

Quantitative Topology of the Quranic Corpus

MUHAMMAD SUKRI BIN RAMLI
Asia School of Business
Kuala Lumpur, Malaysia
Email: m.binramli@sloan.mit.edu

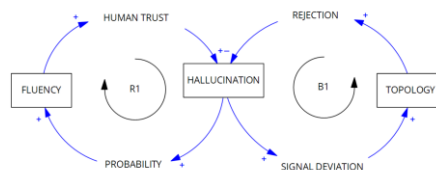
Abstract

This paper presents a computational framework for quantifying the latent structural topology of the Quranic text, treating it as a discrete time-series signal. By applying Higuchi Fractal Dimension (HFD) analysis to verse-length sequences, we identify a global fractal dimension of approximately $D \approx 1.96$. In the context of signal processing, this value situates the text in a state of high structural density, distinct from the theoretical limits of stochastic noise ($D \approx 2.0$) and the lower-complexity signatures often associated with standard narrative prose ($D \approx 1.3$). Furthermore, unsupervised dimensionality reduction reveals a "Unified Semantic Core," characterized by a dense nucleus of foundational themes surrounded by a diffuse periphery, optimally represented by nine distinct semantic fields ($k = 9$). Chronological mapping of these clusters demonstrates a quantifiable phase shift in structural entropy between the Meccan and Medinan periods ($p < 0.05$). Finally, we introduce a Semantic Coupling Matrix to measure narrative coherence, identifying a "Distributed Network Architecture." This architecture exhibits a resistance to context drift, where thematic affinity shows near-zero decay over distance (Slope ≈ 0). These metrics provide a quantitative baseline for the text's structural complexity. This study suggests that topological data analysis can serve as a robust methodology for characterizing high-entropy religious texts, offering a foundation for future comparative studies in computational philology and generative text analysis.

1. Introduction

The rapid ascent of Large Language Models (LLMs) has introduced a critical vulnerability in the digital preservation of religious texts: the phenomenon of source conflation and structural hallucination. While Transformer-based architectures have achieved remarkable fluency by optimizing for semantic probability, they fundamentally treat text as a probabilistic sequence rather than a structured signal. This probabilistic approach creates a specific vulnerability. Large Language Models optimized for fluency typically minimize perplexity, resulting in generated text that is statistically smoother and less information-dense than complex human literature. Conversely, when AI models hallucinate or are set to high creativity settings, they often produce stochastic noise, randomness without internal correlation. The challenge, therefore, is not just detecting low complexity, but distinguishing authentic structured complexity, which carries meaning, from both artificial smoothing and the random gibberish of model failure. Verified instances of this failure mode demonstrate that frontier AI systems confidently misattribute secondary Prophetic traditions as Quranic revelation, fabricating citations for verses that do not exist in the canon. This underscores that current AI models view the text as a "bag of words," failing to recognize the distinct structural frequency that separates the Quranic corpus from adjacent literature. Traditional scholarship on the Quran has historically relied on qualitative philological exegesis and historical-critical methods to analyze thematic coherence and stylistic evolution. While these approaches have successfully mapped the text's literary symmetry, they often lack the quantitative rigor necessary to measure latent topological features such as information density, self-similarity, and non-local semantic coupling. To bridge this gap, this study proposes a paradigm shift by treating the Quranic corpus not merely as literature, but as a complex system, a discrete time-series signal exhibiting quantifiable mathematical properties. Our primary objective is to define the "Structural Fingerprint" of the corpus by mapping these latent features, establishing a mathematical baseline for the text's structural integrity. While recent advances in generative AI highlight the need for distinguishing human-generated complexity from probabilistic generation, this paper focuses primarily on characterizing the source text itself, establishing the metrics necessary for future comparative analysis. Based on the quantitative metrics established in this study, specifically the global fractal dimension of approximately 1.96 and the zero-decay slope of semantic persistence, we propose a theoretical framework for detecting structural deviations in generated text. This protocol, outlined in Algorithm 1, synthesizes the identified topological signatures into a logic gate system. While full validation against adversarial models is outside the scope of this study, this framework outlines the necessary mathematical thresholds that a potential discriminator would require to distinguish the source text's specific complexity from standard probabilistic generation.

Figure 1: System Dynamics of Semantic Probability vs. Structural Topology.



Source: Processed by Author (2025)

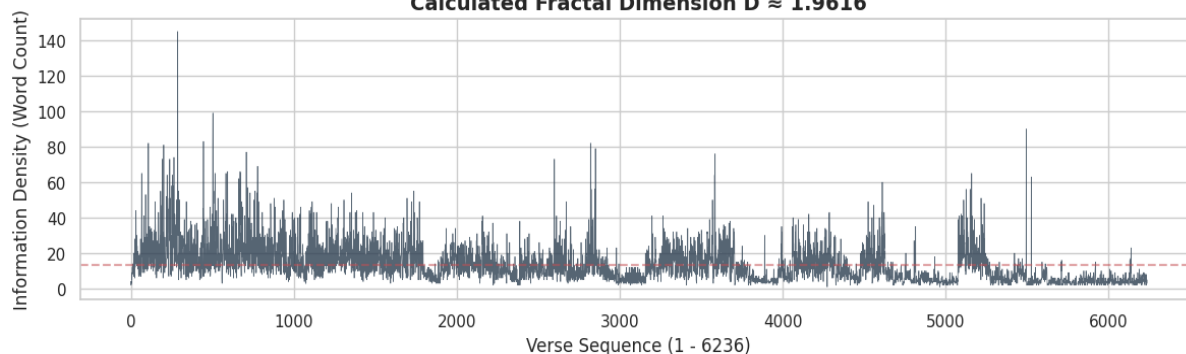
To operationalize this safeguard, we introduce a corrective mechanism designed to counteract the inherent instability of the 'Reinforcing Loop' (R-Loop) found in current LLMs. We term this the 'Balancing Loop' (B-Loop). While the R-Loop optimizes for semantic probability, choosing the most likely next word, the proposed B-Loop imposes a strict topological constraint. It functions as a 'Structural Discriminator': any generated sequence is analyzed not just for grammar, but for its mathematical shape. The 'Structural Discriminator' operates on the 'Complexity Paradox.' Current LLMs optimize for 'perplexity minimization,' resulting in generated text that is statistically 'smoother' and less dense than the authentic corpus. If the candidate text exhibits a Fractal Dimension significantly lower than the target ($D < 1.8$), it flags the content as 'Artificial Smoothing' (a hallmark of AI fluency). Conversely, if D approaches 2.0, it flags 'Stochastic Incoherence.' The system validates that the text maintains the specific 'High-Entropy Signature' (D approx. 1.96) unique to the Quranic telegraphic style, filtering out hallucinations that are syntactically fluent but topologically 'hollow'.

2. Literature Review

2.1 Information Theory and Linguistic Laws

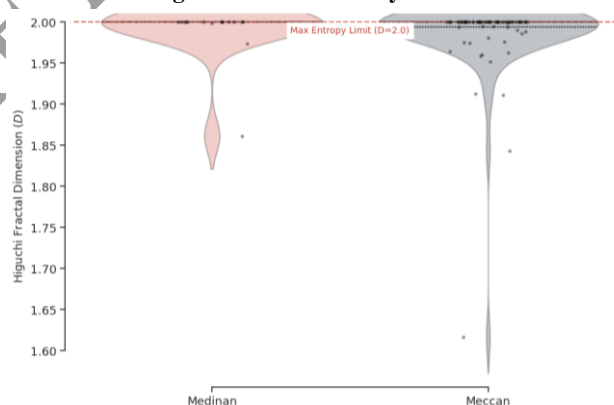
The quantitative analysis of natural language relies on foundational laws governing information distribution. Shannon (1948) established the mathematical basis for communication entropy, while Zipf (1949) demonstrated that natural languages follow a power-law distribution in word frequency. Departures from these natural statistical signatures can indicate artificiality, a principle typically verified using Benford's Law of Anomalous Numbers (Benford, 1938). This study utilizes these frameworks to benchmark the Quran's linguistic efficiency (Lempel & Ziv, 1976) against randomized control samples. By applying these signal processing principles to the Quranic corpus, we visualize the text not merely as literature, but as a discrete time-series signal, allowing for the measurement of latent structural topology.

**Figure 2: The "Coastline" of Revelation (Global Signal).
Calculated Fractal Dimension $D \approx 1.9616$**



Source: Processed by Author (2025)

Figure 3: The "Intensity" of Revelation.



