# Multifractal Detrended Cross-Correlation Analysis for Characterization of Spoken Language – A New Method to explore the Genesis of Languages

**Suparna Panchanan[1], Nilav Darsan Mukhopadhyay[2], Ranjan Banerjee[3], Most Mahabuba Islam[4], Shankha Sanyal[5], Dipak Ghosh[6], Debmalya Mukherjee[7]**

[1]Computer Science and Engineering - CS & DS, Brainware University

Email ID: suparna_mou2k@yahoo.co.in

[2]Computer Science and Engineering - CS & DS, Brainware University

Email ID: contactnilav@gmail.com

[3]Computer Science and Engineering, Brainware University

Email ID: rnb.cse@brainwareuniversity.ac.in

[4]Computer Science and Engineering(AI & ML), Brainware University

Email ID: bwubta21092@brainwareuniversity.ac.in

[5]Sir C. V. Raman Centre for Physics and Music, Jadavpur University

Email ID: ssanyal.sanyal2@gmail.com

[6]Sir C. V. Raman Centre for Physics and Music, Jadavpur University

Email ID: deegee111@gmail.com

[7]Computational Sciences Department, Brainware University

Email ID: dbm.cs@brainwareuniversity.ac.in

## ABSTRACT

This work presents a novel use of chaos-based non-linear techniques coupled with statistical methods for spoken language characterization in the speech signal domain. Our goal is to create a framework that highlights linguistic commonalities. While Multifractal Detrended Cross-Correlation Analysis (MFDXA) assesses long-term cross-correlations between languages using the cross-correlation coefficient as a measure of similarity, Multifractal Detrended Fluctuation Analysis (MF-DFA) is used to examine linguistic correlations among the languages.

Bengali, Assamese, Maithili, Odia, Nepali, Manipuri, Hindi, Urdu, Marathi, Gujarati, Punjabi, Konkani, Tamil, Telugu, Malayalam, Kannada, and Sanskrit are among the seventeen Indian languages that are the emphasis of the present work. The speech corpus, which includes speakers of both sexes, is mostly composed of unplanned conversational material on a variety of subjects, including social welfare, agriculture, and in-person interviews. This model is unique in that it avoids the conventional use of linguistic information. Our findings reveal notable deviations from established linguistic theories in cases such as Bengali-Gujarati, Hindi-Tamil, and Bengali-Kannada, indicating that current language classifications may benefit from re-evaluation. Statistical tools like ANOVA and Mahalanobis Distance have been used to validate the current study. The proposed method offers a valuable approach to investigating the origins of languages within a global framework.

*Keywords: Indian Languages, MFDFA, MFDXA, Cross-Correlation, Language Genesis, ANOVA, Mahalanobis Distance*

## 1. INTRODUCTION

Spoken language is the most practical form of human communication. To capture the diversity of languages and their speakers globally, languages are often grouped into large families. A language family consists of languages that share a common ancestral mother tongue. Although there are over 5,000 languages spoken worldwide, researchers typically classify them into fewer than twenty families.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

In linguistics, languages are categorized in two main ways: typologically and genetically (or genealogically). Typological classification groups languages based on structural features, while genetic classification organizes them by their shared historical origins. Languages within a family often exhibit similarities in vocabulary, sounds, or grammatical structures. Population growth and cross-cultural interactions often lead to a blending of native (deshi) and borrowed (bideshi) words within a mother language. Over time, this fusion contributes to the evolution of new languages from the original mother tongue. In India, it is reported that over 1,500 languages, both official and unofficial, are spoken (Office of the Registrar General & Census Commissioner). These languages are primarily classified into two major language families: Indo-Aryan and Dravidian. In addition to the Indo-Aryan and Dravidian language families, India is also home to Tibeto-Burman and Austro-Asiatic language families. The phonemes in languages within each of these groups are quite similar, and their sequences are often comparable. This highlights the shared characteristics among many Indian languages. Other factors that contribute to linguistic similarities include political influence, regional tourism, and trade connections. Additionally, studying linguistic origins and commonalities provides important information about historical commerce networks, their routes, and international diplomatic ties. Consequently, scholars have shown a great deal of interest in the commonalities between languages. Several important research on language similarity are examined in the section that follows.

Cognates are crucial in assessing language similarity. McMahon and McMahon (McMahon & McMahon, 2005) employed Phylogenetic Analysis to explore relationships between languages. ALINE (Kondrak, 2000) distance-based similarity metric was established by Downey et al. (Downey et al., 2008), the phoneme matching algorithm is used to classify cognates. K. Kettunen et al. (Kettunen et al., 2006) used a file compression method to compare various European languages based on textual content, while Nathaniel Oco et al. (Oco et al., 2013) created a trigram frequency profile for language comparison. While these studies primarily focus on textual data, Bradlow et al. (Bradlow et al., 2010) analyzed language similarity based on phonetic similarity derived from spoken language.

Sengupta et al. (Sengupta & Saha, 2015) examined linguistic similarities between Indian languages using a language verification framework. Their approach integrated two modeling techniques—the Gaussian Mixture Model (GMM) and Support Vector Machine (SVM)—with two feature extraction methods: Mel Frequency Cepstral Coefficients (MFCC) combined with Shifted Delta Coefficients (SDC), and Speech Signal-based Frequency Cepstral Coefficients (SFCC) with SDC. This integration led to the development of four distinct language verification frameworks for assessing language similarities.

Our objective is to develop a generalized framework to assess similarities among spoken languages through a non-linear methodology. This approach effectively captures signal characteristics, recognizing that vocal fold vibrations are greatly influenced by nonlinearity in both tissue properties and airflow dynamics (Titze, 1995). Previous studies on language similarity have largely overlooked the nonlinearity inherent in speech signals. Our current work seeks to address this gap, aiming to make a substantial contribution to this field.

Speech signals and other naturally occurring geometries cannot be accurately described by a single scaling ratio. Different parts of natural systems are subjected to varying scaling ratios to achieve a detailed analysis, leading to a pattern of non-uniform self-similarity throughout the system. Such systems are better referred to as "multifractal." There are two types of multifractal signals: monofractal and multifractal. For voice signal analysis, Multifractal Detrended Fluctuation Analysis (MFDFA) is more effective than Detrended Fluctuation Analysis (DFA) because it can handle both small and large signal fluctuations. This method works well for examining non-stationary signals' multifractal scaling behaviour. While Multifractal Detrended Cross-Correlation Analysis (MF-DXA) (Ghosh et al., 2014; He & Chen, 2011; Jiang & Zhou, 2011; Wang et al., 2013) offers a more thorough investigation of the multifractality in two cross-correlated signals, Detrended Cross-Correlation Analysis (DCCA) can be used to detect long-term cross-correlations between two non-stationary time series (Hedayatifar et al., 2011; Podobnik, Grosse, et al., 2009; Podobnik et al., 2008; Podobnik, Horvatic, et al., 2009; Podobnik et al., 2011; Podobnik & Stanley, 2008; Xu et al., 2010). Autocorrelation and cross-correlation are two crucial characteristics of real-time signals that have been demonstrated to be accurately captured by the DCCA technique (Horvatic et al., 2011). This study examines the commonalities between Indian languages as well as the existence of multifractality. The multifractality of spoken language is measured using MFDFA, while power-law cross-correlations, or the similarity between two languages, are examined using the MF-DXA technique (Zhou, 2008). We've chosen seventeen Indian languages for this. This method is used to find correlation coefficients for every conceivable pair of languages ($\gamma_x$). Importantly, the similarity between musical signals was successfully measured using the cross-correlation coefficient (Sanyal et al., 2016). This suggests that the degree of correlation between the language pairs may be revealed. A smaller value of $\gamma_x$ implies a strong correlation between the data series, while a value of 1 represents uncorrelated data (Ghosh et al., 2014). A negative value of $\gamma_x$ indicates a very high degree of correlation between the data series.

