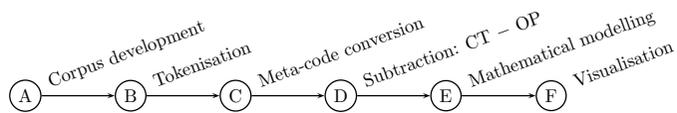


Figure 1: Flowchart of data processing

Table 1: Summary of the contemporary Japanese translations

translation work (year)	pages	manuscript	method
Kaneko (1933)	1105	Teika	word-for-word
Kubota (1960)	1449	Teika	word-for-word
Matsuda (1968)	1998	Teika	not mentioned
Ozawa (1971)	544	Teika	wording changed
Takeoka (1976)	2278	Teika	word-for-word
Okumura (1978)	434	Teika	intention oriented
Kyūsojin (1979)	1260	Teika	words added
Komachiya (1982)	407	Teika	not mentioned
Kojima and Arai (1989)	483	Teika	not mentioned
Katagiri (1998)	3022	Teika	word-for-word

and both of them are separated into tokens using the classical poem tokenizer, `kh` (Yamamoto 2007). We convert the tokens into meta-codes, then using the meta-codes, subtract the elements of the original from the elements of their translations. We examine the length of the portion of meta-codes between the two elements. (Figure 2) As an algorithm for matching the elements of CT and OP, we use Longest Common Subsequence (Traum and Habash 2000). An example of subtraction processing with `code2match.c` (Yamamoto 2005, 2009) is shown in Figure 3.

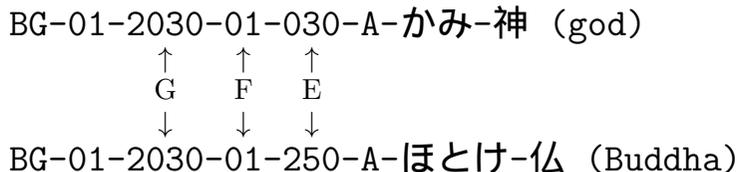


Figure 2: Level of matching elements: group matching (G); field matching (F); exact matching (E); each level is evaluated by the length of corresponding characters of meta-codes from the first letter.

```

+----- pair No.
| +----- value of matching level, exact=17, field=13, group=10
| | +--- POS No.
| | | OP element No.+      +- CT element No.
| | | OP element +      | CT element
| | | | |
1 17 11 *tatsutahime 00 <-> 12 *Tatsutahime (pn.Tatsutahime)
2 17 47      te 04 <-> 25 te (hand)
3 17 47      mukueru 05 <-> 26 mukueru (toward)
4 17 2      kami 06 <-> 32 kami (god)
5 10 61      no 07 <-> 33 ga (SUB)
6 17 47      ari 08 <-> 34 aru (be)
7 10 64      ba 09 <-> 35 kara (because)
8 17 65      koso 11 <-> 36 koso (EM)
9 17 2      aki 12 <-> 38 aki (autumn)
10 17 71     no 13 <-> 39 no (CON)
11 17 2     konoha 14 <-> 40 konoha (leaf of tree)
12 17 2     nusa 19 <-> 45 nusa (present)
13 17 61     to 20 <-> 46 to (CRD)
14 17 47     chiru 21 <-> 49 chiru (fall)
15 13 74     ramu 22 <-> 54 u (CJR)

```

Figure 3: An example of the alignment of the matched elements between OP(298) and CT(298, koma). Each line consists of the matched pair ID number (1), the matching level indicated by the value (17), ID number of POS (11) which indicates a place name, OP element (\*tatsutahime), ID number of OP element, ID number of CT element, CT element (\*Tatsutahime), and the glossary; \* written in different kanji.

Table 2: Result of subtracting the elements of  $OP_{(298)}$  from those of  $CT_{(298, koma)}$ : it indicates the ratio of the ingredients of  $OP_{(298)}$ .

OP (valid number of element)	= 16
E (ratio of exact match)	12/16 = 0.750
F (ratio of field match)	1/16 = 0.062
G (ratio of group match)	2/16 = 0.125
T (ratio of total match)	15/16 = 0.938
U (ratio of unmatched OP)	1 - T = 0.062

### 3 Results

Table 2 indicates a calculation of the components of  $OP_{(298)}$ .  $OP_{(298)}$  refers to a poem by Prince Kanemi.  $CT_{(298, koma)}$ , in turn, refers to the translation of the 298 poem by Teruhiko Komachiya in 1982. 12 elements out of 16 (75 percent) are matched in  $CT_{(298, koma)}$ . One element out of 16 is matched at the field level, and two elements out of 16 are matched at the group level in  $CT_{(298, koma)}$ . One element of  $OP_{(298)}$  does not match the elements of  $CT_{(298, koma)}$ . If we assume that matched elements at all the three levels express in  $CT_{(298, koma)}$ , then 15 elements (94 percent) of  $OP_{(298)}$  express as the elements in  $CT_{(298, koma)}$ . If we assume that matched elements at all the three levels are expressed in  $CT_{(298, koma)}$ , then 15 elements (94 percent) of  $OP_{(298)}$  are expressed as the elements in  $CT_{(298, koma)}$ . The remaining 6 percent of elements of  $OP_{(298)}$  do not match against any elements in  $CT_{(298, koma)}$ . None of the ten modern language translations could be fully expressed with the ancient language. The amount of added information was 80 percent higher than the original.(Table 4) The differences between the theoretical and experimental values were at most 8 percent. Those were rare cases, and in general accounted for around 4 percent.

### 4 Discussion

Based on the analysis of the differences between the two, we assume that translators attempted to express some cultural elements unfamiliar to modern people. Table 3 is an example of a calculation which indicates the components of  $CT_{(298, koma)}$ .  $CT_{(298, koma)}$  uses the same 12 elements as  $OP_{(298)}$ . The total number of elements of  $CT_{(298, koma)}$  is 41; thus 29 percent of  $CT_{(298, koma)}$  is calculated as the component of  $OP_{(298)}$ . The rest of  $CT_{(298, koma)}$ , 71 percent, is considered as added annotated text. Ratio A, however, does not consist only of newly added components: it should be deconstructed into

Table 3: Component of CT in case of KKS 298 by Komachiya (1982):  $\text{fabs}(D-H)$  stands for the function of the absolute value of the practical value, D, minus the theoretical value, H.

