

令和 2 年 6 月 15 日現在

機関番号：13901

研究種目：基盤研究(B) (一般)

研究期間：2016～2019

課題番号：16H03444

研究課題名(和文) 構文部分木頻度の確率的情報に基づく第二言語習得理論構築のための基礎的研究

研究課題名(英文) A fundamental study for constructing second language acquisition theory based on probabilistic information of syntactic tree fragment frequencies

研究代表者

杉浦 正利 (Sugiura, Masatoshi)

名古屋大学・人文学研究科・教授

研究者番号：80216308

交付決定額(研究期間全体)：(直接経費) 11,800,000円

研究成果の概要(和文)：8週間にわたる縦断的学習者コーパス(245名分、1,836エッセイ、約36万語)を構築し、エッセイ評価の上昇に影響する要因を分析し、総語数や節数などが影響を与えることを確認した。コーパスより構文部分木を抽出し、平均統語距離を調べた結果、母語話者に比べ学習者の部分木は統語的複雑性が低いことが明らかになった。構文知識を学習者も持っているかどうかを視線計測実験で調べたところ、二重目的語構文については母語話者と同様に持っていることを示唆する結果が得られた。言語処理の基本単位として構文部分木を使うことにより、第二言語処理能力の発達過程を包括的に説明できる可能性を確認した。

研究成果の学術的意義や社会的意義

約36万語の縦断的学習者コーパスを構築できたことは第二言語習得研究の基礎的データとして研究分野に貢献するものである。また、複雑性を表す言語的特徴としてこれまで平均文長等が使われてきたが、複雑性の結果としての「長さ」ではなく、理論的に複雑性そのものを表す平均統語距離(MSD)を新たに考案した点も重要である。MSDを使い、コーパスから抽出された構文部分木の複雑性を測り、学習者の言語処理が「浅い」ことを実証できた。視線計測実験により、学習者も構文知識を使っていることを実証できたことも意義深い。今後、構文部分木の重要性を生かし、高頻度構文部分木リストが開発されることにより、教育面への貢献も期待できる。

研究成果の概要(英文)：A longitudinal learner corpus (about 360,000 words in 1,836 essays written by 245 learners over the course of 8 weeks) was compiled. An analysis of this corpus found that factors such as the number of words and the number of clauses affected essay ratings. An examination of mean syntactic distances of syntactic tree fragments extracted from the corpus showed that learner tree fragments were lower in syntactic complexity compared to those of native speakers. The results of an eye-tracking experiment showed that regarding double object constructions, learners may have similar syntactic knowledge to native speakers to predict the constructions. By using syntactic tree fragments as the basic units of language processing, the possibility of a comprehensive description of the development of second language processing has been confirmed.

研究分野：第二言語習得論

キーワード：学習者コーパス 第二言語習得論 頻度効果 視線計測 コロケーション 構文文法 平均統語距離

## 様式 C-19、F-19-1、Z-19（共通）

### 1. 研究開始当初の背景

#### (1) 研究の学術的背景

第一言語獲得研究で、Peters (1983)が、幼児は単語を一つずつ覚えるのではなく、覚えた発話の中から共通する部分を分節化することで「語」を認識するとともに、その規則性から「文法」を習得するという仮説をたてた。Saffranら (1996) は、8ヶ月の乳児が発話の中から頻度の高い音連続のパターンを同定できることを心理言語学的な実験で証明し、乳児が音連続の統計的な知識を使って「語」を認識し習得していくことを示した。また、Tomasello (1992, 2000) の動詞の島仮説や、Goldberg (1995, 2006, 2009)の構文（コンストラクション）文法などの用例基盤モデルによる言語獲得研究が盛んになってきている。第二言語習得も用例基盤モデルによる説明が試み出されている。例えば、外国語学習者も「コンストラクション」という単位で第二言語を処理しているかどうかを Gries and Wulff (2005) が検証している。