In summary, we propose an automated scientific algorithm that not only measures complexity or multifractality but also assesses the similarity of 17 Indian languages: Assamese, Bengali, Hindi, Marathi, Gujarati, Punjabi, Urdu, Malayalam, Odia, Konkani, Maithili, Kannada, Manipuri, Nepali, Tamil, Telugu, and Sanskrit. For each language, 10 segments of spoken

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

data, each lasting 10 seconds, are analyzed. Multifractal width measures the multifractality of spoken language using MFDFA. Examine the similarity between two languages, using MF-DXA. Calculate cross-correlation coefficients for all possible pairs of the selected languages to provide a measure of their similarity. Compare the findings with existing theories of language to see if they align or differ, indicating the need to re-examine existing language frameworks. Incorporate one-way ANOVA and Post-hoc Tukey's HSD tests to validate the results. Use Mahalanobis Distance (MD) to find the distance among the languages.
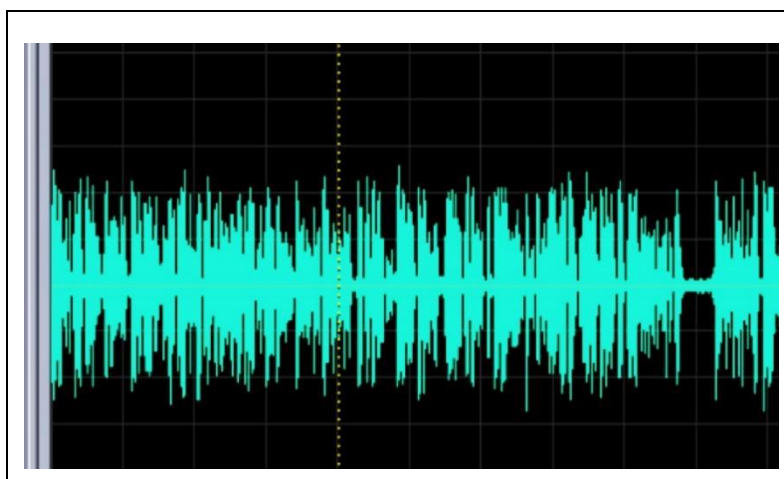
Employ these techniques to trace the origins of languages and determine relationships between them. These objectives outline the complete approach taken in the study to analyze and compare the multifractality and similarity of Indian languages using advanced statistical and mathematical methods. In some cases, the findings are consistent with current theories of language, but in others, they are not. As a result, this non-linear analysis indicates that existing language frameworks need to be re-examined.

Therefore, this technique can be used as a scientific tool to trace the origins of languages and find relationships between them.

Through a systematic methodology and a critical examination of the findings, the rest of the paper demonstrates the viability of this concept.

## 2. DATABASE

The availability of a comprehensive speech database is essential for the development of any speech-based application, and this aspect holds significant importance in the current context. Several factors must be considered when collecting speech data, including the diversity of informants and the inclusion of various environmental elements, such as the surrounding ambiance. For this study, a speech database was compiled for each of the seventeen languages, covering seventeen standard Indian dialects. The data was sourced in mp4 format from YouTube. Com, providing a rich and diverse set of speech samples. The time domain representation of the signal is given in Figure 1.



**Fig.1 Time domain representation of spoken language**

The richness of the speech database is increased by the fact that these dialects come from different language families. The corpus mostly consists of material from live shows, talk shows, interviews, and news bulletins with both male and female presenters. Table 1 gives an overview of the database.

The "Total Audio MP3 Converter®" program was used to extract the audio data from the videos, and it was subsequently converted into .wav format. The conversion technique used 16-bit mono encoding and a sampling rate of 22,050 samples per second. An initial noise reduction was performed using the "Cool Edit 2000®" program. It should be noted that the data retrieved in this manner represents a more realistic environment, as it underwent several compression and decompression processes. Consequently, data gathered in such settings will be more beneficial for developing a functional system compared to data collected in a controlled environment.

**Table 1: Database**

| Name of Language | Abbreviation | Language Family | Duration in Minutes |
|---|---|---|---|
| Bengali | Ben | Indo Aryan Eastern | 39 |

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam,
Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

| Assamese | Asa | Indo Aryan Eastern | 30 |
|---|---|---|---|
| Maithili | Mat | Indo Aryan Eastern | 35 |
| Odia | Odi | Indo Aryan Eastern | 35 |
| Nepali | Nep | Indo Aryan Northern | 42 |
| Manipuri | Man | Tibeto-Burman | 39 |
| Hindi | Hin | Indo Aryan | 40 |
| Urdu | Urd | Indo Aryan Southern | 31 |
| Marathi | Mar | Indo Aryan Southern | 25 |
| Gujarati | Guj | Indo Aryan Western | 30 |
| Punjabi | Pun | Indo Aryan North-western | 40 |
| Konkani | Kon | Indo Aryan Southern | 42 |
| Tamil | Tam | Dravidian | 34 |
| Telugu | Tel | Dravidian | 35 |
| Malayalam | Mll | Dravidian | 31 |
| Kannada | Kan | Dravidian | 32 |
| Sanskrit | San | Indo-European | 32 |

## 3. EXPERIMENTAL DETAILS

For this study, ten segments, each consisting of ten seconds of spoken data, are analyzed from each of the seventeen languages. In speaker recognition, a frame length of 50–100 ms is typically sufficient for multifractal speech signal analysis (DC González, 2012). The Hurst exponent, multifractal width, and auto-correlation are calculated using Multifractal Detrended Fluctuation Analysis (MFDFA), while the cross-correlation between language pairs is assessed using Multifractal Detrended Cross-Correlation Analysis (MF-DXA). The next section provides an overview of the MFDFA and MF-DXA techniques.

## 4. METHODOLOGY

### 4.1. Multifractal detrended fluctuation analysis (MFDFA)

Here, MATLAB is applied to inspect the speech data. The mathematical framework is created on the expressions proposed by Kantelhardt et al. (Kantelhardt et al., 2002). The complete procedure is outlined below:

Step 1: Signals' noise looks alike structure is transformed into a random walk-like signal. It can be denoted as:

$$y_i = \sum (x_k - \bar{x}) \tag{1}$$

$\bar{x}$ is the representation of the mean value of the signal.

Step 2: The speech signal is partitioned into $N_s$ numbers of frames. Each frame contains of a certain number of samples. Therefore, $N_s$ will be the number of frames generated from the original signal with length $N$ and $s$ samples/frame.

$$N_s = \text{int}(\frac{N}{s}) \tag{2}$$

Step 3: The function $F(s,v)$ denotes the variation of local RMS with sample size $s$.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

$$F^2(s,v) = \frac{1}{s}\sum_{i=1}^{s}\{y[(v-1)s+i] - Y_v(i)\}^2$$

(3)

For $v = N_s + 1 \ldots\ldots 2 N_s$, where $Y_v(i)$ is the least square fitted value in the bin of size v. The first-order polynomial (MFDFA-1) is considered to get the least square linear fit. Higher-order extension is possible by fitting higher-order polynomials.

Step 4: Equation (4) is used to find out the overall RMS variation of q-order for various scale sizes

$$F_q(s) = \{\frac{1}{N_s}\sum_{v=1}^{Ns}[F^2(s,v)]^{\frac{q}{2}}\}^{(\frac{1}{q})}$$

(4)

The index $q$ can be any value other than zero as, $\frac{1}{q}$ equal to infinite is not allowed.

Step 5: The scaling nature of the fluctuation function is observed from the log-log graph of $F_q(s)$ Vs. $s$ for different values of $q$.

$$F_q(s) \sim s^{h(q)}$$

(5)

The $h(q)$ is named the generalized Hurst exponent. The Hurst exponent is used in fractal analysis to comprehend a time series' correlation and self-similarity characteristics. The existence of a long-range correlation is measured using $h(q)$. It is unique for all values of $q$ in case of a monofractal time series.

