

# 情報通信端末による日本語入力の語彙の偏りに関する研究

代表研究者 竹原卓真 同志社大学 心理学部 教授

## 1 緒言

昨今の劇的な情報通信網の拡大と、それを利用したコミュニケーションの増加や多様化は、我々の生活に多大な恩恵をもたらしてきた。現在では様々な最新情報を手軽かつ瞬時に得られるだけでなく、その周辺情報さえも連鎖的に収集できる。また、情報通信のイノベーションはコミュニケーションおよび地域活性化のツールとして有効であり、さらにはそれらに関連してビジネスや社会への貢献が著しい。加えて、日常場面だけでなく、災害などの緊急時にも必須のツールとなるなど、情報通信のイノベーションに関するメリットは枚挙に暇がない。

その一方で、発達した情報通信技術のいわゆる影の部分も徐々に明らかになってきている。例えば、電子掲示板での誹謗中傷や、誤った顔文字および絵文字を付与した電子メール送受信における問題などがある(竹原・栗林・水岡・瀧波, 2005a; 竹原・栗林・武川・水岡・瀧波, 2005b)。その他、迷惑メールの問題やウィルス感染による個人情報の漏洩など、情報通信技術がバックボーンとなるインターネットの世界は単に利便性のみで成り立っているものではない。これらの利便性に遮蔽された影の側面は実際のコミュニケーションにおいて重要であるにもかかわらず、科学的に議論されることが少なく、情報通信技術革新至上主義における多数派同調バイアスや正常性バイアス等が関与して軽視されがちであると判断せざるを得ない。そこで、今後も必要となる情報通信の進歩に対し、影の部分について科学的な観点から検証を加え、我々が忘れてはならない問題点を明らかにすることが必要だと考えた。

上述した影の部分をあぶり出す有効な研究対象として、我々が普段から使用している言語がある。一般的に、言語を使用する際に我々は多種多様な言葉や単語を使用し、その使用頻度は特定の単語が平均的になるような正規分布に従うと考えがちである。しかし、認知言語学においては、我々が使用している言語の出現頻度に大きな偏りがあり、正規分布に従わないことが報告されている。Zipf (1949)は英単語の出現頻度とそのランクとの関係を詳細に調べ、ランクが2倍になれば出現頻度が1/2になるというベキ乗則(ベキ乗分布)を発見した。例えば、英単語では”the”が最も高頻度で出現するためランクは1位となり、2位の”of”の出現頻度は1位の”the”の半分ということになる。ベキ乗則の観点から言い換えると、英単語を使う人たちは極めて高頻度で”the”や”of”といった単語を使うが、その一方でほとんど使用しない英単語も膨大な量に及ぶことになる。また、ベキ乗則に従うため、この関係性を両対数グラフで表すと図1に示すように傾きの絶対値がほぼ1.0となる。この単語の出現頻度とランクにおけるベキ乗則の関係はZipf(ジップ)の法則と呼ばれ、日本語やその他の多くの言語でも出現する。なぜこの法則が出現するのかについては明確な根拠が示されていないが、単一化/多様化(unification/diversification)の原理から論じられることが多い(Zipf, 1949)。この原理は単語の送信者と受信者との間における情報通信の労力から説明され、グラフの傾きが1.0周辺の場合に送信者と受信者の労力は拮抗してバランスが良くなる。しかし、傾きが1.0より大きくなれば送信者が使用する単語の語彙が少なく(労力が少なく)、受信者はその少ない情報から状況認識を行わねばならないため多くの労力が必要となる。逆に、傾きが1.0より小さければ送信者が使用する語彙は多く(労力が多く)、受信者は容易に理解できる。

ジップの法則は単語の出現頻度のみならず、様々な対象について発見されてきた。たとえば、横川(2002)によると、都市の人口、文献の参照回数、企業の所得分布、映画の観客動員数などでもジップの法則が当てはまる。無論、これらの現象においては上述の単一化/多様化の原理から説明が難しいのは明白であるが、言語で発見された法則が全く異なる領域で数多く発見されることは非常に興味深いと考えられる。言い換えると、一般的に偏りの少ない正規分布と想定されている現象においても、極めて偏りの大きいベキ乗分布が潜んでいる可能性が示唆される。

