

スペシャルセッション 音声B/音声A/聴覚

[ここまで来た声質変換技術 – 実用可能性の視点からの現状認識と将来展望 –]

歌声インタフェース: 歌声を対象とした信号処理と それに基づくインタフェース構築

中野 倫靖, 後藤 真孝
(産業技術総合研究所)

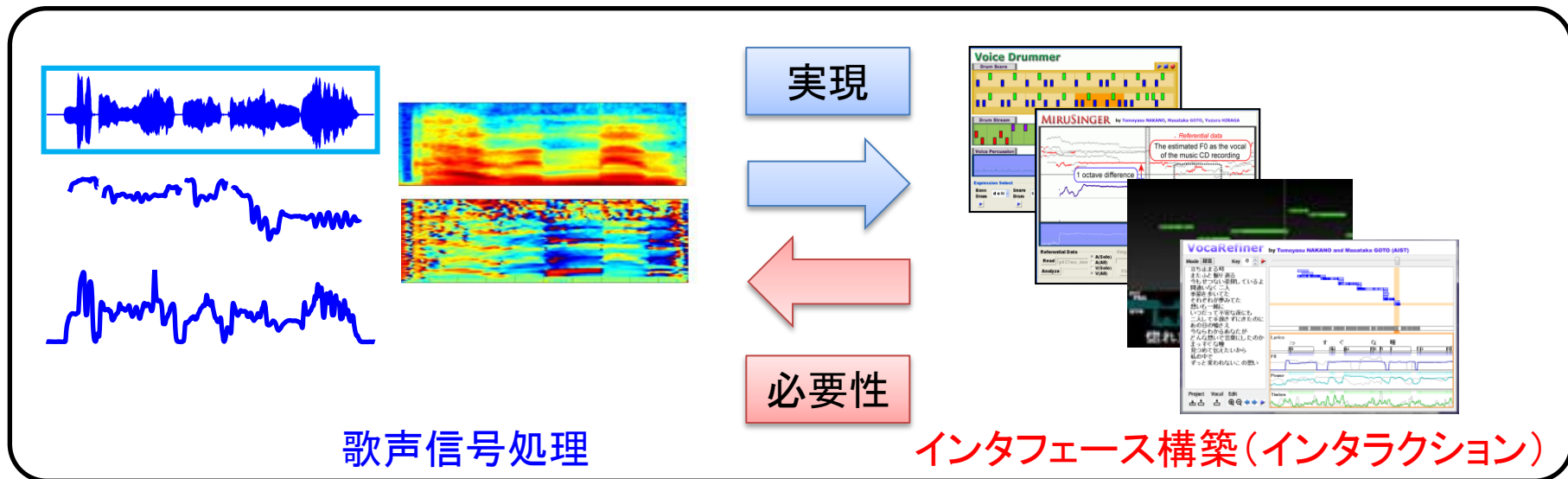
2013年9月27日

日本音響学会 2013年秋期研究発表会(講演番号3-7-3)

ここまで来た声質変換技術 — 実用可能性の視点からの現状認識と将来展望 —

歌声インタフェースとは

□ 歌声信号処理に基づく**インタフェース構築**や
インタラクションによって
 人々の音楽生活を
 より豊かにする研究アプローチ

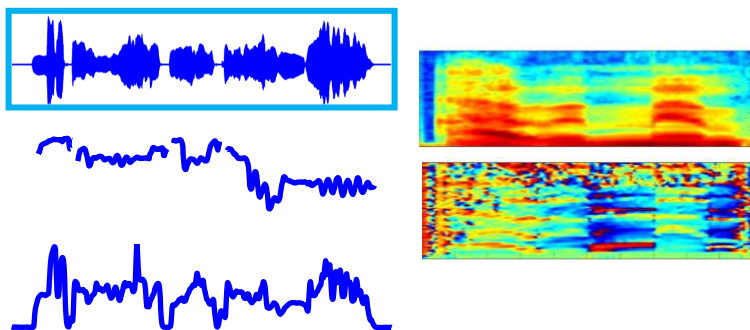


ここまで来た声質変換技術 — 実用可能性の視点からの現状認識と将来展望 —

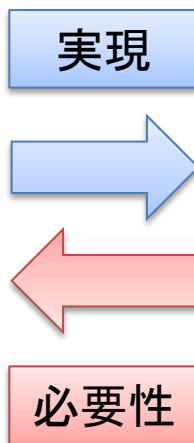
歌声インタフェースとは

□ 歌声信号処理に基づく**インタフェース構築**や
インタラクションによって
 人々の音楽生活を
 より豊かにする研究

有機的な融合



歌声信号処理



インタフェース構築(インタラクション)

インタフェース構築・インタラクションデザイン

□ 様々なユーザに技術を使ってもらうため

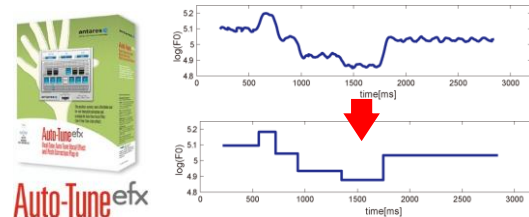
■ 対象ユーザの特性に合わせたインタフェース構築

- ユーザ視点での問題発見(インタラクションデザイン)

□ 歌声インタフェースの現状

■ プロ向けのインタフェースは既に多く存在

- 音高補正ツールAuto-tune 等



■ エンドユーザ向けのインタフェースは少ない

- 一方で、近年はエンドユーザが気軽に歌を発表したり、それを聴いて楽しんだりする文化が広がっている

エンドユーザによる歌声の公開

□ エンドユーザによるWeb上の歌声

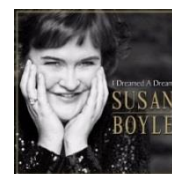
- 歌声合成 (VOCALOID)
- 歌ってみた / Me Singing

□ 多くの視聴者もいる



□ エビデンス

- 音楽CDの販売
- ニコニコ動画



- 27万曲+ (VOCALOID楽曲) 64万曲+ (歌ってみた)
- 10万回+ 再生 (それぞれ3200曲+, 4,650曲+)
- 100万回+ 再生 (それぞれ170曲+, 200曲+)

エンドユーザによる歌声の公開

□ エンドユーザによるWeb上の歌声

- 歌声合成 (VOCALOID)
- 歌ってみた / Me Singing

□ 多くの視聴者もいる



□ エビデンス

- 音楽CDの販売
- ニコニコ動画



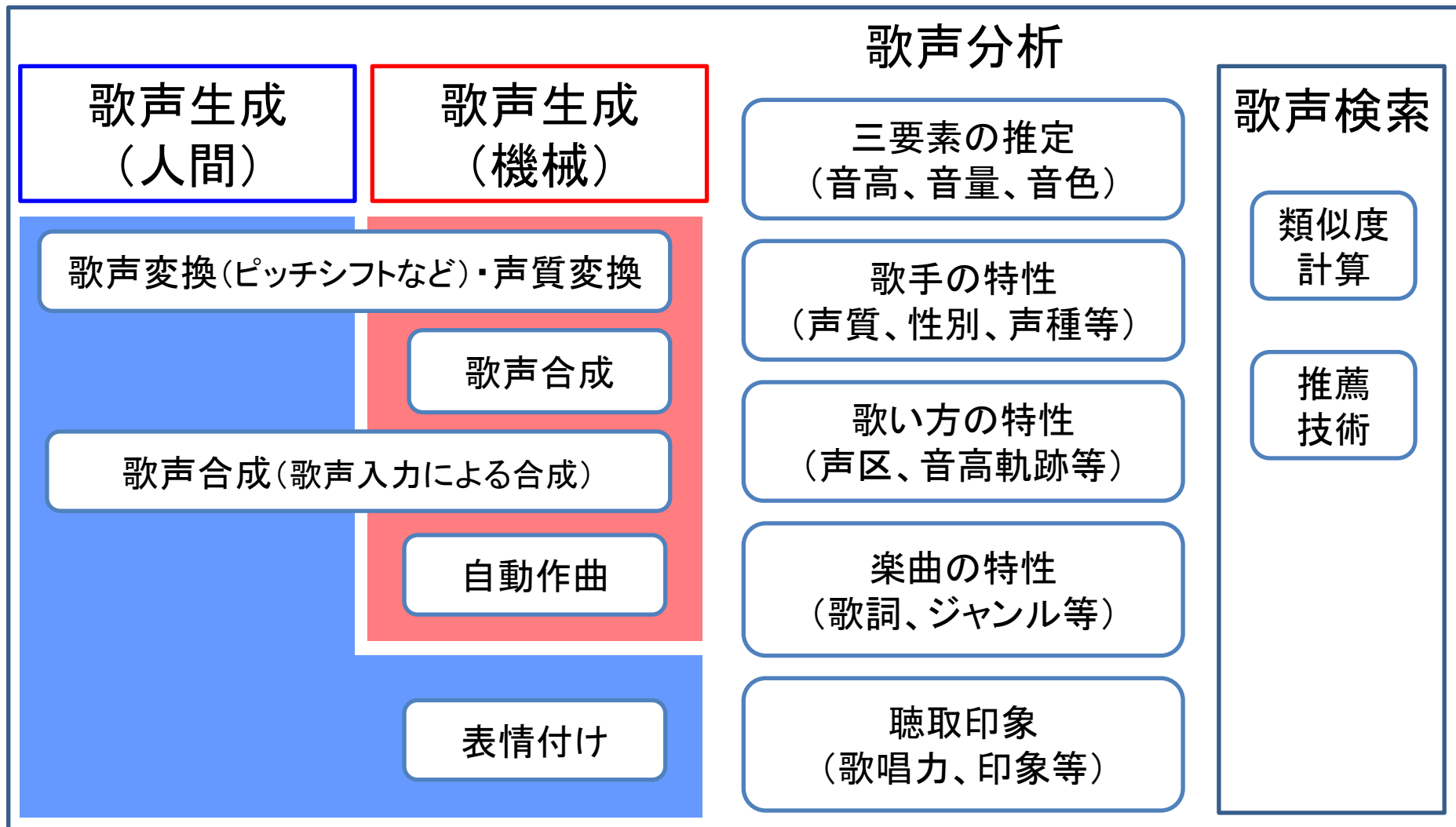
- 27万曲+ (VOCALOID楽曲) 64万
- 10万回+ 再生 (それぞれ3200曲+)
- 100万回+ 再生 (それぞれ170曲+)

「うたスキ動画」
(カラオケ動画のWeb投稿)
などが登場

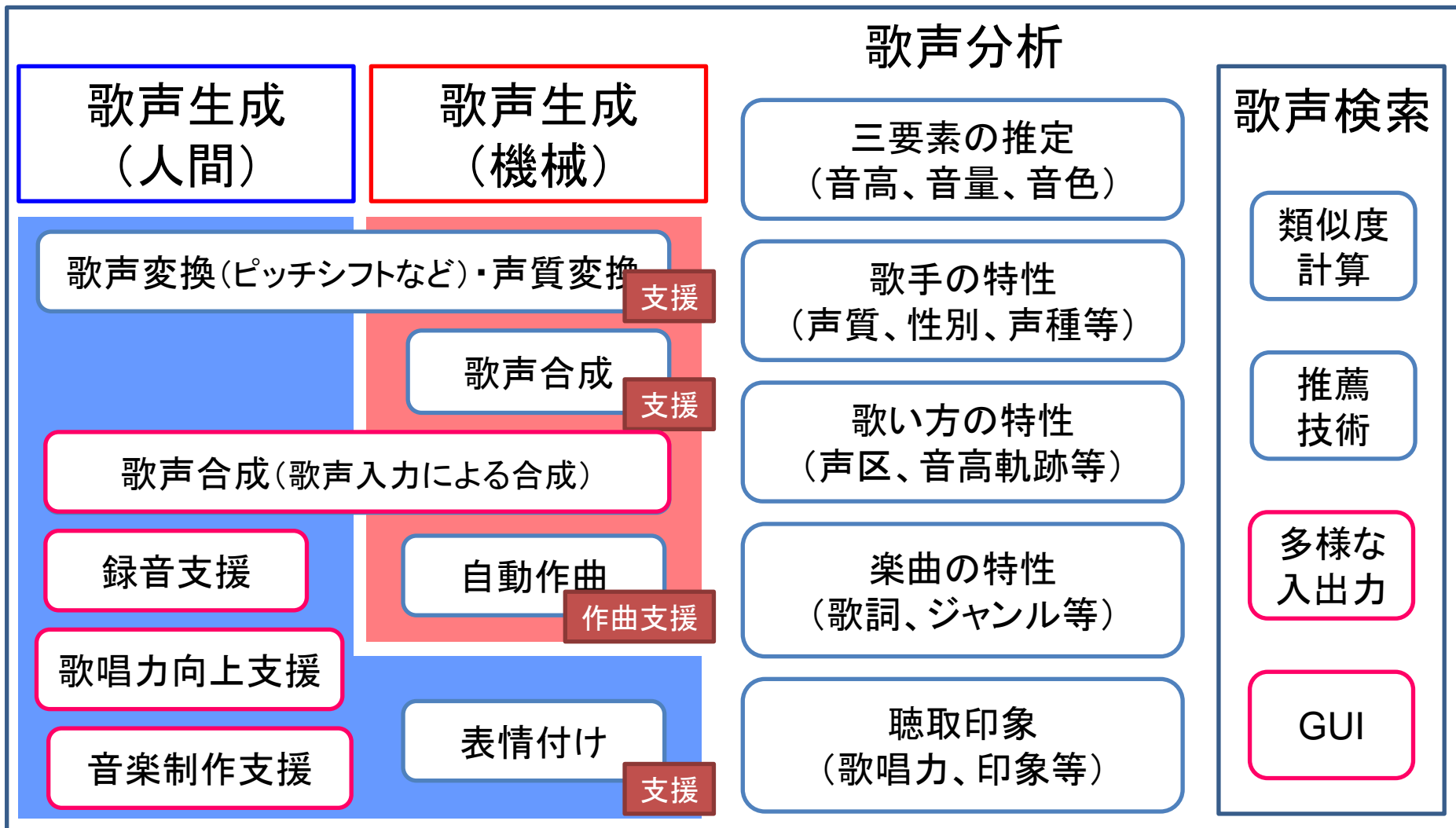
<http://utadou.jp/>

カラオケユーザへ **新たな体験を
提供する「場」**としての機能

歌声信号処理の全体像



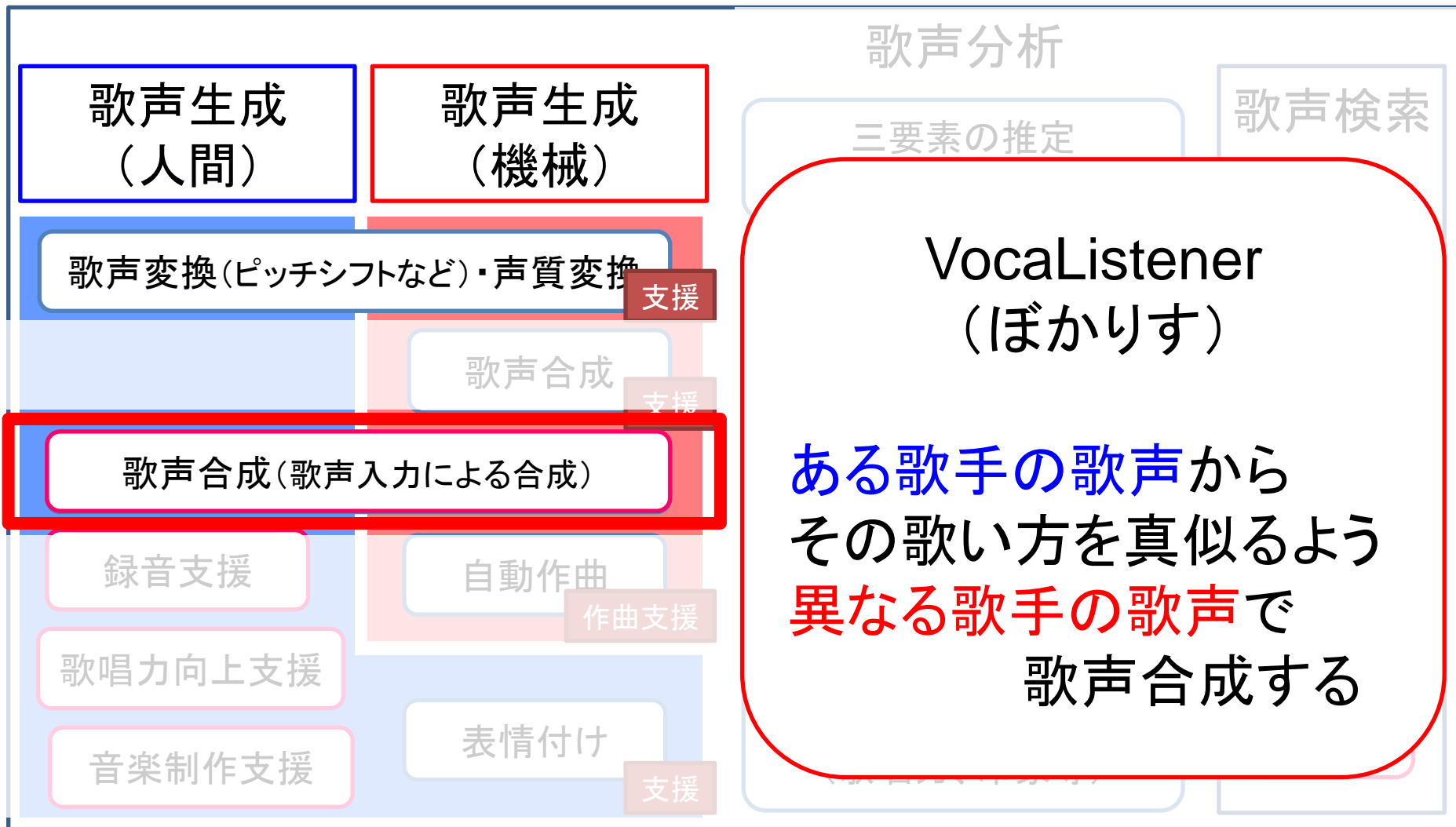
歌声インタフェースの全体像



歌声インタフェースの全体像



歌声インタフェースの全体像



VocaListener(ぼかりす)による合成結果

□ 【初音ミク】 PROLOGUE 【ぼかりす】

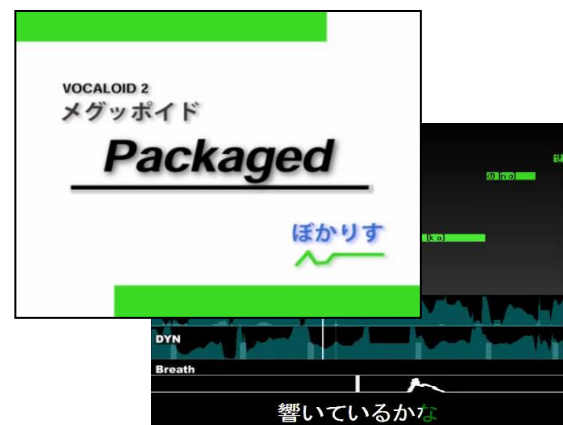
- 2008/04/28 [sm3128145] ~
- [361,055 views](#)