The classical scaling exponent $\tau(q)$ is connected to the generalized Hurst exponent $h(q)$ of MFDFA by the relation

$$\tau(q) = qh(q) - 1$$

(6)

The long-range correlation of a mono-fractal series depends linearly on $\tau(q)$ with a single Hurst exponent H.

In contrast, the multifractal signal possesses multiple Hurst exponents. Here $\tau(q)$ is a non-linear function of $q$ (Ashkenazy et al., 2003).

The relation of singularity spectrum f (α) with $h(q)$ is given below

$$\alpha = h(q) + qh'(q)$$

(7)

where $\alpha$ is the singularity strength.

$$f(\alpha) = q[\alpha - h(q)] + 1$$

(8)

$f(\alpha)$ defines the dimension of a subset series. The multifractal spectrum illustrates the significance of various fractal exponents within a time series. Physically, the width of the spectrum represents the range of these exponents. A least square fitting technique using a quadratic function $f(\alpha)$ (Shimizu et al., 2002) at maximum $\alpha = \alpha_0$, will give the information about the spectra.

$$f(\alpha) = A(\alpha - \alpha_0)^2 + B(\alpha - \alpha_0) + C$$

where C is an additive constant and it is given as $C = f(\alpha_0) = 1$. $B$ indicates the asymmetry of the spectrum. It is zero for the symmetric spectrum.

The constant A plays a crucial role in defining the nature and orientation of the parabola. To measure the spectrum width, the fitted curve is extrapolated to zero. The width $W$ is as follows,

$$W = \alpha_1 - \alpha_2 \text{, with } f(\alpha_1) = f(\alpha_2) = 0$$

(9)

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

$W$ is used as a metric of the multifractal spectrum. Clearly, with increasing width $W$, the multifractality of the spectrum increases. In the case of mono-fractal series, where $h(q)$ does not vary with $q$, the width $W = 0$.

## 4.2. Multifractal Detrended Cross-Correlation Analysis (MF-DXA):

We have followed Zhou's (2008) (Zhou, 2008) prescription to investigate the correlation between different pairs of languages through the study and analysis of cross-correlation.

Here we consider $xx(i)$ and $yy(i)$ are two time series of two different languages of length $N$. Hence, the average value of the two-time series can be indicated as follows

$$x_{avg} = 1/N\sum_{i=1}^{N} xx(i) \quad \text{and} \quad y_{avg} = 1/N\sum_{i=1}^{N} yy(i) \tag{10}$$

$XX(i)$ and $YY(i)$ can be evaluated from equation (10)

$$XX(i) = \left[\sum_{k=1}^{i} xx(k) - xx_{avg}\right] \text{ for i=1.... } N \;; \tag{11}$$

$$YY(i) = \left[\sum_{k=1}^{i} yy(k) - yy_{avg}\right] \text{ for i = 1 .... } N \;; \tag{12}$$

The summation helps to reduce the noise level embedded in data. Each of the time series $XX(i)$ and $YY(i)$ were split up into $N_s$ non-coinciding bins where

$$N_s = \text{int}(\frac{N}{s}) \quad \text{where } s \text{ is the length of the bin.}$$

Now, as $N$ is not an integer multiple of $s$, a small portion of the series is clipped at the end. The entire portion of the signal was mirrored from the opposite side to include the clipped section, ensuring the full signal is analyzed. As a result, the total number of bins becomes $2N_s$. A least squares linear fit was applied within each bin to obtain the fluctuation function, which is denoted as:

$$F(s,v) = 1/s\sum_{i=1}^{s} (YY[(v-1)s+i] - yy_v(i))\{XX[v-1)s+i] - xx_v(i)\}$$

For each bin $v$, $v = 1........N_s$ and

$$F(s,v) = 1/s\sum_{i=1}^{s} \{YY(v-N_s)s+i] - yy_v(i)\}\{XX[N-(v-N_s)s+i] - xx_v(i)\}$$

For each bin $v$, $v = N_s+1.......2N_s$ where $x_v(i)$ and $y_v(i)$ are the least square-fitted values in the bin v. The $q$ th order detrended covariance $F_q(s)$ is obtained after averaging over $2N_s$ bins.

$$F_q(s) = \{1/2N_s\sum_{v=1}^{2N_s}[F(v,s)]^{q/2}\}^{1/q} \tag{13}$$

where $q$ can have any value other than zero. The procedure can be repeated by varying the value of $s$. $F_q(s)$ rises when the value of $s$ increases. If there exists a power correlation over a long-range, $F_q(s)$ can be represented using power law as shown below -

$$F_q(s) \sim s^{\lambda(q)}$$

In this kind of scaling relation, $\log\left(F_q(s)\right)$ linearly depends on log(s), where $\lambda(q)$ is the slope of the straight line. $\lambda(q)$

depicts the degree of the cross-correlation among the data series under consideration. As $F_q$ has a singularity at $q = 0$, $F_q$ can be determined using a logarithmic averaging method.

$$F_0(s) = \{1/4N_s \sum_{v=1}^{2N_s} [F(s,v)]\} \sim s^{\lambda(0)}$$

(14)

The technique takes the form of standard detrended cross-correlation analysis (DCCA) when $q = 2$. The cross-correlation between two time series will be mono-fractal if the scaling exponent $\lambda(q)$ is independent of $q$ and it will be multifractal if the scaling exponent $\lambda(q)$ is dependent on q. Furthermore, when q is greater than zero, the scaling behavior of the segments with large fluctuations is demonstrated by $\lambda(q)$ whereas the slight fluctuations can be described by $\lambda(q)$ with negative $q$. The value λ(q) = 0.5 signifies the lack of cross-correlation. For consistent long-range cross-correlations, it is found that $\lambda(q)$ is greater than 0.5 where a large value in one variable indicates the higher value of another variable, antithetically λ(q) less than 0.5 symbolizes a cross-correlation that is anti-persistent. Here inverse relation is observed between the variables (Movahed & Hermanis, 2008).

If two-time series are fabricated using binomial assessment of the p-model, then as per Zhou (2008), the following can be written (Zhou, 2008):

$$\lambda(q=2) \approx [h_x(q=2) + h_y(q=2)]/2$$

(15)

Podobnik and Stanley (Podobnik, Grosse, et al., 2009) examined this relationship in the context of mono-fractal Autoregressive Fractional Integrated Moving Average (ARFIMA) signals and EEG time series, with the parameter $q = 2$. In this scenario, there is autocorrelation within the individual time series. However, a specific exponent does not show any power-law cross-correlation when the series is modeled using two uncoupled ARFIMA processes (Movahed & Hermanis, 2008). This suggests that while inherent autocorrelation exists within each signal, cross-correlation between them does not follow a power-law distribution unless there is a coupling or interaction between the processes.

The autocorrelation function can be denoted as

$$C(\tau) = \left\langle \left[ xx(i+\tau) - \langle xx \rangle \right] \left[ xx(i) - \langle xx \rangle \right] \right\rangle \sim \tau^{-\gamma}$$

(16)

The cross-correlation function, by definition, takes the following form -

$$C_x(\tau) = < [xx(i+\tau) - < xx >][yy(i) - < yy >] > \sim \tau^{-\gamma_x}$$

(17)

Where $\gamma$ and $\gamma_x$ are the auto correlation coefficient and cross-correlation exponent, respectively.

The autocorrelation exponent is evaluated from the DFA method and it can be denoted as $\gamma = 2 - 2h(q=2)$ (Movahed & Hermanis, 2008). Recently, Podobnik et al. (Podobnik et al., 2008) have established the relation between cross-correlation exponent $\gamma_x$ and scaling exponent $\lambda(q=2)$

$$\gamma_x = 2 - 2\lambda(q=2)$$

(18)

The lower value of $\gamma$ and $\gamma_x$ shows the correlation of signals whereas, $\gamma_x$ is 1 or greater than 1 for uncorrelated data. Predominantly, the presence of multifractality is indicated by the dependency of $\lambda(q)$ on $q$. In this study, we aim to examine the cross-correlation at different time scales among speech signals to quantify the similarities between the languages. By analyzing how the fluctuations evolve across these time scales, we can uncover patterns of self-similarity and multifractality that reflect the distinct or shared characteristics of the languages under investigation.

