

$$-\frac{\partial \text{KL}(\mathbf{y}_{(\text{ICAL})}(t)})}{\partial \mathbf{w}_{(\text{ICAL})}(n)} \cdot \mathbf{W}_{(\text{ICAL})}(z^{-1})^T \mathbf{W}_{(\text{ICAL})}(z)$$

システム情報 猿渡・齋藤研究室 (創造情報 猿渡研究室)の紹介

$$f(x) = x^{k-1} \frac{e^{-x/\theta}}{\Gamma(k)\theta^k}$$

東京大学大学院・情報理工学系研究科

システム情報学・創造情報学専攻 猿渡・齋藤研

(2026年5月)

猿渡・齋藤研(システム情報第一研究室)

教授
猿渡洋



専門分野

- ・教師無し最適化
- ・統計・機械学習
- ・論的信号処理

講師
齋藤佑樹



専門分野

- ・統計的機械学習
- ・音声知覚モデリング
- ・Human Computation

助教
山岡洸瑛



専門分野

- ・多チャンネル信号処理
- ・方位時間差推定
- ・補助関数最適化

特任助教
岡本悠希



専門分野

- ・環境音合成
- ・環境音検出認識
- ・環境音DB構築

特任准教授
高道慎之介



専門分野

- ・音声コミュニケーション拡張
- ・音声言語情報処理

協力教員

- ・北村大地先生(香川高専)
- ・中村友彦先生(産総研)
- ・伊藤信貴先生(産総研)
- ・Wang Ruiさん(特任研究員)

学術専門職員 高宗さん

秘書 丹治さん

学生

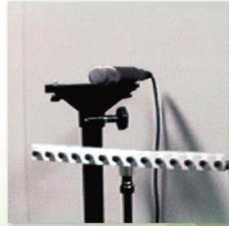
- ・博士課程学生12名
- ・修士課程学生7+9名



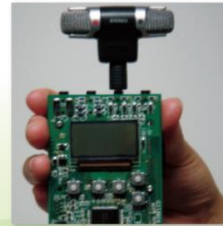
研究俯瞰図

- 音声・音響・音楽メディアに関する信号処理・情報処理
- ヒューマンインターフェイス・コミュニケーションシステムの構築
- 統計的・機械学習論的信号処理、数理最適化問題等を研究

多チャンネル信号処理



- ・教師あり音源分離
- ・統計的信号強調
- ・ロボット聴覚システム



教師なし学習に基づく
ブラインド音源分離

- ・独立線形因子分析
- ・音コミュニケーション拡張



あらゆる声を実現できる
音声合成変換

- ・音声のための深層学習
- ・音声信号処理
- ・人間参加型機械学習

統計信号
処理
音場解析・
合成
音声情報
処理
音楽信号
処理



音空間の解析と合成

- ・音場計測における逆問題
- ・音空間制御
- ・VR/ARのための空間音響



- ・楽音分離・加工
- ・ウェブレット解析
- ・モノラル音源分離

信号处理的深層学習に
基づく楽音分離



マルチモーダル
ヒューマン
インターフェイス

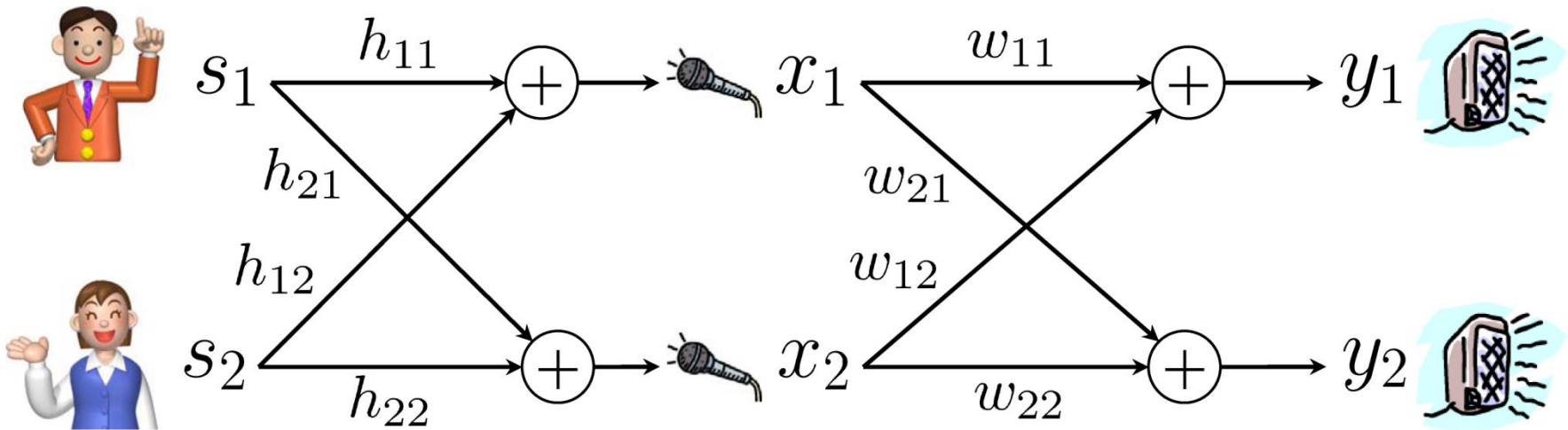
研究紹介1. ブラインド音源分離

- 音の方向・声質・音量など、事前に何も分かっていなくても、瞬時に音を「聞き分ける」ことの出来るシステムを目指す。
- 独自に開発した高速独立成分分析(ICA)、独立低ランク行列分析(ILRMA)という教師無し数理最適化アルゴリズムに基づいて、音を統計的に独立な成分に分解することにより、別々の音声信号を見つける。

スモールデータ
教師無し最適化
低ランクモデル

ブラインド音源分離 (Blind Source Separation): 聞き分けるAI

- 混ぜり合った信号 x_1, x_2 から元の信号を取り出す
- どのように混ぜたかに関する情報 H は利用できない
- 事前トレーニング出来ない \Rightarrow ビッグデータではなく **スモールデータ**



実は上記は**2つのことを同時に推定**している

- [空間] 統計的に独立な音源の分類問題 (分離行列 W の推定)
- [音源] 各音源が属する確率分布 $p(y)$ や構造の推定問題

上記を閉形式で解く方法は存在せず凸問題でもない \Rightarrow **大変困難!**

ILRMA: 音源の独立性と低ランク性に着目したBSS

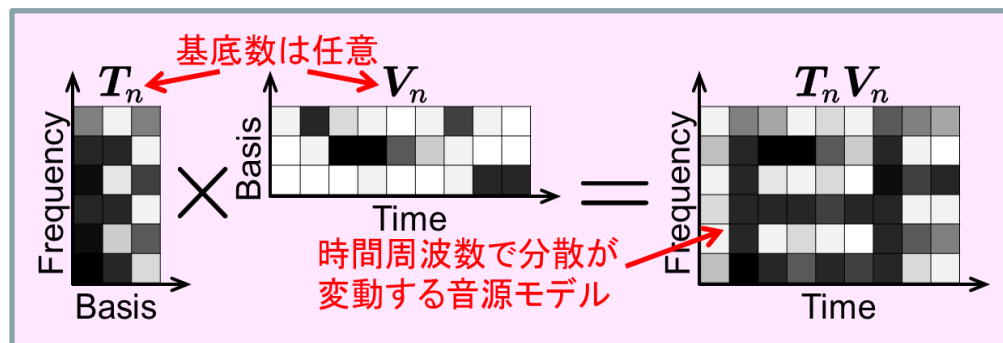
[IEEE Trans. ASLP 2016、IEEE SPS論文賞・ASJ粟屋賞・JSPS育志賞]

- ILRMAのコスト(対数尤度)関数→これを最小化

$$\mathcal{J} = \sum_{i,j} \left[\sum_m \log \sum_k z_{mk} t_{ik} v_{kj} + \sum_m \frac{|y_{ij,m}|^2}{\sum_k z_{mk} t_{ik} v_{kj}} - 2 \log |\det \mathbf{W}_i| \right]$$

音源の低ランク性コスト関数
(音源NMFモデルの推定に寄与)

音源の独立性コスト関数
(空間モデル W の推定に寄与)



$$p(\mathbf{y}) = p(y_1) p(y_2) \dots p(y_m)$$

となる W を推定

両者を交互にMajorization-Minimization(補助関数法)アルゴリズムで反復最小化

- ✓ コスト値の単調減少性を保証(勾配法には無い特徴)
- ✓ 高速かつ安定な求解法を実現(従来の多入力NMFと比較して2ケタ速い)

モデルの多様化・数理解法の開拓

● 音源生成モデルの多様化 IEEE SPS Tokyo Joint Chapter 学生賞!

