

人間の認知特性の基づくキャラクタにあった音声生成の検討

斉藤彰吾^{†1} 大井翔^{†2} 佐野睦夫^{†2}

概要: 現在, 先行研究ではアニメキャラクタに合った声の生成方法を研究し, アニメキャラクタの画像特徴から既存の声を推定していたが, 良い結果は得られなかった. そこで本研究では, キャラクタに合った声の特徴と画像の特徴を関連付けるために, アニメキャラクタにとって違和感のない声の傾向を分析することを考えた. 手法として, 1つのキャラクタイラストに対して複数の音声を提示し, 参加者によって評価を行う. これによって人間の知覚がアニメキャラクタの声の特徴を知覚する傾向を分析する. その後, 得られた傾向データを元に学習データの割合を調整し生成する音声をより自然な音声で生成できないか検討する.

1. はじめに

近年, 電子書籍の市場規模は順調に伸びてきており, 音声読み上げ機能があるものもある事や, 電子書籍にて何か別の物事を行いながら読む事が可能な時代となりつつある[1]. しかし, 音声読み上げのクオリティは低く, 違和感を覚えてしまうことや, 声優が音声を吹き込むオーディオコミックは制作にコストがかかる.

本研究では, 得られた結果を元に音声学習アルゴリズムの改良を考えた. 手法として音声学習を行う際, 学習データを人間の知覚としてキャラクタにあった声の特徴を感じる傾向分析によって得られた傾向データを元に学習データ割合の調整を行う事を提案する. 図1では傾向分析を行った際に仮に音声Aを30%, 音声Bを20%, 音声Cを50%程度の割合が未知キャラクタの声を生成する際に必要な割合だと評価された場合に訓練する方法について示す.

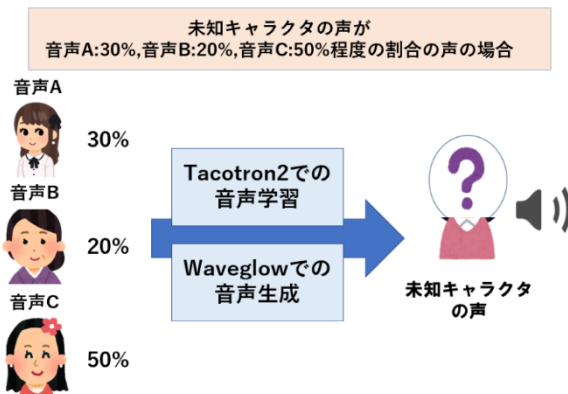


図1 提案アルゴリズムの概要図

2. 関連研究

コミックのキャラクタの特徴から音声を推定する研究[2]やキャラクタの顔画像のパーツに注目し人がキャラクタのどのパーツを見て声を想像しているかのメカニズムを

用いる事で, よりシンプルなキャラクタの画像特徴から音声推定が可能かを検討した研究が行われている[3,4]. この研究では, 人の顔と声には関係性があるという研究を参考に[5,6], 初めに人がキャラクタの顔画像のイラストを見た際に, キャラクタの顔のどのパーツを見て声を想像しているのか調査を行っている. その結果, 人は「目の形」「髪の色」「髪の色」が有力な特徴である事が判明した. この研究では図2に示すようにその三点の中でも特に有力な特徴である「目の形」を抽出している. 画像推定から音声推定を行いキャラクタの画像から人の感性が想像する声に近い音声が生成可能か調査を行った[3,4].

また, 他の研究ではキャラクタの声優のキャスティングに対して調査している研究もあり, キャラクタの声優が担当したいくつかの音声データを収集し, 得られた音声データを元にして音響特徴量を算出する. その後得られた音響特徴量と印象値の関係を学習させ, 活用する事で新しいキャラクタに対しても印象値与える事で図3に示すように適切な音響特徴量を推定するといった学習から音声生成を行った研究もあるが, 音声生成は実際の人の音声で生成されているため, いわゆるアニメ声といった声の生成が難しいと考えられる. 本研究ではアニメ等のキャラクタを演じられている方の音声を用いて音声生成したものをを用いていくためよりキャラクタイラストから生成される音声としての精度が高いものの作成を狙う[7].