The parameters, Hurst exponent, Autocorrelation and Multifractal width under the investigation are compared using the analysis of variance (ANOVA) technique to determine whether their means differ significantly. This comparison is conducted with a 95% confidence level and a p-value threshold of 0.05.

Mahalanobis distance (MD) have been used to measure the spatial distance between each language pair.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

$$MD = \sqrt{(x - \mu)^T \sum{}^{-1} (x - \mu)}$$

(19)

Here, $x$ represents the vector point of interest, $\mu$ is the mean vector of the distribution and $\sum{}^{-1}$ is the inverse of the covariance matrix. $T$ stands for transpose.

## 5. RESULT AND DISCUSSION

This research utilizes multifractal analysis to investigate the multifractal properties of spoken languages and their similarities among seventeen languages. Traditional nonlinear methods, like Lyapunov exponents and correlation dimension, are often noise-sensitive and require stationary conditions. In contrast, the MFDFA technique objectively quantifies the complexity and information content of speech signals. Given that most Indian languages belong to either the Indo-Aryan or Dravidian language families, the analyzed signals exhibit a complex self-similar structure. Various exponents of the power-law spectrum are applied in this analysis, as a single power-law exponent cannot sufficiently capture the self-similar nature of these signals. This method is effective for determining whether a signal exhibits self-similarity. Section 3 provided a detailed description of the technique. To reduce noise, the data sets were modified using Eqs. (2) and (3). The integrated time series were then divided into $N_s$ bins, where $N_s = \text{int}(\frac{N}{s})$ with s being the scale size and $N$ the series length. For values of $q$ ranging from -5 to +5, the $q^{th}$ order detrended covariance $F_q(s)$ was calculated using Eq. (4). Power law scaling of $F_q(s)$ with s is observed for all. To determine the spectrum in MFDFA, the local RMS variation is combined into an overall RMS. The overall RMS is further expanded into a higher-order RMS to differentiate between segments of varying sizes. For positive $q$, the average $F_q(s)$ is subjected by the segments with large variance $F^2(s, v)$. For positive $q$, $h(q)$ tells the scaling behavior of the segments with large fluctuations. The average $F_q(s)$ for $q$ is less than zero will be controlled by the segments with a small variance $F^2(s, v)$. Therefore, the negative values of $q$ have an impact on the segments that exhibit minor fluctuations, and corresponding $h(q)$ describes the scaling behaviour of those segments.

This technique offers significant advantages for analyzing speech data and other non-stationary time series. The multifractal width $(W)$ is a key measure that quantifies the complexity of spoken languages, which can be represented as time series with substantial nonlinearity. A higher value of $W$ indicates greater local variations in the time scale, making it particularly useful for distinguishing and characterizing a specific language. In contrast, a lower $W$ reflects smaller local changes in the temporal scale. Consequently, two signals with similar complexity (i.e., exhibiting similar fluctuations in the local scale) will display comparable $W$ values in the temporal domain. Therefore, the multifractal spectral width is a valuable metric for describing the distinctive characteristics of a language.

For $q = 2$, $h(q)$ becomes H(q). When H(q) equals 0.5, a signal can be thought of as an independent random process. The signal has long-range anti-correlations if $H(q)$ is less than 0.5. Conversely, long-term positive correlations are characterized by 0.5<H(q) <1. The exponent H(q) is equivalent to the well-known Hurst exponent, which is widely used to quantify long-term memory or persistence in a time series.

Our objective in this study is to quantify the variation of multifractality in different spoken languages of India and to re-evaluate the differences. The database includes multiple speakers, featuring content from news bulletins, talk programs, interviews, live shows, and more; however, we were unable to determine the speakers' age ranges. Therefore, we aim to use the nonlinear technique, Multifractal Detrended Fluctuation Analysis (MFDFA), to characterize the languages based on speaker-independent information. This study employs the Hurst index, auto-correlation, and multifractal breadth to describe the languages.
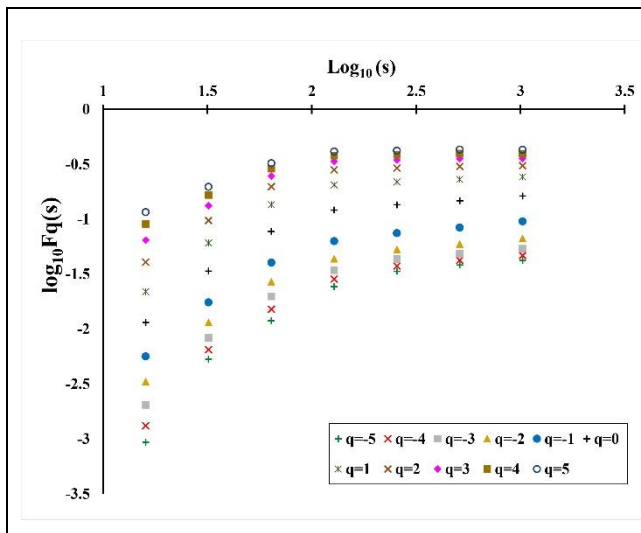
Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

**Fig. 2(a) Log$_{10}$(s) Vs Log$_{10}$F$_q$(s) for Sanskrit**



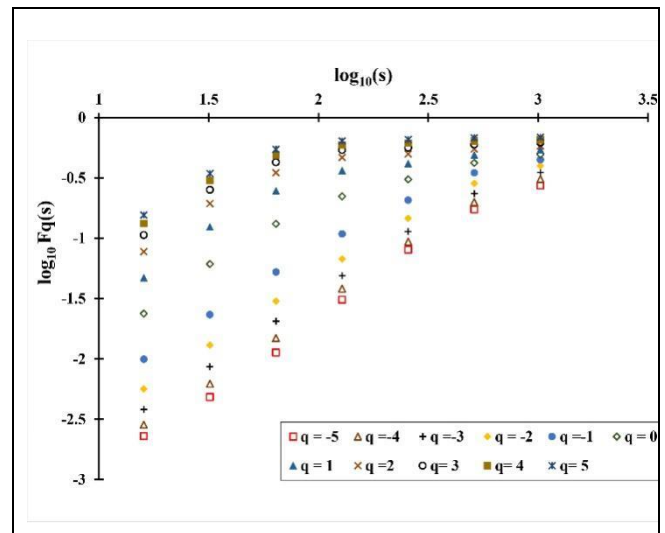**Fig. 2(b) Log$_{10}$(s) Vs Log$_{10}$F$_q$(s) for Bengali**

Figure 2(a) and Figure 2(b) show the regression plot of $\log_{10}(s)$ versus $\log_{10}(F_q(s))$ for different values of q for San and Ben, respectively. The linear relationship observed in these plots indicates a scaling behavior between $\log_{10}(F_q(s))$ and $\log_{10}(s)$. The slope of the linear fit in Figure 2 corresponds to the calculated value of $H_q$.