[original]

□ 【メグツポイド】 Packaged 【ぼかりす】

- 2010/10/4 [sm12320140] ~
- [31,404 views](#)



歌声合成とは

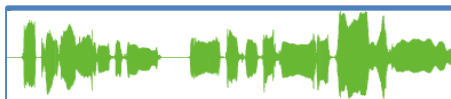
□ コンピュータ上で歌声を作ること

■ 入力

- 楽譜
- 歌詞
- 歌声合成パラメータ(声の高さ、音量)

■ 出力

- 歌声の音響信号



音符を並べて歌詞入力

+

細かなパラメータ調整

楽譜(音符)

歌詞

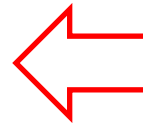
声の高さ
音量

これまでの歌声合成の問題点

□ 音符を並べて入力

問題

音符を並べただけでは
自然性が低い



歌詞

「ためらいもなく飲み込んで
次々と消化してゆく」

従来

VocaListener

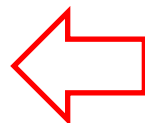


これまでの歌声合成の問題点

□ 音符を並べて入力

問題

音符を並べただけでは
自然性が低い



歌詞

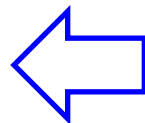
「ためらいもなく飲み込んで
次々と消化してゆく」

従来

VocaListener

問題

品質を高くしようとする
パラメータ調整に数時間
(数十時間)かかる



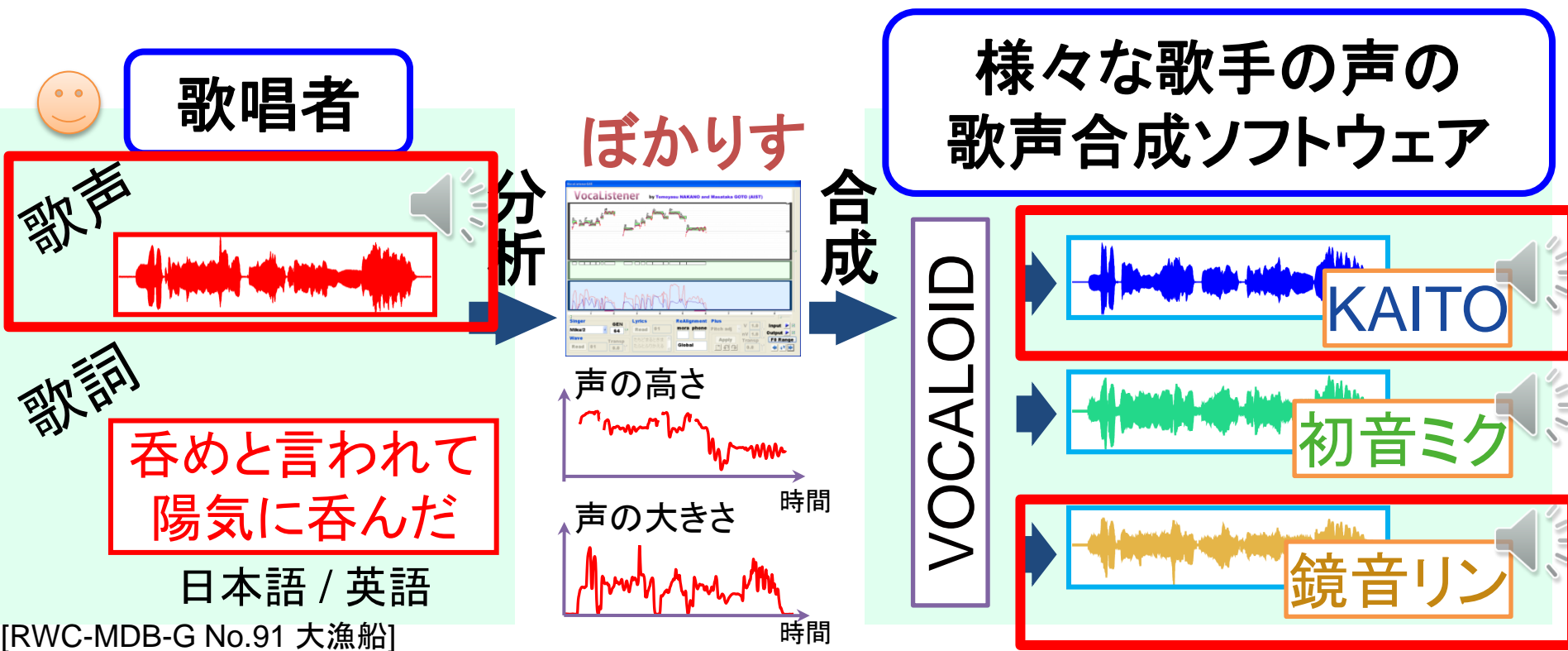
楽譜(音符)

歌詞

声の高さ
音量

VocaListener(ぼかりす) [中野, 後藤, 2008-]

- 歌うだけで**楽譜**と**自然なパラメータ**を自動生成
- 音源(歌手)**を手軽に切り替えて合成できる



これまでの歌声合成の問題点

□ 音符を並べて入力

問題

音符を並べただけでは
自然性が低い

歌詞

「ためらいもなく飲み込んで

歌詞

「呑めと言われて
陽気に呑んだ」

七
姫
の
声

人間
(お手本)

七
姫
の
声

楽譜入力
で合成

七
姫
の
声

VocaListener
(ぼかりす)

ばかりすの実現課題(1)

問題

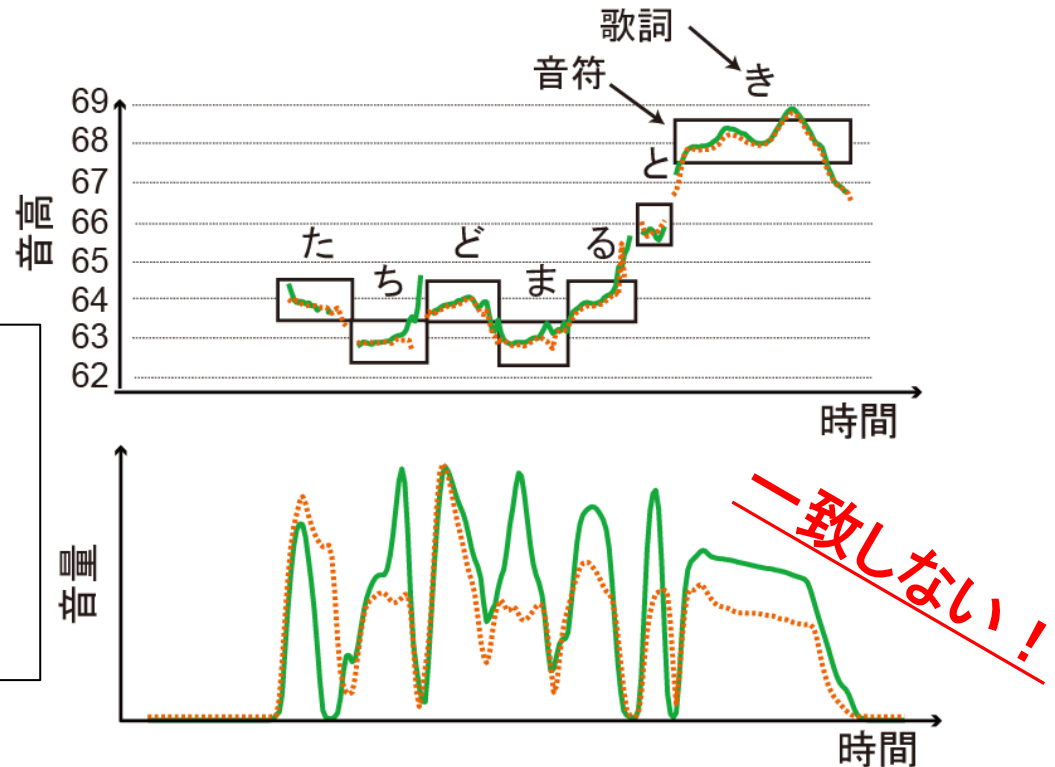
与えたパラメータ通りに合成されない

- 声の高さと声の大きさが歌声合成ソフトに依存

同じパラメータを与えて
二つの歌声合成ソフト

Aと**B**

で合成した結果



ばかりすの実現課題(1)

問題

与えたパラメータ通りに合成されない

解決策

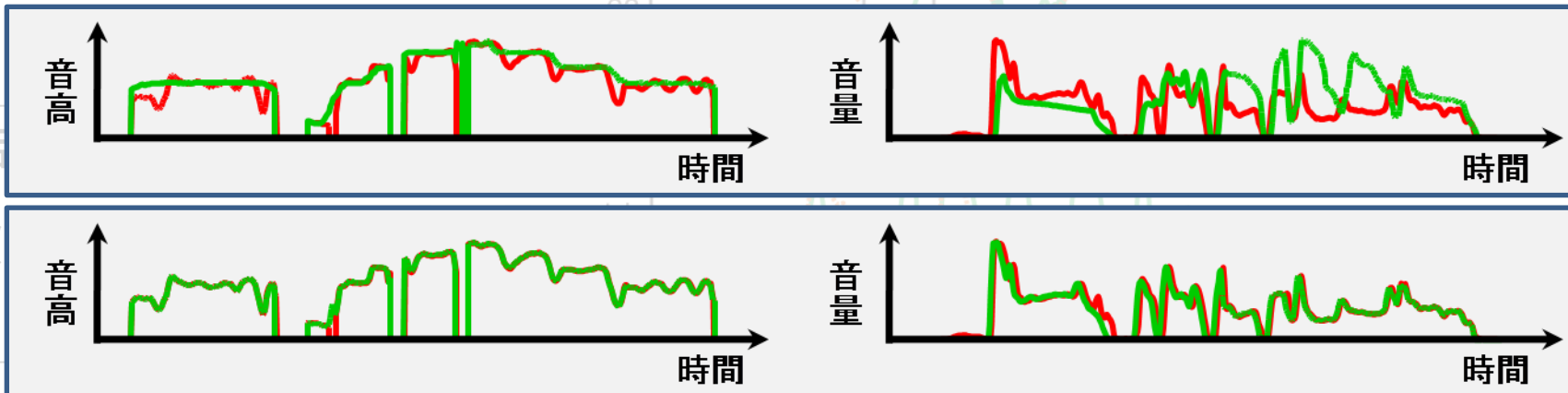
合成結果を分析
して、目標の値に近くなるように
パラメータを反復更新

— 目標
— 合成

反復前



反復後



ばかりすの実現課題(2)

問題

歌詞のどこを、いつ歌っているか分からない

- 歌声の信号と歌詞のテキストだけがある

こんな熱い夢～♪



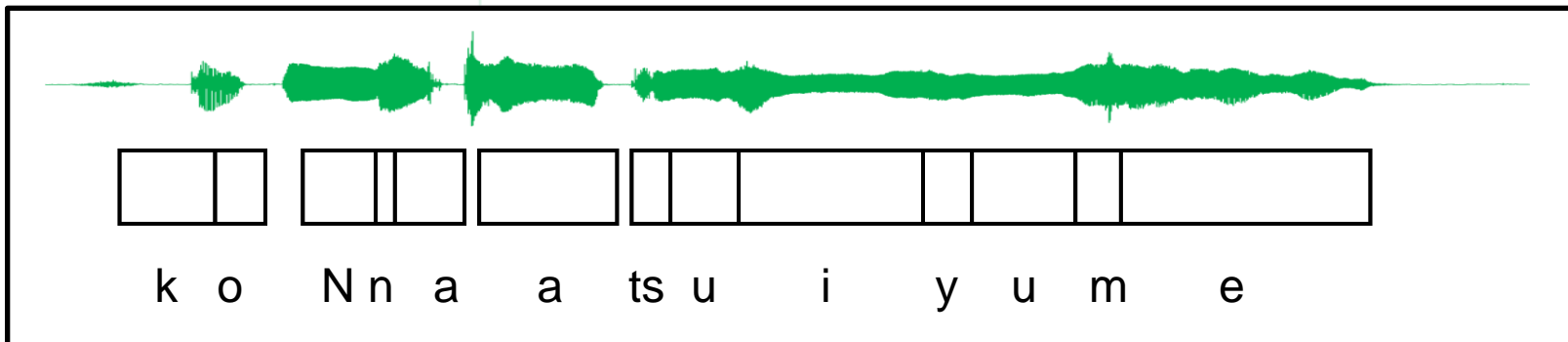
ばかりすの実現課題(2)

問題

歌詞のどこを、いつ歌っているか分からない

解決方策

歌声専用の音響モデル
 (HMM: 声の響き方を学習したモデル)
 を構築して高精度に対応付け



ばかりすの実現課題(2)

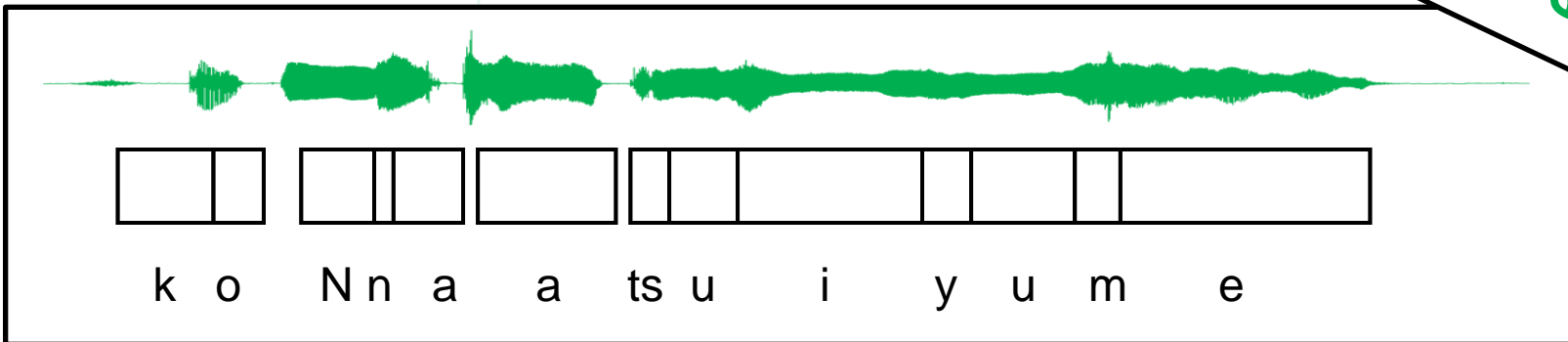
問題

歌詞のどこを、いつ歌っているか分からない

解決方策

歌声専用の音素モデル
 (HMM: 声の響き方を学習)
 を構築して高精度に対応

しかしまだ問題がある

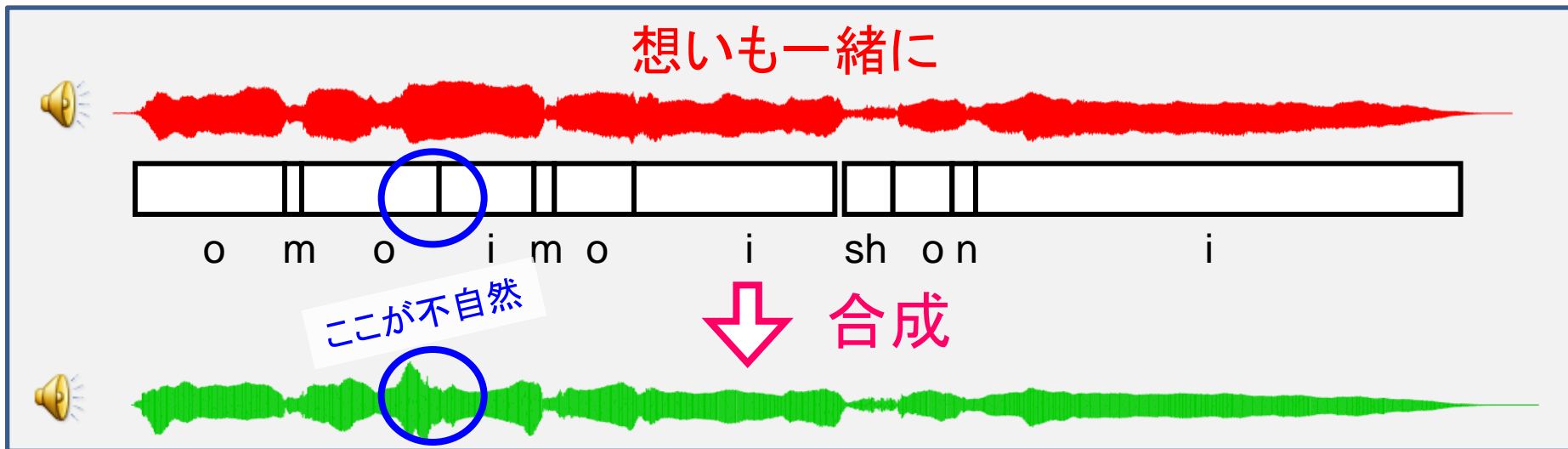


ばかりすの実現課題(3)

問題

合成のためには100%に近い精度が必要

- 高い精度であるが、数%の誤りは不可避



ばかりすの実現課題(3)

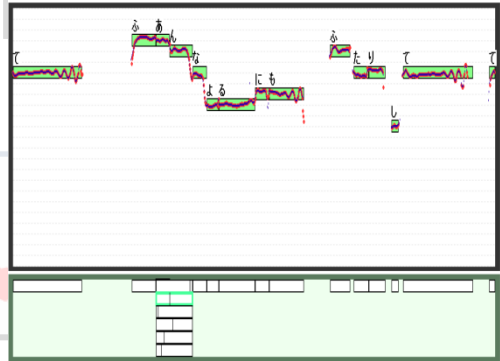
インタラクション

問題

合成のためには100%に近い精度が必要

解決策

ダメ出しインターフェース
(聴いて指摘するだけで訂正)



