

会話を中心とした超短編小説の自動生成

Dialogue-Based Generation of Very Short Stories

高木 大生*¹ 佐藤 理史*¹ 駒谷 和範*²
Daiki Takagi Satoshi Sato Kazunori Komatani

*¹名古屋大学大学院 工学研究科 電子情報システム専攻
Graduate School of Engineering, Nagoya University

*²大阪大学 産業科学研究所
The Institute of Scientific and Industrial Research, Osaka University

This paper describes a scenario that generates very short stories, where a new story is produced from an existing dialogue in a story by replacing some utterances with different utterances in a different story. To realize this scenario, we have implemented a system that identifies a speaker of each utterance of dialogues in Shinichi Hoshi's short stories.

1. はじめに

コンピュータは、人間を楽しませる小説を作ることができるだろうか。人工知能の守備範囲を理性から感性へ拡大する試みとして、グランドチャレンジ「きまぐれ人工知能プロジェクト『作家ですよ』」が2012年秋にスタートした[松原 13]。このプロジェクトでは、ショートショートの大匠である星新一に学び、コンピュータにショートショートを創作させることを目指すプロジェクトである。

ショートショートの自動創作には、多くの機能の自動化が必要である。その中で、最も大きな障害になると予想されるのが、文章生成、すなわち、「人間が読んで意味を取ることができ、かつ、違和感を感じさせない段落以上のテキストを作成すること」である。日本語の文章生成の研究は、解析の研究と比べて全く遅れており、現状では、文生成のツールさえ存在しない。

このような認識に立ち、我々は、作成する小説の長さを、140字、400字、1200字と段階的に長くしていく戦略を取る。最初の140字はツイッター小説と呼ばれ、ツイッターの登場とともに出現した小説形式である。すでに3回のコンテスト(ツイッター小説大賞 <http://twonovel.net/>)が実施され、小説集[内藤 09]も出版されている。

本論文では、140字以内の超短編小説を、会話を中心に組み立てる方法を検討する。会話は、小説を構成する重要な要素の一つであり、ロバート・B・パーカーのスペンサー・シリーズのように、会話が魅力となっている小説も数多く存在する。さらに、第2回ツイッター小説大賞の優秀賞の1つは、下記のように、会話を中核とした作品である。

結婚が決まった私は、母に料理を習うことにした。
「ママ、これ塩多すぎない？ それに油でギトギトだし…。毎日こんなじゃ彼が早死にしちゃうよ」
「あら、ごめんなさい。でも覚えておいて損はないわよ」 (@shortxshort 作. 第2回ツイッター小説大賞ホームページより)

連絡先: 高木大生, 名古屋大学大学院 工学研究科 電子情報システム専攻, 〒4648603 愛知県名古屋市千種区不老町 C3-1(631) IB 電子情報南棟 159, 052-789-4435, d.takagi@nuee.nagoya-u.ac.jp

以下、まず2節で、会話を中心とした超短編小説の生成シナリオを示す。ここでは、会話をゼロから作るのではなく、星新一のショートショートに現れる会話を組み替えて、新たな作品を生成する。次に、3節では、その実現に必要な発話者特定の方法の概要について述べる。4節と5節では、その方法を構成する2つのモジュールについて述べる。最後に6節で、現状と課題をまとめる。

2. 会話を中心とした超短編小説の自動生成法

我々が想定する超短編小説の生成シナリオを、具体例を交えながら示す。このシナリオは、Step 0 から Step 3 までの4つのステップで構成される。

Step 0 準備

既存のショートショートを収集する。それぞれのショートショートに現れる各発話に対して、発話者をあらかじめ特定しておく。

映画や演劇の台本と異なり、小説では各発話の発話者は明示的に記述されない(図1)。しかしながら、発話を組み替えるためには、誰がその発話を発したのかの情報が不可欠である。このため、図2に示すような形で、発話者をあらかじめ特定しておく。

Step 1 利用する作品の決定

超短編小説の作成に利用する2つの作品(作品Xと作品Y)を選択する。

ここでは、作品Xとして星新一の『泉』(星新一ショートショート1001[星98]第1巻pp729-732)を、作品Yとして『小さくて大きな事故』([星98]第1巻pp613-617)を選択する。

