リアルタイムニューラルボコーダにおける学習データ量の影響の調査* 松原圭亮^{1,2},岡本拓磨²,高島遼一¹,滝口哲也¹,戸田智基^{3,2},志賀芳則²,河井恒² ¹神戸大学,²情報通信研究機構,³名古屋大学

1 はじめに

近年,ニューラルネットワークによる音声合成技術 の発展により,テキスト音声合成は人間の声と区別 がつかないほど高品質な音声を合成できるようになっ てきている[1]。これらの発展の大きな転換点として, 直接波形生成モデル WaveNet [2]の登場が挙げられ る。WaveNet の登場以前は,テキストから音響特徴 量を推定する音響モデルにはニューラルネットワーク が用いられていたが,音響特徴量から波形に変換する ボコーダに関してはソースフィルタに基づく手法が 用いられており,ボコーダの処理の部分で肉声感が失 われてしまうという問題があった。そこで WaveNet ボコーダ [3] をはじめとするニューラルボコーダを用 いることによりそれらの問題を解決し,従来手法の 品質を大きく上回るテキスト音声合成を実現した。

WaveNet には合成速度が遅いため応用先が制限さ れるという問題点があったが, WaveNet より高速な 新しいニューラルボコーダモデルが数多く提案され てきており, CPU でもリアルタイム合成が可能なモ デルも登場してきている [4,5]。これらの先行研究で は数時間から数十時間のデータを用いて学習されて いるが,実際の課題解決に応用する場合は大量の学 習データを用意できないことが考えられる。データ 量の少ない話者の例として,障害により発話が不安 定になってしまっている構音障害者が挙げられる。構 音障害者は発話動作による本人の負担が大きいため, 大量の音声データを用意することができない。構音障 害者の発話を支援することは障害による社会活動の 格差を低減し、よりよい社会形成の一助になると考 えられる[6]。ニューラルボコーダに必要な学習デー タ量は WaveNet ボコーダでの調査がされている程度 であり [7], リアルタイムニューラルボコーダに必要 なデータ量の調査は行われていない。

本稿では、少量の学習データしかないような場合で もリアルタイムニューラルボコーダにより高品質な音 声を合成できるかを検証するため、WaveNet, LPC-Net [4], Parallel WaveGAN[5]を 1/8 時間から 9 時 間のデータで学習させ、品質の変化に対する主観評 価を行う。その際、より高品質な合成を実現するため [1,8], 従来はサンプリング周波数 16 kHz であった LPCNet を 24 kHz へと拡張させる。また, LPCNet と Parallel WaveGAN における CPU コア数による 合成速度の変化についても考察する。

2 ニューラルボコーダ

本稿では WaveNet, LPCNet, Parallel WaveGAN の3つのニューラルボコーダを採用した。WavaNet はリアルタイム生成はできないが,後者の2つのモ デルはリアルタイム合成が可能である。

2.1 WaveNet

WaveNet は自己回帰機構を用いて波形を直接推定 するニューラル生成モデルである。一般的な音声信号 は16 bit 整数値であるが, WaveNet では予測を簡単 にするために µ-law 変換を適用し8 bit 整数値に圧縮 している。また波形予測を回帰問題としてではなく, 256 種類の離散型分類問題としている。

WaveNet は高品質な音声を生成可能だが,自己回 帰モデルであることと,数十層の畳み込みニューラル ネットワークを用いた巨大なネットワークであること から生成速度が遅く,応用先が制限されてしまうとい う問題点がある。

2.2 LPCNet

LPCNet は再帰的ニューラルネットワークを用いた WaveRNN[9] ベースのニューラルボコーダであり,線 形予測分析 (linear prediction coding: LPC) した残差 信号をニューラルネットワークで推定する。WaveRNN が Dual Softmax により 16bit 波形を直接生成するの に対して, LPCNet は WaveNet と同様に μ-law 圧縮 した 8bit 信号を生成するが,残差信号であるため量子 化誤差が少なく,結果として WaveRNN よりも少な いモデルパラメータで高品質な音声合成が可能である [4]。Fig. 1 に概略図を示す。LPCNet は,入力の音響 特徴量から特徴量抽出を行う Frame rate network と, LPC で予測された音声波形サンプルから残差信号を 推定する Sample rate network の2つのモデルから 構成される。ここで,サンプリング周波数が16 kHz の場合,入力特徴量は18次元のバーク尺度のケプス トラムとピッチ周期, ピッチ相関である。また, LPC

^{*}Investigation of the effect of amount of training data for real-time neural vocoders by MATSUB-ARA, Keisuke^{1,2}, OKAMOTO, Takuma², TAKASHIMA, Ryoichi¹, TAKIGUCHI, Tetsuya¹, TODA, Tomoki^{3,2}, SHIGA, Yoshinori² and KAWAI, Hisashi² (¹Kobe Univ, ²NICT, ³Nagoya Univ)



Fig. 1 LPCNet structure.