お も い も い しよ に

↑ 😊 ここが間違っている

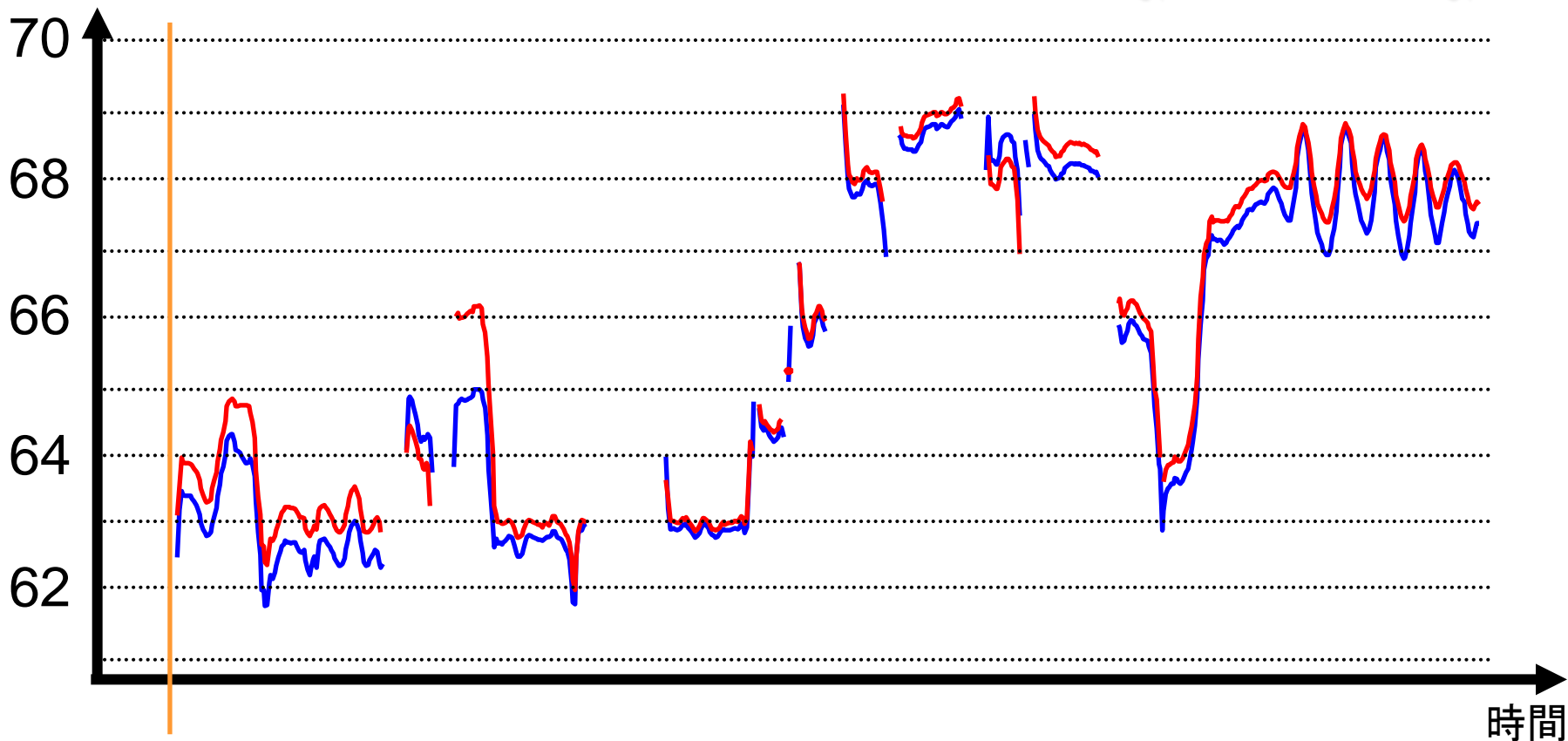
歌唱力補正も可能

音高(声の高さ)

元歌



合成歌唱



歌詞

い ま も せ つ な い

す が た さ が し て い る よ

い ま も せ つ な い

す が た さ が し て い る

よ

VocaListener VOCALOID3 Job Plugin

- 2012/10/19 ヤマハ株式会社より製品化



- 【ぼかりす公式デモ】迷子ライフ【歌：がくっほいどPower】

– 2012/10/22 [sm19181691] ~

– 8,489 views



- 塾講師なので本気出して卒業ソングをつくってみたよ【初音ミク など】

– 2013/3/22 [sm20403611] ~

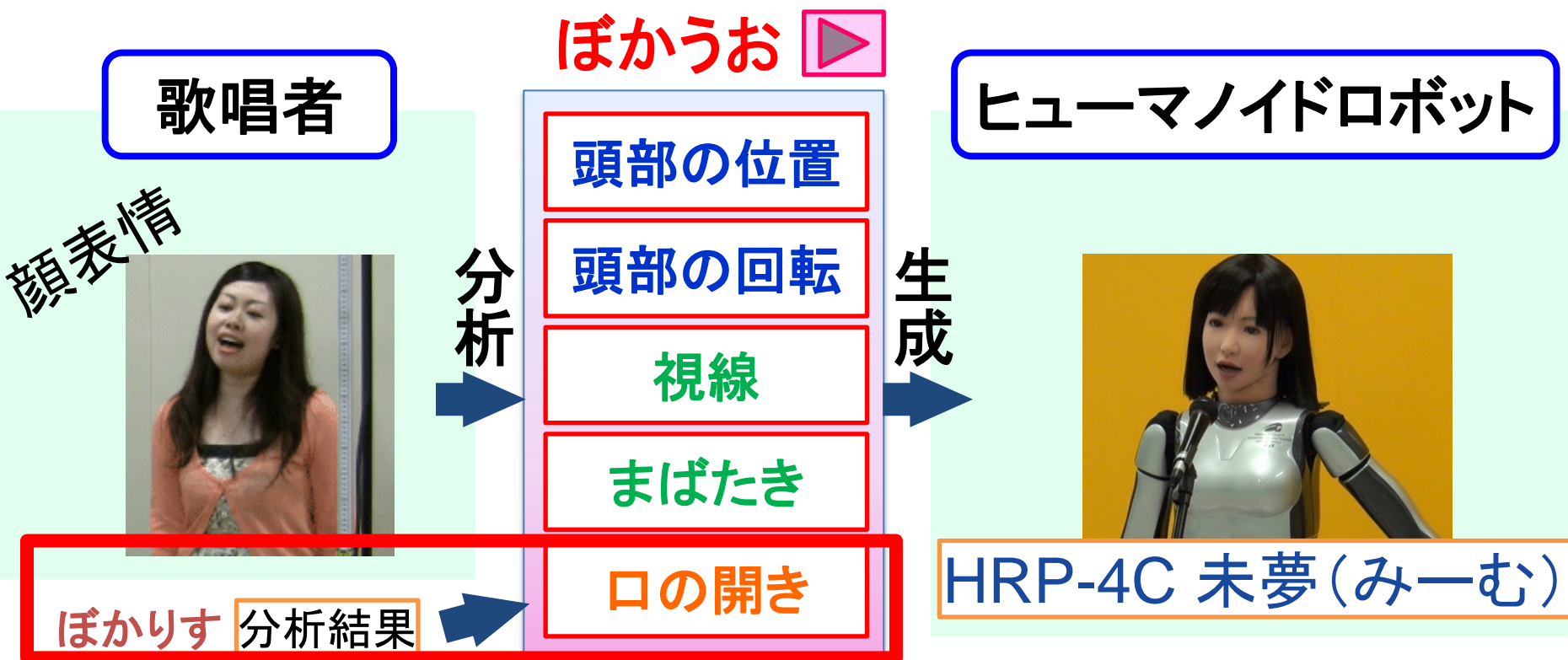
– 111,977 views



VocaWatcher (ぼかうお)

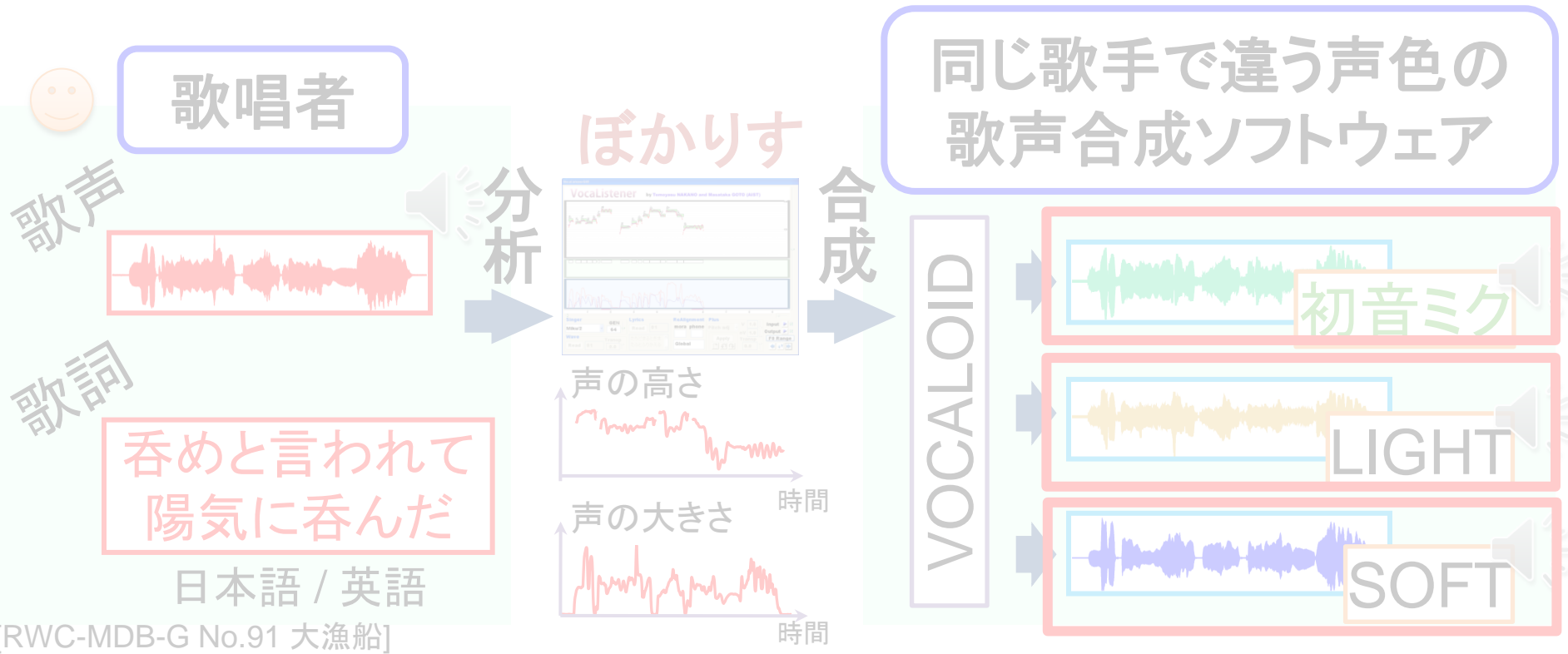
[Kajita, Nakano, Goto+ 2010-]

- 歌うだけでヒューマノイドロボットの
自然な顔動作パラメータを自動生成



複数の声色で歌声合成が可能


- 中間の声色を作ることはできなかった
- ユーザの声色変化を真似ることはできなかった



VocaListener2 (ばかりす2) [中野, 後藤, 2010-]

- **VocaListener**で音高と音量を真似て全て合成
- 声色変化を真似るようにそれらを**混ぜ合わせる**

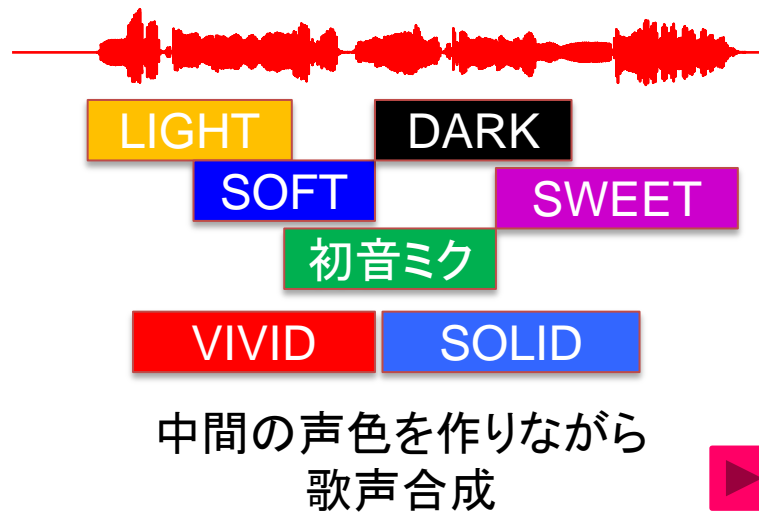
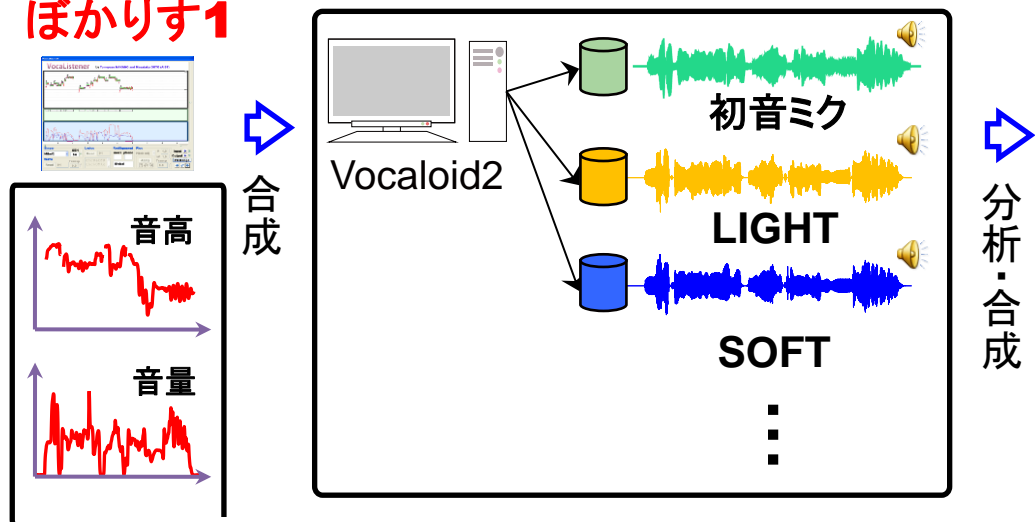
ユーザ歌唱 「呑めと言われて
陽気に呑んだ」



ばかりす2による合成
(イメージ)

分析 初音ミクと初音ミク・アペンド

ばかりす1

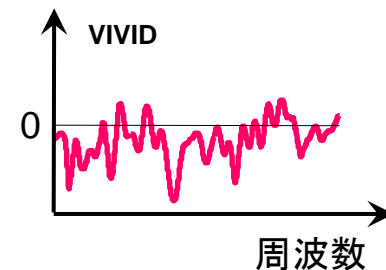
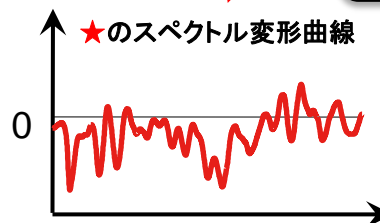
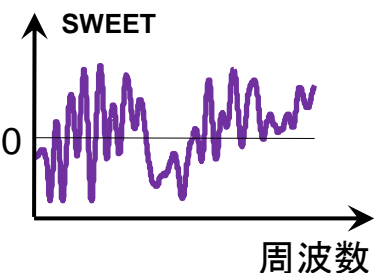
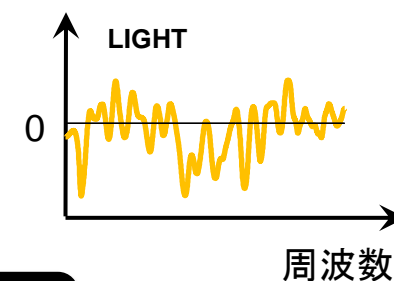
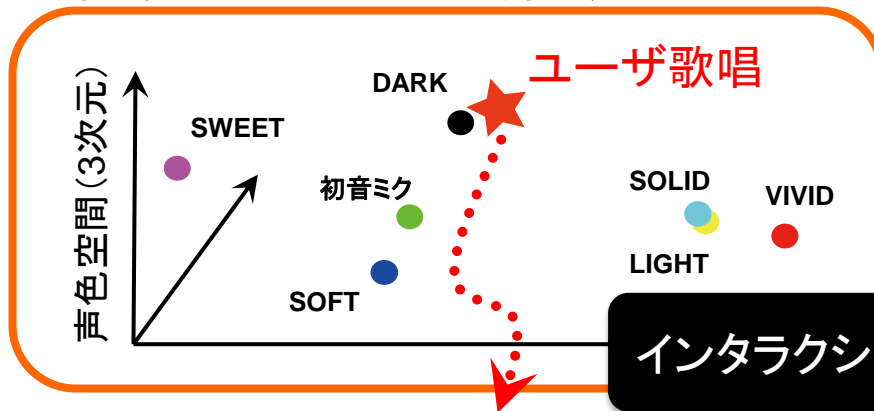
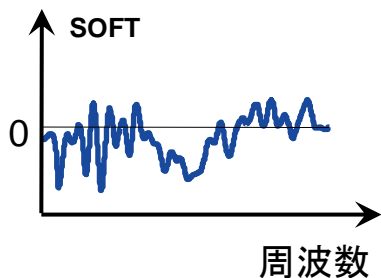
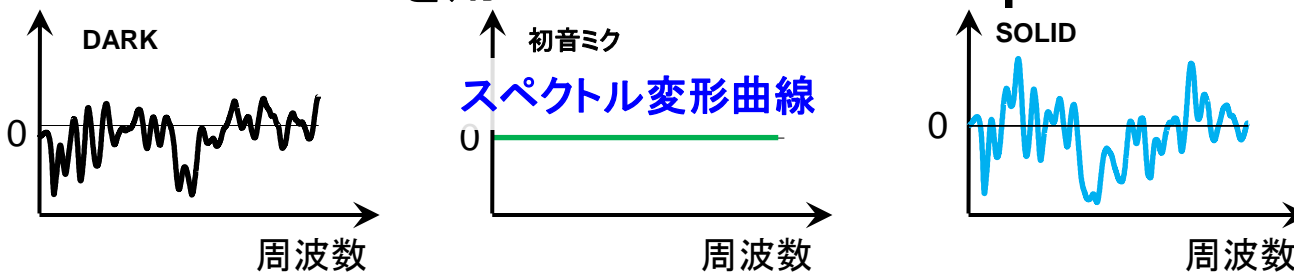