#### (2) 本研究に関連する国内・国外の研究動向及び位置づけ

第二言語習得の普遍性という問題は、Brown (1973)による第一言語の文法形態素の習得順序の発見を、Dulay and Burt (1974) が第二言語の習得において検証し、第二言語習得においても普遍性があると発表したことに始まる。その後、第二言語習得の過程を Pienemann (1998) は、処理可能性理論 (Processability Theory) として精緻化した。処理可能性理論は、言語習得の発達段階を6段階としている。日本語母語英語学習者（大学生7名）を対象に、Sakai (2008) が統語の発達を検証し、全員が6段階の5か6であったと報告している。Eguchi and Sugiura (2015)の報告では、14人の中学生のうち9名が統語面において第5段階に達していた。しかし、こうした第5段階に達したと判定される学習者の第二言語の処理能力は本当に十分に高いといえるであろうか。また、この処理可能性理論の枠組みでは、形態素よりも統語の発達段階の方が高く出る傾向にあることや、定型表現の扱いが無視されているなどの問題点がある。

#### (3) 着想に至った経緯

平成 21～23 年度の研究「第二言語習得における処理単位に関する基礎的研究」で、「動詞＋目的語」に関するコロケーションについて眼球運動を観察したところ、英語母語話者では遷移確率の高い方が有意に速く処理されたのに対し、学習者の場合その差は現れなかった。これにより、遷移確率を基に言語処理能力を説明できるとのではないかと考え、平成 24～27 年度まで「遷移確率に基づく第二言語処理能力発達理論構築のための基礎的研究」を行った。その結果、遷移確率は頻度効果としてまとめられ、その頻度効果は第二言語習得において低頻度でも効果があることが確認されるとともに、単語間や文の構成要素間という、同じレベルでの遷移確率だけではなく、レベルを超えた多層的な言語単位の頻度効果を見る必要があるという着想に至った。

### 2. 研究の目的

本研究は、従来の研究が、文法の発達段階という枠組みを設定しようとしているのに対し、Goldberg がコンストラクションと呼んでいる形式と機能（意味）の結びついたさまざまなレベルの言語表現が、どのような単位と頻度で使用されるのかという観点から、第二言語の処理能力の発達過程の解明のための基礎的研究を行うことを目的としている。

### 3. 研究の方法

本研究では、英文エッセイデータを8週間にわたり収集し縦断的な英語学習者コーパスを構築し、品詞・構文解析をした後に、Discontinuous Data-Oriented Parsingを行う。これにより英文の統語構造としての「ツリー」を構成する全ての「部分木」の頻度を分析することができる。階層構造を含め単語から文全体までその文を構成するすべてのレベルにわたる構成部品が「部分木」となる。8週間にわたる縦断データを分析することで、言語産出単位とその頻度がどのように変化するかを観察することを通し、言語産出における言語処理の発達過程を解明することができる。また、部分木の頻度の差により、言語単位の認知処理にも違いがあるのかどうか、また違うとしたらどのように違うのかということ、視線計測装置を使った実験により明らかにする。コーパスデータの分析も視線計測実験も、英語母語話者データをあわせて収集することにより、母語話者データをベースラインとして、第二言語習得データの特徴を明らかにする。

### 4. 研究成果

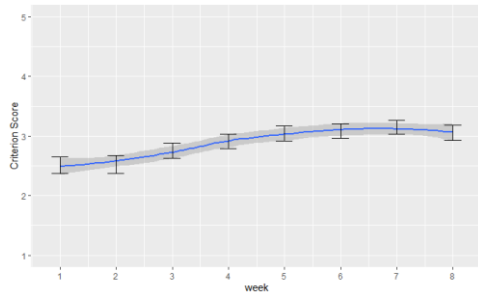
#### (1) 縦断的学習者コーパスの構築

3年間にわたる英文エッセイライティングデータの収集に基づく縦断的学習者コーパスを構築した。第二言語習得の過程を検証するためのコーパスということから「Nagoya Interlanguage Corpus of English for SLA Testbed」（略称「NICEST」）と呼ぶことにする。各学習者が8週にわたり8つのテーマでエッセイを書き、245名から延べ1,836個のエッセイを収集した（一部、収集の過程で学習者の欠席によ

NICEST	全体	第1週	第2週	第3週	第4週	第5週	第6週	第7週	第8週
ファイル数	1,836	226	234	228	233	231	229	224	231
総語数	360,208	35,977	37,540	42,266	45,723	49,992	49,515	48,487	50,708
平均語数	195.8	158.6	158.7	185.2	193.2	218.0	217.1	217.3	218.3
平均異語数	100.1	85.8	88.0	97.1	99.8	103.9	106.7	110.1	109.6
総文数	26,251	2,715	2,892	3,157	3,311	3,425	3,579	3,623	3,549
平均文数	14.3	12.1	12.3	13.9	14.4	14.7	15.8	16.3	15.3
平均文長	13.9	13.5	13.1	13.6	14.0	14.8	14.0	13.6	14.5
総段落数	8,461	965	1,039	1,033	1,063	1,092	1,080	1,084	1,105
平均段落数	4.6	4.3	4.4	4.5	4.6	4.7	4.7	4.8	4.8

