

日本大学 文理学部
マルチメディア情報処理 (2024/07/23)

音声合成・歌声合成が拓く未来

高道 慎之介 (慶應義塾大学)

自己紹介



@forthshinji

名前

高道 慎之介 (たかみち しんのすけ)

現職

東京大学 講師 → 慶應義塾大学 准教授

経歴

離島 → 熊本高専 → 長岡技大 → 奈良先端大
2016年に博士(工学).

兼務

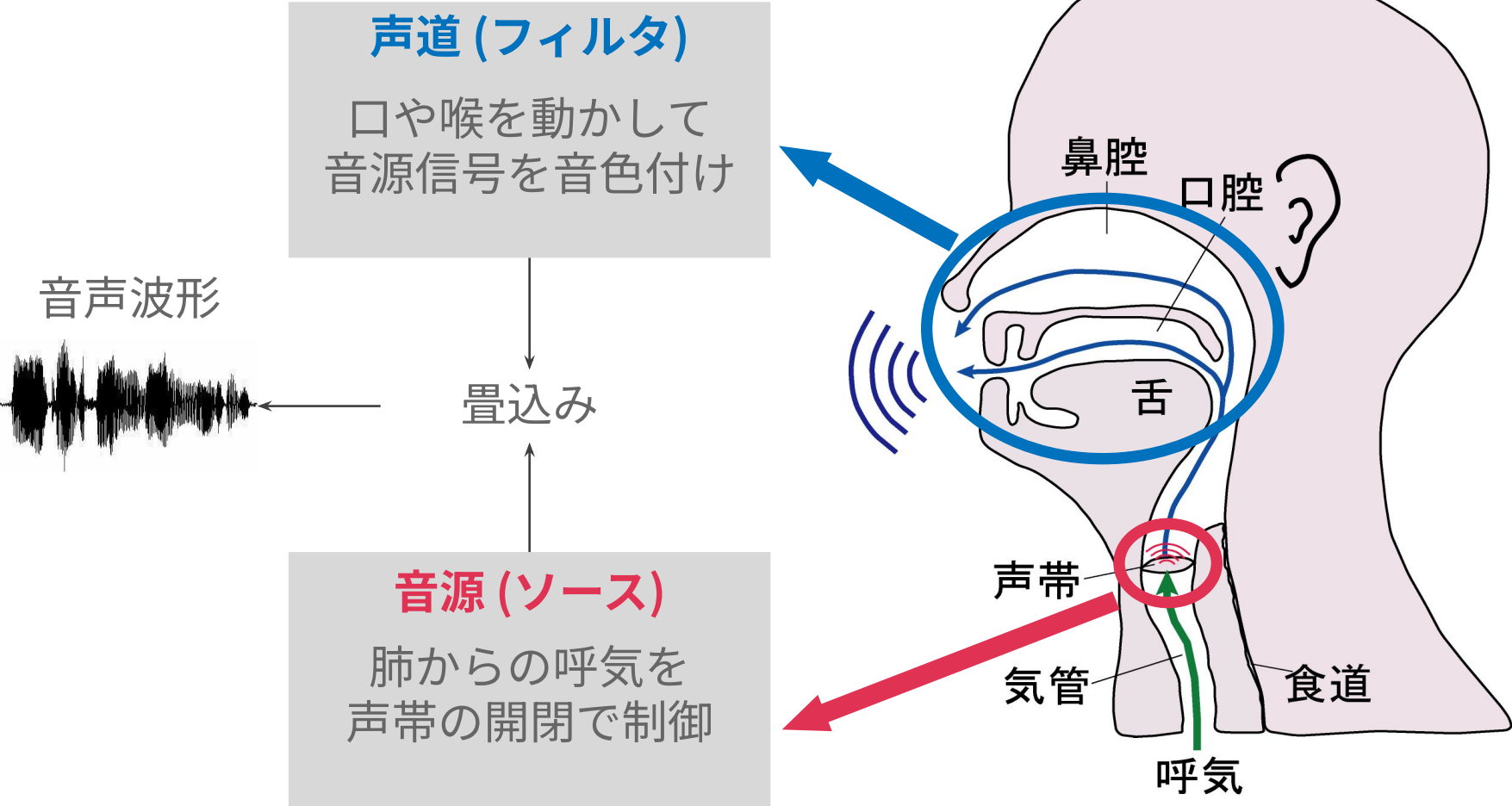
東京大学 特任准教授
理研 客員研究員
他

本講演のテーマ

計算機による声合成はどこまで来たのか？
これからどう進んでいくのか？

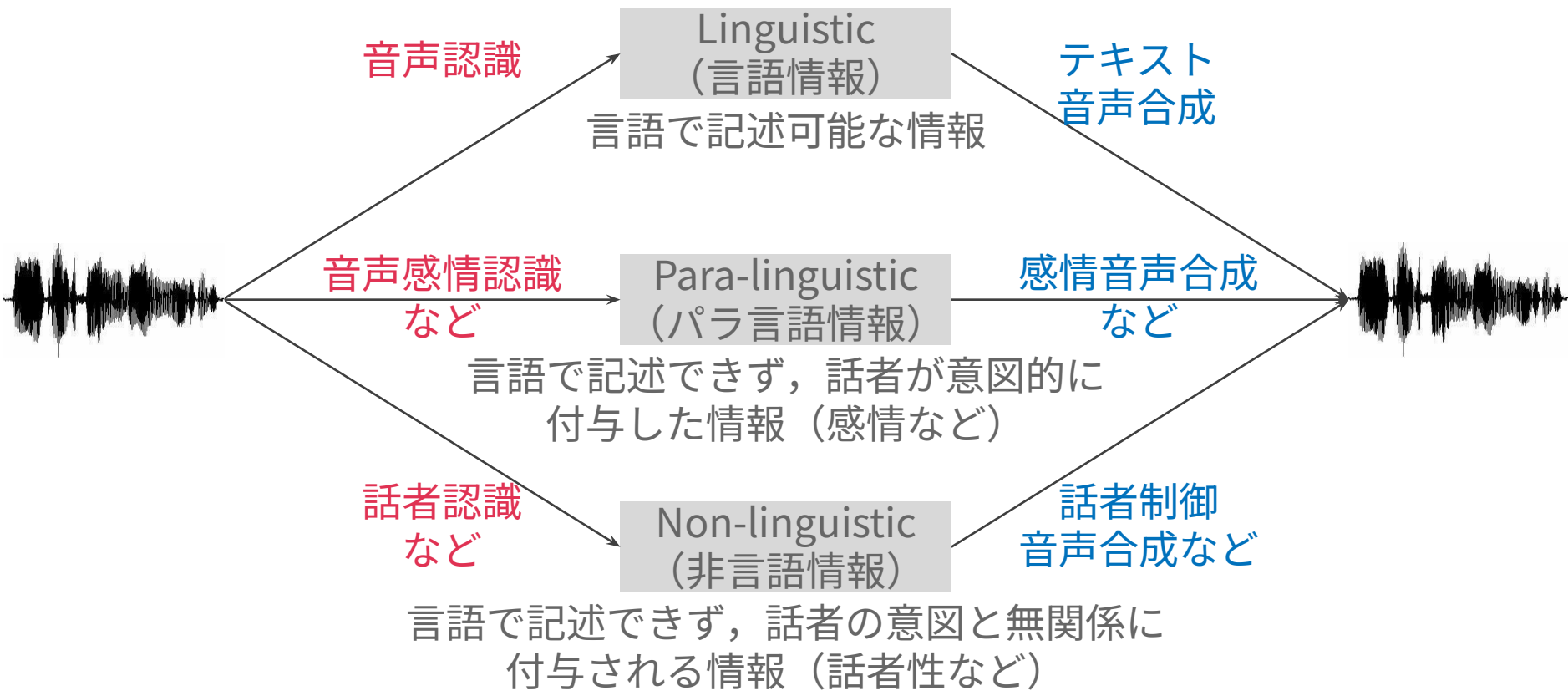
音声情報処理の基礎知識

発声器官のモデル: ソース・フィルタモデル



音声のもつ情報とそれを扱う主な音声技術

物理 → 情報 → 物理

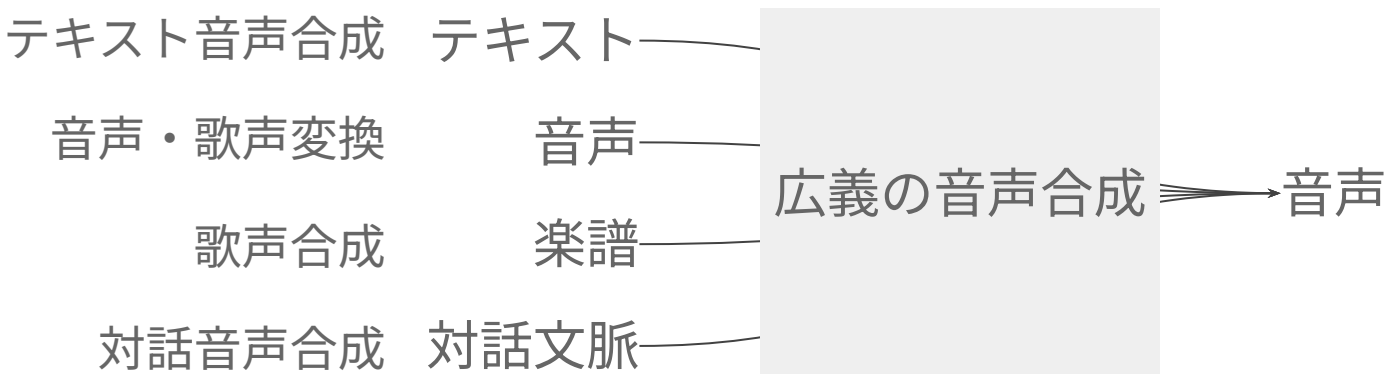


音声合成とは

- 狭義には、テキストから音声を作成する技術
 - = テキスト音声合成



- 広義には、何らかの情報から音声を作成する技術
 - {音声, 動画, 話者情報, 概念} to speech



紹介する研究紹介 20連発！(2023年以降を中心に)

研究事例紹介①：人間の音声表現を拡張する

研究事例紹介②：音声の文化を守る

