

2024 年度後期 修士論文

# 深層学習による口唇音声変換に関する研究

A Study on the Conversion from Lip-Video to  
Speech Using a Deep Learning Technique

2025 年月日

九州大学芸術工学府音響設計コース

2DS23095M

南 汰翼

MINAMI Taisuke

研究指導教員 鎚木 時彦 教授

概要

# 目次

<b>1 序論</b>	<b>1</b>
1.1 背景 . . . . .	1
1.2 目的 . . . . .	3
1.3 本論文の構成 . . . . .	4
<b>2 音声信号処理</b>	<b>5</b>
2.1 音声のフーリエ変換 . . . . .	5
2.2 メルスペクトログラム . . . . .	7
<b>3 深層学習</b>	<b>8</b>
3.1 ニューラルネットワーク . . . . .	8
3.2 学習方法 . . . . .	17
<b>4 動画音声合成モデルの検討</b>	<b>18</b>
4.1 音声合成法 . . . . .	18
4.2 実験方法 . . . . .	22
4.3 結果 . . . . .	27
4.4 まとめ . . . . .	41
<b>5 結論</b>	<b>42</b>
<b>謝辞</b>	<b>43</b>
<b>参考文献</b>	<b>44</b>

# 1 序論

## 1.1 背景

音声は基本的なコミュニケーションの手段であり、人と人とのコミュニケーションの場面において、重要な役割を果たしている。音声は、肺からの呼気流による声帯の振動が音源波を生成し、声道特性に伴ったフィルタリングと口唇からの放射特性に従って生成される。これより、音声の生成には音源を作り出す声帯やその制御のための喉頭、舌や口唇といった調音器官の働きが重要となる。しかし、癌などの重い病気で喉頭を摘出した場合、音源波を生成することができなくなるため、これまで通り発声を行うことが不可能になってしまう。このようなコミュニケーション機能の喪失に対し、現在でも電気式人工喉頭や食道発声、シャント発声といった代用音声手法が存在する。電気式人工喉頭では、専用の発振器を顎下に当てて振動を加えることにより、それを音源とした発声を行う。発振器を用意すれば容易に発声することが可能であるが、生成される音声のピッチが発振器の振動に依存してしまうため、抑揚のない単調な音声になってしまう。食道発声では、まず口や鼻から食道内に空気を取り込み、その空気を逆流させることで食道入口部の粘膜を振動させることによって発声する。電気式人工喉頭と違って道具を必要とせず、ピッチも本人が調節できるが、その習得に長期間の訓練を要する。シャント発声では、手術によって気管と食道を繋ぐ管を設ける。これにより、息を吐き出す際に設けられた喉の穴を手で塞ぐことによって肺からの呼気流が食道に流れる。そのため、食道発声と同様に粘膜の振動を音源とし、発声することが可能となる。習得は容易であり、比較的自由に話すことが可能となるが、設けられた管を交換するための定期的な手術が必要となる。このように、現在用いられている代用音声手法にはそれぞれデメリットが存在する。

そのため、本研究ではビデオカメラで撮影した口唇の動きから音声合成を行うことによる、新たな代用音声手法を検討する。本来、音声は声帯の振動や声道の形状に依存して生成されるものであり、口唇の動きのみから音声波形を直接推定することは困難である。そこで、近年画像や自然言語処理、音声といった分野において成果を上げている深層学習を使用し、データ駆動型の方法で口唇の動きと音声の関係性を学習することで推定を行う。これにより、従来の代用音声手法よりも自然性の高い音声を、訓練や定期的な手術の必要なく提供することを目指す。

これまでの動画音声合成は英語が中心に検討が進んでおり、近年ではYouTube上のデータを収集、処理することによって構築した大規模データセット [1, 2] を用いることで、大規模で表現力の高いモデルが構築可能となっている。これにより、従来行われてきた教師あり学習のみならず、動画と音声の関係性を自己教師あり学習 (Self Supervised Learning; SSL) によって学習し、そのモデルを動画音声合成や、動画からテキストを推定する Visual Speech Recognition (VSR) に FineTuning するアプローチが提案され、その有効性が示されている。自己教師あり学習モデルにもいくつかの種類があり、近年多くの研究で応用例のある AVHuBERT [3] は、動画・音声の入力領域においてマスクされた区間の予測と、予測対象の更新を繰り返して学習を進めていくモデルである。予測対象の更新は5回行われ、1回目は音声波形から計算される MFCC をクラスタリングした結果を利用するが、2回目以降はモデルの中間特徴量をクラスタリング

した結果を新たな予測対象に設定する。更新のたびに再度モデルをランダム初期化して再学習するが、その予測対象の複雑さが増していくことによって学習を促進するようなメカニズムとなっている。また、これに類似した VATLM [4] は、動画と音声のみならずテキストも加えた学習によって、精度改善を達成した。その他、Student と Teacher という二つのネットワークを利用し、Teacher から出力される特徴量を Student がマスクされた入力から予測することによって学習を進める RAVEn [5] や AV-data2vec [6]、RAVEn の改善版として提案された BRAVEn [7] など、多くのモデルが提案されている。

近年の動画音声合成や VSR では、こういった SSL モデルを動画からの特徴抽出器として活用しつつ、さらなる工夫によって精度改善を達成している。動画音声合成について、[8] では予測対象として従来用いられてきたメルスペクトログラムに加え、テキストを予測するマルチタスク学習手法を提案した。損失においては上記の二つに加え、予測したメルスペクトログラムを事前学習済みの ASR (Automatic Speech Recognition) モデルに入力して得られる特徴表現も採用した。音声波形はメルスペクトログラムに対して Griffin-Lim アルゴリズムを適用することで獲得しており、従来のメルスペクトログラムのみを損失とする手法に対して客観評価指標における改善を達成した。これに続き、[9] では前述した手法がテキストアノテーションされたデータのみにはしか用いることができないという課題を解消するため、テキストと同様に言語的な情報を持つと考えられている、音声 SSL モデルの HuBERT[10] から得られた離散特徴量 (HuBERT 離散特徴量とする) を利用する手法を提案した。また、予測されたメルスペクトログラムと HuBERT 離散特徴量の両方を入力とする Multi-input Vocoder、Multi-input Vocoder の学習時にメルスペクトログラムにノイズをかけるデータ拡張を合わせて提案し、客観評価と主観評価の両方で改善を達成した。加えて、ここでは AVHuBERT の転移学習についても合わせて検討が行われ、これによってさらに性能を改善できることを示した。手法 [9] に関連して、上記のようなマルチタスク学習手法以外にも、HuBERT 離散特徴量や HuBERT 連続特徴量 (離散化しない場合を指す) を音声波形までの中間特徴量として扱う手法は提案されている。例えば、[11] ではメルスペクトログラムの推定を行わず、HuBERT 離散特徴量のみを推定して音声波形に変換する手法が提案された。[9] では離散化におけるクラスタ数を 200 にしていたのに対して、[11] ではクラスタ数を 2000 と大きく取っている点で実装が異なっている。メルスペクトログラムを省略する分情報圧縮の程度を軽減することで、音声波形への変換に十分な情報を保持する目的があると考えられる。また、[9] や [11] では AVHuBERT を直接動画音声合成に FineTuning していた一方で、[12] では AVHuBERT を VSR によって FineTuning し、その後重みを固定した上で特徴抽出器として利用するアプローチを提案している。動画音声合成モデルは、VSR で FineTuning した AVHuBERT から得られる動画特徴量を入力とし、HuBERT 特徴量を予測するネットワークを導入して、HuBERT 特徴量のみから音声波形に変換するボコーダを利用することで構築される。ここでは HuBERT 特徴量として連続値および離散値の両方が検討され、連続値を用いる場合の方が客観評価指標が改善することを報告している。検討された離散値のクラスタ数が 100 であったため、[11] の結果と合わせると、HuBERT 離散特徴量のみで音声波形に変換するアプローチを取るのであれば、クラスタ数を十分大きく取る必要があ

ることが予想される。

一方 VSR について、[13] では音声認識を利用して言語情報を格納したメモリを用意し、メモリと動画特徴量の間でアテンションをとることによって、ネットワーク内部で言語情報との関連を考慮する構成を提案した。また、[14] では AVHuBERT が動画あるいは音声のどちらを入力とした場合でもクロスモーダルな特徴量を返すことに着目し、音声認識デコーダに組み合わせる AVHuBERT の Few-shot Learning、Zero-shot Learning による転移学習を検討した。加えて、同様に音声認識デコーダを転移学習するアプローチであるが、事前学習済みモデルの重みを固定し、動画特徴量から音声認識モデルの中間特徴量を予測するネットワークのみを新たに学習することで、両者を合併するようなアプローチ [15] も提案されている。さらに、静止画像と音声から動画を合成するネットワークを構築し、音声認識用のデータセットを用いて VSR の学習データを大量に合成するデータ拡張手法 [16] や、事前学習済みの音声認識モデルによって教師なしデータにラベリングを行うデータ拡張手法 [17]、10 万時間分の教師ありデータを新たに増強した研究 [18] など、大規模な学習データを確保することで精度改善を達成した例も報告されている。

上記の研究は英語データを用いたものであったが、VSR においては英語以外の言語に焦点を当てた研究や、多言語対応モデルの構築も検討が進んでいる。[19] では RAVeN を利用し、英語に加えてスペイン語、イタリア語、ポルトガル語など計 6 種類の言語が含まれるデータセット [20, 21, 22] を用いて多言語モデルの構築を検討した。結果として、教師ありデータの少ない英語以外の言語に対する、多言語モデルの有効性が明らかとなった。また、[23] では英語データで学習された AVHuBERT を用いつつ、特定の言語ごとに構築した音声認識モデルのデコーダを転移学習することで、特定言語ごとにモデルを構築するアプローチを提案した。さらに、[24] では音声認識モデルである Whisper を利用し、教師なしデータへのラベリングによるデータ拡張を行うことで、上記二つのアプローチを超える精度を達成した。

本研究では、近年の主流とも言える英語大規模データセットを用いた実験は計算機のスペックの都合上難しく、世界的に見て日本語での動画音声合成の検討例が少ないことも考慮して、文献 [25, 26] で収録された日本語データを用いて研究を行うこととした。英語データと比較して小規模なデータである分性能に課題を抱えたが、予備実験として英語データで学習された AVHuBERT の FineTuning を検討したところ、スクラッチで構築したモデルと比較して、より高い精度を示すことが明らかとなった。これは、英語データを用いた事前学習済みモデルの多言語対応を検討した先行研究の傾向にも一致する結果であり、日本語においても同様に有効なモデルだと考えられる。しかしながら、それでも依然として合成音声の品質は低く、自然音声に迫る合成音は実現されていないことが課題である。

## 1.2 目的

本研究の目的は、動画音声合成によって得られる合成音声の品質を向上させることである。近年高い精度を達成した手法 [9] では、AVHuBERT の利用および、メルスペクトログラムと音

声 SSL 離散特徴量を利用したマルチタスク学習が採用されている。その他にも近年高い精度を達成したモデルは存在 [11, 12, 27] するが、手法 [9] が採用しているマルチタスク学習の有効性は、テキストを用いた先行研究 [8] でも同様に示されている。これより、このアプローチが現状特に有効そうだと判断し、本研究においてはこの手法をベースラインとして、さらなる改善を狙う形で研究を進めることとした。この手法では、動画を入力としてメルスペクトログラムと音声 SSL 離散特徴量を推定し、これら両方を Multi-input Vocoder に入力することで音声波形へと変換する。しかし、動画と音声の間には、同様の口の動きであっても声道形状の違いによって生じる発話内容の曖昧さや、話者によるパターンの多様さが存在すると考え、推定を動画のみに依存した先行研究の手法ではこういった側面への対処が難しいと考えた。これに対して本研究では、音声 SSL モデルである HuBERT を利用した動画音声合成モデルを提案し、合成音声の推定残差を HuBERT を利用した後処理によって軽減することで、合成音声の品質改善を狙った。HuBERT は、音声波形を畳み込み層を通すことによってダウンサンプリングしつつ特徴量に変換し、ここでマスクをかけた上で Transformer 層を通す。そして、マスクされたフレームにおける予測対象を推定する、Masked Prediction を行うことで学習する。大規模な音声データを用いてこの自己教師あり学習を行うことで、音声のコンテキスト自体をデータそのものから学習することが可能であり、音声認識において有効性が確認されている。本研究では、大規模日本語音声データで事前学習済みの HuBERT を活用し、動画音声合成モデルにおいて生じる推定残差を、音声自体のコンテキストを考慮する形で補うようなアプローチを検討した。

### 1.3 本論文の構成

## 2 音声信号処理

音声にはフォルマントや基本周波数（ピッチ）など、様々な周波数的な特徴が存在している。フォルマントは母音や子音を知覚するため、ピッチはアクセントやイントネーションを表現するために重要なものである。このような音声信号の持つ複雑さから、時間波形のままその特徴を分析することは困難である。これに対し、本節では音声の特徴を捉えやすくするための信号処理について説明する。

### 2.1 音声のフーリエ変換

音声の時間波形に対して、周波数領域での情報を得るためにはフーリエ変換（Fourier Transform）を使用する。特に、音声はマイクロフォンで収録され、コンピュータ内で処理されることが多い。この時、音声はアナログ信号ではなく、サンプリング周波数と量子化ビット数に従って離散化されたデジタル信号として扱われる。このような場合、離散信号に対してのフーリエ変換である離散フーリエ変換（discrete Fourier transform; DFT）が用いられる。また、信号の系列長をゼロパディングして2の冪乗の長さに調整することで、計算量を抑えた高速フーリエ変換（fast Fourier transform; FFT）を用いることができる。

離散信号を  $x[n]$ 、それに対するフーリエ変換を  $X[k]$  とする。ここで、 $n$  はサンプルのインデックス、 $k$  は周波数インデックスである。 $X[k]$  は複素数であり、以下のように極座標表示することができる。

$$X[k] = \text{Re}(X[k]) + j\text{Im}(X[k]) \quad (2.1)$$

$$= |X[k]|e^{j\angle X[k]} \quad (2.2)$$

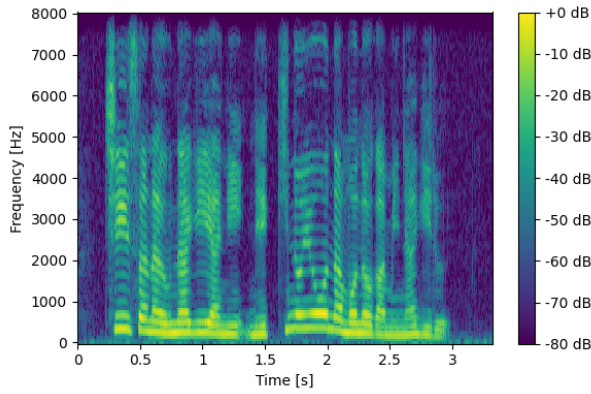
ここで、 $|X[k]|$  は振幅特性（振幅スペクトル）、 $\angle X[k]$  は位相特性（位相スペクトル）であり、以下の式で表される。

$$|X[k]| = \sqrt{\text{Re}(X[k])^2 + \text{Im}(X[k])^2} \quad (2.3)$$

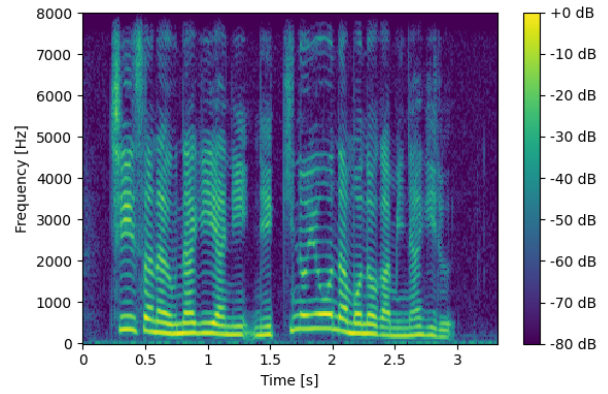
$$\angle X[k] = \tan^{-1} \frac{\text{Im}(X[k])}{\text{Re}(X[k])} \quad (2.4)$$

また、 $|X[k]|^2$  はパワースペクトルと呼ばれる。これにより、信号中にどのような周波数成分がどれくらい含まれているかを調べることができる。しかし、音声はフォルマントやピッチが時々刻々と変化するため、信号全体に対して直接フーリエ変換を適用したとしても有用な結果が得られない。このような音声の持つ非定常性の問題に対して、十分短い時間幅においては信号の定常性が成り立つという仮定のもと、短時間フーリエ変換（short-time Fourier transform; STFT）が用いられる。STFT では、音声信号に対して窓関数による窓処理を適用し、短時間に区切られた信号それぞれに対して DFT を適用する。ここで、窓処理とはある特定の窓関数と音声信号を時間領域でかけ合わせることであり、窓関数の時間幅を窓長という。また、窓関