A constant $H_q$ indicates mono-fractality, whereas a dependence of $H_q$ on q reveals a multifractal nature in the series, as illustrated in Figure 3.
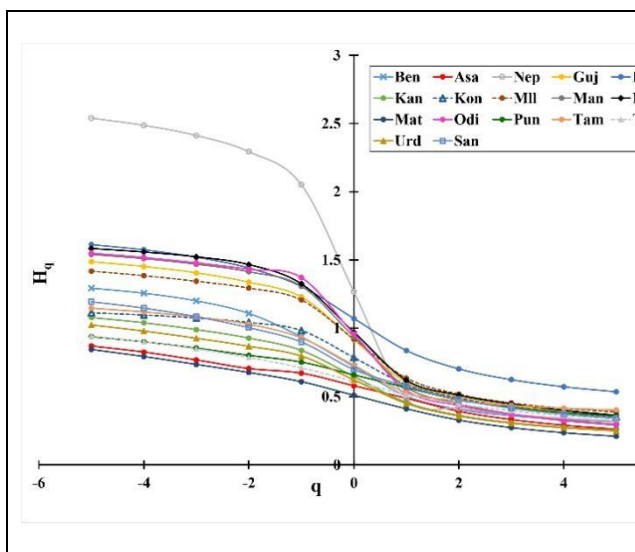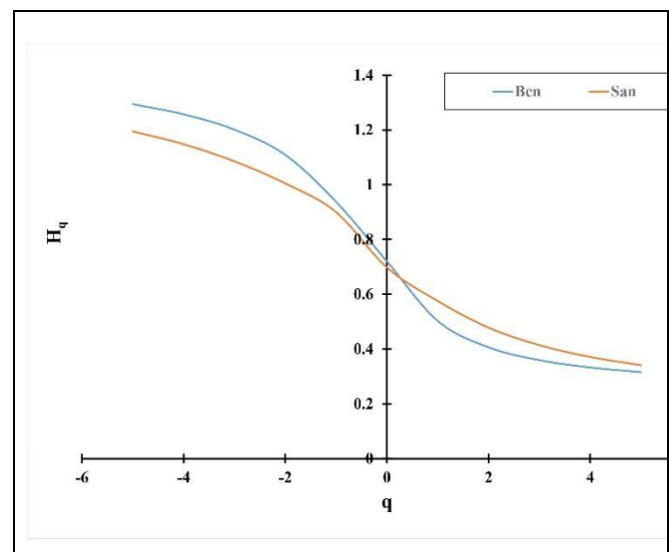


**Fig. 3(a) The variation of q Vs H$_q$**



**Fig. 3(b) The variation of q Vs H$_q$ of Bengali and Sanskrit**

Figure 3(a) illustrates the variations in q and $H_q$ across all seventeen languages, while Figure 3(b) specifically highlights these variations for the Ben and San languages. It is evident from both Figure 3(a) and Figure 3(b) that $H_q$ decreases as q increases.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam,
Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

**Fig. 4(a) Multifractal width for seventeen languages**



**Fig. 4(b) Multifractal width for Bengali**



**Fig. 4(c) Multifractal width for Sanskrit**

Figure 4(a) presents the multifractal widths across seventeen languages, while Figures 4(b) and Figures 4(c) illustrate the widths for Ben and San, respectively. The variation in width for each language indicates significant complexity in the voice signal. For example, the width for Ben ranges from 0.865 to 1.049, whereas for San, it spans from 0.924 to 1.449. In this study, ten sets of spoken data, each lasting ten seconds, are analyzed for each language to calculate an average width value. Figure 5 further shows the multifractal widths across five sections for each of the seventeen languages. The results indicate that complexity varies across different sections within the same language. This variation stems from changes in quasi-periodic and quasi-random signals, quiescent intervals, consonant-vowel (C-V) transitions, and the distribution of pauses within sentences.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee



**Fig. 5 Multifractal width at different parts of the seventeen languages**

The analysis shows that the multifractal width exceeds 0.9 for most languages, indicating a high degree of complexity. Autocorrelation proves to be a valuable measure for assessing the self-similarity of the signal. Table 2 summarizes the computed averages and standard deviations of the Hurst exponent, autocorrelation, and multifractal width across the seventeen languages.

**Table 2. Values of Hurst exponent, Autocorrelation and Multifractal width of seventeen languages**

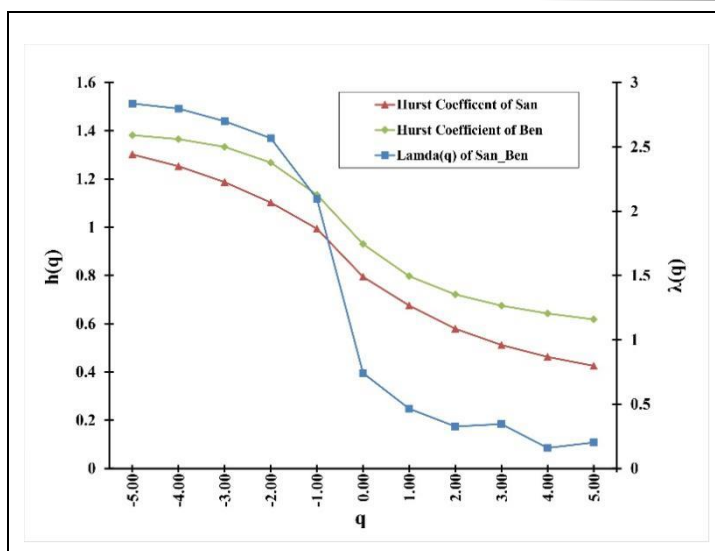| Languages | Hurst Exponent | Autocorrelation | Multifractal Width |
|---|---|---|---|
| Ben | 0.547±0.157 | 0.905±0.313 | 0.926±0.122 |
| Asa | 0.389±0.032 | 1.222±0.064 | 0.864±0.09 |
| Mat | 0.371±0.057 | 1.258±0.114 | 0.974±0.124 |
| Odi | 0.641±0.052 | 0.719±0.104 | 1.126±0.206 |
| Nep | 0.547±0.042 | 0.906±0.084 | 1.174±0.098 |
| Man | 0.587±0.143 | 0.827±0.286 | 1.071±0.236 |
| Hin | 0.704±0.103 | 0.592±0.207 | 1.401±0.106 |
| Urd | 0.462±0.019 | 1.076±0.038 | 1.176±0.114 |
| Mar | 0.51±0.031 | 0.98±0.062 | 1.027±0.14 |
| Guj | 0.369±0.027 | 1.262±0.053 | 1.316±0.264 |
| Pun | 0.47±0.02 | 1.06±0.04 | 0.805±0.112 |
| Kon | 0.493±0.089 | 1.013±0.178 | 1.155±0.194 |
| Tam | 0.310±0.063 | 1.379±0.127 | 1.092±0.127 |
| Tel | 0.499±0.032 | 1.002±0.064 | 1.361±0.131 |
| Mll | 0.396±0.033 | 1.208±0.067 | 0.976±0.071 |
| Kan | 0.604±0.141 | 0.793±0.282 | 2.375±0.319 |
| San | 0.478±0.068 | 1.043±0.137 | 1.176±0.12 |

As explained in Section 3, MF-DXA is used to perform cross-correlation between language pairs. First, undesirable signals are removed from the data using Eqs. 11 and 12. After the time series has been integrated over time, it is split into $N_s$ bins, where $N_s = int(N/s)$, where N is the series length and s is the scale size. Then, using the same procedures as in MFDFA, the qth order detrended covariance $F_q(s)$ is calculated for values of $q$ ranging from -5 to 5. The multifractal cross-correlation in Figure 3(a) and Figure 3(b) illustrates how the power-law cross-correlations change as q increases.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

**Fig. 6 Variation of λ(q) and h(q) for Sanskrit and Bengali**

Figure 6 illustrates the characteristics of the function λ(q) for two cross-correlated signals, San and Ben. This figure also compares the function h(q), derived from MFDFA, for San and Ben, highlighting the multifractal nature of spoken languages.

Previous studies (Berument et al., 2010; C. M. Jones, 1996; Chen, 2009; Deb, 2012; Reboredo et al., 2014) have reported negative cross-correlations, which indicate strong similarities between the analyzed samples.

Among these, only four language pairs show positive cross-correlation. Notably, the Tam-Mll pair has a positive cross-correlation value, suggesting less similarity between them compared to other Dravidian languages. This observation may be attributed to the significant removal of San loanwords and other foreign influences from Tam during the Pure Tam Movement.