以前における言語の使用と比較し、最近になって言語の使用方法そのものに変化が生じてきている。その変化の一つとして情報通信端末における言語の使用があげられ、これは爆発的なスマートフォンの普及が関与していると考えられる。現在我々の日常生活に深く浸透しているスマートフォン等の情報通信端末の日本語入力には、入力の手間を省くための様々なアシスト機能が実装されている。その一つに、学習による予測

変換機能がある。この機能は一度入力した単語変換を端末が記憶し、次の単語入力時に変換候補として提示して入力をアシストするもので、非常に便利である反面、同じ単語ばかりが選択される傾向があり、伝達されるメッセージの語彙低下という問題が予想される。語彙の低下はコミュニケーションや会話の質低下に直結して様々なトラブルの原因になる可能性が高い上、利用者自身の国語能力低下とも関連する可能性がある。昨今声高に叫ばれている若者の国語能力低下に鑑みれば、利用者が将来の日本を背負う若者であればあるほどそれらの問題に取り組む価値があることは自明であろう。加えて、この問題を定量的に明らかにするためには情報通信端末を使用しない従来の紙への筆記による会話や、発話による会話のデータと比較して論じることが重要である。

そこで、本研究では情報通信端末に実装されている予測変換機能にスポットを当て、これらの端末を使用した実際の会話を通じて単語の出現頻度分布の調査を行い、情報通信端末利用者が使用する日本語の語彙においてトラブルの原因となり得る出現頻度の偏りを明らかにすることを主目的とする。具体的には、情報通信端末を利用しない紙への筆記による会話や発話による会話における単語の出現頻度分布の偏りを調査し、情報通信端末を利用した結果と比較する。これによって、便利ではある反面、同じ単語の繰り返し使用の原因となる予測変換機能が送受信者間の通信努力のバランスに影響を与えていることや語彙力の低下等、コミュニケーションに重要な国語能力にネガティブな影響を与えていること等が明らかになる可能性がある。この視点からコミュニケーションと言語との関係について検証された研究は我々が知る限り存在せず、極めて斬新である。また、情報通信端末における予測変換機能、つまり現在発展中である我が国の情報通信技術がもたらす、いわゆる「影」の部分をつまびらかにできることに大きな意義があると考えられる。

## 2 方法

### 2-1 実験参加者

実験参加者は大学生 96 名（男性 26 名、女性 70 名、平均年齢 20.61 歳、標準偏差±0.10 歳）であった。

### 2-2 実験装置

実験装置にはストップウォッチと Apple 製 iPhone4s（以下、iPhone と略記）および、OLYMPUS 製ボイスレコーダー（LINEAR PCM RECORDER LS-P2）を用いた。会話のため、iPhone には LINE アプリ（株式会社 LINE 製；以下 LINE と略記する）をあらかじめインストールした。また、A4 サイズの白紙、および赤色と黒色のボールペンを使用した。

### 2-3 実験デザイン

LINE を使用し、かつ予測変換を使用して会話する条件を予測変換条件、LINE を使用するが予測変換を使用しないで会話する条件を非予測変換条件、筆談での会話をする筆記条件、通常通り発話で会話する発話条件とし、この 4 条件に実験参加者をランダムに割り当てた。

