

平成 30 年 6 月 12 日現在

機関番号：12608

研究種目：基盤研究(B) (一般)

研究期間：2015～2017

課題番号：15H02724

研究課題名(和文) ガウス過程回帰に基づく音声合成技術の確立

研究課題名(英文) Establishment of speech synthesis framework based on Gaussian process regression

研究代表者

小林 隆夫 (Kobayashi, Takao)

東京工業大学・工学院・教授

研究者番号：70153616

交付決定額(研究期間全体)：(直接経費) 10,000,000円

研究成果の概要(和文)：多様で表情豊かな音声合成の実現に向け、統計的パラメトリック音声合成の新たな枠組みであるガウス過程回帰に基づく音声合成(GPR音声合成)技術の確立をめざして研究を行った。ガウス過程回帰に基づいたスペクトルパラメータ生成に加え、基本周波数および音韻継続長予測からなる韻律生成手法を提案し、GPR音声合成システムを構築した。評価実験を通してGPR音声合成手法の有効性を示すとともに、多様な話者性やスタイルによる音声合成への応用、従来手法では合成音声の韻律の自然性が不十分であった声調言語への適用を検討し、提案手法の有用性を示した。

研究成果の概要(英文)：The purpose of the research is to develop a novel statistical parametric speech synthesis framework based on Gaussian process regression (GPR). We have proposed prosody generation techniques including pitch pattern prediction and phone duration prediction as well as the spectral parameter generation technique based on GPR. We developed a GPR-based speech synthesis system and showed its effectiveness through assessment of synthetic speech quality. Furthermore, we examined the proposed framework for generating expressive speech. We also examined it for generating more natural-sounding prosody in speech synthesis of a tonal language.

研究分野：音声情報処理

キーワード：テキスト音声合成 統計的パラメトリック音声合成 韻律生成 ガウス過程回帰 GPR音声合成 HMM音  
声合成 機械学習 深層学習

## 1. 研究開始当初の背景

言語の壁を越えた話し言葉によるユニバーサルコミュニケーションを実現する上で、音声認識、言語翻訳、音響信号処理などと共に、任意の文章テキストから音声を生成する「テキスト音声合成」(Text-to-Speech, TTS, 以下では単に音声合成と記す)が重要な基盤要素技術の一つとなっている。最近では、カーナビゲーションシステム、電子ブック、スマートフォン上の音声対話エージェント、コミュニケーションロボットなど、音声合成技術が実用的に使われ、身近なものとなっている。

このような音声合成の実用化において、近年大きな役割を担っているのが隠れマルコフモデルに基づく音声合成(HMM 音声合成)と呼ばれる統計的パラメトリック音声合成技術の進展である。本研究代表者は、約 20 年前の HMM 音声合成の基本アイデアの提案から現在に至るまで HMM 音声合成の研究に関わり、この 10 年間は平均声モデルに基づく音声合成手法や合成音声のスタイル制御手法など、HMM 音声合成に基づいた多様な表情豊かな音声合成の研究進展に寄与してきた。

その一方で、HMM のモデル構造に起因する根本的な性能の限界が指摘され、この解決策として音声をノンパラメトリックなモデルにより表現する試みや、ディープニューラルネットワーク(DNN)を HMM 音声合成に組み込んだアプローチが提案され、最近では DNN やリカレントニューラルネットワーク(RNN)に基づく音声合成手法の研究が急速に進んでいる。

これに対し本研究代表者らは、既存手法とは異なる、機械学習に基づいた新たなアプローチとしてガウス過程回帰(GPR)に基づくスペクトルパラメータ生成手法を提案し、これまでに提案手法が従来の HMM 音声合成手法を上回る性能を持つことを示した。しかしながら、自然性の高い合成音声を生成するためには、スペクトルパラメータ系列の生成だけでなく、声の抑揚、アクセント、リズムなどを表す韻律、すなわち音声の基本周波数(F0)と音韻継続長の生成が必須であり、研究開始時点でガウス過程回帰に基づく音声合成手法の確立に至っていなかった。