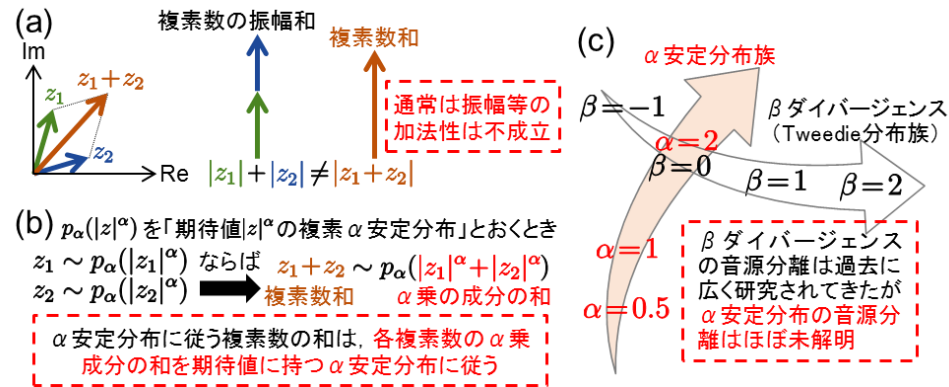
— 複素波形重ね合わせと整合する α 安定分布の導入

⇒ t-ILRMA [EURASIP-JASP2018]

— 複素球状ポアソン分布の導入

⇒ $\beta=1$ -divergence 最小化 ILRMA

[IEICE Trans. 2018]



● 座標降下法におけるバリエーション

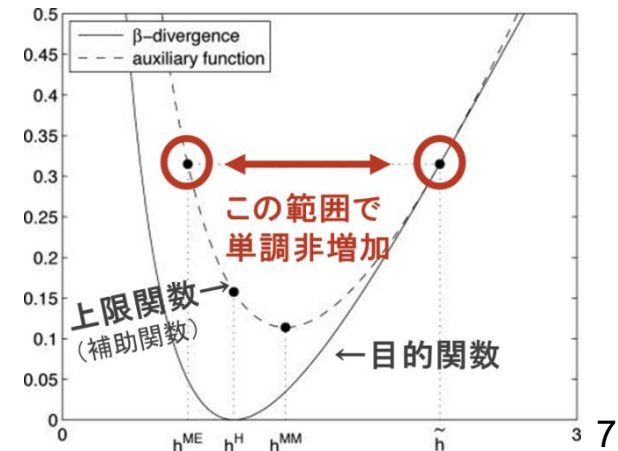
— パラメトリック Majorization-Equalization アルゴリズム による音源・空間最適化の「バランス」化

[CAMSAP2017], [IEEE Trans. ASLP2019]

● 深層学習 (DNN) との融合

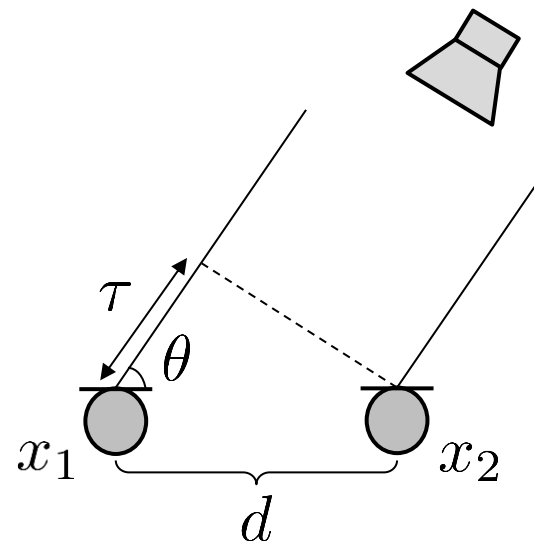
— 独立深層学習行列分析 (IDLMA) の提案

[IEEE Trans. ASLP2019]



補助関数型時間差推定 [Yamaoka+, 2019]

- 音の到来時間差
 - 音の到来方向に対応 ($\tau = d \cos(\theta) / c$)
 - ブラインドに推定する必要あり
 - 音源分離とは相補的な関係

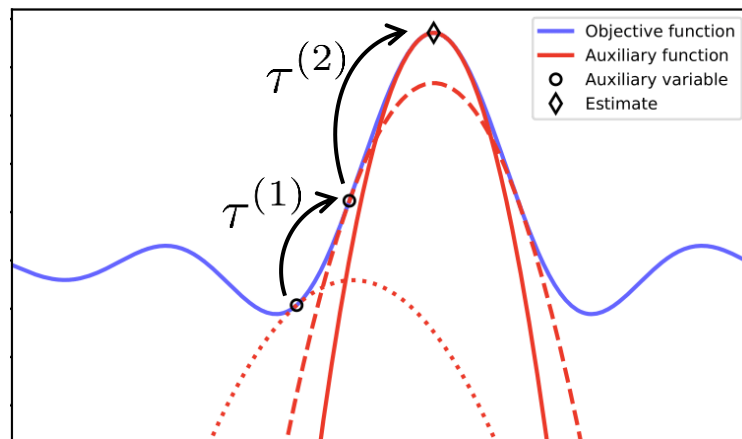


- 補助関数法に基づく推定
 - 目的関数は非凸で閉形式の解はない

$$\arg \max_{\tau} \sum_k x_{1k} x_{2k}^* e^{-j2\pi \frac{k}{K} \tau}$$

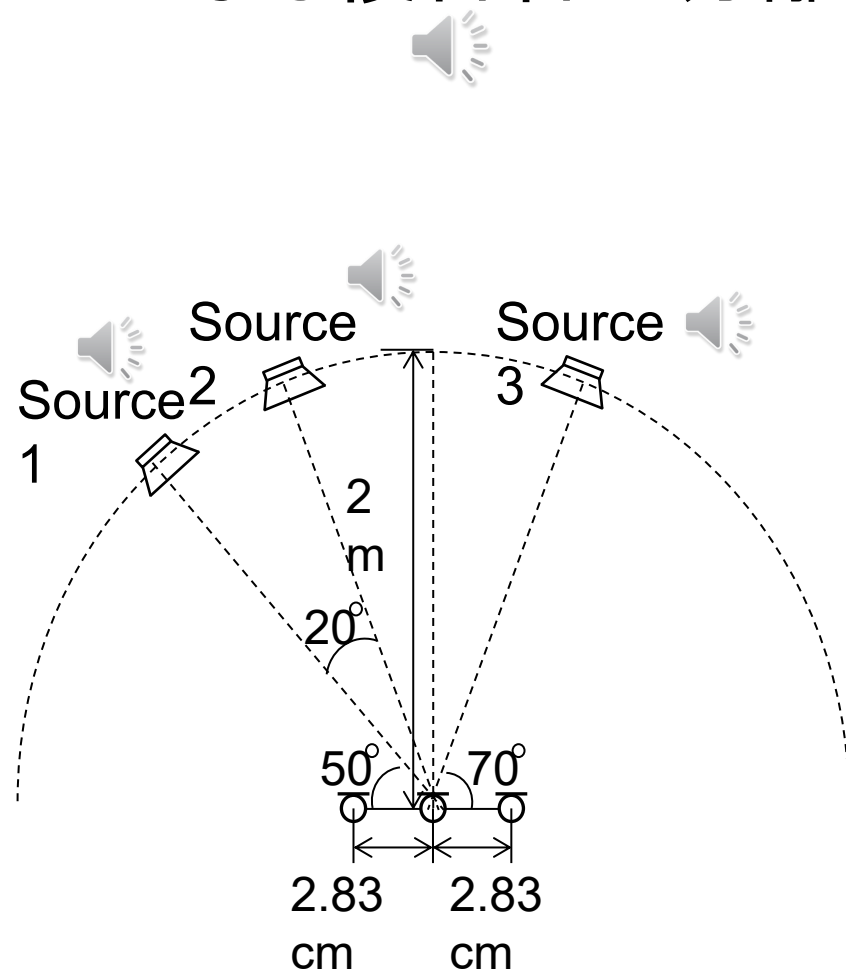
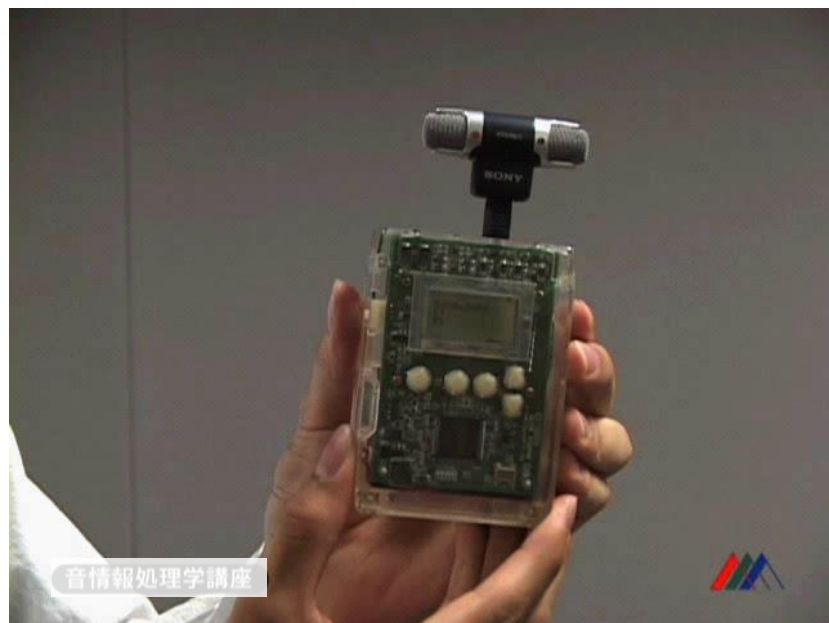
- 二次の補助関数を設計,
閉形式の解による更新則を提案

$$\tau^{(\ell+1)} \leftarrow \tau^{(\ell)} - \frac{\sum_k A_k \operatorname{sinc} \theta_k^{(\ell)} \frac{\theta_k^{(\ell)}}{\omega_k}}{\sum_k A_k \operatorname{sinc} \theta_k^{(\ell)}}$$