 **Voice timbre changes** +
  **Voice timbre** =
  **Synthesized singing**

声色空間に基づくスペクトル変形曲線の推定

- スペクトル包絡の主成分分析により3次元空間へ射影
- Radial Basis Functionを用いたVariational Interpolation

[Turk et al., 2008]



ばかりす2による歌声合成デモ

VocaListener2

1. VocaListener1 による歌声合成結果

初音ミク (Vocaloid2)

音高・音量を真似て歌声合成

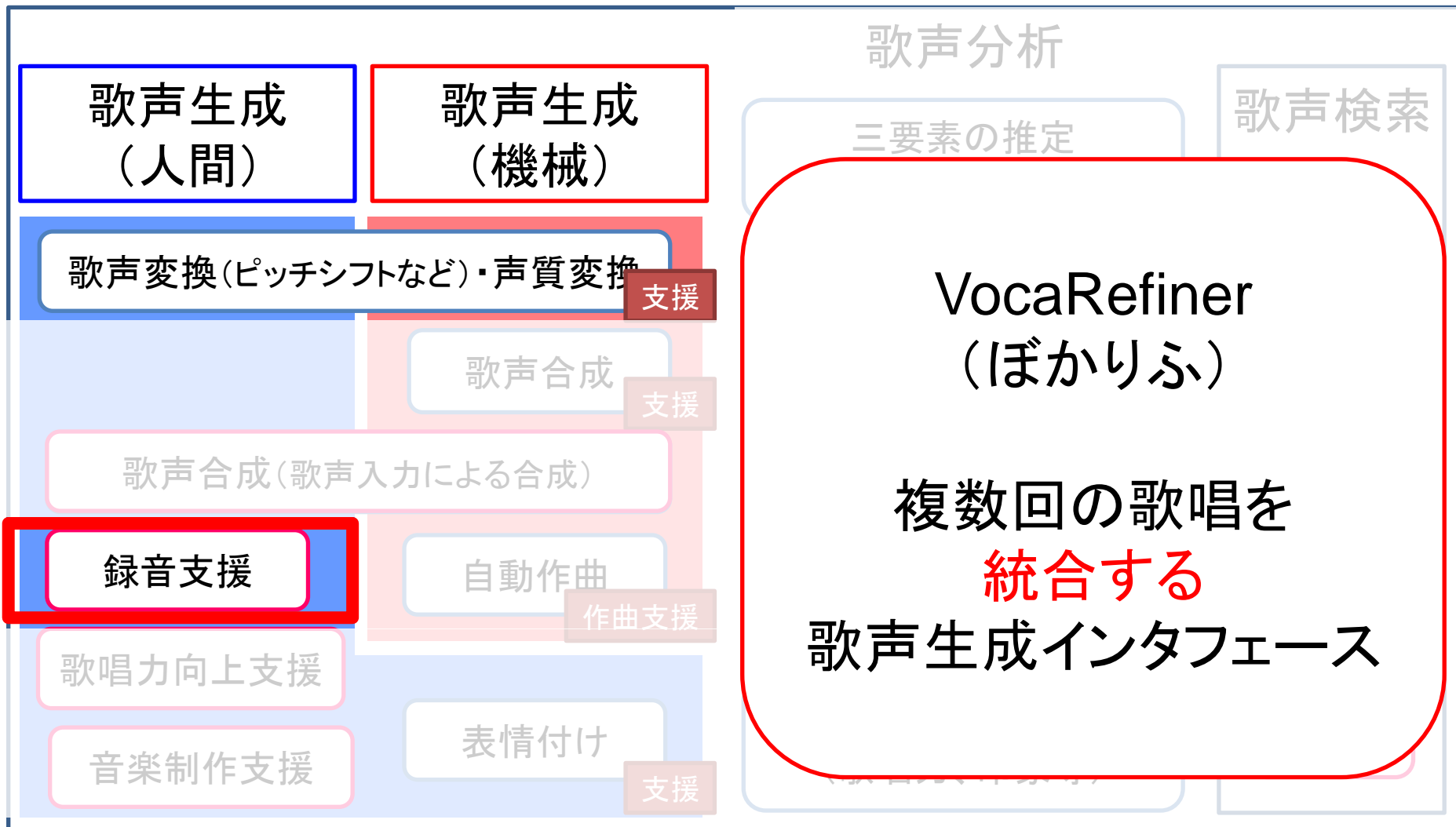
2. VocaListener2 による歌声合成結果

初音ミク (Vocaloid2)

と初音ミク・アペンド (Vocaloid2)

音高・音量・声色変化を真似て歌声合成

歌声インタフェースの全体像



歌をうまく録音したい

- 誰でも気軽に自分の歌を公開して共有する文化



- 歌ってみた(ニコニコ動画)

- 投稿動画数: 64万件以上

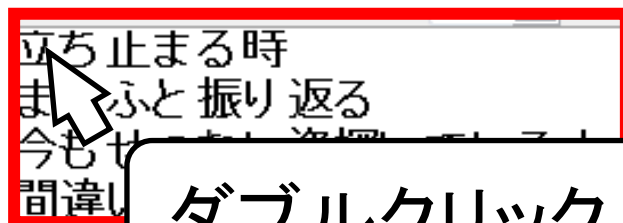
歌を一度で完璧に録音するのは困難

F0適応多重フレーム統合分析法とF0適応F0推定法に基づく [中野, 後藤, 2012]

歌声生成インタフェースVocaRefiner

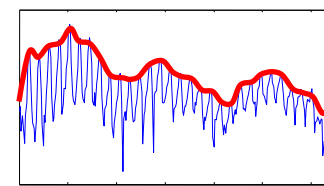
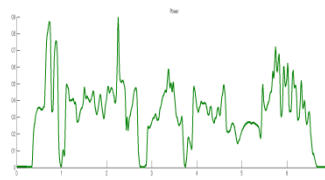
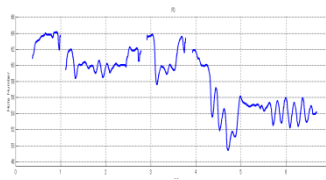
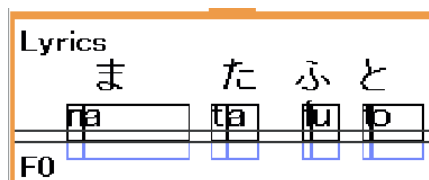
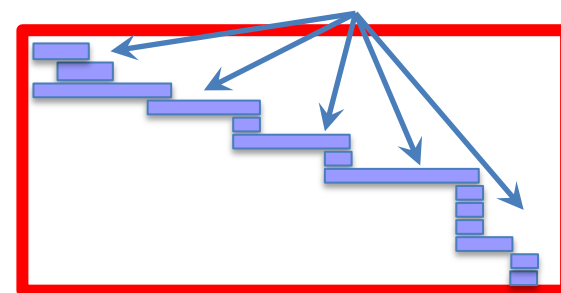
□ 歌詞を使って手軽に録音する

- 効率的に録音できる
- 日本語と英語に対応



□ 歌を歌い直して(複数回録音)統合する

- 音素単位で統合できる
- 音の三要素で統合できる
- 過去に歌った歌を活用する



統合機能

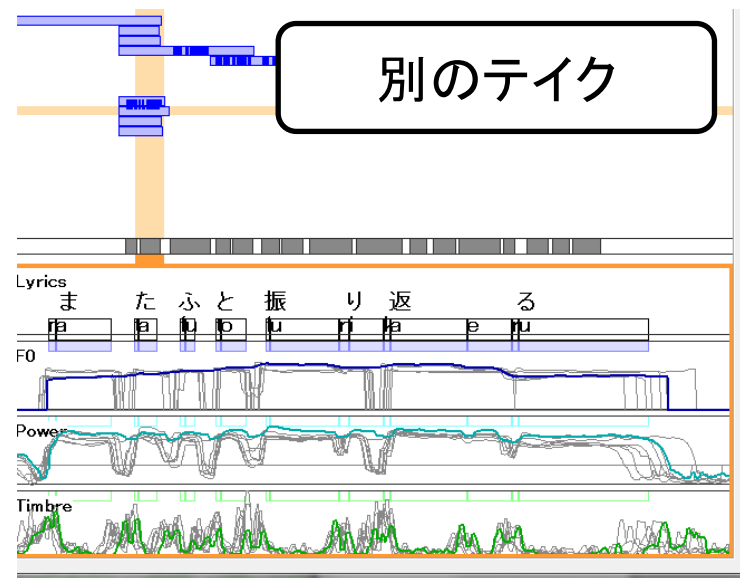
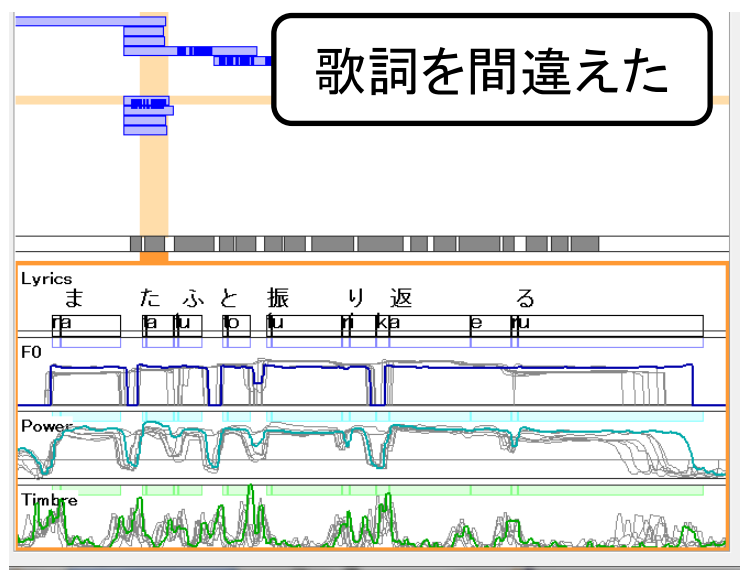
※歌詞を間違えたと想定して
ハミング(ラララ歌唱)で歌った

□ 三要素を組み替えながら統合

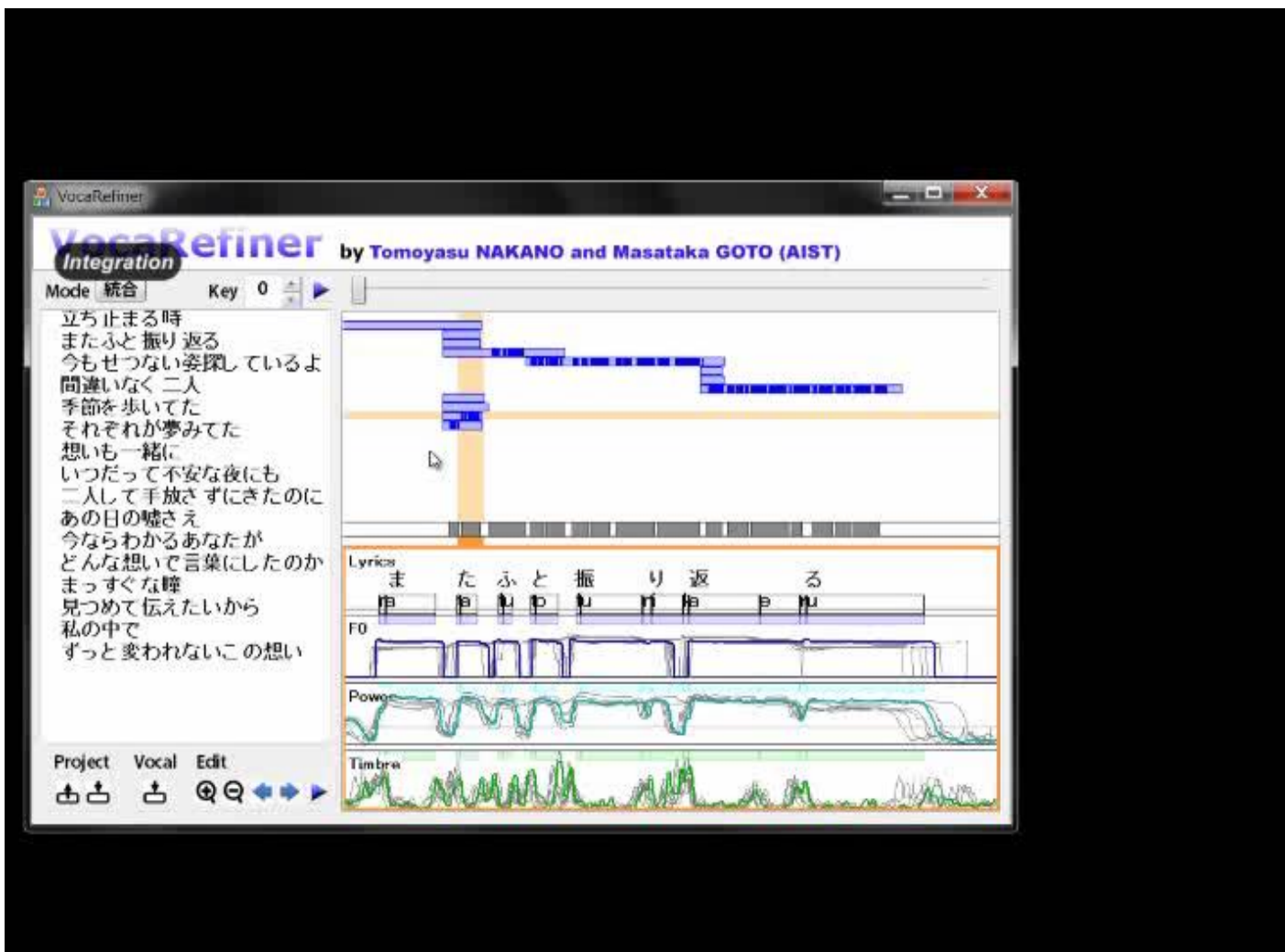
■ 「感情は完璧だったが歌詞を間違えた」場合

- 歌い方(音高・音量)はその録音を使用

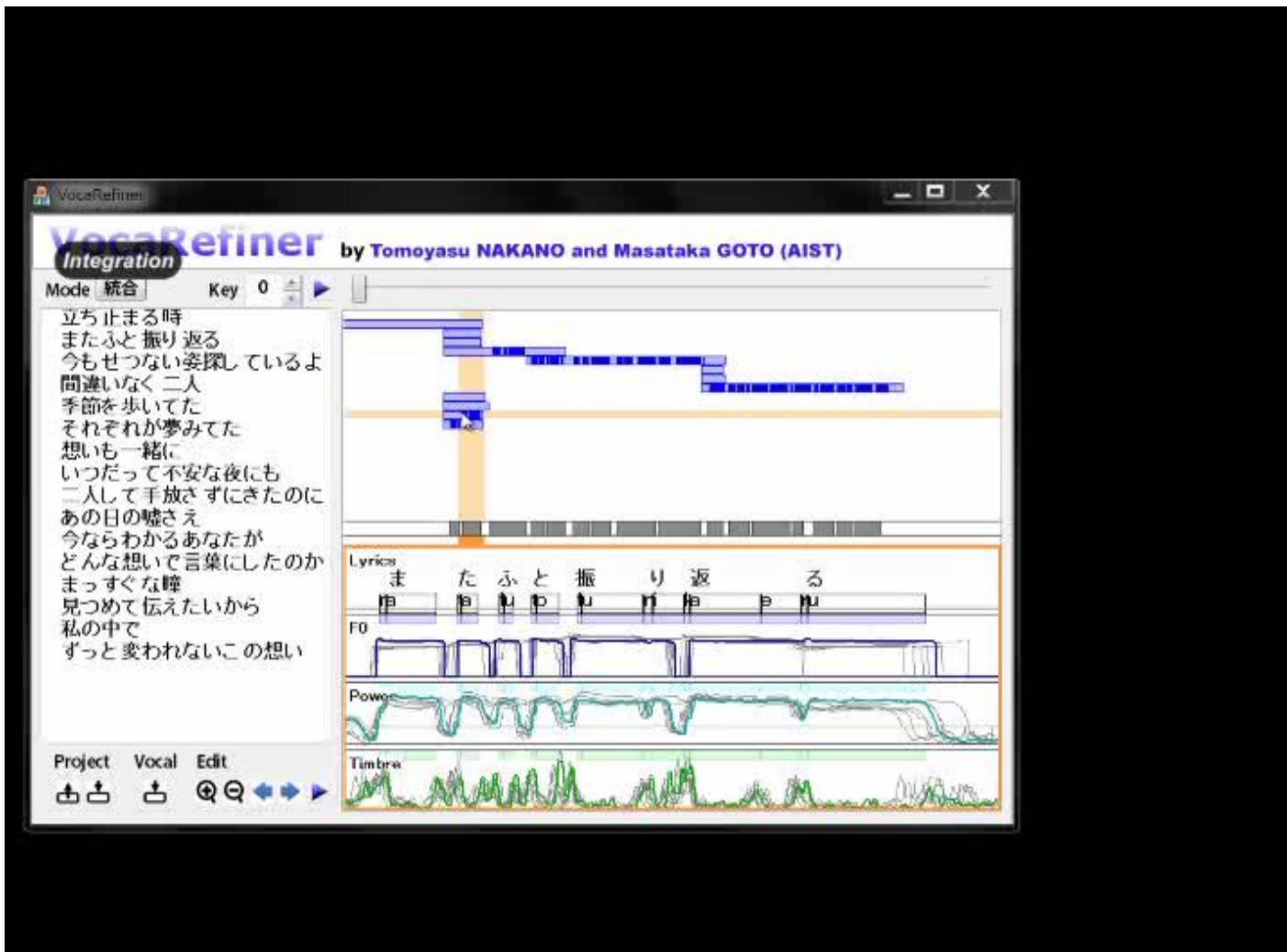
歌詞が正しい**過去録音の音色**を使用



テイク1: 音高は正しいが歌詞を間違えた



テイク2: 歌詞は正しいが音高が誤った



統合(合成): 正しい歌詞と正しい音高

The screenshot displays the VoceRefiner software interface. The main window is titled "VoceRefiner Integration by Tomoyasu NAKANO and Masataka GOTO (AIST)". The mode is set to "統合" (Integration) and the key is "0".