## 2. 研究の目的

本研究ではこれまでの GPR に基づくスペクトルパラメータ生成手法を発展させ、韻律生成を含む新たな音声合成の枠組みである「ガウス過程回帰に基づく音声合成(GPR 音声合成)」技術を開拓・確立することをめざす。具体的には、以下の 4 項目について研究・開発を行うことを目的とする。

### (1) GPR 音声合成手法の確立

ガウス過程回帰に基づくスペクトルパラメータの生成に加え、韻律(F0, 音韻継続長)の生成手法を開発し、従来の HMM 音声合成に

代わる GPR 音声合成手法の枠組みを確立する。

### (2) GPR 音声合成に基づく多様な音声合成

HMM 音声合成で既の実現されている話者性やスタイル(発話様式・感情表現)の制御を、GPR 音声合成においても容易に可能にする手法を開発する。

### (3) 話し言葉音声・表情豊かな音声の合成

GPR 音声合成を自発性の高い話し言葉音声や多様なスタイル音声の合成に適用し、より自然な韻律をもった合成音声を生成するシステムの構築をめざす。

### (4) ユニバーサルコミュニケーションに向けた音声合成

音声合成の基盤技術が十分でなく、韻律の制御が不十分なタイ語を含むいくつかの言語を対象として、多様な音声合成を実現する手法を検討する。

## 3. 研究の方法

### (1) GPR 音声合成手法の確立

既に基本的な検討を終えている GPR に基づくスペクトルパラメータ系列の生成に加えて、合成音声の自然性に大きな影響を及ぼす韻律の生成に提案手法を拡張する。具体的には、F0 パターン生成手法として、まずガウス過程分類の枠組により有声/無声区間の推定処理を実現し、有声区間であれば GPR に基づいて F0 値を生成する手法を開発する。これに、GPR に基づく音韻継続長予測モデルを加えて GPR 音声合成システムを構築する。

### (2) GPR 音声合成に基づく多様な音声合成

HMM 音声合成には、HMM のモデルパラメータを適切に修正するだけで、合成音声の話者性やスタイルを容易に制御できるという大きな特長があった。GPR 音声合成においても多様な音声合成の実現をめざし、特徴量の変換に基づいた適応手法を検討する。

### (3) 話し言葉音声・表情豊かな音声の合成

日常会話といった表情豊かな音声の合成は難しい問題であり、世界的研究レベルでも課題となっている。そこで、平静読上げ調音声に比べ、より表情豊かな音声としてオーディオブックを想定した物語読み聞かせ音声を対象に、提案手法の検討を行う。

### (4) ユニバーサルコミュニケーションに向けた音声合成

韻律が言葉の意味理解に大きく影響を及ぼす声調(トーン)言語では、従来手法による合成音声の韻律の自然性が十分とは言えない。ここでは、音声合成技術基盤が十分整備されていない声調言語としてタイ語を対象とし、GPR 音声合成により品質向上をめざす。また音声合成システムの多言語化に向けた検討を行う。

#### 4. 研究成果

##### (1) GPR 音声合成手法の確立

従来の HMM 音声合成や本研究課題である GPR 音声合成を実現するためには、入力テキストの読み（発声内容）から決まる音韻系列に対応するスペクトルパラメータ系列の生成と、発話の抑揚やリズムに相当する韻律生成（F0 パタン生成と音韻継続長予測）が必要となる。本研究代表者は既に GPR に基づくスペクトルパラメータ系列の生成手法を提案して有効性を示したが、韻律生成は実現していなかった。そこで、本研究における最重要課題を「GPR に基づく韻律生成手法の開発と GPR 音声合成技術の確立」と設定し、研究を進めた。