CT	(valid number of element)	= 41
W	(ratio of original word use)	$12/41 = 0.293$ (E/CT)
A	(ratio of annotation)	$1-0.293 = 0.707$ (1-W)
---breakdown of the annotation---		
P1	(ratio of FG paraphrased)	$(0.62+0.12)/0.707 = 0.073$ (F+G)/A
P2	(ratio of U paraphrased)	$(0.707-0.073)*0.062 = 0.040$ (A-P1)*U
D	(ratio of purely added)	$0.707-(0.073+0.040) = 0.595$ A-(P1+P2)
H	(theoretical value of D)	$1-16/41 = 0.610$ 1-OP/CT
Gap		$\text{fabs}(0.595-0.610) = 0.015$ $\text{fabs}(D-H)$

Table 4: Amount of added information (N=1000)

translator	alignment			subtraction		
	min.	mean (SD)	max.	min.	mean (SD)	max.
1 Kaneko	0.16	0.53 (0.09)	0.80	0.18	0.49 (0.09)	0.73
2 Katagiri	0.21	0.49 (0.08)	0.71	0.16	0.44 (0.08)	0.68
3 Kojima Arai	0.15	0.46 (0.09)	0.74	0.10	0.41 (0.10)	0.69
4 Komachiya	0.12	0.44 (0.08)	0.72	0.11	0.39 (0.08)	0.67
5 Kubota	0.15	0.45 (0.09)	0.77	0.13	0.40 (0.09)	0.72
6 Kyusojin	0.10	0.47 (0.08)	0.73	0.11	0.42 (0.08)	0.69
7 Matsuda	0.00	0.44 (0.09)	0.77	0.07	0.39 (0.09)	0.69
8 Okumura	0.06	0.44 (0.08)	0.75	0.11	0.41 (0.08)	0.72
9 Ozawa	0.10	0.46 (0.08)	0.72	0.20	0.44 (0.07)	0.70
10 Takeoka	0.11	0.42 (0.10)	0.74	0.06	0.38 (0.10)	0.69
mean	0.12	0.46 (0.03)	0.74	0.12	0.42 (0.03)	0.70

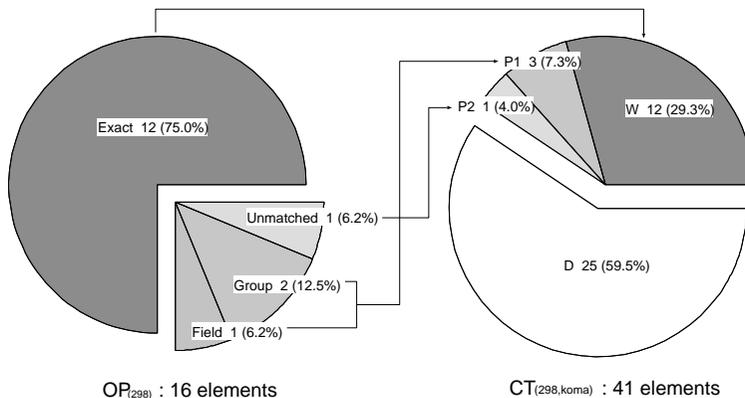


Figure 4: Pie-charts illustrating the components of  $OP_{(298)}$  and  $CT_{(298, koma)}$ : the ratio of purely added components is estimated based on the number of elements in common in  $OP$  and  $CT$ .

three kinds of components: 1) the first level of the paraphrased component,  $P1$ , which can be estimated from the ratio of the elements of the field match  $F$  and the group match  $G$ ; 2) the second level of the paraphrased component,  $P2$ , which can be estimated from the ratio of the unmatched elements, since even unmatched elements are assumed to be somehow translated into  $CT$ ; and 3) the purely added component,  $D$ , which can be estimated from the ratio of the annotation minus  $P1$  and  $P2$ .(Figure 4)

If the estimation from the subtraction of elements of  $OP$  from those of  $CT$  is correct, the practical value,  $D$ , can be close to the theoretical value,  $H$ , and the validity of the operation will be supported. In the case of the values between  $OP_{(298)}$  and  $CT_{(298, koma)}$ , the theoretical value is 0.610, the practical value is 0.595, and their discrepancy is 0.015, which means the two values are very close.

## 5 Conclusion

The current paper discussed the differences between the original poems of the Kokinshū and its translations. We attempted to classify the components of both  $OP$  and  $CT$  to examine whether or not  $CT$  includes added elements, which are the non-literal elements of  $OP$ . After subtracting the matched elements between  $OP$  and  $CT$  from  $CT$ , the presence of a residual indicated that  $CT$  includes newly added elements. It shows that it is impossible to convert the contents in the ancient language into only their equivalents in the modern language.