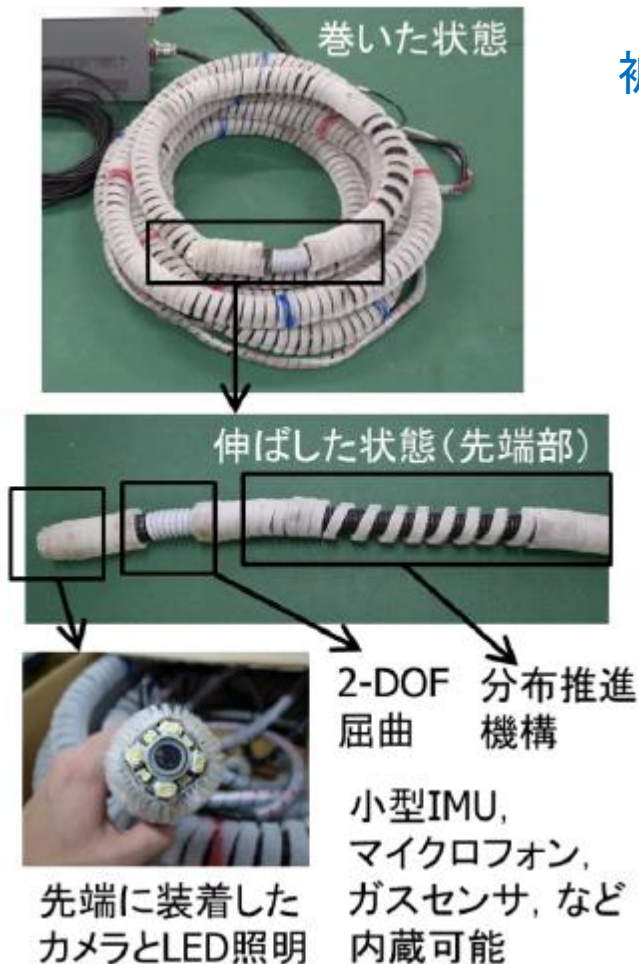
高速ICA、独立低ランク行列分析によるデモ

- リアルタイム音声聞き分け(警察備品に採用)
- ドラム、弦楽器、音声からなる複合音の分離

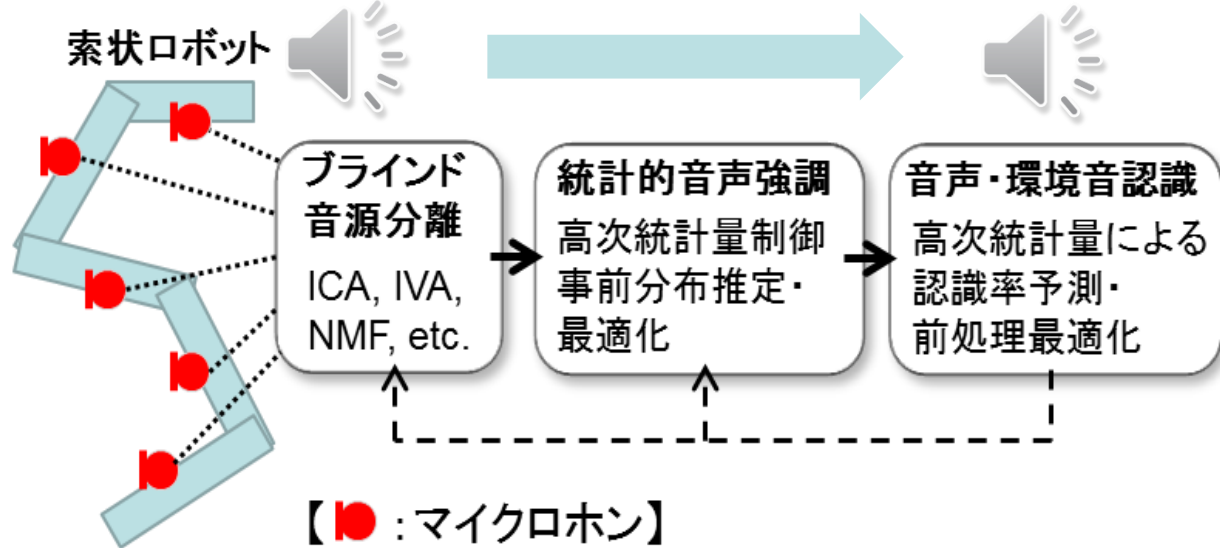


内閣府ImPACT災害対応タフロボット [2016年6月プレスリリース]

- 災害時の倒壊家屋に入り込んで被災者発見
- 環境音認識による状況把握・救助支援



被災者はいらぬのか？



いかなる曲がりくねった形状においても
マイク同士が協調して騒音の中から被災者の声を見つけ出す

オンライン(リアルタイム)ILRMA

- 極端にうるさい環境で使用される音声対話ロボット向けのオンライン(リアルタイム)ILRMA [Ishikawa, Saruwatari+, IEEE Access 2025]
- 全てのパラメータは32 ms以内に最適化される
- 2025年大阪万博のアバター展示会場にて使用



Empathetic listening
dialogue robot
(CommU)

Client (human)



4-ch circular microphone array

【Demo 1】

• Observed



• BSS



【Demo 2】

• Observed



• BSS



独立深層学習行列分析

Independent Deeply Learned Matrix Analysis

(IDLMA: 発音はアイドルエムエー)

[Makishima, Saruwatari+, IEEE-Trans. 2019]

ILRMAにおける問題点：音源の低ランク性？

$$\mathcal{J}_{\text{ILRMA}} = \frac{1}{J} \sum_{i,j,n} \left[\log \sum_l t_{il,n} v_{lj,n} + \frac{|w_{i,n}^H x_{ij}|^2}{\sum_l t_{il,n} v_{lj,n}} \right] - \sum_i \log |\det W_i|^2$$

音源モデル (低ランク性)

空間モデル (音源間が独立)

音源によっては低ランク性が
成り立たない場合がある

音源・マイク位置，部屋の形状，
残響時間などの膨大な物理要因に依存

ならば！

事前に学習データを用いて音源モデルの分散を推定する写像を作る

学習データの用意は非現実的
ブラインドに推定

ILRMAにおける問題点：音源の低ランク性？

$$\mathcal{J}_{\text{ILRMA}} = \frac{1}{J} \sum_{i,j,n} \left[\log \sum_l t_{il,n} v_{lj,n} + \frac{|w_{i,n}^H x_{ij}|^2}{\sum_l t_{il,n} v_{lj,n}} \right] - \sum_i \log |\det \mathbf{W}_i|^2$$

音源モデル (低ランク性)

空間モデル (音源間が独立)

- 深層学習 (DNN) による強力なモデリング能力を活用する！
- 今まで培ってきた「教師あり音源分離 (例：教師ありNMF)」の技術を昇華させる形で研究を発展できる。
- 急速に発展するDNN研究を我々ならではの視点で拡張する。

事前に学習データを用いて音源モデルの分散を推定する写像を作る

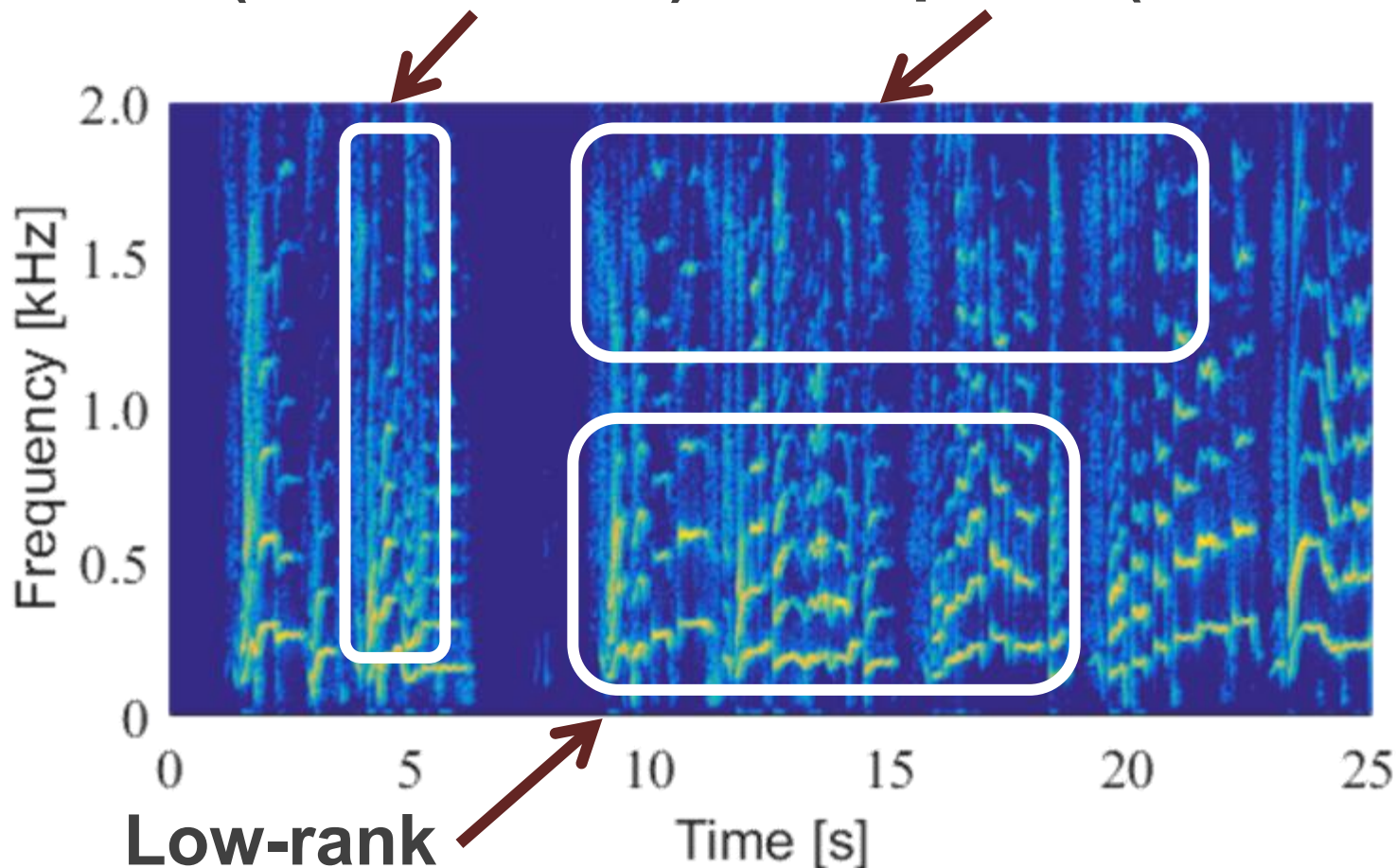
学習データの用意は非現実的
ブラインドに推定

音源の低ランク性？ (例：音声信号)



Dense (not low-rank)

Sparse (not low-rank)



このような複雑な構造を持つ信号はDNNでモデリングする

提案手法：DNN音源モデルによる最尤推定

■ 独立深層学習行列分析 (IDLMA)

