

文字列中からの単語の発見と感覚情報に基づく 単語の意味づけを通じた SIR 名付けゲームによる言語の創発

堀江孝文¹ 谷口彰² 萩原良信³ 谷口忠大²¹ 立命館大学大学院 情報理工学研究科 ² 立命館大学 情報理工学部³ 立命館大学 総合科学技術研究機構

{horie.takafumi,a.taniguchi,yhagiwara,taniguchi}@em.ci.ritsume.ac.jp

概要

相互分節化仮説とは、複数のエージェントによる文のやり取りを通じ、文の分節化と状況の分節化が同時に行われることで、人間の言語が誕生したとする仮説である。本研究では、相互分節化仮説に着想を得て、観測した感覚情報を表現する文字列のやり取りを通じ、エージェントが文字列中から単語を発見し、単語と分類された感覚情報の関係を推論することで、言語が創発する過程のモデル化を目指す。そこで本稿では、教師なしで文を分節化するためのモデル NPYLM と、色や形といった複数の属性を含む感覚情報と単語との結びつきを学習する交差状況学習のモデル CSL-PGM を統合してエージェントを構成し、二体のエージェントを接続したモデル Inter-CSL+NPYLM を構築する。また、エージェント間における推論を、文のやり取りに基づく言語ゲームである SIR 名付けゲームにより実現する。実験では、提案手法が比較手法よりもエージェント間の発話の編集距離を小さくすることが確認された。一方で、提案手法が文の分節化誤りの影響を大きく受けたことで、エージェント間で多くの単語が共有されなかったことも明らかになった。

1 はじめに

相互分節化仮説とは、文を構成する文字列の一部と分節化された状況が対応付けられることで、言語が誕生したとする仮説である [1]。この仮説では、複数のエージェントが共有された状況の中で発話する中で、発話された文の分節化と、状況の分節化が同時に行われていき、言語が創発したと考える。本研究では、この相互分節化仮説に着想を得て、分節化されていない文字列のやり取りを通じ、エージェントが文字列を分節化しつつ、複数の属性を含む感

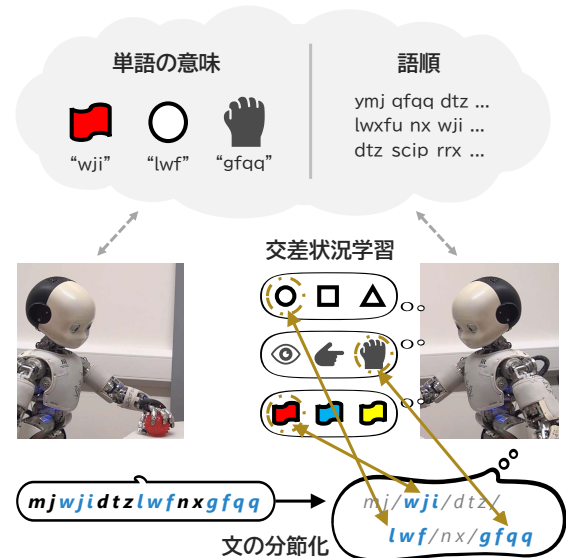


図 1: 本研究で扱う言語の創発の概要。分節化されていない文字列のやり取りを通じ、エージェントが単語を発見しつつ、単語と感覚情報を結びつける。

覚情報のカテゴリゼーションを行い、単語と感覚情報を結びつけていくことで、言語が創発する過程のモデル化を目指す (図 1)。

切れ目のない文字列のやり取りを通じた言語の創発をモデル化するためには、文を分節化すること、状況を分節化すること、分節化によって得られた単語を感覚情報と対応付けることが必要である。これに関し、言語の創発の計算論的なモデル化を扱ったこれまでの研究が多数提案されている [2]。Chaabouni らは、深層強化学習に基づいた言語ゲームを通して、単語と感覚情報が結びつけられていく過程をモデル化した [3]。また萩原らは、エージェントが複数の単語を組み合わせることで観測を説明する言語ゲームを通じて、記号が創発する過程をモデル化した [4]。しかし、これらのモデルは用いる単語の種類を制限しており、また切れ目のない文によるコ

