# Evaluating user preferences in Hindi Text-to-Speech

**Bharat Gupta**
**MeitY. New Delhi, India**
**bharatg@gov.in**

## Abstract

Hindi holds the distinction of being the fourth most extensively spoken first language globally. It serves as an official language in India, encompassing several states within the country. Hindi also has a number of dialects that are scattered over the entire Hindi-speaking region. There hasn't been a phonological comparison of Hindi dialects. Text-To-Speech (TTS) systems can only be evaluated effectively using the mean opinion score (MOS) and degradation mean opinion score (DMOS) as recommended metrics for synthesized speech quality. These subjective metrics are the most widely used to assess speech synthesis. During the evaluation phase, numerous assessors from different locations tend to exhibit a bias towards their prosodic style. They often feel more comfortable when both speaking and listening in their native language with a prosodic manner. In this report, we studied the Hindi region's evaluators' preferences while evaluating the Hindi TTS system due to language's dialects and prosody, the study focuses on the influence of language dialects and prosody on the Degradation Mean Opinion Score (DMOS) and overall system performance. The current research is to discover the preferences and appropriate weightage of the dialects while evaluating the performance of Hindi TTS. Through a comparative analysis and an exploration of the details and recommended weights assigned to different dialects based on preferences, this research examines the variations in scores offered by different evaluators. The current research investigates various patterns of dialects in terms of character and word variations. The words variations have been organized by building Chhattisgarhi's& Haryanvi's word dictionary w.r.t Hindi. The evaluation score has also been analyzed by conducting evaluation testing on Hindi TTS having characters/words variations of the dialects. The W3C SSML (Speech Synthesis Markup Language) has been built for the implementation of various patterns of dialects on the TTS system. The current approach has been made for the development of dialect based TTS that is not currently available in the system.

## 1 Introduction

A system called Text to Speech (TTS) converts text into artificial speech. The purpose of testing and evaluating a TTS system is to determine how well the speech can be understood and how closely it resembles the human voice. TTS voices are typically evaluated using a subjective listening test in which listeners are presented with samples of synthesized speech and asked to rate them along dimensions such as clarity and overall quality of the experience [4,5,6,7,8]. The simplest way for assessing a voice synthesis system's quality is the Mean Opinion Score (MOS). The MOS provides a numerical assessment of the effectiveness of the TTS System's produced synthesized speech. By contrasting the genuine and synthetic voices, the DMOS (Degraded Mean Opinion Score) test evaluates how natural the speech sounds. The evaluation takes into account a number of factors, including word pronunciation, listening effort, speaking rate, naturalness, etc. As there are numerous dialects of Hindi, it is spoken throughout the entire world. According to the literature review, the dialect variances have not been taken into account while evaluating the TTS engine, which will have an impact on the evaluation's findings, which will then have an impact on the system's real performance. As a result, it is necessary to investigate various characteristics, standardize the grading of various terms based on preferences, and analyze the

variations of various dialect variants and by applying the suitable weights to different dialects/region's evaluators while calculating the MOS and DMOS score, the actual performance of Hindi & Dialect based TTS system can be reported. The research also investigates the word replacement mechanism by using the dictionary having dialect variations with respect to Hindi that further helps in building dialect based Text to speech system.

The current research revolves around examining character variations present in various Hindi dialects through the collection of data and the construction of dictionaries. These variations play a pivotal role in unveiling user preferences, elucidated by evaluator scores aligned with specific dialect requirements. Based on this analysis, the research compares the score of different dialects with Khadi Boli as this is the predominant dialect of Hindi language and it is much closer to the modern Hindi. Leveraging artificial intelligence methodologies, this research contributes to sustaining uniformity in the Text-to-Speech (TTS) quality grading system. Moreover, it is instrumental in the development of a TTS system tailored to different Hindi dialects.