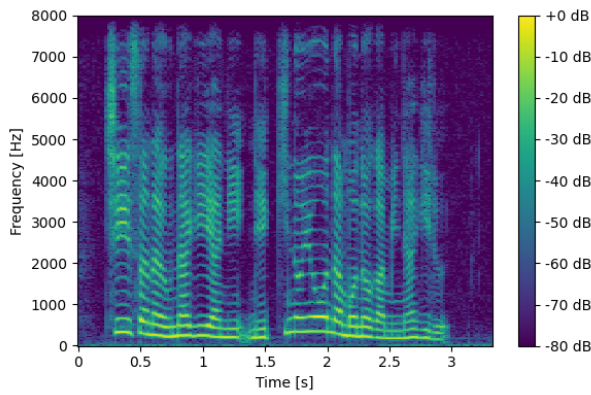




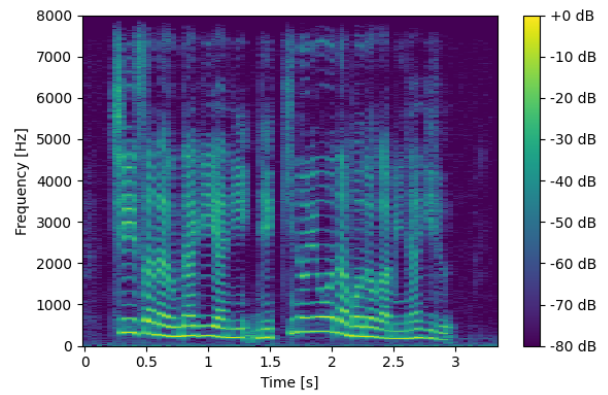
(a) 窓長 12.5ms、シフト幅 5ms



(b) 窓長 25ms、シフト幅 10ms



(c) 窓長 50ms、シフト幅 20ms



(d) 窓長 100ms、シフト幅 40ms

図 2.1: 「小さな鰻屋に、熱気のようなものがみなぎる」と発話した音声から計算された対数パワースペクトログラム

数を時間方向にシフトするときの時間幅をシフト幅という。STFT には時間分解能と周波数分解能の間に不確実性が存在し、両者の間にトレード・オフの関係がある。窓長が長い場合には周波数分解能が向上する一方、時間分解能が低下する。窓長が短い場合にはその逆となる。音声信号  $x[n]$  の STFT を時刻  $j$ 、周波数インデックスを  $k$  として  $X(j, k)$  と表すと、 $X(j, k)$  は時間周波数領域における複素数となる。これを複素スペクトログラムと呼ぶ。また、 $|X(j, k)|$  を振幅スペクトログラム、 $\angle X(j, k)$  を位相スペクトログラム、 $|X(j, k)|^2$  をパワースペクトログラムと呼ぶ。「小さな鰻屋に、熱気のようなものがみなぎる」と発話した音声に対し、窓関数としてハニング窓を用いた上で、複数の窓長・シフト幅によって計算した対数パワースペクトログラムを、図 2.1 に示す。窓長が 100ms と長い場合には周波数分解能が高いが、時間分解能が低下することでスペクトルの時間変化が滑らかでないことがわかる。一方、窓長が 12.5ms と短い場合には時間分解能が高いが、周波数分解能が低下することでスペクトルがぼやけていることがわかる。これが窓長に対する時間分解能と周波数分解能とトレード・オフであり、窓長 25ms や 50ms が程よいパラメータであることがわかる。

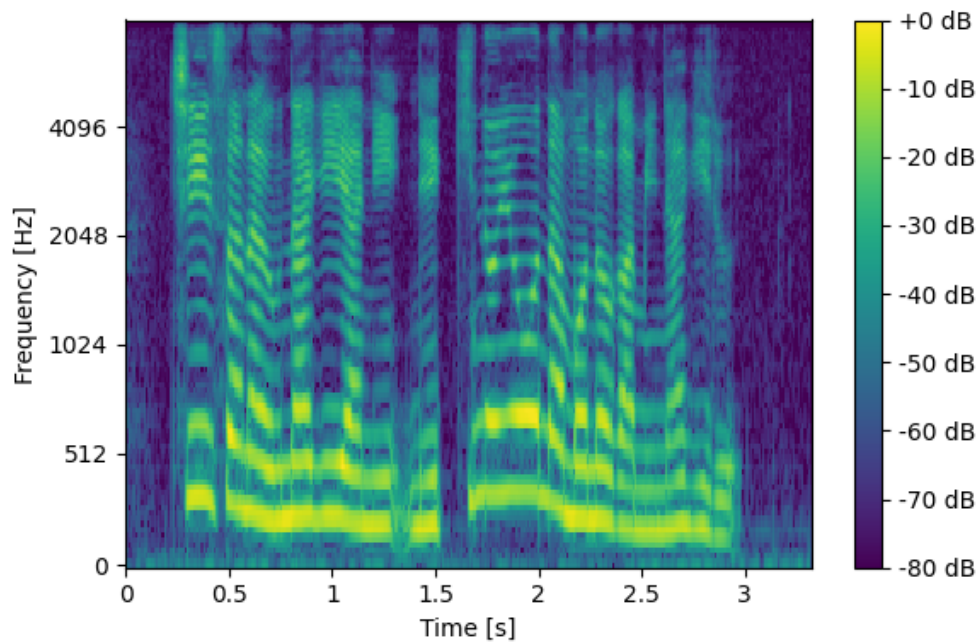


図 2.2: 「小さな鰻屋に、熱気のようなものがみなぎる」と発話した音声に対する対数メルスペクトログラム

## 2.2 メルスペクトログラム

メルスペクトログラムは、振幅スペクトログラムを人間の聴感特性を考慮したメル尺度に変換することによって得られる。周波数軸をメル尺度に変換する際、以下の式を用いる。

$$Mel(f) = 2595 \log_{10} \left( 1 + \frac{f}{700} \right) \quad (2.5)$$

メル尺度は、1000 Hz、40 dB の純音を 1000 mel とする比率尺度である。メル尺度を用いることにより、低い周波数ほど細かく、高い周波数ほど荒い特徴量になる。メルスペクトログラムは、振幅スペクトログラムに対してメルフィルタバンクを適用することによって得られる。メルフィルタバンクの数は任意に決定できるパラメータであり、メルスペクトログラムの周波数方向の次元はこれに一致する。音声合成においては、音声のサンプリング周波数を 16 kHz とするとき、メルフィルタバンクの数を 80 とし、8000 Hz までの帯域に対して適用することが多い。「小さな鰻屋に、熱気のようなものがみなぎる」と発話した音声に対し、窓関数にハニング窓を用い、窓長 25ms、シフト幅 10ms としてパワースペクトログラムを計算した上で、80 次元のメルフィルタバンクを適用して得られた対数メルスペクトログラムを、図 2.2 に示す。

### 3 深層学習

深層学習とは、人間の神経細胞の仕組みを模倣したニューラルネットワークを用いる機械学習手法のことである。特に近年ではその層を深くしたディープニューラルネットワーク（Deep Neural Network; DNN）が用いられ、大量のパラメータを最適化することによって得られる表現力により、自然言語処理や画像処理、音声認識や音声合成など様々な分野で成果を上げている。本章では、本研究で使用するニューラルネットワーク及び、構築した DNN の学習方法について説明する。

#### 3.1 ニューラルネットワーク

##### 3.1.1 全結合層

全結合層は、すべての入力ノードがすべての出力ノードに接続される層である。全結合層の出力は、入力に対して学習可能な重みによる線形変換を適用することで得られる。入力  $\mathbf{x} \in \mathbb{R}^n$  に対し、出力  $\mathbf{y} \in \mathbb{R}^m$  は、

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b} \quad (3.1)$$

で計算される。ここで、 $\mathbf{W} \in \mathbb{R}^{m \times n}$  は重み行列で、各要素  $w_{ij}$  は入力の  $j$  番目の成分と、出力の  $i$  番目の成分の間の重みである。 $\mathbf{b} \in \mathbb{R}^m$  はバイアスベクトルで、各要素  $b_i$  は出力の  $i$  番目の成分のバイアス項である。全結合層は特徴量の次元の変換に用いられ、特に最終層において、ネットワークの特徴量を所望の特徴量の次元に変換する場合に便利である。

##### 3.1.2 畳み込み層

畳み込み層は、入力に対して畳み込み演算を行う層である。一次元畳み込み層について、入力  $\mathbf{x} \in \mathbb{R}^{C_{in} \times T_{in}}$  に対し、出力  $\mathbf{y} \in \mathbb{R}^{C_{out} \times T_{out}}$  は、 $y_k[i]$  を出力テンソルの  $k$  番目のチャンネルの  $i$  番目の成分とすると、

$$y_k[i] = b_k + \sum_{c=0}^{C_{in}-1} \sum_{m=0}^{M-1} x_c \left[ i - \left\lfloor \frac{M}{2} \right\rfloor + m \right] w_{k,c}[m] \quad (3.2)$$

で計算される。ここで、 $x_c \left[ i - \left\lfloor \frac{M}{2} \right\rfloor + m \right]$  が入力テンソルの  $c$  番目のチャンネルの  $\left[ i - \left\lfloor \frac{M}{2} \right\rfloor + m \right]$  番目の成分、 $w_{k,c}[m]$  が出力チャンネル  $k$  と入力チャンネル  $c$  に対応するカーネルの  $m$  番目の成分、 $M$  がカーネルサイズである。上式より、一次元畳み込み層の  $i$  番目の出力は、 $i$  番目の入力を中心とし、カーネルサイズの範囲分が考慮されて得られる値だと解釈できる。一次元畳み込みは、自然言語や音声といった一次元系列に対する特徴抽出のために用いられることが多い。

これに加えて、カーネルを二次元配列とすれば二次元畳み込み層、三次元配列とすれば三次元畳み込み層となる。二次元畳み込み層は主に画像に用いられることが多く、三次元畳み込みは動画に用いられることが多い。

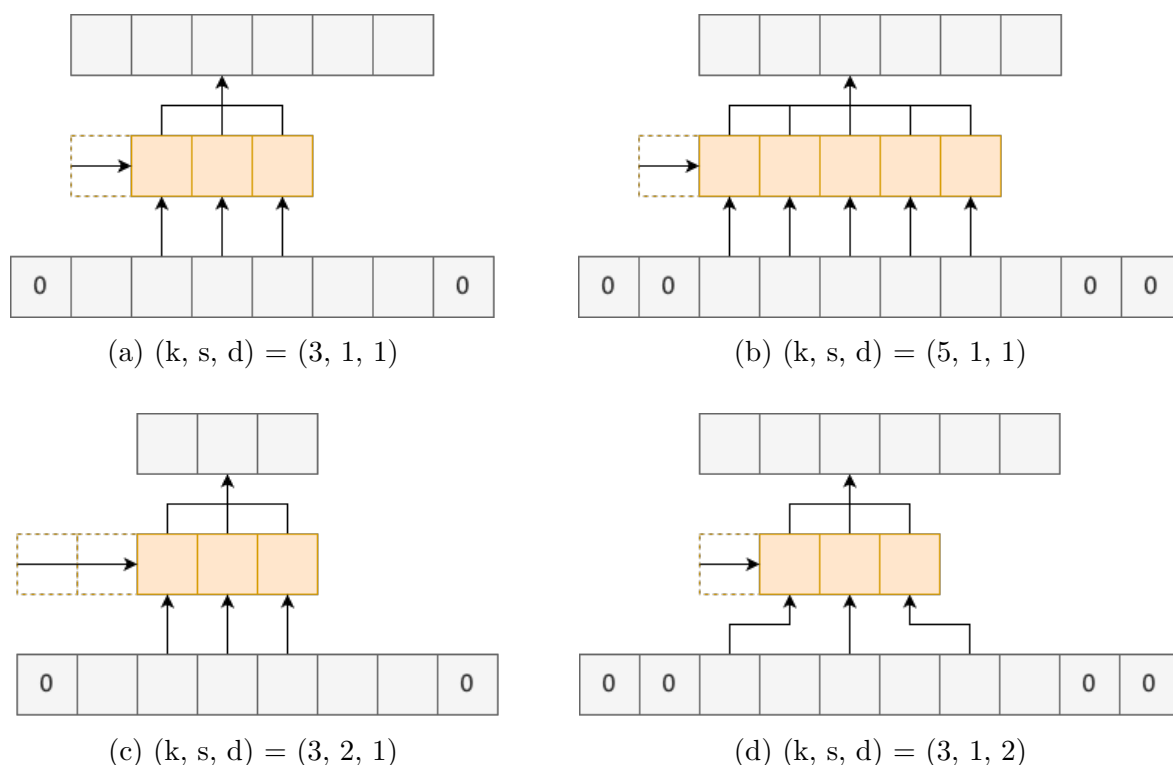


図 3.1: ある入出力チャンネル間における一次元畳み込み層の様子。k はカーネルサイズ、s はストライド、d はダイレーションを表し、図中の 0 はパディング部。

畳み込み層における主要なパラメータは三つある。一つ目は、カーネルサイズである。この値が、畳み込み演算におけるフィルタのサイズを決定しており、考慮できる入力特徴量の範囲が定まる。二つ目は、ストライドである。この値はフィルタのシフト幅にあたり、2 以上の値を設定すると、出力特徴量のサイズは入力特徴量のサイズよりも小さくなる。実際、特徴量をダウンサンプリングしたい場合には、2 以上の値が用いられる。三つ目は、ダイレーションである。これはフィルタ内の要素間の距離を表し、その値を 2 以上とすれば、入力特徴量の飛び飛びの値を考慮した畳み込み演算が行われる。ダイレーションを大きくすることによって、カーネルサイズが同じでも考慮できる範囲が広がるのが特徴であり、系列長の長大な音声波形を扱う場合などに用いられる。また、出力の系列長を入力に対して整数倍に保つためには、上記のパラメータに対して適切なパディング長を指定する必要がある。例えば、カーネルサイズを 3、ストライドとダイレーションを 1 とした場合には、入力の両端に 1 ずつゼロパディングすることで入出力の系列長が保たれる。図 3.1 に、ある入出力チャンネル間における一次元畳み込み層の様子を示す。

### 3.1.3 転置畳み込み層

転置畳み込み層は、畳み込み層の逆の演算を行うような層である。図 3.2 に、ある入出力チャンネル間における一次元畳み込み層の様子を示す。転置畳み込み層では、 $i$  番目の入力とカーネ

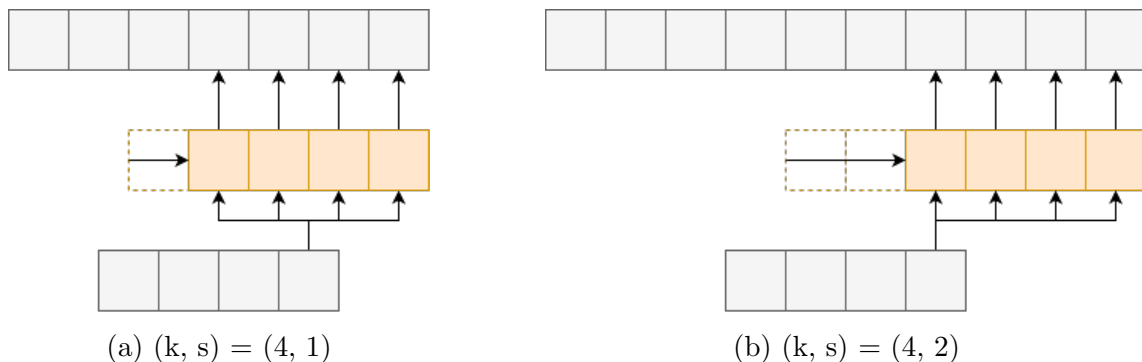


図 3.2: ある入出力チャンネル間における一次元転置畳み込み層の様子。k はカーネルサイズ、s はストライド。

ルの積を計算し、その結果を  $i$  番目から  $i + M - 1$  番目までの出力とする。ここで  $M$  はカーネルサイズを表す。また、 $i$  番目の入力から計算された出力が、 $i - 1$  番目までの入力によってすでに得られている出力とオーバーラップする場合、これらは加算される。図 3.2a は、カーネルサイズを 4、ストライドを 1 とした場合の様子である。転置畳み込み層は特に、入力をアップサンプリングしたい場合に用いられることが多い。その例を図 3.2b に示す。ここでは、カーネルサイズを 4、ストライドを 2 としており、入力系列長が 4 であるのに対して、出力系列長が 10 まで拡大されていることがわかる。ここで、出力系列長を入力系列長の整数倍にするためには、パディングの値を適切に設定する必要がある。転置畳み込み層におけるパディングは、出力の両端を何個落とすかを指定するパラメータである。上述の例においては、パディングを 1 とすることで、出力系列長を入力系列長の 2 倍である 8 に調整することができる。

### 3.1.4 活性化関数

活性化関数は、ニューラルネットワークの各層の出力に非線形性を導入するために用いられる関数である。これにより、ネットワークは単純な線形変換だけでは表現できない複雑な入出力の関係を学習可能になる。以下、代表的な活性化関数を述べる。また、活性化関数とその一階導関数のグラフを図 3.3 に示す。

一つ目は、シグモイド関数である。シグモイド関数は

$$\sigma(x) = \frac{1}{1 + \exp(-x)} \quad (3.3)$$

で与えられ、その一階導関数は

$$\frac{d}{dx}\sigma(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2} \quad (3.4)$$

となる。図 3.3b より、シグモイド関数の一階導関数の最大値は  $x = 0$  における 0.25 であり、入力が 0 から離れるほど微分係数が小さくなることが分かる。DNN の学習では、損失関数の重みに対するの勾配を逆伝播する必要があり、これは合成関数の微分における連鎖律に従って行わ

れる。シグモイド関数は取り得る微分係数の値が小さく、特に層数が深い DNN においては勾配が小さくスケールリングされていくことで、浅い層の重みに対する勾配のノルムが小さくなりすぎる可能性がある。この問題は、勾配消失と呼ばれる。

二つ目は、 $\tanh$  関数である。 $\tanh$  関数は

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \quad (3.5)$$

で与えられ、その一階導関数は

$$\frac{d}{dx} \tanh(x) = \frac{4}{(\exp(x) + \exp(-x))^2} \quad (3.6)$$

$$= \frac{1}{\cosh(x)^2} \quad (3.7)$$

となる。シグモイド関数と比較すると、シグモイド関数の値域が  $[0, 1]$  であったのに対し、 $\tanh$  関数の値域は  $[-1, 1]$  となっている。シグモイド関数や  $\tanh$  関数を実際に活性化関数として利用するニューラルネットワークとして、後述する LSTM や GRU が挙げられる。

三つ目は、ReLU である。ReLU は

$$\text{ReLU}(x) = \max(0, x) \quad (3.8)$$

で与えられ、その一階導関数は

$$\frac{d}{dx} \text{ReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (3.9)$$

となる。ReLU は入力が入力が 0 以上であれば恒等写像となり、0 未満であれば 0 に写す。一階導関数は 0 あるいは 1 のみを取り、特に入力が正の値であれば常に微分係数は 1 となることから、シグモイド関数よりも勾配消失問題が起こりづらくなっている。一方、ReLU は 0 以下の入力に対して微分係数が 0 となるため、勾配も 0 となる。従って、入力が負の値に偏った場合には勾配が 0 になり続け、重みの更新が進まなくなる可能性がある。この問題を、Dying ReLU 問題と呼ぶ。