## References

- Hasumi, Yoko (1991) “Dōitsu jōhō ni motozuku bunshōhyōgen ni tsuite no bunseki / Difference of expressions on the same information”, *Mathematical Linguistics*, Vol. 18, No. 3, pp. 136–144.
- Kaneko, Motoomi (1933) *Kokinwakashū Hyōshaku: Shōwa Shimban*, Tokyo: Meiji-shoin.
- Katagiri, Yoichi (1998) *Kokinwakashū Hyōshaku Jō, Chū, Ge*, Tokyo: Kodansha.
- Kojima, Noriyuki and Eizō Arai (1989) *Kokinwakashū*, Vol. 5 of Shin-Nihon bungaku taikai (A new collection of Japanese literature), Tokyo: Iwanami shoten.
- Komachiya, Teruhiko (1982) *Gendaigo yaku taishō Kokinwakashū (Kokinwakashū with modern Japanese translations)*, Obunsha Bunko Taiyaku Koten Series, Tokyo: Ōbunsha.
- Kubota, Utsubo (1960) *Kokinwakashū Hyōshaku (Vol. 1, 2, 3)*, Tokyo: Tokyodo shuppan.
- Kyūsojin, Hitaku (1979) *Kokinwakashū Zen'yaku Chū (Comprehensive annotations of the Kokinwakashū)*, Vol. 1–5 of Kodansha Gakujutsu Bunko: Kodansha.
- Matsuda, Takeo (1968) *Shinshaku Kokinwakashū Vols.1 and 2*, Tokyo: Kazama Shobo.
- Miyazima, Tatuō (1979) ““Kyōsantō Sengen” no yakugo (Translated terms in the “Communist Manifesto”)”, in Gengogaku Kenkyūkai ed. *Gengo no Kenkyū (Study of language)*, Tokyo: Mugi Shobo, pp. 425–517.
- (1980) ““Jodōshi” to ‘Hojodōshi’ (Auxiliary verbs and subsidiary verbs)”, in Society of Modern Language ed. *Kindaigo kenkyū (Study of contemporary vocabulary)*, Vol. 6, Tokyo: Musashino shoin, pp. 455–468.
- Nakamura, Yasuo, Yoshihiko Tachikawa, and Mayuko Sugita (1999) *Kokubungaku kenkyū shiryōkan dētabēsu koten korekushon “Nijūichidaishū” Shōhobanbon CD-ROM (Database Collection by National Institute of Japanese Literature “Nijūichidaishū” the Shōho edition CD-ROM)*: Iwanami Shoten.
- Okumura, Tsuneya (1978) *Kokinwakashū*, Shinchō Nihon Koten Shūsei, Tokyo: Shinchō sha.
- Ozawa, Masao (1971) *Kokinwakashū*, Vol. 7 of Nihon Koten Bungaku Zenshū, Tokyo: Shōgakkan.

- Suzuki, Tai (1988) “Weirando “Shūshinron” no Kanji (Kanji in the “Elements of moral science” by Francis Wayland)”, in *Gengo no Kenkyū (Study of language)*, Vol. 8 of *Kindai Nihongo to Kanji (Contemporary Japanese and Kanji)*, Tokyo: Meijishoin, pp. 128–164.
- Takeoka, Masao (1976) *Kokinwakashū Zen Hyōshaku Jō Ge (the complete annotated edition of Kokinwakashū, Vols. 1 and 2)*, Tokyo: Yubun Shoin.
- Traum, David and Nizar Habash (2000) “Generation from lexical conceptual structures”, in *NAACL-ANLP 2000 Workshop on Applied interlinguas*, pp. 52–59, Morristown, NJ, USA: Association for Computational Linguistics.
- Yamamoto, Hirofumi Hilo (2005) “A Mathematical Analysis of the Connotations of Classical Japanese Poetic Vocabulary”, Ph.D. dissertation, Australian National University.
- Yamamoto, Hilofumi (2007) “Waka no tame no Hinshi tagu zuke shisutemu / POS tagger for Classical Japanese Poems”, *Nihongo no Kenkyu / Studies in the Japanese Language*, Vol. 3, No. 3, pp. 33–39.
- (2009) “Thesaurus for the Hachidaishu (ca.905–1205) with the classification codes based on semantic principles”, *Nihongo no Kenkyu / Studies in the Japanese Language*, Vol. 5, No. 1, pp. 46–52.

## 第 8 章

# 情報処理学会じんもんこん

### 8.1 概要

### 8.2 発表年

2006, 2007, 2008, 2010, 2011, 2012, 2013, 2014, 2016, 2017

2017 年、ベストポスター賞を共同研究者ボル・ホドシチェク氏（大阪大学）が受賞する。

### 8.3 論文・ポスター・スライド



# 二十一代集シソーラスのための漸近的語彙対応システムの開発

山元 啓史 †\*      村井 源 †      ホドシチェク ボル ‡  
 †東京工業大学   \* カリフォルニア大学サンディエゴ校   ‡ 大阪大学

## 目的

- 本プロジェクトの目的は、和歌用語のシソーラスを開発することである。
- 本研究では、そのための語彙の自動対応づけシステムの開発を行う。
- 和歌とその現代語訳 10 種を用いてパラレルコーパスを作成し、対応づけに用いる。
- 将来的には、これを利用して、客観的な基準による語彙の分類と言語変化の探究を行う。

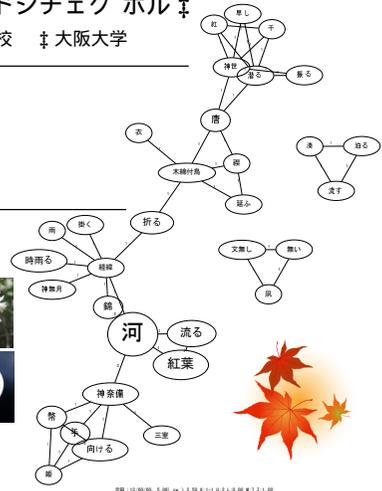
## 問題

### 和歌とは何か？



立田 姫  
 手向くる神の  
 あればこそ  
 秋の木の葉の  
 幣と散るらめ

1. 表記の問題  
 龍田、立田、竜田、たつた
2. 単位の問題  
 卯の花？ 卯の/花？
3. 所属の問題  
 「卯の花」は植物か食品か？
4. 掛詞の問題  
 「海松藻（みるめ）」と「見る目」？



## 方法

材料: 古今和歌集 (905 年頃) とその現代語訳 10 種

Mutual Co-occurrence Rate: 村井 (2010)

$$mcr(o, t) = p(o | t) p(t | o)$$

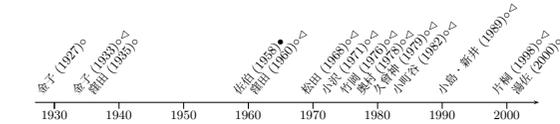


図 1: 古今集の注釈本の出版年。○は現代語訳が掲載されているもの。●は現代語訳が掲載されていないもの。▷は本研究で用いられたもの。

ただし、 $o$  はオリジナル (和歌) の単語。 $t$  は現代語訳の単語。 $mcr(o, t)$  は相互共出現率。 $p(o | t)$  は、オリジナルと現代語訳の対応する 2 つの文に注目した時、単語  $o$  と単語  $t$  が同時に出現した割合。  
 →  $mcr(o, t)$  の値が、十分に大きい時、単語  $o$  と単語  $t$  は、**文脈的に一致している**と推定。