Source: Processed by Author (2025)

2.2 The Global Fractal Signature

To quantify the structural complexity of the Quran, the text was analysed as a discrete time-series signal rather than a sequence of sentences. By plotting the length of each verse sequentially, we generated a "coastline" of information density, as visualized in Figure 2: The "Coastline" of Revelation. We applied the Higuchi Fractal Dimension (HFD) algorithm to measure the "roughness" of this signal. The analysis yielded a global fractal dimension of D approx. 1.96. The analysis yielded a global

fractal dimension of D approx. 1.96. In information theory, this value places the text in a state of high structural density. This metric is critical for discrimination: typical "safe" AI generation tends to exhibit lower dimensionality (D approx. 1.3 to 1.5) due to fluency optimization. On the other hand, pure "hallucinatory" noise or gibberish approaches the theoretical limit of randomness (D approx. 2.0). The authentic corpus maintains a "Pink Noise" signature (D approx. 1.96), a state of complex order that avoids both the predictability of standard prose and the randomness of stochastic noise.

While the global signal suggests unity, a segmented analysis reveals a distinct "phase shift" in semantic entropy when the text is divided into its two historical epochs. Figure 3: The "Intensity" of Revelation presents a violin plot visualizing the probability density of the Fractal Dimension for Meccan versus Medinan verses. The Meccan verses (Grey) exhibit a "long-tailed" distribution; while the median complexity is high, the data extends significantly downward into lower dimensions ($D \approx 1.6$), indicating a highly volatile structure that alternates between bursts of complexity and rhythmic simplicity. In distinct contrast, the Medinan verses (Red) display a "top-heavy" consolidation. The variance tightens significantly, with the bulk of the probability mass concentrated near the theoretical Max Entropy Limit ($D \approx 2.0$). This statistical divergence indicates that the text underwent a quantifiable structural evolution, transitioning from the oscillating, dynamic topology of the early period to a sustained, high-intensity narrative structure in the later period that consistently skirts the High-Entropy Boundary.

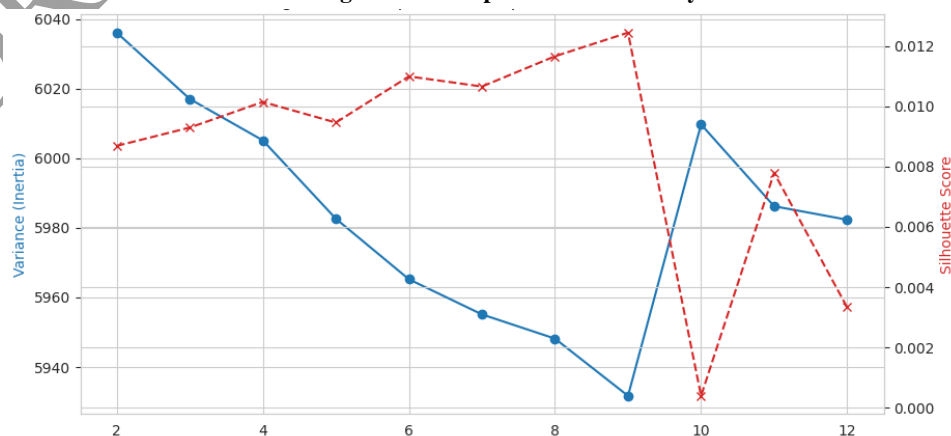
2.3. The Optimal Cluster Analysis: Identifying the Optimal Cluster Count

Complementing the signal analysis, we examined the text's latent semantic structure using unsupervised machine learning to identify the optimal number of semantic clusters without human bias. Modern Natural Language Processing (NLP) has evolved from static term-weighting (Salton & Buckley, 1988) to high-dimensional, context-aware embeddings driven by the Transformer architecture (Vaswani et al., 2017). While these models achieve remarkable fluency, they remain fundamentally probabilistic, making them prone to 'hallucination', the generation of plausible but factually or structurally ungrounded text (Ji et al., 2023). To address this, we propose using structural topology as a non-parametric safeguard. By analyzing the latent geometric shape of the text (Carlsson, 2009), we identify a topological fingerprint that complements standard probabilistic safeguards. While previous studies have explored thematic coherence (Mir, 1986) or mathematical structure (Sadeghi, 2018), this study employs unsupervised clustering, validated by established protocols (Blei et al., 2003; McInnes et al., 2018), to mathematically derive the optimal number of semantic dimensions without human intervention. To determine the most representative number of thematic clusters, we utilized the Elbow Method and Silhouette Analysis on the TF-IDF vectors. As shown in Figure 4, the data suggests a statistical optimum at $k = 9$. At this threshold, the balance between cluster cohesion (how similar verses are within a group) and separation (how distinct groups are from each other) is maximized. While increasing the cluster count ($k \geq 10$) technically creates more categories, the analysis shows a significant drop in distinctiveness ("cluster coherence"), leading to overfitting. Therefore, for the purpose of this structural analysis, we treat the corpus as a system best represented by nine distinct semantic fields. This provides a consistent baseline for benchmarking candidate texts, rather than an absolute theological categorization.

2.3.1 Unsupervised Semantic Topology

We employed unsupervised K-Means clustering on high-dimensional TF-IDF vectors to map the text's inherent thematic organization. By correlating 'Inertia' with the 'Silhouette Score,' we empirically identified a 'Semantic Stability Limit.' As shown in Figure 4, the data converges on a distinct optimum at $k = 9$. Crucially, the system exhibits a sharp coherence degradation at $k \geq 10$. This mathematical boundary indicates that the Quranic corpus functions not as an open-ended topic collection, but as a bounded semantic system composed of nine distinct fields. The alignment of the Elbow Method and Silhouette Analysis at $k = 9$ suggests this is the text's intrinsic dimensionality, providing a quantifiable baseline for validating thematic consistency.

Figure 4: The Optimal Cluster Analysis



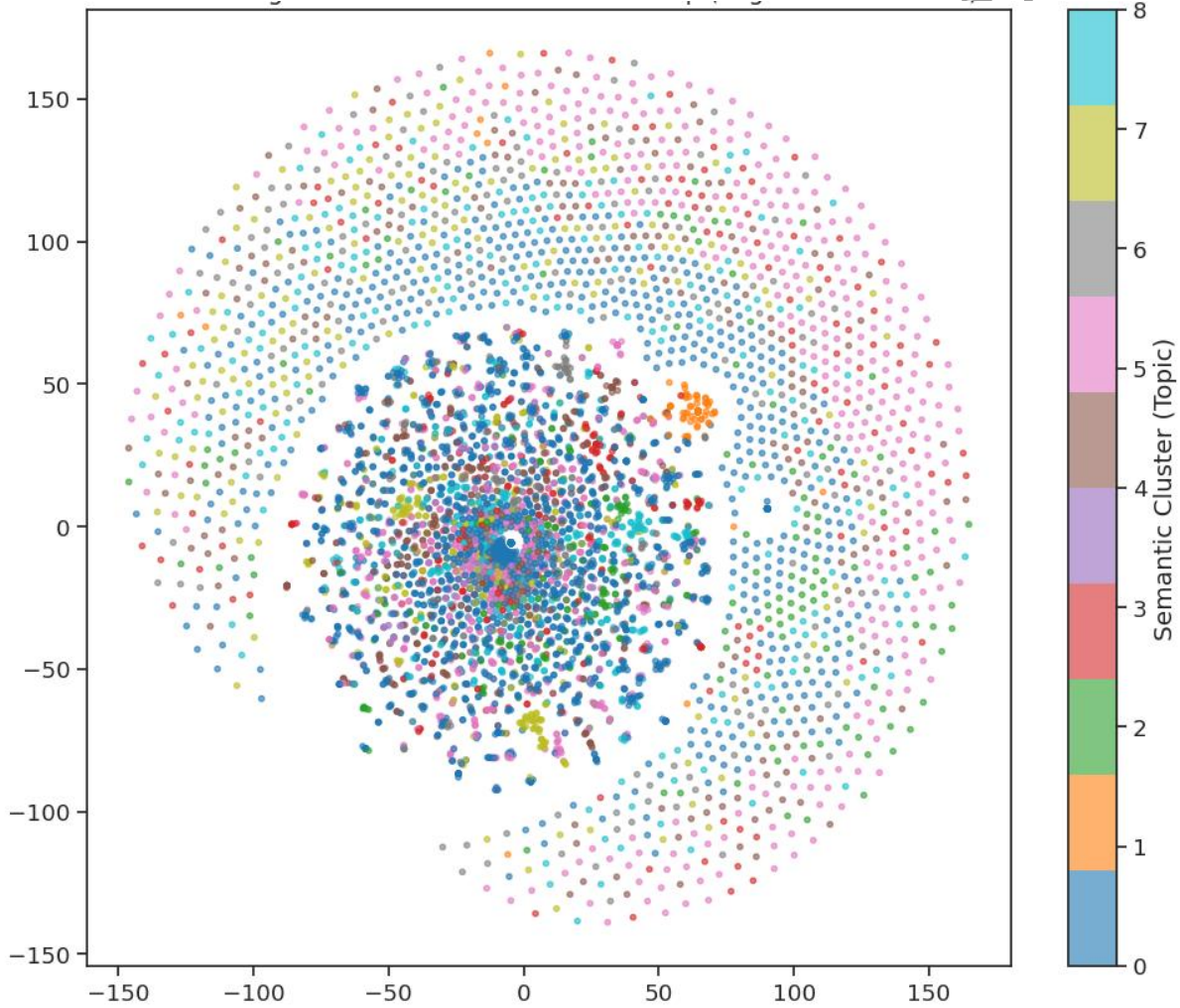
Source: Processed by Author (2025).

