

話者依存型 Recurrent Temporal Restricted Boltzmann Machine を用いた声質変換*

中鹿亘，滝口哲也，有木康雄（神戸大）

1 はじめに

近年，入力音声信号に含まれる音韻情報を残したまま，特定話者の声質のみを変換する技術（声質変換）の研究が盛んに行われている．この背景として，音声認識における話者性の制御による認識精度の向上や発話困難な障がい者の発話支援などへの応用が挙げられる [1, 2, 3]．これまで声質変換は GMM (Gaussian mixture model) を用いた手法 [4, 5] が広く用いられてきた．また近年では，スパース空間への線形射影に基づく手法として，非負値行列因子分解法 (non-negative matrix factorization; NMF) を用いた声質変換法が提案されている [6, 7]．しかしながら，これらの手法による変換式は線形変換をベースとしているため，入力-出力話者の特徴間の詳細な対応付けが困難であるなど，変換精度には限界があった．通常人間の声道形状は非線形的であり，得られる特徴量は単純な線形変換で置き換えられるとは考えにくい．つまり，非線形ベースの変形手法の方が音声特徴量の変換にはより適切であると考えられる．非線形ベースのアプローチの例として，Desai らによる多層 NN (neural network) を用いた声質変換法 [8] や，我々が提案してきた話者依存型 RBM (restricted Boltzmann machine [9]) もしくは DBN (deep belief network) [10] を用いた多層型声質変換法 [11]，Wu らによる CRBM (conditional restricted Boltzmann machine) を用いた非線形声質変換法 [12] が挙げられる．いずれの手法においても，非線形変換に基づくアプローチでは，線形変換ベースの手法と比べて比較的高い精度が得られていることが報告されている．

本稿では，我々の先行研究である話者依存型 RBM を用いた声質変換法 [8] を拡張して，音声信号に含まれる潜在的な時間的依存性を考慮したモデルを提案する．RBM は可視層と隠れ層からなる 2 層ネットワークであり，データの中に潜む潜在的な特徴量を抽出することができる．我々は，先行研究において，話者ごとに RBM を学習させ，その話者固有の潜在的な情報（特徴量）を浮き出させることで，話者性を強調させた潜在特徴量同士を変換する手法を提案してきた．音声信号は時系列データであるため，特徴量間の時間的な関係性が存在する．この性質を捉えることのできるモデルを用いれば，より高い精度でデータを表現

できるため，声質変換の精度にも良い影響を与えらる．本研究では，時系列データの中に潜む時間的依存関係を捉えるグラフィカルモデルとして，RTRBM (recurrent temporal restricted Boltzmann machine) [13] を用いる．RTRBM は RBM を拡張したモデルであり，前フレームの潜在特徴から現時刻の潜在特徴，可視特徴への接続が考慮されている．この接続重みが，潜在空間における時間依存関係を表す．先行研究と同様に，本研究では，特定話者のみの音声信号を用いて学習させた RTRBM を使って潜在特徴量を抽出し，入力話者-出力話者の潜在特徴量間は NN によって変換関数を求める．得られた出力話者の潜在特徴量は RTRBM の後方推論によって MFCC などの音響特徴量への逆変換が可能である．さらに，上記の変換手順（入力話者の音響特徴量から入力話者の潜在特徴への変換，入力話者の潜在特徴から出力話者の潜在特徴への変換，出力話者の潜在特徴から出力話者の音響特徴量への変換）は，一つの RNN (recurrent neural network) として表すこともでき，推定された，話者ごとの RTRBM，話者間の NN の重みを初期値として，入力信号に出力話者の音響特徴量，教師信号に出力話者の音響特徴量を与えて，RNN を再学習させることで，各パラメータの微調整 (fine-tuning) が可能である．

2 RBM ベースの確率モデル

提案手法では，RTRBM (recurrent temporal restricted Boltzmann machine) を用いて，潜在的な時間依存関係を捉えた高次元空間で特徴変換を行う．本節ではまず，基礎技術となる RBM (restricted Boltzmann machine) について述べ，続いて RTRBM について説明する．