\* 「平均」は1エッセイあたりの平均

る欠損があるため 1,960 にはならない)。また、参照用に英語母語話者データとして、20 名 160 個のエッセイもデータとして収集しコーパスデータ化した。8 週間にわたる時系列での変化については、1 回目から 8 回目にかけていずれのグループもエッセイのスコアが上昇していることから、学習者のライティング能力一般は伸びていると判断できる。



## (2) 学習者コーパスの分析

### ①第一言語 (L1) と第二言語 (L2) との相違に関する言語的特徴の分析

英語学習者と英語母語話者のライティングデータを比較し、どのような言語特徴がその二者に違いをもたらしているのかを分析した。書いたエッセイの総語数が影響しないように配慮し何通りかの組み合わせで複数の分析方法で分析をした。NICER データを対象に L1 と L2 との判別にどの言語特徴が影響を与えるかを、機械学習の一方法である Random Forest で分析した結果、MLT (平均 T-unit 長)と語彙の多様性指標 D とで約 90%のデータを正しく判別できた。次に、ICNALE と NICER のデータを対象にロジスティック回帰分析を行った結果、MLS(平均文長)・名詞句の複雑さ・異なり語数・動詞の多様性により約 95%を正しく判別できた。さらに、NICER データを対象に、二次判別分析を行った結果、MLT と語彙の多様性指標 MTL D とで約 96%の判別率を得ることができた。ただし、同じ判別式で ICNALE データを分析したところ判別率は約 90%であった。こうした一連の分析を通し、いずれの場合も、文もしくは文に準じる T-unit の「長さ」が重要な要因であることが分かった。しかし、「長さ」というのは結果的に「長さ」として測定できる性質のものであり、統語的な単位が長くなるということは、統語単位を構成する要素とその結びつき、すなわち統語構造が複雑になることと表れてはならないかと考えられる。

Logistic regression analysis lead to the following formula

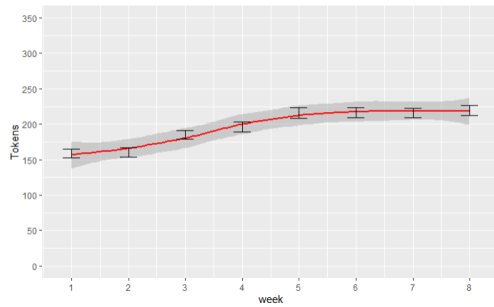
$$Z = .59(MLS) + 1.37(CN/C) + .24(Ndwrz) + .51(Ndwsz) + 1.15(CV1) - 46.34$$

Syntactic "Complexity"
Lexical "Complexity"

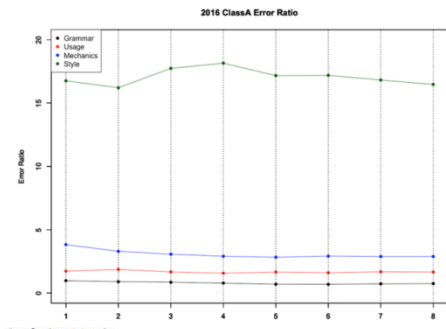
Predicted/Actual	NICER data (.96 accuracy)		ICNALE data (.94 accuracy)	
	L1	L2	L1	L2
L1	39	9	169	12
L2	1	176	26	480

### ②縦断的学習者コーパスの分析

NICEST の 8 週間にわたるエッセイライティングの縦断的データは、エッセイの評価が徐々に上昇していることからライティング能力が上昇していると判断できる。そうしたライティング能力の上昇をもたらしているのはどのような言語的特徴かということ进行分析した。その結果、成績の上昇とエッセイの総語数の上昇とがほぼ同じ傾向にあることから、総語数の影響が強いと判断される。



また、エッセイを評価した Criterion に基づく誤用のタイプの推移を見ると、誤用率はいずれのタイプもほぼ横ばいで変化がないことから、エッセイライティングにおいて、誤用が減ったからスコアが上がったということは言えないと判断される。