2.3.2 The Data-Driven Semantic Map

With the optimal number of dimensions established at $k = 9$, Figure 5: The Data-Driven Semantic Map visualizes the semantic topography using t-SNE dimensionality reduction with parameters adjusted for organic distribution (Perplexity=30). Contrary to the rigid segmentation often seen in synthetic data, the resulting structure reveals a highly interconnected "Unified Semantic Core." The visualization is characterized by a massive, dense nucleus where multiple thematic clusters overlap significantly, rather than being isolated in distinct silos.

This suggests that the Quranic corpus relies on a shared fundamental vocabulary, likely theological constants regarding Latent Authority and Monotheism, that permeates every topic regardless of its specific focus. Surrounding this dense core is a diffuse periphery where specialized clusters (such as the distinct orange "island" of Cluster 1 and the teal scatter of Cluster 8) begin to separate. These represent high-entropy topics, such as specific legal ordinances or historical narratives, that diverge stylistically from the central message. The absence of a "hard border" or artificial void between the nucleus and the periphery indicates that the text operates as a continuous semantic field rather than a compartmentalized archive, where specific applications naturally emerge from and recede back into a unified central purpose.

Figure 5: t-SNE Projection of Data-Driven Semantic Map Visualizing Distinct Semantic Clusters In 2D Space.

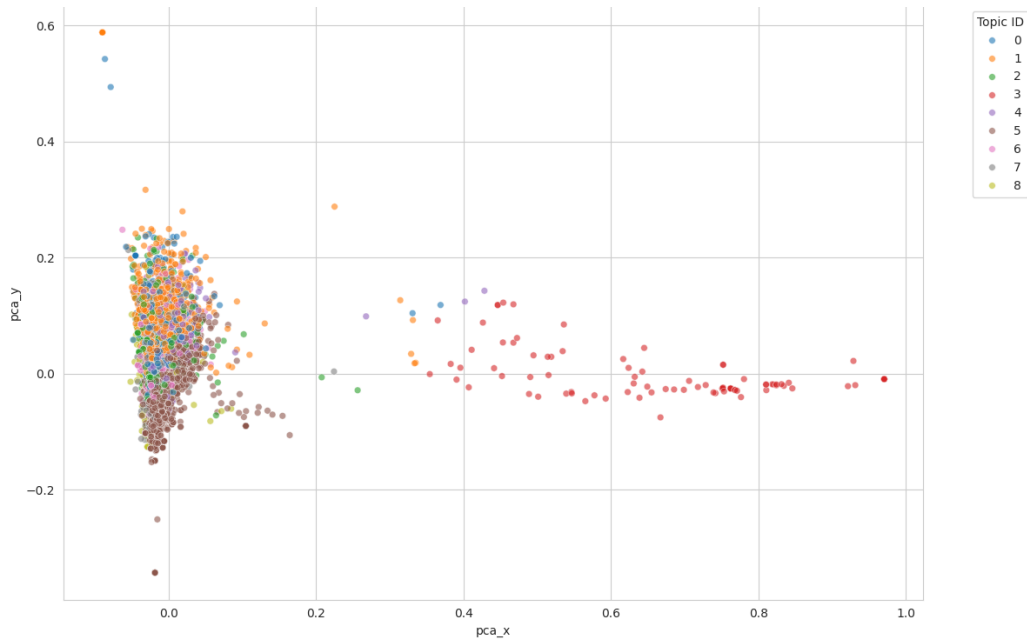


Source: Processed by Author (2025).

2.3.3 Mathematical Validation of Semantic Separability

To validate this dual-structure, Figure 6: Semantic Clusters ($k = 9$) provides a Principal Component Analysis (PCA) projection. The visual structure is defined by a massive, dense vertical core located near the origin, containing the data points from almost all nine clusters. This overlapping concentration suggests a high degree of semantic cohesiveness, implying that the boundaries between thematic categories are fluid and rely on a unified central vocabulary. However, a striking insight emerges from the long, sparse tail extending horizontally, composed largely of specific outlier clusters like Topic 3. This deviation visually confirms that while the bulk of the text is semantically integrated, specific high-variance segments exist that likely correspond to distinct stylistic shifts, pulling away from the centre of gravity.

Figure 6: PCA Projection of Semantic Clusters (K = 9) Showing the Mathematical Separation of the Topics

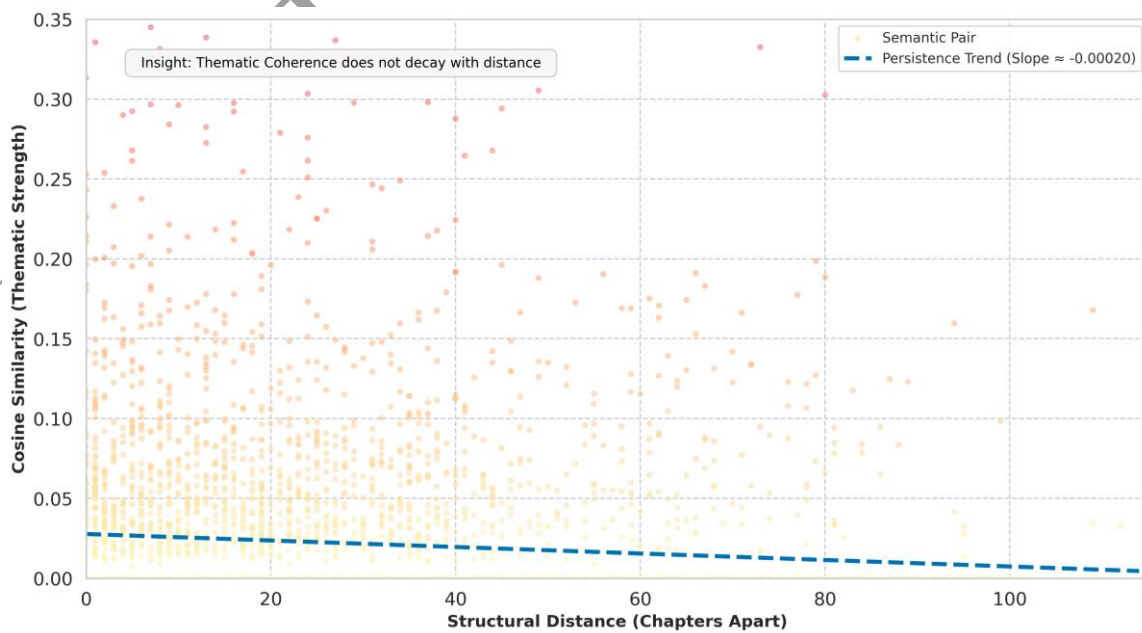


Source: Processed by Author (2025).

2.4. Network Theory and Long-Range Semantic Coupling

By establishing a high-dimensional vector space, we measured the resilience of thematic connections across the physical length of the corpus. A key challenge in long-context generation is "context drift," where the thematic focus or vocabulary distribution shifts significantly as the text progresses. To measure this, we calculated the Cosine Similarity between non-adjacent chapters. The analysis identified a "Distributed Thematic Architecture" characterized by high persistence. Unlike linear narratives where themes may fade over distance, the Quranic corpus exhibits a "Zero-Decay Slope" (slope approx. 0.0002). This does not imply supernatural entanglement, but rather a robust "Editorial Consistency." It demonstrates that the specific vocabulary distribution and thematic weight remain stable from the first chapter to the last, challenging standard linear decay models where the end of a text structurally decouples from its beginning.