Nep exhibits a strong relationship with several Indo-Aryan languages, particularly Mat and San (Jain & Cardona, 2007). The cross-correlation between Nep and San is significantly high, while that between Nep and Mat is even higher. However, the cross-correlation value for Nep and Mat remains negative, indicating a level of comparability between the two languages. Table 3 provides the cross-correlation coefficients for these language pairs, organized in a confusion matrix.

The following observations are drawn from the confusion matrix and distribution curve:

- The last row of the diagonal matrix (Table 3) is empty because the cross-correlation values for San with the other languages are provided in the upper part of the matrix.

- Negative cross-correlation is observed across the majority of language pairs. A consistent negative cross-correlation was initially reported by S.S. Chen (Chen, 2009).

- According to the frequency distribution curve, the Ben-Hin, Ori-Hin, and Hin-Kan pairs exhibit strong correlations. Hin is a prominent language in Jharkhand and Chhattisgarh, located on the northern and western borders of Odisha. Therefore, the influence of Hin on the languages of Odisha cannot be overlooked.

- Kan and Hin are influenced by San in addition to sharing similar vocabularies. This influence, along with lexical similarities, explains the notable resemblance between these two languages, despite their lack of genealogical affinity.

- Additionally, both Hin and Ben belong to the Indo-Aryan language family and share linguistic roots in San, further contributing to their linguistic similarities.

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam, Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

**Table 3: Cross-Correlation coefficient among the language pairs**

| | Asa | Mat | Ori | Nep | Man | Hin | Urd | Mar | Guj | Pun | Kon | Tam | Tel | Mll | Kan | San |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| San | -1.2515 | -0.5936 | -0.4783 | -1.3227 | -1.0935 | -1.2701 | -0.5936 | -0.9253 | -1.0341 | -0.6646 | -0.8660 | -1.1060 | -0.2534 | -0.9521 | -0.5842 | **-1.5814** |
| Kan | -1.2651 | -1.1020 | -1.0618 | -2.0722 | -1.7382 | -1.7109 | -2.2124 | -1.3160 | -1.5938 | -0.8935 | -1.2992 | -1.6275 | -0.7187 | -1.4911 | **-1.1354** | |
| Mll | -0.9563 | -0.2688 | -0.2631 | -1.1878 | -0.8280 | -1.0837 | -1.3758 | -0.6446 | -0.6428 | -0.0166 | -0.6298 | -0.7442 | 0.0133 | **-0.5916** | | |
| Tel | -1.2048 | -0.5012 | -0.4868 | -1.4926 | -1.1563 | -1.5328 | -1.6799 | -0.8041 | -0.9181 | -0.5295 | -0.9599 | -0.9608 | **-0.4114** | | | |
| Tam | -0.5570 | 0.0978 | 0.1741 | -0.8980 | -0.5877 | -0.6453 | -1.0144 | -0.1985 | -0.3929 | 0.2122 | -0.2495 | **-0.2645** | | | | |
| Kon | -1.2097 | -0.5427 | -0.6995 | -1.5716 | -1.1062 | -1.4348 | -1.8182 | -1.0807 | -1.2779 | -0.5343 | **-0.9650** | | | | | |
| Pun | -1.1416 | -0.4216 | -0.4201 | -1.4200 | -0.9280 | -1.2485 | -1.6159 | -0.7955 | -0.8941 | **-0.4522** | | | | | | |
| Guj | -1.3596 | -0.0415 | -0.1398 | -1.0084 | -0.6158 | -0.8581 | -1.3322 | -0.3670 | **-0.6142** | | | | | | | |
| Mar | -1.3516 | -0.4951 | -0.7247 | -1.6485 | -1.1750 | -1.2810 | -1.7455 | **-0.8823** | | | | | | | | |
| Urd | -1.1837 | -0.3517 | -0.5099 | -1.4391 | -1.0045 | -1.2750 | **-1.5567** | | | | | | | | | |
| Hin | -2.5151 | -1.2257 | -1.2981 | -2.3281 | -1.9871 | **-2.1163** | | | | | | | | | | |
| Man | -1.7436 | -1.0774 | -1.0134 | -1.9475 | **-1.4124** | | | | | | | | | | | |
| Nep | -1.3859 | -0.7580 | -0.7271 | **-1.7442** | | | | | | | | | | | | |
| Ori | -1.8480 | -1.1432 | **-1.1444** | | | | | | | | | | | | | |
| Mat | -0.8058 | **-0.1274** | | | | | | | | | | | | | | |
| Asa | **-1.3842** | | | | | | | | | | | | | | | |

| Ben | | | | | | | | | | | | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | Ben | Asa | Mat | Ori | Nep | Man | Hin | Urd | Mar | Guj | Pun | Kon | Tam | Tel | Mll | Kan | San |

Given that its cross-correlation values with the majority of other languages fall within a comparable range, San is a crucial language in this category. It is commonly acknowledged that contemporary Indic languages, including Ben, Hin, Pun, and Guj, are descended from the San language (Dutt, 2013; Strazny, 2013). Numerous studies have also emphasized San's importance in the formation of Indo-Aryan languages (Krishnamurti, 2006) and its influence on the Dravidian language family (Kolipakam et al., 2018; Witzel, 1999). Additionally, it has been noted that San and Mll are closely related languages (Academi).

To get more accurate insight one ANOVA was carried out on all 17 languages. The three parameters obtained from the multifractal namely Hurst exponent (Table 4), Autocorrelation (Table 5) and Multifractal width (Table 6) were used to test the statistical significance of the genesis of languages.

**Table 4: One-way ANOVA parameters on Hurst exponent**

| | DF | SS | MS | F Value | Prob>F |
|-------|-----|---------|---------|----------|--------|
| **Model** | 16 | 1.79159 | 0.11197 | 18.02413 | 0 |
| **Error** | 153 | 0.95051 | 0.00621 | | |
| **Total** | 169 | 2.7421 | | | |

**Table 5: One-way ANOVA parameters on Autocorrelation**

| | DF | SS | MS | F Value | Prob>F |
|-------|-----|---------|--------|----------|--------|
| **Model** | 16 | 7.16636 | 0.4479 | 18.02412 | 0 |
| **Error** | 153 | 3.80204 | 0.02485 | | |
| **Total** | 169 | 10.9684 | | | |

**Table 6: One-way ANOVA parameters on Multifractal width**

| | DF | SS | MS | F Value | Prob>F |
|-------|-----|----------|---------|---------|----------|
| **Model** | 19 | 33.82974 | 1.78051 | 3.37142 | 1.03E-05 |
| **Error** | 180 | 95.06155 | 0.52812 | | |
| **Total** | 199 | 128.8913 | | | |

The result was obtained with a 95% significance level. Following (George Cardona, 2007; Sardesai, 2019; Strazny, 2013), a post-hoc analysis using Tukey's HSD comparison test was conducted to determine the significant difference between the two languages. The significant values indicate substantial differences between the Asa, Ben, Guj, Hin, Mat, Mar, and Nep languages. The significance value of Autocorrelation and the Hurst exponent can be used to assess how different the languages are. Although Asa-Guj, Mat-Asa, Mat-Guj, Mar-Guj, and Mar-Asa have some overlap. These linguistic groups cannot be distinguished from one another using all of the three criteria. The languages that overlap the most are Odi, Pun, Kon, Urd, San, Kan, Mll, and Tel. It is evident that Tam differs greatly from the other languages, and that multifractal width is enough to set Tam apart from the others.

The analysis reveals that Ben is most closely aligned with Kan, while Asa exhibits the smallest MD relative to Mat. In case of Kan the smallest MD, is 1.122102, for Hin, highlighting their strong linguistic affinity. Additionally, the geographical proximity of Manipur and Nepal is reflected in their linguistic similarity, as Nep and Man emerge as the closest pair based

on the MD.