## 結果

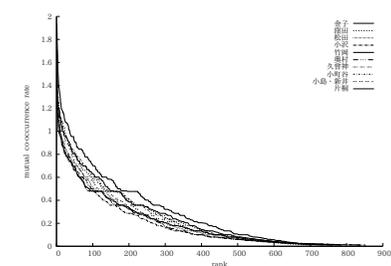


図 2: 古今集と現代語訳 10 種中の任意の単語対の mcr 値の分布。

表 1: 「古今和歌集」とその現代語訳 10 種の単語の対応処理を実施し、その推定値の上位、中位、下位ごとの単語の対応例

no.	上位対 (1.3 以上)	中位対 (0.16 より)	下位対 (0.01 以下)
1	鳴く 鳴く	老ゆ 年老いる	異なり あの
2	風 風	乱る 乱れる	雫 どうして
3	世の中 世の中	来 いらっしゃる	此の この
4	人 人	問ふ 問う	随に まま
5	春 春	問ふ 訪ねる	匂ふ 美しい
6	秋 秋	名 噂	見る せい
7	時鳥 時鳥	変はる 変る	連れ つく
8	時鳥 ほととぎす	燃ゆ 燃える	立ち返る 言う
9	散る 散る	濡づ 濡れる	有り つく
10	見る 見る	難し むずかしい	有り まさしく

## 考察

1. 動詞: 「落つ 落ちる (自動詞)」 「捨つ 捨てる (他動詞)」
2. 名詞: 「木綿付鳥 (ゆふつけどり) 鶏」「朝な朝な 毎朝」
3. 古語固有: 機械的に、古語「名」に対して「噂」が推定できた。
4. 推定範囲: 単語対の推定が可能なのは 0.2 あたりまで。

## おわりに

和歌・現代語訳のパラレルコーパスを利用した単語対の推定法により、新規シソーラスに追加すべき語の決定はほぼ計算手続きによってできるだけでなく、比喩、含意、掛詞などの和歌特有の対応づけも可能なことがわかった。

## 文献

- 村井源 (2010) 「漸近的対応語彙推定法に基づく翻訳文の解釈的特徴の抽出: 日本語翻訳聖書の計量的比較」, 『情報知識学会誌』, 第 20 巻, 第 3 号, 293-310.

# 歌ことば「橘」「梅」「桜」における関連対の抽出

ホドシチェック ボル  
大阪大学  
bor@lang.osaka-u.ac.jp

山元 啓史  
東京工業大学  
yamagen@ila.titech.ac.jp

## 材料

八代集 データベース 山元 (2007) で下記資料を収集 シソーラスに異なる表記を登録、正規化した

- ・新編国歌大観CD-ROM版の二十一代集に相当するデータ
- ・国文学研究資料館編集二十一代集データベース (中村樸 1999)
- ・新日本古典文学大系本二十一代集に相当する書籍
- ・新潮日本古典集成の新古今集 (5,6冊 1979)
- ・ヴァージニア大学日本語テキストインシアティブ監修 (http://www.iaj.org/iaj-tdb/tdb/tdbmain.html)

## 方法



「主役語」  
八代集  
データベース



ネットワークコミュニティ検出  
・Rでlinkcommを使用  
・類似度の高い(類似度)の抽出  
・ノードが2つ以上のコミュニティに所属できる  
・コミュニティクラスタ

表1: 橘のサブクラスター1番から10番までを抽出。average法, mcquity法, single法のクラスターングを用いた。括弧内はそれぞれMaximum partition densityを示す。

No.	node	edge	node	edge	node	edge
1	昔	7	昔	7	昔	5
2	心ふ	6	心ふ	6	心ふ	4
3	風	5	今年	4	夢	4
4	夢	5	辺	4	香る	3
5	今年	4	待つ	4	今年	3
6	辺	4	風	4	初む	3
7	待つ	4	夢	4	五月雨	3
8	香る	3	初む	3	折	3
9	五月雨	3	香	3	折	3
10	初む	3	間	3	思寝	3

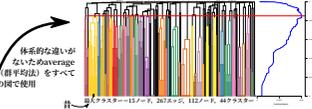


表2: 「梅」のサブクラスターから抽出された「顔し」「恋」「折」に関する歌

No.	歌番号	作者	歌
1	2000	藤原	顔しにのほほほほの梅。色こそ見ゆ。ゆめゆめと
2	2002	藤原	もろく。おのほほほほ。梅。色こそ見ゆ。ゆめゆめと
3	2003	藤原	顔しにのほほほほの梅。色こそ見ゆ。ゆめゆめと
4	2004	藤原	顔しにのほほほほの梅。色こそ見ゆ。ゆめゆめと
5	2005	藤原	顔しにのほほほほの梅。色こそ見ゆ。ゆめゆめと
6	2006	藤原	顔しにのほほほほの梅。色こそ見ゆ。ゆめゆめと

## 考察

### 関連対

### 直接的⇔間接的

- ・連想しやすい関連対よりも間接的な関連対が存在する
- ・和歌の短い文脈において主役語が同時に現れない
- ・関連対と類似対

### 可視化からみた問題点

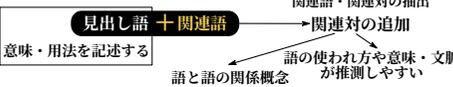
単語単位: 同じ文に使われた \* 連語が成立していた  
重複したエッジ: 同一語形のため多義語が一義的に扱われてしまう

### 結論

歌ことば辞典の開発にあたり、従来の「見出し語とその解説」による記述に加え、「見出し語—関連語」の形式による記述を提案し、「橘」「梅」「桜」のコミュニティ分析によって、関連対の抽出を行い、それらは実際の和歌において関連した語であることが確認できた。

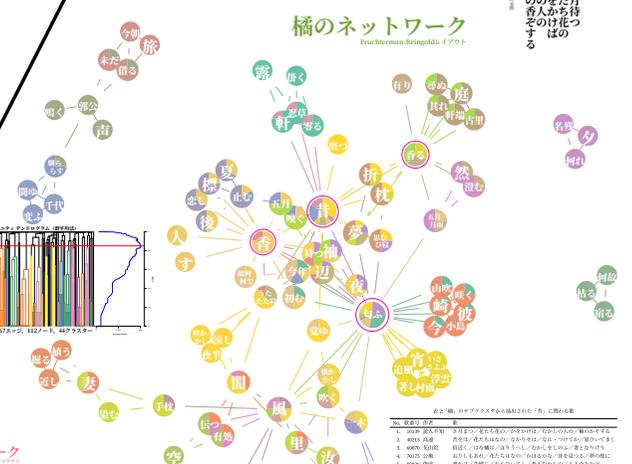
## 問題

過去のものとなってしまった古代語の意味記述  
従来の見出し語とその解説に加え...