On the left, the lyrics are displayed in Japanese:

立ち止まる時
 またふと振り返る
 今もせつない姿探しているよ
 間違はなく二人
 季節を歩いた
 それぞれが夢みてた
 想いも一緒に
 いつだって不安な夜にも
 二人して手放さずにきたのに
 あの日の嘘さえ
 今ならわかるあなたが
 どんな想いで言葉にしたのか
 まっすぐな瞳
 見つめて伝えたいから
 私の中で
 ずっと変わらないうの思い

A large black circular overlay in the bottom left corner contains the text: **Recompose (synthesis) & Playback**.

The main display area shows a pitch contour plot with blue bars representing notes. Below the plot, there are four tracks: "Lyrics" (showing the Japanese text), "F0" (fundamental frequency), "Power", and "Timbre".

In the top right corner, a small window titled "C:\Users\Yn..." displays a list of files or settings.

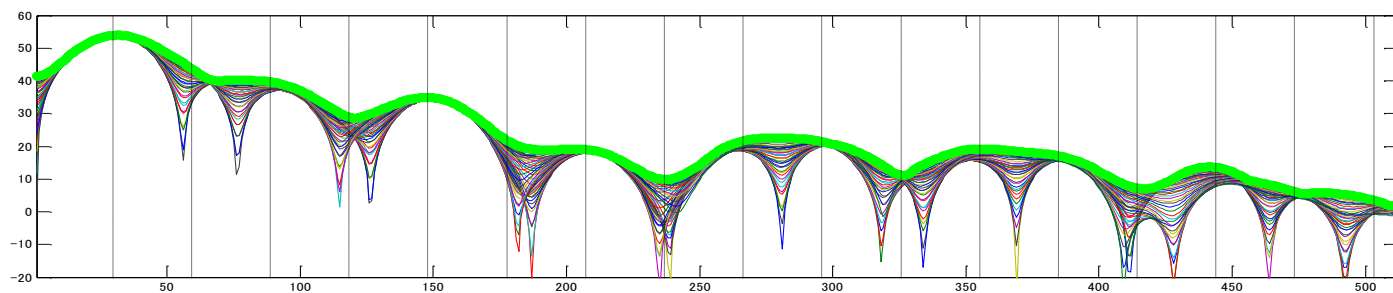
話声と歌声の分析合成のための F0適応多重フレーム統合分析法

[中野, 後藤, 2012]