## 2  Literature Survey

It is worthwhile to conduct studies evaluating the performance of the synthesizers with human listeners because the ultimate goal of deploying the synthetic speech is to make it usable to applications. According to [9], in order to more accurately evaluate speech synthesizers, it was required to include perception aspects in the synthetic speech evaluation rather than just gauging intelligibility [10] evaluated the listener's perception in a comprehension assignment to determine how well the listeners could comprehend the synthetic speech produced by the synthesizers. [11]. Through comprehension challenges, different voice synthesizers' performance can also be assessed. Intelligibility can be evaluated via comprehension evaluation, according to several researchers [12, 13]. A single score for MOS, which incorporates scores for a variety of speech impairments, serves as a general indicator of speech quality. The MOS score shouldn't be used only to describe speech quality because it is a generic statistic that takes numerous aspects into account. By contrasting the natural and synthetic voices, the DMOS test evaluates how natural the speech sounds. To prevent biased scoring when using the DMOS approach, evaluators must listen to both synthetic and natural voice samples in random order without knowing beforehand whether they are synthetic or natural [14].

For the proper evaluation of the TTS system, it is important to consider all the aspects that affect the MOS and DMOS scores of the synthetic voices. There is a need to examine and consider all the parameters for the evaluation of the TTS that is developed for particular language say Hindi that has varieties of dialects. Hindi is basically used by the citizens of northern and central states of India [15].

The list of States with the highest percentage of native Hindi speakers is shown below [16].

| State | Hindi Speaking rates |
|---|---|
| **Bihar** | 76.16% |
| **Uttar Pradesh** | 91.40% |
| **Haryana** | 87.56% |
| **Rajasthan** | 91.03% |
| **Himachal Pradesh** | 89.02% |
| **Uttarakhand** | 88.05% |
| **Chhattisgarh** | 82.76% |
| **Jharkhand** | 57.64% |
| **Madhya Pradesh** | 87.20% |

Table 1: State wise Native Hindi Speakers

It is recommended that Hindi should be one of the UN's official languages because it has such a big native speaker population. The Indian government is aggressively addressing this issue in this regard [17]. There are numerous varieties of Hindi, some of which are considered to be "proper" dialects. Below is a list of state-specific dialects:

| State | Dialect | Region |
|---|---|---|
| **Uttar Pradesh** | Awadhi, Bhojpuri, Bagheli, Brajbhasha, Bundeli, Kannauji, Khadiboli | Northwestern part of state |
| **Bihar** | Bihari | |
| **Haryana** | Haryanvi | Northern state of Haryana, Delhi |
| **Rajasthan** | Rajasthani | Rajasthan and neighboring states of Gujarat, Haryana and Punjab |
| **Himachal Pradesh/ Uttarakhand** | Pahari, Pahari (Kumaoni, Garhwali) | Himachal Pradesh & Uttarakhand |
| **Chhattisgarh** | Chhattisgarhi | Chhattisgarh , adjacent areas |

| | | of Madhya Pradesh, Orissa and Jharkhand |
| --- | --- | --- |
| **Madhya Pradesh** | Bundelkhandi/ Bundeli | Madhya Pradesh and southern parts of Uttar Pradesh |

Table 2:  State wise Hindi Dialects

There are numerous dialects of Hindi as mentioned above that has phonetic variations that leads to the different written form in some regions. But as per the literature survey, the complete information of the variations pertaining to different dialects is not available on the web resources. Although some offline resources discuses about the phonetic representation of the dialects of different regions. The current study collate variations occur in the major Hindi dialects that helps in evaluating the performance of TTS engines.

## 3   Methodology

The following research methodology has been adopted in the present research:
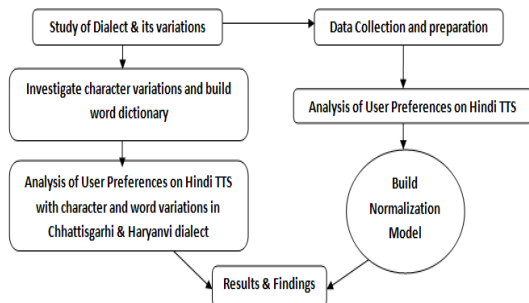


Fig 1:  Methodology of current research

## 4   Data Collection and preparation

The current study examined 20 different sentences from 5 different text-to-speech systems using male and female voices. The data on Hindi dialect variants has been gathered based on the goals of the current study. The evaluation sheet has been constructed to include information about the evaluators' backgrounds, including their education, age, and addresses. The 10 evaluators have been chosen to assess the various TTS systems.

| No. of Evaluators | No. of Sentences | Voice | No. of TTS System for evaluation |
| --- | --- | --- | --- |
| 10 | 20 | Male/Female | 5 |

Table 3:  Data information for evaluation

The evaluators have been given the various sentences in order to grade MOS and DMOS. The entire assessment sheet, which includes the grades assigned for each sentence together with the average MOS and DMOS computation, has been generated. The full results sheet has been created in order to compare the ratings of various evaluators depending on many variables, such as geography, variations, etc.