## 結果

### 橘のネットワーク



### 梅のネットワーク

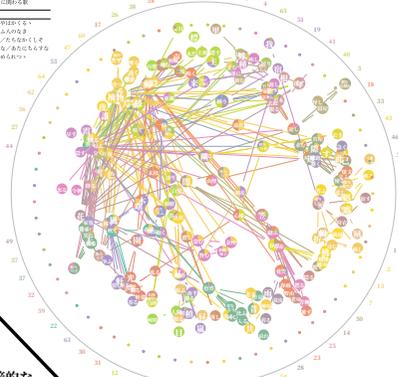


表3: 「桜」のサブクラスターから抽出された「高嶺 (1,2)」「曙 (2,3)」「春風 (4,5,6)」に関する歌

No.	歌番号	作者	歌
1	50036	公実	白雲とノをものぬねのノ見つるはノ心まとはノさくらなりけり
2	80113	後鳥羽院	みよし野のたかねのさくらノ散にけりノ風もしるノ春のあけほの
3	80114	後醍醐	又かみノつた野のみノ桜ノ花の香るノ春のあけほの
4	10058	賀正	たしなもノとてあまふつ。智恵ノちかむとてノ自の御を
5	30042	元輔	春かすみノちかへたてそノ花さかりノみてたにあかぬノさくらを

### 桜のネットワーク







# 八代集「桜の花」歌における作者の分類



山元 啓史 (東京工業大学) ホドシチエク ボル (大阪大学)

**意義:**

よみ人しらず歌は**すぐれた作品**ばかり。

**問題:**

1. よみ人しらず歌の作者はわからない。

2. 歌データは**1首31文字**のみである。

**目的:**

作者は特定できなくても、コミュニティ、つまり、

類似度により、誰の作品に近いかぐらいは推定したい。

**方法:**

トピックを桜に限定し、内容語も分析に含める。

機能語は大まかな一致、たとえば「なむ」「らむ」は同じ。

名前を伏せ、6歌人の相互関係を力学モデルで表現する。



1.

紀貫之



2.

凡河内躬恒



3.

遍照



4.

素性



5.

小野小町



6.

伊勢

**結果は後ほど...**



# 八代集「桜の花」歌における作者の分類

山元 啓 史 (東京工業大学)

ホドシチェク ボル (大阪大学)



大阪大学  
OSAKA UNIVERSITY

## はじめに

**意義** よみ人しらず歌はすぐれた作品ぞろいであり、これらの歌の分類、特徴推定は有意義である。

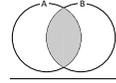
**問題** 1. よみ人しらず歌の作者はわからない。2. 歌データは1首31文字のみである。

**目的** 和歌の類似度によりコミュニティ(作品グループ)を推定できるかどうか。

## 方法

1. トピックを桜に限定<sup>1</sup>し、内容語も分析に含める<sup>2</sup>。
2. 機能語<sup>3</sup>は大きな一致、たとえば「なむ」「らむ」は同じ。
3. 有名6歌人による歌を材料として用い、類似度計算を行う。
4. 6歌人の相互関係をデンドログラムと力学モデル<sup>4</sup>で描画する。
5. モデルから、6歌人の相互関係が説明できるかどうかを考察する。

$$Jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$$



$$Dice(A, B) = \frac{2|A \cap B|}{|A| + |B|}$$

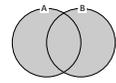


TABLE 1: 六歌人の桜歌: 小町は桜を意味する歌を選んだ。

No	生没年	歌番号	歌
1a.	紀貫之	古今 049	今年より / 春しりそむる / 桜花 / ちるといふことは / ならはさらなむ
b.	877-945	古今 058	たれしかも / とめておりつる / 春露 / たちかくすらん / 山の桜を
2a.	凡河内躬恒	古今 086	雪とのみ / ふるたにあるを / 桜花 / いかちれれとか / 風の吹らむ
b.	857-927	古今 358	山たかみ / 雲ふにみゆる / 桜花 / 心のゆきて / あらぬ日そなき
3a.	遍昭法師	古今 394	山風に / 桜吹まき / みたれなん / 花のまきせに / 立とまるへく
b.	816-890	古今 091	花の色は / 露にこめて / みせずとも / かをたにぬすめ / 春の山かせ
4a.	素性法師	古今 055	見てのみや / 人にかたらん / 桜花 / てことありて / いへつとにせむ
b.	???-910	古今 056	見わたせば / 柳桜を / こきませて / 都を春の / 跡なりける
5a.	小野小町	古今 113	花の色は / うつりにけりな / いたづらに / 我身世にふる / なかめせしまに
b.	???-???	古今 797	いるみえて / うつるふ物は / 世中の / 人の心の / 花にそありける
6a.	伊勢	古今 061	桜花 / 春くははれる / としたにも / 人の心に / あかれやはせぬ
b.	872-938	古今 068	みるひと / なき山嵐の / 桜はな / ほかのちりなん / のちそきかまし

Fig. 1: Methods of Similarity calculations

P1: 1 今年 2 より 3 春 4 知る 5 初む 6 桜 7 花 8 散る 9 と 10 言ふ  
11 こと 12 は 13 言ふ 14 ず 15 なむ..15 語

P2: 1 雪 2 と 3 のみ 4 ふる 5 だに 6 ある 7 を 8 桜 9 花 10 いか  
11 散る\*と 12 か 13 風 14 の 15 吹く 16 らむ..16 語