Step 2 利用する会話文の決定

作品Xから、2人の登場人物AとBが交互に発話している4発話を抜き出す。この4発話を、 $A_1B_1A_2B_2$ と記述する。

発話の抜き出し方には、色々な方法がある。ここでは、『泉』の先頭から4発話を抜き出す。これを図3に示す。

「ねえ、ちょっと。起きてよ」
 男は妻にゆり起こされた。彼が、都心ちかくに新しく建てられた、さほど大きくはないがスマートなビルの管理人として雇われ、その地階の一室に夫婦でねとまりするようになってから、何日目の真夜中のことだった。
 「どうしたんだ」
 「いま、あたしがお便所から出ようとした時にね、だれかに背中をぼんとたたかれたような気がしたのよ」
 「泥棒かな」

図 1: 『泉』の冒頭部

妻:「ねえ、ちょっと。起きてよ」
 男は妻にゆり起こされた。彼が、都心ちかくに新しく建てられた、さほど大きくはないがスマートなビルの管理人として雇われ、その地階の一室に夫婦でねとまりするようになってから、何日目の真夜中のことだった。
 男:「どうしたんだ」
 妻:「いま、あたしがお便所から出ようとした時にね、だれかに背中をぼんとたたかれたような気がしたのよ」
 男:「泥棒かな」

図 2: 『泉』の冒頭部 (発話者特定後)

妻:「ねえ、ちょっと。起きてよ」
 男:「どうしたんだ」
 妻:「いま、あたしがお便所から出ようとした時にね、だれかに背中をぼんとたたかれたような気がしたのよ」
 男:「泥棒かな」

図 3: 『泉』の冒頭部の 4 発話

彼女:「ねえ。お話があるの……」
 夫:「なんだ。言ってみろ」
 彼女:「いいかげんで、あたしと別れてくださらない……」
 夫:「なんの話かと思ったら、そんなことか。ばかばかしい。おまへはおれにとって、遊んで暮すための大事な金づるだ。別れてやるわけにはいかないぜ」

図 4: 『小さくて大きな事故』の 4 発話

「ねえ、ちょっと。起きてよ」
 「どうしたんだ」
 「いいかげんで、あたしと別れてくださらない……」
 「なんの話かと思ったら、そんなことか。ばかばかしい。おまへはおれにとって、遊んで暮すための大事な金づるだ。別れてやるわけにはいかないぜ」

図 5: 生成結果 1

「だれだとは、なんです。あなたのご要望にこたえて出現した悪魔ですよ」
 「これはありがたい。みごとに成功したようだな。こううまくゆくとは……」
 「そうですよ。理屈もなにもない無茶な願いを、真剣になってとねえた。そこがわたしの気に入った点です。まったく、そういう人が少なくなつた。そういう人を相手にするのが、わたしの働きがいなのに。ことが不合理であればあるほど、悪魔のほうもやって楽しんでますよ」
 「もちろんでございます。ぜひ、わたしをお助け下さい」

図 6: 生成結果 2

Step 3 発話の入れ替え

作品 X から抜き出した 4 発話中の 2 発話を、作品 Y の 2 発話で組み替える。このとき、組み替える発話には、似たような属性を持つ発話者の発話を選ぶ。

発話の組み替え方として、当面、次の 2 種類を考える。作品 Y に現れる登場人物を C と D とするとき、

人物の組み替え $A_1B_1A_2B_2$ のうち、同一人物の発話を組み替える。すなわち、 $A_1D_1A_2D_2$ 、または、 $C_1B_1C_2B_2$ を作成する。

後半部の組み替え $A_1B_1A_2B_2$ のうち、後の 2 発話を組み替える。すなわち、 $A_1B_1C_2D_2$ を作成する。

図 4 に、作品 Y 『小さくて大きな事故』に現れる連続した 4 発話を示す。これを $C_1D_1C_2D_2$ とする。Step 2 で作成した $A_1B_1A_2B_2$ の後半部の A_2B_2 を、 C_2D_2 で組み替えると、図 5 に示す新しい会話を得られる。文字数は 118 文字である。