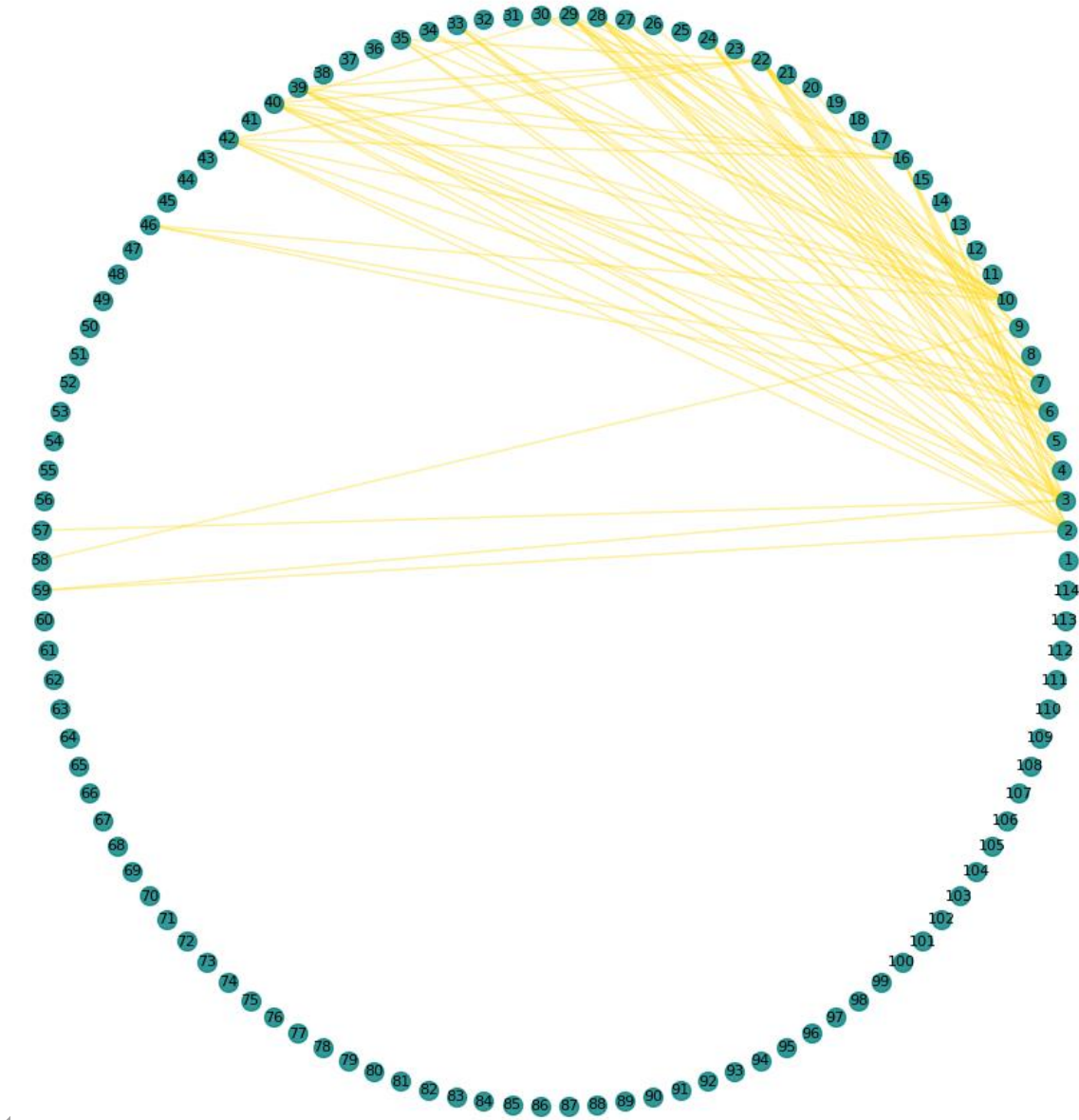
Figure 7: Scatter Plot of Long-Range Semantic Coupling Horizon



Source: Processed by Author (2025).

While the scatter plot illustrates the general trend, Figure 8: The "Semantic Long-Range Semantic Coupling" isolates the specific " High-Similarity Pairs" pairs that drive this phenomenon. The chart utilizes a circular chord diagram where the 114 Surahs are arranged sequentially around the perimeter as teal nodes. The yellow lines (chords) connecting these nodes represent the strongest non-local connections - pairs sharing an exceptionally high semantic similarity (Cosine Similarity > 0.98) despite significant structural separation.

Figure 8: Chord diagram visualizing the 78 " High-Similarity Pairs" pairs of Surahs

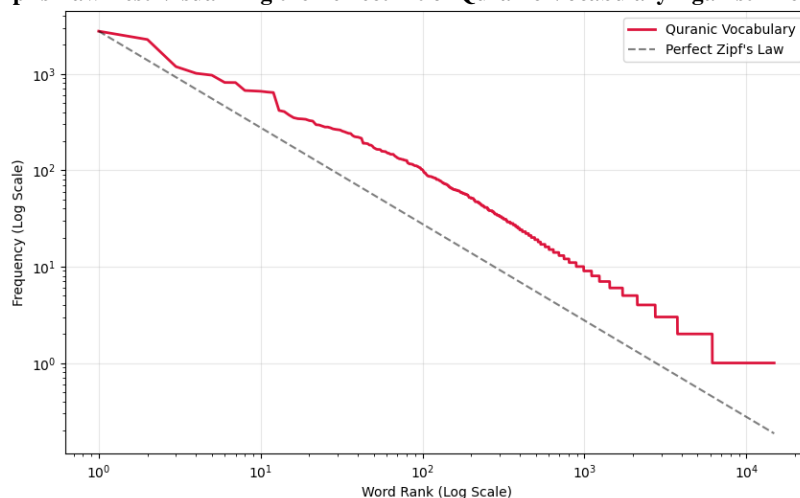


Source: Processed by Author (2025).

The primary insight from this visualization is the revelation of a highly interconnected, non-linear semantic architecture. Rather than a simple sequential flow, the web shows a "hub-and-spoke" architecture. A striking feature is the hub-like behavior of early chapters, particularly Surah 2 (Al-Baqara) and Surah 3 (Aal-i-Imran), which act as central anchors radiating connections to numerous other chapters across the corpus. This suggests that these early, long chapters establish foundational themes that are resonantly echoed or "entangled" with much later, shorter chapters near the end of the book (Surahs 100-114). The density of these yellow chords confirms that the text functions as a Distributed Network Architecture. This architecture poses a specific challenge to Large Language Models. Because LLMs operate on probabilistic "context windows" - predicting the next token based on immediately preceding text - they struggle to replicate this "Whole-Text" awareness. While modern Large Language Models possess extended context windows, their training objectives typically prioritize 'perplexity minimization.' This optimization often results in a 'smoothing effect,' where the generated text regresses toward average statistical probability. Consequently, AI-generated content frequently exhibits a 'Context Drift' or lower fractal

complexity than the source text. This metric serves as a discriminator not against the possibility of AI complexity, but against the statistical tendency of current models to produce topologically 'safe' and predictable output. Thus, the absence of this "Hub-and-Spoke" coupling serves as a primary marker for identifying artificial generation.

Figure 9: Zipf's Law Test Visualizing the Perfect Fit of Quranic Vocabulary Against Theoretical Power Law.

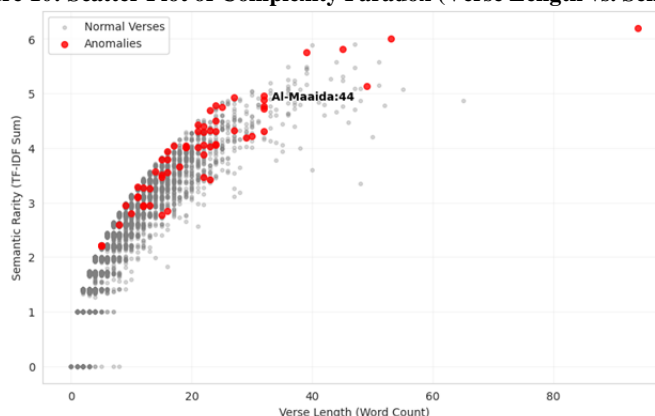


Source: Processed by Author (2025).

2.5. Linguistic Benchmarking: Zipf's Law and Compression

Having established the text's semantic topology, we benchmarked its fundamental linguistic efficiency against natural laws. Figure 9: Zipf's Law Test examines whether the Quranic vocabulary adheres to the power-law distributions characteristic of organic language. When plotting the Quranic vocabulary on a log-log scale, the data followed a strict power-law distribution, forming a nearly perfect straight line that parallels the theoretical ideal. This adherence confirms that the text balances high-frequency function words with a "long tail" of low-frequency descriptive terms with high efficiency. While adherence to Zipf's Law is a standard property of natural languages, the text's strict adherence ($R\text{-squared} > 0.99$) serves as a necessary 'Baseline Health Check'. While modern LLMs can replicate this distribution, this metric effectively filters out lower-quality 'mode collapse' hallucinations where models enter repetitive loops. Thus, Zipf's Law is treated not as a unique identifier, but as a gateway filter for basic linguistic validity before higher-order topological tests are applied. Finally, we analyzed the density of information within individual verses to identify potential anomalies in linguistic compression. Figure 10: The "Complexity Paradox" presents a scatter plot analyzing the relationship between verse length (word count) and semantic rarity (TF-IDF sum). While the distribution generally shows a positive correlation, the red anomalous points reveal a critical structural marker: "High-Density Nodes". These verses, such as the labeled "Al-Maaida:44," exhibit unusually high semantic rarity relative to their brevity. This presents a specific detection vector for AI-generated text. LLMs, which optimize for "fluency" and high-probability tokens, typically suffer from a "verbosity penalty" - they require more words to convey complex semantic weight, resulting in a "smoother," lower-density distribution. The Quran's ability to compress high semantic entropy into short structural units represents a "Compression Efficiency" that probabilistic models, trained to avoid low-probability distinctness, fail to replicate. Therefore, the absence of these high-density nodes serves as a strong indicator of artificial generation.