CS: 1. と、 2. 桜、桜 3. 花、花 4. 散る、散る 5. なむ、らむ..5 語

$$Jaccard(p1, p2) = \frac{5}{15 + 16 - 5} \quad (1)$$

$$= .19 \quad (2)$$

$$Dice(p1, p2) = \frac{2 \times 5}{15 + 16} \quad (3)$$

$$= .32 \quad (4)$$

## 結果

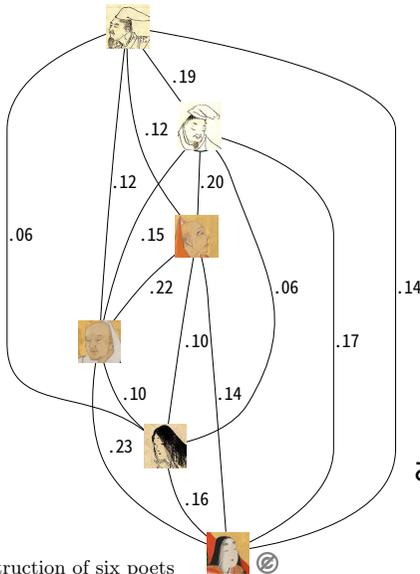


Fig. 2: Construction of six poets

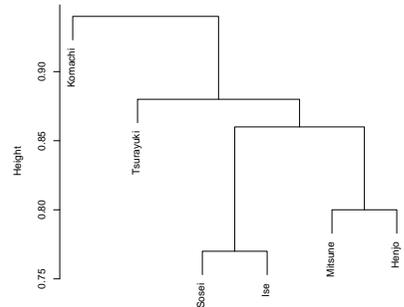


Fig. 3: Dendrogram of six poets; hclust, as.dist(d) complete.

## おわりに

いずれの類似度計算においても、局所的には素性と伊勢が最も近く、貫之と小町が最も遠くなった。力学モデルでは、古今集編者の2名(貫之、躬恒) 僧侶・親子の2名(遍昭、素性) 女性2名(小町、伊勢)の3つのグループにわかれた。

## 文献

- Ellson, J., E. R. Gansner, E. Koutsofos, S. C. North, and Gordon W. (2004) "Graphviz and Dynagraph: Static and Dynamic Graph Drawing Tools", in M. Jünger and P. Mutzel eds. *Graph Drawing Software*, Berlin Heidelberg New York: Springer-Verlag, pp. 127-148.
- Fortunato, Santo (2010) "Community detection in graphs", *Physics Reports*, Vol. 486, No. 3, pp. 75 - 174.
- Jaccard, Paul (1912) "The distribution of the flora in the alpine zone", *New Phytologist*, Vol. 11, p. 3750.
- Kamada, Tomihisa and Satoru Kawai (1989) "An algorithm for drawing general undirected graphs", *Information Processing Letters (Elsevier)*, Vol. 31, No. 1, p. 715.

## 第9章

# Language Classes and Book list

### 9.1 辞書・専門書

1. Japanese: *A Comprehensive Grammar*, Routledge, Tyler and Francis, 2012 (p. 77)
2. 通時コーパスによる言語の研究, 「コーパスと日本語史研究」ひつじ書房 2015 (p. 77)
3. 第4章コーパスから抽出した複合辞, シリーズ: 講座日本語コーパス7 「コーパスと辞書」朝倉書店 2019 (p. 77)
4. *A Gradual Approach to Technology Based Instruction, Learning Japanese in the Network Society*, University of Calgary Press March 2003 (p. 77)

### 9.2 日本語: 大学院

1. 大学院 (入門) テーマ別日本語 T51、Strategic Japanese (p. 77)
2. 大学院 (初級) テーマ別日本語 23,24、マイクロビット・ジャパニーズ (p. 77)
3. 大学院 (上級) 日本語「セミナーを開く」(p. 77)
4. 大学院 (上級) 日本語「プレスリリースを読む」(p. 77)
5. 大学院 (上級) 日本語「ニュースを聞く」(p. 77)

### 9.3 日本語: 学部

1. 学部 (上級) 日本語 Advanced Japanese (p. 77)
2. 教養特論 (学部): ライティング・スキル (p. 77)

## 9.4 日本文化関連

1. 学部: 日本文化演習「ことばの文化: とんちんかんのちんぷんかんぷん」(p. 77)
2. 学部: 日本文化演習「現代文化論と日本語」(p. 77)
3. テーマ別日本語 T51 レポート集 (p. 77)

## 9.5 言語学関連

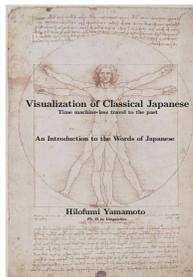
1. 学部: 言語と文化 (p. 78)
2. 学部: 言語学 A (p. 78)
3. 学部: 言語学 B (p. 78)
4. 大学院: 言語学特講 (p. 78)

## 9.6 MOOC TokyoTechX

1. Ten Sentences A Day for Eight Weeks: Dictation Every Day, Volume 1. (p. 78)
2. Ten Sentences A Day for Eight Weeks: Dictation Every Day, Volume 2. (p. 78)
3. Ten Sentences A Day for Eight Weeks: Dictation Every Day, Volume 3. (p. 78)
4. Idiomatic Japanese: The Secret of Advanced Japanese, Volume 1. (p. 78)
5. Idiomatic Japanese: The Secret of Advanced Japanese, Volume 2. (p. 78)
6. Idiomatic Japanese: The Secret of Advanced Japanese, Volume 3. (p. 78)
7. Kanji Vocabulary: The Secret of Advanced Japanese, Volume 1. (p. 78)
8. The RONBUN SHINAN: Writing The Basics (p. 78)
9. 科学者の人生: 彼らはどう自分の人生を生きたか (p. 78)
10. 日本語読本: これから言語学を楽しむ皆様のために (p. 78)

## 9.7 その他

1. 目で見てわかる今の日本語・昔の日本語: 中学生版 (ワークブック)(p. 78)
2. 目で見てわかる今の日本語・昔の日本語: 中学生版 (英語版)(p. 78)
3. 言語学ノート「ことばのおと」(p. 78)
4. 目で見てわかる今の日本語・昔の日本語: 小学生版 (ワークブック)(p. 78)
5. 創造性育成プログラムポスター: 2016年(p. 80), 2017年(p. 81), 2018年(p. 82)