研究事例紹介③：コンピュータが心をもったときに

研究事例紹介④：歌の可能性を広げる

研究事例紹介⑤：音声に関する感性を定量化する

研究事例紹介⑥：環境音を人工的に作り出す

研究事例紹介①

人間の音声表現を拡張する技術

ハンターハンター ボイスチェンジャー

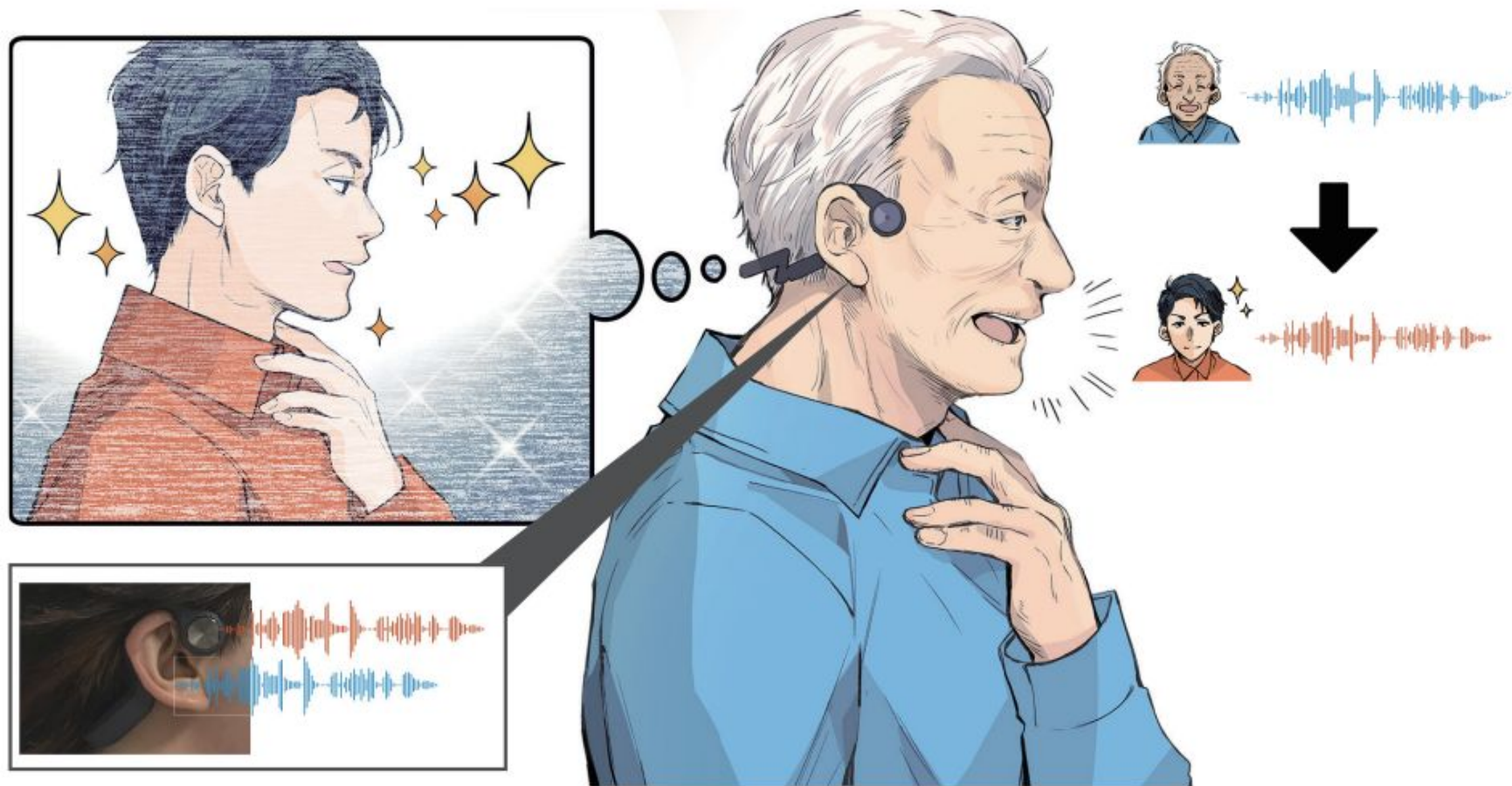


“自分の声が『HUNTER×HUNTER』のキャラクターの声に変換されてスピーカーから流れ、なりきりアフレコ体験ができます”

スマホで動くリアルタイムボイスチェンジャー



Digital speech makeup : 音声の「鏡」を作る



自分で聞く自分の声が変わると、自己の認識が変わる．認識が変わると行動が変わる

エモーションキャンセリング： 音声への「盾」を作る



話し手の音声と聞き手の知覚感情を明らかにし、聞き手の心的負担を軽減する

Spatial voice conversion: 空間上の狙った声だけ変える

Two persons are speaking to me (woman from left, man from right.)



**Spatial voice
conversion**

Converts desired speaker's voice (left in demo) without changing other info.