図 6 に、人物を組み替える方法で生成した作品を示す。この作品は『魔法の大金』([星 98] 第 1 巻 pp1410-1411) と、『条件』([星 98] 第 3 巻 pp65-68) から作成されている。

図 5 や図 6 の会話は、面白いかどうかはさておき、いちおう、会話として成立している。すなわち、異なる作品の発話を組み替えることにより、新たな会話 (作品) をつくり出すことは可能である。

当然のことながら、人間は、超短編小説の創作に、このような方法を用いない。なぜならば、このままでは、あきらかに剽窃となるからである。しかしながら、我々は、機械による小説生成は、まず、このレベルの習作からスタートする必要があると考える。この後、発話の一部を変更したり、地の文を挿入し

たりすることによって、剽窃の度合を薄めていくことが可能である。

3. 発話者特定システム

前節のシナリオを実現するためには、準備の Step 0 で、すべての発話に対して、あらかじめ発話者を特定しておく必要がある。なぜならば、Step 3 の発話の組み替えを盲目的に行なった場合、意味のある会話を得られる可能性は低いからである。先に示した例では、作品 X の「妻 (A)」と「男 (B)」、作品 Y の「彼女 (C)」と「夫 (D)」において、 $A \approx C$ 、 $B \approx D$ という関係が成立しているため、発話の組み替え後も会話として成立するのである。

図 7 に、今回作成した発話者特定システムの構成図を示す。本システムは、登場人物が 2 名である作品のみを対象とする。この図に示した通り、システムは、登場人物抽出モジュールと発話者特定モジュールから構成されている。これら 2 つのモジュールの詳細を、後述する 2 つの節で述べる。

4. 登場人物抽出モジュール

発話者特定の最初のステップでは、入力された作品のテキストに登場する人物 2 名を指し示す表現 (登場人物ラベル) を抽出する。

4.1 アルゴリズム

星新一の作品において、登場人物を表す表現は、大きく 4 種類に大別できる。

1. 人名の固有名詞
2. 人称代名詞
3. 人を表す名詞 (「博士」など)

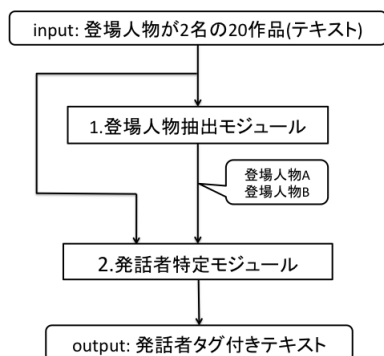


図 7: 発話者特定システムの構成

表 2: 登場人物抽出の結果の一部

タイトル	登場人物 (正解)		抽出した結果	
	A	B	A	B
悪魔	エス氏	悪魔	エス氏	悪魔
暑さ	巡査	男	巡査	男
妖精	女の子ケイ	妖精	彼女	妖精
愛の通信	男	宇宙人の女性	*彼女	女
告白	花屋の女	青年	女	青年
帰路	私	通信相手	私	*—
再認識	社長	私	社長	私
初夢	男	サービスマン	男	相手
泉	妻	夫	妻	男
宝船	エヌ氏	福の神	エヌ氏	福の神

表 1: 登場人物検出用辞書

中項目	分類項目	具体例
人間	神仏・精霊	悪魔, 妖精
	男女	男, 女
	老少	子供, 老婆
	人物	先生, 博士
家族	夫婦	妻, 夫
	親・先祖	お父さん, お母さん
	子・子孫	孫
	兄弟	兄, 妹
仲間	親戚	いとこ, めい
	相手・仲間	相手, 同僚
	友・なじみ	友達, 恋人
成員	主客	主人, 主婦
	専門的・技術的職業	医者, 女優
	管理的・書記的職業	大臣, 助手
	実業・商業	起業家, 経営者
	農林水産業	漁師, セールスマン
	運輸業	パイロット, 運転士
	職人	コック, パーティンダー
	保安サービス	巡査, 消防士
	サービス	家政婦, モデル
	学徒	大学生, 高校生
	長	社長, 会長
その他の仕手	リーダー, 目撃者	

表 3: 登場人物抽出の正解率

登場人物の総数に対する正解率	42/46	91%
作品の総数に対する正解率	20/23	87%

4.2 評価

星新一のショートショート 23 作品に対して、登場人物抽出モジュールを実行した結果の一部を、表 2 に示す。この表で“*”が付与されているものは、抽出誤りまたは抽出失敗を表す。