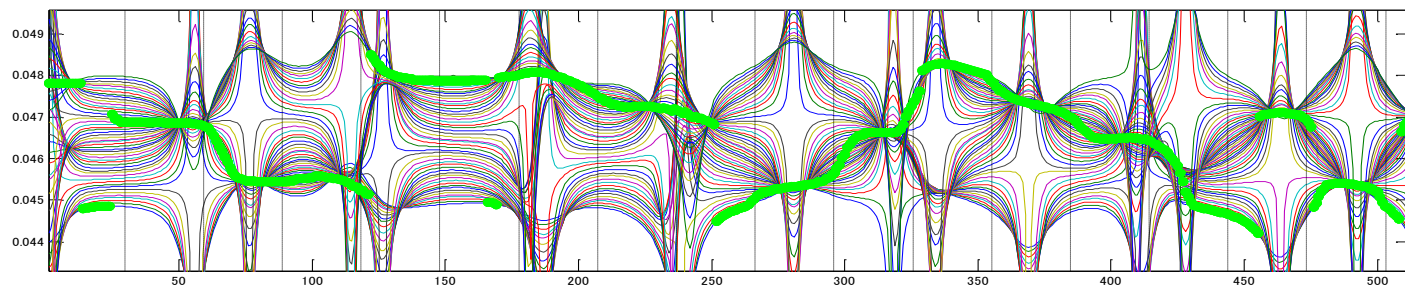
□ 高い精度と時間分解能で推定可能

- スペクトル包絡: **音素、声質、声色**
- 群遅延(位相の周波数微分): **波形の形、声の質感**

F0適応スペクトル
と
最大包絡

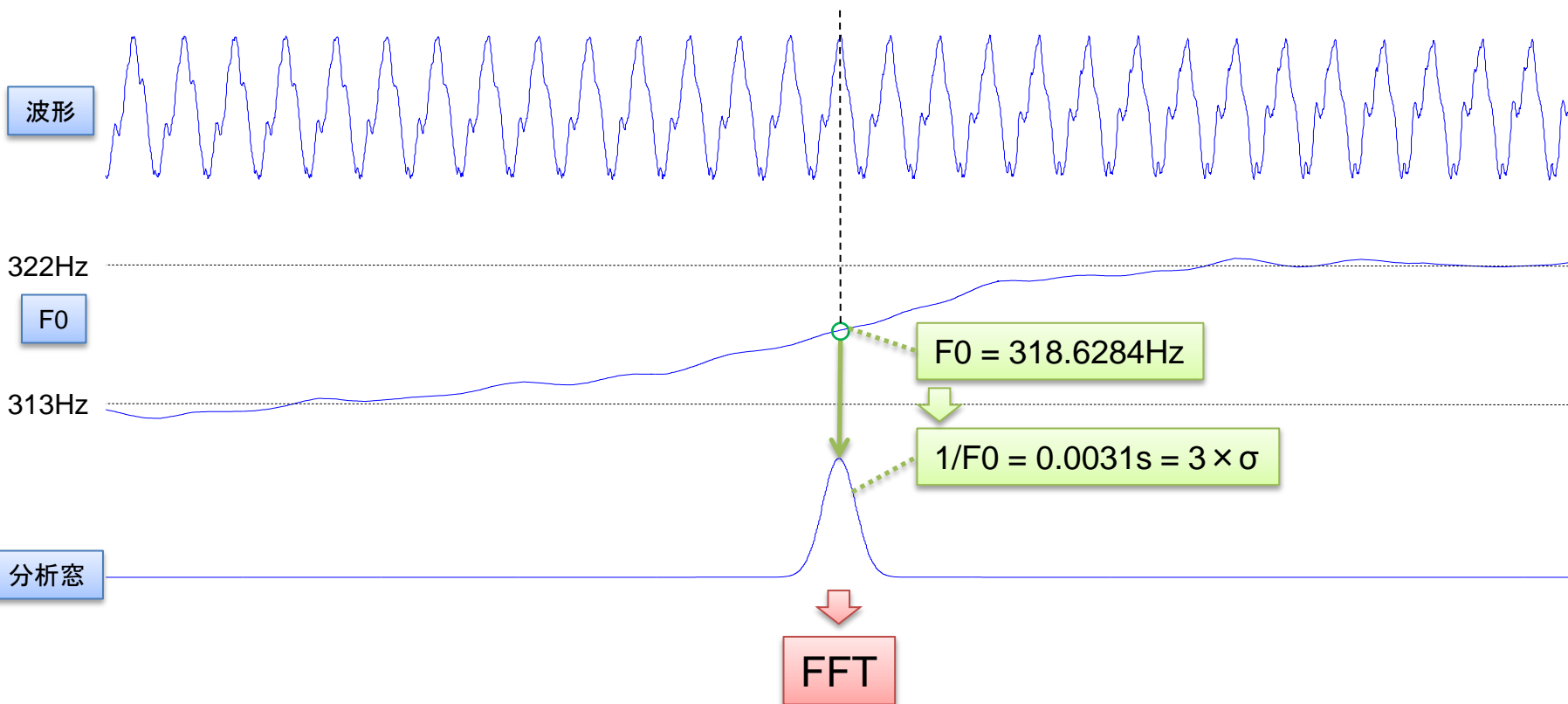


最大包絡
に対応する
群遅延



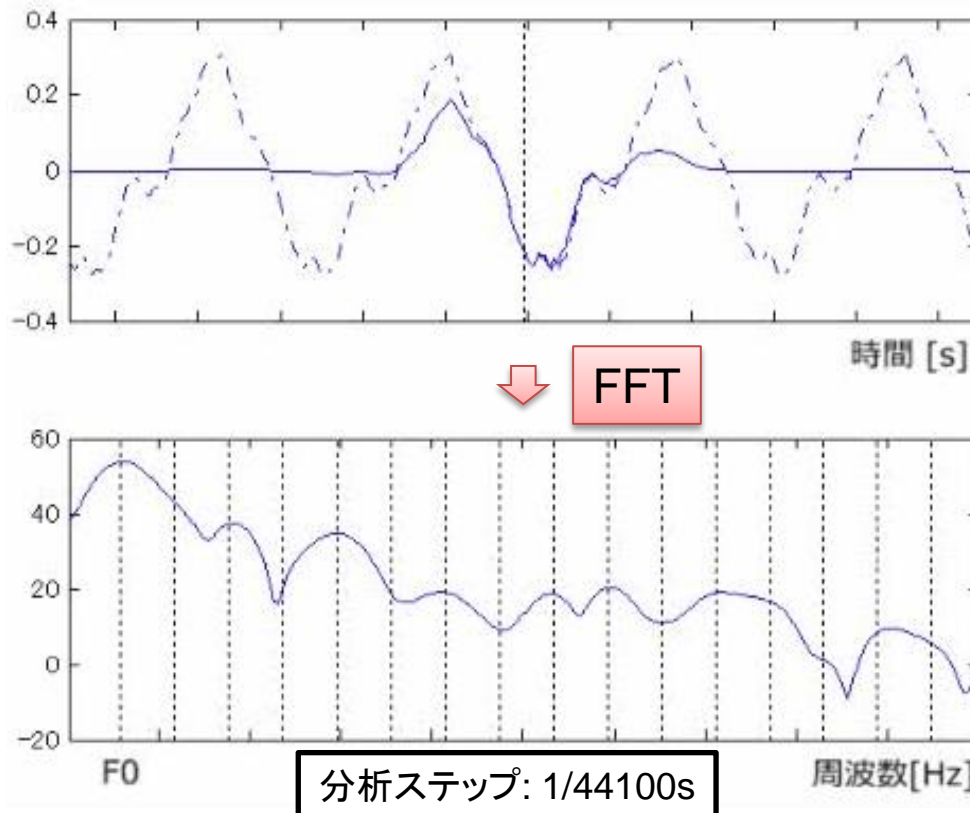
F0適応分析:F0に適応した長さの短い窓

- F0に適応した長さの短い窓で
全時間(全サンプリング点)に対してFFT



F0適応分析：全時間に対してFFT

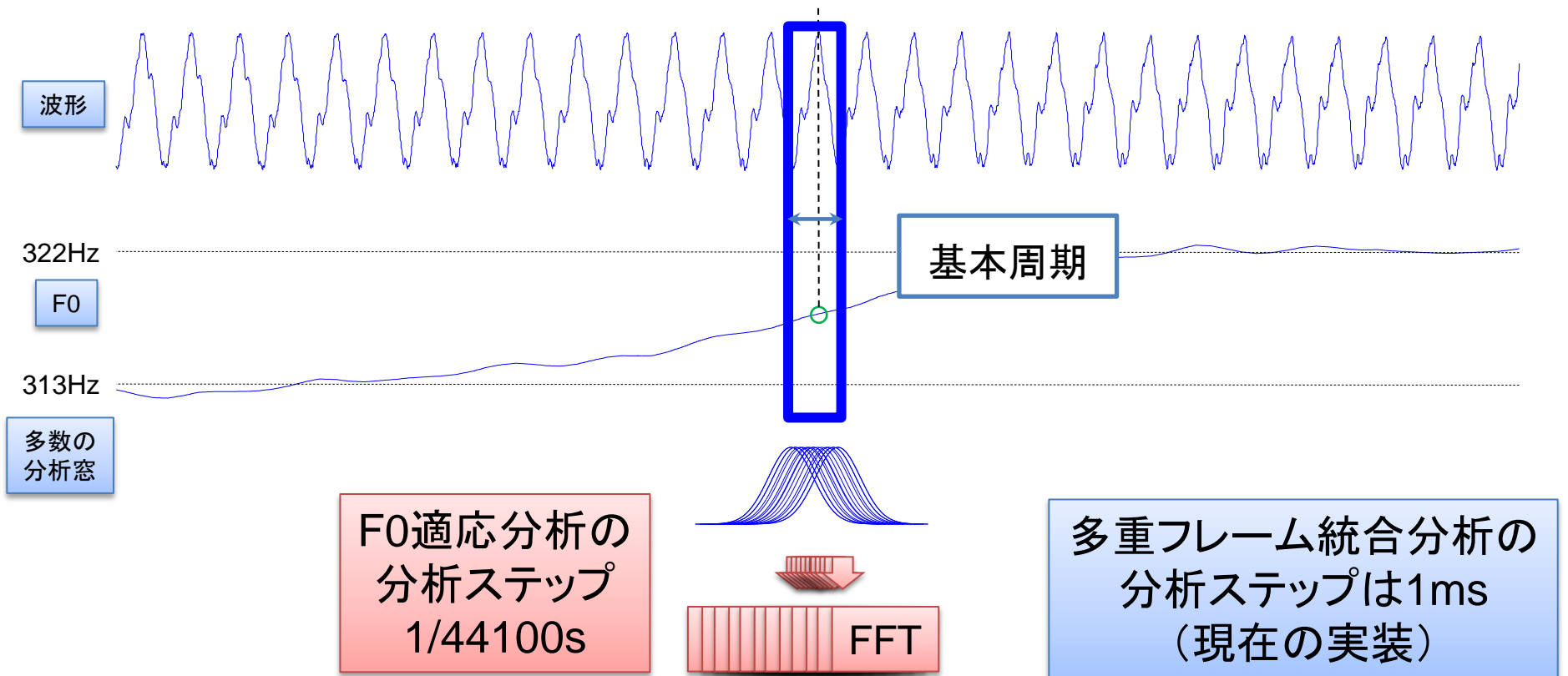
- F0に適応した長さの短い窓で
全時間(全サンプリング点)に対してFFT



多重フレーム統合分析:基本周期の範囲で統合

□ 基本周期の範囲で

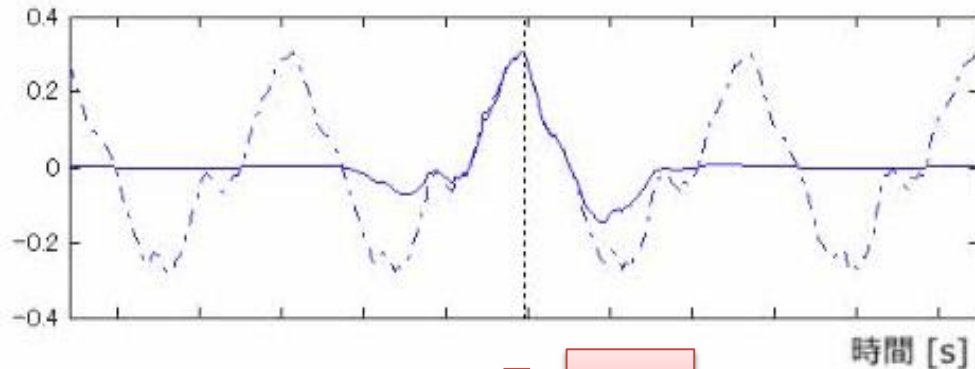
F0適応スペクトルを統合して変動を消去



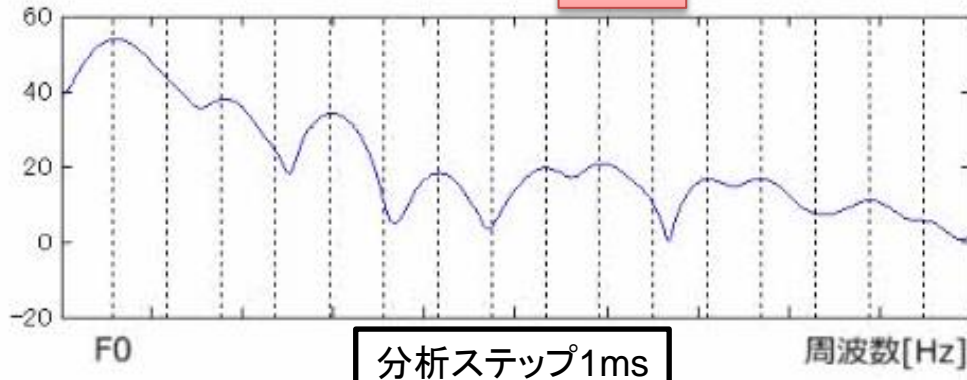
多重フレーム統合分析:基本周期の範囲で統合

□ 基本周期の範囲で

F0適応スペクトルを統合して変動を消去



↓ FFT



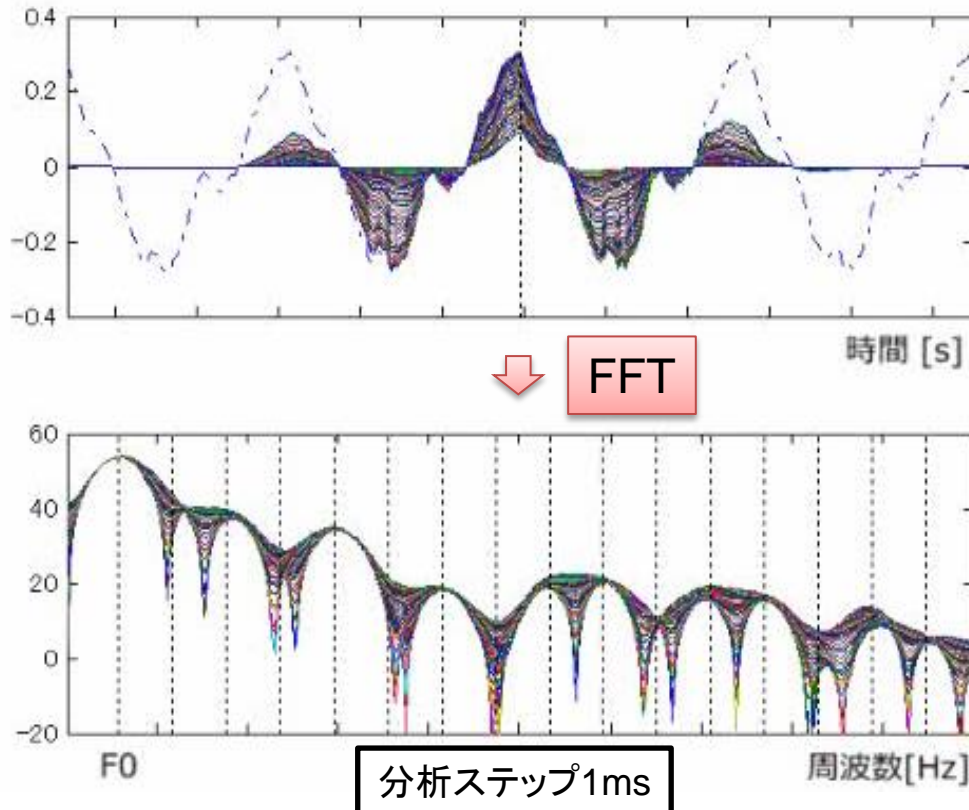
統合を
行わない場合

谷が発生

多重フレーム統合分析:基本周期の範囲で統合

□ 基本周期の範囲で

F0適応スペクトルを統合して変動を消去



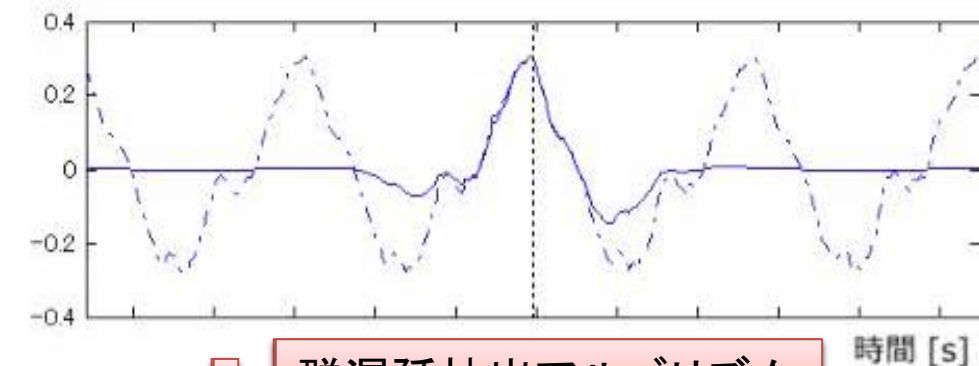
統合(重畳)
により

どこで分析しても
谷が埋まる

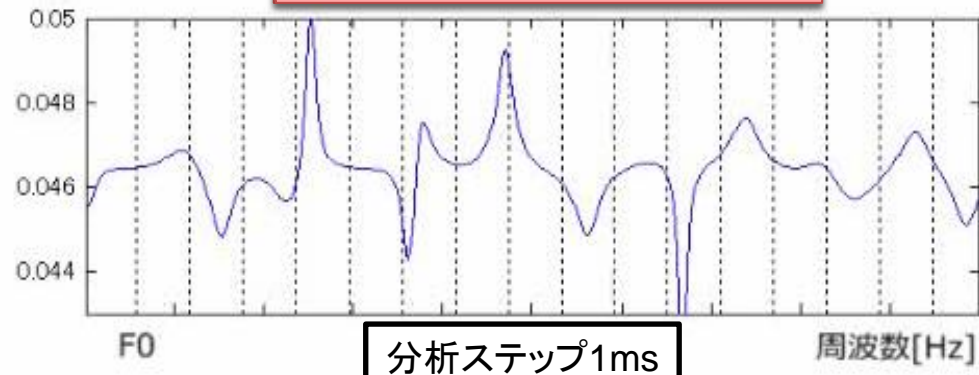
多重フレーム統合分析:基本周期の範囲で統合

□ 基本周期の範囲で

F0適応群遅延を統合して変動を消去



群遅延抽出アルゴリズム



統合を
行わない場合

ピークが発生

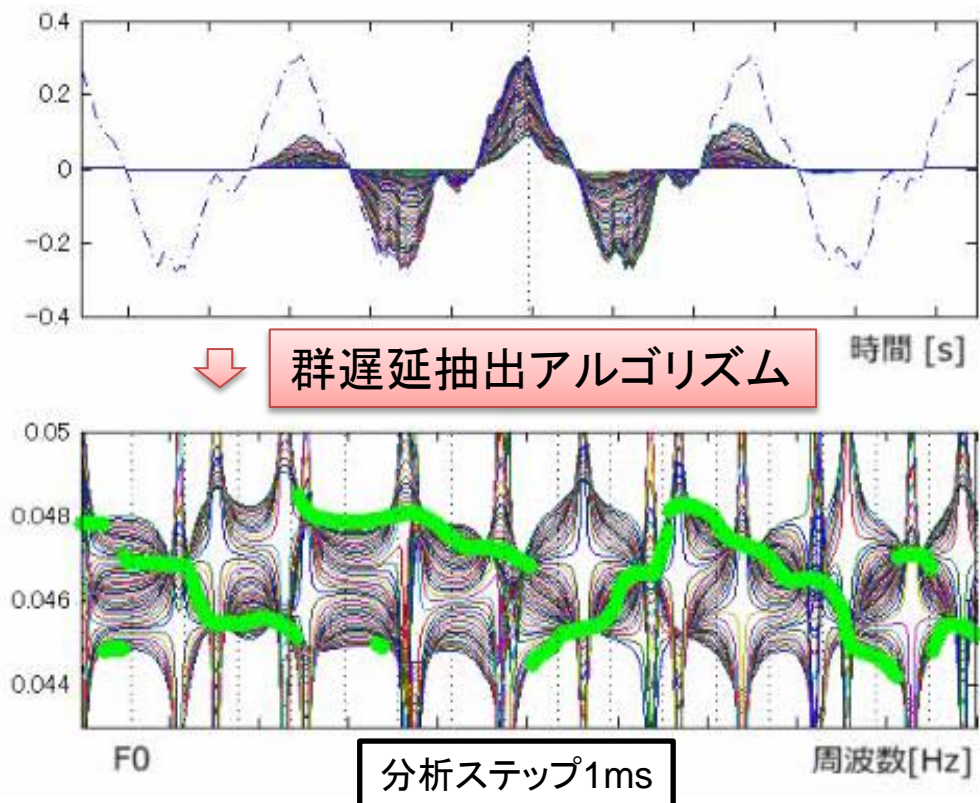
多重フレーム統合分析:基本周期の範囲で統合

□ 基本周期の範囲で

F0適応群遅延を統合して変動を消去

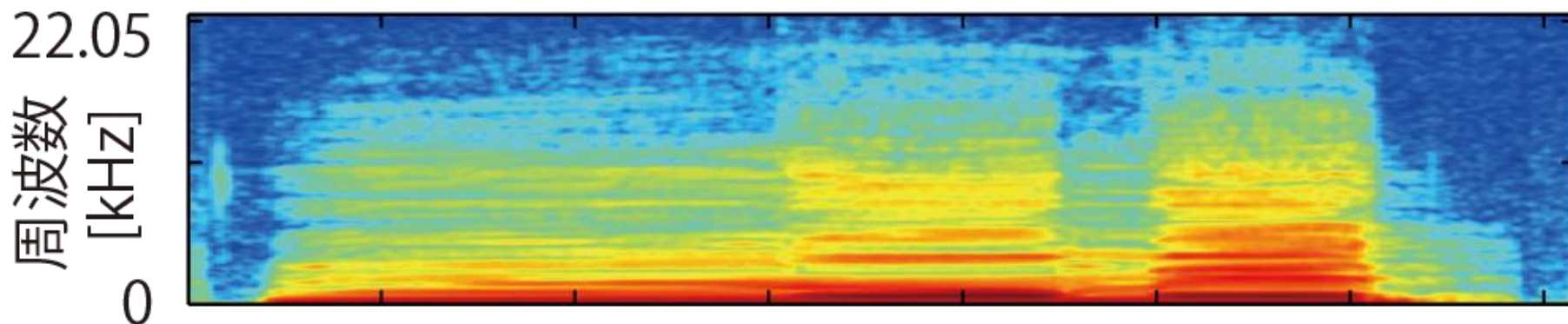
統合(重畳)
により

どこで分析しても
相対的に同じ形状

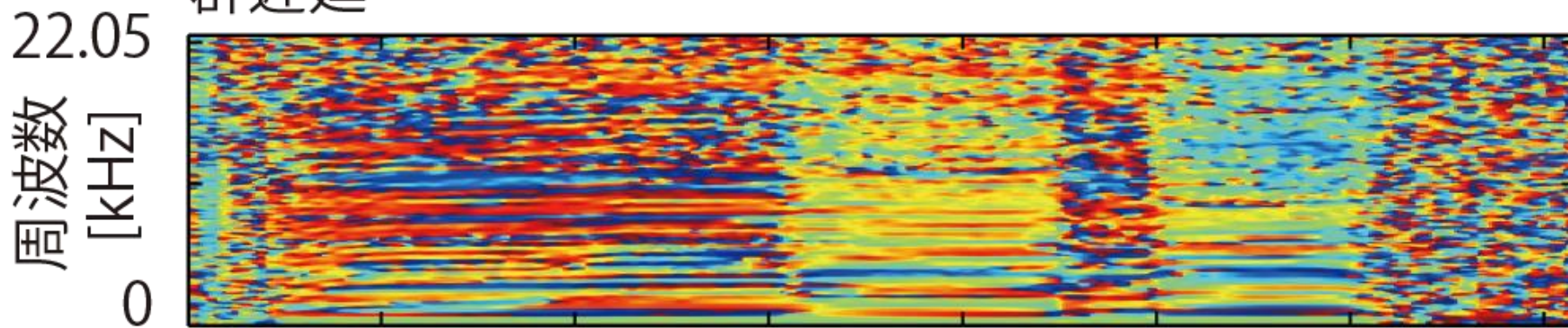


歌声の例(分析結果)

スペクトル包絡

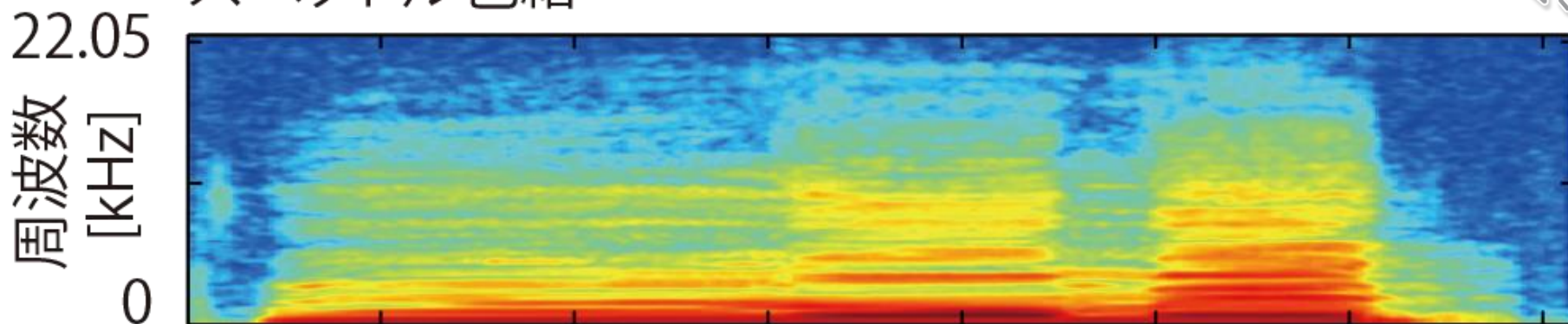


群遅延

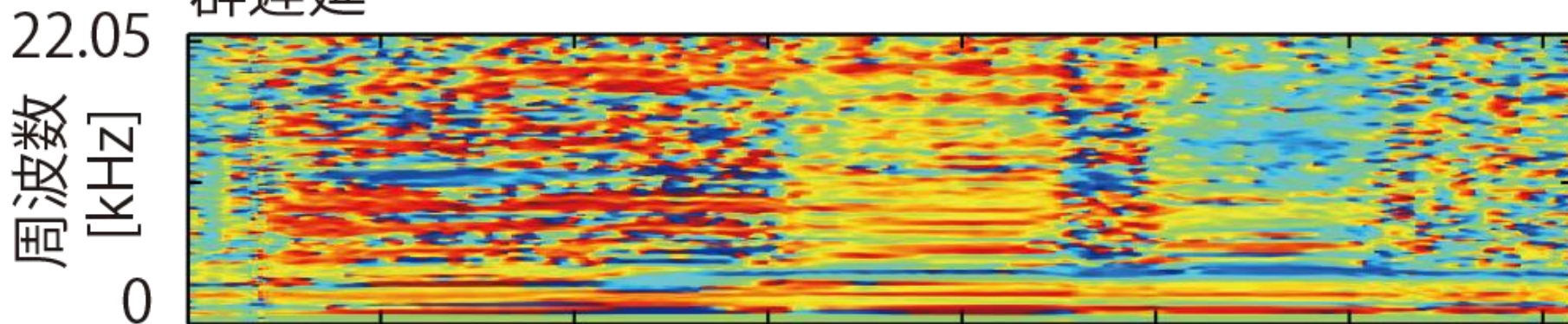


歌声の例（再合成結果）：音高・時間軸も変更可能

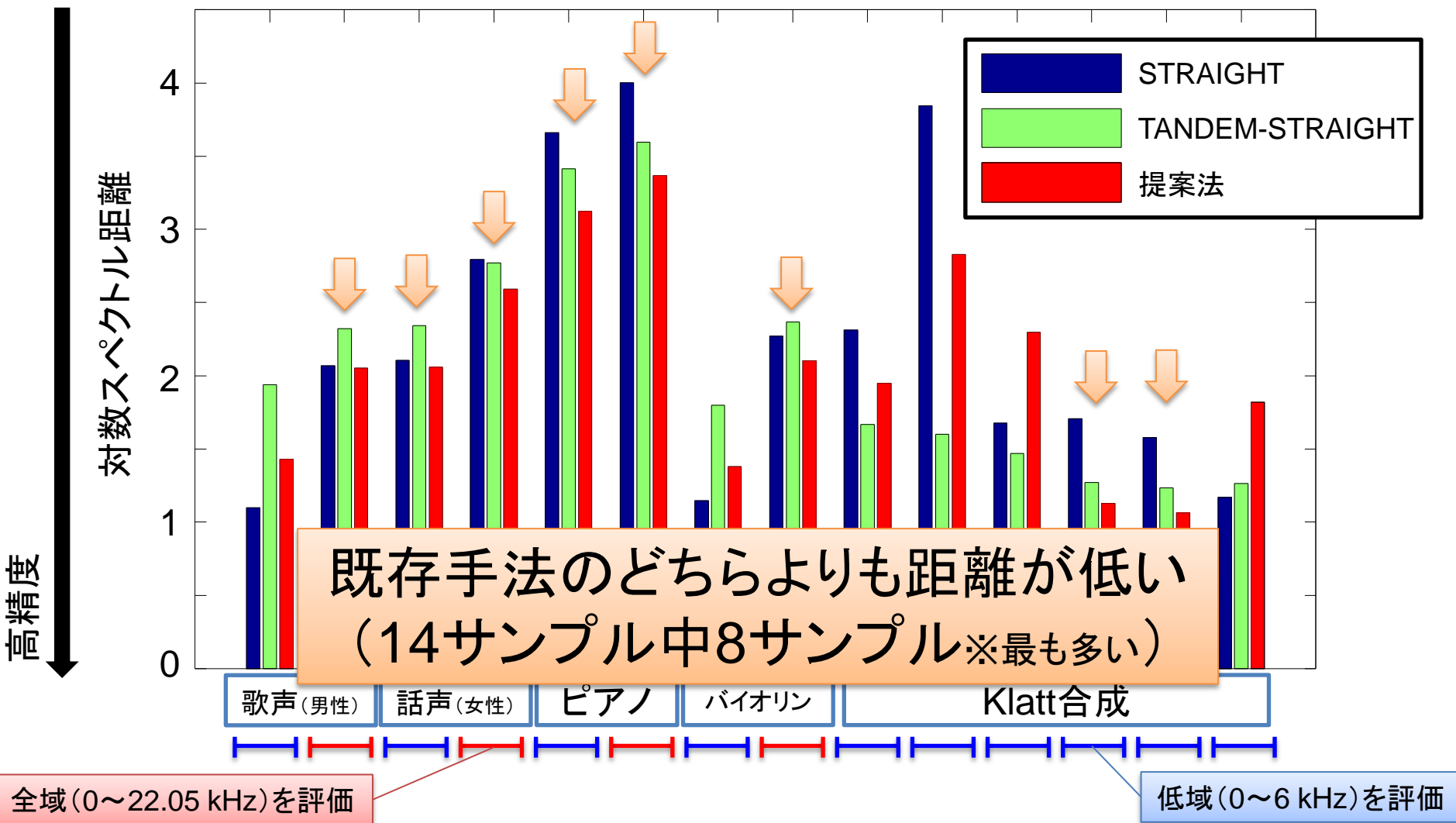
スペクトル包絡



群遅延



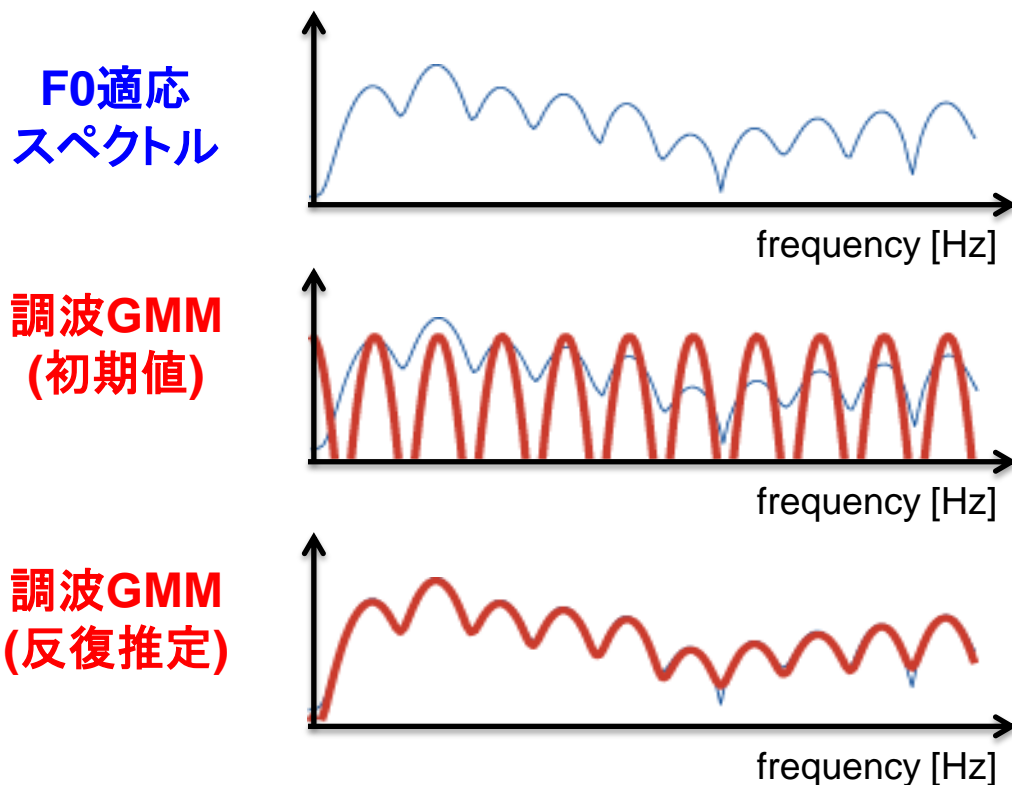
スペクトル包絡の推定精度



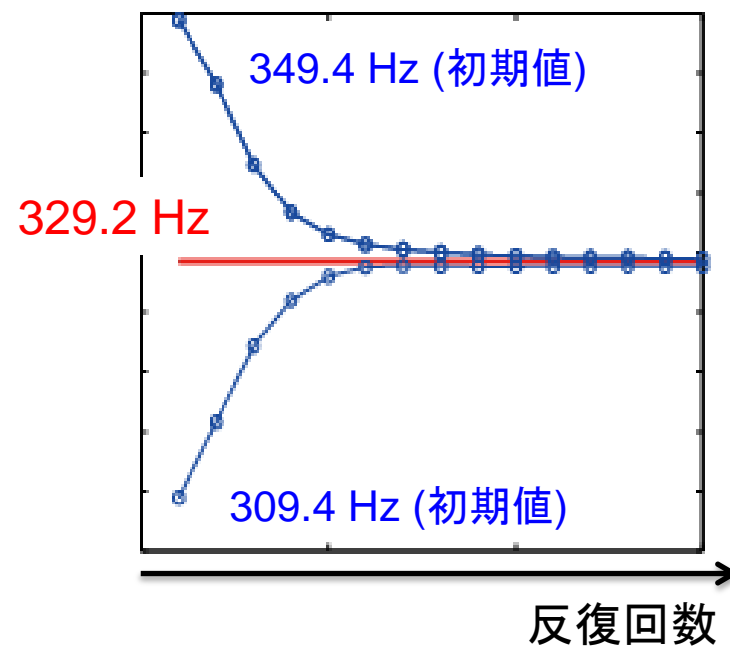
歌声の音高(F0)推定手法

[Nakano and Goto, 2013]