2.1 RBM

RBM は Fig. 1(a) のような 2 層ネットワークであり，可視素子 v と隠れ素子 h の確率変数分布を表現する無向グラフィカルモデルである [9]．連続値の入力をサポートした改良型 GB (Gaussian-Bernoulli)-RBM [14] (以下，この改良型 RBM を単に RBM とする) では，連続値の可視素子 $v = [v_1, \dots, v_I]^T, v_i \in \mathbb{R}$ と隠れ素子 $h = [h_1, \dots, h_J]^T, h_j \in \{0, 1\}$ の同時確率

*Voice conversion using speaker-dependent recurrent temporal restricted Boltzmann machine. by Toru NAKASHIKA, Tetsuya TAKIGUCHI, Yasuo ARIKI (Kobe University)

$p(\mathbf{v}, \mathbf{h})$ は、以下のように表される。

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})} \quad (1)$$

$$E(\mathbf{v}, \mathbf{h}) = \left\| \frac{\mathbf{v} - \mathbf{b}}{2\sigma} \right\|^2 - \mathbf{c}^T \mathbf{h} - \left(\frac{\mathbf{v}}{\sigma^2} \right)^T \mathbf{W} \mathbf{h} \quad (2)$$

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (3)$$

ここで、 $\mathbf{W} \in \mathbb{R}^{I \times J}$ 、 $\sigma \in \mathbb{R}^{I \times 1}$ 、 $\mathbf{b} \in \mathbb{R}^{I \times 1}$ 、 $\mathbf{c} \in \mathbb{R}^{J \times 1}$ はそれぞれ可視層-隠れ層間の重み行列、可視素子の偏差とバイアス、隠れ素子のバイアスを示しており、いずれも推定すべきパラメータである。

RBM では可視素子間、または隠れ素子間の接続は存在しないため（つまり、それぞれの可視素子、隠れ素子は互いに条件付き独立であるため）、それぞれの素子の条件付き確率 $p(v_i = 1 | \mathbf{h})$ 、 $p(h_j = 1 | \mathbf{v})$ は以下のような単純な関数で表現される。

$$p(v_i = 1 | \mathbf{h}) = \mathcal{N}(v_i; b_i + \mathbf{W}_{i \cdot} \mathbf{h}, \sigma_i^2) \quad (4)$$

$$p(h_j = 1 | \mathbf{v}) = \mathcal{S}(c_j + \mathbf{W}_{:j}^T \left(\frac{\mathbf{v}}{\sigma^2} \right)) \quad (5)$$

ここで、 $\mathbf{W}_{i \cdot}$ と $\mathbf{W}_{:j}$ は \mathbf{W} の第 i 列ベクトル、第 j 行ベクトルを表す。また、 $\mathcal{N}(x; \mu, \sigma)$ は平均 μ 、偏差 σ の正規分布、 $\mathcal{S}(x)$ はシグモイド関数を表す（すなわち $\mathcal{S}(x) = \frac{1}{1+e^{-x}}$ ）。

それぞれの RBM のパラメータは、 N 個の観測データを $\{\mathbf{v}_n\}_{n=1}^N$ とするとき、この確率変数の対数尤度 $\mathcal{L} = \log \prod_n p(\mathbf{v}_n)$ を最大化するように推定される。この対数尤度をそれぞれのパラメータで偏微分すると、

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{ij}} = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial b_i} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}} \quad (7)$$

$$\frac{\partial \mathcal{L}}{\partial c_j} = \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}} \quad (8)$$

が得られる。ただし、 $\langle \cdot \rangle_{\text{data}}$ と $\langle \cdot \rangle_{\text{model}}$ はそれぞれ、観測データ、モデルデータの期待値を表す。しかし、一般に後者の期待値に関しては計算困難であるため、代わりに式 (5)(4) によって得られる再構築したデータの期待値 $\langle \cdot \rangle_{\text{recon}}$ が用いられる（CD 法: Contrastive Divergence 法 [10]）。それぞれのパラメータは式 (6)(7)(8) から、確率的勾配法を用いて繰り返し更新される。

2.2 RTRBM

RTRBM は RBM の拡張モデル [13] であり、時系列データを取り扱うことに適している。RBM における可視層-隠れ層間の無向グラフに加えて、RTRBM では過去 P フレーム前から現時刻 t までの隠れ素子集合 $\{\mathbf{h}^{(p)}\}_{p=t-P}^t$ から、現時刻 t における可視素子

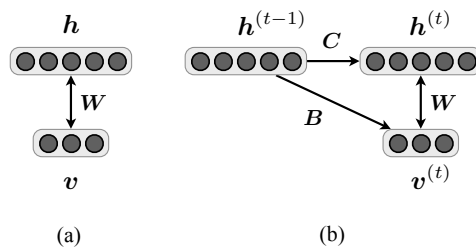


Fig. 1 Graphical representation of (a) an RBM and (b) an RTRBM.