さらに、詳細な分析として、エッセイスコアの上昇に、どの言語的特徴が寄与しているかについて、文の複雑性の一要因である関係節の使用頻度と、文章構成に関係する談話標識の頻度も含めて、順序ロジスティック回帰を使って分析した。その結果、総語数だけでなく、節数、従属節数、複雑な名詞句の数、平均 T-unit 長、T-unit あたりの動詞句の数、低頻度語のタイプ数とトークン数が有意に影響するという結論を得た。こうした言語的特徴により約 58%のスコアを正しく予測できた。

Confusion Matrix

Observed	Estimated				
	1	2	3	4	5
1	0	8	18	0	0
2	0	9	150	9	0
3	0	10	376	104	0
4	0	0	52	116	2
5	0	0	1	3	1

Over-estimated  
Under-estimated

Overall accuracy: 502/859 = .58

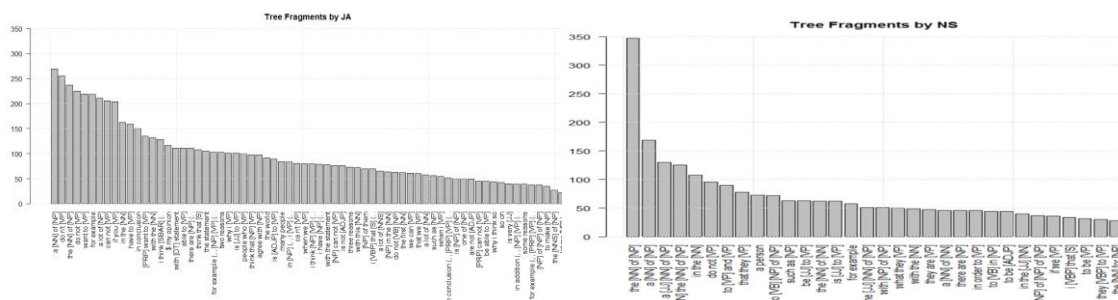
### ③構文部分木の頻度分析

NICEST より、各プロンプトあたりのデータ量 (総語数) がほぼ同じ (11,000 語) になるように、JA と JB の学習者 2 グループ (プロンプトの提示順序に A とその逆順の B の二種類があるため) と母語話者 (NS) 1 グループの、3 グループからなるデータセットを作成した。JA と JB については 1 週目から 8 週目までの時系列順と、時系列に関係なくプロンプトの種類別にデータを組み合わせさせた二通りの分析を行った。このデータをもとに構文部分木の頻度一覧表を作成した。構文部分木の総タイプ数はおよそ 6 万項目となった。学習者データ全体と母語話者データ全体の上位 500 の構文部分木の対数化頻度の相関を分析し相関係数 0.77 という結果を得た。すなわち、学習者と母語話者とでは使用する構文部分木はおおむね共通す

Ratio of Tree Fragments including Content Words

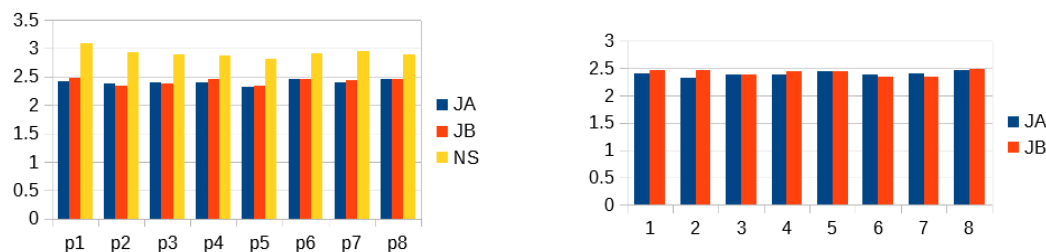
	L1	L2 (A)	L2 (B)
Tree Fragments including Content Words	3	32	30
Total number of Tree Fragments	32	66	74
Ratio	9%	48%	41%

ることが確認された。逆に、学習者と母語話者との間の頻度差が大きい部分木に見られた特徴として、母語話者は冠詞や前置詞など機能語を含む木を多く使用しており、学習者は人称代名詞や動詞+名詞句から成る動詞句を多く使用していることが観察された。