コミュニケーションを扱っていなかった。

本研究では相互分節化仮説における文の分節化と状況の分節化を同時にモデル化するため、NPYLM [5] と Cross-situational learning with Bayesian probabilistic generative model (CSL-PGM) [6] という二つのモデルを統合する。NPYLM は、教師なしで文字列を分節化することができる言語モデルである。本研究では、例えば図 1 で文を “mj gfqq dtz lwf nx wji” のように単語へと分割するために用いる。CSL-PGM は、交差状況学習に基づき、感覚情報を分類しつつ単語と結びつく属性・感覚情報を推論するモデルである。ここで交差状況学習とは、単語と感覚情報の共起を学習することで、単語と結びついている属性（例: “wji” は色属性, “lwf” は形属性など）と、感覚情報のカテゴリ（例: “wji” は赤, “lwf” は球など）を推論することである。

本稿では、CSL-PGM と NPYLM を統合したエージェントを構成し、そのエージェント二体を接続したモデル **Inter-CSL+NPYLM** を提案した。また、Inter-CSL+NPYLM 上での推論を、観測した感覚情報に結び付いた文をやり取りする言語ゲーム Sample Importance Resampling 名付けゲーム (**SIR 名付けゲーム**, **SIRNG**) により定式化した。

2 提案手法: Inter-CSL+NPYLM with SIRNG

2.1 Inter-CSL+NPYLM

図 2 に Inter-CSL+NPYLM のグラフィカルモデルを示す。また、表 1 に各変数のうち後述の SIR 名付けゲームに関係するものの定義を示す。その他の変数の定義については Appendix を参照されたい。なお、 $* \in \{A, B\}$ は任意のエージェントを示す記号である。

本研究では、隠れ単語列 w_d^A, w_d^B および文字列 (character sequence) ℓ_d の生成過程を、unigram rescaling [7] を用いて以下のように近似する。なお以下では、CSL-PGM に対応する部分の変数のうち、 w_d^* 以外のものに関する集合を Ω^* と表記する。

$$\begin{aligned} & P(\ell_d, w_d^A, w_d^B \mid \mathbf{v}^A, T^A, \mathbf{v}^B, T^B, \Omega^A, \Omega^B) \\ \approx & \frac{P(\ell_d \mid w_d^A, \mathbf{v}^A) P(w_d^A \mid T^A) P(w_d^A \mid \Omega^A)}{\prod_{d,s} P(\ell_{d,s}) \prod_{d,n} P(w_{d,n}^A)} \\ & \times \frac{P(\ell_d \mid w_d^B, \mathbf{v}^B) P(w_d^B \mid T^B) P(w_d^B \mid \Omega^B)}{\prod_{d,n} P(w_{d,n}^B)} \quad (1) \end{aligned}$$

この近似により、 ℓ_d, w_d^A, w_d^B の生成確率を各エー

ジェントの CSL-PGM 部分・NPYLM 部分の確率によって計算することができる。なお、 $P(\ell_d \mid w_d^*, \mathbf{v}^*), P(w_d^* \mid T^*), P(w_d^* \mid \Omega^*)$ の具体的な計算式については Appendix を参照されたい。

ℓ_d, w_d^A, w_d^B 以外の変数のうち、CSL-PGM 部分に対応する変数の生成・推論過程は、オリジナルの CSL-PGM と同様である [6]。また、NPYLM 部分に対応する変数は、階層 Pitman-Yor 過程 [5] により生成される。

2.2 SIR 名付けゲーム

Inter-CSL+NPYLM では、文のやり取りを通じたエージェント間の推論を、SIR 名付けゲームとしてモデル化する。これは、確率的生成モデルを接続する SERKET フレームワーク [8,9] における、SIR 法に基づいた手続きである。なお以下では、エージェント A, B のうち話し手を Sp, 聞き手を Li とする。

SIR 名付けゲームは、式 (1) を重点サンプリングと SIR 法により近似する。SIR 名付けゲームの具体的な手続きを以下に示す。

- (i) 話し手側エージェントは、CSL-PGM 部分により、 w_d^{Sp} の仮のサンプル $\bar{w}_d^{\text{Sp}(r)}$ ($r = 1, \dots, R^{\text{Speak}}$) を得たのち、NPYLM 部分、および文字ユニグラム確率によって各サンプルの重み

$$\mathcal{W}_d^{(r)} = \frac{P(\bar{w}_d^{\text{Sp}(r)} \mid T^{\text{Sp}(r)})}{\prod_{d,s} P(\bar{\ell}_{d,s}^{(r)}) \prod_{d,n} P(\bar{w}_{d,n}^{(r)})} \quad (2)$$