<p>教養特論: 言語と文化 にほんご ワークブック</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>	<p>Linguistics: Truth of Daily Language An Introduction to the Study of Words</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>言語学: ことばの研究史 私たちが使う日本語のルーツ 授業の進め方とテキストシート</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>	<p>2nd Edition Diachronic Linguistics Language, Time, and Space</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>
<p>Ten Sentences A Day for Eight Weeks Dictation Everyday Volume 1 Let's learn Japanese through Listening Practices</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>Ten Sentences A Day for Eight Weeks Dictation Everyday Volume 2 Let's learn Japanese through Listening Practices</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>Ten Sentences A Day for Eight Weeks Dictation Everyday Volume 3 Let's learn Japanese through Listening Practices</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>論文指導: 書き方の道筋</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>
<p>Idiomatic Japanese The Secret of Advanced Japanese Volume 1 Let's learn Japanese through Listening of Japanese</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>Idiomatic Japanese The Secret of Advanced Japanese Volume 2 Let's learn Japanese through Listening of Japanese</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>Idiomatic Japanese The Secret of Advanced Japanese Volume 3 Let's learn Japanese through Listening of Japanese</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>Kanji Vocabulary The Secret of Advanced Japanese Volume 1 Let's learn Kanji through Reading of Kanji</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>
<p>科学者の人生 科学者としての人生を学ぶ ワークブック</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>	<p>日本語読本 二千年の日本語を学ぶための読本 第1巻 ワークブック</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>		
<p>目で見てわかる昔の日本語と今の日本語 タイムマシンに乗ってわかる昔の日本語 ワークブック</p>  <p>山光 啓史 Hi. Sh. in Japanese</p>	<p>Visualization of Classical Japanese Time machine has opened to the past An Introduction to the Words of Japanese</p>  <p>Hilofumi Yamamoto Hi. Sh. in Japanese</p>	<p>言語と文化 ことばのおと 山光 啓史 Hi. Sh. in Japanese</p> 	<p>目で見てわかる昔の日本語と今の日本語 タイムマシンに乗ってわかる昔の日本語 ワークブック 2018.12.26 山光 啓史 Hi. Sh. in Japanese</p> 



平成28年度 創造性育成科目

# 「言語学A・B・C」・「言語と文化」

平川八尋・山元啓史・赤間啓之 (リベラルアーツ研究教育院)



## 授業の目的: 言語を客観的に眺める!

**言語学の流儀:** 日常言語の疑問を教師が問いかける!  
総勢80名によるディスカッション!  
1人1枚全80枚のポスター大会!

手書きの  
ポスター  
迫力抜群

言語学の  
各論を  
ディスカ  
ッションで

それぞれの「...ている」の意味が異なります。  
どう異なり、なぜそうなるのかを考えてください。  
(1) 太郎は 今 走っている  
(2) この釘は 曲 がっている

Q: 英語モードのSiriで  
「掘った羊、いじるな!」  
といったら、どうなる?

### 「言語学とはどんな科学か」

1. ウィトゲンシュタインと言語ゲーム
2. チューリングとチューリングマシン
3. ジップとジップの法則
4. ダニエル・ジョーンズの18の基本母音
5. ソシュールと記号論
6. フィルモアと格文法
7. シュラムとコミュニケーション理論

### 「言語学と言語の仕組み」

8. 音声 (連濁)
9. 語形成 1
10. 語形成 2
11. 統語
12. 意味
13. 母語習得

[ポスター発表] 「はじめての言語研究」

### ミスチルの歌詞

には掛詞や韻が  
たくさん使われている。  
音声学・音韻論  
の立場からそれらを  
分析せよ。

あなたは言語のやりとりを見て育ちました。  
何を見てことばが話せるようになったのでしょうか。

言語学は  
なぜ始まった?

学生の活動: 自分で考えてディスカッション!  
教員の役割: 考える場を作る! 答えはない!

今週の日曜日に何をしたかを  
ことばを使わずに、グループの  
メンバーに伝えてみなさい。

この授業(言語学)では、「言語的直感」をもたらす知識を「文法」とよびます。いままでのように、現代国語や古典の文法での「暗記」とは異なります。

先生も  
学生と  
いっしょに悩  
む

時には、先生同士が  
疑問をぶつけあう、  
ツッコミを入れるバトルあり。

先生にも  
わからないこと  
があるんだ!



# 言語学：自分のテーマで考える仕組み

東京工業大学  
リベラルアーツ  
研究教育院

言語と文化・言語学A・言語学B 担当：赤間啓之・平川八尋・山元啓史

平成29年度教育推進室創造性育成科目委員会平成28年度教育革新センターの助成を得た。

## 概要

学生が自主的におこなったことは必ず**フィードバック**が得られ、自分自身の行動が**モニター**できるようにすることが不可欠である。しかし、大人数クラスではフィードバックは難しい。

大人数クラスであっても**ポスター発表**を可能にし、そのポスターの管理、コメントの整理、フィードバックの送付を自動的に行い、内容の充実を促進するシステムを開発した。例年、創造性育成科目に指定されている言語学は昨年度は80名を超え、今年度は115名となった。このシステムにより、昨年度は**80枚のポスター管理**、**80名の学生からの80名へのコメント返却が自動的に**行えるようになった。

