

くずし字翻刻のための古文用音声認識システムの検討

研究代表者

文学研究科

情報理工学研究所

戸塚史織

ZHANG Yutao

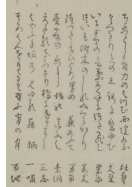
1. 研究背景

くずし字

- 古典籍、古文書などの前近代資料に使われてきた文字
- 日本には数百万の古典籍・古文書があり、近年多くの機関がインターネットを通じてこれを公開し始めているが、くずし字で記されているため現代の多くの人は読めない

くずし字翻刻

- くずし字を現在で使用されている文字に変換する



ちよつとした角方は、はつむり顔際
くつきり色のしろき髪はひ
いきまの官裏にくきにつきまし
つとも附木のしめくからり
備つきもとり算の業あふて
ほほの光に耀の温ふ也
えた顔でどつきり横に連りわたし
そはんをきてきてきて書み書み月

百穂 一穂 三穂 五穂 七穂 九穂 十一穂 十三穂 十五穂 十七穂 十九穂 二十一穂 二十三穂 二十五穂 二十七穂 二十九穂 三十一穂 三十三穂 三十五穂 三十七穂 三十九穂 四十一穂 四十三穂 四十五穂 四十七穂 四十九穂 五十一穂 五十三穂 五十五穂 五十七穂 五十九穂 六十一穂 六十三穂 六十五穂 六十七穂 六十九穂 七十一穂 七十三穂 七十五穂 七十七穂 七十九穂 八十一穂 八十三穂 八十五穂 八十七穂 八十九穂 九十一穂 九十三穂 九十五穂 九十七穂 九十九穂 百穂

- 負荷が大きい
- くずし字解読では近年AIやOCR(光学的文字認識)技術が注目されている

現状

- くずし字解読を得意とする人々の多くはAIやOCR技術を活用したくずし字解読技術を使用するよりも、使い慣れたアナログ的な方法の方がより効率的に翻刻できる
- 翻刻に慣れた人が過去に自分の手元に手書きで残した翻刻のテキストデータ化が求められる

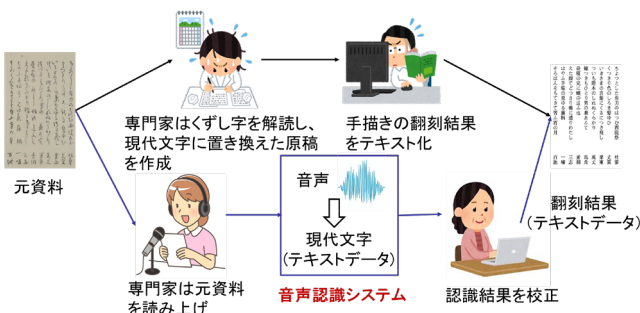
翻刻作業の負荷を減らし、特に翻刻のテキストデータ化を効率化させることが現在の重要課題

2. 古文用音声認識システム

くずし字翻刻の音声入力システム

- 負荷を減らし、翻刻のテキストデータ化効率化させるため
 - 画像だけではなく、音声からくずし字翻刻をサポートする方法を検討
- 翻刻に音声入力システムの導入を検討

高精度な古文認識を実現できる音声認識システムが求められる



3. 古文用音声認識システムの構築

学習データの収集

- 古文テキスト学習データの収集
 - 既存データセットに含まれている作品
 - 江戸時代洒落本(15作品) 文字数:約17万
 - 独自の学習データ
 - 役者評判記の形態素解析データ
 - ◆ 「役者白虎通 江戸の巻」(先頭5000語)(昨年度)
 - ◆ 「役者白虎通」のより詳細な正誤判定を行った解析データ(京都の巻先頭2000語・江戸の巻先頭2000語の計4000語解析修正データ)(今年度)
- 古文音声データセットの収録(今年度)



気伝導マイクと皮膚密着型マイク(NAM)を使用して同時に収録

- 話者: 古文を流暢に読める8人の話者(女性6名、男性2名)
- 読み上げ作品: 形態素情報が付与された洒落本作品の翻刻テキスト
- 収録場所: ARCスタジオの録音ブース
- 収録音声の長さ: 各話者が40分間読み上げ、合計で約5時間のラベル付き音声データを収集

古文用音声認識システムの構築

- DNN-HMMベースの音声認識システム(昨年度)
- 商用音声認識APIを試す(今年度)
 - 古文の認識精度を向上させるため、カスタム言語モデル機能を利用(ドメイン固有の音声の文字起こしの精度を向上させるように設計される)
 - 古文学習データを用いて言語モデルをチューニング
- End-to-End 古文用音声認識モデルを構築(今年度)
 - ラベルあり古文音声データが限られているため、自己教師あり学習モデルwav2vec 2.0を活用

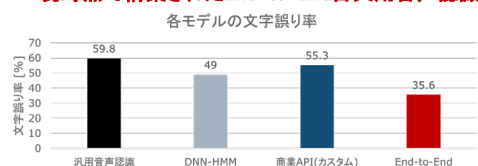
4. 評価実験

実験目的

- 現時点で構築された古文音声認識システムの性能評価

実験手順

- テストデータ
 - 洒落本の読み上げデータ 内容:「聖遊廓」の一部
- 認識テストに使用する音声認識システム
 - 汎用日本語音声認識システム
 - DNN-HMMベースの音声認識システム
 - 商用音声認識API(カスタム言語モデル使用)
 - 現時点で構築されたEnd-to-End古文用音声認識モデル



- 昨年度より、古文音声認識システムの認識性能を向上させることができた