を計算する。

- (ii) $\mathcal{W}_d^{(r)}$ に比例する確率で $\bar{w}_d^{\text{Sp}(r)}$ をリサンプリングし、発話される単語列のサンプル $\bar{w}_d^{\text{Sp}(q)}$ ($q = 1, \dots, Q$) を得る。なお、 $\bar{\ell}_d^{(r)}$ は $\bar{w}_d^{(r)}$ を \mathbf{v}^{Sp} によって変換し得られる文字列である。
- (iii) 話し手側の単語列 $w_d^{\text{Sp}(q)}$ を $\mathbf{v}^{\text{Sp}(q)}$ により変換し、切れ目のない文字列 ℓ_d のサンプル $\ell_d^{(q)}$ ($q = 1, \dots, Q$) を得る。
- (iv) 聞き手側エージェントは $\bar{\ell}_d^{(1)}, \dots, \bar{\ell}_d^{(Q)}$ をそれぞれ分節化し w_d^{Li} の仮のサンプル $\bar{w}_d^{\text{Li}(q,r)}$ ($r = 1, \dots, R^{\text{segment}}$) を得たうえで、各 $\bar{w}_d^{\text{Li}(q,r)}$ に対し重み

$$\begin{aligned} \mathcal{W}_d^{(q,r)} = & \frac{P(\bar{w}_d^{\text{Li}(q,r)} \mid \bar{T}^{\text{Li}(q)})}{P(\bar{w}_d^{\text{Li}(q,r)} \mid \bar{\ell}_d^{(q)}, \bar{\mathbf{v}}^{\text{Li}(q)}, \bar{T}^{\text{Li}(q)})} \\ & \times \frac{P(\bar{w}_d^{\text{Li}(q,r)} \mid \Omega^{\text{Li}(q,r)})}{P(\bar{w}_d^{\text{Li}(q,r)} \mid \bar{\ell}_d^{(q)}, \bar{\mathbf{v}}^{\text{Li}(q)}, \bar{T}^{\text{Li}(q)})} \quad (3) \end{aligned}$$

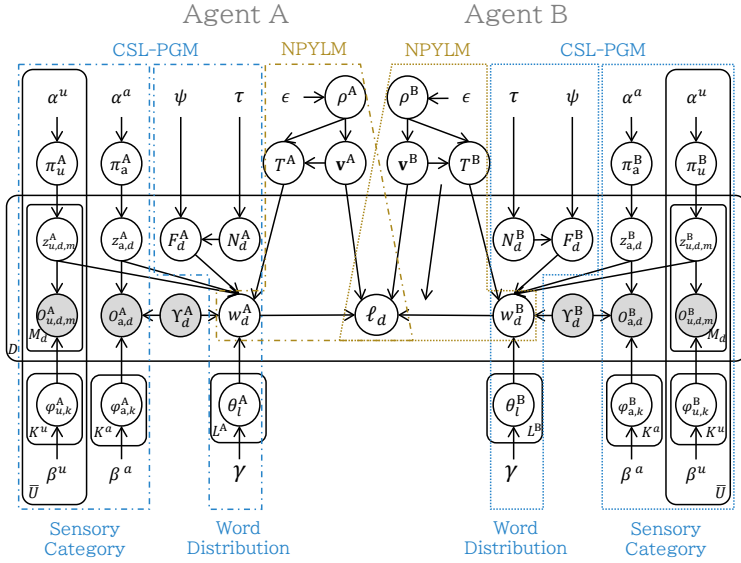


図 2: Inter-CSL+NPYLM のグラフィカルモデル

を計算する。ここで、 $\Omega^{\text{Li}(q,r)}$ は $\bar{w}_d^{\text{Li}(q,r)}$ をすべて使い、 Ω^{Li} のうち $F_d^{\text{Li}}, \theta_l^{\text{Li}}$ を Gibbs sampling で学習して得られる、 Ω^{Li} のサンプルである。