四つ目は、LeakyReLU である。LeakyReLU は

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{if } x \leq 0 \end{cases} \quad (3.10)$$

で与えられ、その一階導関数は

$$\frac{d}{dx} \text{LeakyReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ a & \text{if } x \leq 0 \end{cases} \quad (3.11)$$

となる。ここで、 $a$  は negative slope と呼ばれるパラメータである。ReLU と比較すると、0 以下の入力に対しても 0 でない値を取り、また微分係数も 0 にならない点が異なっている。この工夫により、勾配が 0 になり続ける恐れがなくなるため、前述した Dying ReLU 問題に対処することが可能である。具体的な利用ケースとしては、敵対的生成ネットワーク（Generative adversarial network; GAN）における Generator や Discriminator が挙げられる。

五つ目は、PReLU である。これは四つ目に述べた LeakyReLU と非常に似た活性化関数であるが、LeakyReLU の negative slope を学習可能なパラメータに変更したという点で異なっている。

六つ目は、GELU である。GELU は

$$\text{GELU}(x) = x\Phi(x) \quad (3.12)$$

で与えられる。ここで、 $\Phi(x)$  は標準正規分布の累積分布関数である。GELU の一階導関数は、

$$\frac{d}{dx}\text{GELU}(x) = \Phi(x) + \frac{x}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \quad (3.13)$$

となる。GELU は入力に対して標準正規分布に基づく確率値を利用したスケーリングを行う活性化関数であり、後述する Transformer など大規模なモデルにおいて採用される場合がある。

### 3.1.5 再帰型ニューラルネットワーク

再帰型ニューラルネットワーク（Recurrent Neural Network; RNN）は、自身の過去の出力を保持し、それをループさせる再帰的な構造を持ったネットワークである。

近年よく用いられる RNN として、長・短期記憶（Long Short-Time Memory; LSTM） [28] がある。LSTM では入力ゲート、忘却ゲート、出力ゲートの 3 つを持ち、これらゲートによってネットワーク内部の情報の取捨選択を行うことで、長い系列データからの学習を可能にした。LSTM のネットワーク内部で行われる計算を以下に示す。

$$\mathbf{f}_t = \sigma(\mathbf{W}_f[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_f) \quad (3.14)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_i[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_i) \quad (3.15)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_h[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_h) \quad (3.16)$$

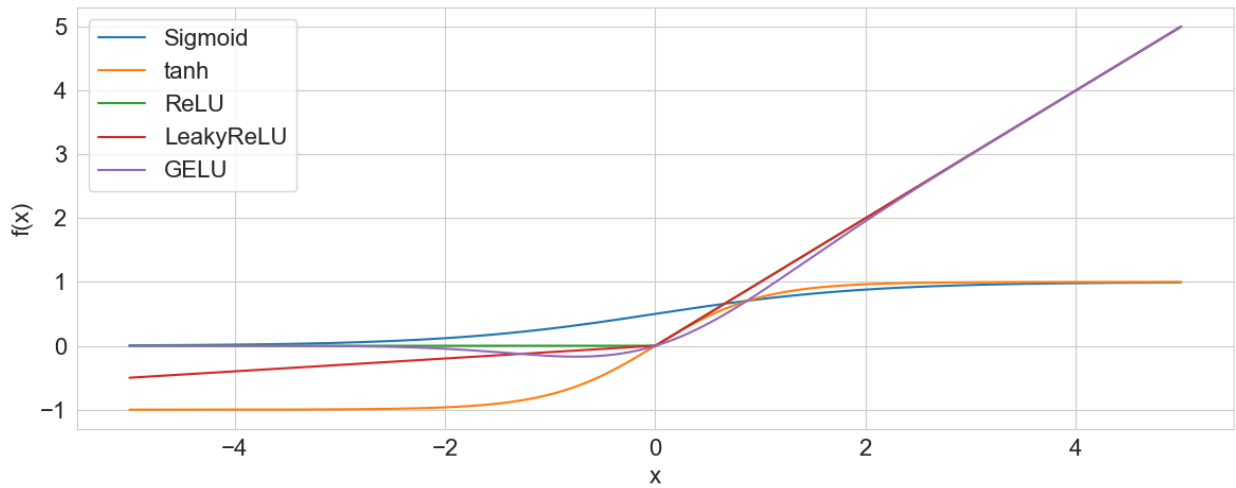
$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \quad (3.17)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_o) \quad (3.18)$$

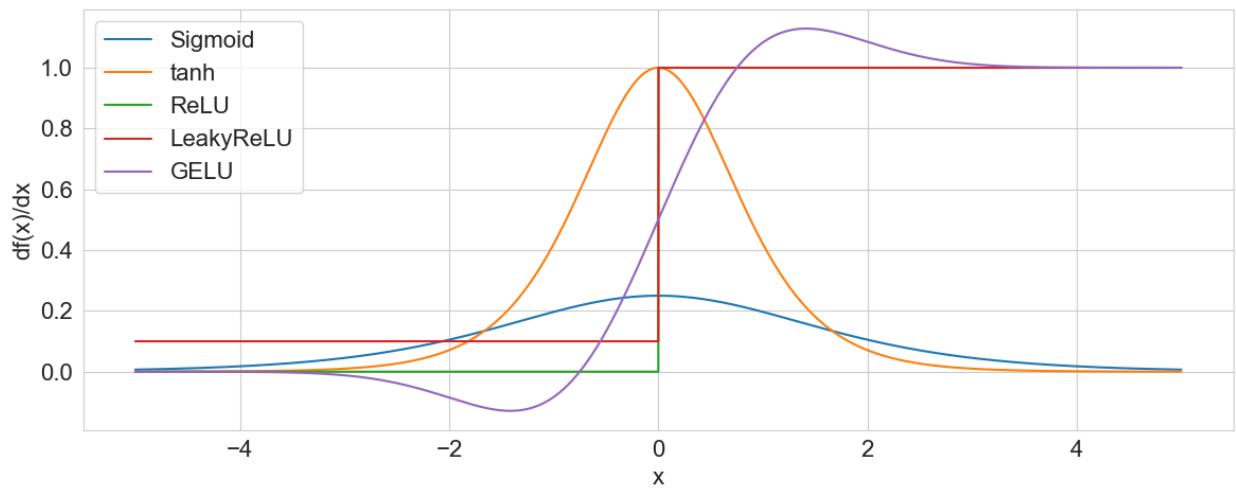
$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (3.19)$$

ここで、 $\mathbf{x}_t$  は時刻  $t$  の入力、 $\mathbf{f}_t$  は忘却ゲートの出力、 $\mathbf{i}_t$  は入力ゲートの出力、 $\mathbf{c}_t$  が時刻  $t$  におけるセルの状態、 $\mathbf{o}_t$  が出力ゲートの出力、 $\mathbf{h}_t$  が時刻  $t$  における隠れ状態を表し、それぞれベクトルである。また、 $\mathbf{W}$ 、 $\mathbf{b}$  はそれぞれ学習可能な重み行列、バイアスベクトルである。忘却ゲートの出力  $\mathbf{f}_t$  は、前時刻のセル状態  $\mathbf{c}_{t-1}$  との要素積に用いられ、これによってネットワーク内





(a) 活性化関数



(b) 活性化関数の一階導関数

図 3.3: 活性化関数の例

に保持されていた長期記憶に対して忘却すべき情報を与え、取捨選択を行う。入力ゲートの出力  $i_t$  は時刻  $t$  のセル状態の候補となる  $\tilde{c}_t$  との要素積に用いられ、これによって新たにネットワークが長期記憶すべき情報が選択され、セル状態  $c_t$  として保持される。出力ゲートの出力  $o_t$  は時刻  $t$  のセル状態  $c_t$  に  $\tanh$  を適用した値との要素積に用いられ、これによって時刻  $t$  において出力すべき情報を選択する。最後に、隠れ状態  $h_t$  は、時刻  $t$  の出力と  $\tanh$  を適用したセル状態  $c_t$  から決定される。このように、 $h_t$  及び  $c_t$  を入力に対して適応的に変化させながら情報を取捨選択することで、長い系列長のデータに対しての学習を可能にした。

また、LSTM では3つのゲートを必要とするが、ゲートを2つに減らすことでネットワークの軽量化を図ったネットワークとして、ゲート付き回帰型ユニット (Gated Recurrent Unit; GRU) [29] がある。GRU ではリセットゲートと更新ゲートの2つのゲートを用いて隠れ状態  $h_t$  を再



帰的に更新する。また、LSTM と違ってセル状態  $c_t$  を省いており、よりシンプルな構造となっている。GRU のネットワーク内部で行われる計算を以下に示す。

$$\mathbf{z}_t = \sigma(\mathbf{W}_z[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_z) \quad (3.20)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_r[\mathbf{x}_t, \mathbf{h}_{t-1}] + \mathbf{b}_r) \quad (3.21)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h[\mathbf{x}_t, \mathbf{r}_t \odot \mathbf{h}_{t-1}] + \mathbf{b}_h) \quad (3.22)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (3.23)$$

ここで、 $\mathbf{x}_t$  が時刻  $t$  における入力、 $\mathbf{z}_t$  が更新ゲートの出力、 $\mathbf{r}_t$  がリセットゲートの出力、 $\mathbf{h}_t$  が時刻  $t$  における隠れ状態を表し、それぞれベクトルである。また、 $\mathbf{W}$  は学習可能な重み行列、 $\mathbf{b}$  はバイアスベクトルである。更新ゲートの出力  $\mathbf{z}_t$  はネットワーク内で新たに記憶すべき情報の割合を決定する役割を果たす。一方、リセットゲートの出力  $\mathbf{r}_t$  は前時刻の隠れ状態  $\mathbf{h}_{t-1}$  との要素積に用いられ、これによって忘却すべき情報が選択される。リセットゲートの出力によって処理された内部表現から、時刻  $t$  の隠れ状態の候補となる  $\tilde{\mathbf{h}}_t$  が計算される。最後に、前時刻の隠れ状態  $\mathbf{h}_{t-1}$  と隠れ状態の候補  $\tilde{\mathbf{h}}_t$  に対して更新ゲートの出力  $\mathbf{z}_t$  を用いた重み付け和を計算することで、隠れ状態  $\mathbf{h}_t$  を決定する。

### 3.1.6 Transformer

Transformer [30] は、自己注意機構 (Self-Attention) を用いて、入力系列全体に渡る依存関係を捉えることができるニューラルネットワークである。特に、再帰的な計算を必要とする RNN と比較して、Transformer は並列計算のみ行うため、GPU による計算の高速化が可能である。Transformer の内部構造について説明するため、入力系列  $\mathbf{X} \in \mathbb{R}^{T \times d_{model}}$  をとる。ここで、 $T$  は入力系列の系列長、 $d_{model}$  は入力系列の次元である。また、Transformer 層の構造を図 3.4 に示す。

まず、Self-Attention について、これは以下の二つの要素からなる。

1. クエリ、キー、バリューの計算
2. Attention スコアの計算

まず、クエリ、キー、バリューの計算は、それぞれを  $\mathbf{Q} \in \mathbb{R}^{T \times d_k}$ 、 $\mathbf{K} \in \mathbb{R}^{T \times d_k}$ 、 $\mathbf{V} \in \mathbb{R}^{T \times d_v}$  とおくと、

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q \quad (3.24)$$

$$\mathbf{K} = \mathbf{X}\mathbf{W}_K \quad (3.25)$$

$$\mathbf{V} = \mathbf{X}\mathbf{W}_V \quad (3.26)$$

で与えられる。ここで、 $d_k$  はキーとクエリの次元、 $d_v$  はバリューの次元、 $\mathbf{W}_Q \in \mathbb{R}^{d_{model} \times d_k}$  はクエリに対する重み行列、 $\mathbf{W}_K \in \mathbb{R}^{d_{model} \times d_k}$  はキーに対する重み行列、 $\mathbf{W}_V \in \mathbb{R}^{d_{model} \times d_v}$  は

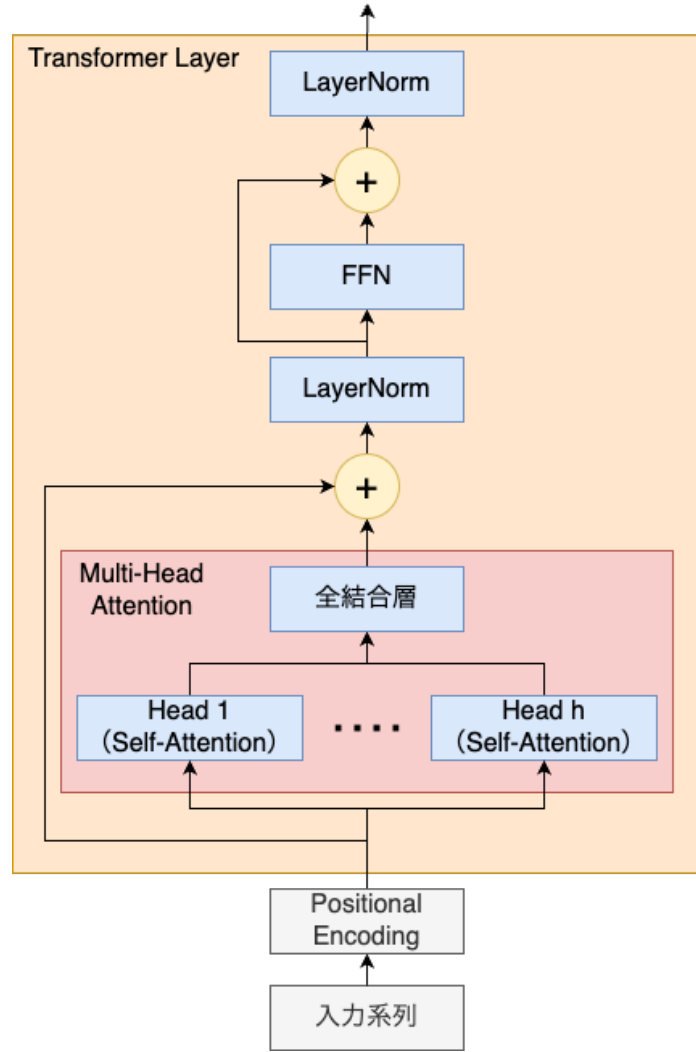


図 3.4: Transformer 層の構造

バリューに対する重み行列である。次に、ここで得られたクエリ、キー、バリューを利用した Attention スコアの計算は、

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (3.27)$$

で与えられる。式中の  $\mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{T \times T}$  について、この行列の  $(i, j)$  成分  $(\mathbf{Q}\mathbf{K}^\top)_{i,j}$  は、クエリベクトル  $Q_i \in \mathbb{R}^{d_k}$  と、キーベクトル  $K_j \in \mathbb{R}^{d_k}$  の内積であり、 $Q_i$  と  $K_i$  の関連度を表していると解釈できる。また、softmax 関数は列方向に適用されるため、アテンション重み  $\text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right)$  の  $(i, j)$  成分  $\left( \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \right)_{i,j}$  は、

$$\left( \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \right)_{i,j} = \frac{\exp \left( \frac{(\mathbf{Q}\mathbf{K}^\top)_{i,j}}{\sqrt{d_k}} \right)}{\sum_{j=1}^T \exp \left( \frac{(\mathbf{Q}\mathbf{K}^\top)_{i,j}}{\sqrt{d_k}} \right)} \quad (3.28)$$

となる。softmax 関数を適用することでアテンション重みの行ベクトルは確率値として扱うこ

とができるようになり、これは  $i$  番目のクエリからのキーの全要素に対する注意度が、確率分布として表されていると解釈できる。最後にアテンション重みをバリュー  $V$  にかけることにより、各クエリからの注意度をもとにバリューベクトル  $V_j$  の重み付け和が計算され、出力が得られる。

Transformer では、Self-Attention をただ行うだけではなく、これを複数のヘッドで並列に計算し、各ヘッドの出力を結合して最終出力を得る Multi-Head Attention が行われる。ヘッド数を  $h$  とすると、各ヘッドにおけるクエリ、キー、バリューの計算は、

$$Q^i = XW_Q^i \quad (3.29)$$

$$K^i = XW_K^i \quad (3.30)$$

$$V^i = XW_V^i \quad (3.31)$$

で与えられる。ここで、 $W_Q^i \in \mathbb{R}^{d_{model} \times \frac{d_k}{h}}$  は  $i$  番目のヘッドのクエリに対する重み行列、 $W_K^i \in \mathbb{R}^{d_{model} \times \frac{d_k}{h}}$  は  $i$  番目のヘッドのキーに対する重み行列、 $W_V^i \in \mathbb{R}^{d_{model} \times \frac{d_v}{h}}$  は  $i$  番目のヘッドのバリューに対する重み行列である。よって、 $i$  番目のヘッドのクエリ、キー、バリューの形状はそれぞれ、 $Q^i \in \mathbb{R}^{T \times \frac{d_k}{h}}$ 、 $K^i \in \mathbb{R}^{T \times \frac{d_k}{h}}$ 、 $V^i \in \mathbb{R}^{T \times \frac{d_v}{h}}$  となる。 $i$  番目のヘッドにおける Attention スコアの計算は

$$\text{Attention}(Q^i, K^i, V^i) = \text{softmax} \left( \frac{Q^i (K^i)^\top}{\sqrt{d_k/h}} \right) V^i \quad (3.32)$$

で行われる。その後、すべてのヘッドからの Attention スコアを結合し、さらに線形変換を適用することで Multi-Head Attention の出力が得られる。これは、

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}^1, \dots, \text{head}^h) W_o \quad (3.33)$$