まず、GPR に基づくスペクトルパラメータ系列手法に倣い、アクセントなどの韻律情報に対してフレームレベルのコンテキストを定義することにより、GPR に基づくフレームレベルの F0 パタン生成手法を提案した（雑誌論文）。提案手法では、音素弁別特性に関する従来のフレームコンテキストに加え、アクセントの基本的な特徴であるモーラの高低、句頭のピッチ上昇（句頭音調）、アクセント核などを用い、それに伴う新たなカーネル関数を定義した。

一方、F0 パタンには値の存在する有声音のフレームと、値が定義できない無声音のフレームが存在することから、F0 パタンの生成には有声/無声を考慮したモデルが必要となる。ここでは、有声/無声推定モデルとしてガウス過程分類（GPC）を導入した（雑誌論文）。F0 モデルには対数 F0 を特徴量とする GPR を使用し、学習データの有声フレームを用いて学習する。合成時には、まず有声/無声推定モデルを用いて合成する文の各フレームの有声/無声を決定する。そして、F0 モデルから有声フレーム区間における F0 の予測分布を求め、その平均系列を生成 F0 としている。

GPR に基づくフレームレベル音響モデリングを用いたスペクトル特徴量および F0 パタン生成モデルに、GPR に基づく音素継続長の予測モデルを加えた GPR 音声合成システムを構築した（雑誌論文）（学会発表）。表 1 に音声合成システムを構成するモデルの種類を、図 1 にシステムの構成ブロック図を示す。

システムの実現にあたって、フレームレベルの音響特徴量の予測の際に、計算量削減の

表 1 GPR 音声合成システムを構成するモデル

モデル	出力変数	予測	単位
マルチケプストラム	多次元連続	GPR	フレーム
有声/無声	2値	GPC	フレーム
F0	1次元連続	GPR	フレーム
非周期性指標	多次元連続	GPR	フレーム
音素継続長	1次元連続	GPR	音素

表 2 対比較試験による主観評価結果[%]  
（雑誌論文）

HMM	DNN	GPR	Neutral	p 値
15.5	55.7		28.9	$< 10^{-10}$
14.0		52.4	33.6	$< 10^{-10}$
	19.9	23.5	56.5	0.261

ためグラム行列をブロック対角行列と低ランク行列の和で近似する PIC (partially independent conditional) 近似を用いた。この近似により、ブロック境界において音声パラメータが不連続になる場合が生じるが、HMM 音声合成で利用されている動的特徴量を用いた音声パラメータ生成手法を GPR 音声合成の枠組みに導入することでこの問題を解決した。

構築した GPR 音声合成システムの性能を、従来の HMM 音声合成システムと、HMM 音声合成に代わる手法として世界的規模で研究が進みつつある DNN 音声合成システムを用いて比較評価した。話者は ATR 日本語音声データベース B に含まれる女性と男性の各 2 名、計 4 名で、学習データ量は各話者 450 文章と比較的少量であるが、従来手法の HMM 音声合成では最低限必要な学習データ量を満たしている。HMM 音声合成には 5 状態のスキップなし隠れセマルコフモデル (HSMM) を使用した通常構成とした。DNN 音声合成では、予備実験により客観評価結果が良くなる構成を選んだ。

合成音声の自然性を対比較試験で主観的に評価した結果を表 2 に示す（雑誌論文）（学会発表）。表中の数字は合成音声より自然と感じられた手法の割合を示し、差がないと感じられた場合は Neutral を選べるものとした。この結果より、HMM 音声合成に対し、DNN 音声合成と提案手法である GPR 音声合