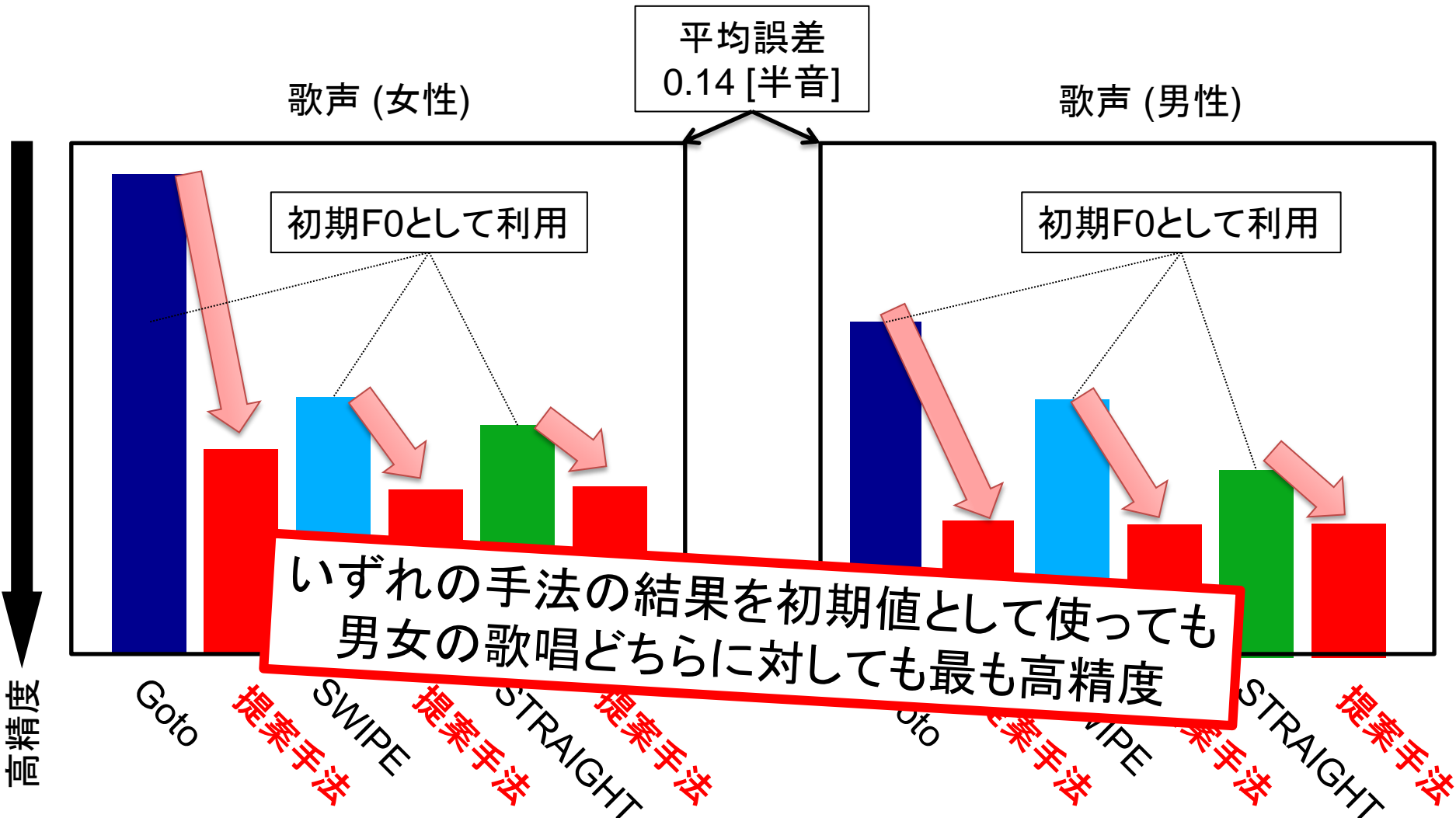
- F0情報からF0適応スペクトルを算出し、調波GMMを用いて**高い精度と時間分解能でF0を再推定**する



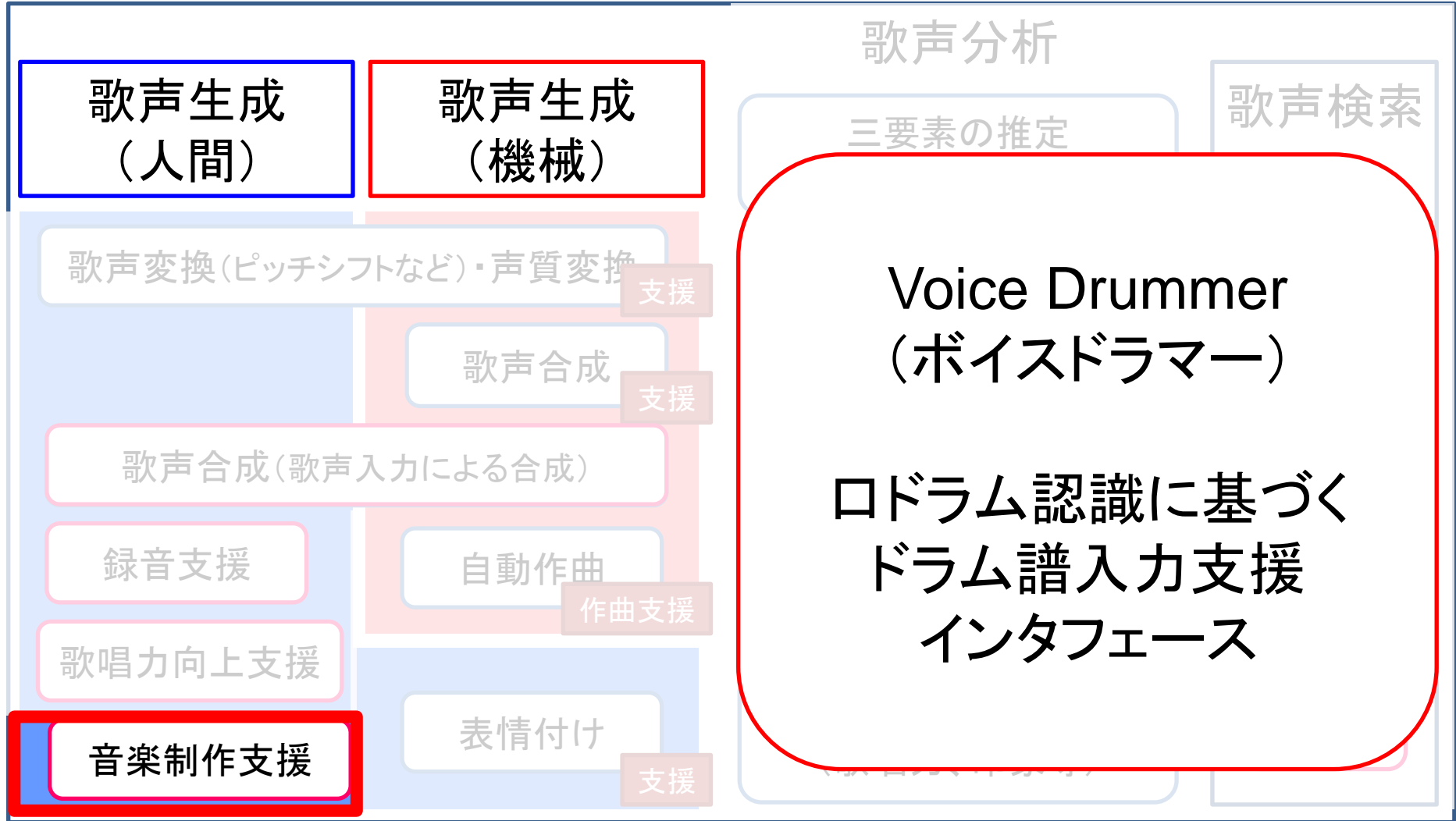
初期値の誤りにもロバスト



既存手法よりも推定誤差を低くできる



歌声インタフェースの全体像



歌声生成
(人間)

歌声生成
(機械)

歌声分析

歌声検索

三要素の推定

歌声変換(ピッチシフトなど)・声質変換 支援

歌声合成 支援

歌声合成(歌声入力による合成)

録音支援

自動作曲 作曲支援

歌唱力向上支援

表情付け 支援

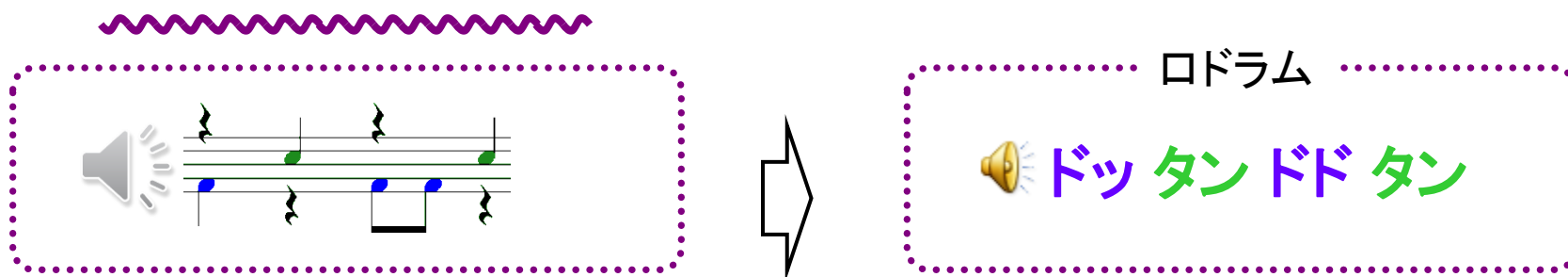
音楽制作支援

Voice Drummer
(ボイスドラマー)

ロドラム認識に基づく
ドラム譜入力支援
インタフェース

ロドラムとは

□ ドラム音（楽器音）の擬音語表現（リズムカルな歌声）



□ ロドラム認識の重要性

- 楽譜入力の新しい方式の実現
 - ・ ロドラムによる直感的で手軽な楽譜入力
- 音楽情報検索へ応用可能
 - ・ 「ラララー」等のハミングでの検索はこれまでであったが、ドラムパートを検索する目的では使えなかった

Voice Drummer

[Nakano, Goto,
Ogata, Hiraga, 2005-]

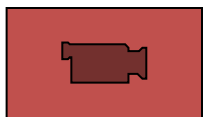
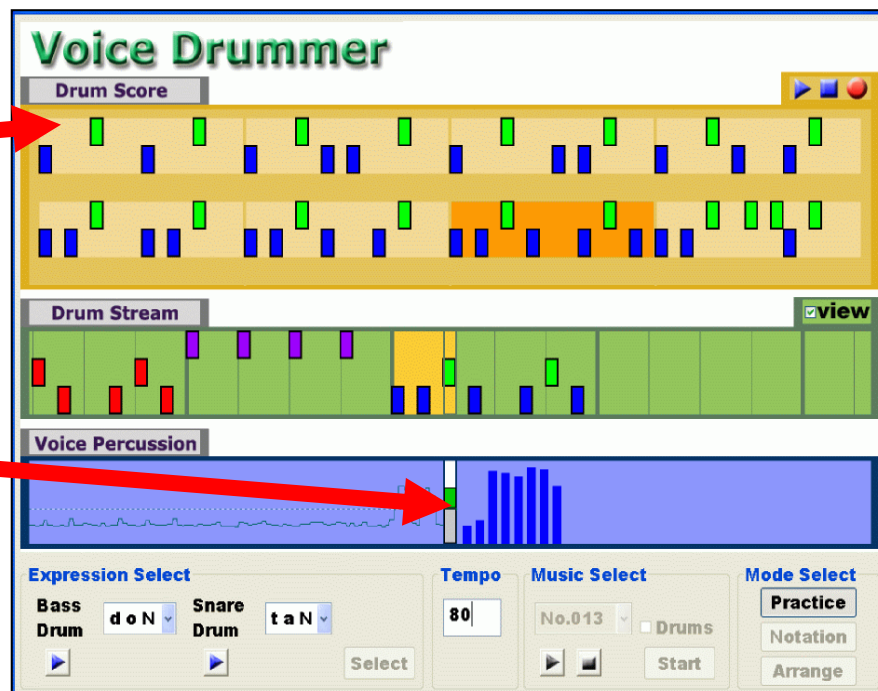
□ □(くち)ドラムによるドラム譜入力システム

- ロドラム(ボイスパーカッション)によってドラムパターンを検索
 - ・ ドラム音を真似た「ドンタンドタン」のような発声
- 練習することで認識モデルを歌唱者へ適応
- 既存の楽曲のドラムパートだけを差し替えて編曲

インタラクション

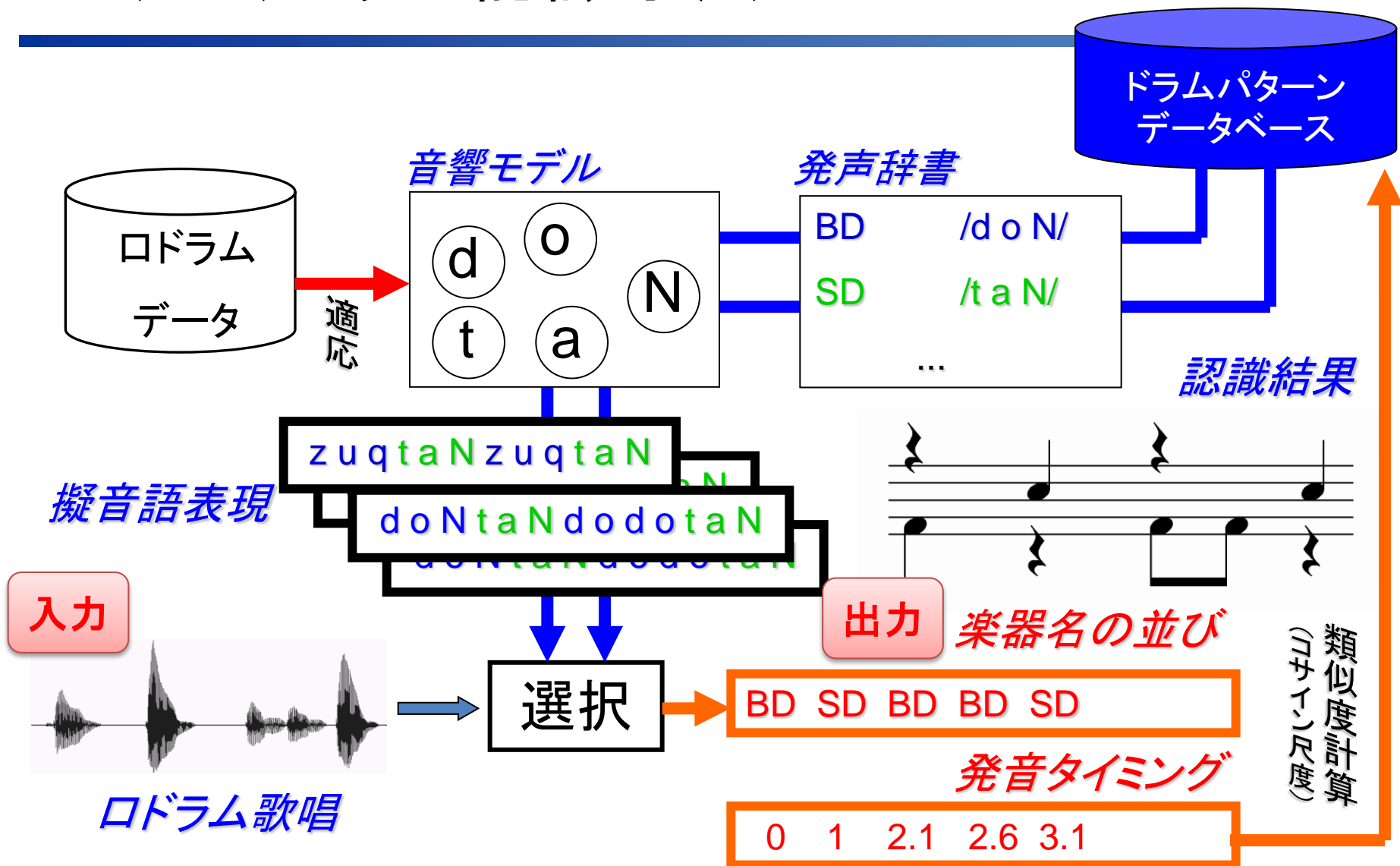
認識結果

ロドラム入力
「ドンタンドタン」



口(くち)ドラム認識手法

[Nakano, Goto, Ogata, Hiraga, 2004-]