による予測波形は以下の式で表される。

$$p_t = \sum_{k=1}^M a_k s_{t-k} \tag{1}$$

ここで p_t は t ステップ目の予測波形サンプル, a_k は k次 LPC 係数, s_{t-k} はt-kステップ目の実際の波 形である。LPC係数は入力特徴量のバーク尺度のケ プストラムから計算される。具体的には,まずバーク ケプストラムから通常の周波数尺度のパワースペク トル密度を計算する。次にパワースペクトル密度を 逆フーリエ変換することで自己相関関数を計算する。 最後に自己相関関数からレビンソンダービン法を用 いて LPC 係数を計算する。LPC の予測波形サンプル と,音響特徴量を入力した Frame rate network の出 力,1ステップ前の実際の音声波形サンプルと予測残 差波形サンプルを入力として,現在の予測残差波形サ ンプルを Sample rate network で予測する。Sample rate network は Gated Recurrent Unit(GRU) を用い た自己回帰モデルであるため,本来は WaveNet と同 様に生成に時間を要するが, Sparse coding と呼ばれ るネットワークの重み行列のうち値が小さいものをゼ ロに置き換える手法を用いることで高速化を行って おり, モバイル CPU でもリアルタイム生成が可能と なっている。

2.2.1 サンプリング周波数 24 kHz への拡張

LPCNet は従来サンプリング周波数 16 kHz の音声 を合成可能なニューラルボコーダとして提案されて いるが,更に高品質な 24 kHz の音声を合成するため に入力特徴量の修正を行う。従来手法では 18 次元の バークケプストラムとピッチ周期,ピッチ相関の計 20 次元を入力特徴量としていたが,ここではバーク ケプストラムの次元数を 20 に拡張する。

バーク尺度でのフィルタバンク分割は, 音声圧縮手 法の Opus[10] に準拠して行う。Fig. 2 にフィルタバ



Fig. 2 Opus and Bark band layout.

ンク構造を示す。この手法では低周波領域ではバーク 尺度の大きさによらず一定の間隔で分割していき,高 周波領域ではバーク領域での間隔が一定となるように 分割している。16 kHz でサンプリング可能な8 kHz までをカバーするには低次18次元,24 kHz でサンプ リング可能な12 kHz までカバーするには低次20次 元までを用いる。

2.3 Parallel WaveGAN

Parallel WaveGAN は, Parallel WaveNet を Generator とする敵対的生成ネットワーク (Generative Adversarial Network: GAN) をベースとするニュー ラルボコーダである。入力はホワイトノイズと音響特 徴量であり, Generator である WaveNet が全てのサ ンプルを同時に生成する。Discriminator は Generator が生成した音声波形を偽と判別できるように学習 される。また通常の GAN で用いられる Adversarial loss に加えて,出力波形と目標波形に短時間フーリエ 変換 (STFT) を行った場合の誤差である STFT loss を導入して Generator の学習を補助している。STFT は時間解像度と周波数解像度がトレードオフの関係 であり , 一方を細かく分析しようとするともう一方の 分析が粗くなってしまう問題点があるが,ここでは解 像度の異なる複数の STFT での損失を導入すること で時間解像度と周波数解像度の両方を担保している。

Parallel WaveGAN は自己回帰機構を持たないた め、一度に複数のサンプルを同時に生成することで高 速生成を可能にしており、CPUでもリアルタイム生成 が可能である¹。また、Parallel WaveNet と比べ、知 識蒸留なしで直接並列生成モデルを学習可能である。

3 実験

3.1 実験条件

WaveNet, LPCNet, Parallel WaveGAN ボコーダ における学習データ量と品質の関係を調査するため, サンプリング周波数 24 kHz の音声を用いた分析合成 を行う。データには JSUT コーパス [11] より日本人女 性話者による 7697 文(約10時間)の音声を使用した。 また大語彙連続音声認識エンジンである Julius[12] を 用いて音素アライメントを行い,音声データの前後に