$\mathbf{v}^{(t)}$ 、隠れ素子 $\mathbf{h}^{(t)}$ への有向グラフを考慮したモデルとなっている。簡単のため、ここでは $P = 1$ とする（このときの RTRBM を Fig. 1(b) に示す）。このモデルでは、Fig. 1(b) のように、3 種類の推定すべき重みパラメータ： $\mathbf{W} \in \mathbb{R}^{I \times J}$ ($\mathbf{v}^{(t)}$ と $\mathbf{h}^{(t)}$ 間の無向重み行列)、 $\mathbf{B} \in \mathbb{R}^{I \times J}$ ($\mathbf{h}^{(t-1)}$ から $\mathbf{v}^{(t)}$ への有向重み行列)、 $\mathbf{C} \in \mathbb{R}^{J \times J}$ ($\mathbf{h}^{(t-1)}$ から $\mathbf{h}^{(t)}$ への有向重み行列) が存在する。これらの重みパラメータは、RBM と同様に、CD 法を用いて推定される。RTRBM では、過去の素子集合 $\mathcal{A}^{(t)} = \{\mathbf{v}^{(\tau)}, \mathbf{h}^{(\tau)} | \tau < t\}$ が与えられたときの \mathbf{v} の条件付き確率密度は以下のように表される。

$$p(\mathbf{v}^{(t)} | \mathcal{A}^{(t)}) = \frac{1}{Z} \sum_{\mathbf{h}^{(t)}} e^{-E(\mathbf{v}^{(t)}, \mathbf{h}^{(t)} | \mathcal{A}^{(t)})} \quad (9)$$

ただし、 Z は正規化項を表す。 E は以下のエネルギー関数を示している。

$$\begin{aligned} E(\mathbf{v}^{(t)}, \mathbf{h}^{(t)} | \mathcal{A}^{(t)}) &= \left\| \frac{\mathbf{v}^{(t)} - \mathbf{b}^{(t)}}{2\sigma} \right\|^2 - \mathbf{c}^{(t)T} \mathbf{h}^{(t)} - \left(\frac{\mathbf{v}^{(t)}}{\sigma^2} \right)^T \mathbf{W} \mathbf{h}^{(t)} \end{aligned} \quad (10)$$

$$\mathbf{b}^{(t)} = \mathbf{b} + \mathbf{B} \mathbf{h}^{(t-1)} \quad (11)$$

$$\mathbf{c}^{(t)} = \mathbf{c} + \mathbf{C} \mathbf{h}^{(t-1)} \quad (12)$$

RBM の場合と同様に、観測データの対数尤度 $\mathcal{L} = \log \prod_t p(\mathbf{v}^{(t)} | \mathcal{A}^{(t)})$ をそれぞれのパラメータで偏微分すると、

$$\frac{\partial \mathcal{L}}{\partial \mathbf{B}_{ij}} = \left\langle \frac{v_i^{(t)} \hat{h}_j^{(t-1)}}{\sigma_i^2} \right\rangle_{\text{data}} - \left\langle \frac{v_i^{(t)} \hat{h}_j^{(t-1)}}{\sigma_i^2} \right\rangle_{\text{model}} \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{C}_{j'j}} = \langle \hat{h}_{j'}^{(t-1)} \hat{h}_j^{(t)} \rangle_{\text{data}} - \langle \hat{h}_{j'}^{(t-1)} \hat{h}_j^{(t)} \rangle_{\text{model}} \quad (14)$$