### 2-4 手続き

まず全ての条件に対して共通の手続きを述べる。実験参加者募集にあたり、仲の良い同性の友人であることを条件に設定した。その理由は、事前に 139 名の男女を対象にして LINE で最も会話が弾む人数を報告させたところ、約 55%にあたる 76 名が 2 人と報告したからである。実験参加者を上述の 4 条件にランダムに割り付け、2 人 1 組で実験室に入室させ、実験室の中央に設置した机に向かい合って座らせた。なお、予測変換条件、非予測変換条件において、机にはアイコンタクト等を遮蔽する目的で衝立を設置し、実験参加者が互いに姿を確認できないようにした。そして、実験内容について説明を行い、同意を得た。その後、条件毎に使用する装置を手渡し、全ての条件において自由に 15 分間会話させるセッションを 2 回、セッション間に 3 分ほど休憩時間を設けて行った。つまり、実験参加者には合計で 30 分間、会話を行わせた。この 2 回のセッションにはそれぞれランダム順に「アルバイト」および「日本の県民性」という 2 種類のテーマを設定した。テーマ設定は会話の導入を統一することを目的としており、会話内容については一切制限を設けず、テーマから会話が逸れても良いと教示した。なお、各条件における個別の手続きは以下の通りである。

#### （1）予測変換条件

実験者は実験参加者に iPhone を手渡し、個人情報等についての情報が一切取得されないことを確認させた後、同様に LINE には実験参加者ペアの端末のみが友達登録されていることを確認させた。その後、LINE 使用のときには顔文字、スタンプ、絵文字などの明確に言語化できないやりとりになり得るアイコンの使用を禁止すると教示した。上述の 2 セッションの会話終了後、報酬を手渡して実験室を退出させた。

## (2) 非予測変換条件

原則的に、予測変換条件と同じ手続きを用いた。ただし予測変換条件と異なり、教示の際に予測変換の使用を禁止する旨の説明を加えた。また、実験参加者が別人に変わるたびに iPhone を初期出荷状態に戻した。

## (3) 筆記条件

実験者は参加者のペアごとに A4 用紙と赤ペン、黒ペンを 1 本ずつ手渡した。ペアの片方が黒ペン、片方が赤ペンを使用して A4 用紙に横書きで会話を筆記してその紙を渡し合う形で筆談を行わせた。会話中はアイコンタクトと私語を禁止することを教示し、会話終了後、報酬を手渡して実験室を退出させた。

## (4) 発話条件

実験者は参加者に日常で行うような通常の会話をするよう教示し、ボイスレコーダーを参加者ペアの間に設置したテーブルの中央に置いて会話を録音した。また、他の条件と同様に 15 分間の会話のセッションを 2 回行わせた。会話終了後、報酬の受け渡しを行い、退出させた。

## 3 結果

実験終了後、全条件における全会話データを電子化した。具体的には、LINE アプリを使用した 2 条件ではスマートフォンに残っている会話データをそのまま言語データとして取り出し、筆記条件では紙に書かれた文章をコンピュータ上で文字起こしし、発話条件では IC レコーダーに記録された音声データを逐語録としてコンピュータ上で文字起こしした。このようにして電子化された会話データは、形態素解析サイト「茶まめ」に貼り付け、条件ごとに全ての出現単語の形態素単語とその出現頻度を算出した。そして、算出された単語に対し出現頻度が高いものから順にランクを付けて、ジップの法則で利用された横軸をランク、縦軸を出現頻度とする両対数グラフにデータをプロットした。さらに、各条件において回帰分析を行って回帰曲線の傾きおよび説明率 ( $R^2$  値) を算出した。

図 1 に予測変換条件のグラフを示す。単語の出現種類の総数は 2733 種類となり、出現頻度で 1 位は「笑」(普通名詞) で 570 回、2 位は「！」(補助記号) で 396 回、3 位は「た」(助動詞) で 353 回、4 位は「？」(補助記号) で 353 回、5 位は「の」(終助詞) で 303 回であった。回帰曲線は傾きの絶対値は 1.14、説明率は 99.49% となった。