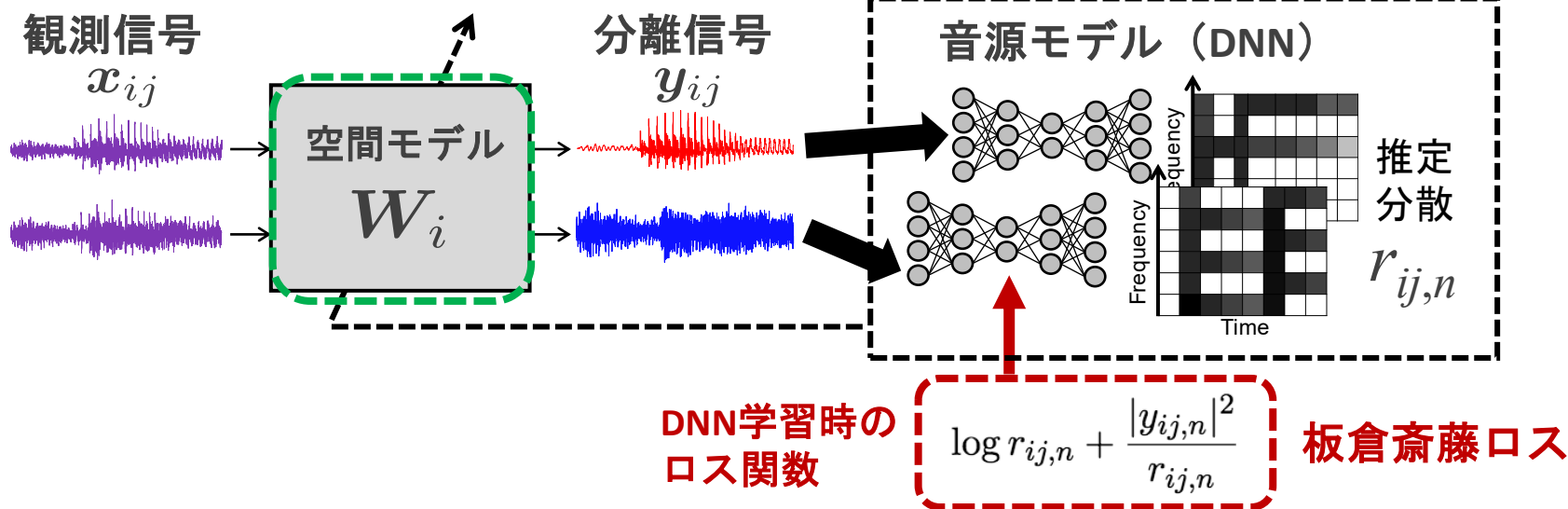
$$y_{ij,n} = \mathbf{w}_{i,n}^H \mathbf{x}_{ij}$$

$$\mathcal{J} = \frac{1}{J} \sum_{i,j,n} \left(\log r_{ij,n} + \frac{|\mathbf{w}_{i,n}^H \mathbf{x}_{ij}|^2}{r_{ij,n}} \right) - \sum_i \log |\det \mathbf{W}_i|^2$$

音源モデル (DNN)

交互に最適化

空間モデル (音源間が独立)



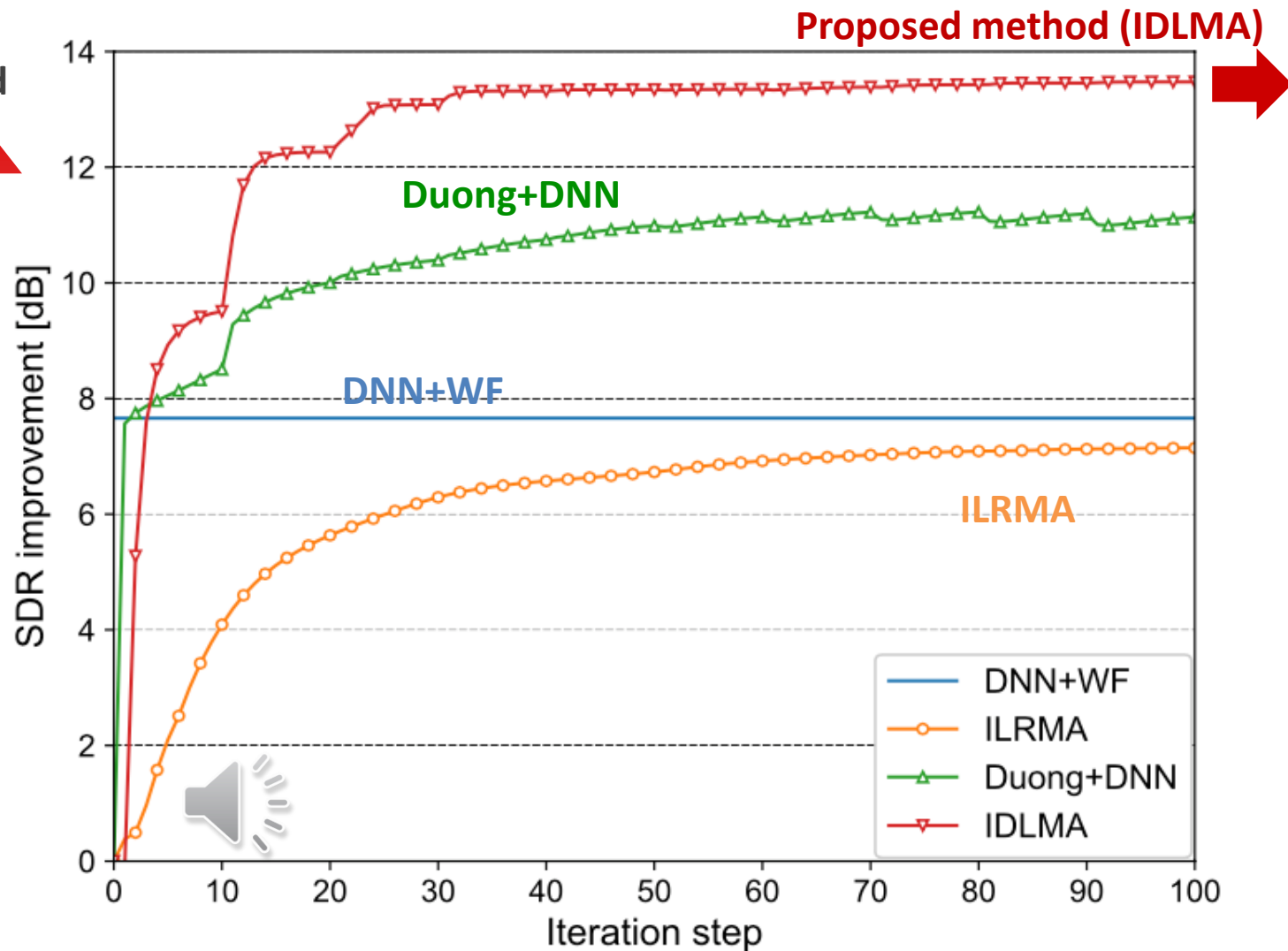
■ 空間モデル：各音源が統計的に独立となる分離行列を推定

■ 音源モデル： \mathcal{J} を最小化するような分散 $r_{ij,n}$ を推定するDNNを各音源ごとに構成

実験結果例（反復最適化）



Good



Vo.



Ba.



研究紹介2. 統計的音声合成

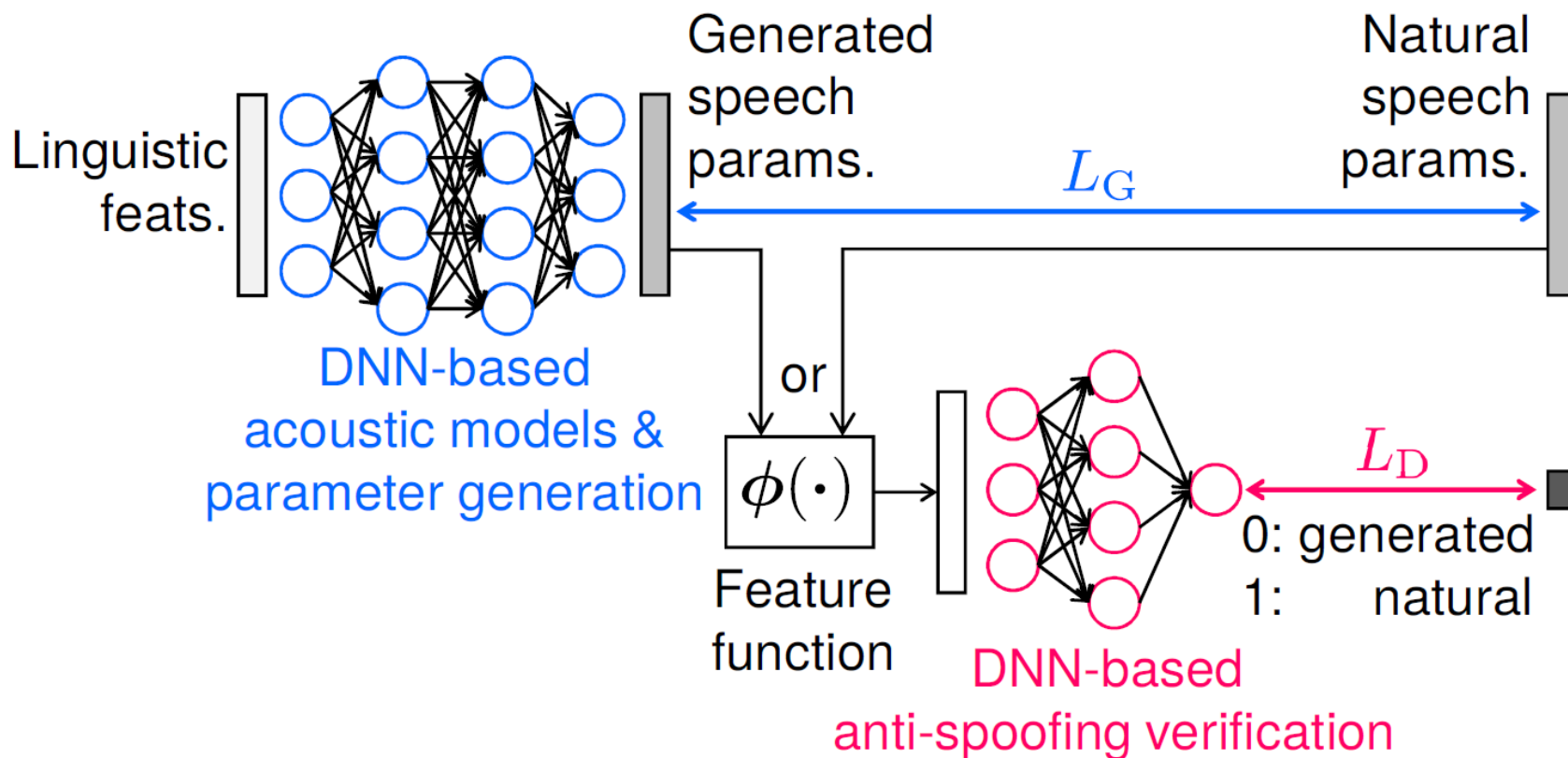
- テキストもしくは自分の声を入れるだけで、誰の声でも、どんな訛りでも、何語でも、喋ることが出来るような統計的音声合成システムを実現する。
- 深層学習 (Deep Neural Net) の枠組みを活かし、「AIオレオレ詐欺師」と「AI防犯課刑事」を対決させて、お互いに精度を高めるAnti-Spoofing 敵対学習理論を独自に提唱