が得られる。ただし、 $\hat{\mathbf{h}}^{(t-1)}$ は $\mathbf{h}^{(t-1)}$ の期待値を表す。無向グラフに関するパラメータの偏微分に関してはそれぞれ式 (6)(7)(8) と同様にして導出される。

一度パラメータが推定されれば、RTRBM の前方推論 ($\mathbf{v}^{(t)}$ と $\mathbf{h}^{(t-1)}$ が与えられたときの $\mathbf{h}^{(t)}$ の条件付き確率) と後方推論 ($\mathbf{h}^{(t)}$ と $\mathbf{h}^{(t-1)}$ が与えられた

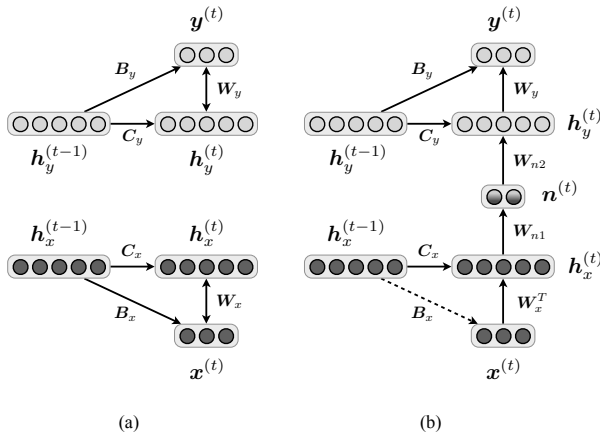


Fig. 2 (a) RTRBMs for a source speaker (below) and a target speaker (above), (b) our proposed voice conversion architecture, which combines two speaker-dependent RTRBMs with an NN.

ときの $v^{(t)}$ の条件付き確率) は, それぞれ以下のよう
に計算される .

$$p(h_j^{(t)} = 1 | v^{(t)}, h^{(t-1)}) = \mathcal{S}(c_j^{(t)} + \mathbf{W}_{:j}^T (\frac{\mathbf{v}^{(t)}}{\sigma^2})) \quad (15)$$

$$p(v_i^{(t)} = 1 | h^{(t)}, h^{(t-1)}) = \mathcal{N}(v; b_i^{(t)} + \mathbf{W}_{i \cdot} h^{(t)}, \sigma_i^2) \quad (16)$$

3 提案手法

提案手法による声質変換の概要を Fig. 2 に示す . 我々のアプローチでは, 特定話者の音声のみを含むデータを用いて RTRBM の特定話者モデルを学習させておく (Fig. 2 (a)) . 変数 $x^{(t)}$ と $y^{(t)}$ は時刻 t での, 入力話者と出力話者の RTRBM における可視素子であり, 例えば MFCC などの音響特徴量ベクトルを表す . 音声データは時系列データであるので, 前後フレームの特徴間には強い相関がある . RTRBM を用いれば, この潜在的な相関関係を捉えることにより, 音声データをより適切に表現できると考えられる . 提案手法では, このようにして (式 (15) より) 得られた特徴量同士の変換 ($h_x^{(t)}$ から $h_y^{(t)}$ へ) に, NN を用いる (便宜上 L を NN の隠れ層の数とし, $0 \leq L \leq 1$ とする) . この NN は, 入力話者音響特徴量を入力話者 RTRBM で射影した高次特徴量 $h_x^{(t)}$ を入力, 出力話者音響特徴量を出力話者 RTRBM で射影した高次特徴量 $h_y^{(t)}$ を教師信号として学習を行う . NN のパラメータ $\{W_l, d_l\}_{l=0}^L$ は典型的な NN の枠組みと同様に, NN の出力ベクトル $\eta(h_x^{(t)})$ と教師ベクトル $h_y^{(t)}$ の差を最小化するように推定される . 一度パラメータが推定されれば, 入力ベクトル $h_x^{(t)}$ は以下のよう

に出力話者の高次特徴量へ変換される .

$$h_y^{(t)} \approx \eta(h_x^{(t)}) = \bigodot_{l=0}^L \mathcal{S}(W_l h_x^{(t)} + d_l) \quad (17)$$

ただし, $\bigodot_{l=0}^L$ は $L+1$ 個の合成関数を表す (例えば隠れ層を 1 つ持つ NN の場合, $\eta(z) = \mathcal{S}(W_1 \mathcal{S}(W_0 z + d_0) + d_1)$) . NN の出力ベクトルから出力話者の音響特徴量へ逆射影するには, 式 (16) による RTRBM の後方推論によって計算される .