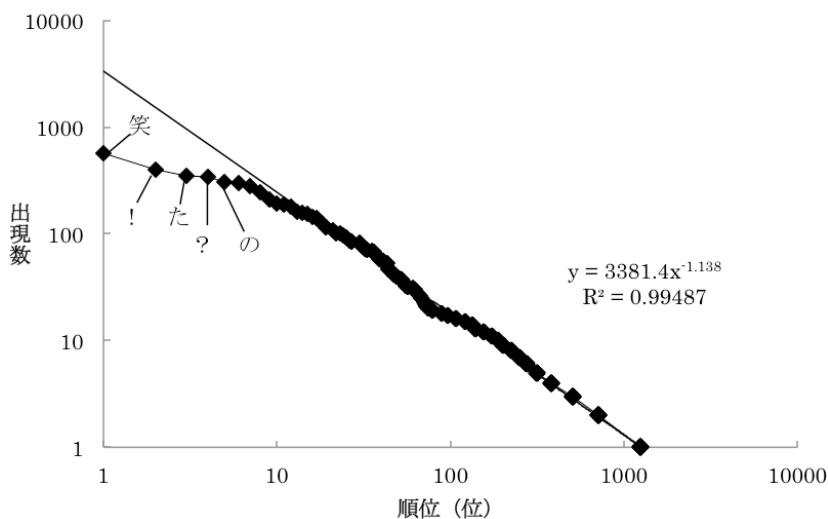


図1. 予測変換条件の分布図

図 2 に非予測変換条件のグラフを示す。単語の出現種類の総数は 2890 種類となり、出現頻度で 1 位は「？」(補助記号) で 460 回、2 位は「！」(補助記号) で 426 回、3 位は「た」(助動詞) で 326 回、4 位は「の」(格助詞) で 315 回、5 位は「は」(係助詞) で 285 回であった。回帰曲線は傾きの絶対値は 1.15、説明率は 99.07% となった。

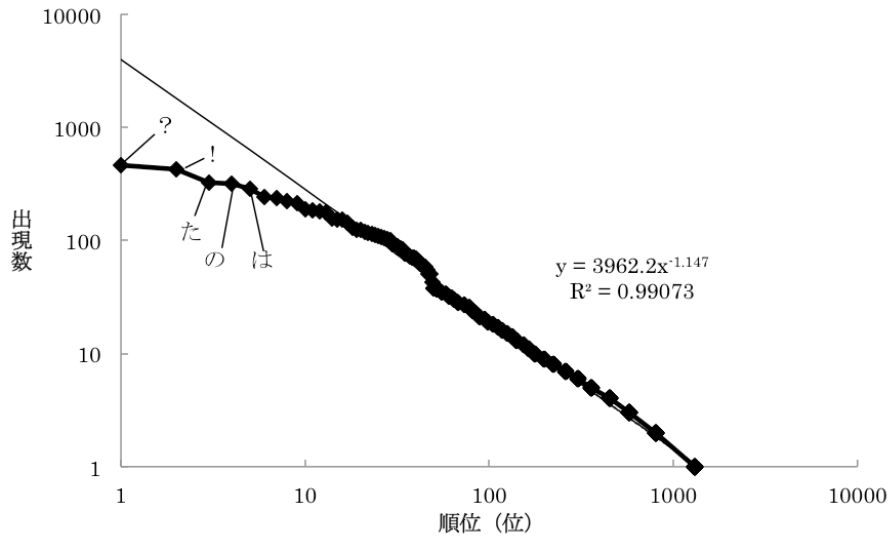


図2. 非予測変換条件の分布図

図3に筆記条件のグラフを示す。単語の出現種類の総数は1659種類となり、出現頻度で1位は「。」(補助記号)で286回、2位は「!」(補助記号)で222回、3位は「ば」(係助詞)で221回、4位は「か」(副助詞)で205回、5位は「の」(格助詞)で204回であった。回帰曲線は傾きの絶対値は1.18、説明率は98.91%となった。

