

小説テキストを用いた雑談対話コーパスの自動構築とその分析に関する研究

○岩本 和真¹ 安藤 一秋² (香川大学¹ 創発科学研究科,² 創造工学部)

1. 概要

- 雑談対話コーパスの構築コストを削減するため、小説テキストから会話グループを抽出し発話者を特定・付与する**雑談対話コーパス（小説対話コーパス）の自動構築手法**を提案
- 会話グループ抽出として台詞間の**発話応答関係モデル**を構築し、**F1値80%**の性能を確認
- 発話者特定として、ルールベース、大規模言語モデル（LLM）、口調の特徴と**複数の手がかりを用いた特定手法**を提案し**Precision78%, F0.5値70%**の性能を確認
- 提案手法を用いて小説対話コーパスを構築し、その特徴を分析
 - コーパスの規模分析** ▶ **会話文脈の長い会話や複数人による会話を構築可能**
 - 口調に関する分析** ▶ コーパスには**多様な口調や人間関係を考慮した会話**が存在
 - 対話モデルを用いた分析** ▶ **親密性が高く、口調特性を反映した会話モデル**の構築が可能

「太郎くん！次郎くん！おはよう！」
明るいうちに振り返ると、そこには花子が立っていた。
太郎と次郎の幼なじみで、いつも3人で一緒にいる仲良しグループだった。
「花子、おはよう。今日から高校生だね」
「うん！楽しみだよ！新しい制服、似合ってる？」
花子はぐるりと一回りして見せた。
「似合ってるよ」
「似合ってるって...よかった花子」
次郎はにやりと笑った。
「もう、次郎ったら、からかわないでよ！」
花子は頬を膨らませた。
教室に入ると、窓際の席に一人の女の子が座っていた。
「あの子、誰だろう？」

小説テキスト

花子：太郎くん！次郎くん！おはよう！
太郎：花子、おはよう。今日から高校生だね
花子：うん！楽しみだよ！新しい制服、似合ってる？
太郎：似合ってるよ
次郎：似合ってるって...よかった花子
花子：もう、次郎ったら、からかわないでよ！

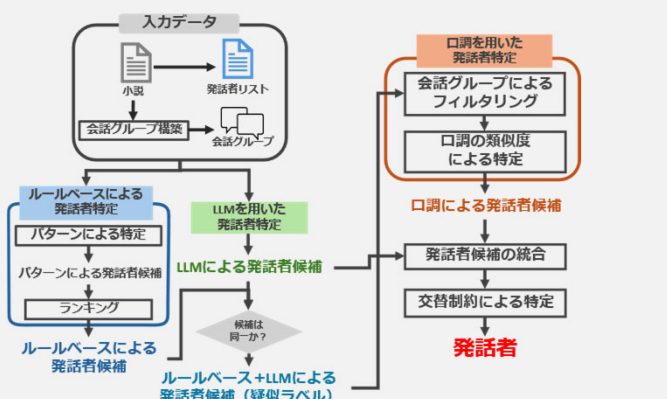
雑談対話コーパス

2. 会話グループ抽出

- 小説内の連続した台詞を1会話として抽出した場合、1会話に含まれる平均発話数が2.96と**少ない**
- 文脈を考慮するためBERTを用いた発話応答関係判定モデル**を構築し、その判定結果を会話グループの抽出に利用
- 提案モデルはBERTの複数の隠れ層を用いたモデル構造を採用
- 性能評価の結果、ルールベースがF1値64%の性能に対して提案モデルはF1値**80%**の性能であることを確認

3. 発話者特定

- ルールベース、LLM、口調ベクトルと複数の手がかりを用いることで、**発話者を正確に特定**することが目的
- 評価結果から、Precision**78%**とF0.5値**70%**の性能で特定できることを確認



4. コーパスの規模分析

- 「小説家になろう」の恋愛ジャンル20タイトルで構築
- 7,500発話から1会話平均5発話含む会話が1,300件構築可能
- 各小説で**10発話以上含む会話**と**ターン数が多い会話**を構築
- 3人以上**で構成されている会話など**複数人による会話**が存在

	発話数	会話数	1会話内の発話数	特定発話数	特定会話数
合計	149,855	26,301	-	98,550	8,317
平均	7,492.7	1315.0	5.73	4927.5	415.9

7. 応用活用例と今後の展望

活用例1：小説を用いたキャラクターAIの開発

- 著者側：技術的要素、構築コストから実用化が困難
 - ▶ **自動構築システムにより著者自身がAIを構築可能**
- エンジニア側：構築コスト、著作権により商用化が困難
 - ▶ **基盤システムの提供により、実用化を促進**

活用例2：エンタメ業界における創作支援

- ゲーム作成におけるNPCの会話文生成
- キャラクターの人物像やシナリオに応じた会話文の創作支援

活用例3：教育における応用

- キャラクターAIを介したコミュニケーション練習



小説の著者

小説テキスト

雑談対話コーパスの自動構築システム

システム提供

技術者



雑談対話コーパス



キャラクターAIの構築システム

キャラクターAI

今後の展望

- 小説から周辺情報や人間関係、感情情報を抽出して、小説対話コーパスを拡張する手法の検討
- より柔軟にユーザと雑談できる対話モデルの検討
- 構築手法の性能向上による小説対話コーパスの質向上