ビッグデータ
深層学習
敵対学習



Anti-Spoofingと敵対する音響モデル学習理論

[IEEE SPS Young Author Best Paper賞・IEEE SPS] 学生論文賞他]



人間の声に似せようと努力

ウソ(合成音)に騙されまいと攻防



リアルタイムDNN声質変換



2019年3月 日経xTECHプレスリリース



さらに柔軟な音声合成へ



- 松任谷由実+AI荒井由実「Call me back」人工声の提供(2022年の紅白にて紹介)



出典:YouTube 松任谷由実 – Call me back/松任谷由実with 荒井由実

- フィラー挿入付き自発音声合成(言いよどむAI)

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

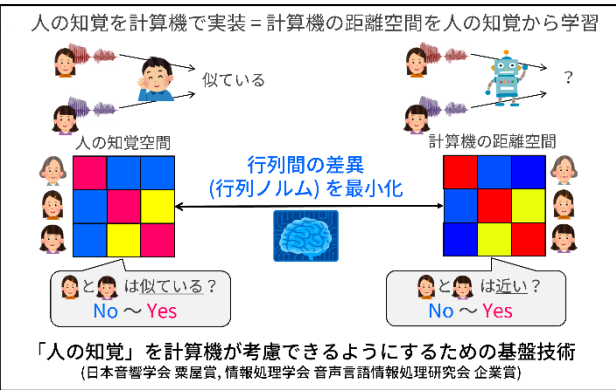


ざっくりいうと、(アノ)先ほど(アノ)
少し(アノ)お話ししましたけども、
戦後のそういうサブカルチャーの
イメージという...

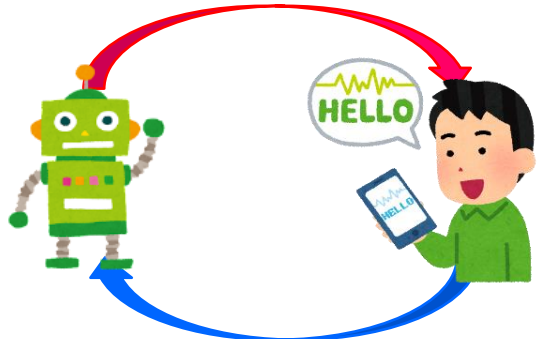


...

人と機械の相互作用を通じた音声合成・変換の最適化



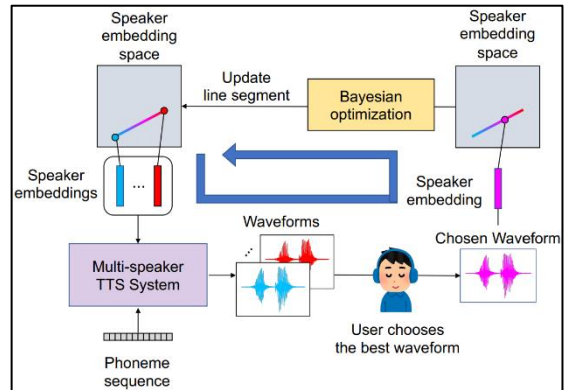
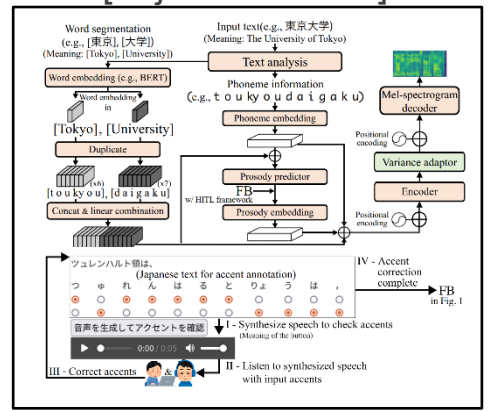
機械による
音声情報処理能力の拡張



人の知覚を導入した最適化

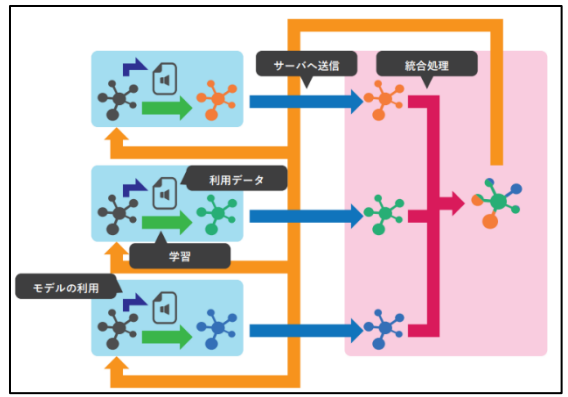
話者知覚に基づく
音声表現学習 [Saito+TASLP21]

アクセント訂正FB音声合成
[Fujii+APSIPA22]



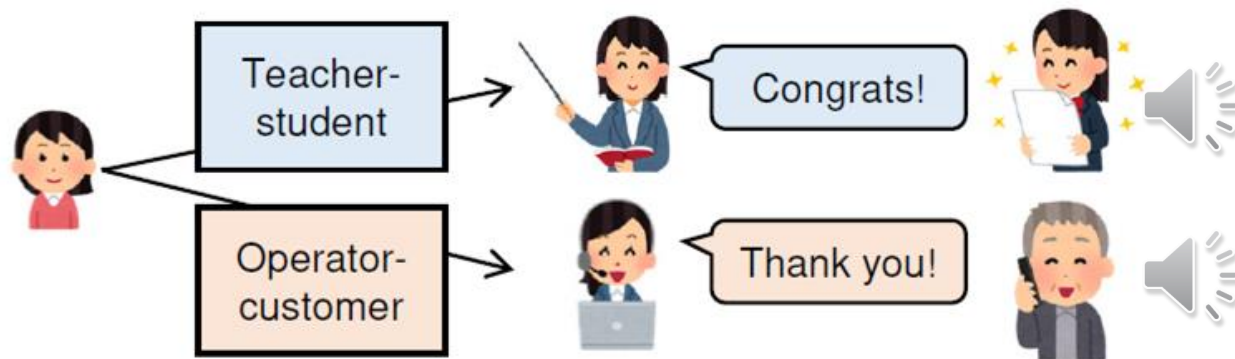
Human-in-the-loop
話者適応 [Udagawa+IS22]

Federated 声質変換
[平井+SLP23]

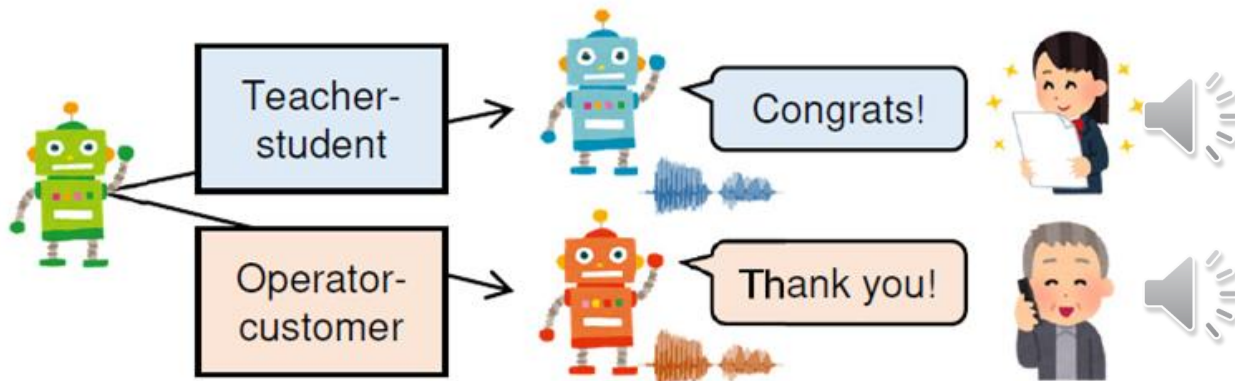


人と機械が協調して情報伝達を行うための基盤技術を構築

研究例: 多ドメイン共感的対話音声合成



人間同士のコミュニケーション:
様々な対話ドメインで相手と共感的に会話可能



人間・ロボット間のコミュニケーションでもこれを再現

ChatGPT-EDSS: 大規模言語モデル (LLM) と対話して 発話スタイルを制御する音声合成

対話相手
(人間)



Hi, teacher!

Oh, did you get
a good score?

Bingo!!

聞き手
(AI)



Speaker: Hi, teacher!

Listener: Oh, did you get
a good score?

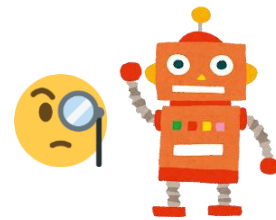
Speaker: Bingo!!

Listener: Congrats!!

対話履歴を考慮し、
相手にどう応答すべきかを回答

LLM に
「どう応答すべきか?」
を質問

対話
アドバイザー
(LLM)



喜んで、
祝うように

Congrats!!