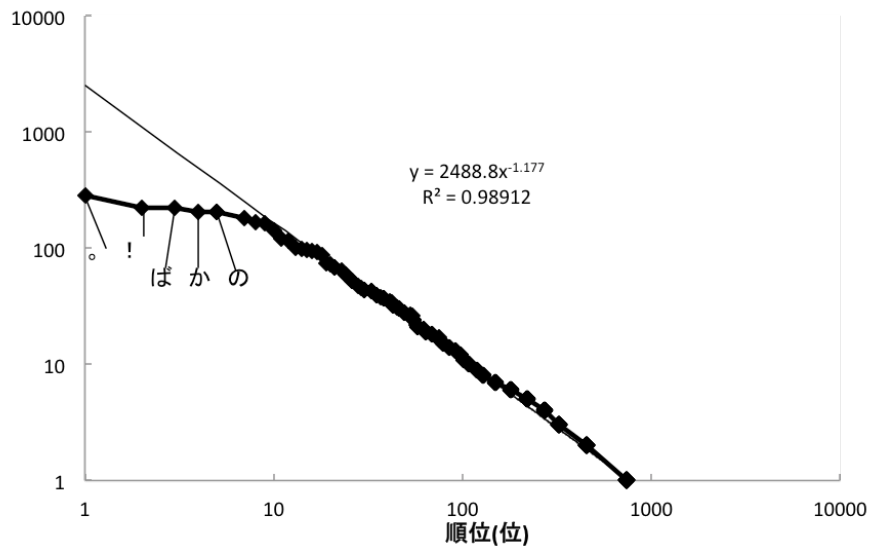


図3. 筆記条件の分布図

図4に発話条件のグラフを示す。単語の出現種類の総数は5998種類となり、出現頻度で1位は「か」(副助詞)で2123回、2位は「うん」(感動詞)で1720回、3位は「て」(接続助詞)で1651回、4位は「た」(助動詞)で1537回、5位は「さん」(接尾辞)で1471回であった。回帰曲線は傾きの絶対値は1.41、説明率は98.72%となった。

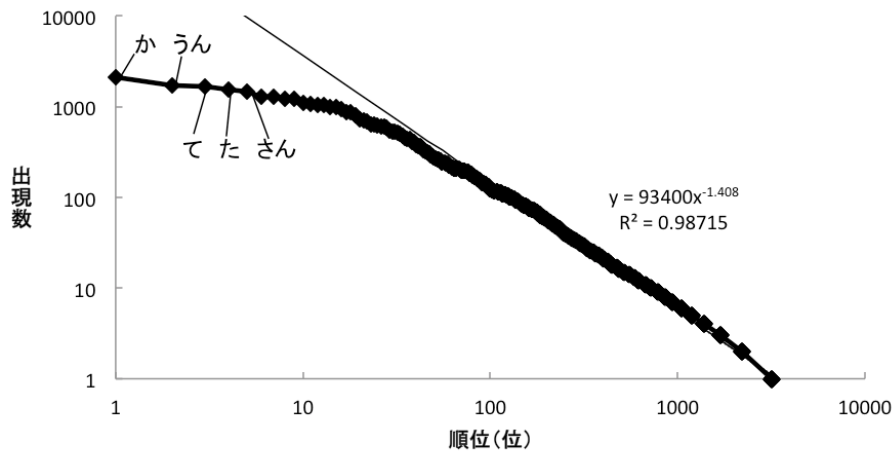


図4. 発話条件の分布図

#### 4 考察

近年の爆発的な情報通信端末の普及とその眼を見張るようなイノベーションから我々は様々な恩恵を受けている。ハードウェアや科学技術の進歩は我々の生活に必要な不可欠なものになっている反面、それらの利用方法においては従来から色々な問題点や問題となる行動が指摘されてきた(e.g., 竹原ら, 2005a, 2005b)。また、若者の国語力低下が叫ばれ始めて久しいが、我々が普段使用している語彙は、実は極めて偏った状態にあることが判明しており、これはジップの法則と呼ばれる(Zipf, 1949)。ここで、当初は情報通信端末の利便性を高める目的で考案・実装された文字入力の予測変換機能は一度使用した単語が変換候補として必ず表示されるため、情報通信端末を利用した会話においてジップの法則をより偏らせている可能性が考えられる。そこで、本研究では情報通信端末に実装されている予測変換機能が実際には会話等のコミュニケーションに暗い影を落としているかもしれないと考え、メッセージの送受信者間の通信努力のバランスが偏っているのか、また、予測変換機能により利用する単語の種類が低下し、その結果として語彙力が低下しているかを調べた。