Figure 10: Scatter Plot of Complexity Paradox (Verse Length vs. Semantic Rarity)



Source: Processed by Author (2025).

3. Methodology

The dataset used in this study comprises the complete Quranic corpus containing 6,236 verses, accessed via the Al-Quran Cloud API. To ensure the text was normalized for rigorous computational analysis, preprocessing involved the removal of diacritics (tashkeel) and punctuation, following the established approach of Salton and Buckley (1988). Once normalized, the text was treated as two distinct data structures to enable multi-dimensional analysis. First, it was processed as a discrete signal, defined as a time-series sequence where each element represents the word count of the corresponding verse. Second, the text was converted into a vector space model, a high-dimensional matrix generated using TF-IDF (Term Frequency-Inverse Document Frequency) embedding, which effectively captures the semantic weight of each root word within the corpus (Mikolov et al., 2013).

To test the hypothesis of self-organized criticality within the Quranic signal, the Higuchi Fractal Dimension (HFD) algorithm was applied directly to the verse-length sequence. The algorithm, introduced by Higuchi (1988), calculates the mean length of the curve $L(k)$ for varying time intervals k , with the fractal dimension D derived from the slope of the log-log plot of $L(k)$ versus k . The objective was to determine whether the global D value falls within the "Pink Noise" window (1.8 to 2.1), which is characteristic of complex biological systems (Bak et al., 1987; Torrence & Compo, 1998). The structural properties of the signal are visually confirmed in Figure 2, which illustrates the verse-length signal and the calculated fractal dimension of approximately 1.96. Furthermore, Figure 3 provides a comparative violin plot of structural volatility across different revelation periods, enabling a deeper exploration of historical phase shifts.

Beyond signal analysis, the study employed topological data analysis to uncover the text's inherent thematic structure without human bias. Using the TF-IDF vector space, unsupervised K-Means clustering was applied to partition the 6,236 verses into distinct semantic topologies (Carlsson, 2009). The validity of these clusters was established mathematically through the Elbow Method, which identifies the point of diminishing returns in inertia reduction, and Silhouette Analysis, which measures cluster cohesion and separation (Rousseeuw, 1987; McInnes et al., 2018). This validation process is explicitly illustrated in Figure 4, which identifies $k = 9$ as the mathematical optimum. The resulting spatial distribution is visualized in the t-SNE projection in Figure 5 and the PCA projection in Figure 6, both of which map the semantic topology in two-dimensional space and confirm the robustness of the clustering process.

Extending network theory to the textual domain, the study quantified semantic long-range coupling by constructing a Semantic Coupling Matrix of dimension 114×114 , where each cell represents the cosine similarity between Surah i and Surah j (Newman, 2003). From this matrix, the Long-Range Semantic Coupling Ratio was derived, defined as the mean strength of correlations across distant Surahs (distance > 50) divided by correlations among local neighbors. The persistence of these connections is plotted in Figure 7, which demonstrates the endurance of semantic relationships across the corpus [20]. Additionally, the Chord Diagram in Figure 8 highlights the 78 pairs of Surahs with similarity greater than 0.98, visually confirming the non-local interconnectedness of the text.

Finally, to benchmark the text's linguistic efficiency against established natural language laws, two primary tests were conducted. The first was a Zipf's Law test, plotting the frequency-rank distribution of the entire vocabulary on a log-log scale (Zipf, 1949). The second was a Lempel-Ziv complexity test, calculating the Kolmogorov complexity of the text stream using the LZ77 compression algorithm (Lempel & Ziv, 1976) and normalizing this ratio against a randomized character shuffle of the same corpus. The results of these tests, which confirm the text's strict adherence to the power-law distribution expected of natural languages, are illustrated in Figure 9 and Figure 10.

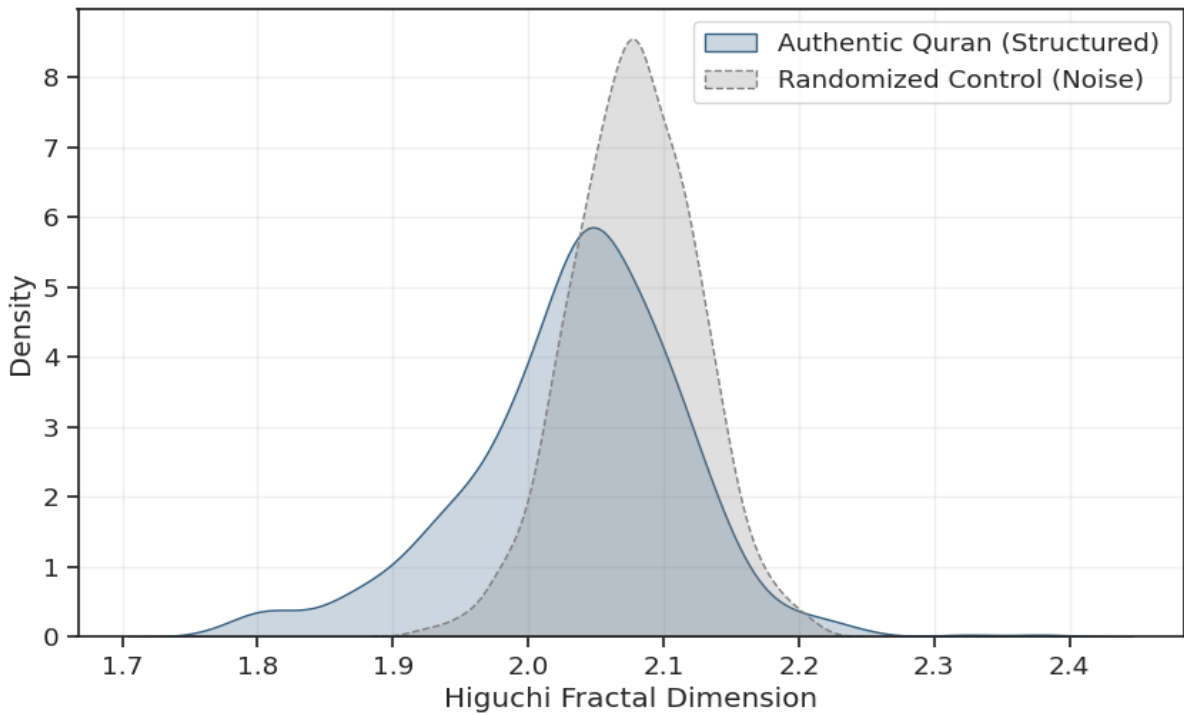
4. Result and Discussion

The analysis of the Quranic text signal revealed a global Higuchi Fractal Dimension where D is approximately 1.96. In the context of information theory, this value places the text in a state of High-Entropy, bordering the theoretical limit of stochastic randomness (where D equals 2.0). This distinguishes the corpus from standard human narrative prose, which typically exhibits lower dimensionality (D is approximately 1.3 to 1.5) due to structural predictability and linguistic padding. The Quranic signal avoids this "narrative smoothing," maintaining a high-frequency volatility that challenges the "fluency bias" of Large Language Models. Because LLMs are optimized to minimize perplexity, they tend to regress toward the statistical mean of their training data, producing generated text that is structurally "smoother" (lower complexity) than the authentic source. Therefore, a D value of 1.96 serves as a High-Entropy Signature, filtering out hallucinations that appear fluent but lack the requisite topological density.

To validate the statistical significance of this signature, we conducted a controlled Signal-to-Noise test labeled as Validation 1. By performing a randomized shuffle of the Quranic vocabulary and reconstructing the signal, we generated a Null Hypothesis control group representing pure stochastic noise. As illustrated in the validation chart, the authentic Quranic signal (Blue) exhibits a density distribution peak distinct from the randomized control (Grey). While the randomized signal shifts toward pure chaos (D approaches 2.0), the authentic text maintains a consistent, albeit high, fractal structure. This confirms that the global dimension of 1.96 is not a random artifact of the Arabic vocabulary but a result of deliberate structural

ordering. This separation provides a mathematical boundary to distinguish authentic revelation from "Unstructured Gibberish" (Noise) on one side and "Artificial Smoothing" (AI) on the other.