(v) 求めた $\mathcal{W}_d^{(q,r)}$ に比例する確率で $\bar{\ell}_d^{\text{Li}(q)}$, $\bar{w}_d^{\text{Li}(q,r)}$ をリサンプリングし、最終的なサンプル $\ell_d^{\text{Li}(q)}, w_d^{\text{Li}(q)}$ ($q = 1, \dots, Q$) を得る。

聞き手は、以上の手続きで得たサンプル $\{\ell_d^{\text{Li}(q)}, w_d^{\text{Li}(q)}\}_{q=1}^Q$ を用いて CSL-PGM 部分、NPYLM 部分のパラメータを Gibbs sampling で学習する。

本稿では、各エージェントにおいて ℓ_d^*, w_d^* を a から z のアルファベットで構成されるランダムな文字列で初期化し、 w_d^* と感覚情報を用いてあらかじめ CSL-PGM 部分を学習させたうえで凍結する。そのうえで、両エージェントが聞き手と話し手の立場を入れ替えながら SIR 名づけゲームのステップ (i)-(v) と凍結されていない部分の Gibbs sampling を繰り返すことで、モデル全体の学習を実現する。

3 実験

本研究では、各エージェントが感覚情報と単語を対応付けられているか、またエージェントの間で単語とその語順が共有されているかを検証するため、実験を行った。エージェント間で単語・語順が共有されているかについては、同じ状況に対する発話の類似度を評価した。具体的には、編集距離を比較対象のうち長い方の文の文字数で割って得られる値（以下、文字あたり編集距離）と、各エージェントが発話した単語列における tri-gram 確率の間の JS ダイ

表 1: Inter-CSL+NPYLM の各変数の定義

N_d^*	ℓ_d 中の単語の数
ℓ_d	d 番目の観測に対応する文字列 $(\ell_{d,1}, \dots, \ell_{d,s}, \dots, \ell_{d, \ell_d })$
w_d^*	ℓ_d を生成する、単語のインデックスの列 $(w_{d,1}^*, \dots, w_{d,n}^*, \dots, w_{d,N_d}^*)$
F_d^*	w_d^* に対応する属性の列 $(F_{d,1}^*, \dots, F_{d,n}^*, \dots, F_{d,N_d}^*)$
θ_l^*	感覚情報カテゴリ l での単語確率
T^*	w_d^* を生成する n-gram 確率 $\{T_{t',t}^*\}_{v_{t',t}}$. $T_{t',t}^*$ は単語列 t' の直後における単語 t の確率
v^*	単語辞書 $\{v_i\}_{v_i}$. v_i は単語インデックス i に対応する文字列
ρ^*	T^*, v_i^* を生成するパラメータの集合

バージェンスを用いた。

3.1 実験条件

実験においては、萩原らの先行研究 [4] で用いられていた感覚情報のデータセットを用いる。このデータセットには、感覚情報とエージェント間の共同注視の情報が含まれている。感覚情報は、iCub シミュレータを用いて取得された、各属性（ロボットの触覚と体性感覚: Action, 物体の位置情報: Position, 物体の形状情報: Object, 物体の色情報: Color）の情報で構成されている。

実験では、三つの異なる手法を比較した。なおモデルのパラメータの設定については Appendix を参照されたい。

- **SIRNG**: 本研究の提案手法である Inter-CSL+NPYLM with SIRNG.
- **no communication**: Inter-CSL+NPYLM with SIRNG において、各エージェントが相手の発話を受け取らず、代わりに自分の発話を受け取る手法。エージェントがコミュニケーションを行っていない場合に相当する。
- **no listener resampling** Inter-CSL+NPYLM with SIRNG において、Speaker エージェントが相手の発話を受け取った際、文を単語へ分節化した後の resampling を行う手法。聞き手が相手の発話をそのまま受け入れる場合に相当する。