まず、結果の中で注目すべきことは、設定した全ての会話条件で語彙の出現頻度分布がほぼ100%の説明率でベキ乗則に従ったということである。出現頻度分布がベキ乗則に従うということは、ほんの少数の語彙だけが極めて高頻度で使用される一方で、極めて多数の語彙がほとんど使用されないことを意味する(Barabási, 2003)。言い換えると、我々がどのような手段であれ会話する場合は、特定の語彙だけを偏って使用していることになる。人間の心に関わる数多くの指標においてはバランスの取れた正規分布が一般的であるが(Barabási, 2003)、普段から利用している語彙はその例に当てはまらないことが判明し、極めて重要な発見となった。正規分布する測定指標には明確な代表的数値、つまり平均値が存在するが、ベキ乗則に支配される測定指標には明確な代表値が存在せず、還元論的な解析を適用することができない。なぜなら、平均値が意味を持たないうえに分布がフラクタル分布となるため、いくら結果を還元しても単純な構造に帰着しないからである(Buchanan, 2002)。本研究のベキ乗則の出現はサンプル数が多いこともあり、非常にロバストな法則に支配されていることにほかならない。

続いて、本研究結果がジップの法則に従うかどうかを検証する。序論でも言及したように、情報の送受信者間における労力バランスが均衡するのが、図1における回帰直線の傾きの絶対値が1.0周辺である(Zipf, 1949)。この値と照らし合わせてみると、本研究では予測変換条件と非予測変換条件では約1.14と1.15と、1.0よりも数値が大きい。この値は同程度の数値において送受信者間の通信労力バランスが崩れていると結論したTakehara, Ochiai, & Suzuki (2015)の研究結果と類似しており、本研究の予測変換および非予測変換条件、つまりスマートフォンを用いた入力による会話では情報の送信者側の労力が軽減され、逆に受信者の労力負担が大きくなる方向に働くと考えられる。加えて、会話で使用される語彙が低減していることの証

左となり、スマートフォン等の情報通信端末を用いることによって語彙が貧弱になっている可能性が示唆される。しかし、予測変換機能の有無における分布の傾き値がほぼ同じことから、語彙の低下に予測変換機能が影響を及ぼしたわけではない。これらの結果から推測できることは、予測変換機能の有無にかかわらず、情報通信端末の文字入力機能を使用して会話を行うだけで、メッセージ送信者の語彙が低下してしまうということである。オンラインでのメッセージ送信は推敲を重ねることが少ないため(竹原・佐藤, 2003)、使用する語彙が特定化されてしまうことが原因であるかもしれない。