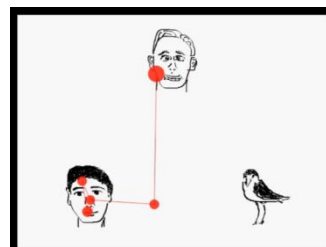
#### ④構文部分木の平均統語距離 (Mean Syntactic Distance)

構文部分木を構成するノード間の距離の合計を、ノードをつなぐ枝の数で割ることで「平均統語距離 (Mean Syntactic Distance)」を算出する方法を今回新たに考案した。上記の JA データ、JB データおよび NS データの MSD を、プロンプトごとに算出して比較したところ、JA・JB 間にはほぼ同じであるのに対し、NS とでは差があることが観察された。JA と JB について、時系列で観察したところ、JA と JB いずれも第 1 週から 8 週にかけて上昇は見られなかった。これは、統語的複雑さに増加は見られなかったことを示していると考えられる。



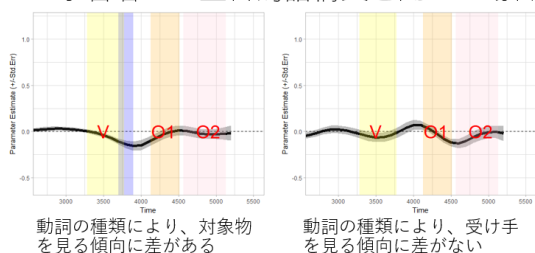
#### (3) 視線計測実験

英語の構文として広く研究が行われている二重目的語構文と交代可能な前置詞を伴う与格構文とを分析対象として視線計測実験を行った。Gries & Stefanowitsch (2004) のコーパス分析データをもとに与格構文よりも二重目的語構文をとる傾向が強い動詞 (give, show)、その逆で与格構文をとる傾向の強い動詞 (bring, take)、および、いずれも同程度の動詞 (get, send) を使い、学習者 29 名と母語話者 14 名を対象に、Visual World Paradigm という絵と音声を使った視線計測実験を行った。例えば、「John will give Taro a bird.」という英文を聞く時に、どのように視線が移動するかを計測する。その際に、give を使った文の場合、「John will give」を聞いた段階で、give は二重目的語構文をとる傾向が強いため、give の次に「人」が来ることを予測し、視線が鳥ではなく「Taro」の方へ移動すると考えられる。この「予測して視線を移動させる」ということが、その動詞がどのような構文をとる傾向が強いかという構文に関する確率的知識を持っている証拠となる。

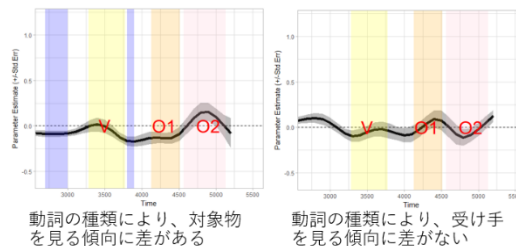


実験の結果、二重目的語構文を使った場合、学習者・母語話者いずれも、動詞の種類によって動詞の直後の予測に差があることが確認された。逆に、与格構文の場合には、学習者・母語話者いずれも、動詞の種類による差は確認されなかった。