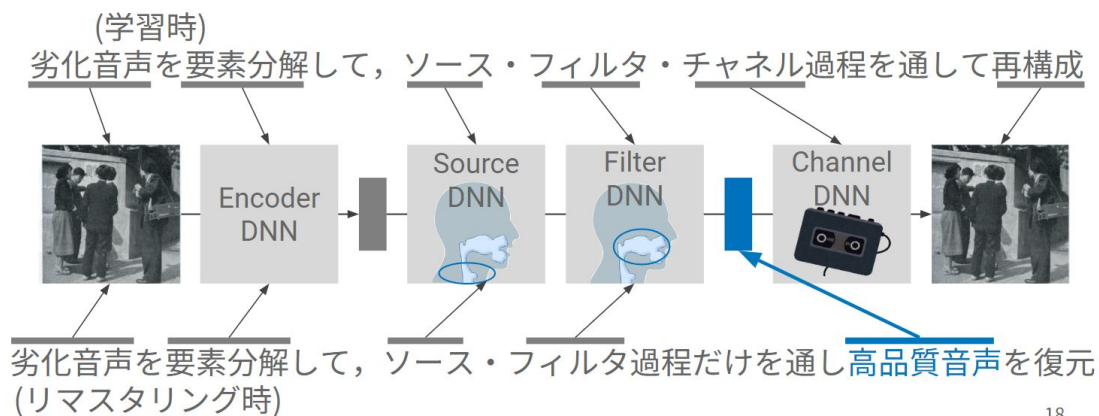
研究事例紹介②

音声の文化を守る

Neural speech restoration : 音声文化を復元する



デジタルメディアの前の録音機器は劣化する。廃棄される



18

音声録音過程を模擬した自己教師あり音声復元。歴史音声だけから現代音質への復元を学習できる



1960~70年代にオープンリールに録音された東北方言昔話をデジタル化(200時間)



現代品質に近づける

音声文化を復元する→新たな動的コンテンツを作る

雪女郎

(第1章 途中まで)

音声文化を受動鑑賞する時代を終え，生成 AI と組み合わせた能動鑑賞の時代へ

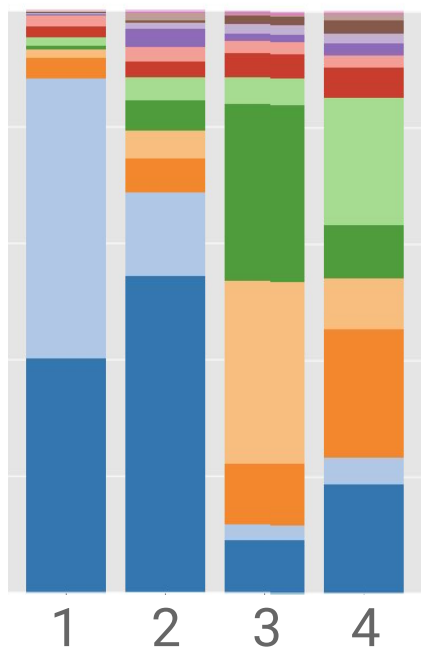
研究事例紹介③

コンピュータが心をもったときに

自然に間違える音声AI エージェント： 「あー」「えー」のフィラーとその個人性



フィラー語の使用頻度



ざっくりいうと、(アノ)先ほど(アノ)少し(アノ)お話ししましたけども、戦後のそういうサブカルチャーのイメージという...

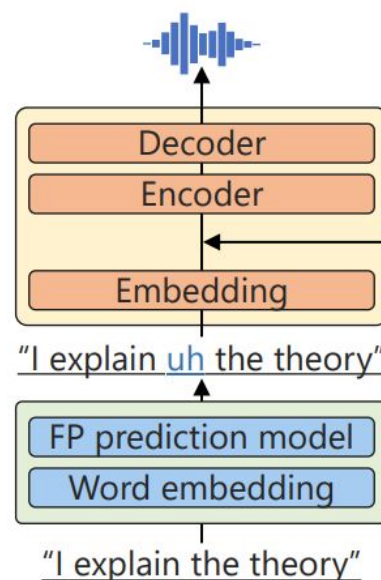
Synthetic speech including FP

Seq2seq model trained using spontaneous speech

FP-included text

FP prediction model trained using FP-annotated texts

Fluent text

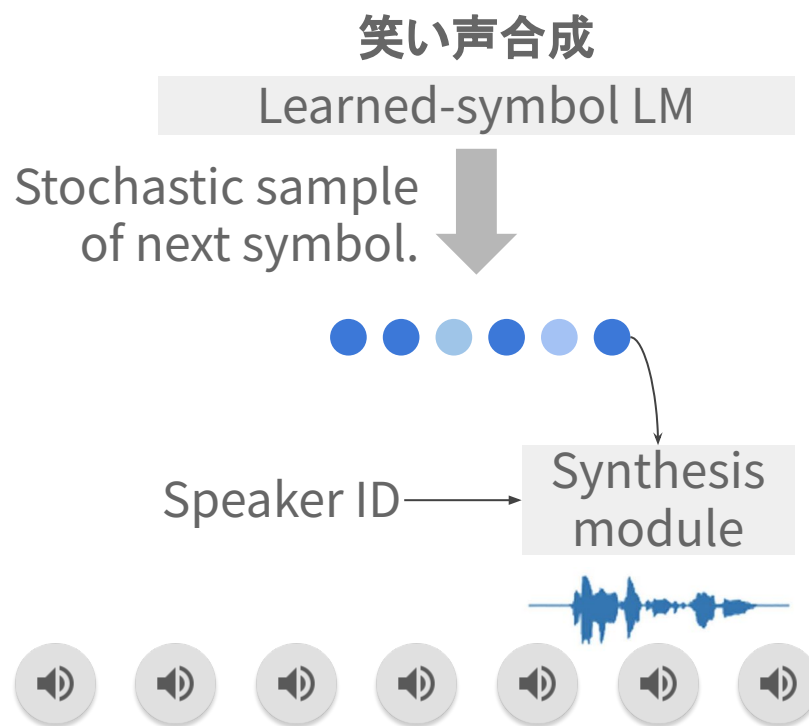
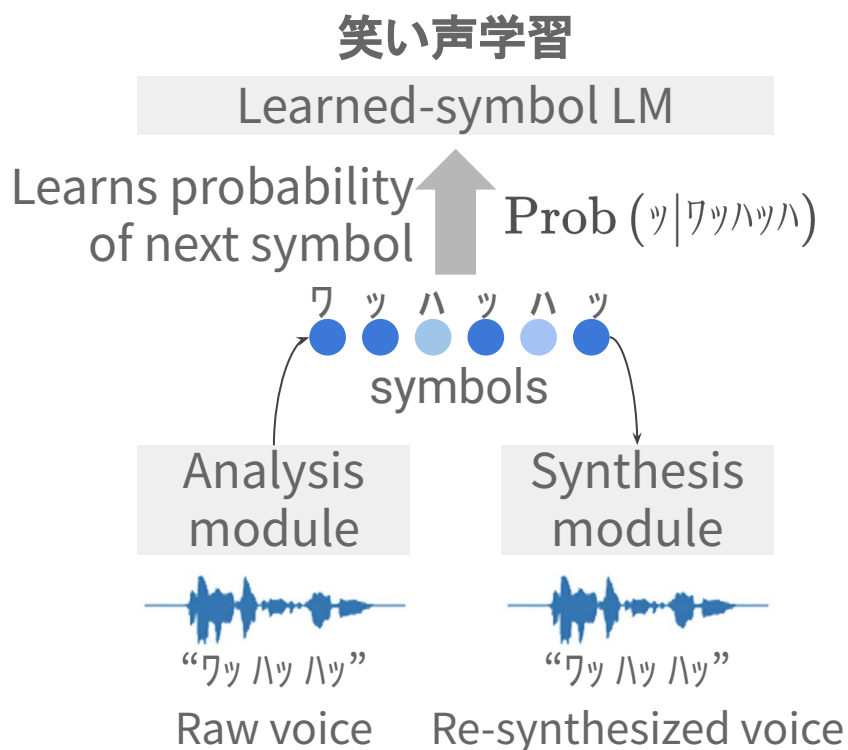


ざっくりいうと、先ほど少しお話ししましたけども、戦後のそういうサブカルチャーのイメージという...