と表される。ここで、 $\text{head}^i \in \mathbb{R}^{T \times \frac{d_v}{h}}$  が各ヘッドから出力された Attention スコア、 $\text{Concat}(\text{head}^1, \dots, \text{head}^h) \in \mathbb{R}^{T \times d_v}$  が全ヘッドの Attention スコアを結合した特徴量、 $W_o \in \mathbb{R}^{d_v \times d_{model}}$  が最終的な線形変換に用いられる重み行列である。ヘッドを分割することで複数パターンのアテンションが可能になり、これが入出力間の複雑な関係性を考慮するのに役立っていると考えられる。Multi-Head Attention 後には、残差結合と後述するレイヤー正規化を適用する。これは、出力系列を  $Y \in \mathbb{R}^{T \times d_{model}}$  とすると、

$$Y = \text{LayerNorm}(\text{MultiHead}(Q, K, V) + X) \quad (3.34)$$

で与えられる。

Multi-Head Attention によって系列全体の依存関係を考慮した後は、各フレームごとに独立して全結合層を適用する。これは、

$$\text{FFN}(Y) = \text{ReLU}(Y W_1 + b_1) W_2 + b_2 \quad (3.35)$$

で与えられる。ここで、 $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}$ 、 $W_2 \in \mathbb{R}^{d_1 \times d_{model}}$  は学習可能な重み行列、 $b_1 \in \mathbb{R}^{d_{ff}}$ 、 $b_2 \in \mathbb{R}^{d_{model}}$  はバイアスベクトルである。 $d_{ff}$  は  $d_{model}$  の 4 倍とされることが多い。全結合層適用後は、Multi-Head Attention 後と同様に、残差結合とレイヤー正規化を適用する。

上述した Multi-Head Attention とフレームごとに適用される全結合層を合わせて、Transformer 層と呼ぶ。実際には、Transformer 層を多層積み重ねて用いることが多い。

最後に、Transformer では RNN と違い、並列計算によって系列全体を一度に処理することが可能であるが、それと引き換えに入力 of 順序情報を考慮することができなくなる。この問題に対し、入力系列に対して位置情報を与えるために行われるのが、Positional Encoding である。Positional Encoding は sin 関数と cos 関数に基づいて計算される値であり、

$$\text{PE}_{pos,2i} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (3.36)$$

$$\text{PE}_{pos,2i+1} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (3.37)$$

で与えられる。ここで、 $pos$  は入力系列の位置を表し、 $i$  は入力系列の次元を表す。

## 3.2 学習方法

### 3.2.1 損失関数

### 3.2.2 確率的勾配降下法

### 3.2.3 誤差逆伝播法

### 3.2.4 最適化手法

### 3.2.5 学習率のスケジューリング

### 3.2.6 特徴量の正規化

DNN の学習過程では学習の進行に伴って重みが変わるため、その度に各層への入力の分布が変わってしまう。これは内部共変量シフト (Internal Covariate Shift) と呼ばれ、ネットワークの学習を不安定にする原因となる。これに対し、バッチ正規化 (Batch Normalization) が有効である。バッチ正規化は、ミニバッチ学習を行う際にそのミニバッチ内における平均と標準偏差を計算し、特徴量を標準化する。

### 3.2.7 正則化

### 3.2.8 学習の安定化

### 3.2.9 自己教師あり学習

## 4 動画音声合成モデルの検討

### 4.1 音声合成法

提案手法の構築手順は3段階に分かれる。ネットワークの構造を図4.1に示す。一段階目では、動画を入力として、メルスペクトログラムとHuBERT離散特徴量、HuBERT中間特徴量を推定するネットワークAを学習する(図4.1のA)。ここで、HuBERT離散特徴量はHuBERT Transformer層から得られる特徴量をk-means法によってクラスタリングすることで離散化した値、HuBERT中間特徴量はHuBERTにおける畳み込み層出力で、HuBERTの事前学習時にマスク対象となる値のことを指す。図4.2にこれらの取得位置を示す。第一段階では、AVHuBERTを動画からの特徴抽出に利用した。これにより、動画の空間情報は完全に圧縮され、768次元の一次元系列となる。その後、事前学習済みの話者識別モデル[31]によって音声波形から得られる256次元の話者Embeddingを、各フレームでチャンネル方向に結合する。これによって特徴量は1024次元に拡張され、全結合層によって再度768次元に圧縮する。その後、畳み込み層と全結合層からなるConvDecoder(図4.1のConvDecoder)を通すことによって、話者Embeddingを結合した特徴量に対する変換を施した。これにより、特にメルスペクトログラムにおいて話者性が正しく反映されることを狙った。複数話者モデルであっても、入力である動画の見た目から話者性を判別できる可能性があったが、将来的な未知話者対応への拡張性も考慮して、本研究では補助特徴量として入力することとした。ConvDecocerは残差結合を利用したブロック単位で構成され、各ブロックに2層の畳み込み層を設けた。各畳み込み層のチャンネル数は768、カーネルサイズは3であり、3ブロック積み重ねた。最後に全結合層を通し、所望の次元に変換することで予測対象を得た。ネットワークAの役割は、続くネットワークBの入力であるHuBERT中間特徴量を提供することである。これに対し、メルスペクトログラムとHuBERT離散特徴量の推定を同時に行った理由は、先行研究においてマルチタスク学習の有効性が確認されていることを考慮し、ネットワークAでもマルチタスク学習を採用しておこうと考えたからである。

二段階目では、一段階目に学習されたネットワークAの重みを固定した状態でHuBERT中間特徴量を推定し、それを入力としてメルスペクトログラムとHuBERT離散特徴量を推定する、HuBERT Transformer層を中心としたネットワークBの学習を行う(図4.1のB)。HuBERT Transformer層出力はAVHuBERT出力と同じ768次元の特徴量となるため、これに対してネットワークAと同様に話者Embeddingを結合し、ConvDecoderを通すことで予測値を得た。ネットワークBの役割は、音声波形への変換に必要なメルスペクトログラムとHuBERT離散特徴量の予測である。HuBERT Transformer層の転移学習を検討した狙いについて、HuBERTは自己教師あり学習時、畳み込み層出力にマスクを適用し、Transformer層を通すことによってマスクされた部分を推定しようとする。これにより、音声の文脈を考慮するのに長けた学習済み重みが、特にTransformer層で獲得されると仮定した。これに基づき、本研究ではHuBERT Transformer層を動画音声合成にFine Tuningすることにより、動画を入力としたAVHuBERTを中心とするネットワークAにおける推定残差を、音声自体の文脈を考慮することによって軽

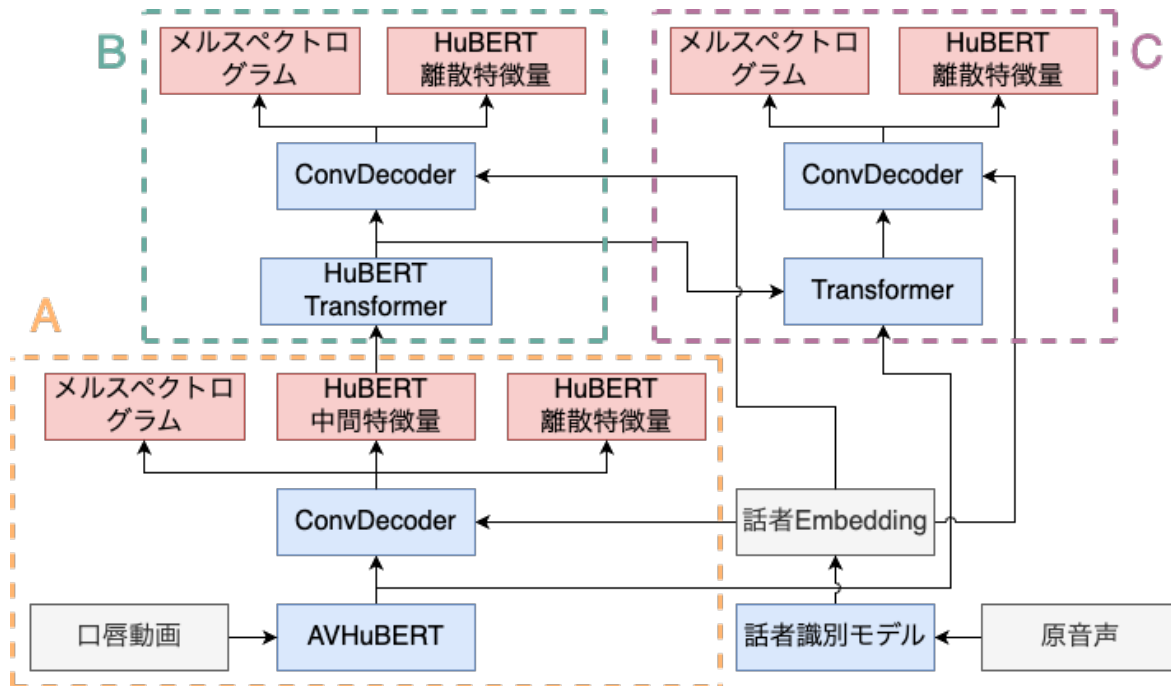


図 4.1: 提案するネットワークの構造

減し、動画から直接推定しきれなかった部分を補うことでの精度改善を狙った。

三段階目では、二段階目までに学習されたネットワーク A とネットワーク B の重みを固定した状態で、AVHuBERT から得られる特徴量と、HuBERT Transformer 層から得られる特徴量の二つを結合し、それらを入力として再びメルスペクトログラムと HuBERT 離散特徴量の予測を行うネットワーク C を学習した (図 4.1 の C)。ネットワーク C では、はじめに前述した二つの特徴量をチャンネル方向に結合することで、1536 次元の入力特徴量を得る。これに対して全結合層を施すことで再度 768 次元に圧縮し、4 層の Transformer 層を通すことで系列全体を考慮した特徴抽出を改めて行った。その後、ネットワーク A,B と同様に話者 Embedding を結合し、ConvDecoder を通すことによって予測値を得た。ここで、ネットワーク C の Transformer 層におけるパラメータについては、AVHuBERT や HuBERT と同様にチャンネル数を 768、ヘッド数は 12 とした。ネットワーク C の役割は、ネットワーク B と同様に音声波形への変換に必要な特徴量の予測である。ここでの狙いについて、まず、AVHuBERT から得られる特徴量と HuBERT Transformer 層から得られる特徴量は、どちらも ConvDecoder への入力となる点で同じである。一方、AVHuBERT は動画を入力、HuBERT Transformer 層は HuBERT 中間特徴量を入力とするため、これら特徴量の元となる入力は異なっている。ここでは、概ね同じ予測対象のために利用される二つの特徴量 (ネットワーク A では HuBERT 中間特徴量の予測も行っているため、全く同じではない) が、入力の違いに依存して内部の Self Attention により注意される部分が変化し、何らかの異なった情報を持っている可能性があると仮定した。この仮定に基づき、両方の特徴量を考慮して単一特徴量への依存を解消することで、汎化性能向上による予測精度の改善を狙った。

以上が提案手法の全体像であるが、今回ベースラインとする先行研究 [9] に基づいたマルチタ

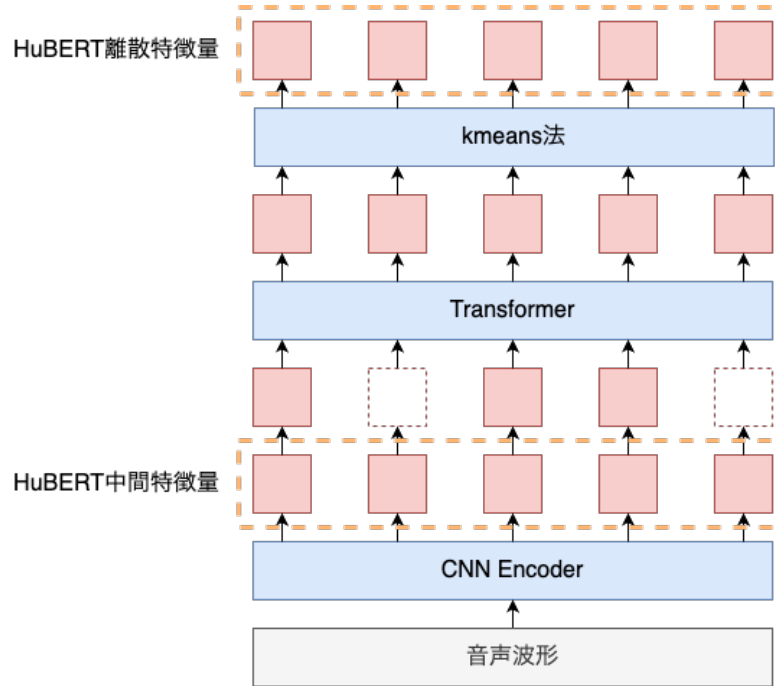


図 4.2: HuBERT 中間特徴量と HuBERT 離散特徴量の取得位置

スク学習手法は、本研究におけるネットワーク A で、HuBERT 中間特徴量を推定しないものに当たる。AVHuBERT 以降の ConvDecoder について、これは先行研究と異なる構成である。先行研究では、メルスペクトログラムに対して畳み込みをベースとしたデコーダ、HuBERT 離散特徴量に対して全結合層を利用したデコーダを別々に設けており、話者 Embedding はメルスペクトログラムの推定のみを用いていた。この構成についても検討したが、ConvDecoder の方がベースライン、提案手法問わず全体的に性能が改善したため、本研究ではこちらを採用した。

以上のモデルにより、動画からメルスペクトログラムと HuBERT 離散特徴量が推定可能となる。その後、先行研究 [9] に基づく Multi-input Vocoder を用い、メルスペクトログラムと HuBERT 離散特徴量を入力として音声波形に変換することで、最終的な合成音声を得た。Multi-input Vocoder は HiFi-GAN [32] をベースとしたモデルであり、音声波形を生成する Generator と、Multi-Period Discriminator (MPD) および Multi-Scale Discriminator (MSD) という二つの Discriminator によって構成される。

Generator の構造を図 4.3 に示す。左の特徴量予測モデルは、動画からメルスペクトログラムと HuBERT 離散特徴量を予測する、本研究において主な検討対象となる部分を表す。Generator の内部構造について、まず、前処理層はメルスペクトログラムと HuBERT 離散特徴量を入力として受け取り、その後のレイヤーに入力するための形状に変換する役割を持つ。メルスペクトログラムに対しては、時間方向に隣接したフレームを次元方向に縦積みすることによって、100 Hz・80 次元の特徴量から 50 Hz・160 次元の特徴量に変換した後、全結合層によって 128 次元の特徴量に変換する。一方、HuBERT 離散特徴量は 50 Hz のインデックス系列であり、インデックスから 128 次元のベクトルへと変換する。その後、これらをチャンネル方向に結合することで 256 次元の特徴量を構成し、これをその後のレイヤーへの入力とする。この特徴量は、転置

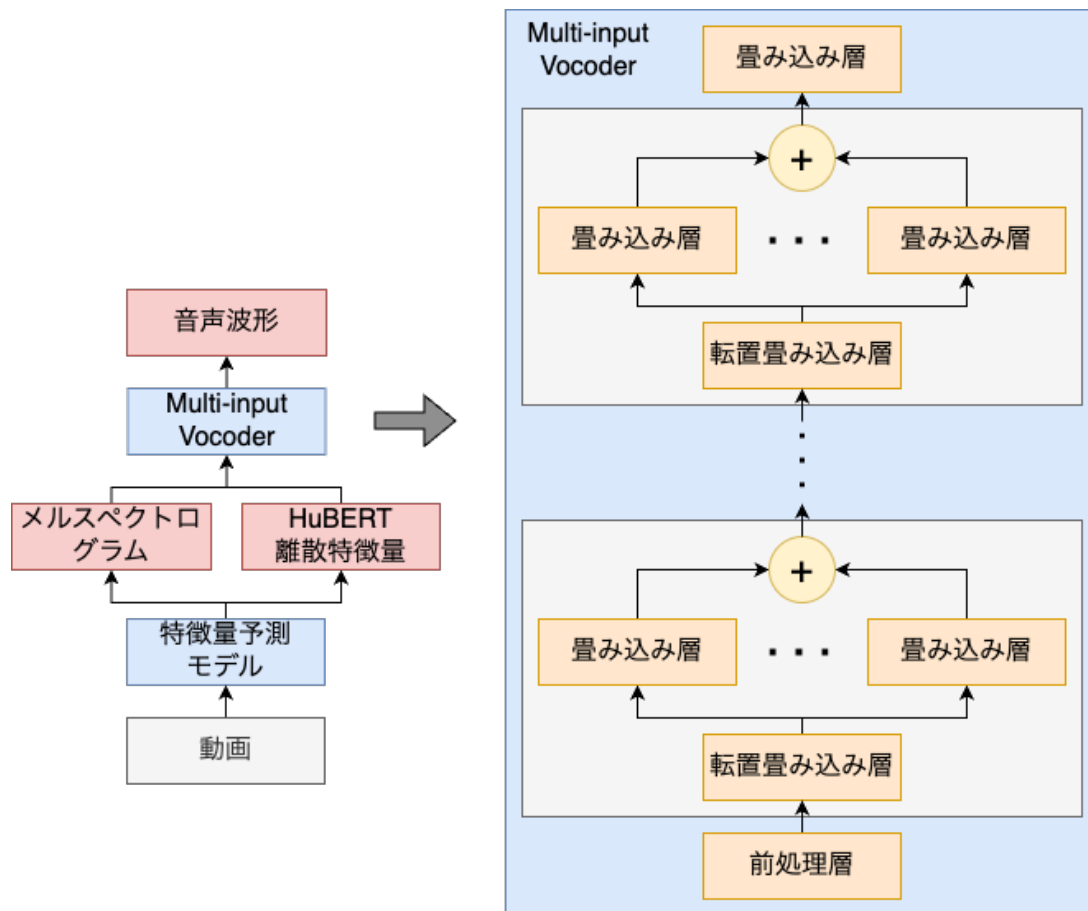


図 4.3: Multi-input Vocoder の構造

畳み込み層と複数種類の畳み込み層から構成されるブロックを通過していく。各ブロックについて、まず、転置畳み込み層は特徴量を時間方向にアップサンプリングする役割を果たす。実際、本研究では 50 Hz の入力特徴量から 16 kHz の音声波形まで、時間方向に 320 倍のアップサンプリングを行う必要がある。Generator では、これを複数のブロックを通して段階的に行っている。また、各ブロックの転置畳み込み層に積まれた複数種類の畳み込み層は、そのカーネルサイズとダイレーションが全て異なっている。複数の時間的な受容野を持つ畳み込み層からの出力をすべて加算することで、アップサンプリング後の特徴量からの特徴抽出を行う仕組みとなっている。表 4.1 に、Generator の各ブロックにおけるパラメータとブロックごとの出力特徴量の形状を示す。カラム名の右にある括弧がきが数値の意味を表しており、K はカーネルサイズ、S はストライド、D はダイレーション、C は次元（チャンネル数）、T は系列長である。畳み込み層についてはカーネルサイズとダイレーションを集合として表記しているが、実際はこれらの直積の元、すなわち (3, 1) や (3, 3)、(3, 5) をパラメータとする畳み込み層が存在することを表す。