一般に、登場人物の記述は、作品中で一貫しているわけではない。そのため、作品中で、その記述により登場人物が一意に定まる場合は正解とした。例えば、作品『妖精』で、「女の子ケイ」の抽出結果は「彼女」であるが、作品中で「彼女」は、「女の子ケイ」を意味するので正解とした。

登場人物の総数に対する成功率と、作品に対する成功率を表 3 に示す。この表に示すように、このモジュールの精度は、90%程度である。この結果を言い換えるならば、星新一の作品では、発話の前後の地の文に、発話者が明示的に書かれることが多いということである。

現在のモジュールは、照応解析を行っていないため、「彼女」のような人称代名詞が実際に誰であるかを特定できない。また、対象とする作品は、登場人物が 2 名である作品に限定されている。これらの問題を解決するためには、作品の流れに沿って登場人物を発見し、かつ、人物記述表現の同一性を同定していく登場人物トラッキングを実現する必要がある。

4. 人間以外を表す名詞 (「悪魔」など)

これらのうち、1 は、形態素解析 (MeCab+IPAdic) の出力より検出することができる。

残りの 3 種類の人物記述表現を検出するために、登場人物検出用辞書を作成した。この辞書は、分類語彙表^{*1}から、『体』の『主体』に属する分類のうち、中項目、分類項目を用いて作成した。登場人物検出用辞書の一部を表 1 に示す。

この辞書を用いた、登場人物抽出アルゴリズムを以下に示す。

1. 作品に含まれる全ての発話の直前・直後の地の文を抽出し、形態素解析する。
2. 形態素解析結果と登場人物検出用辞書により人物記述表現を検出し、その直後が「は」「が」のいずれかであった場合に、その数をカウントする。
3. 出現回数の上位 2 つを登場人物を表すラベルとして出力 (採用) する。

このアルゴリズムでは、人物記述表現の誤検出を避けるために、検出対象を発話の直前・直後の地の文に限定し、かつ、後続形態素が「は」「が」である場合に限定している。

*1 <http://www.ninjal.ac.jp.archives/goihyo>

5. 発話者特定モジュール

発話者特定モジュールでは、それぞれの発話に、発話者ラベルを付与する。

5.1 アルゴリズム

本アルゴリズムは、まず、発話者の特定が容易な発話に、発話者ラベルを付与する。これには、後述する 4 つのルールを用いる。次に、付与された発話者ラベルを使って、前後の (発話者が同定されていない) 発話に、発話者ラベルを付与する。これを発話者ラベルの伝搬と呼ぶ。原則として、2 名の会話では、発話者は交互に交替するので、この性質を利用する。この性質は、発話が連続する場合には、ほぼ間違いなく成立する。しかしながら、2 つの発話間が離れていくにつれて、成り立たなくなる。そのため、発話間に距離を導入し、その距離が短いものから優先して、発話者ラベルを付与する。なお、例外的に、2 つの発話間に地の文が 1 文存在し、その文が「つづけて言った」という表現を含む場合は、2 つの発話に、同一の発話

表 4: 発話者ラベルの付与状況

	対象発話数	タグの付与数	正解数	Precision
ルール 1	653	22	21	95%
ルール 2		15	14	93%
ルール 3	65	34	33	97%
ルール 4	78	48	48	100%
ラベル伝搬	534	482	368	76%
全体	653	601	484	81%

者ラベルを付与する。

上記のアルゴリズムでは、発話者ラベルを伝搬させる前に、できるだけ多くの発話の発話者を正しく同定しておくことが重要である。これには、以下の4つのルールを用いる。なお、 S は、登場人物抽出モジュールで抽出された発話者ラベルを示す。

ルール 1 「 S は言った」などの、明示的に発話者を示す地の文があり、かつ、その前後どちらか一方のみが発話である場合、その発話に発話者ラベル S を付与する。

ルール 2 発話の次の地の文に、『 S は「...」と、言った』『「...」と、 S は』という表現があるときには、発話「...」に発話者ラベル S を付与する。