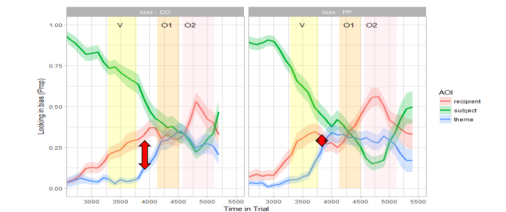
#### 学習者が二重目的語構文を聞いた場合



#### 母語話者が二重目的語構文を聞いた場合

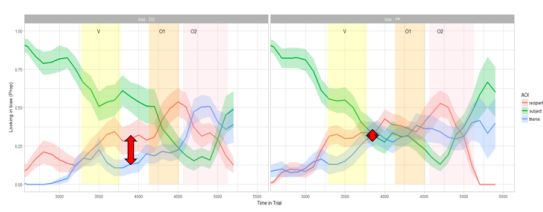


## 学習者が二重目的語構文を聞いた場合



二重目的語構文をとる傾向のある動詞の場合、動詞直後で受け手（赤線）を見る傾向がある  
与格構文をとる傾向のある動詞の場合、動詞直後で差（赤線と青線間）がない

## 母語話者が二重目的語構文を聞いた場合



二重目的語構文をとる傾向のある動詞の場合、動詞直後で受け手（赤線）を見る傾向がある  
与格構文をとる傾向のある動詞の場合、動詞直後で差（赤線と青線間）がない

これにより、学習者も、母語話者と同様に、二重目的語をとる傾向が高い動詞 (give, show) については、動詞を聞いた段階で、次に受け手である「人」が来ることを予測していることが観察された。すなわち、構文に関する確率的知識を持っていると結論付けられる。また、二重目的語構文で使われる傾向が低い動詞 (bring, take) については、母語話者も学習者も動詞の直後で仮説としては「モノ」が来ることを予測すると思われたが、実験結果としては、いずれも、「モノ」と「人」との間で差が見られなかったということから、構文に関する確率的知識がない、もしくは、働いていなかったと結論付けられる。

与格構文をとる動詞に関して仮説に反する結果とはなったが、いずれにせよ、学習者と母語話者と同じ結果になったということは、学習者と母語話者が同じ確率的知識を持っているのではないかということを示す結果となった。

しかし、この結果については、本研究で使用した構文が、非常に基本的な動詞に関する構文であったことから、学習者も典型的な構文に慣れ親しんでいたため、母語話者との間で差がみられなかったという可能性も指摘できる。むしろ、差がある構文にはどのような構文があるか、を今後探索していく必要があるといえる。

### (4) 第二言語処理能力の発達理論

学習者コーパスの分析より、L1 と L2 との違いは、統語的複雑さ (平均文長や名詞句の複雑さ) と語彙的複雑さ (語彙の多様性、動詞の多様性) によるところが大きいと考えられる。統語的複雑さのうち平均文長 (もしくは平均 T-unit 長) という「長さ」の指標は、「長さ」と「複雑性」の関係からすると、複雑な結果長くなるという因果関係にあると考えるのが合理的である。すなわち、統語的に複雑であるがゆえに、結果的に文が長くなる、ということである。であれば、指標としては、「長さ」ではなく統語的複雑性そのものを代表する指標の方が望ましい。

そこで、平均統語距離 (MSD) という新たな統語的複雑性指標を開発し、NICEST のデータを使い構文部分木の MSD を比較したところ、学習者と母語話者とで差があることが観察された。学習者の方が MSD の値が低い。これは学習者の使用する構文部分木の統語的複雑さが低いことを表している。これは、第二言語では深い統語処理ができないという Clahsen らの「Shallow Structure Hypothesis」の主張に沿うものである。また、NICEST の 8 週にわたるエッセイライティングで、スコアは上昇しているにもかかわらず MSD は上昇していないことから、NICEST の実験参加者である大学生レベルでは、統語能力の発達は観察されなかったといえる。これは、平均文長の伸びも観察されなかったということとも整合性がある。第二言語での統語能力の発達を観察するには、習得初期からのデータを分析する必要があるが、それだけでは第二言語処理能力の発達は説明できない。

構文部分木という言語処理単位を想定することで、様々なレベルの言語処理単位を包括的に取り扱い議論することができるようになる。興味深い言語現象でありながら、その定義が定まらず、それゆえに議論も収束しない各種コロケーションやレキシカル・バンドルも統語範疇と語彙項目とを含んだ構文部分木として一元的に扱えるようになる。また、構文文法が形式と意味が対応するあらゆる言語単位がコンストラクションであると言う場合の単位も構文部分木としてとらえることができる。言語処理上の基本単位として構文部分木をとらえ、その習得と運用を観察することが第二言語処理能力の発達研究の発展へとつながる。

本研究の視線計測実験で、統語処理が自動的に行われる証拠として構文の予測能力を調べたところ、母語話者と同様に学習者でも予測をしていることが観察された。しかし、今回の実験項目は非常に基本的な動詞であったために学習者も予測できたとも考えられる。今後は逆に、母語話者は予測できるのに学習者はできない項目を探すことが、第二言語処理能力固有のメカニズムの解明につながると考えられる。

内容語さえわかれば第二言語の理解はある程度できるが、逆に、内容語だけを並べても第二言語の産出はままならない。高頻度で使われる機能語を含んだ構文部分木の習得が第二言語処理能力の発達を効率よく促すと考えられる。そうした高頻度構文部分木のリスト作りが今後の第二言語処理能力の発達研究の発展の礎の一つとなるであろう。