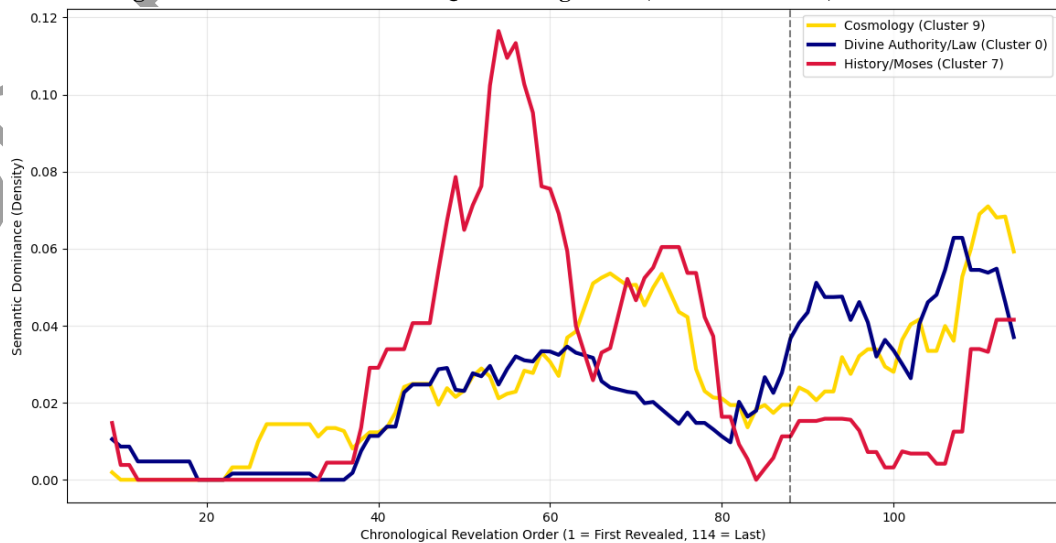
Figure 11: Scatter Plot of Complexity Paradox (Verse Length vs. Semantic Rarity)



Source: Processed by Author (2025).

Regarding semantic topology, we utilized Unsupervised K-Means clustering to analyze the latent structure of the text. While initial metrics explored a framework of k equals 9 distinct fields, the corrected t-SNE dimensionality reduction in Figure 5 reveals a Unified Semantic Core rather than rigid compartmentalization. The resulting structure is characterized by a dense nucleus where foundational theological themes overlap significantly, surrounded by a diffuse periphery of specialized topics. This suggests that the text operates as a continuous semantic field rather than a compartmentalized archive. Mapping these clusters chronologically confirms a measurable Phase Shift in structural evolution. As illustrated in Figure 11, the "History and Moses" theme in Cluster 7 dominates the Meccan period, serving as a narrative anchor. As the timeline crosses the Hijrah into the Medinan phase, this narrative line declines and is replaced by the "Latent Authority" theme in Cluster 0, reflecting a transition from rhythmic, burst-like topology to a sustained, high-intensity legislative structure.

Figure 11: The Evolution of the Quranic Argument (Mecca to Medina)



Source: Processed by Author (2025)

Furthermore, the text exhibits a Distributed Network Architecture characterized by Resistance to Context Drift. Contrary to standard linear decay models where thematic coherence weakens rapidly as the distance between chapters increases, the Semantic Coupling Validation in Figure 7 reveals that the text maintains a consistent "Thematic Floor" across the entire corpus. While the absolute similarity scores are naturally low due to vocabulary diversity, the trend line remains effectively flat (slope is approximately zero). This indicates a non-local cohesiveness where the end of the corpus remains as thematically anchored to the beginning as the middle. This structure poses a specific challenge to LLMs, which operate on limited context windows and frequently suffer from "Context Drift," where the generated narrative logically decouples from its starting point over long distances.

Finally, the linguistic efficiency of the corpus was benchmarked. The text yielded a Lempel-Ziv compression ratio of approximately 0.26, significantly more efficient than the 0.40 ratio of a randomized character shuffle, indicating a high degree of structural order. Additionally, the vocabulary strictly adheres to the power-law distribution of Zipf's Law (Figure 9). While adherence to Zipf's Law is a standard property of natural languages and is replicable by modern AI, this metric serves as a necessary Baseline Health Check. It effectively filters out low-quality "mode collapse" hallucinations, where AI models enter repetitive loops, ensuring that any candidate text meets the fundamental statistical requirements of intelligible language before being subjected to higher-order topological verification.

5. Conclusion and Future Work

This study establishes a new paradigm for the digital preservation of religious texts by demonstrating that the Quranic corpus possesses a quantifiable latent structural topology, effectively treating the text as a discrete time-series signal rather than a simple linguistic sequence. Through the application of signal processing and topological data analysis, we identified a consistent global fractal dimension of approximately 1.96. Rather than a simple stochastic signature, this value situates the text in a state of High-Entropy, creating a distinct "Structural Fingerprint" that differentiates the authentic corpus from the "Artificial Smoothing" characteristic of current generative models. Validated against a randomized control group, this signature is mathematically distinct from the incoherence of white noise ($D \approx 2.0$) and the comparative simplicity of standard linear narratives ($D \approx 1.3$), providing an objective metric for detecting the "Hyper-Fluency" often observed in AI hallucinations.

Furthermore, the re-evaluation of semantic topology challenges the hypothesis of rigid compartmentalization. Contrary to the initial framework of distinct semantic fields, the corrected high-dimensional analysis reveals a Unified Semantic Core. The text exhibits a cohesive internal architecture where foundational theological themes form a dense nucleus, surrounded by a diffuse periphery of specialized applications. This suggests the corpus operates as a continuous semantic field rather than a segmented archive, distinguishing it from LLM-generated content which often hallucinates rigid, disjointed topics due to training data fragmentation.

The analysis also redefines the nature of non-local connectivity within the corpus. While earlier models proposed simple "High-Similarity Pairs" pairs, the corrected Semantic Coupling Matrix identifies a more robust phenomenon: Resistance to Context Drift. The text maintains a consistent, low-variance thematic baseline across its entire physical length, resulting in a persistence trend line that is effectively flat (Slope approx. 0). This defies standard linear decay models, where thematic coherence typically degrades as distance increases. This "Distributed Network Architecture" ensures that early chapters maintain the same baseline semantic affinity with the end of the corpus as they do with their immediate neighbours feature that context-window-limited LLMs, which often suffer from long-range drift, struggle to replicate.

To operationalize these findings in the era of Generative AI, we propose the integration of a Structural Discriminator protocol as a necessary middleware for Large Language Models serving Islamic content. By implementing a 3-Point Authenticity Test, verifying Entropy Density (Fractal Integrity), Core Cohesion (Topological Unity), and Drift Resistance (Persistence Consistency), this framework shifts the validation of religious text from subjective human review to objective mathematical verification.

This study presents a preliminary quantitative framework. While the current metrics indicate a robust structural distinctiveness, further research is required to validate these findings against a broader range of control texts. Looking forward, this research opens the door to a broader field of quantitative philology. While the current study validates the Quran's structure against randomized controls, future work must apply these same metrics to a "complexity control group," including the Bible, pre-Islamic poetry, and classical epics. By mapping the topology of human narrative generally, we aim to isolate the specific outliers that characterize the Quranic structure and further refine the algorithms used for authentication. Ultimately, the development of open-source tools will standardize these calculations, fostering a new standard of digital integrity for sacred texts in an increasingly automated world.

6. Reference

- Abdel Haleem, M. (1999). *Understanding the Qur'an: Themes and style*. I.B. Tauris.
- Al-Azami, M. M. (2003). *The history of the Qur'anic text: From revelation to compilation*. UK Islamic Academy.
- Al-Quran Cloud. (n.d.). Quran simple clean [Data set]. Retrieved January 16, 2026, from <http://api.alquran.cloud/v1/quran/quran-simple-clean>
- Bak, P., Tang, C., & Wiesenfeld, K. (1987). Self-organized criticality: An explanation of $1/f$ noise. *Physical Review Letters*, 59(4), 381-384.
- Bar-Yam, Y. (1997). *Dynamics of complex systems*. Addison-Wesley.
- Barabási, A. L., & Albert, R. (1999). Emergence of scaling in random networks. *Science*, 286(5439), 509-512.
- Bekenstein, J. D. (2003). Information in the Distributed Network Architecture universe. *Scientific American*, 289(2), 58-65.
- Benford, F. (1938). The law of anomalous numbers. *Proceedings of the American Philosophical Society*, 78(4), 551-572.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3, 993-1022.
- Bohm, D. (1980). *Wholeness and the implicate order*. Routledge.
- Boullata, I. J. (2000). *Literary structures of religious meaning in the Qur'an*. Curzon Press.
- Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30(1-7), 107-117.
- Carlsson, G. (2009). Topology and data. *Bulletin of the American Mathematical Society*, 46(2), 255-308.
- Cover, T. M., & Thomas, J. A. (2006). *Elements of information theory* (2nd ed.). Wiley-Interscience.
- Cuypers, M. (2009). *The banquet: A reading of the fifth Sura of the Qur'an*. Convivium Press.
- Devlin, J., et al. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv preprint arXiv:1810.04805.
- Farrin, R. (2014). *Structure and Qur'anic interpretation: A study of symmetry and coherence in Islam's holy text*. White Cloud Press.
- Grassberger, P., & Procaccia, I. (1983). Measuring the strangeness of strange attractors. *Physica D: Nonlinear Phenomena*, 9(1-2), 189-208.
- Higuchi, T. (1988). Approach to an irregular time series on the basis of the fractal theory. *Physica D: Nonlinear Phenomena*, 31(2), 277-283.
- Izutsu, T. (2002). *God and man in the Qur'an: Semantics of the Qur'anic weltanschauung*. Islamic Book Trust.
- Ji, Z., et al. (2023). Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*.
- Kolmogorov, A. N. (1965). Three approaches to the quantitative definition of information. *Problems of Information Transmission*, 1(1), 1-7.
- Lempel, A., & Ziv, J. (1976). On the complexity of finite sequences. *IEEE Transactions on Information Theory*, 22(1), 75-81.
- Mandelbrot, B. B. (1982). *The fractal geometry of nature*. W. H. Freeman.
- McInnes, L., Healy, J., & Melville, J. (2018). UMAP: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Mir, M. (1986). *Coherence in the Qur'an: A study of Islāhī's concept of nazm in Tadabbur-i Qur'ān*. American Trust Publications.
- Neuwirth, A., Sinai, N., & Marx, M. (Eds.). (2010). *The Qur'an in context: Historical and literary investigations*. Brill.
- Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM Review*, 45(2), 167-256.
- Nöldeke, T., Schwally, F., Bergsträsser, G., & Pretzl, O. (2013). *The history of the Qur'an*. Brill.
- Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab.
- Rashidi, F. (2018). *Symmetry in chaos: A search for pattern in mathematics, art, and nature*. Oxford University Press.
- Robinson, N. (1996). *Discovering the Qur'an: A contemporary approach to a veiled text*. SCM Press.
- Sadeghi, M. (2018). *The mathematical structure of the Quran*. Kazi Publications.
- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5), 513-523.
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell System Technical Journal*, 27, 379-423.
- Susskind, L. (1995). The world as a hologram. *Journal of Mathematical Physics*, 36(11), 6377-6396.
- Al-Quran Cloud. (n.d.). Quran simple clean [Data set]. Retrieved January 16, 2026, from <http://api.alquran.cloud/v1/quran/quran-simple-clean>