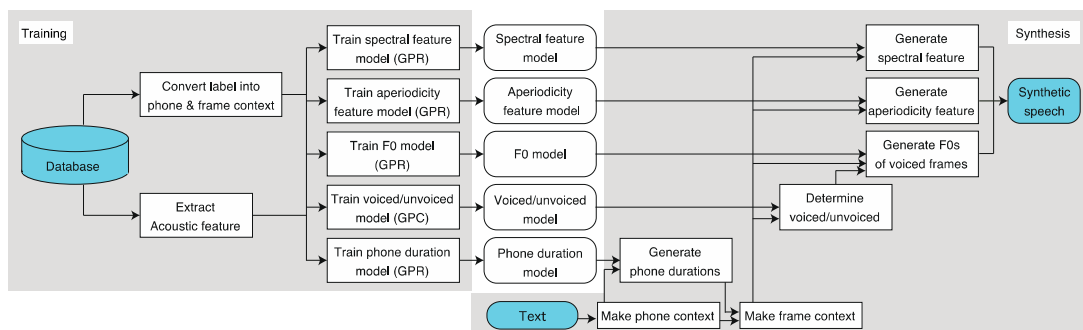


図 1 構築した GPR 音声合成システムの構成（雑誌論文）

成は有意にスコアが高いことを示した。DNN と GPR を比較すると、有意な差は見られなかったものの、GPR 音声合成は DNN 音声合成に比べてわずかに高いスコアが得られた。

以上の検討により、基本的な GPR 音声合成手法を確立することができ、合成音声の評価結果でも従来の HMM 音声合成を有意に上回り、比較的学習データが少量の条件において DNN 音声合成と同等以上の品質評価結果が得られることを示した。なお、提案した GPR 音声合成手法では PIC 近似を用いており、この際に木構造によるブロック分割を用いていることから、性能が木構造による予測性能に依存するという点が課題として残った。

これに対し、確率的勾配降下法により効率的な学習が可能な確率的変分ガウス過程 (SVGP) と、コンテキストの特徴抽出器としての DNN を組み合わせたハイブリッド手法 (GP-DNN ハイブリッド音声合成) を提案した [学会発表]。提案手法では、従来法と同様に音素弁別特性などの特徴量やフレームの相対位置をコンテキストとして、コンテキストの類似度に基づくカーネル関数をガウス過程の共分散関数として用いる。このとき、コンテキストを直接カーネル関数への入力に使用せず、コンテキストを DNN によって変換して得られた特徴スペクトルをカーネル関数への入力に使用する。

提案 GP-DNN ハイブリッドモデルと従来の DNN 音声合成、および DNN 音声合成に比べてより性能が高いとされる LSTM-RNN (long short-term memory RNN) と対比較試験に基づく主観性能評価を行った結果を表 3 に示す。ここでは、女性話者 1 名の 1593 文章 (約 119 分) を学習データとして用いた。GP-DNN と同様に深層構造に基づく手法である DNN と比べると、GP-DNN は DNN に比べ有意にスコアが高く、LSTM-RNN と同等かわずかに高いスコアが得られた。

GP-DNN ハイブリッドモデルでは深層構造に DNN を用いた。これに対し、ガウス過程 (GP) 自体を深層構造にすることが考えられる。そこで、深層ガウス過程を用いた音響モデリング手法を提案し、GPR 音声合成における深層構造の有用性について検討した [学会発表]。

GP-DNN ハイブリッドモデルを用いた音声合成の評価と同じく、女性話者 1 名の 1593 文章 (約 119 分) を学習データとして用いて提案深層 GP モデルに基づいた音声合成手法 (DeepGP) の評価を行った。比較対象として前述の GP-DNN ハイブリッド音声合成、従来手法として DNN 音声合成 (DNN) と双方向 LSTM に基づく RNN 音声合成 (LSTM) を用いた。5 段階 MOS 評価による合成音声の自然性の主観評価結果を図 2 に示す。その結果、危険率 5% で有意に DNN 音声合成の MOS 値が低かった。また、GP-DNN ハイブリッドモデルと深層 GP は、モデルにリカレント構造を有していないにも関わらず、LSTM-RNN に基づく音声合成

表3 GP-DNNハイブリッドモデルに基づいた合成音声の対比較テストによる主観評価結果[%] [学会発表]