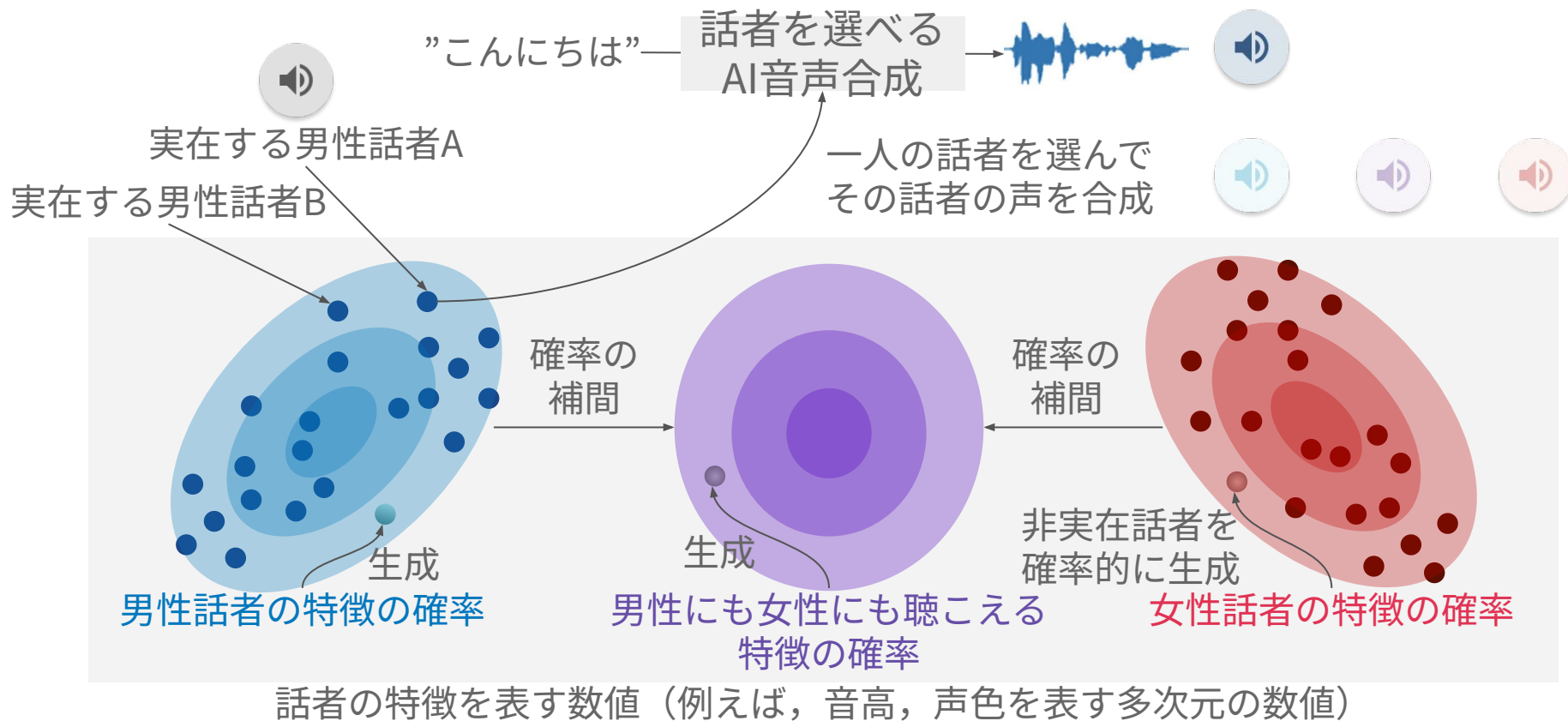
Neural laughter generation : 笑い声を合成する

- コンピュータも笑ってほしい

- 笑い声はテキストで書けない。だから、テキストの代わりに「自動で見つけた音の塊」を使う。それを確率モデルで予測。

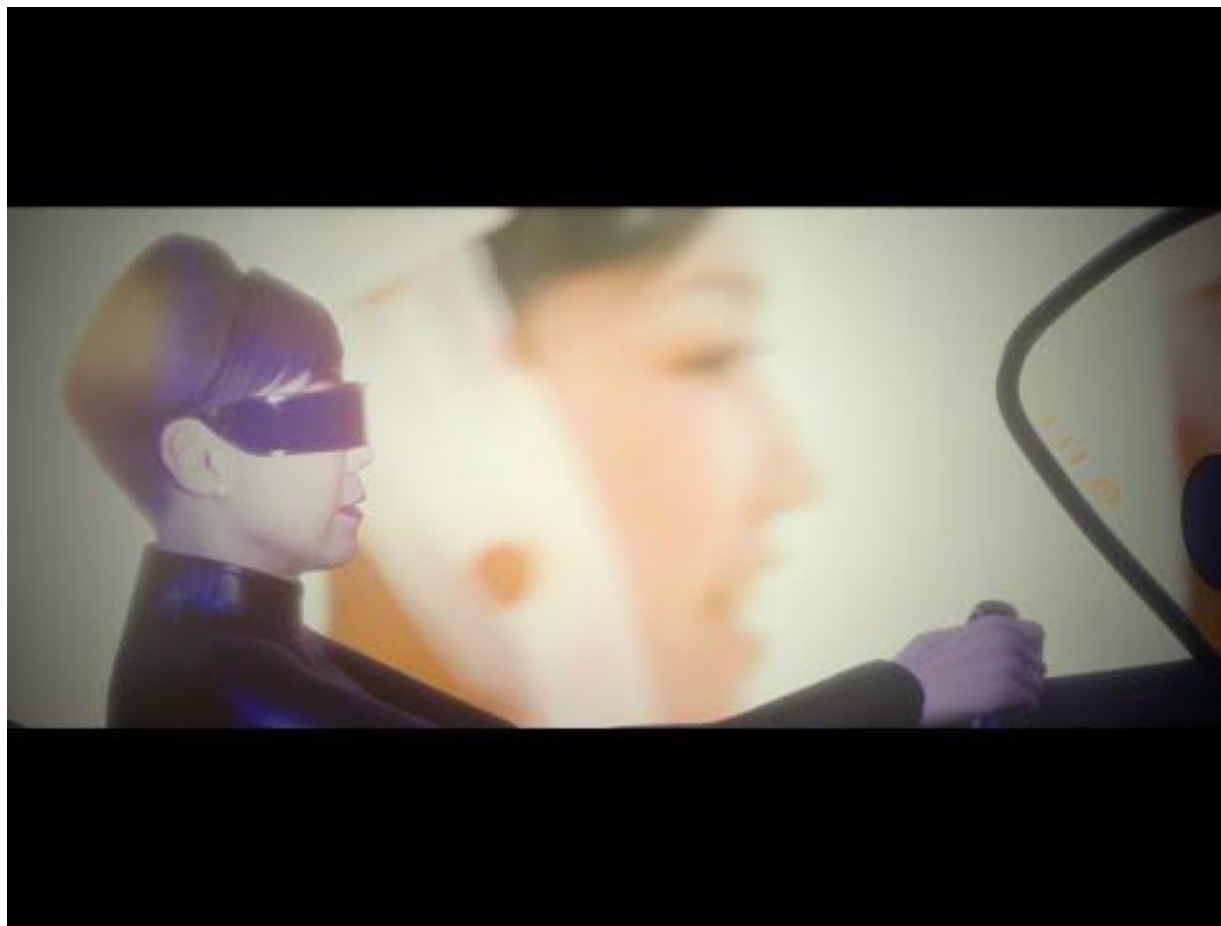


Mid-attribute speaker generation : バイナリの属性を超えた非実在話者



研究事例紹介④ 歌の可能性を広げる

時を超えて蘇る50年前の歌声 ～スモールデータを用いたタスク混合深層学習による歌唱再現～



“**歌手の松任谷由実氏が50年前にデビューした当時の歌声を人工再現する技術**を開発しました。(中略) 当時の声色と歌唱表現を忠実に再現することに成功しました” 23

歌声合成



The advertisement features a central illustration of the character Megpoid, a young girl with vibrant green hair styled in pigtails, wearing a blue visor and an orange jacket over a yellow top. To her left are two product boxes for the software, each displaying the character and the name 'megpoid'. A pink circular badge with the word 'NEW' is positioned between the boxes and the character. The background is a light blue gradient with the word 'INTERNET' in a stylized font at the top right.