## 5   Statistics of user preferences

The completed study demonstrates the assessors' technical educational backgrounds by demonstrating their ability to appropriately and adaptively assign grades to voice data. The various TTS systems have been evaluated using MOS and DMOS methodologies. The comparison was conducted using the mean MOS and DMOS values reported by specific TTS system evaluators. The complete charts of the user preferences w.r.t evaluation scores have been prepared.

## 6   Implementation to maintain the uniformity in the evaluation of Hindi TTS

This strategy results in the consistency of the ratings systems provided by various users from various places. Therefore, in the evaluation of Hindi text to speech system, the findings on the selected dataset for Hindi text to speech clearly shows the variations of MOS and DMOS. Therefore it is essential to identify the weight age in order to maintain the desired grading for the Hindi text to speech system.

### 6.1  Build Normalization model by AI techniques for unified Hindi text to speech system

The section presents the adopted experimental study of the different dialects by using linear regression technique. The technique figures out the appropriate weight for different Hindi dialects and proposed the normalized value in order to maintain the uniformity in the user scores while evaluating of the TTS system.

The supervised learning approach has been adopted to visualize the score of different dialects with reference

to the Khadi boli dialect as this is the predominant dialect of Hindi language and it is much closer to the modern Hindi.

The experimental work has been taken by comparing the scores of different dialects with Khadi Boli Dialect of the same sentences. Based on the dataset, the regression line has been plotted and figure out the slope i.e. a and intercept value i.e. b of the equation y= ax+b that can adjust the line to fit the data.

Linear regression model has been used to fit and train the data through which intercept and regressor coefficient has been inspected. Further, the score of dialects other than Khadi Boli has been predicted.

The regression coefficient (a) is the slope of the regression line which is equal to the average change in the dependent variable (b) for a unit change in the independent variable.

By applying the above defined model the predicted values have been examined initially for Haryanvi & Bundelkhandi dialects (x axis) with reference to Khadi Boli dialect (y axis):
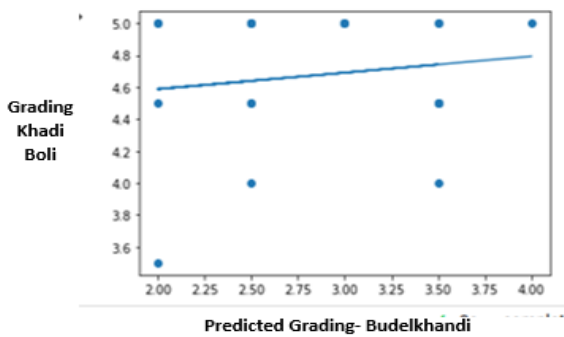
## 6.2 Results

**Bundelkhandi**



Fig 3: Regression: Khadi Boli-Bundelkhandi

b = regressor.intercept = 4.75297619

a = regressor.coef = -0.01190476

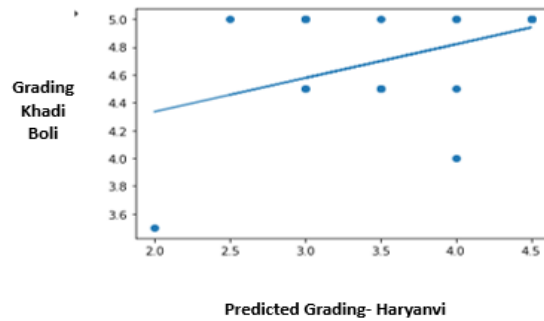|   | Actual | Predicted |
|---|--------|-----------|
| 0 | 3.5 | 4.729167 |
| 1 | 5.0 | 4.711310 |
| 2 | 4.5 | 4.711310 |
| 3 | 5.0 | 4.729167 |