In Karnataka, Maharashtra, and Kerala, Kon, the official language of Goa, is also acknowledged as a minority language (George Cardona, 2007). There will inevitably be linguistic similarities between these areas due to their geographic spread. There have also been reports of contentious linguistic ties between Kon and Mar (Sardesai, 2019; Strazny, 2013), underscoring their intimate bond. The apparent parallels between Kon, Mar, and Kan are probably explained by these characteristics. Table 7 displays the confusion matrix based on the MD.

**Table 7: Confusion Matrix with Mahalanobis Distance as parameter**

|  | Asa | Ben | Guj | Hin | Kan | Kon | Mll | Man | Mar | Mat | Nep | Odi | Pun | San | Tel | Tam | Urd |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Asa | 0.00 | 1.64 | 2.42 | 9.41 | 2.17 | 2.99 | 1.77 | 2.40 | 5.20 | 1.17 | 6.28 | 5.91 | 3.82 | 3.18 | 6.27 | 2.22 | 4.97 |
| Ben | 1.64 | 0.00 | 2.66 | 4.41 | 0.68 | 1.49 | 1.57 | 0.82 | 0.88 | 1.57 | 2.26 | 1.41 | 1.23 | 2.35 | 3.61 | 2.37 | 2.41 |
| Guj | 2.42 | 2.66 | 0.00 | 4.66 | 2.40 | 2.03 | 1.96 | 2.28 | 5.32 | 1.67 | 5.11 | 7.27 | 5.31 | 2.42 | 4.64 | 1.87 | 4.12 |
| Hin | 9.41 | 4.41 | 4.66 | 0.00 | 1.12 | 3.47 | 9.41 | 2.28 | 4.94 | 7.66 | 3.98 | 2.14 | 8.04 | 3.88 | 3.16 | 7.63 | 5.03 |
| Kan | 2.17 | 0.68 | 2.40 | 1.12 | 0.00 | 1.17 | 2.06 | 0.60 | 0.95 | 2.20 | 0.65 | 0.87 | 1.33 | 1.53 | 1.11 | 2.83 | 1.42 |
| Kon | 2.99 | 1.49 | 2.03 | 3.47 | 1.17 | 0.00 | 2.37 | 0.95 | 0.85 | 2.39 | 1.19 | 2.43 | 2.44 | 0.28 | 1.53 | 2.80 | 0.61 |
| Mll | 1.77 | 1.57 | 1.96 | 9.41 | 2.06 | 2.37 | 0.00 | 2.03 | 4.31 | 0.61 | 5.77 | 5.71 | 3.05 | 2.43 | 5.55 | 1.77 | 4.35 |
| Man | 2.40 | 0.82 | 2.28 | 2.28 | 0.60 | 0.95 | 2.03 | 0.00 | 0.81 | 2.17 | 0.65 | 0.71 | 1.97 | 1.34 | 1.68 | 2.65 | 1.30 |
| Mar | 5.20 | 0.88 | 5.32 | 4.94 | 0.95 | 0.85 | 4.31 | 0.81 | 0.00 | 3.28 | 1.72 | 3.21 | 2.68 | 1.51 | 2.52 | 4.24 | 2.04 |
| Mat | 1.17 | 1.57 | 1.67 | 7.66 | 2.20 | 2.39 | 0.61 | 2.17 | 3.28 | 0.00 | 4.61 | 5.55 | 2.55 | 2.39 | 4.58 | 1.35 | 3.38 |
| Nep | 6.28 | 2.26 | 5.11 | 3.98 | 0.65 | 1.19 | 5.77 | 0.65 | 1.72 | 4.61 | 0.00 | 2.19 | 4.63 | 1.80 | 2.09 | 5.62 | 2.84 |
| Odi | 5.91 | 1.41 | 7.27 | 2.14 | 0.87 | 2.43 | 5.71 | 0.71 | 3.21 | 5.55 | 2.19 | 0.00 | 4.79 | 3.63 | 4.34 | 7.00 | 5.42 |
| Pun | 3.82 | 1.23 | 5.31 | 8.04 | 1.33 | 2.44 | 3.05 | 1.97 | 2.68 | 2.55 | 4.63 | 4.79 | 0.00 | 3.25 | 4.68 | 3.65 | 3.45 |
| San | 3.18 | 2.35 | 2.42 | 3.88 | 1.53 | 0.28 | 2.43 | 1.34 | 1.51 | 2.39 | 1.80 | 3.63 | 3.25 | 0.00 | 1.86 | 2.83 | 0.66 |
| Tel | 6.27 | 3.61 | 4.64 | 3.16 | 1.11 | 1.53 | 5.55 | 1.68 | 2.52 | 4.58 | 2.09 | 4.34 | 4.68 | 1.86 | 0.00 | 4.99 | 2.02 |
| Tam | 2.22 | 2.37 | 1.87 | 7.63 | 2.83 | 2.80 | 1.77 | 2.65 | 4.24 | 1.35 | 5.62 | 7.00 | 3.65 | 2.83 | 4.99 | 0.00 | 4.98 |
| Urd | 4.97 | 2.41 | 4.12 | 5.03 | 1.42 | 0.61 | 4.35 | 1.30 | 2.04 | 3.38 | 2.84 | 5.42 | 3.45 | 0.66 | 2.02 | 4.98 | 0.00 |

Odi is an Indo-Aryan language though it is strongly correlated with Kan which is a Dravidian language. A similar type of

observation is observed from MD. It is reported that ancient Orissa was an inhabitant of Dravidian-speaking people (Paulsen et al., 2019). Kan and Tel both languages are from the Dravidian language family, hence 60% of Tel is misclassified as Kan (Khadabadi & Akādamī, 1997) and it is also supported by MD.

The classification of the Man language remains contentious. Traditionally, it's categorized under the Tibeto-Burman language family (Singh, 1996); however, some scholars contend that it might belong to the Indo-Aryan family. This argument arises from observed linguistic similarities between Man and well-known Indo-Aryan languages like Ben, Mar, Hin, Mat, and Pun. The observation can be supported by analyzing the confusion matrices obtained from methods like MD.

## 6. CONCLUSION

This research delves into the multifractal properties of spoken languages in India using Multifractal Detrended Fluctuation Analysis (MFDFA). The study focuses on the complexity and self-similarity in speech signals, showcasing how MFDFA surpasses traditional nonlinear techniques in capturing intricate variations in non-stationary time series data like speech. By quantifying complexity through multifractal width (W), the analysis reveals that languages with higher W values exhibit greater temporal variations, highlighting their unique characteristics.

Comparative analysis of Ben and San displays significant differences in their multifractal spectra, emphasizing spoken languages' self-similar yet diverse nature. The cross-correlation of language pairs indicates notable linguistic similarities, with mostly negative cross-correlation values suggesting shared structural features despite genealogical differences.

The study identifies San as a foundational language influencing other Indo-Aryan and Dravidian languages, evidenced by its strong correlation with Ben and other Indic languages. ANOVA results highlight significant differences across languages for parameters such as the Hurst exponent, autocorrelation, and multifractal width.

Post-hoc Tukey's HSD tests reveal considerable overlap among closely related languages, particularly within groups like Asa, Guj, and Mar. Therefore, this study provides the relationships between geographically and linguistically proximate languages.

## REFERENCES

[1] Academi, K. S. KERALA SAHITYA AKADEMI. http://www.keralasahityaakademi.org/

[2] Ashkenazy, Y., Baker, D. R., Gildor, H., & Havlin, S. (2003). Nonlinearity and multifractality of climate change in the past 420,000 years. *Geophysical research letters*, *30*(22).

[3] Berument, M. H., Ceylan, N. B., & Dogan, N. (2010). The impact of oil price shocks on the economic growth of selected MENA1 countries. *The Energy Journal*, *31*(1), 149-176.

[4] Bradlow, A., Clopper, C., Smiljanic, R., & Walter, M. A. (2010). A perceptual phonetic similarity space for languages: Evidence from five native language listener groups. *Speech communication*, *52*(11-12), 930-942.