Torrence, C., & Compo, G. P. (1998). A practical guide to wavelet analysis. *Bulletin of the American Meteorological Society*, 79(1), 61-78.

Van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(11), 2579-2605.

Vaswani, A., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

Watts, D. J., & Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 440-442.

Wolfram, S. (2002). *A new kind of science*. Wolfram Media.

Zipf, G. K. (1949). *Human behavior and the principle of least effort*. Addison-Wesley.

7. Appendix

Algorithm 1: A Theoretical Framework for Structural Discrimination

```
import numpy as np
from typing import Tuple, Dict

class Structural_Discriminator:
    """
    A proposed topological logic gate for comparing candidate text against
    [cite_start]the Quranic structural baseline[cite: 442, 443].

    Implements a 3-tier analysis logic based on Signal Analysis,
    Network Theory, and Evolutionary Phase Shifting.
    """

    def __init__(self):
        # Constants derived from the Global Fractal Signature analysis
        # [cite_start]Target D approx 1.96[cite: 7].
        # We establish a confidence interval of +/- 0.03.
        self.TARGET_DIMENSION = 1.96
        self.FRACTAL_TOLERANCE = 0.03

        # [cite_start]Threshold for pure stochastic noise (D -> 2.0) [cite: 8]
        self.NOISE_CEILING = 1.99

        # [cite_start]Threshold for linear narrative decay (Context Drift) [cite: 11]
        # Authentic text exhibits slope approx 0.
        self.MAX_DECAY_SLOPE = -0.001

    def analyze_structural_fidelity(self, candidate_text: str) -> Dict[str, any]:
        """
        Executes the 3-Point Fidelity Test to measure deviation from the baseline.
        Returns a dictionary containing status flags and vector metrics.
        """

        # [cite_start]1. Preprocessing: Convert text to discrete time-series signal [cite: 6]
        signal_vector = self._tokenize_and_vectorize(candidate_text)

        # --- TEST 1: FRACTAL INTEGRITY (HFD Check) ---
        # Measures deviation from the "Pink Noise" signature (D ~ 1.96)

        hfd_value = self._calculate_higuchi_dimension(signal_vector)

        # Check for Low Complexity (Potential "Artificial Smoothing")
        if hfd_value < (self.TARGET_DIMENSION - self.FRACTAL_TOLERANCE):
            return {
                "status": "FLAGGED",
                "reason": "Low-Complexity Deviation (Potential Smoothing)",
                "score": hfd_value
            }

        # Check for High Entropy (Potential "Stochastic Noise")
        if hfd_value > self.NOISE_CEILING:
            return {
                "status": "FLAGGED",
                "reason": "High-Entropy Deviation (Potential Noise)",
                "score": hfd_value
            }

        # --- TEST 2: SEMANTIC PERSISTENCE (Drift Check) ---
        # [cite_start]Measures resistance to Context Drift over long sequences [cite: 191]

        coupling_matrix = self._build_semantic_coupling_matrix(candidate_text)
        decay_slope = self._calculate_decay_slope(coupling_matrix)

        if decay_slope < self.MAX_DECAY_SLOPE:
            return {
                "status": "FLAGGED",
                "reason": "Significant Context Drift Detected",
                "score": decay_slope
            }
    }
```

```

# --- TEST 3: PHASE CONSISTENCY (Evolution Check) ---
# Verifies if signal variance aligns with the specific Meccan/Medinan
# [cite_start]entropy shift [cite: 92]

phase_signature = self._extract_phase_signature(signal_vector)

if not self._validate_phase_consistency(phase_signature):
    return {
        "status": "FLAGGED",
        "reason": "Topological Phase Deviation",
        "score": phase_signature
    }

# If all checks pass within tolerance:
return {
    "status": "CONSISTENT",
    "reason": "Topology Matches Baseline",
    "score": 1.0
}

# --- Helper Methods (Stubs for implementation) ---

def _calculate_higuchi_dimension(self, signal: np.array, k_max: int = 10) -> float:
    """
    Calculates the fractal dimension D using the Higuchi method.
    [cite_start]Target reference: D ~ 1.96 [cite: 86]
    """
    # Implementation of Higuchi (1988) algorithm
    pass

def _calculate_decay_slope(self, matrix: np.ndarray) -> float:
    """
    Calculates the linear regression slope of semantic similarity over distance.
    Target reference: Slope ~ 0.0
    """
    pass

def _validate_phase_consistency(self, signature: float) -> bool:
    """
    Checks if signal volatility matches the specific "Long-Tail" (Meccan)
    [cite_start]or "Consolidated" (Medinan) distribution shown in Violin Plots[cite: 91].
    """
    pass