Fig 4: Regression: Khadi Boli-Haryanvi

b = regressor.intercept = 4.56060606

a = regressor.coef = 0.07575758

|   | Actual | Predicted |
|---|--------|-----------|
| 0 | 3.5 | 4.712121 |
| 1 | 4.5 | 4.825758 |
| 2 | 5.0 | 4.863636 |
| 3 | 4.0 | 4.863636 |

## 7 Implementation of dialect characters & Word variations on Hindi text to speech system

The following methodology has been adopted in order to identify the different dialect patterns and initially character and word variations of Chhattisgarhi & Haryanvi dialects has been tested that helps in building dialect based TTS system:

**Identification of different patterns of Hindi dialects**

Based on the study and survey of different Hindi dialects, the following dialect variations and their patterns identified with respect to Hindi [23][24]:

| Chhattisgarhi | |
|---|---|
| **Patterns used in Chhattisgarhi** | **Examples** |
| ''ज' is replaced and pronounced by Hindi word ' झ' | जन - झन |
| 'द' is replaced and pronounced by 'ध'and 'क'replace by 'ख' | दौड़ –धौड़<br><br>इलाका — इलाखा |

| | |
|---|---|
| In an another variation 'चह'is replaced by Hindi word 'छे'and similar way 'स' is replaced by 'छ' | कचहरी– कछेरी<br><br>सीता – छीता |
| 'ब' is replaced and pronounced by Hindi word 'प' | शराब- शराप<br><br>खराब - खराप |
| 'म' is replaced and pronounced by Hindi word 'व' | नाम - नाव |
| 'ग' is replaced and pronounced by Hindi word 'क' | बंदगी- बंदकी |
| 'ए' is replaced and pronounced by Hindi word 'ये' | आए- आये |

### Awadhi

| Different Patterns | |
|---|---|
| 'ण' is replaced by 'न' | गुण - गुन<br><br>लक्ष्मण - लक्ष्मन |
| 'व' is pronounced like consonant Hindi Word 'ब' | वाहन - बाहन<br><br>व्याकुल- ब्याकुल<br><br>वन - बन |
| श' is pronounced like consonant Hindi Word 'स' | शंकर– संकर<br><br>शाम– साम<br><br>शेर– सेर |
| vowel 'इ'will be addition in 'स्' consonant | स्कूल - इस्कूल<br><br>स्त्री – इस्त्री |
| 'ऋ'word is replaced by 'र' , Hindi word 'ड़' replaced by 'र' and 'द' is replaced by word 'ड' | ऋषि - रिसि<br><br>किवाड़ - किवार<br><br>दंड - डंड |

### Hariyanvi

| Different Patterns | Examples |
|---|---|
| Single consonant is replaced by Dual Consonant. | भीतर –भित्तर, राजा -राज्जा, मित्र – मित्तर, गाड़ी – गाड्डी, आँसू-आँस्सू, आलू-आल्लू . ऊँचा-ऊँच्चा, चाचा-चाच्चा, ढीला-ढील्ला, हँसी-हँस्सी, |
| Hindi word 'थ'consonant is replaced by 'त' Consonant. | हाथ – हात<br><br>साथ – सात |
| 'न' is replaced with ण, 'श' is replaced with'स', 'ष' is replaced with 'स', 'इ' and'ल' is replaced with 'र' | अपने – अपणे, खाना – खाणा, चना – चणा, आकाश - आकास , सृष्टि– सिरस्टी, काला - कारा, कीड़ी – कीरी |
| **Augmentation of vowels** | अ : लग्न- लगन<br><br>इ : खजूर - खिजूर, जब –जिव, जवान - जिवान<br><br>उ : गवाही -गुवाही, जवाब -जुवाब, सपना –सुपना |
| Augmentation of consonants | च : ओछा –ओच्छा<br><br>त : चौथा –चौत्था<br><br>व : आना -आवणा, पीना –पीवणा |
| Disappearance of characters | अ : अठारह - ठारह, अनाज - नाज, अमावस—मावस , अहीर - हीर<br><br>इ : इकट्ठा -कट्ठा, कलियुग - कलयुग, घिसना -घसणा, जाति – जात<br><br>उ : उठाना -ठाणा, धातु –धात<br><br>व : ध्वजा -धजा, पाँव –पाँ<br><br>स : इकतीस-इकती, उन्नीस -उनी, छब्बीस –छब्बी<br><br>ह : अठारह –अठारा<br><br>अ is replaced with आ: अगला - आगला, ककड़ी - |