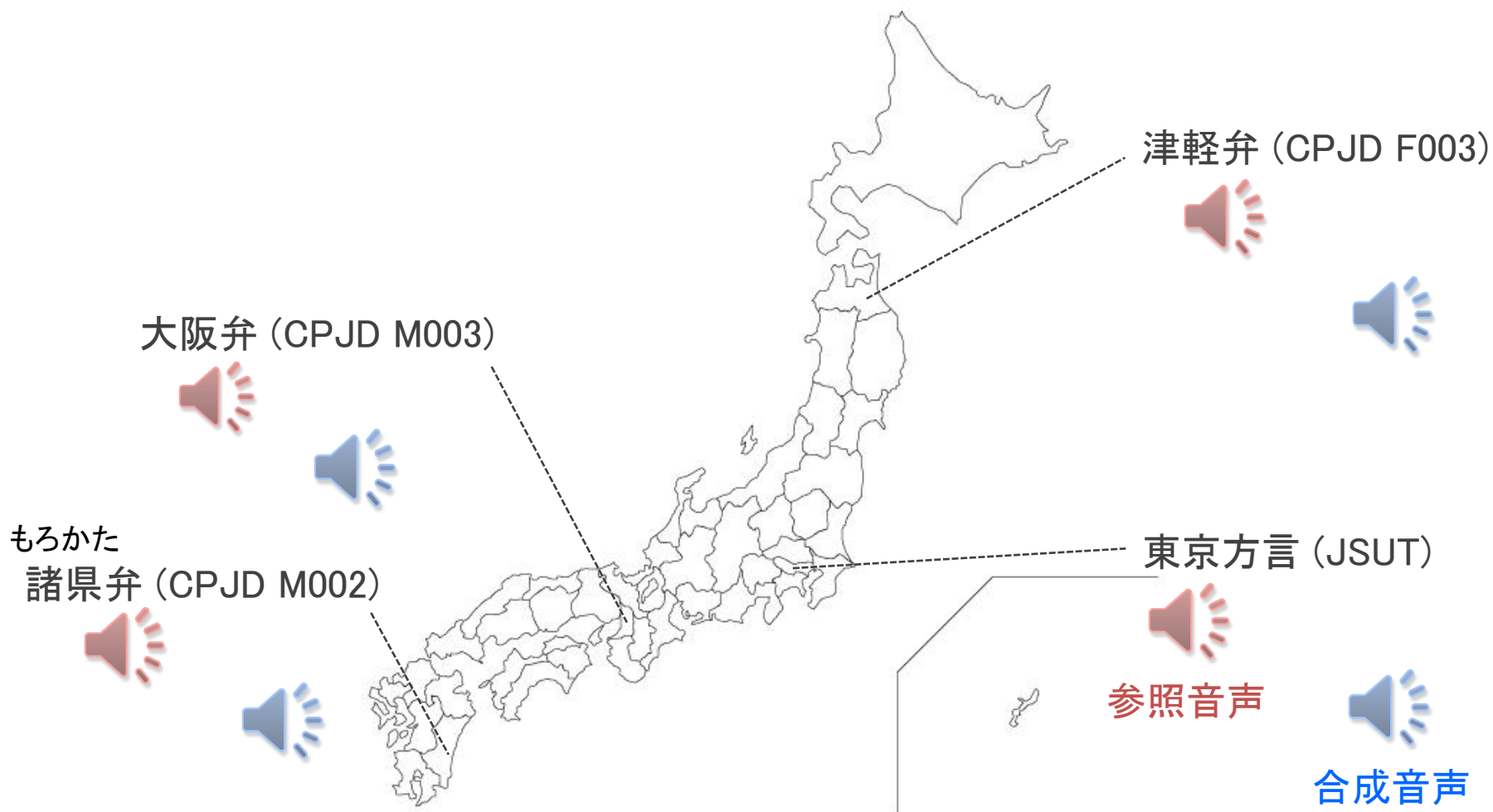


合成 w/ 感情ラベル

合成 w/ 対話履歴

合成 w/ 対話履歴 + LLM

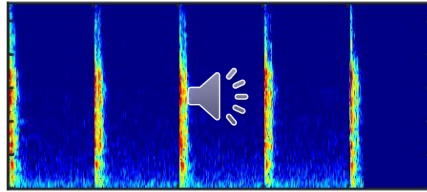
Cross-dialect TTS: 方言をまたいで喋らせる音声合成



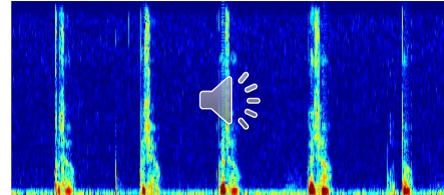
Voice-to-Foley: 声真似に基づく環境音合成



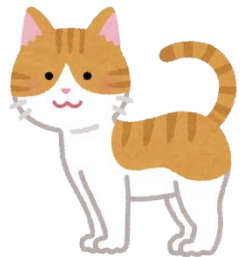
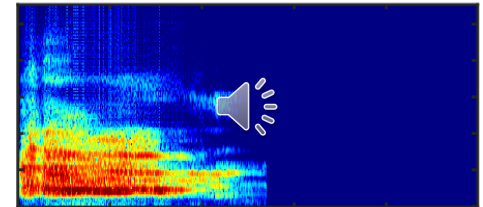
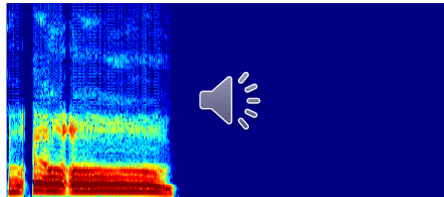
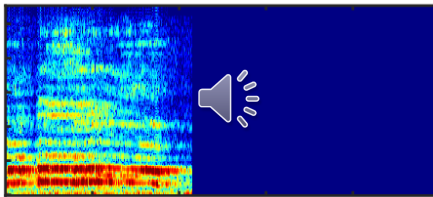
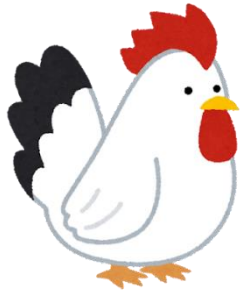
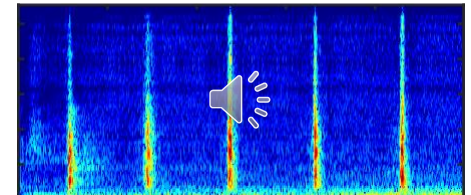
目標音



声真似



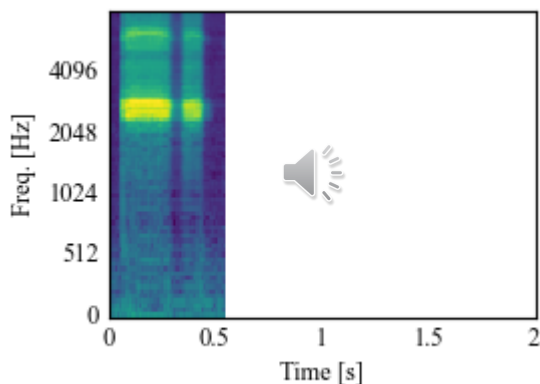
合成音



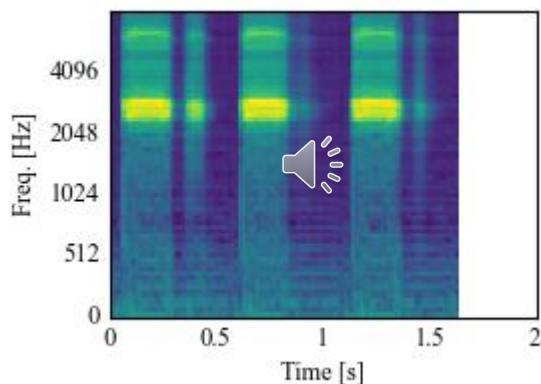
Visual Onoma-to-Wave: 画像を考慮した環境音合成

オノマトペ (擬音) 文字画像からの合成

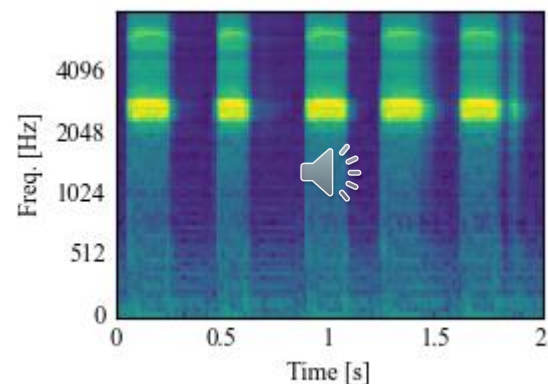
ビッ



ビッビッビッ



ビッビッビッビッビッ



環境音イメージ画像で条件付けした合成

キーン



カッ



ポーン



研究紹介3. 音楽信号解析

- 様々な楽器がまじりあった音楽信号の中から、自分の好きな楽器を見つけ出し、自分の好みのリミックス版を製作する。
- 非負値行列因子分解 (NMF) や深層学習 (DNN) という教師有リアルゴリズムを用い、音を事前に学習したパターンに基づいて分解することにより、信号を解析する。

スパース分解
低ランクモデル
教師有り学習

多重解像度深層分析(1/2): 動機

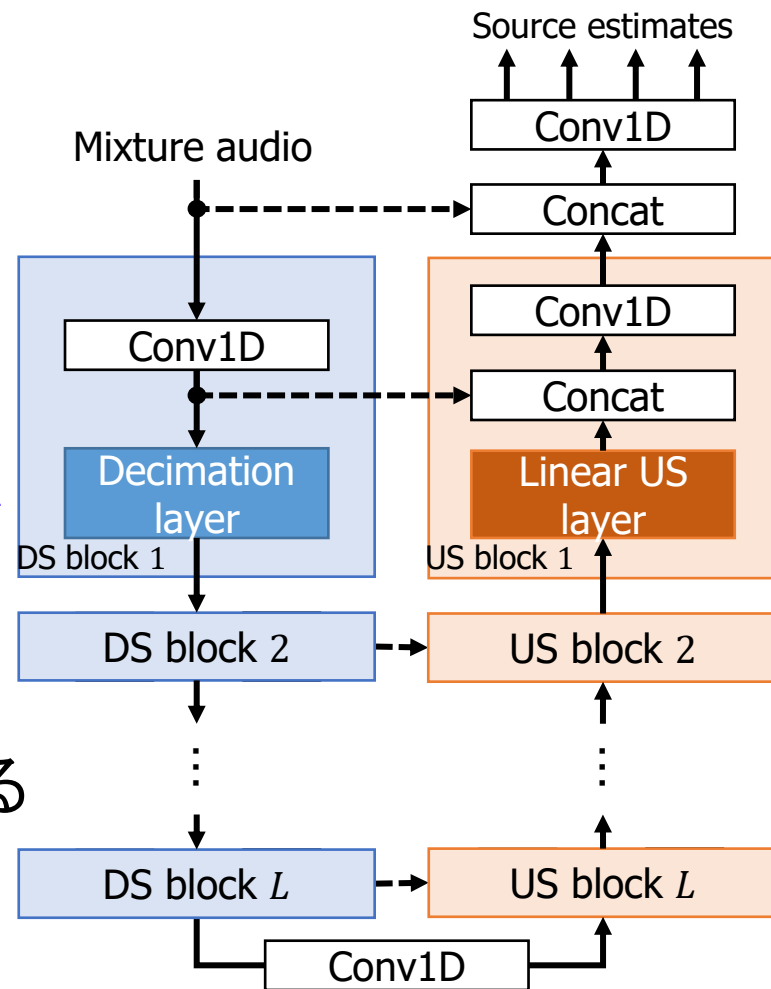
[Nakamura+ IEEE Trans. ASLP2021]