Generator の学習時に用いられるのが、MPD および MSD という二つの Discriminator である。これらの概要を図 4.4 に示す。MPD では、一次元の音声波形を指定した周期をもとに Reshape することで二次元に変換し、これに対して二次元畳み込みを適用することで、入力された音声



表 4.1: Generator の各ブロックにおけるパラメータ

	転置畳み込み層 (K, S)	畳み込み層 (K, D)	出力特徴量の形状 (C, T)
1	(11, 5)	({3, 5, 7, 9, 11}, {1, 3, 5})	(1024, 250)
2	(8, 4)	({3, 5, 7, 9, 11}, {1, 3, 5})	(512, 1000)
3	(4, 2)	({3, 5, 7, 9, 11}, {1, 3, 5})	(256, 2000)
4	(4, 2)	({3, 5, 7, 9, 11}, {1, 3, 5})	(128, 4000)
5	(4, 2)	({3, 5, 7, 9, 11}, {1, 3, 5})	(64, 8000)
6	(4, 2)	({3, 5, 7, 9, 11}, {1, 3, 5})	(32, 16000)

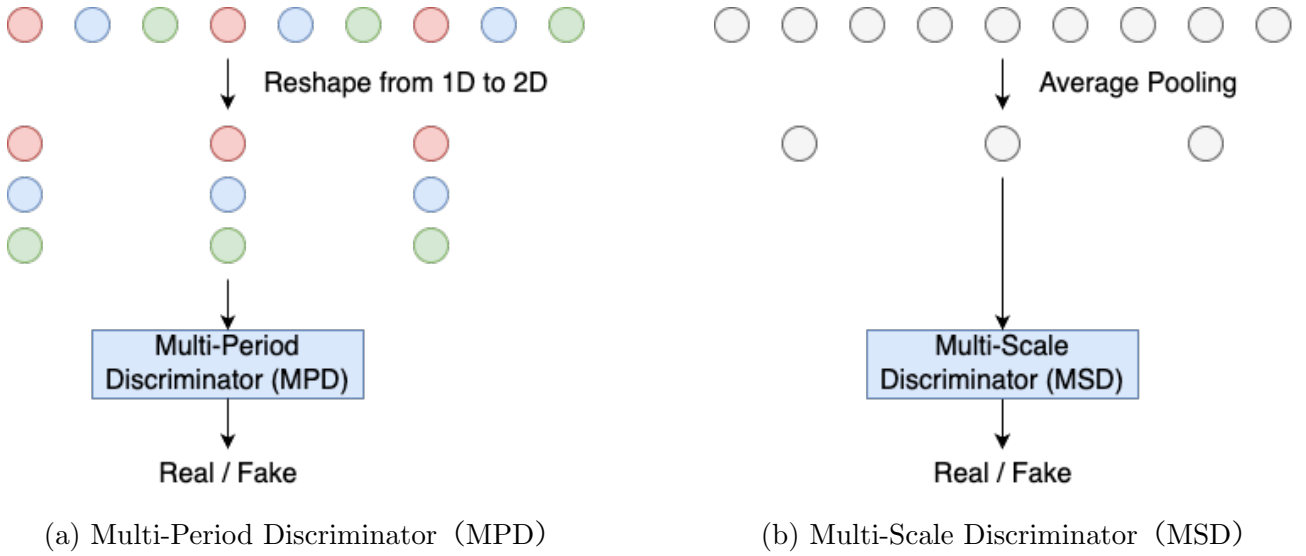


図 4.4: Multi-Period Discriminator (MPD) と Multi-Scale Discriminator (MSD) の概要

波形の原音声らしさを判定する Discriminator である。異なる周期を設定した Discriminator を複数利用することで、時間的な特徴を考慮できるように構成されている。一方、MSD では、音声波形に対して Average Pooling を適用することでダウンサンプリングし、これに一次元畳み込みを適用することで入力された音声波形の原音声らしさを判定する Discriminator である。MSD は Average Pooling によって系列長を抑えつつ、時間方向の連続的な特徴を考慮する狙いがある。これら二つの Discriminator については、HiFi-GAN と同様のパラメータで用いた。

## 4.2 実験方法

### 4.2.1 利用したデータセット

動画音声データセットには、男女二人ずつから収録した合計 4 人分のデータセット [25, 26] を用いた。これは ATR 音素バランス文 [33] から構成され、全話者共通で A から H セットを学習データ、I セットを検証データ、J セットをテストデータとして利用した。各分割ごとの文章数を表 4.2 に示す。

表 4.2: 利用したデータセットの文章数

	学習	検証	テスト
動画音声データセット	1598	200	212
Hi-Fi-Captain	37714	200	200
JVS	10398	1299	1300

Multi-input Vocoder の学習に利用する音声データセットには、Hi-Fi-Captain（日本語話者二名分）[34] と JVS（parallel100 と nonpara30）[35] を利用した。Hi-Fi-Captain は train-parallel および train-non-parallel を学習データ、val を検証データ、eval をテストデータとして分割した。各分割ごとの文章数を表 4.2 に示す。JVS には話者に対して 1 から 100 まで番号が割り振られており、本実験では 1 から 80 番の話者を学習データ、81 番から 90 番の話者を検証データ、91 番から 100 番までの話者をテストデータとした。各分割ごとの文章数を表 4.2 に示す。また、JVS には読み上げ音声の parallel100 および nonpara30 と、裏声の falset10、囁き声の whisper10 が含まれる。本研究では、parallel100 と nonpara30 のみを利用した。

#### 4.2.2 データの前処理

動画データは 60 FPS で収録されたものを ffmpeg により 25 FPS に変換して用いた。その後、手法 [36] により動画に対してランドマーク検出を適用した。このランドマークを利用することで口元のみを切り取り、画像サイズを (96, 96) にリサイズした。モデル入力時は動画をグレースケールに変換し、各フレームに対する正規化および標準化を適用した。正規化では、グレースケールが 0 から 255 までの値を取るため、最大値の 255 で割り、標準化では、AVHuBERT のプログラムで利用されていた平均値 0.421 と、標準偏差 0.165 をそのまま利用して、平均値を引いた後に標準偏差で割った。全体として、今回は事前学習済みの AVHuBERT の転移学習を行うため、そこでの前処理に合わせている。学習時は、ランダムクロップ、左右反転、Time Masking（一時停止）をデータ拡張として適用した。ランダムクロップは、(96, 96) で与えられる画像から (88, 88) をランダムに切り取る処理である。検証およびテスト時は、必ず画像中央を切り取るよう実装した。左右反転は、50% の確率で左右が反転されるよう実装した。Time Masking は、連続する画像の時間平均値を利用することによって、一時停止させるような効果を与えるデータ拡張手法である。動画 1 秒あたり 0 から 0.5 秒の間でランダムに停止区間を定め、その区間における動画の時間方向平均値を計算し、区間内のすべてのフレームをこの平均値で置換した。

音声データは 16 kHz にダウンサンプリングして用いた。それから、窓長 25 ms のハニング窓を用いて、シフト幅 10 ms で STFT を適用することでフレームレート 100 Hz のスペクトログラムに変換した。さらに、振幅スペクトログラムに対して 80 次のメルフィルタバンクを適用し、メルスペクトログラムを得た上で対数スケールに変換した。また、Multi-input Vocoder の学習に利用した Hi-Fi-Captain と JVS については、無音区間のトリミング (-40 dBFS 未満かつ

500 ms 継続する区間を 100 ms までカット) を適用した。なぜなら、Multi-input Vocoder の学習時は、全体から 1 秒分をランダムサンプリングするように実装しており、元データに存在する無音区間を除去することによって、学習の安定化および得られる音声の品質改善に繋がったからである。また、話者 Embedding の取得には事前学習済みの話者識別モデル [31] を利用した。動画音声データセット、Hi-Fi-Captain、JVS とともに、各話者学習データの中から 100 文章をランダムサンプリングし、各発話に対して得られたベクトルの平均値を用いた。この値を学習・検証・テストで一貫して用いるため、学習外の検証データやテストデータには非依存な値となっている。モデルに入力する際には、ベクトルの大きさを割って正規化した。

HuBERT は、HuggingFace に公開されている ReazonSpeech というデータセットによって学習されたモデル [37, 38] を利用した。ReazonSpeech は約 19000 時間の日本語音声からなるデータセットであり、日本語音声のコンテキストを大量のデータから学習したモデルである。今回用いるデータセットが日本語であることから、本研究の検討対象としては日本語音声に関する事前知識を豊富に有するモデルが適していると考え、このモデルを選択した。本研究で用いる HuBERT 中間特徴量および HuBERT 離散特徴量について、HuBERT 中間特徴量は、音声波形に対して畳み込み層を適用した出力で、HuBERT の事前学習時にはマスクの対象となる特徴量を利用した。一方、HuBERT 離散特徴量は、HuBERT Transformer 層の 8 層目出力に、k-means 法によるクラスタリングを適用することで得た。あえて 8 層目出力を選択した理由は、HuBERT のレイヤーごとの特徴量について、音素の One-hot ベクトルおよび単語の One-hot ベクトルとの相関を、Canonical Correlation Analysis (CCA) によって調べた先行研究 [39] より、8 層目出力がそのどちらとも相関が高く言語的な情報に近いと判断したからである。クラスタ数については決め打ちとなるが、今回は離散化した結果が言語的な情報を持つことを狙い、比較的少ない 100 とした。k-means 法の学習には動画音声データセットにおける学習用データ全てを利用し、これを用いて動画音声データ、Hi-Fi-Captain、JVS に対するクラスタリングを実施した。

#### 4.2.3 学習方法

一段階目について、損失関数はメルスペクトログラムの MAE Loss  $L_{mel}$  と HuBERT 離散特徴量の Cross Entropy Loss  $L_{ssl^d}$ 、HuBERT 中間特徴量の MAE Loss  $L_{ssl^i}$  の重み付け和とした。それぞれの重み係数を  $\lambda_{mel}$ ,  $\lambda_{ssl^d}$ ,  $\lambda_{ssl^i}$  とすると、

$$L = \lambda_{mel} * L_{mel} + \lambda_{ssl^d} * L_{ssl^d} + \lambda_{ssl^i} * L_{ssl^i} \quad (4.1)$$

となる。最適化手法には AdamW [40] を利用し、 $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ ,  $\lambda = 0.01$  とした。学習率に対するスケジューラには、Cosine Annealing with Warmup を利用した。開始時の学習率は  $1.0 \times 10^{-6}$  として、最大エポック数の 10% に至るまでは学習率を  $1.0 \times 10^{-3}$  まで線形に増加させ、その後のエポックでは cosine 関数に基づいて  $10 \times 10^{-6}$  まで減少させた。バッチサイズはメモリの都合上 4 としたが、学習の安定化のため、Gradient Accumulation によって各イテレーションにおける勾配を累積させ、8 イテレーションに一回重みを更新するようにした。モ

デルに入力する動画の秒数は10秒を上限とし、それを超える場合はランダムにトリミング、それに満たない場合はゼロパディングした。ゼロパディングした部分は損失の計算からは除外した。勾配のノルムは3.0を上限としてクリッピングすることで、過度に大きくなることを防止した。最大エポック数は50とし、10エポック連続して検証データに対する損失が小さくならない場合には、学習を中断するようにした (Early Stopping)。また、学習終了時には検証データに対する損失が最も小さかったエポックにおけるチェックポイントを保存し、これをテストデータに対する評価に用いた。

第二段階について、損失関数はメルスペクトログラムの MAE Loss と HuBERT 離散特徴量の Cross Entropy Loss の重み付け和とした。これは式 (4.1) において、 $\lambda_{mel} = 1.0$ ,  $\lambda_{ssl} = 0.0$  と固定した場合に相当する。最適化手法には AdamW を利用し、 $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ ,  $\lambda = 0.01$  とした。学習率に対するスケジューラには、Cosine Annealing with Warmup を利用した。開始時の学習率は  $1.0 \times 10^{-6}$  として、最大エポック数の 10% に至るまでは学習率を  $5.0 \times 10^{-4}$  まで線形に増加させ、その後のエポックでは cosine 関数に基づいて  $10 \times 10^{-6}$  まで減少させた。学習率の最大値は第一段階の  $1/2$  となっているが、これは値を半減させることによって学習を安定させることができたためである。その他のパラメータは、第一段階における値と同じである。

第三段階について、学習率のみ第一段階と同様に  $1.0 \times 10^{-3}$  としたが、それ以外は第二段階と同様である。

Multi-input Vocoder の学習では、Hi-Fi-Captain と JVS を用いた。はじめに Hi-Fi-Captain のみを用いて学習させ、その後学習済みモデルを JVS によって再学習した。損失関数は HiFi-GAN と同様である。Hi-Fi-Captain は男女一人ずつの文章数が豊富なデータセットであるため、高品質なモデルを構築可能であった。しかし、学習できる話者数が少ない分、学習外話者に対する合成音声の品質が低かった。そのため、一人当たりの文章数は 100 文章程度と少ないながらも、100 人分の話者からなる JVS を利用して再学習することによって、学習外話者に対する合成音声の品質を向上させた。最適化手法には AdamW を利用し、 $\beta_1 = 0.8$ ,  $\beta_2 = 0.99$ ,  $\lambda = 1.0 \times 10^{-5}$  とした。学習率は  $2.0 \times 10^{-4}$  から開始し、1 エポック経過するごとに 0.99 かけて徐々に減衰させた。バッチサイズは 16 とし、ここでは Gradient Accumulation は利用しなかった。モデルへの入力 は 1 秒を上限とし、それを超える場合はランダムにトリミング、それに満たない場合はゼロパディングした。勾配のノルムは 3.0 を上限としてクリッピングすることで、過度に大きくなることを防止した。最大エポック数は 30 とし、ここでは Early Stopping は適用しなかった。また、学習終了時には検証データに対する損失 (メルスペクトログラムに対する L1 Loss) が最も小さかったエポックにおけるチェックポイントを保存し、これをテストデータに対する評価に用いた。また、Multi-input Vocoder の提案された先行研究 [9] においては、学習時にあえてメルスペクトログラムにノイズをかけることによって、合成音声に対する汎化性能を向上させる学習方法が提案されている。本研究では、動画から推定されるメルスペクトログラムと HuBERT 離散特徴量の推定精度向上に焦点を当てたため、Multi-input Vocoder の学習は原音声から計算される特徴量そのもので行い、ボコーダ自体の汎化性能向上による精度改善は追求しなかった。

実装に用いた深層学習ライブラリは PyTorch および PyTorch Lightning である。GPU には NVIDIA RTX A4000 を利用し、計算の高速化のため Automatic Mixed Precision を適用した。

#### 4.2.4 比較手法

比較手法は、以下の五つである。

1. メルスペクトログラムと HuBERT 離散特徴量のマルチタスク学習手法（ベースライン）
2. 提案手法のネットワーク B で、事前学習済み重みを読み込まずに HuBERT Transformer 層を用いる手法
3. 提案手法のネットワーク C で、事前学習済み重みを読み込まずに HuBERT Transformer 層を用いる手法
4. 提案手法のネットワーク B で、事前学習済み重みを読み込んで HuBERT Transformer 層を用いる手法
5. 提案手法のネットワーク C で、事前学習済み重みを読み込んで HuBERT Transformer 層を用いる手法

手法 1 が先行研究において有効性が確認された手法であり、今回の実験においてベースラインとなる。これに対する改善案として、手法 2 から手法 5 が提案手法である。手法 2 と手法 3 及び手法 4 と手法 5 の違いは事前学習済み重みを読み込むか否かであり、これらを比較することで初期値を事前学習済み重みとして転移学習することの有効性を調べた。

#### 4.2.5 客観評価

合成音声の客観評価では、二種類の指標を用いた。一つ目は、音声認識の結果から算出した Word Error Rate (WER) である。音声認識には Whisper [41] を利用し、出力される漢字仮名交じり文に対して MeCab を用いて分かち書きを行った上で、jiwer というライブラリを用いて算出した。Whisper は Large モデルを利用し、MeCab の辞書には unidic を利用した。WER の値は 0% から 100% であり、この値が低いほど音声認識の誤りが少ないため、より聞き取りやすい音声であると判断した。二つ目は、話者 Embedding から計算したコサイン類似度である。話者 Embedding の計算は、モデルへの入力値を計算したものと同様の話者識別モデルを利用し、対象音声と原音声のペアでコサイン類似度を計算した。今回構築するモデルは 4 人の話者に対応するモデルとなるため、原音声に似た声質の合成音声を得られているかをこの指標で評価した。値は 0 から 1 であり、高いほど原音声と類似した合成音声だと判断できる。

## 4.3 結果

### 4.3.1 客観評価

まず、損失関数 (4.1) の重み係数  $\lambda_{ssld}$  を変化させた時の、客観評価指標の全テストデータに渡る平均値を表 4.3 に示す。各手法ごとに 0.0001 から 1.0 まで 10 倍刻みで 5 段階検討し、各手法の客観指標ごとに最も優れた値を下線で示している。最良エポックは検証データに対する損失が最小となったエポックであり、テストデータの合成にはこのエポックにおけるチェックポイントを利用した。また、 $L_{mel}$ 、 $L_{ssld}$ 、 $L$  は最良エポックにおける検証データに対する損失の平均値である。また、これ以降の比較のために、最適だと考えられる  $\lambda_{ssld}$  の値を選択しており、選択された行を太字で表している。