INTERNET

Synthesizer V AI

アーティスト名

megpoid

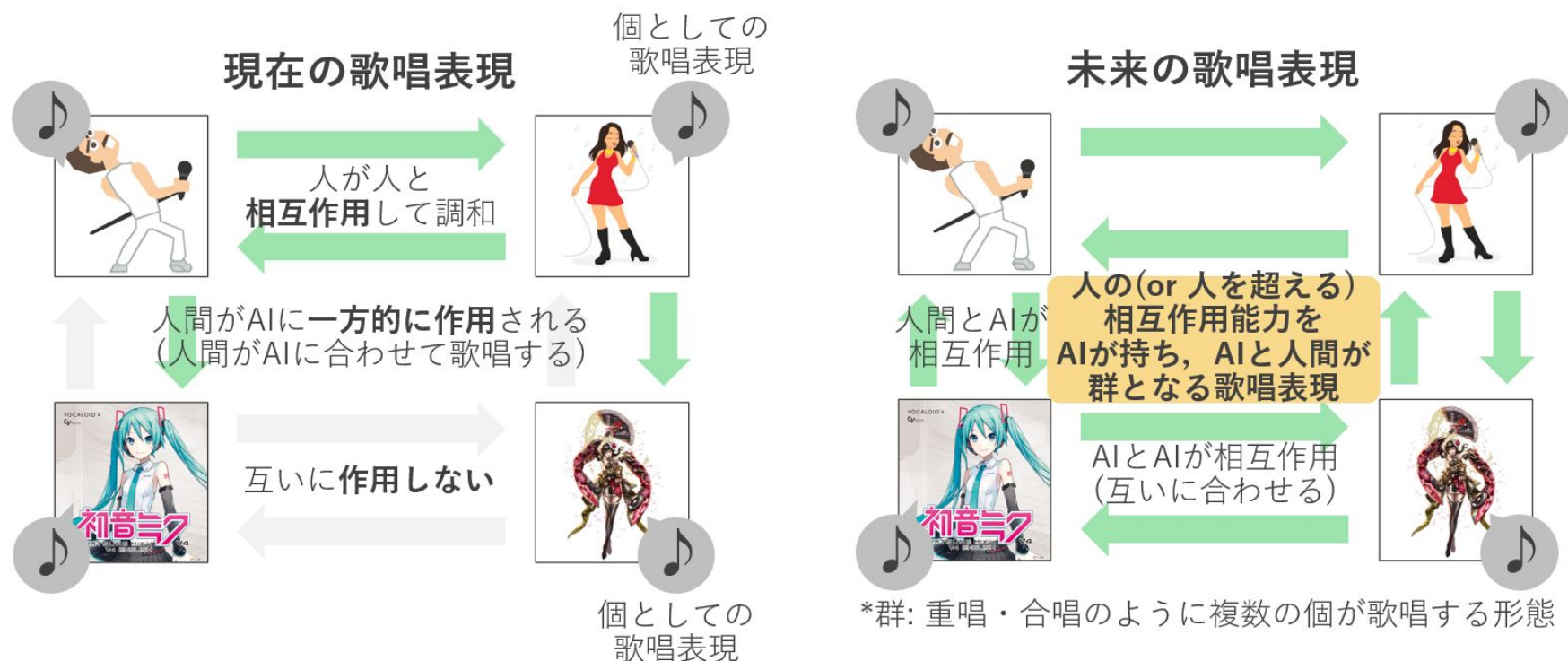
メグリポイド

NEW

2023.12.20登場

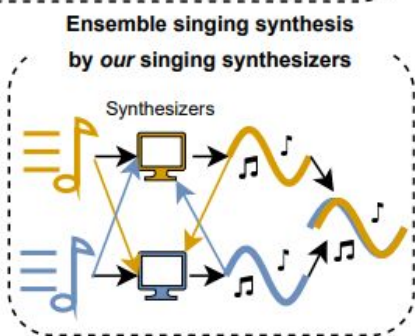
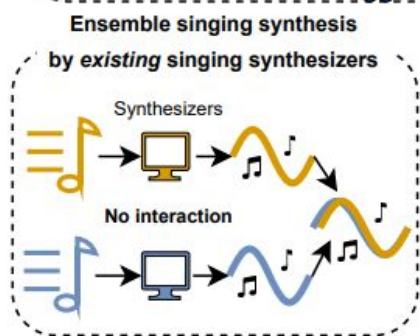
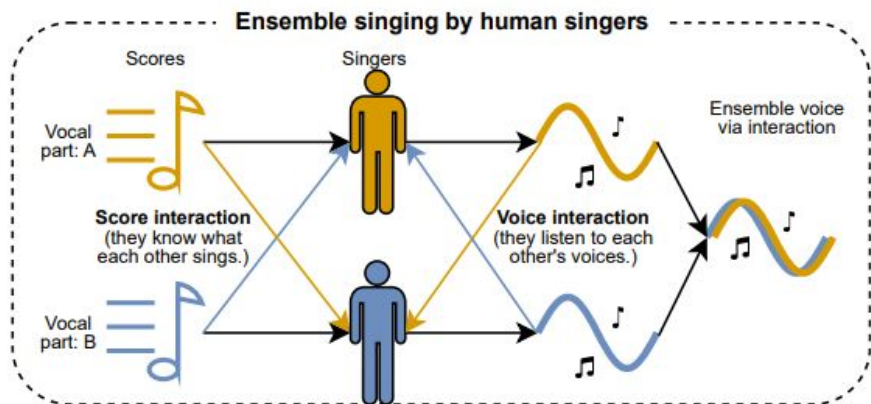
jaCappella : 他者と相互作用するAIシンガー & アカペラDB

- 人間歌手とAI歌手が相互作用する社会へ
 - 歌声信号レベルでの相互作用の例がアカペラ
- 相互作用を分析・再現する技術
 - 相互作用により生まれる一体感とは？
 - 計算機は一体感を生み出せるか？