[5] C. M. Jones, G. K. (1996). Oil and the Stock Markets. *The Journal of Finance*, 915-932. https://doi.org/ https://doi.org/10.1111/j.1540-6261.1996.tb02691.x

[6] Chen, S.-S. (2009). Oil price pass-through into inflation. *Energy economics*, *31*(1), 126-133.

[7] DC González, L. L. L., F Violaro. (2012). Analysis of the multifractal nature of speech signals. Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications: 17th Iberoamerican Congress, CIARP 2012, Buenos Aires, Argentina.

[8] Deb, D. (2012). On case marking in assamese bengali and oriya. *International Journal of Applied Linguistics & English Literature*, *1*(2), 102.

[9] Downey, S. S., Hallmark, B., Cox, M. P., Norquest, P., & Lansing, J. S. (2008). Computational feature-sensitive reconstruction of language relationships: Developing the ALINE distance for comparative historical linguistic reconstruction. *Journal of Quantitative Linguistics*, *15*(4), 340-369.

[10] Dutt, S. (2013). India in a globalized world. In *India in a globalized world*. Manchester University Press.

[11] George Cardona, D. J. (2007). *The Indo-Aryan Languages*. Routledge, London.

[12] Ghosh, D., Dutta, S., & Chakraborty, S. (2014). Multifractal detrended cross-correlation analysis for epileptic patient in seizure and seizure free status. *Chaos, Solitons & Fractals*, *67*, 1-10.

[13] He, L.-Y., & Chen, S.-P. (2011). Multifractal detrended cross-correlation analysis of agricultural futures markets. *Chaos, Solitons & Fractals*, *44*(6), 355-361.

[14] Hedayatifar, L., Vahabi, M., & Jafari, G. (2011). Coupling detrended fluctuation analysis for analyzing coupled nonstationary signals. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, *84*(2), 021138.

[15] Horvatic, D., Stanley, H. E., & Podobnik, B. (2011). Detrended cross-correlation analysis for non-stationary

time series with periodic trends. *Europhysics Letters*, *94*(1), 18007.

[16] Jain, D., & Cardona, G. (2007). *The Indo-Aryan Languages*. Routledge.

[17] Jiang, Z.-Q., & Zhou, W.-X. (2011). Multifractal detrending moving-average cross-correlation analysis. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, *84*(1), 016106.

[18] Kantelhardt, J. W., Zschiegner, S. A., Koscielny-Bunde, E., Havlin, S., Bunde, A., & Stanley, H. E. (2002). Multifractal detrended fluctuation analysis of nonstationary time series. *Physica A: Statistical Mechanics and its Applications*, *316*(1-4), 87-114.

[19] Kettunen, K., Sadeniemi, M., Lindh-Knuutila, T., & Honkela, T. (2006). Analysis of EU languages through text compression. International Conference on Natural Language Processing (in Finland),

[20] Khadabadi, B., & Akādamī, P. t. B. (1997). Studies in Jainology, Prakrit literature, and languages: a collection of select 51 papers. *(No Title)*.

[21] Kolipakam, V., Jordan, F. M., Dunn, M., Greenhill, S. J., Bouckaert, R., Gray, R. D., & Verkerk, A. (2018). A Bayesian phylogenetic study of the Dravidian language family. *Royal Society open science*, *5*(3), 171504.

[22] Kondrak, G. (2000). A new algorithm for the alignment of phonetic sequences. 1st Meeting of the North American Chapter of the Association for Computational Linguistics,

[23] Krishnamurti, B. (2006). *The Dravidian languages* (1st ed.). Cambridge University Press.

[24] McMahon, A., & McMahon, R. (2005). *Language classification by numbers*. Oxford University Press.

[25] Movahed, M. S., & Hermanis, E. (2008). Fractal analysis of river flow fluctuations. *Physica A: Statistical Mechanics and its Applications*, *387*(4), 915-932.

[26] Oco, N., Syliongka, L. R., Roxas, R. E., & Ilao, J. (2013). Dice's coefficient on trigram profiles as metric for language similarity. 2013 International Conference Oriental COCOSDA held jointly with 2013 Conference on Asian Spoken Language Research and Evaluation (O-COCOSDA/CASLRE),

[27] Office of the Registrar General & Census Commissioner, I. *Office of the Registrar General & Census Commissioner, India*. Office of the Registrar General & Census Commissioner, India. Retrieved 10/04/2019 from https://censusindia.gov.in/census.website/

[28] Paulsen, J., Bergh, K., Chew, A., Gugerty, M. K., & Anderson, C. L. (2019). Wheat value chain: Bihar. *Gates Open Res*, *3*(593), 593.

[29] Podobnik, B., Grosse, I., Horvatić, D., Ilic, S., Ivanov, P. C., & Stanley, H. E. (2009). Quantifying cross-correlations using local and global detrending approaches. *The European Physical Journal B*, *71*, 243-250.

[30] Podobnik, B., Horvatic, D., Ng, A. L., Stanley, H. E., & Ivanov, P. C. (2008). Modeling long-range cross-correlations in two-component ARFIMA and FIARCH processes. *Physica A: Statistical Mechanics and its Applications*, *387*(15), 3954-3959.

[31] Podobnik, B., Horvatic, D., Petersen, A. M., & Stanley, H. E. (2009). Cross-correlations between volume change and price change. *Proceedings of the National Academy of Sciences*, *106*(52), 22079-22084.

[32] Podobnik, B., Jiang, Z.-Q., Zhou, W.-X., & Stanley, H. E. (2011). Statistical tests for power-law cross-correlated processes. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, *84*(6), 066118.

[33] Podobnik, B., & Stanley, H. E. (2008). Detrended Cross-Correlation Analysis: A New Method<? format?> for Analyzing Two Nonstationary Time Series. *Physical review letters*, *100*(8), 084102.

[34] Reboredo, J. C., Rivera-Castro, M. A., & Zebende, G. F. (2014). Oil and US dollar exchange rate dependence: A detrended cross-correlation approach. *Energy economics*, *42*, 132-139.

[35] Sanyal, S., Banerjee, A., Patranabis, A., Banerjee, K., Sengupta, R., & Ghosh, D. (2016). A study on improvisation in a musical performance using multifractal detrended cross correlation analysis. *Physica A: Statistical Mechanics and its Applications*, *462*, 67-83.

[36] Sardesai, M. (2019). *Mother tongue blues*. Retrieved 5/4/2019 from M. Sardesai

[37] Sengupta, D., & Saha, G. (2015). Study on similarity among Indian languages using language verification framework. *Advances in Artificial Intelligence*, *2015*(1), 325703.

[38] Shimizu, Y., Thurner, S., & Ehrenberger, K. (2002). Multifractal spectra as a measure of complexity in human posture. *Fractals*, *10*(01), 103-116.

[39] Singh, C. M. (1996). *A history of Manipuri literature*. Sahitya Akademi.

[40] Strazny, P. (2013). *Encyclopedia of linguistics*. Routledge.

[41] Titze, I. R. (1995). *Workshop on acoustic voice analysis: Summary statement*. National Center for Voice and

Suparna Panchanan, Nilav Darsan Mukhopadhyay, Ranjan Banerjee, Most Mahabuba Islam,
Shankha Sanyal, Dipak Ghosh, Debmalya Mukherjee

Speech.

[42] Wang, F., Liao, G.-p., Zhou, X.-y., & Shi, W. (2013). Multifractal detrended cross-correlation analysis for power markets. *Nonlinear Dynamics*, *72*, 353-363.

[43] Witzel, M. (1999). Substrate languages in old indo-Aryan.(Ṛgvedic, middle and late vedic). *Electronic Journal of Vedic Studies*, *5*(1), 1-67.

[44] Xu, N., Shang, P., & Kamae, S. (2010). Modeling traffic flow correlation using DFA and DCCA. *Nonlinear Dynamics*, *61*, 207-216.

[45] Zhou, W.-X. (2008). Multifractal detrended cross-correlation analysis for two nonstationary signals. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, *77*(6), 066211.