GP-DNN	DNN	LSTM	Neutral	$p$ 値
63.3	13.3		23.3	$< 10^{-4}$
32.5		25.8	41.7	0.34

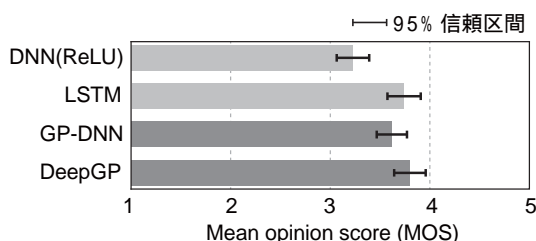


図2 深層 GP に基づく合成音声の MOS スコア [学会発表]

と同程度のスコアが得られた。

## (2) GPR 音声合成に基づく多様な音声合成

HMM 音声合成は平均声モデルと話者適応により、多様な話者性を持つ音声合成を容易に実現できる利点があった。HMM 音声合成の品質を上回る音声出力が可能で GPR 音声合成において同様な枠組みが実現できれば、多様な音声合成が可能になり有用であると考えられる。これに対し、GPR 音声合成における話者適応手法として、音響特徴量空間における線形変換に基づく手法を提案した [雑誌論文] [学会発表]。

HMM 音声合成における話者適応では、モデルパラメータであるガウス分布の平均および共分散行列を線形変換する手法が用いられている。しかし GP はノンパラメトリックモデルであるため、この手法を直接適用することができない。そこで、特徴量空間における適応手法 (fMLLR) に倣って特徴量に対しアフィン変換を行い、変換後の全話者の音響特徴量が同一の GP に従って生成されると仮定して音響特徴量の同時分布を求める。そして、適応データに対する条件付き確率を最大にすることにより、話者毎に特徴量変換パラメータを求めている。

提案手法 (GPR-SA) の有効性を調べるために、HMM 音声合成 (HMM-SA) との比較評価を行った。ATR 日本語音声データベースセット B に含まれる男性話者 2 名による各話者 450 文章、計 900 文章を用いて話者適応前のモデルを学習し、これとは異なる男性話者 2 名を目標話者とした。各話者 10 文章を適応データとして用いたときの、XAB 試験による合成音声の話者類似性の評価結果を図 3 に示す。提案手法は、HMM 音声合成を上回る品質を得ることができた。

同様の手法は話者適応だけでなくスタイル適応にも適用可能であり、GPR 音声合成における特徴量変換に基づくスタイル適応手法を



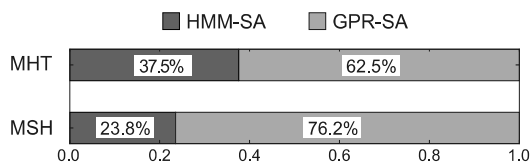


図3 話者適応に基づくGPR音声合成における話者類似性の評価(目標話者MHT, MSH) [雑誌論文]

提案し、HMM音声合成との比較評価を通して有効性を示した(学会発表)。

### (3) 話し言葉音声・表情豊かな音声の合成

構築したGPR音声合成システムを、表情豊かな音声の一例としてオーディオブック音声に適用し、提案手法により生成された表情豊かな合成音声の評価と問題点の検討を行った(学会発表)。その結果、従来法のHMM音声合成に比べて、韻律の変動が大きいオーディオブック音声に対して、スペクトル歪、F0歪、音素継続長歪が共に減少することを示した。しかし、オーディオブック音声を持つ表現の豊かさを十分に表現するまでには至っておらず、適切なコンテキストの検討が今後の課題として残された。

この他にも、GPR音声合成の枠組みを歌声音声合成に適用し、その可能性を検討した(学会発表)。

### (4) ユニバーサルコミュニケーションに向けた音声合成

提案GPR音声合成の枠組みによる多言語の合成に向けて、韻律の自然性の再現が難しいとされる声調言語の一つであるタイ語を対象として、詳細な検討を行った。

まず、GPR音声合成システムをタイ語音声合成に適用したときの基本性能の評価を行い、学習音声データ量が1ないし2時間程度の場合、従来のHMM音声合成やDNN音声合成に比べて客観評価、主観評価ともにより高い性能を持つことを明らかにした(雑誌論文)。