以上の議論をまとめると, 過去の特徴ベクトル $x^{(t-1)}$, $y^{(t-1)}$ を観測したときに, 時刻 t の入力話者音声の音響特徴量 $x^{(t)}$ から出力話者音声の音響特徴量 $y^{(t)}$ へ変換する, 提案法による声質変換式は以下のように表すことができる .

$$y^{(t)} = \underset{y^{(t)}}{\operatorname{argmax}} p(y^{(t)} | x^{(t)}, h_x^{(t-1)}, h_y^{(t-1)}) \\ = \mathbf{a}_{L+2}^{(t)} + \mathbf{W}_{L+2} \bigodot_{k=0}^{L+1} \mathcal{S}(\mathbf{a}_k^{(t)} + \mathbf{W}_k x^{(t)}) \quad (18)$$

ここで, $\mathbf{a}_{(k)}^{(t)}$ と $\mathbf{W}_{(k)}$ はそれぞれ

$$\mathbf{a}^{(t)} = \{\mathbf{a}_k^{(t)}\}_{k=0}^{L+2} = \{c_x^{(t)}, d_0, \dots, d_L, b_y^{(t)}\} \quad (19)$$

$$\mathbf{W} = \{\mathbf{W}_k\}_{k=0}^{L+2} = \{\mathbf{W}_x^T, \mathbf{W}_0, \dots, \mathbf{W}_L, \mathbf{W}_y\} \quad (20)$$

の要素を表す . ただし, $c_x^{(t)}$, $b_y^{(t)}$ はそれぞれ式 (12), 式 (11) で計算される, 入力話者 RTRBM, 出力話者の RTRBM の動的バイアスを表す .

式 (18) で表せられる変換式は, $L+4$ 個の層を持つ RNN を示唆している . すなわち, 2 つの異なる RTRBM と中間の NN を合わせて 1 つのネットワークとみなすことができ, BPTT (back-propagation through time) により音響特徴量の平行データを用いてそれぞれのパラメータを微調整することが可能である .

4 評価実験

本実験では ATR 研究用日本語音声データベース (A Set) [15] を用いた声質変換を行い, 提案手法の効果を確認した . このデータベースから, 入力話者として男性話者 1 名 (MMY) を, 出力話者として女性話者 1 名 (FTK) を選んだ . 学習・評価用のデータは STRAIGHT パラメータから抽出した 0 次元を除く $D = 24$ 次元の MFCC 特徴を用いた . 学習用のデータとして, 216 単語の音声から動的計画法を用いて作成した平行データのうち, ランダムに選んだ 5,000, 10,000, 20,000, 40,000 フレームを使用した . 提案手法 (SD-RTRBM) の評価方法として, 平均 MCD (mel cepstral distortion) による客観評価値と,

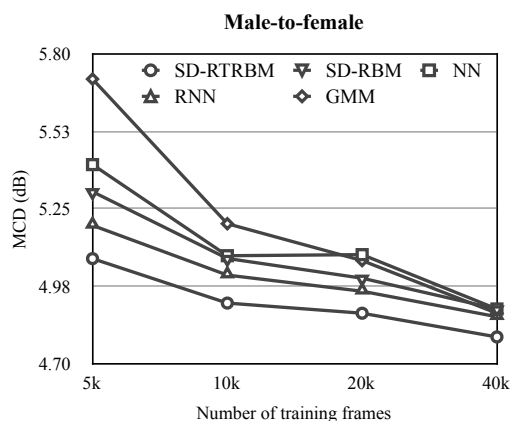


Fig. 3 Average MCD for each method with varying amounts of training data.

Table 1 MOS w.r.t. similarity for each method.