| | काकड़ी, लड्डी –लाड्डी |
| --- | --- |
| | अ is replaced with इ: अब -इब, अलावा - इलावा |
| | आ is replaced with ई :अब - ईब, अजगर - ईजगर |
| | उ is replaced with ऊ :कठपुतली -कठपूतली, चुंगी –चूंगी |
| | ओ is replaced with ऊ:गोंद -गूँद, क्यों –क्यूँ |
| | ऐ is replaced with इ :ऐसा - इसा, जैसा – जिसा |

| KhadiBoli | |
| --- | --- |
| **Different patterns** | **Examples** |
| Single consonant is replaced by Dual Consonant. | बेटा – बेड्डा, भेजा – भेज्जा, बड़ा – बड्डा, छोटा – छोट्टा, रानी – रान्नी, सादी – साद्दी, रोटी – रोट्टी |
| before Mahaprandhvani (महाप्राण ध्वनि( placing AlapapranDhvani (अल्पप्राणध्वनि( makes places. | देखा- देक्खा भूखा – भूक्खा |
| In the characteristic tendency of Balaghat, the preceding short vowel (Hrasv), diminishes or disappears in the form of dissipation | असाढ़- साढ़, उठाना-ठाना, इलाज-लाज, खूश्बू-खसबू |

| BrajBhasa | |
| --- | --- |
| **Different patterns** | |
| ड'and 'ल' sound is replaced by 'र' | पड़ेपर -काला – कारा, कीड़ी – कीरी , बिजली – बिजुरी , तले – तरे , सड़क – सरक , बल – बर |

| Bihari | |
| --- | --- |
| ड' sound is replaced by 'र'and vice versa | लड़का – लरका, कड़ी- करी, करेगा-कड़ेगा, कर- कड़, रात- ड़ात |

Table 4: Dialect Characters Variations

# 8 Implementation of dialect character & Words variations for Dialect based TTS system through SSML (Speech Synthesis Markup Language) Technology

In the evaluation process, the metadata of the evaluators should be prepared that covers evaluators from different regions of the different states where particular dialect exists.

SSML is a component of a larger collection of markup specifications for voice browsers created by the W3C using open processes. It is intended to offer a comprehensive, XML-based markup language to support the creation of synthetic speech in Web and other applications. The primary function of the markup language is to provide authors of synthesizable content with a uniform means of controlling various characteristics of speech output across various synthesis-capable platforms, including pronunciation, loudness, pitch, pace, etc. The pronunciation of different characters variations occur in a particular dialect should be reflected in the pronunciation element in building Hindi and dialect based TTS engines. Specific rules that give the content developer express control over pronunciation should be defined in the SSML elements[25][26][27].

# 9 Evaluations results & Findings

The SSML with all the identified variances has been implemented in the Hindi TTS and tested through the evaluation page of the text to speech system with Chhattisgarhi & Haryanvi dialect variations [28][29]. The different sentences have been developed that covers all the identified patterns for the appropriate results and findings. The dictionary of the dialects Chhattisgarhi and Haryanvi [30,31]initially of approx 1500 & 2100 words w.r.t. Hindi also developed for analysis of TTS systems. The metadata of the evaluators has been recorded in the database. The evaluators provide the grading on the speech generated by both the TTS system without knowing the type of the TTS. The main intension of this strategy is to record the scores given by same evaluators on the Hindi TTS and Hindi Text to speech system with dialect characters and word patterns in order to analyze the user preferences.

The following consolidate results cover the actual grading given by various evaluators and the variations among the grading given by different evaluators of different dialects have been generated in the Evaluation Process:

| Evaluator | Evaluation Score : Chhattisgarhi | | | Evaluation Score: Haryanvi | | |
|---|---|---|---|---|---|---|
| | Score : Hindi TTS system | Character variations | Word Variations using dictionary | Score : Hindi TTS system | Character variations | Word Variations using dictionary |
| Evaluator1 | 3.9 | 4.3 | 4.1 | 3.9 | 4.2 | 4.0 |
| Evaluator2 | 3.7 | 3.9 | 3.8 | 3.7 | 4.0 | 3.9 |
| Evaluator3 | 3.8 | 4.5 | 4.3 | 3.8 | 4.3 | 4.2 |
| Evaluator4 | 4.3 | 5 | 4.7 | 4.3 | 4.5 | 4.4 |
| Evaluator5 | 3.6 | 4.2 | 3.9 | 3.6 | 3.9 | 3.8 |
| Evaluator6 | 3.9 | 4.0 | 4.0 | 3.9 | 3.9 | 4.1 |
| Evaluator7 | 3.8 | 4.1 | 4.2 | 4.0 | 4.0 | 4.0 |
| Evaluator8 | 3.6 | 3.9 | 3.9 | 3.8 | 4.1 | 4.1 |
| Evaluator9 | 4.0 | 4.1 | 3.9 | 3.7 | 4.3 | 4.2 |
| Evaluator10 | 4.1 | 4.2 | 4.0 | 3.8 | 4.0 | 3.9 |

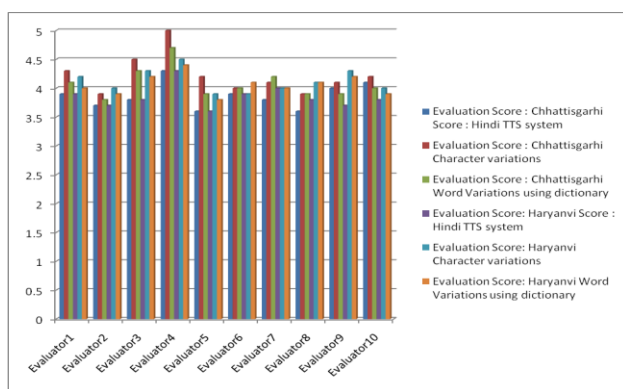Table 5: Evaluators Gradings



Fig 5: Grading Comparison: different dialects

## 10 Conclusion

The above result shows the differences in the scores of both the TTS system given by the Evaluators of the Chhattisgarhi & Haryanvi dialect with respect to the most prominent dialect of Hindi say KhadiBoli. Through this study, we examine that the preferences of the users have been changed in both the TTS systems based on their native place and dialects and that leads to the variations in the scores which are not feasible. So it is required to maintain the uniformity in the grading system by adopting the weightage and differences in the evaluation process of TTS system. Further, Based on the dictionary of Chhattisgarh and Haryanvi dialect as discussed in the above sections has also being created for the development of dialect based text to speech system. All the patterns have been implemented in the dialect based TTS and corresponding scores have been analyzed.

The findings will be used in order to maintain the uniformity in the evaluation of the Hindi TTS system. Also, the Chhattisgarhi & Haryanvi dialect based TTS system has been originated by using this approach. This approach can be adopted for the development of other dialect based TTS system. By covering the wider range of words corpora, better recording and evaluation, the naturalness and performance of the dialect based TTS system can be improved.

## References

1. Rob Clark, Hanna Silen, Tom Kenter, and Ralph Leith. 2019. Evaluating Long-form Text-to-Speech: Comparing the Ratings of Sentences and Paragraphs. arXiv:1909.03965 [cs, eess] (Sept. 2019). http://arxiv.org/abs/1909.03965 arXiv: 1909.03965.

2. Mirjam Wester, Cassia Valentini-Botinhao, and Gustav Eje Henter. 2015. Are We Using Enough Listeners? No!—An Empirically-Supported Critique of Interspeech 2014 TTS Evaluations. In Proc. Interspeech 2015. https://www.isca-speech.org/archive/interspeech_2015/papers/i15_3476.pdf

3. Introduction to Hindi Language, Yale University, 2022, https://hindi.yale.edu/language

4. Alan W Black and Keiichi Tokuda. 2005. The Blizzard Challenge - 2005: Evaluating corpus-based speech synthesis on common datasets. In Proc. Interspeech 2005. 77–80.

5. Maxine Eskenazi, Gina-Anne Levow, Helen Meng, Gabriel Parent, and David Suendermann. 2013. Crowdsourcing for speech processing: Applications to data collection, transcription and assessment. John Wiley & Sons.