- Wave-U-Net [Stoller+2018]

- 時間信号を直接入力し分離音を得る
End-to-end DNNモデル
- 繰り返しダウン・アップサンプリングを行うU-Net構造をもつ

- しかし, 信号処理の観点から見ると
ダウンサンプリングが問題...

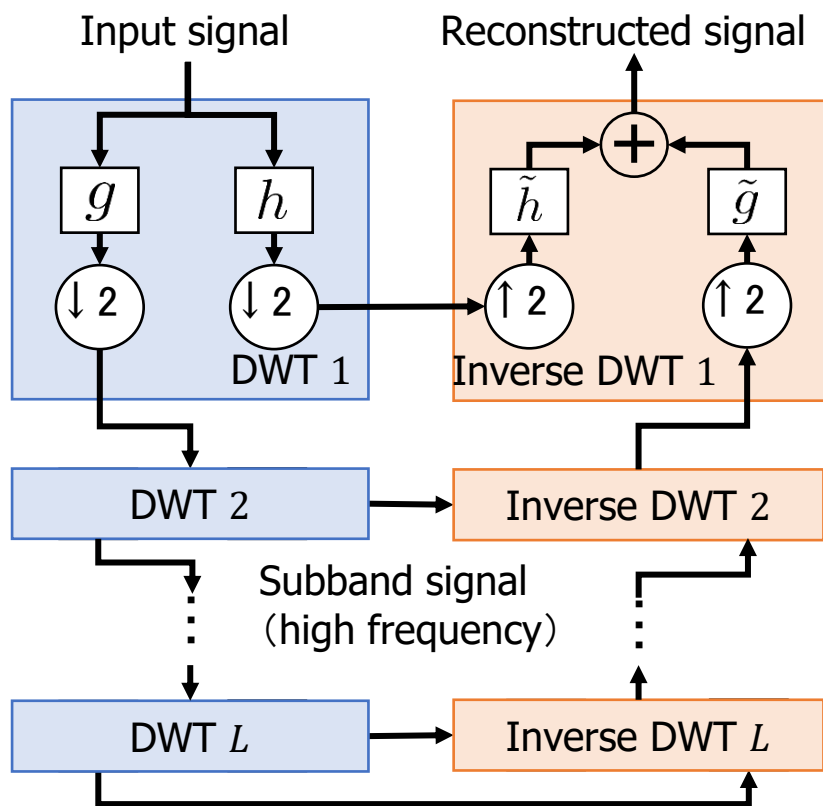
- 特徴量ドメインで**エリアシング**が発生
- ダウンサンプリングで**情報が欠落**する
⇒ **分離性能の低下**を招く 🤔



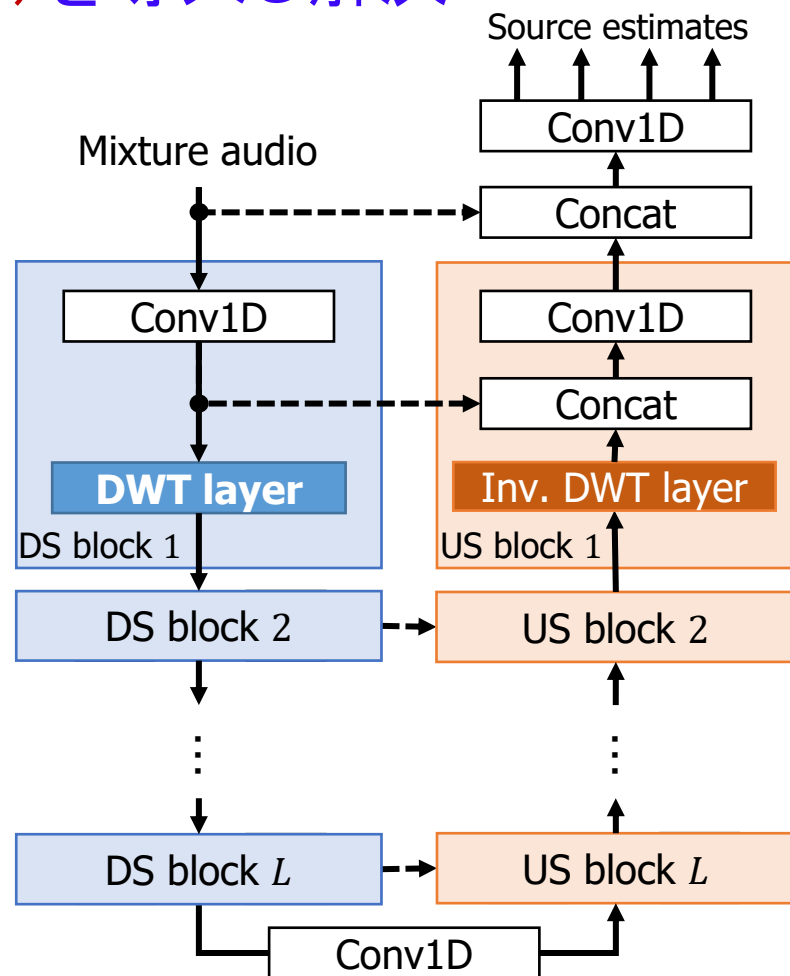
多重解像度深層分析(2/2): 解決方法

[Nakamura+ IEEE Trans. ASLP2021]

- Wave-U-Netと多重解像度解析の構造の類似性に着目
⇒ 離散ウェーブレット変換(DWT)を導入し解決!!



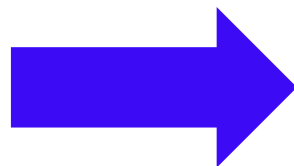
多重解像度分析



多重解像度深層分析(MRDLA)

多重解像度深層分析による楽音分離デモ

- Vocal, bass, drums, guitarの音に分離



分離音

正解音

Vocal: 



Bass: 



Drums: 

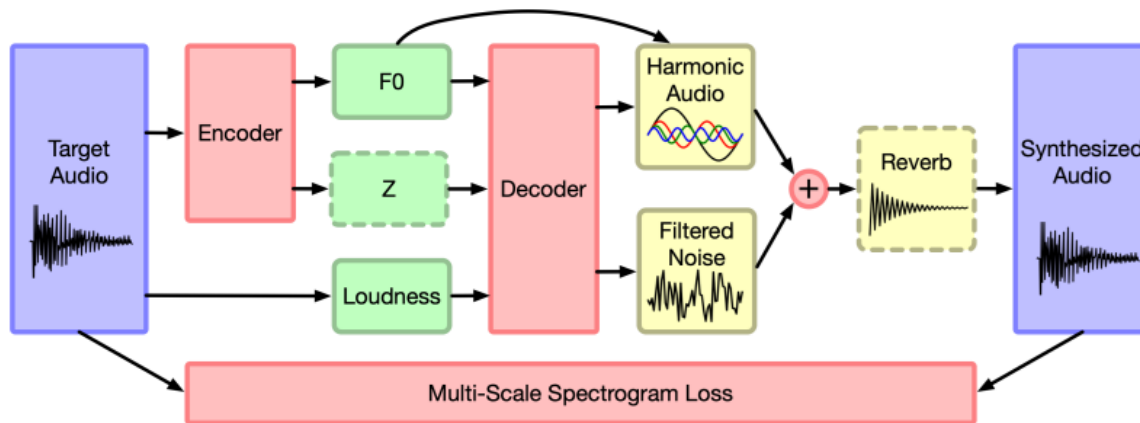


Guitar: 



信号処理を援用した楽器音合成用DNN

- 可微分信号処理 (DDSP) [Engel+2020]
 - シンセサイザーなどで使われる信号処理技術をDNNのモジュールとして利用
- DDSP自己符号化器: 入力音響信号をシンセサイザを用いて再構成