SD-RTRBM	SD-RBM	NN	GMM
2.86	2.80	2.77	2.14

MOS (mean opinion score) による主観評価値を用いて、従来手法である GMM と NN, 我々の先行研究である話者依存型 RBM (SD-RBM) と比較した。提案手法は隠れ素子数 $J = 72$, 中間 NN の隠れ層数 $L = 0$ とした。比較手法である NN は隠れ素子数 $J = 24$, 隠れ層数 $L = 2$ とした。GMM では 64 混合を用いた。いずれのハイパーパラメータも、先行実験で最もよい変換精度 (MCD 基準) が得られたものを用いている。客観評価に用いたデータは、学習データには含まれない 15 文の音声から作成された。主観評価では、このうちランダムに選んだ 5 文を 7 人の被験者が聞き、5 段階評価 (5: excellent; 4: good; 3: fair; 2: poor; and 1: bad) を付け、その平均を他の手法と比較した。また、主観評価実験では参考のため、提案手法と同じ構造を持つが初期値がランダムである RNN とも比較した。

客観評価、主観評価による実験結果をそれぞれ Fig. 3, Table 1 に示す。これらの図表より、総じて提案手法が他の手法と比較して高い変換精度を示したことが確認できる。

5 おわりに

本稿では、話者ごとに学習を行った RTRBM から得られた高次特徴量同士を変換させることで、出力話者の音響特徴量を得る声質変換法を提案し、実験によりその効果を確認した。今後は、最適な隠れ層数や素子数を自動的に決定する手法を提案していきたい。

参考文献

- [1] A. Kain and M. W. Macon: "Spectral voice conversion for text-to-speech synthesis", IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 285–288 (1998).
- [2] K. Nakamura, T. Toda, H. Saruwatari and K. Shikano: "Speaking-aid systems using GMM-based voice conversion for electrolaryngeal speech", Speech Communication, **54**, 1, pp. 134–146 (2012).
- [3] A. Kunikoshi, Y. Qiao, N. Minematsu and K. Hirose: "Speech generation from hand gestures based on space mapping", Proc. Interspeech, pp. 308–311 (2009).
- [4] Y. Stylianou, O. Cappé and E. Moulines: "Continuous probabilistic transform for voice conversion", IEEE Transactions on Speech and Audio Processing, **6**, 2, pp. 131–142 (1998).
- [5] T. Toda, A. W. Black and K. Tokuda: "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory", IEEE Transactions on Audio, Speech, and Language Processing, **15**, 8, pp. 2222–2235 (2007).
- [6] R. Takashima, T. Takiguchi and Y. Ariki: "Exemplar-based voice conversion in noisy environment", 2012 IEEE Spoken Language Technology Workshop (SLT), pp. 313–317 (2012).
- [7] R. Takashima, R. Aihara, T. Takiguchi and Y. Ariki: "Noise-robust voice conversion based on spectral mapping on sparse space", SSW8, pp. 71–75 (2013).
- [8] S. Desai, E. V. Raghavendra, B. Yegnanarayana, A. W. Black and K. Prahallad: "Voice conversion using artificial neural networks", IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3893–3896 (2009).
- [9] Y. Freund and D. Haussler: "Unsupervised learning of distributions of binary vectors using two layer networks", Computer Research Laboratory (1994).
- [10] G. E. Hinton, S. Osindero and Y.-W. Teh: "A fast learning algorithm for deep belief nets", Neural computation, **18**, 7, pp. 1527–1554 (2006).
- [11] T. Nakashika, R. Takashima, T. Takiguchi and Y. Ariki: "Voice conversion in high-order eigen space using deep belief nets", Proc. Interspeech, pp. 369–372 (2013).
- [12] Z. Wu, E. S. Chng and H. Li: "Conditional restricted boltzmann machine for voice conversion", IEEE China Summit and International Conference on Signal and Information Processing (ChinaSIP) (2013).
- [13] I. Sutskever, G. Hinton and G. Taylor: "The recurrent temporal restricted boltzmann machine.", NIPS, Vol. 19, pp. 1601–1608 (2008).
- [14] K. Cho, A. Ilin and T. Raiko: "Improved learning of gaussian-bernoulli restricted boltzmann machines", Artificial Neural Networks and Machine Learning, pp. 10–17 (2011).
- [15] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara and K. Shikano: "ATR Japanese speech database as a tool of speech recognition and synthesis", Speech Communication, **9**, 4, pp. 357–363 (1990).