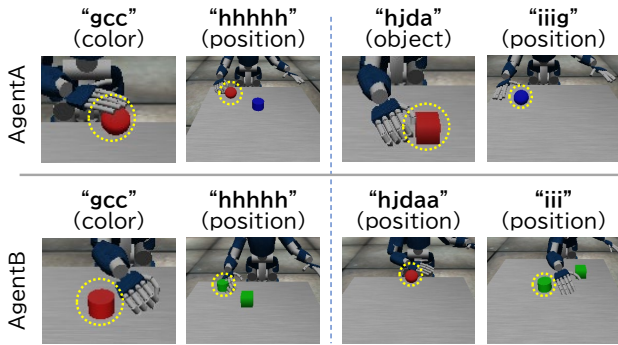


図 3: 各エージェントの発話に含まれていた単語の例。()内は単語の属性を示す。一部の単語が共有されていた一方、似た綴りの語が異なる感覚情報と結び付けられる例が多く見られた。

3.2 実験結果と考察

3.2.1 感覚情報と単語の対応・共有

図 3 に、提案手法において学習後のエージェントが用いていた単語と、それに対応する状況の例を示す。実験では、図の“gcc”や“hhhhh”のように、両エージェントが一部の単語を共有し、同じ感覚情報を指す言葉として学習していることが確認された。しかし、大部分の単語は共有されず、同じ状況に対して各エージェントが全く異なる単語で会話する現象がみられた。

エージェント間で共有されていない単語に注目すると、片方のエージェントでつかわれている単語とよく似ているにもかかわらず、異なる感覚情報と結び付けられている単語が多数確認された。具体的には、図 3 の“hjda”と“hjadaa”のように異なる属性の感覚情報を表す語として学習されているものや、“iiig”と“iii”のように同じ属性の異なる感覚情報を表す語として学習されているものが存在した。

これは、Inter-CSL+NPYLM のうち、感覚情報と単語を結びつける CSL-PGM 部分において、単語の文字列としての類似度を考慮していなかったためだと考えられる。CSL-PGM 部分では、単語同士の文字列としての類似度を考慮せず、一字でも違う単語は全く別の単語として扱っていた。しかし、文のやり取りにおいては、聞き手の分節化誤りなどが発生し、本来同じである単語が異なる文字列で表現されることが考えられる。ゆえに、違う綴りの単語を全く異なる単語として扱ったことが、エージェント間における単語の意味の共有を妨げたと考えられる。

表 2: エージェントの発話の類似度に関する定量評価。“編集距離”はエージェント間の文字あたり編集距離、“JS 距離”は文字 tri-gram の JS ダイバージェンスを示している

method	編集距離↓		JS 距離↓	
	mean	t-test	mean	t-test
SIRNG	0.758	-	0.346	-
w/o listener resampling	0.775	n.s.	0.346	n.s.
w/o communication	0.843	**	0.347	n.s.

3.2.2 エージェント間での単語・語順の共有

表 2 に、エージェント間における発話の文字あたり編集距離、および単語 tri-gram 確率の間の JS ダイバージェンスを示す¹⁾。

表 2 を見ると、比較手法に比べ、提案手法の方が文字あたり編集距離が小さいことが分かる。このことから、聞き手が話し手の文をそのまま受け入れるのではなく、SIR 法に基づきリサンプリングすることが言語の共有を促すと考えられる。

一方で、どの手法も文字あたり編集距離が 0.7 を上回っていることや、エージェント間の単語 tri-gram 確率に関する JS divergence がほとんど同じになっていることも確認できる。これは 3.2.1 節での議論の通り、エージェント間で単語の意味の共有が十分にできていないためだと考えられる。

4 結論

本稿では、切れ目のない文字列によるコミュニケーションを通じ、エージェントが単語を発見しつつ感覚情報を分類し、単語と感覚情報を結びつけることで、言語が創発する過程をモデル化した。具体的には、CSL-PGM と NPYLM に基づきモデル Inter-CSL+NPYLM を構築した。また、エージェント間での推論を文のやり取りとして実現する SIR 言語ゲームを提案した。

今後の展望としては、エージェント間で単語が共有されない問題を解決するため、単語と感覚情報を結び付ける際に単語の文字列としての類似度を考慮することが挙げられる。また、相互分節化仮説 [1] による言語の創発をシミュレーションするため、エージェントによる感覚情報の分類を、コミュニケーションと同時に進めさせることも課題である。

1) mean はそれぞれ、5 回の実験で得た値の平均を示している。t-test は提案手法との t 検定の結果を示しており、p 値が 0.05 以上である場合を n.s.、0.05 未満である場合を *、0.01 未満である場合を ** と示している。