6. Simon King. 2014. Measuring a decade of progress in text-to-speech. Loquens 1, 1 (2014), 006.

7. Petra Wagner, Jonas Beskow, Simon Betz, Jens Edlund, Joakim Gustafson, Gustav Eje Henter, Sébastien Le Maguer, Zofia Malisz, Éva Székely, Christina Tånnander, and Jana Voße. 2019. Speech Synthesis Evaluation — State-of-the-Art Assessment and Suggestion for a Novel Research Program. In Proc. 10th ISCA Speech Synthesis Workshop. 105–110. https://doi.org/10.21437/SSW.2019-19

8. Cambre Julia, et.al, Choice of Voices: A Large-Scale Evaluation of Text-to-Speech Voice Quality for Long-Form Content, 2020, In CHI Conference on Human Factors in Computing Systems,

https://dl.acm.org/doi/fullHtml/10.1145/3313831.3376789

9. C. Stevens, et al., "On-line experimental methods to evaluate text-to-speech (TTS) synthesis: effects of voice gender and signal quality on intelligibility, naturalness and preference," Computer Speech and Language, vol. 19, pp. 129-146, 2005.

10. D. B. Pisoni, et al., "Perception of synthetic speech generated by rule," in Proceedings of the IEEE, 1985, pp. 1665-1676.

11. H. A. Sydeserff, et al., "Evaluation of speech synthesis techniques in a comprehension task," Speech Communication, vol. 11, pp. 189-194, 1992.

12. ChangYu-Yun, Evaluation of TTS Systems in Intelligibility and Comprehension Tasks, aclanthology, https://aclanthology.org/O11-1004.pdf

13. K. Yorkston, et al., "Comoprehensibility of dysarthric speech: Implications for assessment and treatment planning," American Journal of Speech-Language Pathology, vol. 5, pp. 55-66, 1996.

14. K. C. Hustad, "The relationship between listener comprehension and intelligibility scores for speakers with dysarthria," Journal of Speech, Language, and Hearing Research, vol. 51, pp. 562-573, 2008

15. K. Yorkston, et al., "Comoprehensibility of dysarthric speech: Implications for assessment and treatment planning," American Journal of Speech-Language Pathology, vol. 5, pp. 55-66, 1996.

16. K. C. Hustad, "The relationship between listener comprehension and intelligibility scores for speakers with dysarthria," Journal of Speech, Language, and Hearing Research, vol. 51, pp. 562-573, 2008

17. Paper from IIT Madras. Measuring Quality for Text-to-Speech Systems using degraded MOS scale and Word Error Rate

18. Hindi Language Spoken States in India, WorldListMinia, 2021, https://www.worldlistmania.com/hindi-speaking-states-in-india/

19. List of Hindi speaking states of India, https://www.learnsabkuch.in/2017/08/list-of-hindi-speaking-states-of-india.html

20. https://start.mgkvp.ac.in/Uploads/Lectures/18/6715.pdf

21. https://mdu.ac.in/UpFiles/UpPdfFiles/2020/Jan/bhasavigyan-hindi%20bhasa-final.pdf

22. https://start.mgkvp.ac.in/Uploads/Lectures/18/6715.pdf https://start.mgkvp.ac.in/Uploads/Lectures/18/6715.pdf

23. https://dl.acm.org/doi/fullHtml/10.1145/3313831.3376789#BibPLXBIB0050

24. Requirements in PLS, SSML and SRGS Standard—Hindi as a Case Study, 2021, Smart Systems: Innovations in Computing pp 113–120, Springer series

25. Indian Languages Requirements for String Search/comparison on Web, 2021, AIST 2021: Artificial Intelligence and Speech Technology pp 210–214, Springer

26. W3CPLS1.0: https://www.w3.org/TR/pronunciation-lexicon/#S1

27. W3CSSML: https://www.w3.org/TR/speech-synthesis11/

28. Evaluation Test(Chhatisgarhi: http://tdil-dc.in/ttsapi/tts_evaluation_c/index.php)

29. Evaluation Test Hindi: https://tdil-dc.in/ttsapi/tts_evaluation_h/index.php

30. Chattisgarhi Dictionary: https://scert.cg.gov.in/pdf/

31. Haryanvi & Hindi Dictionary: https://www.exoticindiaart.com/book/details/haryanvi-hindi-dictionary-rzz882/