次に、より正確な音素継続長予測をめざし、音節単位の継続長が自然性に影響を及ぼすことを考慮して、まず音節継続長を予測するモデルを作り、次にこれから得られた予測値を新たなコンテキストとして加えて音素継続長を予測するという、2段構成のガウス過程回帰に基づく継続長予測手法を提案した(雑誌論文)。さらに、音素継続長予測を2段構成にするアプローチを従来のHMM音声合成やDNN音声合成の枠組みにも適用し、それぞれの音声合成の枠組みを用いて詳細な予測性能の比較評価を行った。その結果、提案手法が客観評価、主観評価とも従来手法に比べて有意に評価結果が上回ることを明らかにした。評価結果の例として、図4に従来の音素継続長予測(Single-level)を用い

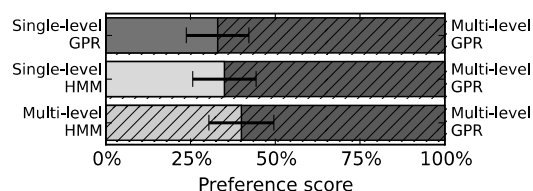


図4 タイ語音声合成における2段音素継続長予測の評価結果(対比較試験)(雑誌論文)

るHMM音声合成と、GPR音声合成に提案2段構成音素継続長予測(Multi-level)を適用したときの、合成音声の対比較試験の結果を示す。

さらに、複数レベルの継続長を考慮した音素継続長予測手法として、音節継続長と音素継続長の二つの予測モデルの積により表されたモデルから最適な音素継続長を予測する手法を提案した(雑誌論文)。また、基本周波数(F0)パタン生成時に、音節単位の韻律特徴を考慮した手法を提案した(雑誌論文)。

タイ語の他にも、インドネシア・スラバヤ工科大学の研究者の協力を得て、インドネシア語を対象とした音声合成システムの基礎的な検討を行った他、英語音声合成合成についてGPR音声合成が容易に適用可能であることを確認した。また、英語・日本語のクロスリンガル音声合成についてもHMM音声合成の枠組みで詳細な検討を行った(雑誌論文)上で、GPR音声合成への適用可能性を検討した。

## 5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

[雑誌論文](計25件)

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, GPR-based Thai speech synthesis using multi-level duration prediction, 査読有, Speech Communication, Vol.99, pp.114-123, DOI: 10.1016/j.specom.2018.03.005, 2018.

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, Enhanced F0 generation for GPR-based speech synthesis considering syllable-based prosodic features, 査読有, Proc. APSIPA Annual Summit and Conference, APSIPA ASC 2017, 4 pages, DOI: 10.1109/APSIPA.2017.8282285, 2017.

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, Duration prediction using multiple Gaussian process experts for GPR-based speech synthesis, 査読有, Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2017, pp.5945-5948, DOI: 10.1109/ICASSP.2017.7953207, 2017.

長濱大樹, 能勢 隆, 郡山知樹, 小林隆夫, クロスリンガル音声合成のための共有決定木コンテキストクラスタリングを用いた話者適応, 査読有, 電子情報通信学会論文誌 D, Vol.J100-D, No.3, pp.385-393, DOI: 10.14923/transinfj.2016PDP0020, 2017.

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, Unsupervised stress information labeling using Gaussian process latent variable model for statistical speech synthesis, 査読有, Proc. 17th Annual Conference of the International Speech Communication Association, INTERSPEECH 2016, pp.1591-1595, DOI: 10.21437/Interspeech.2016-273, 2016.