謝辞

本研究は JSPS 科研費 JP21H04904 および JP23H04835 の助成を受けたものです。

参考文献

- [1] Kazuo Okanoya and Bjorn Merker. Neural substrates for string-context mutual segmentation: A path to human language. In **Emergence of communication and language**, pp. 421–434. Springer, 2007.
- [2] Lukas Galke, Yoav Ram, and Limor Raviv. Emergent communication for understanding human language evolution: What’s missing? In **Emergent Communication Workshop at ICLR 2022**, 2022.
- [3] Rahma Chaabouni, Eugene Kharitonov, Diane Bouchacourt, Emmanuel Dupoux, and Marco Baroni. Compositionality and generalization in emergent languages. In **Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics**, pp. 4427–4442, Online, July 2020. Association for Computational Linguistics.
- [4] Yoshinobu Hagiwara, Kazuma Furukawa, Takafumi Horie, Akira Taniguchi, and Tadahiro Taniguchi. Symbol emergence as interpersonal cross-situational learning: the emergence of lexical knowledge with combinatoriality. **arXiv preprint arXiv:2306.15837**.
- [5] Daichi Mochihashi, Takeshi Yamada, and Naonori Ueda. Bayesian unsupervised word segmentation with nested pitman-yor language modeling. In **Proceedings of the ACL-IJCNLP**, pp. 100–108, 2009.
- [6] Akira Taniguchi, Tadahiro Taniguchi, and Angelo Cangelosi. Cross-situational learning with bayesian generative models for multimodal category and word learning in robots. **Frontiers in neurobotics**, Vol. 11, p. 66, 2017.
- [7] Daniel Gildea and Thomas Hofmann. Topic-based language models using em. In **Sixth European Conference on Speech Communication and Technology**, pp. 2167–2170, 1999.
- [8] Tomoaki Nakamura, Takayuki Nagai, and Tadahiro Taniguchi. Serket: an architecture for connecting stochastic models to realize a large-scale cognitive model. **Frontiers in neurobotics**, Vol. 12, p. 25, 2018.
- [9] Tadahiro Taniguchi, Tomoaki Nakamura, Masahiro Suzuki, Ryo Kuniyasu, Kaede Hayashi, Akira Taniguchi, Takato Horii, and Takayuki Nagai. Neuro-serket: development of integrative cognitive system through the composition of deep probabilistic generative models. **New Generation Computing**, Vol. 38, No. 1, pp. 23–48, 2020.

表 3: Inter-CSL+NPYLM の各変数の定義 (補足)

Ξ_d^*	エージェントが d 番目の観測で注視している物体の番号
$O_{a,d}^*$	d 番目の観測における Action 属性の情報
$O_{u,d,m}^*$	d 番目の観測における m 番目の物体についての属性 u の感覚情報
$z_{u,d,m}^*$	d 番目の観測における m 番目の物体での、属性 u のカテゴリのインデックス
π_u^*	$z_{u,d,m}^*$ を生成するカテゴリ分布のパラメータ
$\varphi_{u,k}^*$	カテゴリ k , 属性 u の感覚情報を生成するガウス分布のパラメータ
$\alpha^u, \beta^u, \gamma, \epsilon, \tau, \psi, \phi$	ハイパーパラメータ
D	観測の数
M_d	d 番目の観測に含まれる物体の数
K^u	属性 u におけるカテゴリの集合
\bar{U}	物体に関する感覚情報の属性の集合
L^*	すべての属性に対するカテゴリの集合

A 付録

A.1 Inter-CSL+NPYLM with SIRNG の補足

Inter-CSL+NPYLM の変数のうち、表 1 で省略したものの定義を表 3 に示す。また、SIRNG の流れをまとめたものを図 4 に示す。

A.2 近似式 (1) の詳細

式 (1) における $P(\ell_d | w_d^*, \mathbf{v}^*), P(w_d^* | T^*), P(w_d^* | \Omega^*)$ は、以下のように計算される:

$$P(\ell_d | \mathbf{v}^*, w_d^*) = \prod_{n=1}^{N_d^*} \prod_{s=1}^{S_{d,n}^*} \delta(\ell_{d,S_{d,(n-1)}^*+s}^*, v_{w_{d,n}^*,s}^*) \quad (4)$$

$$P(w_d^* | \Omega^*) = \prod_{n=1}^{N_d^*} P(w_{d,n}^* | \Omega^*) \quad (5)$$

$$P(w_d^* | T^*) = \prod_{n=1}^{N_d^*} \text{Cat}(w_{d,n}^* | T_{w_{d,n-1:n-\zeta-1}, w_{d,n}^*}^{\text{word}}). \quad (6)$$