一方、情報通信端末を利用しない古典的な会話ではどうだろうか。語彙の出現頻度分布における回帰直線の傾き値を見る限りでは、筆記による会話で 1.18、音声による会話で 1.41 と、情報通信端末を利用した会話と比較するとほぼ同じか明らかに大きい数値となった。この結果から、筆記による会話ではメッセージ送信者の語彙が低下して受信者の認知処理負荷が増加しており、情報通信端末と同様にメッセージ送信者の語彙低下を招くことが見て取れる。筆記による会話はいわゆる筆談であり、本研究では互いの顔が見えないように実験参加者間に衝立を立てて実験を行った。この状態は相手の感情状態や気分状態、あるいはバーバルおよびノンバーバル行動を一切観測できない。言い換えれば、筆談してはいるものの、環境的には情報通信端末を利用した会話と同質である。従って、情報通信端末を利用した会話において出現した傾き値 1.14 や 1.15 に類似しているとしても不思議ではない。他方、音声による会話である発話条件では傾き値が 1.41 と明らかに他の 3 条件とは異なる値となった。つまり、音声会話では使用語彙が極めて偏っており、メッセージ送信者の労力は少なく、逆に受信者の労力は非常に多くなる。一見すると、古来最も一般的な会話環境である音声会話が確実に語彙低下を招いており、国語力低下の原因であると考えられるかもしれない。しかしながら、音声会話条件では実験参加者間に IC レコーダーが設置されただけで、互いに顔の表情やノンバーバル行動を視認することができる。つまり、会話を行ううえで利用できる他者の情報が他の条件よりも多い。利用可能情報が多いということは、文字を介して伝達しなければならない情報のアシストとなるため、情報の送信者の語彙が低下したものと考えられる。同様に、傾き値こそ 1.41 と大きい、情報受信者も利用可能情報が多いため、メッセージの理解に特に負荷がかかっているわけではないのかもしれない。その他、発話条件のグラフの傾きが他 3 条件に比べ大きい値を示している理由として、会話の送信中にも受信者の相槌などの反応を得やすく、会話における間を考慮しない点も考えられる。他 3 条件においてはメッセージ送信してからその反応を得るまでの空白の時間があり、また発信するメッセージを視認することができるため、発話条件に比べて受信者の負担を考慮している可能性も考えられる。加えて、小説や新聞などの媒体での単語出現頻度分布の傾きがほぼ 1.0 に近くなるというのは、メッセージを効率的に伝達し、受信者の解読に要する負荷を大きくしないように配慮しているためではないかとも捉えられる。

本研究では、これまで明らかにされなかった情報通信端末を用いた語彙量の測定を行い、ベキ乗分布に従うことや情報通信端末が我々の語彙に落とす影の部分の部分を明らかにすることができた。言語に関するジップの法則を適用した研究は英語を対象にしたものがほとんどであり(e.g., Corominas-Murtra, Fortuny, & Sole, 2011; Ferrer-i-Cancho & Sole, 2003; Thurner, Hanel, Liu, & Corominas-Murtra, 2015)、日本語を対象にしたものはほとんどない。従って、今後は日本語の様々な側面におけるジップの法則を検証することが重要になるだろう。加えて、英語では単語単位でジップの法則は出現せず、フレーズ単位で出現するという報告があるため(Williams, Lessard, Desu, Clark, Bagrow, Danforth, & Dodds, 2015)、日本語のフレーズでもそれが再現されるかどうかを検証することが必要になるだろう。

## 5 結言

本研究では情報通信端末に実装されている予測変換機能によって会話におけるメッセージの送受信者間の通信努力バランスが偏っているのか、また、予測変換機能により利用する単語の種類が低下し、その結果として語彙力が低下しているかを調べる実験を行った。本研究では大学生とその友人をペアの実験参加者として募集し、スマートフォンと LINE アプリを用いた予測変換条件と非予測変換条件、および情報通信端末を使用しない筆記条件、一般的な発話条件の 4 条件を設定し、それぞれの条件下でテーマを定めて自由に会話を行わせた。その会話記録を電子データ化し、形態素解析を行って出現頻度分布を作成し、その分布の形状と回帰直線の傾きを調べた。その結果、すべての条件においてほぼ 100%の説明率でベキ乗分布に従うことが判明し、使用語彙の偏りが明らかとなった。加えて、回帰直線の傾きが予測変換条件では 1.13、非予測変換条件では 1.14、筆記条件では 1.17、発話条件では 1.41 となり、すべての条件下においてジップの法則にお