手法 1 では、 $\lambda_{ssld}$  の値が 0.0001 の時に WER が最も低く、0.01 の時に話者類似度が最も高くなった。現状 WER の高さが特に課題であり、話者類似度は 0.004 とわずかな違いでもあるため、今回は WER が最小であることを優先して 0.0001 が最適であると判断した。 $\lambda_{ssld}$  の値による評価指標の変化について、WER は  $\lambda_{ssld}$  の値を大きくするのに伴って単調に増加していることがわかる。一方、話者類似度は  $\lambda_{ssld}$  の値が 0.1 以上となった時に、0.01 以下であった場合と比較して顕著に低下することがわかる。次に、図 4.5 に手法 1 における学習曲線の結果を示す。横軸がエポック数、縦軸が損失の値を表す。損失の値は各エポックにおける平均値である。実線は検証データに対する損失、点線は学習データに対する損失を表しており、線の色は  $\lambda_{ssld}$  の違いを表す。また、丸いマーカーは表 4.4 に示した最良エポック時における損失の値を表す。学習曲線より、 $\lambda_{ssld}$  の値を変化させることによって、特に  $L_{ssld}$  の傾向が変化していることがわかる。具体的には、 $\lambda_{ssld}$  の値を 0.0001 から 1.0 へと増加させるのに伴って、学習初期における損失の上がり方が急峻になっており、達する最小値自体が小さくなっていることがわかる。また、 $\lambda_{ssld}$  の値が 0.1 以上の場合、検証データに対する  $L_{ssld}$  は早いうちから増加傾向に転じている。これに伴い、今回は検証データに対する  $L$  の値を監視し、Early Stopping の適用と最良エポックの決定を行なったため、 $L_{mel}$  が下がり切らない状態で学習が中断される結果となった。客観評価指標において、特に  $\lambda_{ssld}$  の値が 0.1 以上となると、0.01 以下の場合と比較して話者類似度の低下や、WER の上昇傾向が見られていたが、学習曲線の挙動より、 $L_{mel}$  を下げきれなくなっていたことが原因として考えられる。最適な  $\lambda_{ssld}$  の値は客観評価指標から 0.0001 としたが、0.01 以下ではそれと概ね同程度の品質であったことを考えると、手法 1 では検証データに対する  $L_{mel}$  の値を十分小さくすることのできる  $\lambda_{ssld}$  が適していると考えられる。

手法 2 では、 $\lambda_{ssld}$  の値が 0.1 の時に WER が最も低く、0.0001 の時に話者類似度が最も高くなった。ここでは  $\lambda_{ssld}$  の値が 0.1 の時に、ベースラインで選択された最適なケースと比較して WER が 9.1% 低下しており、話者類似度についても 0.003 高くなっていることから、0.1 が最適だと判断した。 $\lambda_{ssld}$  の値による評価指標の変化について、 $\lambda_{ssld}$  を 0.1 以上とすることで、0.01 以下の場合と比較して WER が低下する傾向が見られた。しかし、1.0 まで大きくすると 0.1 の場合よりも WER が大きくなっており、単調に減少していないこともわかる。話者類似度については、 $\lambda_{ssld}$  の値を 0.1 としたときに、0.01 以下の場合と比較して値が少し低下し、さらに 1.0 ま

で大きくすることで0.1以下の場合と比較して顕著に低下していることがわかる。次に、図 4.6 に手法2における学習曲線の結果を示す。手法1と同様に、 $\lambda_{ssl^d}$  の値を0.0001から1.0へと増加させるのに伴って、学習初期における  $L_{ssl^d}$  の下がり方が急峻になっており、達する最小値自体が小さくなっていることがわかる。また、 $\lambda_{ssl^d}$  が1.0の場合に、 $L_{mel}$  を下げきれなくなる傾向が見られる。加えて、手法2では  $\lambda_{ssl^d}$  が0.1の場合における  $L_{ssl^d}$  の増加が緩やかであり、 $L_{mel}$  も十分下げられていることがわかる（ただし、表 4.3 の  $L_{ssl^d}$  より、手法1の  $\lambda_{ssl^d}$  が0.1の場合と比較して、損失の値自体は大きい）。さらに、 $\lambda_{ssl^d}$  が0.1の場合、 $L_{mel}$  が達する最小値自体が、 $\lambda_{ssl^d}$  が0.01以下の場合と比較して小さくなっていることがわかる（具体的な値は表 4.3 の  $L_{mel}$  の値を参照）。これは手法1では見られなかった新たな傾向であった。最適な  $\lambda_{ssl^d}$  の値は客観評価指標から0.1としたが、学習曲線の挙動より、この時他の値の場合と比較して  $L_{mel}$  と  $L_{ssl^d}$  の両方をバランスよく下げられていたことがわかった。よって、手法2においては  $L_{mel}$  と  $L_{ssl^d}$  の両方をバランスよく下げられるような、程よい大きさの重みが適していると考えられる。

手法3における最適値の選択理由は手法2と同様であり、 $\lambda_{ssl^d}$  の値による評価指標の変化についても同様であった。次に、図 4.7 に手法3における学習曲線の結果を示す。傾向は手法2と概ね同様であるが、手法3では  $\lambda_{ssl^d}$  の値が0.1の場合における  $L_{ssl^d}$  の増加が早く、学習が早期に中断されている点は異なる。しかし、この時  $L_{mel}$  も十分下げられており、 $L_{mel}$  が達する最小値自体が、 $\lambda_{ssl^d}$  が0.01以下の場合と比較して小さくなっていることがわかる（具体的な値は表 4.3 の  $L_{mel}$  の値を参照）。以上より、細かな違いはあるものの手法3は手法2と同様の挙動を示し、 $\lambda_{ssl^d}$  は  $L_{mel}$  と  $L_{ssl^d}$  の両方をバランスよく下げられるような、程よい大きさの重みが適していると考えられる。

手法4及び手法5についても、最適値の選択理由は手法2と同様である。しかし、これらは  $\lambda_{ssl^d}$  の値を0.1としたときの話者類似度の低下の度合いが手法2、手法3と比較して大きい点で、傾向が異なっていた。次に、図 4.8 に手法4、図 4.9 に手法5における学習曲線の結果を示す。傾向については、 $L_{ssl^d}$  の増加の程度について細かな違いはあるが、概ね手法2と同様であり、 $L_{mel}$  と  $L_{ssl^d}$  の両方をバランスよく下げられるような、程よい大きさの重みが適していると考えられる。また、手法4と手法5において、 $\lambda_{ssl^d}$  が0.1の場合における話者類似度は、手法2および手法3と比較して低くなっていたが、 $L_{mel}$  や  $L_{ssl^d}$  の値からは一貫した傾向が見られない。これより、平均値からこの違いを分析することは難しいが、実際の細かな誤差の出方が異なっていて、これが影響を与えているのではないかと考える。

次に、最適なチューニングをした場合における、手法ごとの客観評価指標の全テストデータに渡る平均値を表 4.4 に示す。分析合成は、原音声から計算した特徴量を入力として、Multi-input Vocoder で逆変換した合成音声であり、本実験下において合成音声により達成され得る上限値を表す。手法1から手法5については、表 4.3 において太字としたもの、すなわち最適なチューニングだと判断されたものを選択している。また、分析合成、原音声を除いた合成音声（手法1から手法5）の中で、最も優れた値を下線で示している。これより、WER、話者類似度ともに手法3が最も優れていることが分かる。特に、今回のベースラインである手法1と比較すると、

WER が 9.5% 低下し、話者類似度は 0.01 高くなっていることから、提案手法がベースラインよりもより聞き取りやすく、話者性を反映した音声を合成できたと考えられる。また、HuBERT の事前学習済み重みを初期値とすることの効果について、手法 2 と手法 4 を比較すると、手法 2 の方が WER が 1.2% 低く、話者類似度が 0.017 高いことがわかる。よって、本研究では HuBERT の事前学習済み重みを初期値とした転移学習が有効である可能性を検討したが、むしろ事前学習済み重みを初期値とせず、ランダム初期化したモデルの方が優れていたと考えられる。これは仮説に反した結果であったが、事前学習済み重みによって与えられる初期値が、より良い局所解への収束には繋がらなかったことが原因だと考えられる。一方、アンサンブル手法の有効性について、手法 2 と手法 3 を比較すると、アンサンブル手法を導入することによって WER が 0.4% 低下し、話者類似度が 0.007 高くなっていることから、わずかではあるが改善していることがわかる。しかし、手法 4 と手法 5 を比較すると、WER は変化せず、話者類似度は 0.017 低下していることから、特に話者類似度の観点で悪化したことがわかる。よって、アンサンブル手法は必ずしも改善につながるわけではなく、改善するとしても顕著な変化はもたらさないと考えられる。

次に、最適なチューニングをした場合における手法ごとに、発話文章ごとの WER の比較を行った結果を図 4.10 に示す。ここで、縦軸は発話文章を表し、横軸はベースラインである手法 1 とその他手法の間で WER の差を計算し、発話文章ごとに平均した値を表す。負の値を取っているとき、手法 1 に対してより低い WER を達成したと解釈できる。図より、表 4.4 における平均値のみを見れば、手法 2 から手法 5 の全てが手法 1 に対してより低い WER を達成していたが、発話文章ごとに見れば、手法 1 に対して WER が高くなっているサンプルもあることがわかる。例えば、「ATR503\_j11」では、手法 2 から手法 5 の全てがベースラインよりも 20% 前後高い WER となっていることがわかる。また、全体的な傾向として手法ごとに WER はばらけており、いかなる場合においても最良となるようなモデルは構築できていないことがわかる。これより、現状のいかなるモデルも発話文章に対する汎化性能が不十分であり、さらなる改善が必要だと考えられる。

次に、話者ごとの WER の比較を行った結果を図 4.11 に示す。ここで縦軸は話者を表し、横軸はベースラインである手法 1 とその他手法の間で WER の差を計算し、話者ごとに平均した値を表す。負の値を取っているとき、手法 1 に対してより低い WER を達成したと解釈できる。これより、手法 2 から手法 5 は平均的にいかなる話者に対しても手法 1 より低い WER を達成していることがわかるが、一方でその改善の度合いは話者によって異なることがわかる。特に、「F02\_kablab」はその他三人の話者と比較して改善の度合いが小さい。また、例えば「F01\_kablab」では手法 2 と手法 3 がより有効である一方、「F02\_kablab」では手法 4 と手法 5 がより有効となっており、有効な手法が話者によって異なることもわかる。これより、話者に対する性能の依存があると考えられ、さらなる汎化性能の向上が必要だと言える。

次に、話者ごとの話者類似度の比較を行った結果を図 4.12 に示す。ここで縦軸は話者を表し、横軸はベースラインである手法 1 とその他手法の間で話者類似度の差を計算し、話者ごとに平均した値を表す。正の値を取っているとき、手法 1 に対してより高い話者類似度を達成したと



解釈できる。これを見ると、手法4と手法5はいかなる話者に対しても話者類似度を悪化させたことがわかる。手法4と手法5は事前学習済み重みを初期値としたモデルであったため、この影響が出ている可能性がある。今回検討した HuBERT は音声の言語情報を考慮する点で強みがあり、実際音声認識ではその転移学習の有効性が確認されている。しかし、話者性を考慮するという点で特に有効でなく、話者性を軽視するような局所解への収束に繋がった可能性が考えられる。ただ、表 4.3 より手法4と手法5においても  $\lambda_{ssl^d}$  が 0.01 以下であれば、比較的高い話者類似度を達成していたため、損失の重み付けにも依存していると言える。また、手法3はいかなる話者に対しても類似度を向上させており、この点で優れていたと言える。手法2では「F02\_kablab」のみ話者類似度が悪化しており、手法3と比較して話者に対する汎化性能が低かったと考えられる。

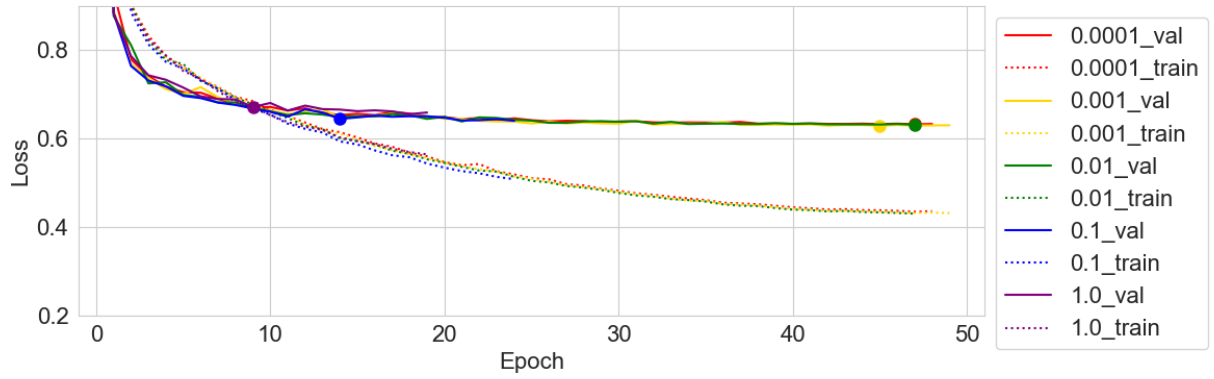
以上のことから、提案手法である手法2から手法5はいずれも手法1に対して平均的に WER を低下させ、特に手法2及び手法3は話者類似度についても若干の改善を達成したことから、提案手法によるベースラインからの改善が達成できたと考える。また、HuBERT の事前学習済み重みを用いることは仮説に反して有効ではなく、アンサンブル手法についてもその効果は顕著なものではなかった。さらに、特に手法2から手法5は損失関数の重み係数である  $\lambda_{ssl^d}$  による性能の変化が著しく、ここで最適値を選択できなければベースラインである手法1から改善しないこともわかった。よって、提案手法におけるベースラインからの改善に寄与した可能性のあるポイントは、以下でまとめられる。

1. HuBERT 中間特徴量を入力とし、メルスペクトログラムと HuBERT 離散特徴量を推定するネットワークを導入したこと
2. 上記のネットワーク構造を HuBERT Transformer としたこと
3. ネットワークをランダム初期化した上で学習させたこと
4. 最適な  $\lambda_{ssl^d}$  の値を発見できたこと

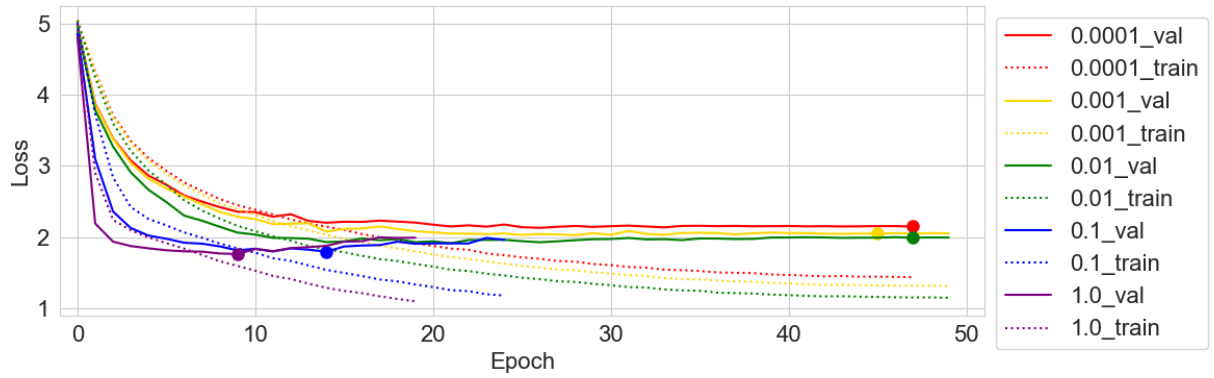
本実験からは、HuBERT Transformer をランダム初期化する必要性と、 $\lambda_{ssl^d}$  の値のチューニングを行う必要性が明らかとなった。しかし、入力特徴量やネットワーク構造については HuBERT を用いることを前提とした実験条件であったことから、検討できていない。実験結果から、もはや HuBERT に依存する必要性は無くなっているため、入力が HuBERT 中間特徴量でなくても良いし、ネットワークも任意に選択できる。入力については、例えば最終予測値であるメルスペクトログラムや HuBERT 離散特徴量を採用しても実装は可能であるが、これによって性能が変化するのであれば、HuBERT 中間特徴量が良い入力特徴量であると考えられる。ネットワークについては、Transformer 自体が多くの場合有効であるため現状で適切なものとなっている可能性もあるが、検討の余地はある。こういった点が、今後の検討課題として挙げられる。