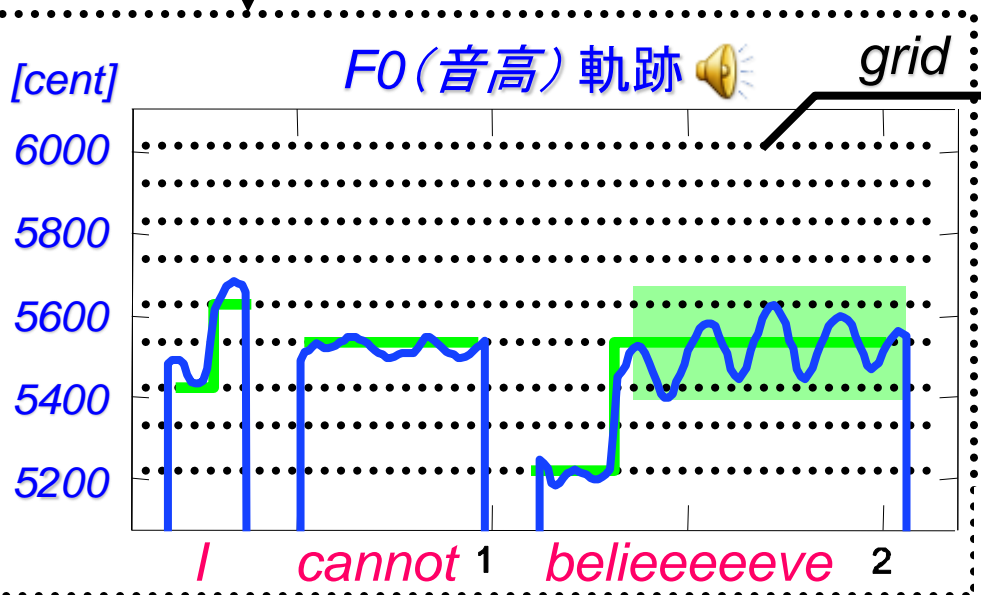
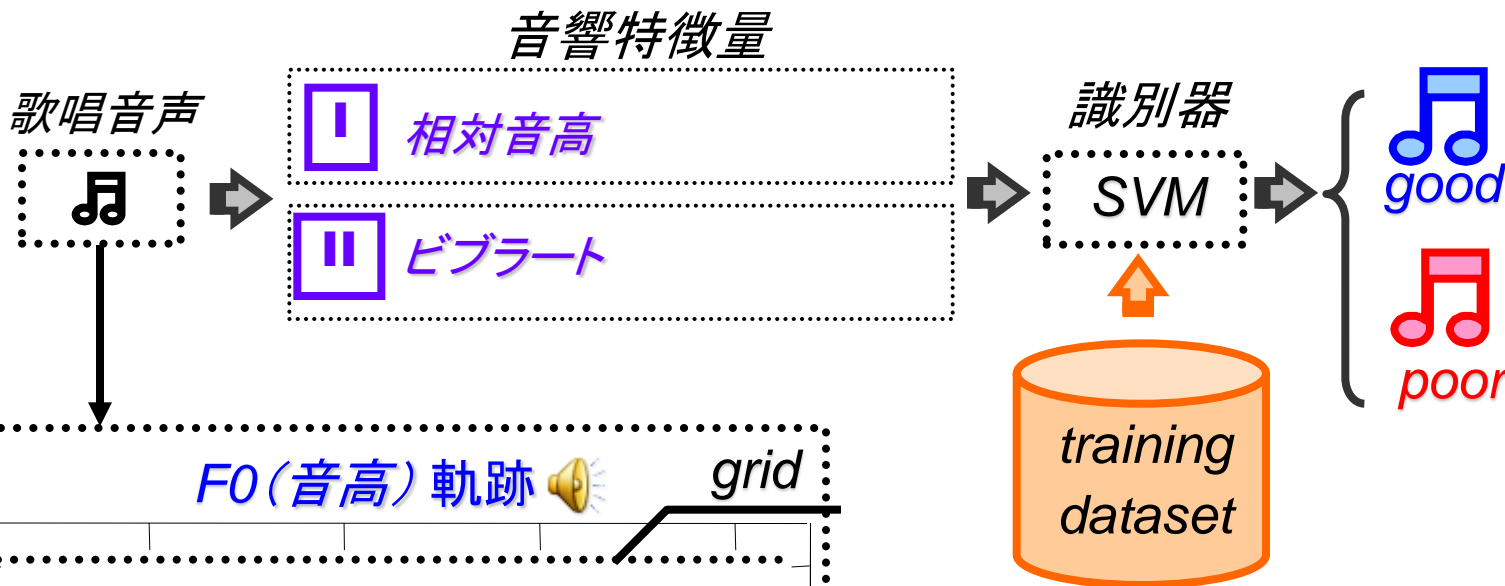


歌声インタフェースの全体像



歌唱力自動評価（楽譜なし）

[Nakano, Goto, Hiraga, 2006-]



相対音高
絶対的な音高 (F0) の正しさではなく、
相対的な音高変化に着目 (半音: 100cent)

※ centは対数周波数
1200 cent = 1 octave

ビブラート
音高 (F0) を周期的、意図的に揺らすテクニック
この区間を自動的に推定する

$$\text{Grid Frequency}_n = n \times 100 + F \quad (0 \leq F < 100)$$

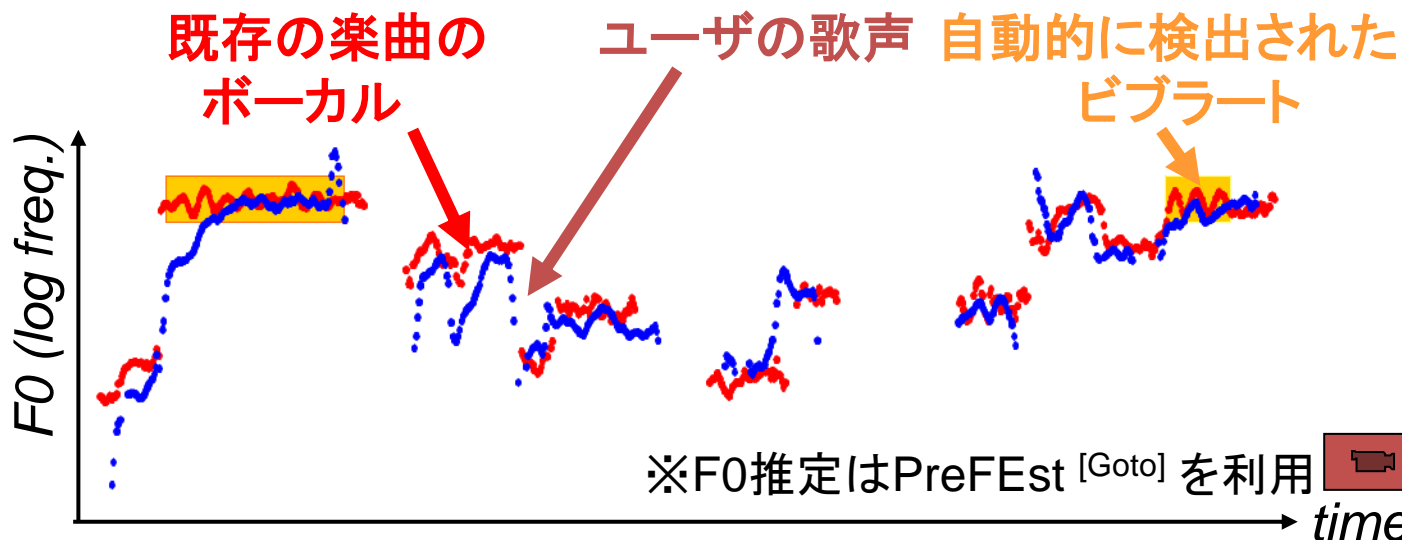
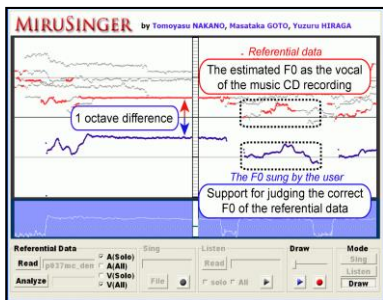
MiruSinger

[Nakano, Goto,
Hiraga, 2007-]

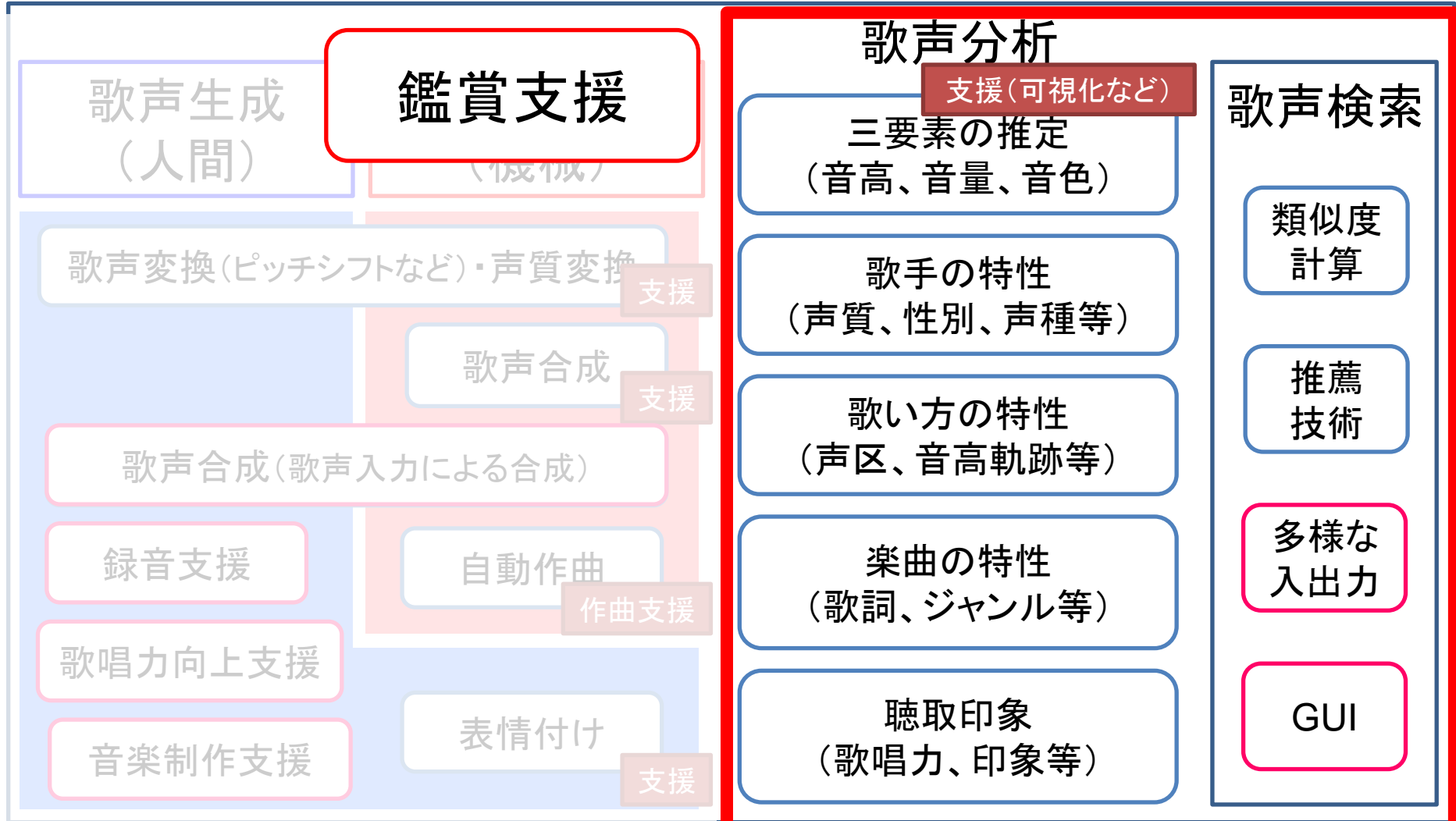
□ 歌唱力向上支援インタフェース

インタラクション

- 既存曲のボーカルの歌い方に忠実に歌いたい！
- 混合音中のボーカルを分析して可視化して修正可能
- それに合わせてユーザの歌声も比較表示
- リアルタイムに音高(F0)が可視化され、ビブラート区間も表示



歌声インタフェースの全体像



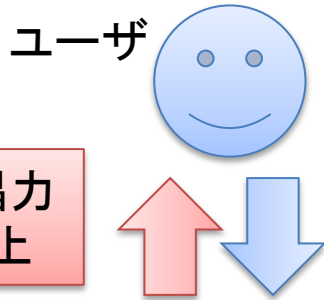
まとめ

- 「歌声インタフェース」と名付けた研究アプローチ
 - 歌声信号処理に基づくインタフェース構築やインタラクションによって人々の音楽生活をより豊かに

- 技術が世の中で広く活用されるために
 - 対象ユーザーの特性に合わせたインタフェース構築
 - ユーザ視点での問題発見（インタラクションデザイン）
 - VocaRefiner: 1度で完璧に歌えないユーザー
 - VocaListener, Voice Drummer: 楽譜に不慣れなユーザー
 - MiruSinger: 自分の好きな歌手の歌で練習したいユーザー

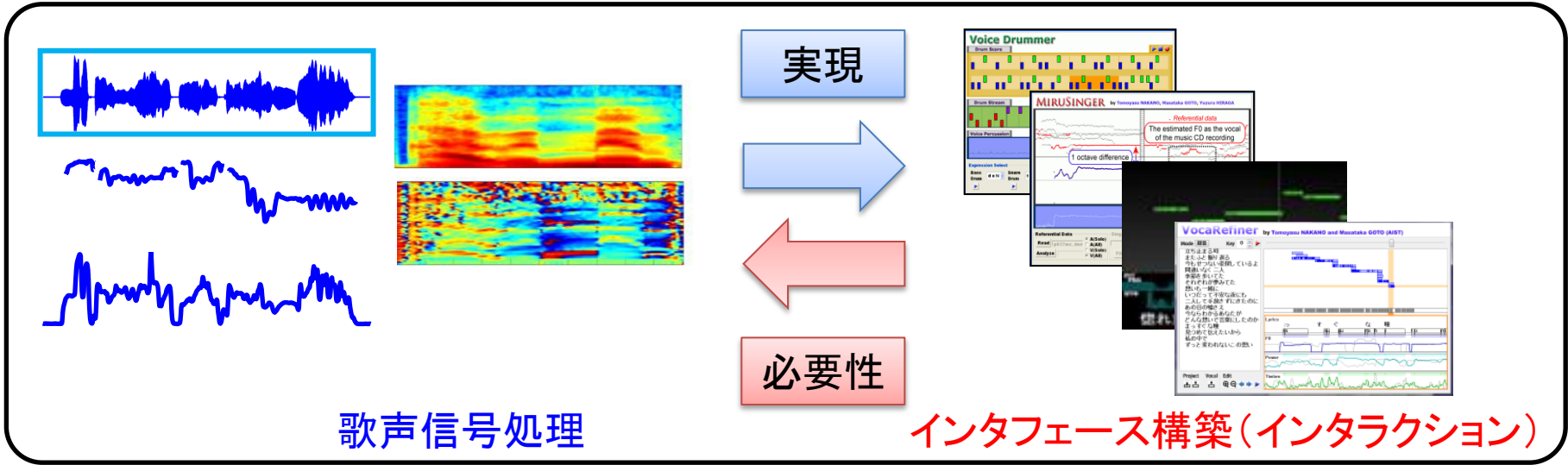
歌声インタフェースの未来

□ 人間と技術・インタラクションが
相互に成長できる未来



ユーザ自身の歌唱力
や表現力等が向上

必要な機能の実現: 歌声インタフェース自体
の機能を豊かで高度に



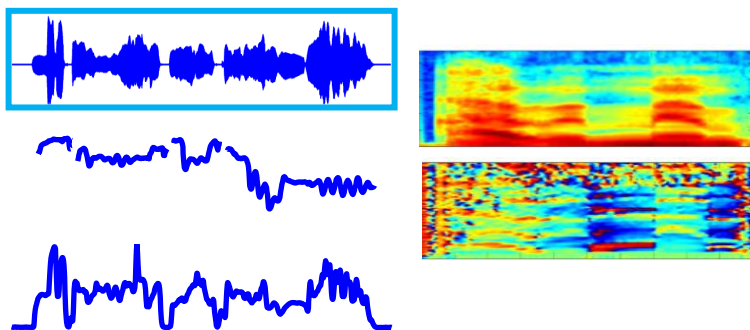
歌声インタフェースの未来

新たな歌唱表現を追求したり
音楽に関する理解が深まったり
できるようになる

単に技術のみが
高度化していく
のではない未来
が切り拓ける

ユーザ自身の歌唱力
や表現力等が向上

必要な機能の実現: 歌声インタフェース自体
の機能を豊かで高度に

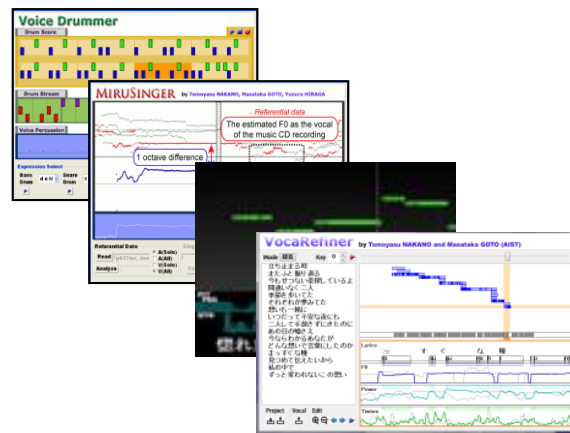


歌声信号処理

実現



必要性



インタフェース構築(インタラクション)

スペシャルセッション 音声B/音声A/聴覚

[ここまで来た声質変換技術 – 実用可能性の視点からの現状認識と将来展望 –]

歌声インタフェース: 歌声を対象とした信号処理と それに基づくインタフェース構築

中野 倫靖, 後藤 真孝
(産業技術総合研究所)

2013年9月27日

日本音響学会 2013年秋期研究発表会(講演番号3-7-3)