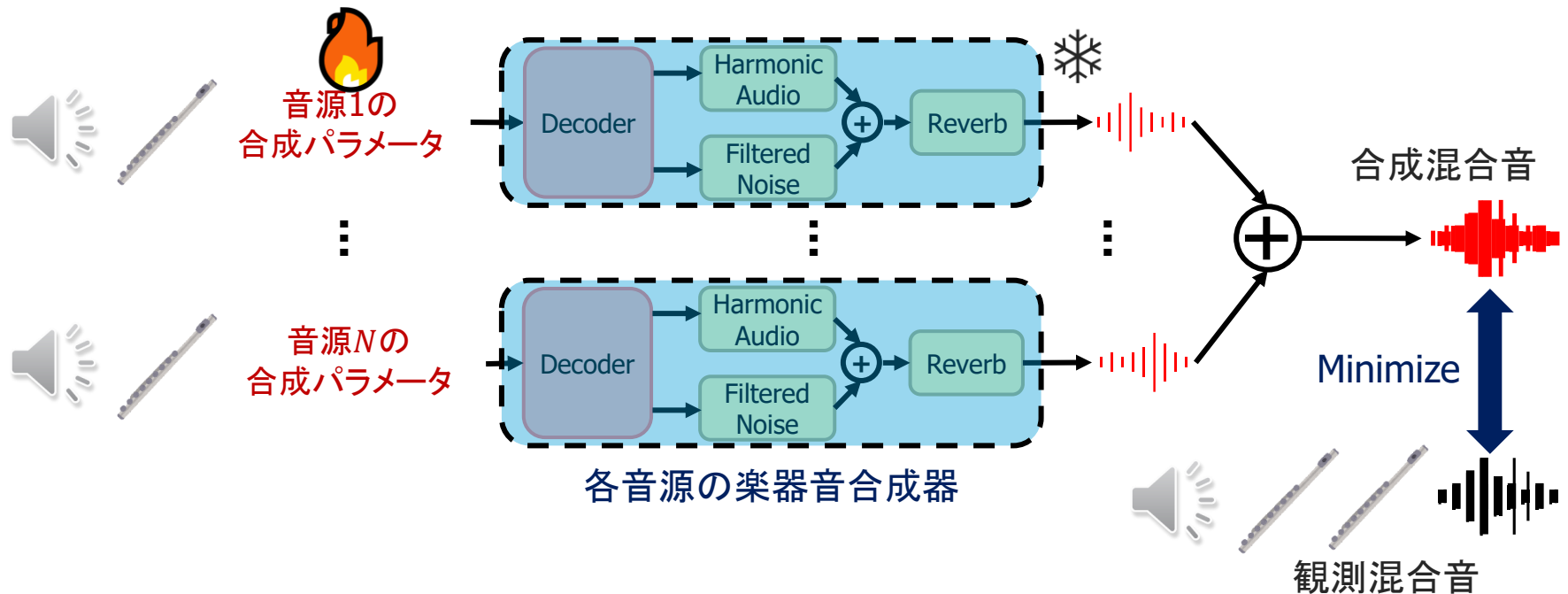
楽器変換
(歌声 to Violin)



DDSP自己符号化器は単旋律の楽器音しか扱えない 🙄

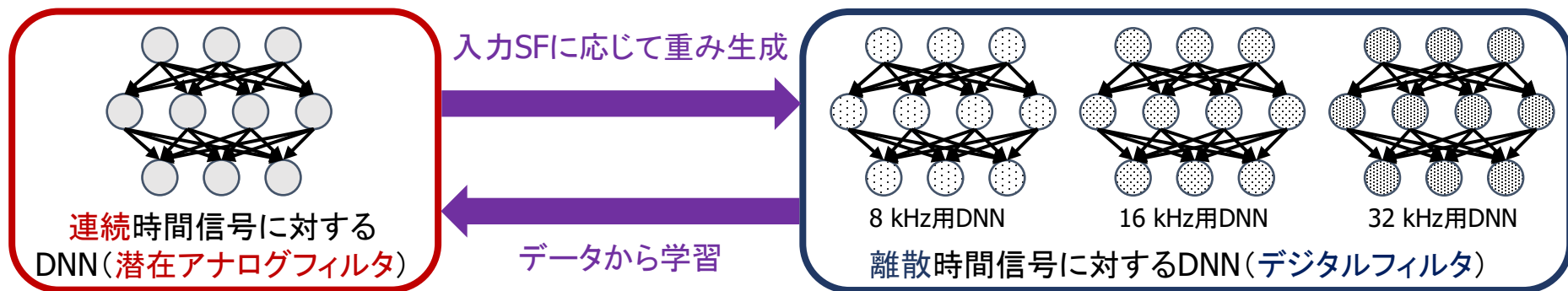
DDSP混合モデルを用いた音源分離 [Kawamura+2022]

- 混合音から各音源の基本周波数, 音色パラメータ, ラウドネスを推定
 - 事前学習したDDSP自己符号化器の一部を楽器音合成器として利用
- 楽譜を援用することで, 同一楽器同士での分離でも高性能に動作!!



標本化周波数非依存深層学習 [Saito+2022]

- DNNにとって、標本化周波数(SF)は関係ない！
 - 同じ現象から生じた音でも異なるSFで標本化されたデータは別物
 - 個別のSF毎にDNNを学習しなおさなければならない
- **連続時間/周波数領域(SF非依存)で定義された関数**から、入力SFに応じて重みを生成する過程を導入 ⇒ **未学習のSFも対応可能!!**
 - 重み生成はアナログフィルタからのデジタルフィルタ設計と解釈可能
 - ⇒ **デジタルフィルタ設計手法を援用可能**



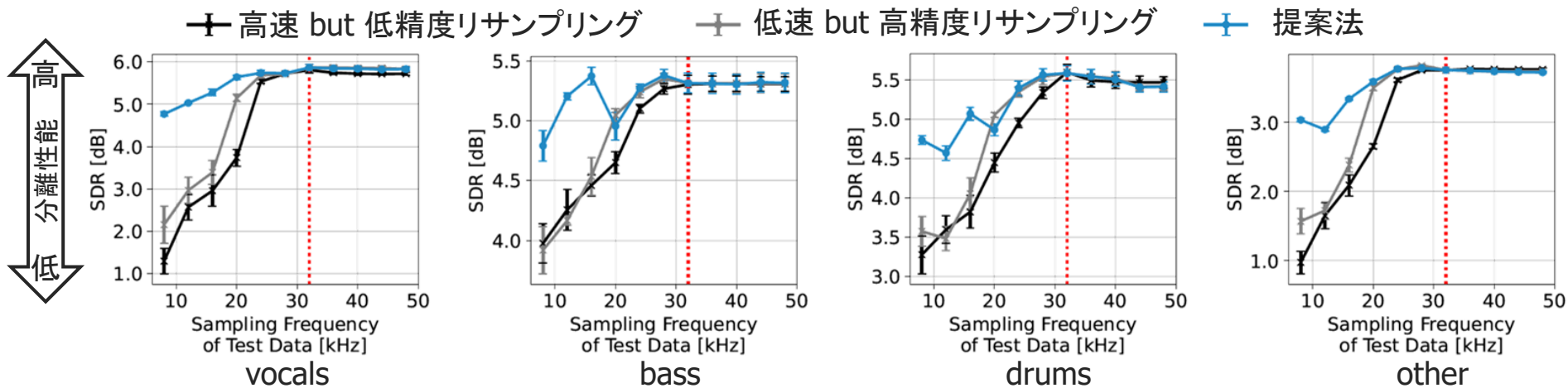
提案アプローチでの学習対象

従来の学習対象

- 様々な拡張も可能: 非整数ストライドへの対応 [Imamura+2023], ニューラルアナログフィルタ [Imamura+2024, in press], etc.

提案法 vs リサンプリング

- 4音源の楽音分離タスクで、提案法をリサンプリングと比較
 - MUSDB18 [Rafii+2018] の学習・テストデータを利用
 - 学習時のSFは32 kHzのみ, テスト時は8~48 kHz(4 kHz毎)のSFを使用
- 提案法は, 未学習のSFでリサンプリングよりも高い分離性能を達成
 - 未学習のSFへの適応方法のみ異なる(i.e., 学習済みモデルは同一)



研究プロジェクト・国内外研究者コラボ

[2026年～]

- ・サイバネティックアバター(内閣府ムーンショット) with 京大河原・阪大石黒先生
- ・楽音信号分解(ヤマハ研究開発センター) with 北村先生・高橋さん・近藤さん
- ・劣決定音源分離(NTT-CS研) with 池下さん・中谷さん
- ・安心声変換(Beyond AIソフトバンク)
- ・次世代音声生成基盤(科研基盤A分担) with 名工大徳田先生
- ・音声言語情報処理の基盤モデル構築(産総研) with 産総研 深山さん・緒方さん
- ・絵文字ベース音声合成(科研費若手)
- ・AIセーフティ(産総研)
- ・音声合成・評価の統合的最適化(JST BOOST)
- ・最適チャンネル選択(科研費若手)
- ・分散アレイデータセット(TAF) with 電気系 高木先生
- ・補聴器のための音声強調(科研費基盤B) with 電気系 高木先生
- ・時変伝達系信号処理(SCAT) with 技科大 若林先生
- ・環境音合成の自動評価(科研費若手)
- ・マルチモーダル環境音認識(SCAT研究助成)
- ・Video-to-audio generation(LINEヤフー) with LINEヤフー橘さん
- ・光ファイバーセンシング(NEC) with NEC 砺波さん

積極的な外部交流を推奨しています！

とにかく音が好き人は集まれ！

2014年以降の研究教育実績

- ・原著論文：Top論文誌IEEE, ASA, EURASIP, ISCA 53編, 他28編
- ・国際・国内会議発表：**数えきれないので略**
- ・教員・指導学生が2014年以降に**149件の学術賞**を受賞



音の情報処理に興味がある人

統計的信号処理の数理に興味がある人

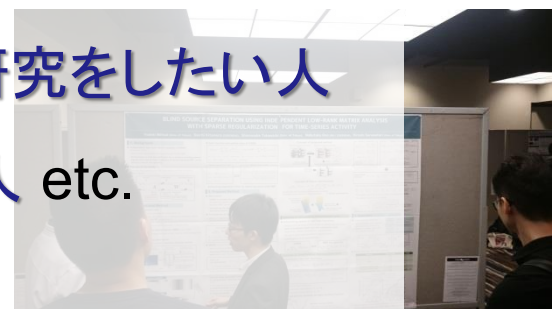
機械学習を使って生のデータを扱いたい人



一期一会なスモールデータの研究をしたい人

波動現象を数理的に考えたい人 etc.

そんなあなたにお勧めです。



猿渡・齋藤研に配属を希望するには？

- 大学院入試受験の際に、以下のどちらかを選択してください。
 - システム情報学専攻へ進学し、猿渡・齋藤研究室を希望
 - 創造情報学専攻へ進学し、猿渡研究室を希望
- **どちらのパスを選択しても、指導体制は変わりません。**
 - 例えば、齋藤先生は創造情報学専攻の教員ではないですが、創造情報学専攻へ進学しても齋藤先生の研究テーマ(音声合成・変換)を希望できます。
 - **ただし、大学院入試の試験科目は異なります。**
(詳細は各専攻の入試案内書を参照)