# List of Figures

# List of Tables

# List of Acronyms

**AAE** African American English

**AF** Amplitude of frication

**AH** Amplitude of aspiration

**AIR** All India Radio

**ANN** Artificial Neural Network

**ANOVA** Analysis of Variance

**ATR** Advanced Telecommunications Research

**AV** Amplitude of voicing

**A1** Amplitude of first formant

**A2** Amplitude of second formant

**BW** Bandwidth

**BFGS** Broyden-Fletcher-Goldfarb-Shanno

**CPP** Cepstral Peak Prominence

**CTD** Central Thai Dialect

**CLSP** Center for Language and Speech Processing

**CMU** Carnegie Melon University

**DCT** Discrete Cosine Transform

**DFW** Dynamic Frequency Warping

**DMOS** Differential Mean Opinion Score

**DNN** Deep Neural Network

**DRT** Diagnostic Rhyme Test

**DTW** Dynamic Time Warping

**EGG** Electro Glottograph

**EM** Expectation Maximization

**EMA** Electromagnetic Acoustic Imaging

**FFT** Fast Fourier Transform

**F0** Fundamental Frequency (Hz)

**F1** First Formant Frequency (Hz)

**F2** Second Formant Frequency (Hz)

**F3** Third Formant Frequency (Hz)

**FT** Formant Transformation

**GMM** Gaussian Mixture Model

**GSS** General Smoothing Splines

**G2P** Grapheme-to-Phoneme

**GVV** Glottal Volume Velocity

**H1** Amplitude of first harmonic

**H2** Amplitude of second harmonic

**HMM** Hidden Markov Model

**HTS** HMM based Text to Speech System

**HNR** Harmonic to Noise Ratio

**HTK** HMM-based Toolkit

**IVRS** Interactive Voice Response System

**IT** Information Technology

**IITG** Indian Institute of Technology, Guwahati

**IITG-HTS** HTS for standard Assamese built at IITG

**IPA** International Phonetic Alphabet

**JHU** John Hopkins University

**KLD** Kull-Back Leibler Divergence

**LLT** Local Linear Transformation

**LMR** Linear Multivariate Regression

**LP** Linear Prediction

**LPC** Linear Prediction Coefficients

**LSF** Line Spectral Frequency

**LSP** Line Spectrum Pairs

**L2S** Letter to Sound

**MA** Moving Average

**MCD** Mel Cepstral Distortion

**MCEP** Mel Cepstral Coefficients

**MDL** Minimum Distance Length

**MFC** Mel Frequency Cepstrum

**MFCC** Mel Frequency Cepstral Coefficients

**MGC** Mel Generalized Coefficients

**MIT** Massachusetts Institute of Technology

**ML** Maximum Likelihood

**MLSA** Mel Log Spectral Approximation

**MOS** Mean Opinion Score

**MRI** Magnetic Radio Imaging

**MUSHRA** Multiple Stimuli with Hidden Reference and Anchor

**NAL** Nalbaria

**NTD** Northern Thai Dialect

**OBR** Onset of Burst Release

**OQ** Open Quotient

**OV** Onset of Voicing

**PARCOR** Partial Correlation

**PID** Pitch Intensity Duration

**PSD** Power Spectral Density

**PSNR** Peak Signal to Noise Ratio

**RAPT** Robust Algorithm for Pitch Tracking

**RCILTS** Resource Centre for Indian Language Technology Solutions

**RMSE** Root Mean Square Error

**SDS** Spoken Dialogue Systems

**SPTK** Speech Processing Toolkit

**SPSS** Statistical Parametric Speech Synthesis

**SSANOVA** Smoothing Spline Analysis of Variance

**Stem-ML** Soft Template Mark-up Language

**SWIPE** Sawtooth Waveform Inspired Pitch Estimator

**S2S** Speech-to-Speech

**TD** Temporal Decomposition

**TP-FT** Text Prompts for the Formant Transformation module

**TP-TTS** Text Prompts for the TTS

**TP-VC** Text Prompts for the Voice Conversion module

**TTS** Text to Speech

**TU** Tezpur University

**TU-TTS** TTS for standard Assamese built at TU

**UCLA** University of California Los Angeles

**UV** Unvoiced

**V** Voiced

**VoQ** Voice Quality

**VQ** Vector Quantization

**VS** Voice Source

**VFS** Vowel Formant Space

**VFS-AIR** Vowel Formant Space for AIR

**VFS-NAL** Vowel Formant Space for NAL

**VFS-TTS** Vowel Formant Space for TTS

**VOT** Voice Onset Time

**VT** Voice Transformation

**WAE** White American English