ける1.0よりも数値が大きく、送受信者間の通信労力バランスが崩れていることが示唆された。本研究の最初の着眼ポイントは、同じ単語を繰り返し使用する原因となる予測変換機能が送受信者間の通信努力のバランスを崩していることや、我々の語彙力低下等、コミュニケーションに重要な国語能力にネガティブな影響を与えていること等を明らかにし、ひいては、情報通信端末における予測変換機能、つまり現在発展中である我が国の情報通信技術がもたらす功罪をつまびらかにすることであった。そのポイントに立脚した上で研究結果から結論できることは、3つある。1つ目は、すべてのコミュニケーション条件において通信バランスが受信者に高負荷となるように崩れていることである。送信者側が高負荷となるような崩れ方になっていないことから、言語を介したコミュニケーションはそもそも受信者が送られてくる言語情報に様々な類推を加えて成立している可能性がある。しかしながら、コミュニケーション自体は特に問題なく成立しているため、特に受信者側の負荷が深刻な問題になることはないだろう。2つ目は、予測変換条件と非予測変換条件のグラフの傾きがほぼ同じであったことから、予測変換の使用は語彙低下に直結するわけではないということである。これは情報通信技術における「功」の部分である。3つ目は、ノンバーバル情報の重要性である。会話条件では衝立などを使用せず、実験参加者同士が互いの顔を見ることができた。従って、言語情報のみならず、ノンバーバル情報を同時に送信することとなってコミュニケーションが円滑にすすんだ結果、大きなグラフの傾き値、つまり送信者の負荷がかなり低くなったと考えられうる。言い換えると、ノンバーバル情報のおかげで特定の語彙しか使用しなくても受信者に対して十分な情報を送信することができたのであろう。総じて、情報通信端末の予測変換機能は語彙量を低下させずにコミュニケーションを実現させる有用なツールであって、その力は従来の手紙等によるコミュニケーションとほぼ同等となり、情報通信技術における「功」の部分が改めて浮き彫りとなった。

## 【参考文献】

- Barabási, A. L. (2003). *Linked*. NY: Penguin Books.
- Buchanan, M. (2002). *Nexus: Small worlds and the groundbreaking science of networks*. NY: W. W. Norton & Company.
- Corominas-Murtra, B., Fortuny, J., & Solé, R. V. (2011). Emergence of Zipf's law in the evolution of communication. *Physical Review E*, *83*, 036115.
- Ferrer-i-Cancho, R. F., & Solé, R. V. (2003). Least effort and the origins of scaling in human language. *Proceedings of the National Academy of Sciences of the United States of America*, *100*, 788-791.
- 竹原卓真・佐藤直樹 (2003). 顔文字の有無によるメッセージの印象の違いについて 日本顔学会誌, *3*, 83-87.
- 竹原卓真・栗林克匡・水岡郁美・瀧波恵美子 (2005a). 顔文字の多用は逆効果！ —謝罪状況時に付加する顔文字の個数および種類と印象形成の関係— 日本顔学会誌, *5*, 21-32.
- 竹原卓真・栗林克匡・武川直樹・水岡郁美・瀧波恵美子 (2005b). メッセージの感情と矛盾した顔文字の付加効果 日本顔学会誌, *5*, 75-89.
- Takehara, T., Ochiai, F., & Suzuki, N. (2015). Scaling laws in emotion-associated words and corresponding network topology. *Cognitive Processing*, *16*, 151-163.
- Thurner, S., Hanel, R., Liu, B., & Corominas-Murtra, B. (2015). Understanding Zipf's law of word frequencies through sample-space collapse in sentence formation. *Journal of The Royal Society Interface*, *12*, 20150330.
- Williams, J. R., Lessard, P. R., Desu, S., Clark, E. M., Bagrow, J. P., Danforth, C. M., & Dodds, P. S. (2015). Zipf's law holds for phrases, not words. *Scientific Reports*, *5*, 12209.
- 横川壽彦 (2002). ジップの法則 日本ファジィ学会誌, *14*, 604.
- Zipf, G. K. (1949). *Human behavior and the principle of least effort*. Addison-Wesley, Cambridge.

〈発表資料〉

題名	掲載誌・学会名等	発表年月
Scaling laws in emotion-associated words and corresponding network topology.	Cognitive Processing	2015年5月
A big data study on how Zipf's law governs multiplex emotion-associated words	International Non-linear Science Conference	2017年4月