ルール 3 登場人物が男性と女性の場合、女性の役割語を持つ発話に女性の登場人物ラベルを付与する。

ルール 4 登場人物が男性と女性の場合、男性の役割語を持つ会話文に男性の登場人物ラベルを付与する。

ルール 3 の女性の役割語としては、文献 [中村 13] を参考に、文末の 2 文字のいずれかが「わ」である、文末の「のよ」「かしら」「そうね」、発話中の「あたし」を用いた。ルール 4 の男性の役割語としては、文献 [小川 06] を参考に、文末の「だ」「ぞ」「ぜ」「ものか」「だろう」「だい」「のさ」「くれ」「だが」、発話中の「ぼく」「おれ」を用いた。

5.2 評価

発話者特定モジュールを、登場人物の抽出に成功した 20 作品に対して実行した。表 4 に発話者ラベルの付与状況を示す。

表 4 から、ルール 1 からルール 4 は、それぞれ、高い精度で発話者のタグを付与できていることが分かる。発話の役割語を利用するルール 3 とルール 4 は、ルール 1 とルール 2 に比べ、より多くの発話に、発話者ラベルを付与できる。しかし、これらのルールは、男性と女性の 2 人が登場する作品に対してのみしか機能しない。

図 8 に、20 作品のそれぞれの recall 値の分布を示す。

表 4 と図 8 から分かることは 2 つある。1 つは発話者特定モジュールの recall は 74% (484/653) であることである。2 つ目は、個々の作品の recall の分布を見ると、高い recall で発話者を特定できる作品と、低い recall でしか発話者を特定できないものの差が激しいことである。著しく recall の低い作品が、全体の正解率を低下させている。

発話者ラベルの伝搬は、局所的な発話群に対しては正しく動作するが、段落をまたいだ発話間には、正しく動作しない場合が多い。このため、アルゴリズムの最初の段階で用いるルールを強化し、より多くの発話の発話者を同定する必要がある。

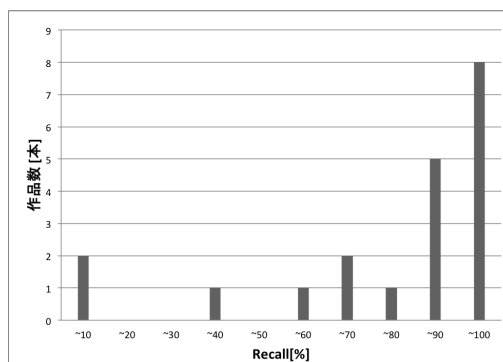


図 8: Recall の分布

6. まとめ

本論文では、会話を中心とした超短編小説の生成シナリオを示し、その実現に必要な発話者特定法を示した。会話を組み替えてストーリーを作る方法は、文章生成の問題を回避しつつ、140 字程度の作品を作成できる可能性がある。組み替えた会話が、会話として成立しやすいのは、発話が独立した単位であることと、地の文よりも前後との整合性に関する制約が緩いことに、その理由があると考えられる。

謝辞

本研究では、星新一ショートショート 1001(新潮社) のテキストデータを使用した。データを提供していただいた、星ライブラリ、および新潮社に感謝します。

参考文献

- [星 98] 星新一: 星新一ショートショート 1001, 新潮社 (1998)
- [松原 13] 松原 仁, 佐藤 理史, 赤石 美奈, 角 薫, 迎山 和司, 中島 秀之, 瀬名 秀明, 村井 源, 大塚 裕子: コンピュータに星新一のようなショートショートを創作させる試み, JSAI2013 論文集 (2013)
- [内藤 09] 内藤 みか, 安達 瑠 乙, 新城 カズマ, 小林 正規, 渡辺 やよい, 吉井 春樹, 泉 忠司, 黒崎 薫, 柊野 浩一, 円城 塔: twitter 小説集 140 字の物語, ディスカバートウエンティワン (2009)
- [中村 13] 中村 桃子: 翻訳が作る日本語, 白澤社 (2013)
- [小川 06] 小川 早百合: 話し言葉の終助詞の男女差の実際と意識 -日本語教育での活用へ向けて-, 日本語ジェンダー学会 (編), 日本語とジェンダー, ひつじ書房 (2006)