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, Tone modeling using Gaussian process latent variable model for statistical speech synthesis, 査読有, Proc. 8th International Conference on Speech Prosody, SPEECH PROSODY 2016, pp.1014-1018, DOI: 10.21437/SpeechProsody.2016-208, 2016.

Tomoki Koriyama, Syohei Oshio, Takao Kobayashi, A speaker adaptation technique for Gaussian process regression based speech synthesis using feature space transform, 査読有, Proc. IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2016, pp.5610-5614, DOI: 10.1109/ICASSP.2016.7472751, 2016.

Decha Moungsri, Tomoki Koriyama, Takao Kobayashi, Duration prediction using multi-level model for GPR-based speech synthesis, 査読有, Proc. Annual Conference of the International Speech Communication Association, INTERSPEECH 2015, pp.1591-1595, [https://www.isca-speech.org/archive/interspeech\\_2015/papers/i15\\_1591.pdf](https://www.isca-speech.org/archive/interspeech_2015/papers/i15_1591.pdf), 2015.

Tomoki Koriyama, Takao Kobayashi, A comparison of speech synthesis systems based on GPR, HMM, and DNN with a small amount of training data, 査読有, Proc. Annual Conference of the International Speech Communication Association, INTERSPEECH 2015, pp.3496-3500, [https://www.isca-speech.org/archive/interspeech\\_2015/papers/i15\\_3496.pdf](https://www.isca-speech.org/archive/interspeech_2015/papers/i15_3496.pdf), 2015.

Tomoki Koriyama, Takao Kobayashi, Prosody generation using frame-based Gaussian process regression and classification for statistical parametric speech synthesis, 査読有, Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2015, pp.4929-4933, DOI: 10.1109/ICASSP.2015.7178908, 2015.

郡山知樹, GPR 音声合成における深層ガウス過程の利用の検討, 電子情報通信学会・日本音響学会音声研究会, 2018.

郡山知樹, GPR 音声合成のための深層構造の利用の検討, 日本音響学会春季研究発表会, 2018.

郡山知樹, GP-DNN ハイブリッドモデルに基づく統計的音声合成の検討, 電子情報通信学会・日本音響学会音声研究会, 2018.

小林隆夫, 表現豊かな音声合成に向けた多様な話者性とスタイルによる音声合成への取組み(招待講演), 第 19 回音声言語シンポジウム, 2017.

郡山知樹, ガウス過程回帰に基づく歌声合成の検討, 日本音響学会 2017 年秋季研究発表会講演論文集, 2017.

津野駿幸, GPR 音声合成に基づいたオーディオブック音声の合成, 日本音響学会 2017 年春季研究発表会, 2017.

前野雄也, GPR 音声合成における区分線形特徴量変換を用いたスタイル適応の検討, 日本音響学会 2016 年秋季研究発表会, 2016.

岡元伶洋, 多様なスタイルによる GPR 音声合成の検討, 日本音響学会 2016 年春季研究発表会, 2016.

押尾翔平, GPR 音声合成における話者適応手法の検討, 日本音響学会 2015 年秋季研究発表会, 2015.

郡山知樹, ガウス過程回帰に基づく音声合成システムの評価, 日本音響学会 2015 年秋季研究発表会, 2015.

[ その他 ]

ホームページ等

<http://www.kbys.ip.titech.ac.jp/>

## 6 . 研究組織

### (1) 研究代表者

小林 隆夫 (KOBAYASHI, Takao)  
東京工業大学・工学院・教授  
研究者番号 : 70153616

### (2) 研究分担者

郡山 知樹 (KORIYAMA, Tomoki)  
東京工業大学・工学院・助教  
研究者番号 : 50749124

### (3) 連携研究者

### (4) 研究協力者

Decha Moungsri (MOUNGSRI, Decha)  
長濱 大樹 (NAGAHAMA, Daiki)  
能勢 隆 (NOSE, Takashi)  
Dhany Arifianto (ARIFIANTO, Dhany)