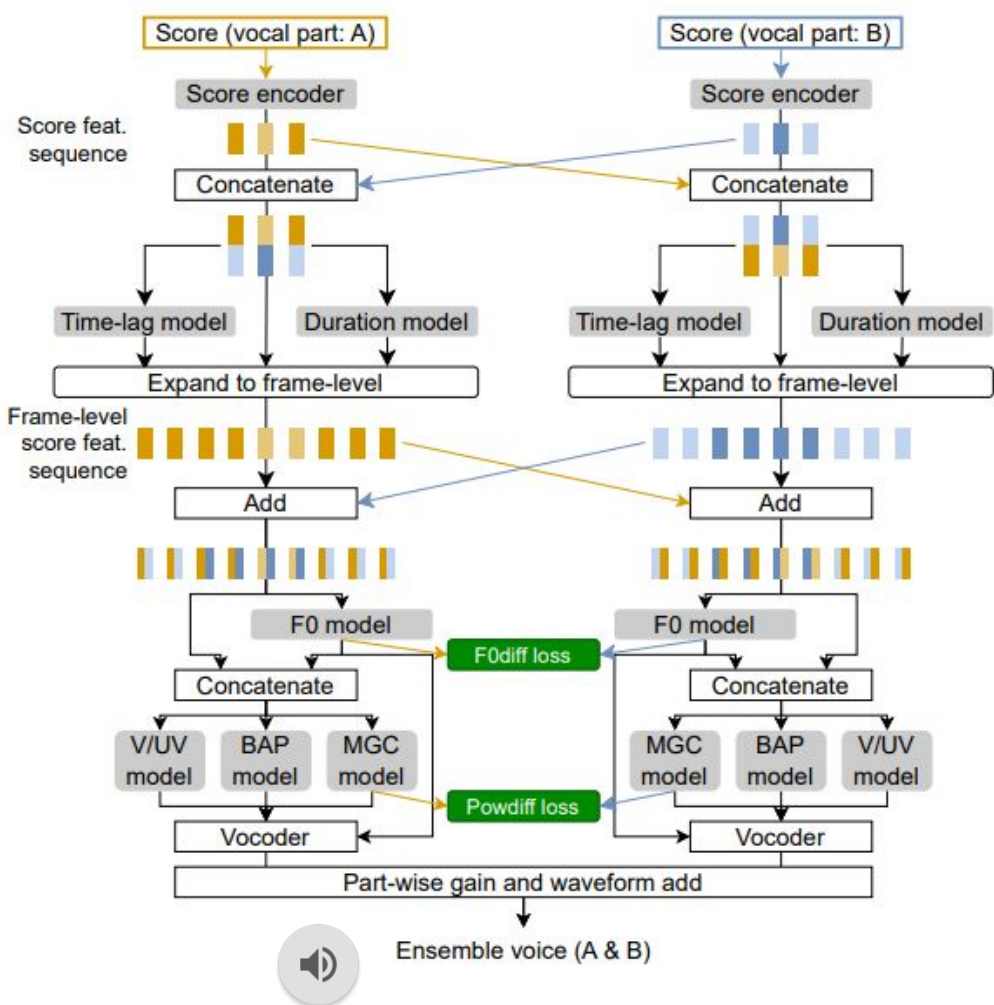


歌唱者間相互作用を模擬するアカペラ合成

人間歌唱者同士は、楽譜レベル・音声レベルの相互作用を起こす

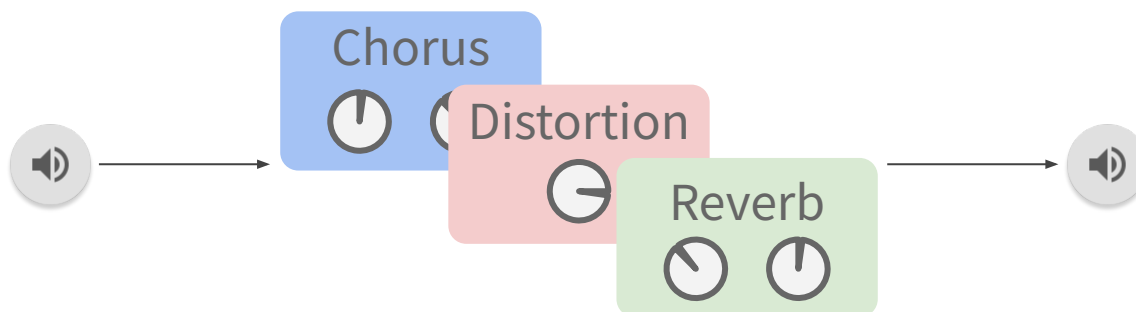


AI歌唱者同士で、楽譜レベル・音声レベルの相互作用を起こす

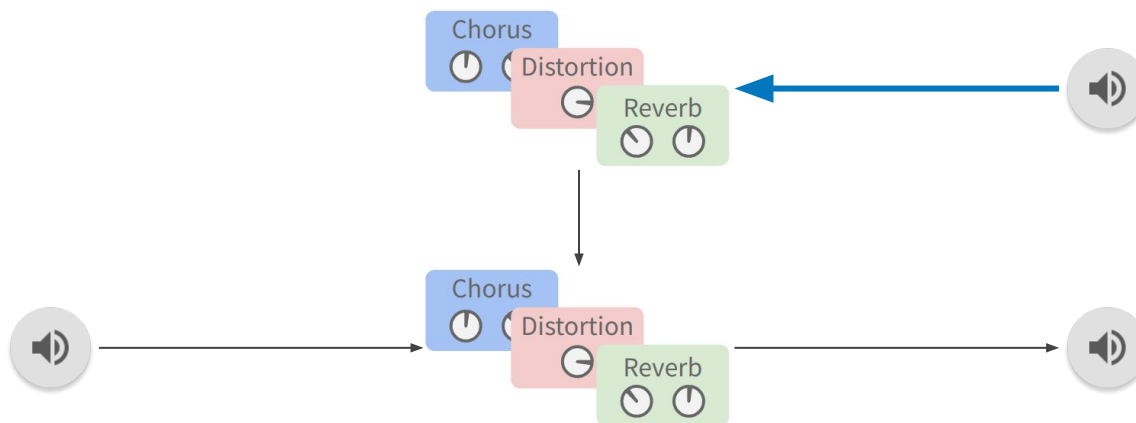


Audio effect chain 推定

- Audio effect をかけて音楽制作・ライブの芸術性を高める



- Audio effect chain 推定は、**処理後の信号 (wet signal) から effect を推定**する。その effect を他の音にも適用できる。



研究事例紹介⑤

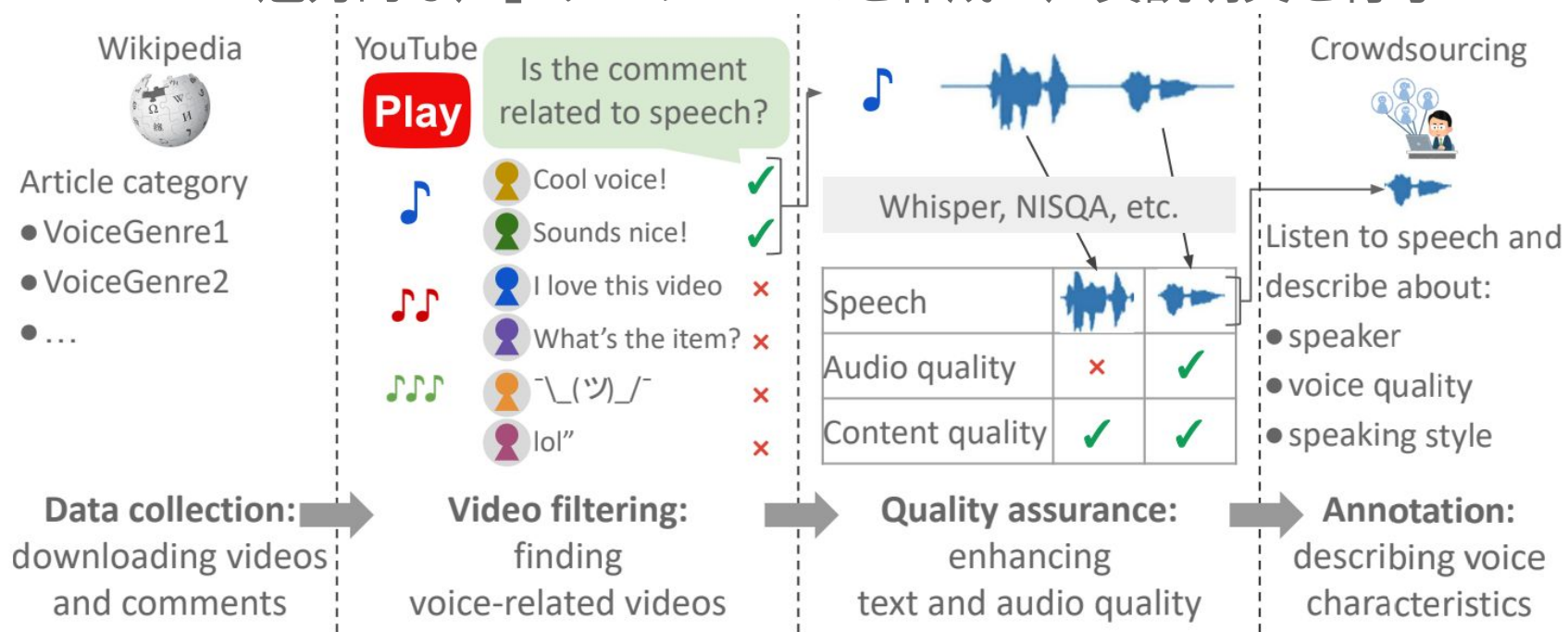
音声に関する感性を定量化する

Coco-Nut:

「魅力的な声」データベース

• 動画サービスにおける「魅力的な声」とは？

- いろんなユーザが声質について言及している声！
- → 「魅力的な声」データベースを作成&声質説明文を付与

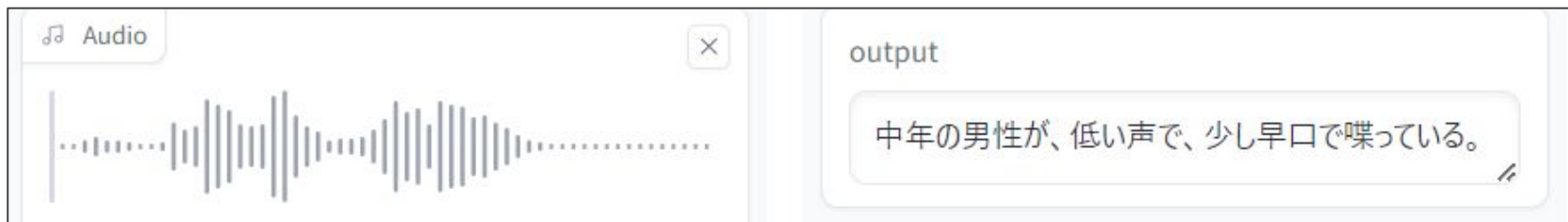


	30代くらいの男性の声。ゆっくりと穏やかな話し方でした。苦悩に満ちた、けだるそうな声でした。
	明るい中年の女性のはきはきとした声で楽しそうに喋っている。

Cococap-beta : 声質キャプションニング技術

- **Cococap-beta**

- <https://huggingface.co/spaces/sarulab-speech/CoCoCap-beta>

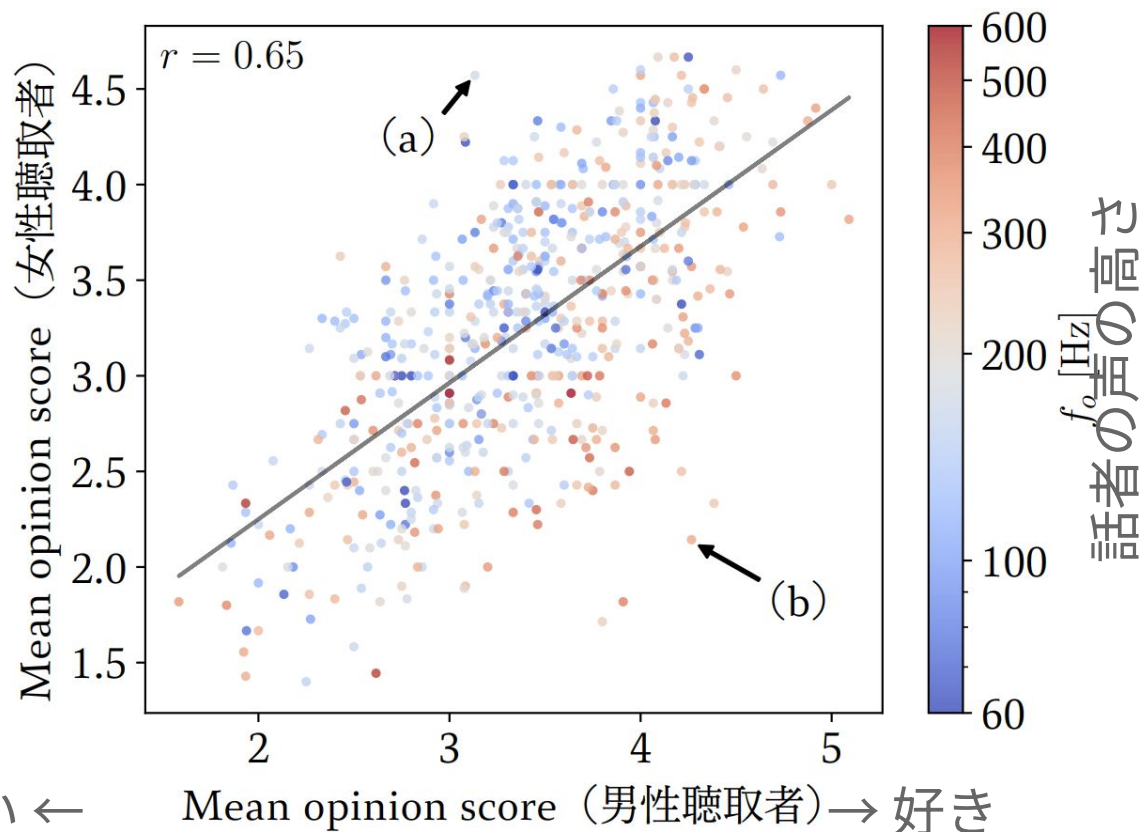


- **これ以外にも、声質表現文と音声の相互変換が可能になる！**
 - 声質表現文 → 音声 (自然言語による音声検索)
 - 声質表現文 → TTSシステム (自然言語による音声AIデザイン)
 - それ以外にも：感情キャプションニング