5. 主な発表論文等

〔雑誌論文〕 計0件

〔学会発表〕 計12件（うち招待講演 2件 / うち国際学会 9件）

1. 発表者名 Sugiura, M.
2. 発表標題 At the intersection of corpus linguistics and second language research
3. 学会等名 New Methods and Data in Second Language Learning Research Meeting (招待講演) (国際学会)
4. 発表年 2019年

1. 発表者名 Sugiura, M., Nishimura, Y., Abe, D.
2. 発表標題 Which linguistic features contribute to essay ratings and when?: A preliminary study employing ordinal logistic regression analysis
3. 学会等名 The International Conference on Foreign Language Education & Technology (FLEAT VII) (国際学会)
4. 発表年 2019年

1. 発表者名 Komori, S., Sugiura, M., Li, W.
2. 発表標題 Examining MDD and MHD as Syntactic Complexity Measures with Intermediate Japanese Learner Corpus Data
3. 学会等名 SyntaxFest 2019 (International Conference on Dependency Linguistics: Depling 2019) (国際学会)
4. 発表年 2019年

1. 発表者名 Sugiura, M.
2. 発表標題 Toward an integrated theory of SLA using tree fragments
3. 学会等名 ESRC-AHRC UK-Japan SSH Connection Grants Seminar/ LCSAW (Learner Corpus Studies in Asia and the World) 4th Meeting (招待講演) (国際学会)
4. 発表年 2019年

1. 発表者名 Sugiura, M., Abe, D., and Nishimura, Y.
2. 発表標題 Linguistic traces of L2 processing difficulties found in learner corpus data: A quest of discriminant function between L1 and L2
3. 学会等名 The 39th Annual Conference of the International Computer Archive for Modern and Medieval English (ICAME39) (国際学会)
4. 発表年 2018年

1. 発表者名 Abe, D. and Sugiura, M.
2. 発表標題 The development of English noun phrase complexity in EFL students
3. 学会等名 The 39th Annual Conference of the International Computer Archive for Modern and Medieval English (ICAME39) (国際学会)
4. 発表年 2018年

1. 発表者名 杉浦正利, 西村嘉人, 阿部大輔
2. 発表標題 ライティング自動評価システムを活用したライティング指導 縦断的学習者コーパス構築に向けて
3. 学会等名 外国語教育メディア学会 (LET) 第58回 (2018年度) 全国研究大会
4. 発表年 2018年

1. 発表者名 Sugiura, M., Nishimura, Y. and Abe, D.
2. 発表標題 Distinguishing L1 and L2 Using Three Linguistic Aspects: A Logistic Regression Model Study
3. 学会等名 The Fourth Asia Pacific Corpus Linguistics Conference (APCLC 2018) (国際学会)
4. 発表年 2018年

1. 発表者名 Sugiura, M. and Abe, D.
2. 発表標題 On Japanese EFL Learners' Sensitivity to Verb Biases in the Dative Alternation Constructions: A Visual World Paradigm Study
3. 学会等名 American Association For Applied Linguistics 2019 Annual Conference (AAL2019) (国際学会)
4. 発表年 2019年

1. 発表者名 杉浦正利、阿部大輔、西村嘉人
2. 発表標題 縦断的英語学習者コーパス構築の試み
3. 学会等名 外国語教育メディア学会 (LET) 第89回春季中部支部研究大会
4. 発表年 2017年

1. 発表者名 杉浦正利、阿部大輔、西村嘉人
2. 発表標題 英文エッセイライティングにおけるトピックの影響 縦断的学習者コーパスの分析
3. 学会等名 平成29 (2017) 年度大学英語教育学会 (JACET) 中部支部大会
4. 発表年 2017年

1. 発表者名 Sugiura, M., Abe, D., and Nishimura, Y.
2. 発表標題 What Kind of Linguistic Features Distinguish Second Language Learners' Texts from Those of Native Speakers, and Why?
3. 学会等名 The 4th Learner Corpus Research Conference (LCR 2017) (国際学会)
4. 発表年 2017年



〔図書〕 計0件

〔産業財産権〕

〔その他〕

Nagoya Interlanguage Corpus of English for SLA Testbed (NICEST)  
<https://nicest-sugiura.blogspot.com/>

6. 研究組織

	氏名 (ローマ字氏名) (研究者番号)	所属研究機関・部局・職 (機関番号)	備考
--	---------------------------	-----------------------	----