## 実装

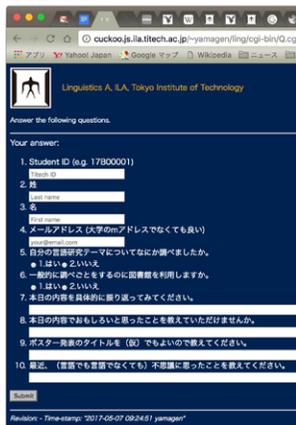


図1 通常授業の振り返り入力

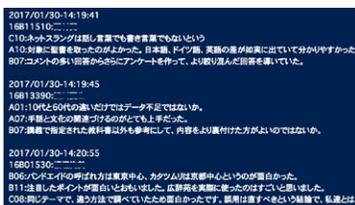


図4 コメント・フィードバック



図2 振り返りの出力・ポスターテーマ

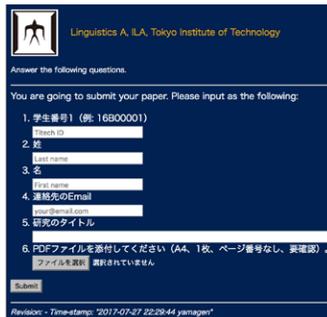


図5 ポスターの提出

2017.07.27: 本日の宿題のリンクはありません。

さて、本日のポスター発表2日目です。Cグループ、Dグループの発表です。たのしくやりましょう。

### ポスター発表会

1. ポスター発表の時のポスター（大きさはA3縦1枚）です。板ダンボールはこちらで用意します。
2. ポスター発表の時のポスターテンプレートは [pptx版](#) [pdf版](#) でこれにたがって作ること。
3. ポスターの内容、論文はインターネットからの引用はできません。
4. ポスターでのデータは自分で作ったデータに限られます。
5. ポスター発表は朝7時20分(7月24日)と7月27日(日)を使って実施します。
6. ポスター発表では、発表者だけでなく聴衆をしてのみさんのコメントも単位必須条件となりますので、ご注意ください。
7. 発表は2分くらいで、結論(成果)、明らかにした仕組み(仕組み)、意義の順で話す短い時間で議論を入れるでしょう。
8. 聴衆のみさんは、1グループ5作の発表を選んで、コメントと質問をGoogle Formで提出してください。
9. コメント・質問にしたがって、ポスターの内容を修正し、ポスターをA4のpdfでアップロードしてください。
10. アップロードしたポスターと質疑・出席・ディスカッションの内容を整合して成績付けます。

図3 ポスター発表実施方法の説明



図6 ポスター・リストと冊子生成

## 結果

ポスター実施システムを利用した結果、次のような3種類の教師・学生間、学生・学生間のインタラクションが見られた。**1. 教員への学生のコメント**（通常授業）、**2. ポスター発表者への学生・教員からのコメントと質問**（ポスター発表時）、**3. 学生から学生へのコメント**（ポスター発表後）。また、ポスター冊子をPDF形式で配布したことも振り返りの機会となった。ポスター冊子とは**修正版ポスターをWeb投稿により回収し、それらを自動的に編集（目次他）処理したものである。**

## 考察

昨年度の約80名、今年度約60名のクラスにしては、多くの意見交換を行うシステムができた。このシステムは言語学だけでなく、**ライティングスキル**においても利用されており、言語学以外のクラスでも利用できる。ただし、現時点のシステムを利用するには管理者による設定が不可欠で、広く使うには、設定、マニュアル、運営上の事情による柔軟な変更・追加機能が必要である。これらの問題はシステム開発によってだけでなく、授業の目的や進め方などとも同時に協調的に改善されるべきである。

## 結語

大人数クラスでの学生自身・学生相互の活動量を増やせた。ポスターに対するコメントが個々の学生に速やかにフィードバックできた。授業のみよりも学生のレポートの質も上がった。



# 「言語学A・B・C」「言語と文化」



言語学に  
数学を使う?  
もちろん!

問 英語モードのSiriで  
掘った芋、いじるな!  
といったら、どうなる?

授業の目的: **言語を客観的に眺める!**

言語学の流儀: **日常言語の疑問を教師が問いかける!**  
**学生全員によるディスカッション! ポスター大会!**

それぞれの「...ている」の意味が異なります。  
どう異なり、なぜそうなるのかを考えてください。  
(1) 太郎は今走っている  
(2) この釘は曲がっている

言語学の  
不思議を  
ディスカッション

### 「言語学とはどんな科学か」

1. ウィトゲンシュタインと言語ゲーム
2. チューリングとチューリングマシン
3. ジップとジップの法則
4. ダニエル・ジョーンズの18の基本母音
5. ソシュールと記号論
6. フィルモアと格文法
7. シュラムとコミュニケーション理論

### ミステルの歌詞

には掛詞や韻が  
たくさん使われている。  
音声学・音韻論  
の立場からそれらを  
分析せよ。

### 「言語学と言語の仕組み」

8. 音声 (連濁)
9. 語形成 1
10. 語形成 2
11. 統語
12. 意味
13. 母語習得

あなたは言語を明示的に教わるこ  
ともなく話せるようになりました。  
何を見てことばが話せるように  
なったのでしょうか。

「ホッピー」はホップの  
ビールのような飲料な  
のに、なぜ「ホッピー」  
と叫ぶのか?



言語学は  
なぜ始まった?



この授業(言語学)では、「言語的直感」  
をもたらす知識を「文法」とよびます。い  
ままでのように、現代国語や古典の文  
法での「暗記」とは異なります。

### ポスター発表「はじめての言語研究」

問いだらけ、というより、問いしか載っていない教科書

### 学生の活動

自分で考えて  
ディスカッション!  
教員の役割  
考える場を作る!  
答えはない!



時には、先生同士が  
疑問をぶつけあう、  
ツッコミを入れるバトルあり。

教師も  
学生と  
いっしょに悩む



ボーッとしていると  
マイクが飛んでくる!

先生にも  
わからないこと  
があるんだ!



## 裏表紙について



写真の日時計にはラテン語で“Carpe Diem”（カルペ・ディアム）と彫ってあります。英語では“Seize the day”、日本語では「その日を摘め」と訳されています。そこには「その日を楽しみ、精一杯いきること」という意味があります。紀元前1世紀の古代ローマの詩人ホラティウスの詩に登場する句で、映画“Dead Poets Society”（1989年、邦題「いまを生きる」ロビン・ウィリアムズ主演）にも出てきます。



クイントゥス・ホラティウス・フラックス  
Quintus Horatius Flaccus  
BC.65.12.8–BC.8.11.27  
古代ローマ時代の南イタリアの詩人

言語と文化、東京工業大学

山元研究室カタログ

2018年 4月 12日 第1版

2018年 10月 15日 第2版

2018年 12月 26日 第3版

2019年 1月 1日 第4版

著者: 山元啓史

©2018, Hilofumi Yamamoto



CARPE DIEM



東京工業大学  
Tokyo Institute of Technology

科研費  
KAKENHI