図2 目の抽出方法

^{†1} 大阪工業大学大学院

^{†2} 大阪工業大学

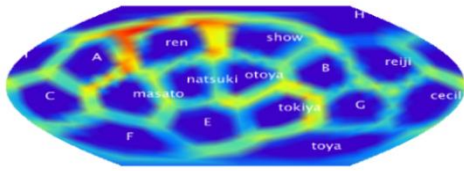


図3 音響特徴量と印象値の学習結果

3. 提案手法

本研究の目的はキャラクタの画像に対して音声推定や生成する研究をより精度が高いものにする事を目指す。そのため、本稿ではキャラクタイラストと声の関係性を調査する事を目的とする。従来研究[2]ではキャラクタの顔画像の中で特に目の形が声を想像する際に有力だと判明しており、目の形の特徴量を用いて音声推定、音声学習が行われている。しかし、それらはキャラクタイラストのみを提示しその画像を見てどういった印象を持つかを調査し得られたもので、実際の音声を用いて調査が行われたものではない。そのため、本研究ではキャラクタイラストと音声生成を行った音声を一対比較法により重複しないランダムな組み合わせを提示し、評価を頂くことでより精度の高い音声学習の特徴モデルの生成を試みる。

本研究でキャラクタイラストと声の関係性を調査する流れは以下の図4のとおりである。



図4 提案手法の流れ

(1) 自動生成された音声データとキャラクタイラストを一対比較法を用いて提示する。また、この際一度提示した組み合わせは再度表示されないようにする。

(2) 1~5の数値でキャラクタと音声とがどれだけ合致していると感じるのか実験参加者に評点を行ってもらう。

(3) 得られた評点データから画像と音の特徴モデルの生成を狙う。

この時、(2)の評点は低い点数である程キャラクタと音声の組み合わせが一致していないといった評価になり、高い点数程キャラクタと音声の組み合わせが一致していると感じているとする。続いて、実験に用いるデータに関しての説明を行う、用いたキャラクタイラストは以下の図のものを用いた。また、キャラクタイラストに区別を行うために今回は左上のキャラクタから右に向かって male1, male2, male3, male4 とし、左下のキャラクタから右に向かって female1, female2, female3, female4. と名称付け

る。



図5 使用したキャラクタイラスト
cre8tiveAI (<https://cre8tiveai.com/sc>)

続いて、音声は8種類用いて実験を行い先行研究[2]のシステムを用いて音声生成された音声を用いた。音声の内容は発話内容によるアンケート結果の差異を生み出さないため「こんにちは」の音声のみを用いて行った。

更に本研究では得られた傾向データから人間の知覚として違和感の無い音声を生成する為音声学習を行う際の割合を推定し学習を行う。

4. 実験

本研究の実験は、一対比較法を用いたキャラクタイラストと音声の提示を行い、それらを実験参加者にどれ程一致しているのかを直感的に評価してもらう事で、人間の知覚としてキャラクタイラストにあった声の特徴を感じる傾向分析が可能だと考える。また、得られた傾向から今後の画像と音の特徴モデルの生成を行っていく。

本実験で用いた音声は、声優の「田村ゆかり」「松風雅也」「沢城みゆき」「鈴木健一」「花澤香菜」「逢坂良太」「平田広明」「小野友樹」の音声を Tacoron2 を用いて学習させ Waveglow にて音声生成を行ったものを用いる。実験ではキャラクタイラスト8種類と音声8種類の組み合わせを重複無しでランダムに提示したため64種類の組み合わせで提示を行った。

本実験を行う際に提示用プログラムを制作した。実験では「押下するとランダム画像表示」のボタンを押下し提示されたキャラクタイラストと音声の組み合わせを評価する。その後、1~5のラジオボタンを選択し選択後に押下する「評価値選んでね」のボタンを押すことで評価値が決定される仕様とした。また本研究では、評点値1を選んだ時キャラクタ画像と音声とが一致していると全然感じない、2の時は感じない、3の時はどちらでもない、4の時は感じる、5の時はとても感じるとした。

以下の図に実験で用いた提示する際に用いた実験用プログラム画面の表示を行う(図6)。



図6 実験用プログラムの画面

本実験は、21~24歳の男子大学生・大学院生7名を対象に行った、実験では評点に加え、「どういった基準でその評価を選定したか、教えてください」の質問をアンケート終了後に行った。

続いて、実験結果を以下の図に提示する。実験結果はそれぞれのキャラクターイラストごとに得られた1~5の評価値の割合を示したものとなる。

実験結果の考察を行うと **male1** の画像に対して松風雅也さんと小野友樹さんの音声を用いて音声生成を行った音声が顔画像と一致していると見受けられる。逆に田村ゆかりさんの音声を用いて音声生成を行った音声とはあまり一致していないと見受けられる。続いて、**male2** の画像に対しては逢坂良太さんの音声を用いて音声生成を行った音声が顔画像と特に一致していると見受けられる。逆に田村ゆかりさんの音声を用いて音声生成を行った音声とは一致していないと見受けられる。続いて、**male3** の画像に対して平田広明さんと小野友樹さんの音声を用いて音声生成を行った音声が顔画像と一致していると見受けられる。逆に田村ゆかりさんの音声を用いて音声生成を行った音声とは一致していないと見受けられる。