表 4.3: 損失関数の重み係数  $\lambda_{ssl^d}$  による客観評価指標の比較

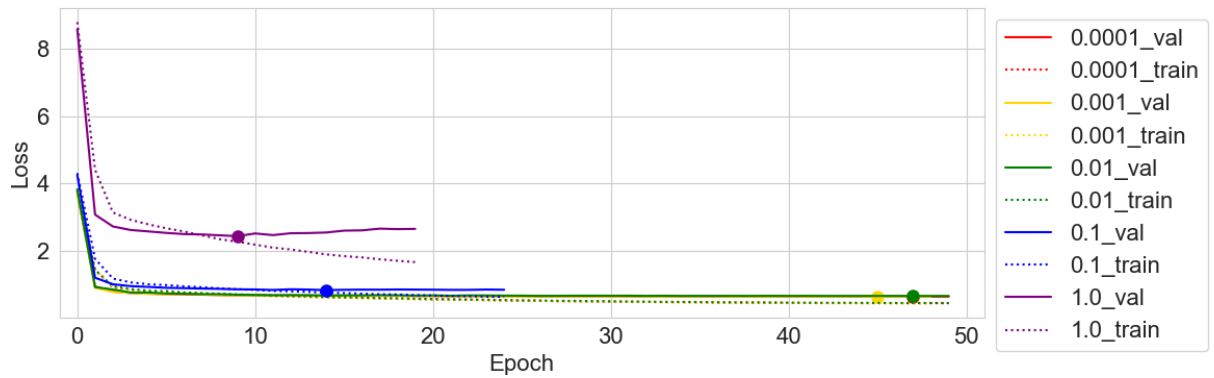
手法	$\lambda_{ssl^d}$	WER [%]	話者類似度	最良エポック	$L_{mel}$	$L_{ssl^d}$	$L$
<b>1</b>	<b>0.0001</b>	<b><u>55.3</u></b>	<b>0.841</b>	<b>47</b>	<b>0.632</b>	<b>2.147</b>	<b>0.633</b>
1	0.001	55.6	0.841	45	0.629	2.049	0.631
1	0.01	56.7	<u>0.845</u>	47	0.631	1.990	0.651
1	0.1	57.7	0.764	14	0.645	1.793	0.824
1	1.0	62.1	0.699	9	0.672	1.760	2.431
2	0.0001	58.9	<u>0.860</u>	46	0.630	2.521	0.630
2	0.001	58.3	0.857	45	0.631	2.382	0.634
2	0.01	56.7	0.859	48	0.631	2.104	0.652
<b>2</b>	<b>0.1</b>	<b><u>46.2</u></b>	<b>0.844</b>	<b>32</b>	<b>0.620</b>	<b>1.942</b>	<b>0.814</b>
2	1.0	50.3	0.706	9	0.651	1.713	2.364
3	0.0001	58.8	0.860	22	0.632	2.584	0.632
3	0.001	58.1	0.854	22	0.633	2.448	0.636
3	0.01	56.5	<u>0.865</u>	45	0.630	2.077	0.651
<b>3</b>	<b>0.1</b>	<b><u>45.8</u></b>	<b>0.851</b>	<b>18</b>	<b>0.622</b>	<b>2.041</b>	<b>0.826</b>
3	1.0	48.6	0.755	14	0.634	1.735	2.370
4	0.0001	58.2	0.849	20	0.632	2.593	0.632
4	0.001	57.4	0.847	24	0.633	2.496	0.635
4	0.01	56.0	<u>0.853</u>	30	0.630	2.106	0.651
<b>4</b>	<b>0.1</b>	<b><u>47.4</u></b>	<b>0.827</b>	<b>17</b>	<b>0.622</b>	<b>1.929</b>	<b>0.815</b>
4	1.0	48.4	0.720	10	0.648	1.746	2.394
5	0.0001	57.8	0.857	22	0.632	2.547	0.632
5	0.001	57.5	0.853	27	0.635	2.449	0.637
5	0.01	58.1	<u>0.865</u>	44	0.632	2.100	0.653
<b>5</b>	<b>0.1</b>	<b><u>47.4</u></b>	<b>0.810</b>	<b>6</b>	<b>0.633</b>	<b>1.981</b>	<b>0.831</b>
5	1.0	47.7	0.748	16	0.637	1.786	2.423



(a) メルスペクトログラムの MAE Loss (式 (4.1) の  $L_{mel}$ )

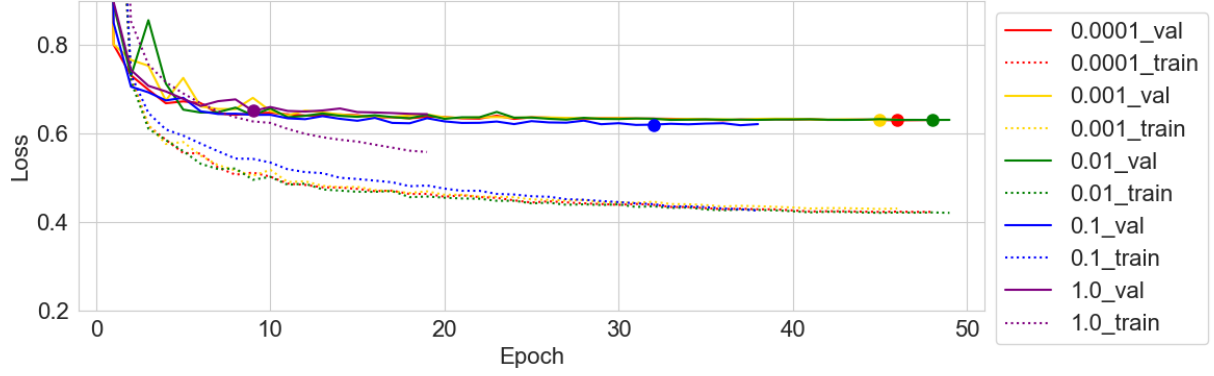


(b) HuBERT 離散特徴量の Cross Entropy Loss (式 (4.1) の  $L_{ssld}$ )

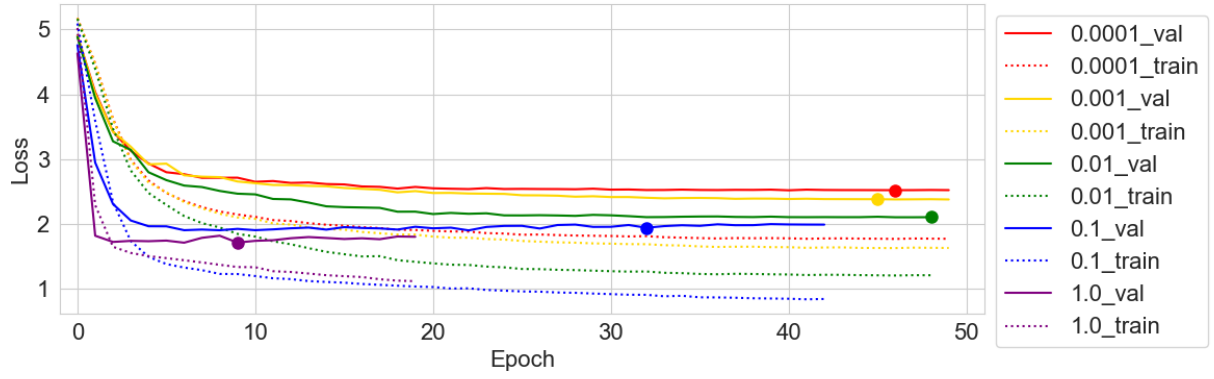


(c) 損失の合計値 (式 (4.1) の  $L$ )

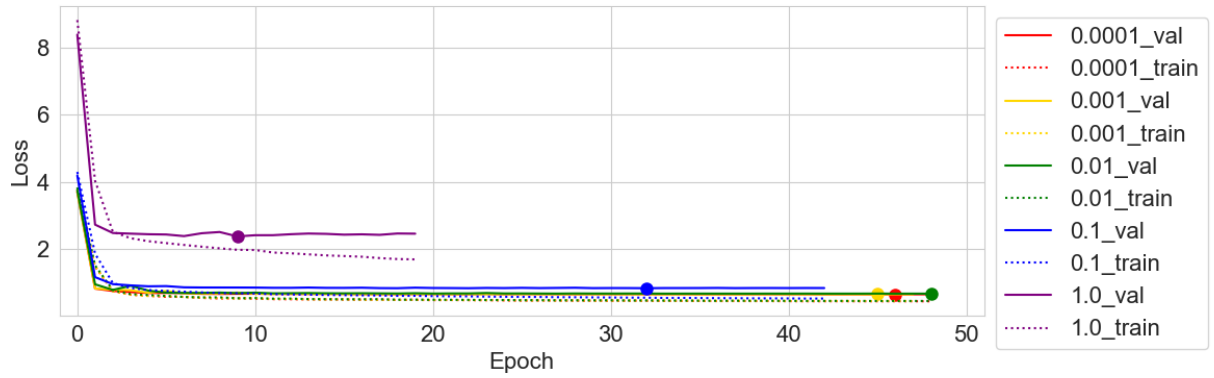
図 4.5: 手法 1 における学習曲線



(a) メルスペクトログラムの MAE Loss (式 (4.1) の  $L_{mel}$ )

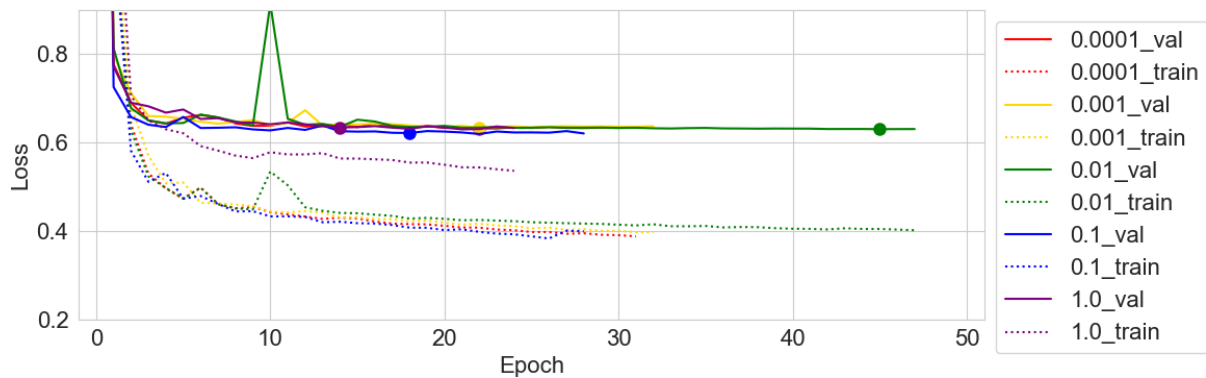


(b) HuBERT 離散特徴量の Cross Entropy Loss (式 (4.1) の  $L_{ssld}$ )

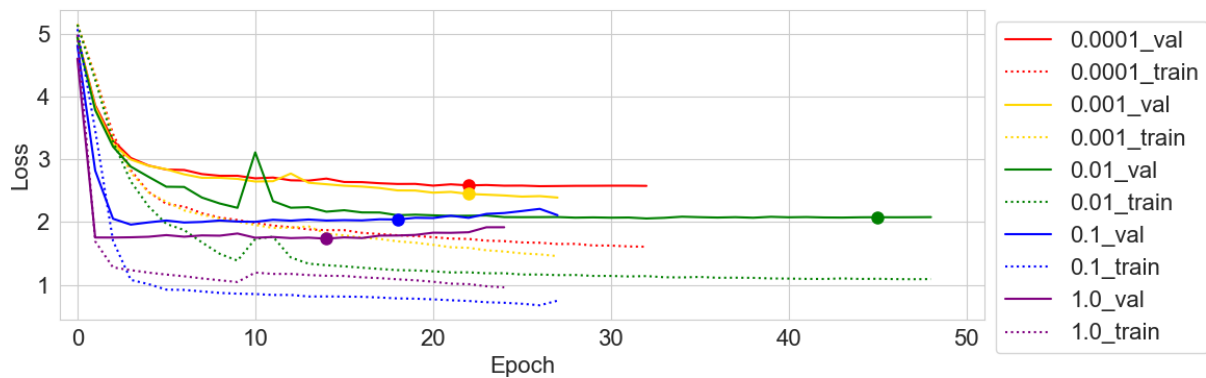


(c) 損失の合計値 (式 (4.1) の  $L$ )

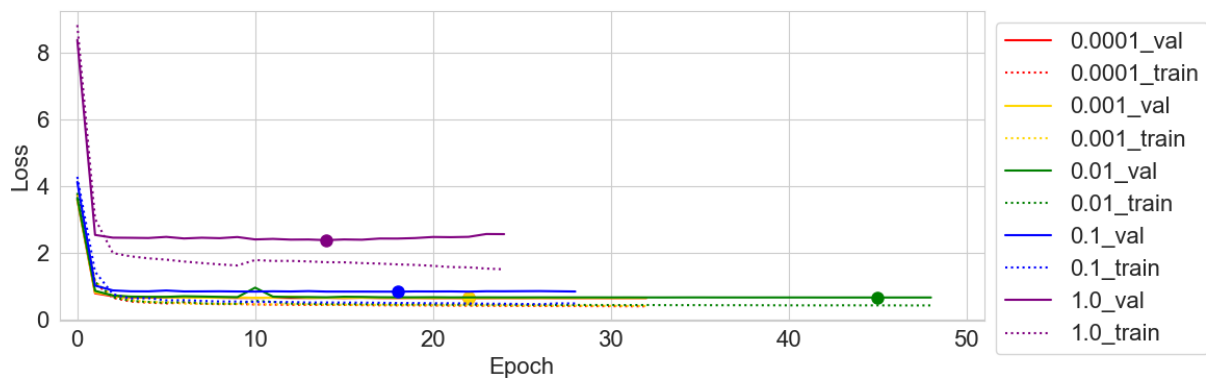
図 4.6: 手法 2 における学習曲線



(a) メルスペクトログラムの MAE Loss (式 (4.1) の  $L_{mel}$ )

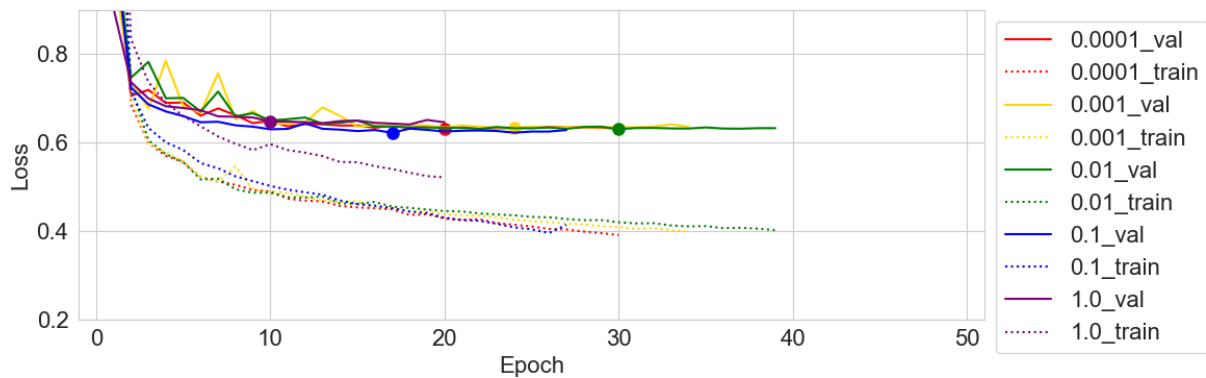


(b) HuBERT 離散特徴量の Cross Entropy Loss (式 (4.1) の  $L_{ssld}$ )

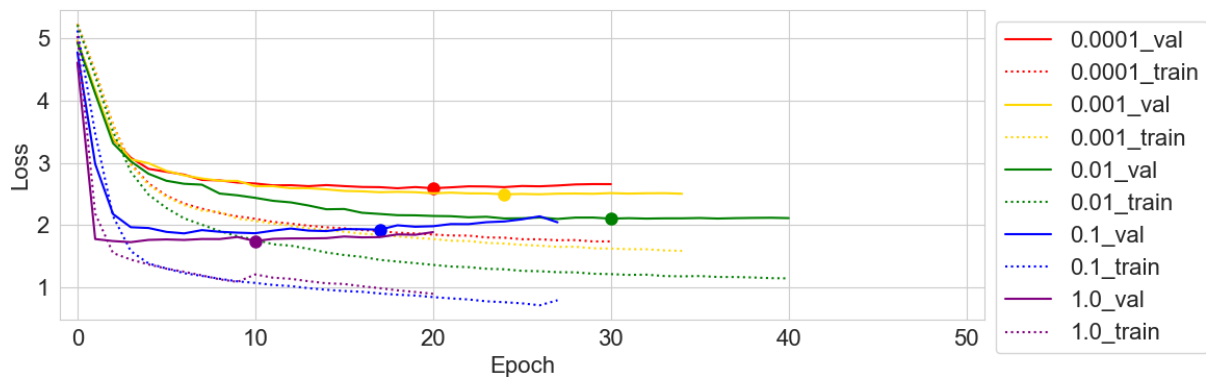


(c) 損失の合計値 (式 (4.1) の  $L$ )

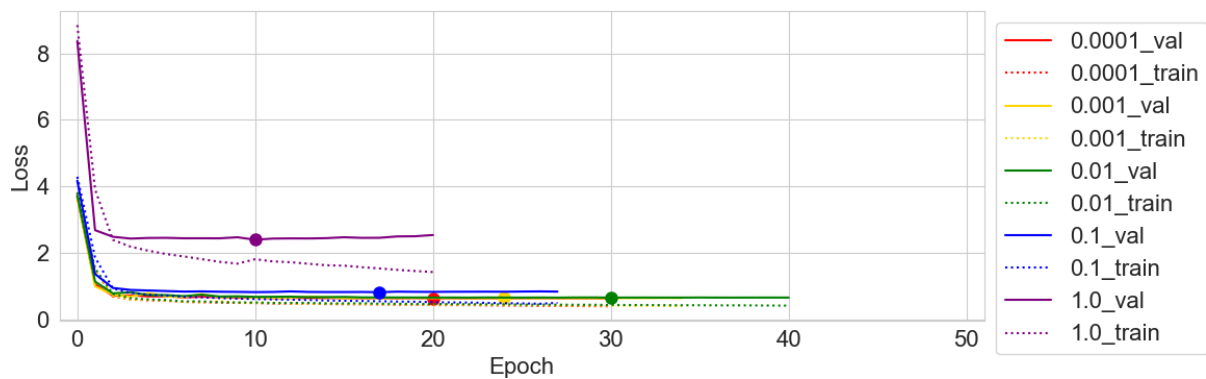
図 4.7: 手法3における学習曲線



(a) メルスペクトログラムの MAE Loss (式 (4.1) の  $L_{mel}$ )

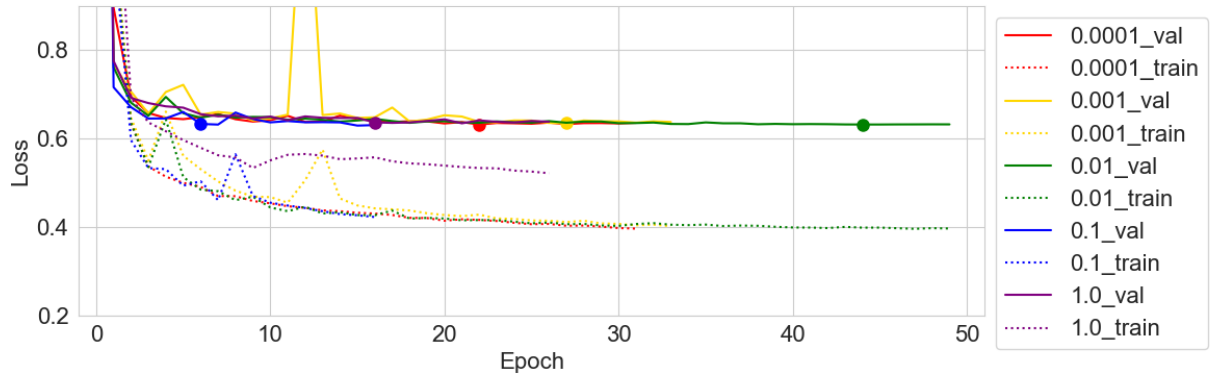


(b) HuBERT 離散特徴量の Cross Entropy Loss (式 (4.1) の  $L_{ssld}$ )