¹https://github.com/kan-bayashi/ParallelWaveGAN

含まれる無音区間の除去を行った。テストセットと検 証セットに 100 文ずつ,学習セットには最大 7497 文 (約9時間)を使用した。学習データ量の種類として は,LPCNet,Parallel WaveGAN は9,5,3,1,1/2, 1/4,1/8時間の7種類,WaveNet は9時間と1/8時 間の2種類を評価した。また比較として,WORLD ボ コーダ [13](D4C edition [14])と9時間のデータ量で 学習させた WaveGlow[15] ボコーダも評価に加えた。

WaveNet ボコーダは [3] と同様のモデルを使用した。また予測誤差による雑音成分を抑えるため,時不変ノイズシェーピング法 [16] を適用した。入力特徴量にはメルスペクトログラム 80 次元を使用した。メルスペクトログラムの計算では 42.7 ms の Hann 窓を用いて,フレームシフトを 5 ms とした。

LPCNet は [4] と同様のモデルを使用した。入力特 徴量にはバーク尺度の 20 次元ケプストラムとピッチ 周期, ピッチ相関を用いた。バークケプストラムの計 算に際しては, 20 ms の Vorbis 窓を用い, フレーム シフトを 10 ms としてスペクトル分析を行い, バー ク尺度を用いたフィルタバンクを適用した後に離散コ サイン変換を行った。ピッチの計算にはオープンルー プの相互相関関数をベースとする手法を用いた。

Parallel WaveGAN は [5] と同様のモデルを使用し た。入力特徴量にはメルスペクトログラム 80 次元を 使用した。メルスペクトログラムの計算では 20 ms の Hann 窓を用いて,フレームシフトを 12.5 ms とし た。また周波数帯域を 80 から 7600 Hz に制限した。 WaveGlow は [18] と同様のモデルを使用した。

パラメータの更新は, WaveNet, LPCNetはAdam を使用し, Parallel WaveGAN と WaveGlow では RAdam[17]を用いた。学習には NVIDIA V100 の GPUを使用し, WaveNet, LPCNet, Parallel Wave-GAN, WaveGlowの学習にそれぞれ4枚,8枚,2枚, 4枚使用した。

主観評価には,聴取実験による平均オピニオン評点 テストを行った。実験参加者は健常な聴覚である18 人の成人日本語母語話者で,テストセット20文に対 して18条件と原音の計19条件の380文をヘッドホ ン聴取により評価した。

3.2 実験結果

Table 1 に WaveNet, LPCNet, Parallel Wave-GAN, WaveGlow の生成速度を示す。また, Fig. 3 に LPCNet と Parallel WaveGAN における CPU コア数 による生成速度の変化を示す。WaveNet は NVIDIA の TeslaV100 の GPU を, LPCNet と Parallel Wave-GAN は Intel の Xeon6152 の CPU を用いて合成し た。Table 1 の結果より, LPCNet はシングル CPU で RTF=0.24 の高速な生成が可能であることが分か



Fig. 3 Result of real-time factors for inference by increasing the number of CPU cores.

Table 1 Result of real-time factors for inference using an NVIDIA Tesla V100 or Intel Xeon 6152. (*) denotes number of CPU cores or GPUs.

| model | RTF-CPU | RTF-GPU |
|------------------|----------|---------|
| WaveNet | - | 196(1) |
| LPCNet | 0.24(1) | 0.22(1) |
| Parallel WaveGAN | 0.41(16) | 0.02(1) |
| WaveGlow | - | 0.07(1) |

った。また Parallel WaveGAN はシングル CPU で RTF=2.38 となり, リアルタイム生成に届かない結 果となった。しかし Fig. 3 の結果より CPU コア数 を増やすことで生成速度が改善され, 16 コアで最大 RTF=0.41 まで高速化することができた。LPCNet は CPU コア数による生成速度の改善は見られなかった。 これは, Parallel WaveGAN が複数の波形サンプルを 同時に生成できるのに対し, LPCNet は自己回帰モ デルであるため, CPU コア数を増やしても並列処理 が出来ないからであると考えられる。