続いて、**male4** の画像に対して逢坂良太さんの音声を用いて音声生成を行った音声が顔画像と特に一致していると見受けられる。逆に田村ゆかりさんの音声を用いて音声生成を行った音声とは特に一致していないと見受けられる。続いて、**female1** の画像に対して田村ゆかりさんと沢城みゆきさんの音声を用いて音声生成を行った音声が顔画像と一致していると見受けられる。逆に逢坂良太さんと小野友樹さんの音声を用いて音声生成を行った音声とは特に一致していないと見受けられる。続いて、**female2** の画像に対して田村ゆかりさんと花澤香菜さんの音声を用いて音声生成を行った音声が顔画像と一致していると見受けられる。逆に平田広明さんの音声を用いて音声生成を行った音声とは特に一致していないと見受けられる。続いて、**female3** の画像に対して田村ゆかりさんの音声を用いて音声生成を行った音声が顔画像と特に一致していると見受けられる。逆に逢坂良太さんと平田広明さんの音声を用いて音声生成を行った音声とは特に一致していないと見受けられる。続いて、**female4** の画像に対して花澤香菜さんと沢城みゆきさんの音声を用いて音声生成を行った音声が顔画像と一致していると見受けられる。逆に小野友樹さんの音声を用いて音声生成を行った音声とは一致していないと見受けられる。

これらの結果から考察を行うと、男性キャラクタを用いた傾向調査では小野友樹さんと逢坂良太さんの音声を用いて音声生成を行った音声は人間が感じる傾向として適しているといった傾向が見受けられ、田村ゆかりさんの音声を元にして音声生成を行った音声は男性キャラクタの音声には適していない傾向が見受けられた。続いて女性キャラクタを用いた傾向調査では田村ゆかりさん、沢城みゆきさん、花澤香菜さんの3名の音声を元として生成を行った音声が良い結果が出たが、その中でも田村ゆかりさんの音声を元とした音声が傾向として良い結果が現れた。また、女性キャラクタの画像と小野友樹さんの音声は適していないといった傾向が見受けられた。