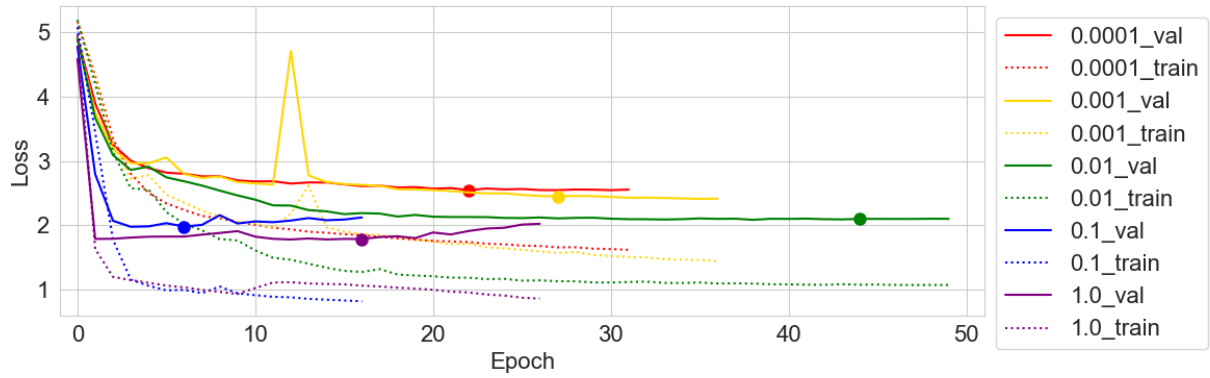


(c) 損失の合計値 (式 (4.1) の  $L$ )

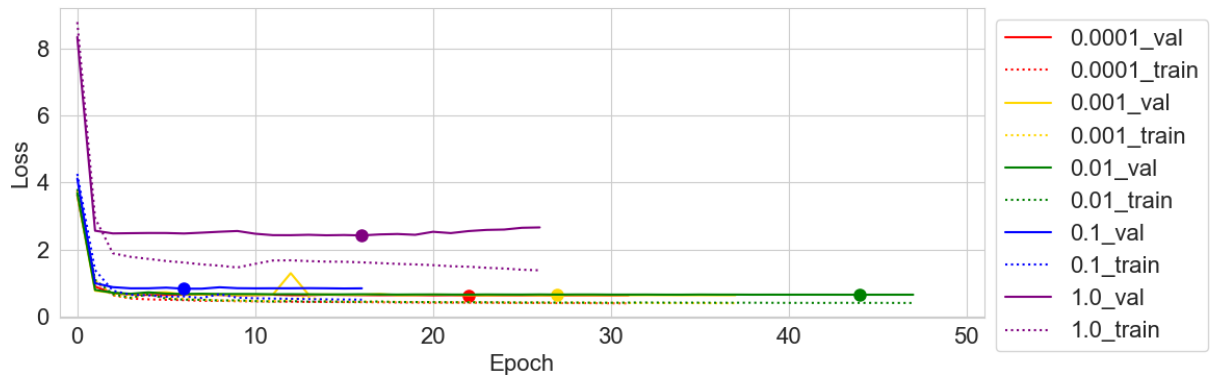
図 4.8: 手法 4 における学習曲線



(a) メルスペクトログラムの MAE Loss (式 (4.1) の  $L_{mel}$ )



(b) HuBERT 離散特徴量の Cross Entropy Loss (式 (4.1) の  $L_{ssld}$ )



(c) 損失の合計値 (式 (4.1) の  $L$ )

図 4.9: 手法 5 における学習曲線

表 4.4: 最適なチューニングをした場合における手法ごとの比較

手法	詳細	WER [%]	話者類似度
1	ベースライン	55.3	0.841
2	ランダム初期化 Transformer	46.2	0.844
3	ランダム初期化 Transformer + アンサンブル	<u>45.8</u>	<u>0.851</u>
4	事前学習あり Transformer	47.4	0.827
5	事前学習あり Transformer + アンサンブル	47.4	0.810
6	分析合成	4.8	0.944
7	原音声	4.5	1.000



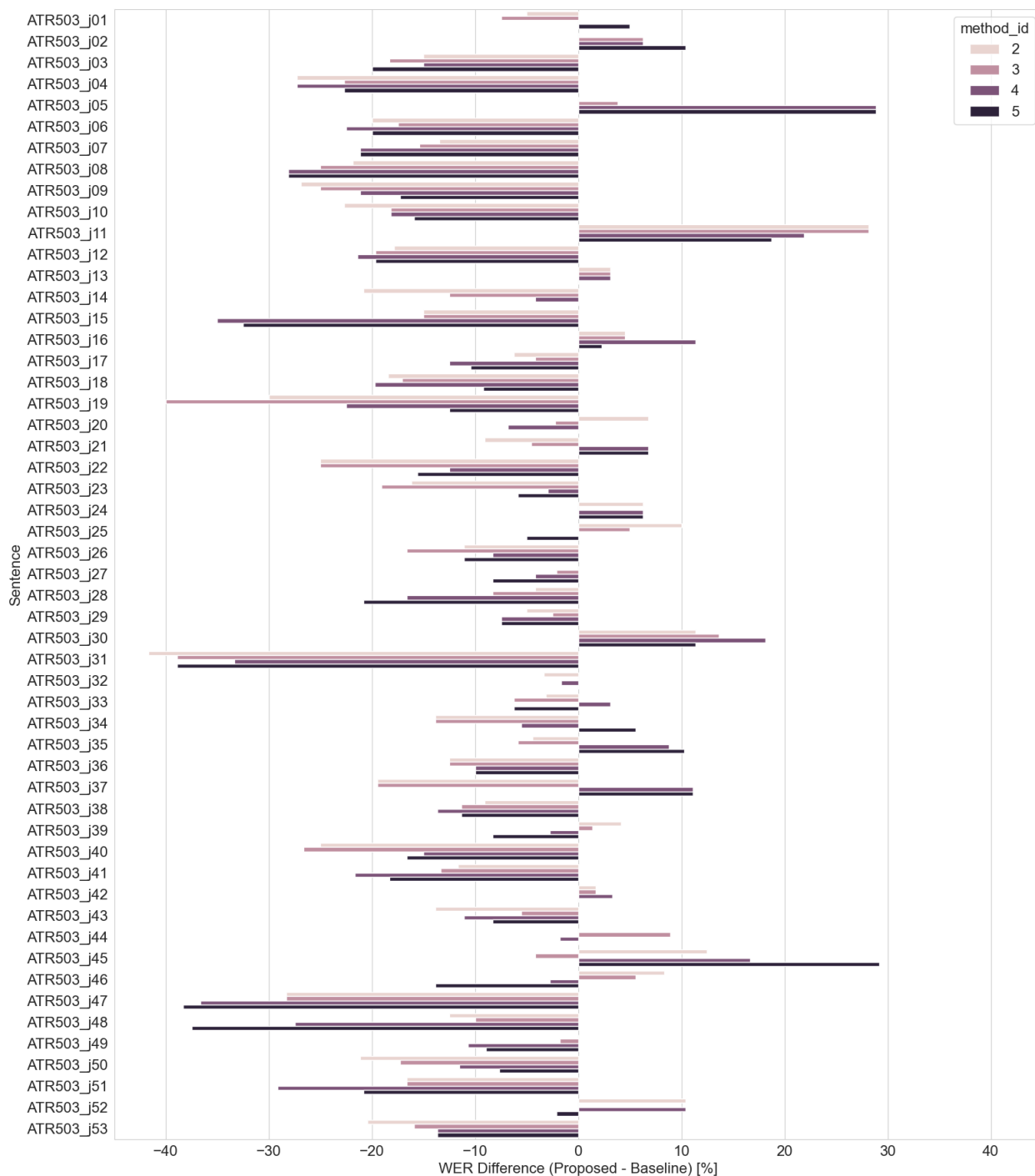


図 4.10: テストデータにおける発話文章ごとの WER の比較

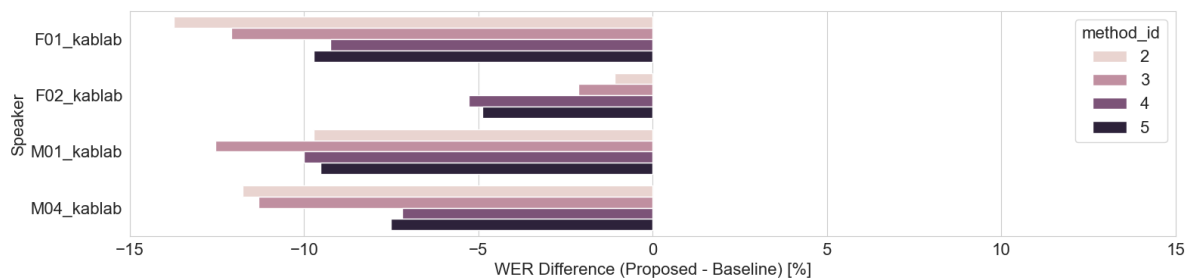


図 4.11: テストデータにおける話者ごとの WER の比較

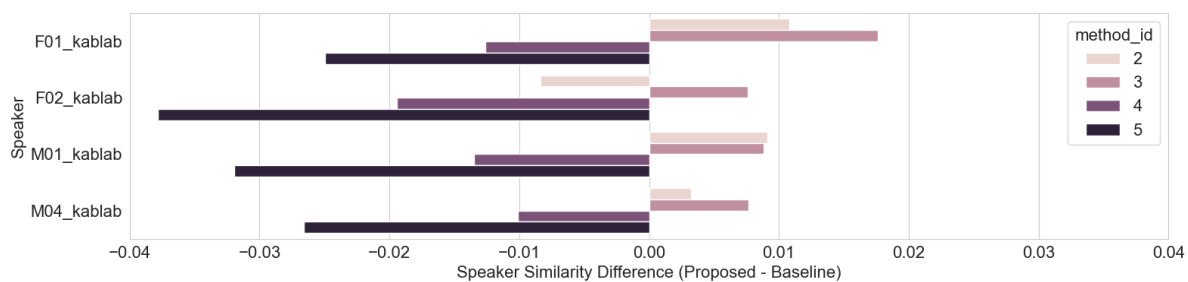


図 4.12: テストデータにおける話者ごとの話者類似度の比較

#### 4.3.2 主観評価

## 4.4 まとめ

## 5 結論

## 謝辭

## 参考文献

- [1] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Lrs3-ted: a large-scale dataset for visual speech recognition. *arXiv preprint arXiv:1809.00496*, 2018.
- [2] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*, 2018.
- [3] Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdelrahman Mohamed. Learning audio-visual speech representation by masked multimodal cluster prediction. *arXiv preprint arXiv:2201.02184*, 2022.
- [4] Qiushi Zhu, Long Zhou, Ziqiang Zhang, Shujie Liu, Binxing Jiao, Jie Zhang, Lirong Dai, Daxin Jiang, Jinyu Li, and Furu Wei. Vatlm: Visual-audio-text pre-training with unified masked prediction for speech representation learning. *IEEE Transactions on Multimedia*, 2023.
- [5] Alexandros Haliassos, Pingchuan Ma, Rodrigo Mira, Stavros Petridis, and Maja Pantic. Jointly learning visual and auditory speech representations from raw data. *arXiv preprint arXiv:2212.06246*, 2022.
- [6] Jiachen Lian, Alexei Baevski, Wei-Ning Hsu, and Michael Auli. Av-data2vec: Self-supervised learning of audio-visual speech representations with contextualized target representations. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 1–8. IEEE, 2023.
- [7] Alexandros Haliassos, Andreas Zinonos, Rodrigo Mira, Stavros Petridis, and Maja Pantic. Braven: Improving self-supervised pre-training for visual and auditory speech recognition. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 11431–11435. IEEE, 2024.
- [8] Minsu Kim, Joanna Hong, and Yong Man Ro. Lip-to-speech synthesis in the wild with multi-task learning. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- [9] Jeongsoo Choi, Minsu Kim, and Yong Man Ro. Intelligible lip-to-speech synthesis with speech units. *arXiv preprint arXiv:2305.19603*, 2023.
- [10] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM transactions on audio, speech, and language processing*, Vol. 29, pp. 3451–3460, 2021.

- [11] Wei-Ning Hsu, Tal Remez, Bowen Shi, Jacob Donley, and Yossi Adi. Revise: Self-supervised speech resynthesis with visual input for universal and generalized speech regeneration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18795–18805, 2023.
- [12] Neha Sahipjohn, Neil Shah, Vishal Tambrahalli, and Vineet Gandhi. Robustl2s: Speaker-specific lip-to-speech synthesis exploiting self-supervised representations. In *2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 1492–1499. IEEE, 2023.
- [13] Jeong Hun Yeo, Minsu Kim, Jeongsoo Choi, Dae Hoe Kim, and Yong Man Ro. Akvsr: Audio knowledge empowered visual speech recognition by compressing audio knowledge of a pretrained model. *IEEE Transactions on Multimedia*, 2024.
- [14] Xize Cheng, Tao Jin, Linjun Li, Wang Lin, Xinyu Duan, and Zhou Zhao. Opensr: Open-modality speech recognition via maintaining multi-modality alignment. *arXiv preprint arXiv:2306.06410*, 2023.
- [15] Yasser Abdelaziz Dahou Djilali, Sanath Narayan, Haithem Boussaid, Ebtessam Almazrouei, and Merouane Debbah. Lip2vec: Efficient and robust visual speech recognition via latent-to-latent visual to audio representation mapping. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 13790–13801, 2023.
- [16] Xubo Liu, Egor Lakomkin, Konstantinos Vougioukas, Pingchuan Ma, Honglie Chen, Ruiming Xie, Morrie Doulaty, Niko Moritz, Jachym Kolar, Stavros Petridis, et al. Synthvsr: Scaling up visual speech recognition with synthetic supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18806–18815, 2023.
- [17] Pingchuan Ma, Alexandros Haliassos, Adriana Fernandez-Lopez, Honglie Chen, Stavros Petridis, and Maja Pantic. Auto-avsr: Audio-visual speech recognition with automatic labels. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- [18] Oscar Chang, Hank Liao, Dmitriy Serdyuk, Ankit Shahy, and Olivier Siohan. Conformer is all you need for visual speech recognition. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 10136–10140. IEEE, 2024.
- [19] Andreas Zinonos, Alexandros Haliassos, Pingchuan Ma, Stavros Petridis, and Maja Pantic. Learning cross-lingual visual speech representations. In *ICASSP 2023-2023 IEEE*



- International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- [20] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. *arXiv preprint arXiv:1804.03619*, 2018.
  - [21] Elizabeth Salesky, Matthew Wiesner, Jacob Bremerman, Roldano Cattoni, Matteo Negri, Marco Turchi, Douglas W Oard, and Matt Post. The multilingual tedx corpus for speech recognition and translation. *arXiv preprint arXiv:2102.01757*, 2021.
  - [22] Ya Zhao, Rui Xu, and Mingli Song. A cascade sequence-to-sequence model for chinese mandarin lip reading. In *Proceedings of the 1st ACM International Conference on Multimedia in Asia*, pp. 1–6, 2019.
  - [23] Minsu Kim, Jeong Hun Yeo, Jeongsoo Choi, and Yong Man Ro. Lip reading for low-resource languages by learning and combining general speech knowledge and language-specific knowledge. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15359–15371, 2023.
  - [24] Jeong Hun Yeo, Minsu Kim, Shinji Watanabe, and Yong Man Ro. Visual speech recognition for low-resource languages with automatic labels from whisper model. *arXiv preprint arXiv:2309.08535*, 2023.
  - [25] 田口史郎. ”深層学習を用いたデータ駆動型調音・音声間変換に関する研究”. 九州大学大学院芸術工学府芸術工学専攻 博士論文, 2021.
  - [26] 江崎蓮. ”深層学習を用いた口唇動画・音声変換に関する調査”. 九州大学大学院芸術工学府芸術工学専攻 修士論文, 2022.
  - [27] Ji-Hoon Kim, Jaehun Kim, and Joon Son Chung. Let there be sound: Reconstructing high quality speech from silent videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38, pp. 2759–2767, 2024.
  - [28] S Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.
  - [29] Kyunghyun Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
  - [30] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

- [31] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno. Generalized end-to-end loss for speaker verification. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4879–4883. IEEE, 2018.
- [32] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in neural information processing systems*, Vol. 33, pp. 17022–17033, 2020.
- [33] Yoshinori Sagisaka, Kazuya Takeda, M Abel, Shigeru Katagiri, Tetsuo Umeda, and Hisao Kuwabara. A large-scale japanese speech database. In *ICSLP*, pp. 1089–1092, 1990.
- [34] T Okamoto, Y Shiga, and H Kawai. Hi-fi-captain: High-fidelity and high-capacity conversational speech synthesis corpus developed by nict, 2023.
- [35] Shinnosuke Takamichi, Kentaro Mitsui, Yuki Saito, Tomoki Koriyama, Naoko Tanji, and Hiroshi Saruwatari. Jvs corpus: free japanese multi-speaker voice corpus. *arXiv preprint arXiv:1908.06248*, 2019.
- [36] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Proceedings of the IEEE international conference on computer vision*, pp. 1021–1030, 2017.
- [37] Yukiya Hono, Kentaro Mitsui, and Kei Sawada. rinna/japanese-hubert-base.
- [38] Kei Sawada, Tianyu Zhao, Makoto Shing, Kentaro Mitsui, Akio Kaga, Yukiya Hono, Toshiaki Wakatsuki, and Koh Mitsuda. Release of pre-trained models for the Japanese language. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 13898–13905, 5 2024. Available at: <https://arxiv.org/abs/2404.01657>.
- [39] Ankita Pasad, Bowen Shi, and Karen Livescu. Comparative layer-wise analysis of self-supervised speech models. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [41] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.