ただし、

$$P(w_{d,n}^* | \Omega^*) = \begin{cases} \text{Cat}(w_{d,n}^* | \theta_x) & (\text{if } F_{d,n} = x) \\ \text{Cat}(w_{d,n}^* | \theta_{l=z_{F_{d,n}^*, d, \Xi_d^*}^*}) & (\text{otherwise}), \end{cases} \quad (7)$$

である。ここで、 ζ は NPYLM 部分の単語 n-gram オーダー、 $S_{d,n}$ は d 番目の観測に対応する文字列における n 番目の単語境界の位置、 $\delta(\cdot)$ はクロネッカーのデルタ、 $\text{Cat}(\cdot)$ はカテゴリカル分布を示して

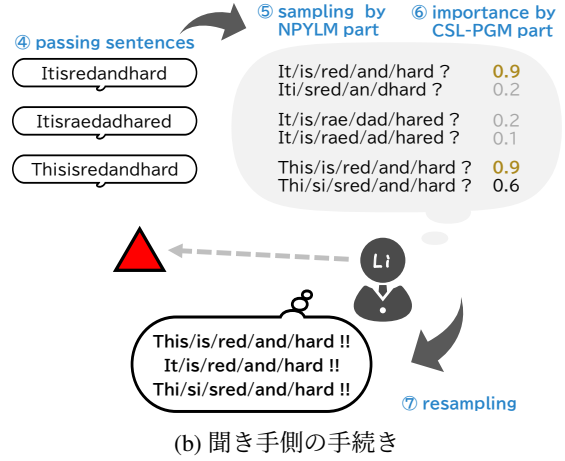
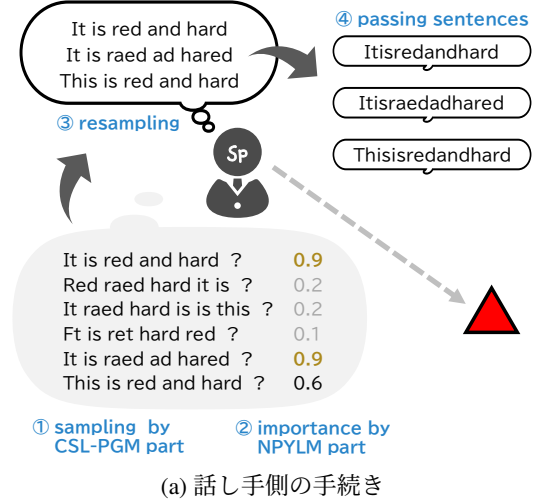


図 4: SIR 名付けゲームの流れ。話し手が発話する文を決定する際、および聞き手が分節化結果を評価する際に、SIR 法 [8,9] による近似が行われている。

いる。また、 x は、対応する単語が機能語もしくは分節化誤りであることを示す属性である。

A.3 実験におけるモデルのパラメータ

実験において、CSL-PGM 部分のハイパーパラメータについては全手法で $\alpha^* = 1.0$, $\gamma = 0.1$, $m_0^* = \mathbf{0}_{\text{dim}^*}$, $\kappa_0^* = 0.001$, $V_0^* = 0.01I_{\text{dim}^*}$, $v_0^* = \text{dim}^* + 2$ とする。ただし、 dim^* は属性 * の感覚情報の次元数であり、 $\mathbf{0}_{\text{dim}^*}$ は dim^* 次元の零ベクトル、 I_{dim^*} は dim^* 次の単位行列である。属性列 F_d の定義域 λ については、集合 $\{a, p, o, c, x\}$ の元からなる長さ N_d の重複を許した順列全てとする。NPYLM 部分では、文字モデルを可変長 n -gram、単語モデルを trigram (3-gram) とする。また SIR 名付けゲームにおいては、毎回のイテレーションで文中の単語の数 N_d を 3~7 の一様乱数で設定している。