CocoNut-Humoresque : 「魅力的な声」 データベースと嗜好分析

各個人にとって「魅力的な声」ってなんだろう？

1800話者の声に対する900聴取者の「この声すきスコア」を集めてみた



魅力は男女間で相関する傾向。一方で片方の性別のみに好かれる場合も。

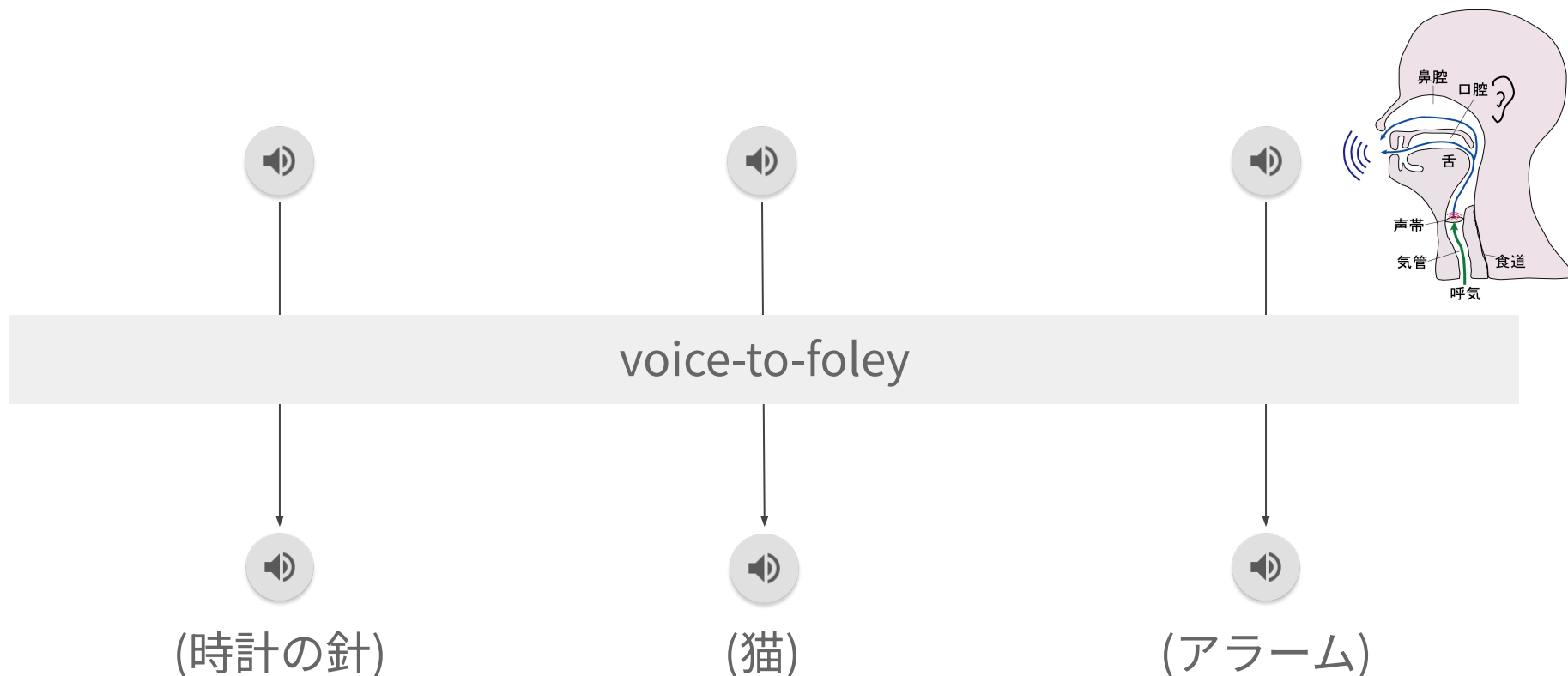
(a) 「若い男性が、はきはきした低い声で怒ったように喋っている」

(b) 「10代の少女が、かわいらしい声でまったりとした口調で喋っている」

研究事例紹介⑥ 環境音を人工的に作り出す

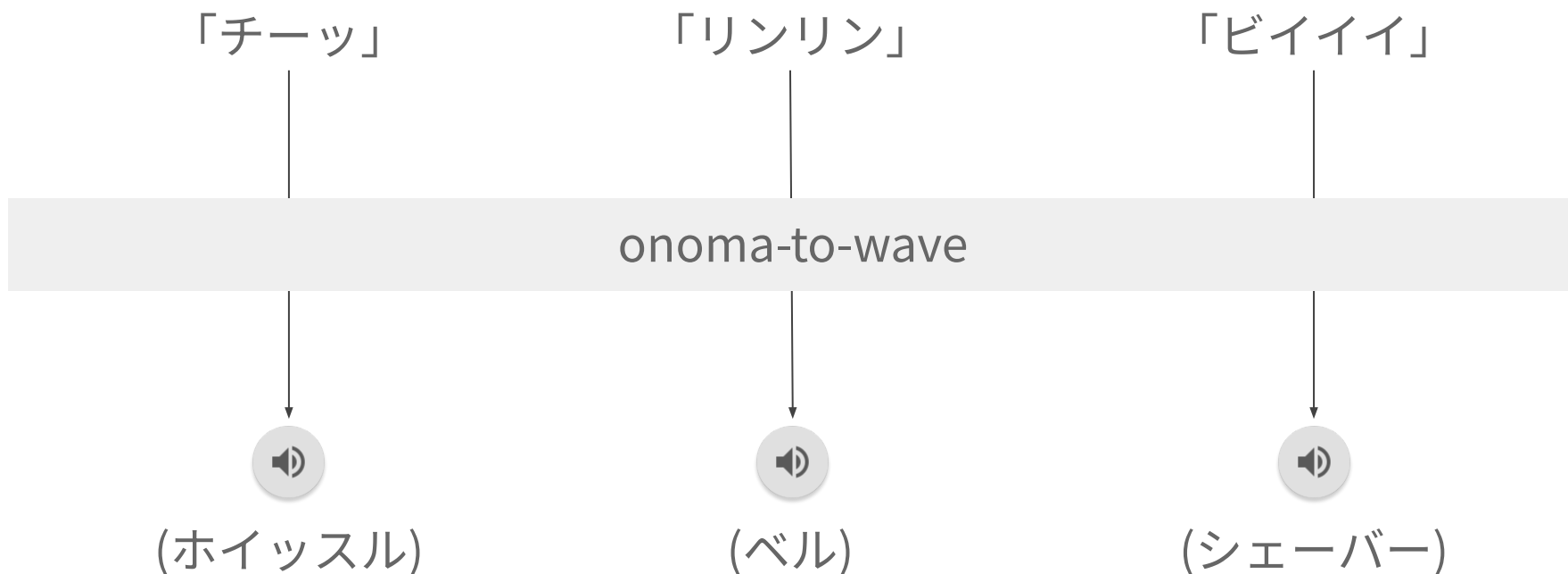
voice-to-foley : 音声模倣からの環境音合成

- 音声模倣：非音声を声まねしたもの
 - 直感的に環境音を合成できる方法
 - → 音声模倣から環境音を合成できるのでは？



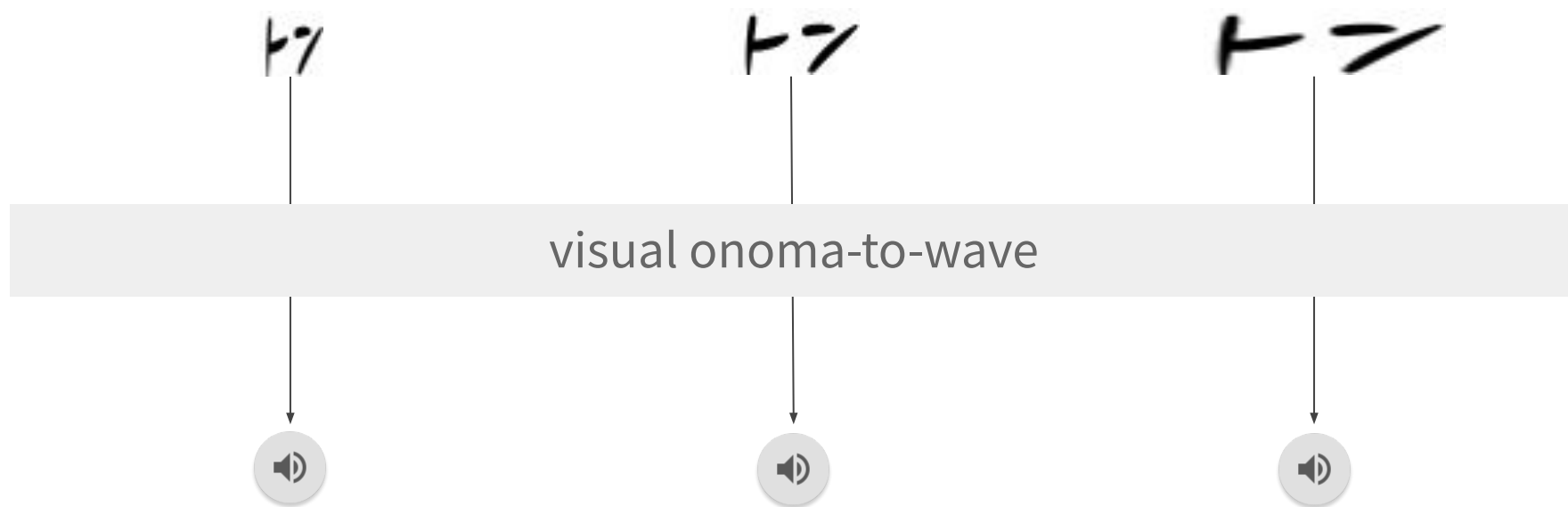
onoma-to-wave : オノマトペからの環境音合成

- **オノマトペ：音を文字で模倣したもの**
 - 環境音の時間構造を文字で表してくれる。
 - → オノマトペ文字から環境音を合成できるのでは？



visual onoma-to-wave : オノマトペ画像からの音声合成

- 画像オノマトペ：漫画などで用いられる音表現
 - オノマトペだけでなく，画像エフェクトでも情感を伝える
 - → オノマトペ画像から環境音を合成できるのでは？



まとめ

まとめ

研究事例紹介①：人間の音声表現を拡張する

研究事例紹介②：音声の文化を守る

研究事例紹介③：コンピュータが心をもったときに

研究事例紹介④：歌の可能性を広げる

研究事例紹介⑤：音声に関する感性を定量化する

研究事例紹介⑥：環境音を人工的に作り出す

音声合成・歌声合成 = 音を正確に再現する、だけではありません。
価値観をどう更新できるか、人間をどう幸せにできるかを考えてください¹³⁷