聴取実験の結果を Fig. 4 に示す。学習データ量が 9時間の場合は WaveNet が最も音質が良いという結 果になった。学習データ量を減らしていくと,LPC-Net, Parallel WaveGAN ともに音質が劣化していが, t 検定の結果,1時間から9時間までは結果に有意な 差は見られなかった。1/8 時間では, Parallel Wave-GAN は WORLD の品質より悪くなったが, LPCNet はWORLDより良い品質を維持していることが分かっ た。結果として, LPCNet と Parallel WaveGAN は 1時間程度の学習データ量で十分学習が可能と分かっ た。LPCNet が Parallel WaveGAN や WaveGlow よ り音質が良い理由としては,自己回帰モデルのため 過去の波形情報も使えることや,予測残差は白色化 しているため,結果として時不変ノイズシェーピン グ法 [16] と同様に予測誤差によるノイズが聴感上で 抑制されていることなどが考えられる。また,今回 の実験では WaveGlow の方が Parallel WaveGAN よ り若干音質が高かったが(t検定により有意差あり),



Fig. 4 Result of MOS test with 18 listening subjects. Confidence level of the error bars is 95 %.

Parallel WaveGAN は adversarial loss のハイパーパ ラメータを調整することにより品質が向上する可能 性があり,これは今後の課題とする。

4 まとめ

WaveNet, LPCNet, Parallel WaveGAN ボコーダ を少量の学習データで学習させ品質を評価した。結果 として, LPCNet と Parallel WaveGAN は1時間程度 の学習データで十分に学習ができることが示された。 また生成速度の実験結果から, LPCNet と Parallel WaveGAN はシングル CPU あるいはマルチ CPU で リアルタイム生成が可能であることが示された。

参考文献

- J. Shen *et al.*, "Neural TTS synthesis by conditioning WavaNet on mel spectrogram predictions," in *Proc. ICASSP*, Apr. 2018, pp. 4779– 4783.
- [2] A. van den Oord *et al.*, "WaveNet: A generative model for raw audio," in *Proc. SSW9*, Sept. 2016, p. 125.
- [3] A. Tamamori *et al.*, "Speaker-dependent WaveNet vocoder," in *Proc. Interspeech*, Aug. 2017, pp. 1118–1122.
- [4] J. Valin, J Skoglund, "LPCNet: Improving Neural Speech Synthesis Through Linear Prediction," in *Proc. ICASSP*, May. 2019, pp. 5891–5895.
- [5] R. Yamamoto *et al.*, "Parallel WaveGAN:A fast Waveform generation model based on generative adversarial networks with multi-resolution spectrogram," *arXiv:1910.11480*, 2019.
- [6] 南阪ら、"構音障害者の少量データを用いた深層
 学習による音声合成の検討",音講論, pp. 1011– 1014, Sept. 2019.
- [7] 林ら、"WaveNet ボコーダにおける学習データ 量の影響に関する調査",音講論, pp. 249-250, Mar. 2018.

- [8] A. van den Oord *et al.*, "Parallel WaveNet: Fast high-fidelity speech synthesis," in *Proc. ICML*, July 2018, pp. 3915–3923.
- [9] N. Kalchbrenner *et al.*, "Efficient Neutal Audio Synthesis," in *Proc. ICML*, July 2018, pp. 2415–2424.
- [10] J. Valin *et al.*, "High-quality, low-delay music coding in the Opus codec," in *Proc. 135th AES Convention*, Oct. 2013.
- [11] S. Takamichi *et al.*, "JSUT and JVS: free Japanese voice corpora for accelerating speech synthesis research," *Acoust. Sci. Tech.*, (in press).
- [12] A. Lee, T. Kawahara. "Recent Development of Open-Source Speech Recognition Engine Julius," in *Proc. APSIPA ASC* Oct. 2009, pp. 131–137.
- [13] M. Morise *et al.*, "WORLD: a vocoder-based high-quality speech synthesis system for realtime applications," *IEICE trans. Inf. Syst.*, vol. E99-D, no. 7, pp. 1877–1884, 2016.
- [14] M. Morise, "D4C, a band-aperiodicity estimator for high-quality speech synthesis," Speech Communication, vol. 84, pp. 57–65, Nov. 2016.
- [15] R. Prenger *et al.*, "WaveGlow: A flow-based generative network for speech synthesis," in *Proc. ICASSP*, May. 2019, pp. 3617–3621.
- [16] K. Tachibana *et al.*, "An investigation of noise shaping with perceptual weighting for WaveNet-based speech generation," in *Proc. ICASSP*, Apr. 2018, pp. 5664–5668.
- [17] L. Liu *et al.*, "On the Variance of the Adaptive Learning Rate and Beyond," *Proc. ICLR*, Apr. 2020.
- [18] T. Okamoto *et al.*, "Tacotron-based acoustic model using phoneme alignment for practical neural text-to-speech systems," in *Proc. ASRU*, Dec. 2019, pp. 214–221.