また、本研究では上記の傾向の中から数種類の傾向を元

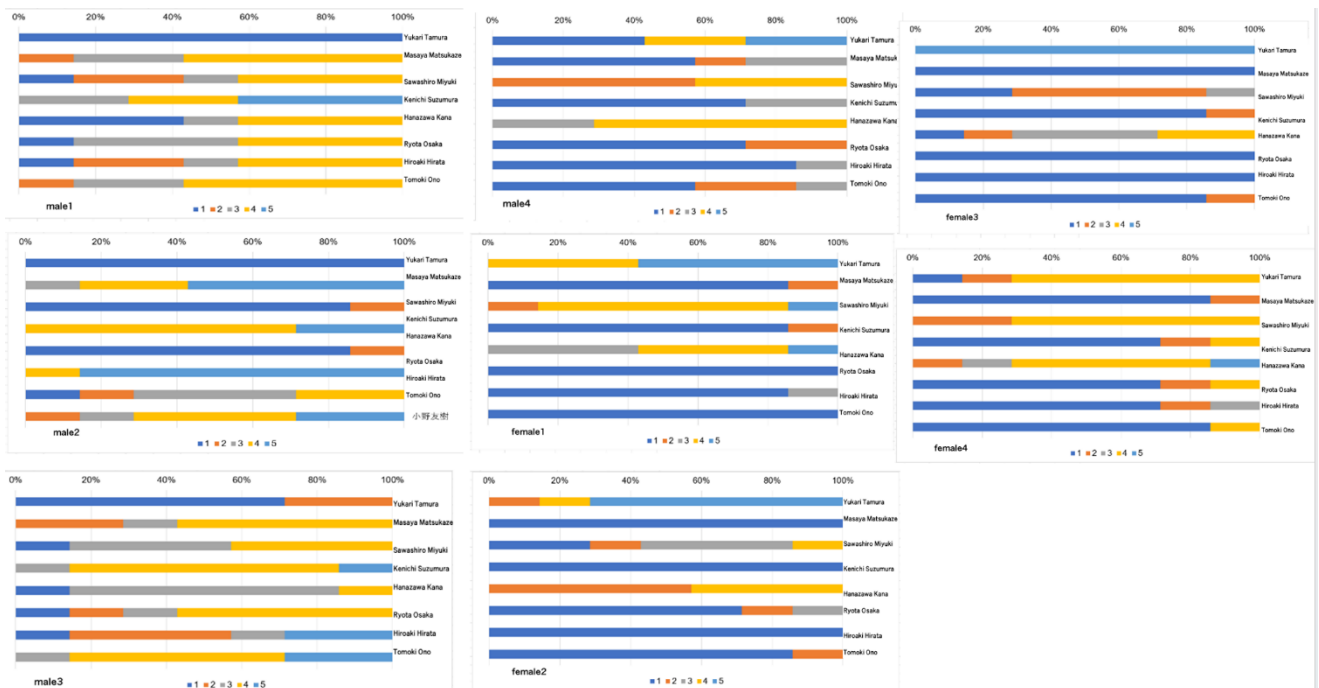


図7 傾向分析結果

に学習データの調整を行う。今回は図 8 に図示する male2 の結果を用いて学習を行う。



図 8 male2 の結果

学習を行う際 male2 の結果から音声の割合配分を制定する必要がある為評点 1~5 をそれぞれ 1~5 点として扱いその点数を元に学習データ割合配分を行った。上記の male2 の場合田村ゆかり:7 点, 松風雅也:29 点, 沢城みゆき:8 点, 鈴木健一:30 点, 花澤香菜:8 点, 逢坂良太:34 点, 平田広明:20 点, 小野友樹:27 点の結果が出ている為, 各点数/総合点数 $\times 100$ にて割合計算を行うと田村ゆかりの音声データを 4%, 松風雅也の音声データを 18%, 沢城みゆきの音声データを 5%, 鈴木健一の音声データを 18%, 花澤香菜の音声データを 5%, 逢坂良太の音声データを 21%, 平田広明の音声データを 12%, 小野友樹の音声データを 17% となった為, その配分で音声学習を行い音声の出力を行う。

5. おわりに

本研究では, 複数枚のキャラクターのイラストに対して複数の音声と組み合わせたものを提示した, これによって人間の知覚としてキャラクターにあった声の特徴を感じる傾向分析を行った。本実験では 7 名の実験参加者に実験を行い 8 種類のキャラクターと音声の組み合わせによって実験を行い, 人間の知覚としての特徴データが得られた。しか

し, 実験を行うにあたって傾向分析するデータ量が少なく十分な調査結果が得られなかったという点が考えられる。また, 実験終了時に評価基準の選定方法についての調査を行った結果キャラクターの雰囲気やイメージ, 性別で声を選定した声が多いという結果が得られた。また, 得られた傾向データから音声学習の割合を調整し音声学習を行う事を検討した。今後は調査に使用する音声の種類を増やし音声の種類を細分化する事で特徴モデルをより深め, より人間の知覚として違和感無く感じるキャラクターの音声生成を行う事を目指す。

参考文献

- [1] 一般社団法人 電子出版政策・流通協議会 "平成 30 年度 電子書籍等の情報アクセシビリティの現状等に関する調査研究・報告"(参照 2022-06-18).
- [2] Wang, Yujia, et al., "Comic-guided speech synthesis," ACM Transactions on Graphics (TOG) 38. 6 (2019): 1-14.
- [3] 大道昇, 大井翔, 佐野睦夫, "オーディオブック自動生成のための 2 次元キャラクター特徴と声の関係性の調査," 情報処理学会関西支部, 支部大会, 2021.
- [4] 大道昇, 大井翔, 佐野睦夫, "オーディオブック自動生成のための 2 次元キャラクター特徴と声の関係性の調査," 情報処理学会 インタラクシオン 2021, (2021).
- [5] Smith, Harriet MJ, et al., "Concordant cues in faces and voices: Testing the backup signal hypothesis," Evolutionary Psychology 14. 1 (2016): 1474704916630317.
- [6] 大杉 康仁, 齋藤 大輔, 峯松 信明, "Eigenvoice と CLNF を用いた顔から声への統計的対応付けの検討," 研究報告音声言語情報処理 (SLP) 2017. 3 (2017): 1-6.
- [7] 酒井えりか, 伊藤彰教, 伊藤貴之, "ゲームキャラクターと声質の傾向分析," (可視化, キャラクターアニメーション, 映像表現・芸術科学フォーラム 2016). "映像情報メディア学会技術報告 40. 11. 一般社団法人 映像情報メディア学会, (2016).