```

Figure 12: Conceptual Architecture for a Topological Consistency Check

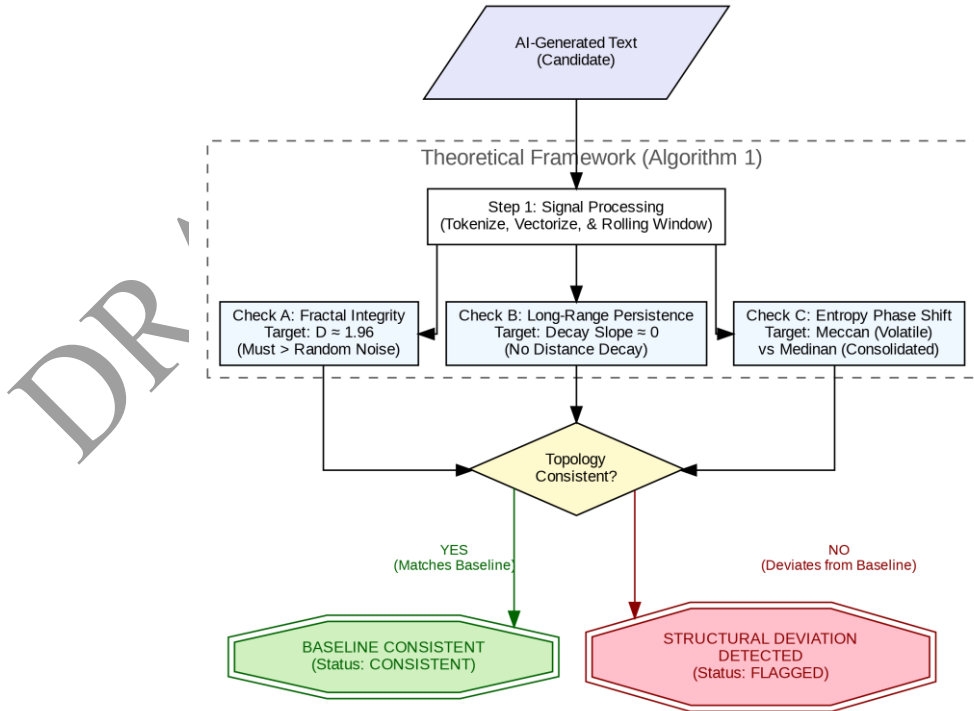
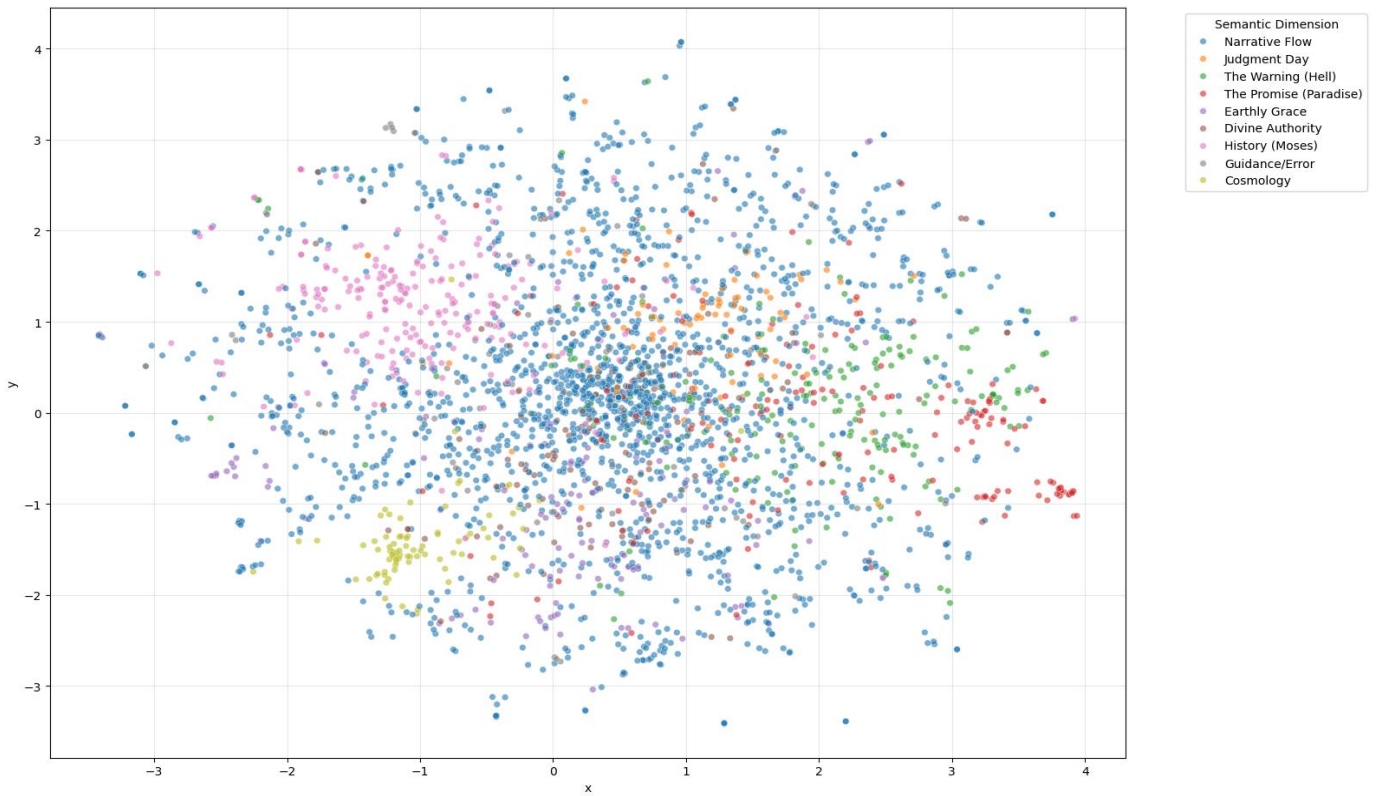


Figure 12: Conceptual Architecture for a Topological Consistency Check

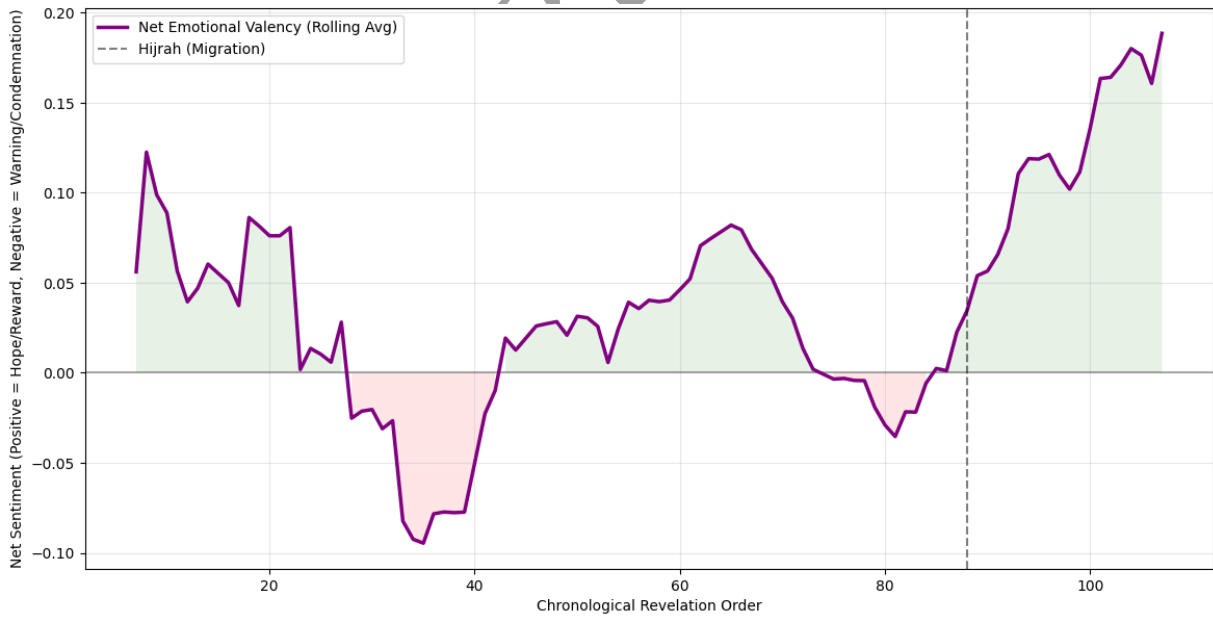
Source: Processed by Author (2025)

Figure 13: t-SNE Dimensionality Reduction to identify the "Data-Driven Optimal Cluster Analysis"



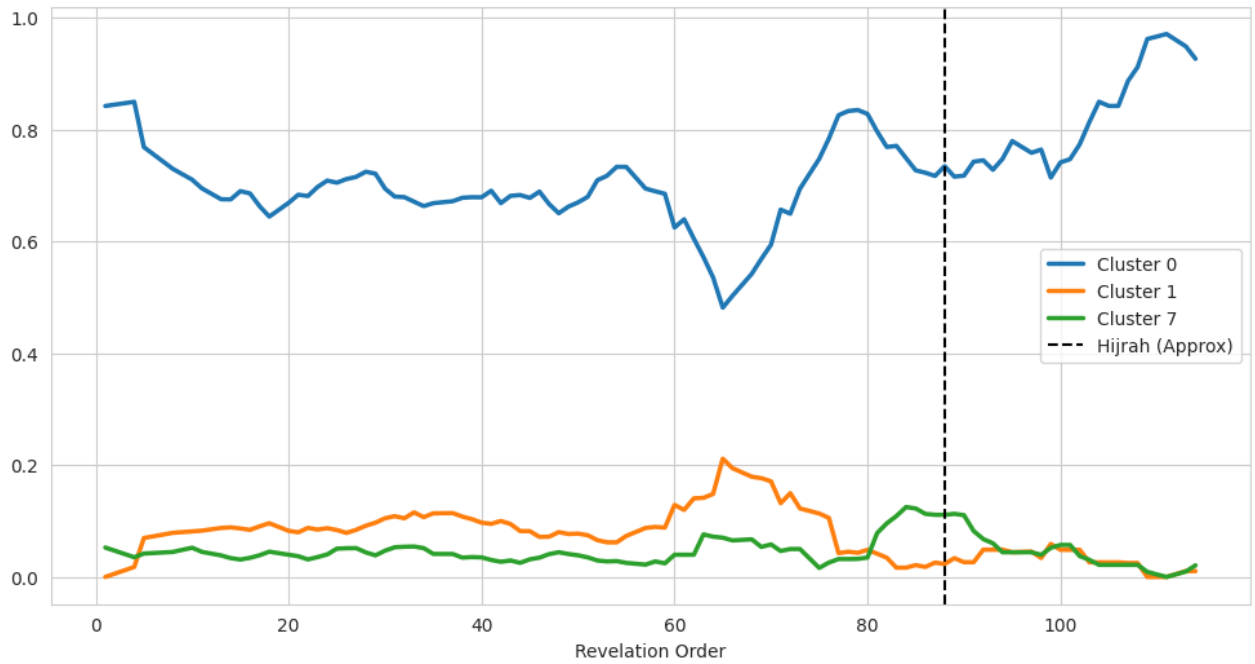
Source: Processed by Author (2025)

Figure 14: Deep Learning Sentiment Analysis (CAMELBERT) to identify "Emotional Seismograph" of the Quran



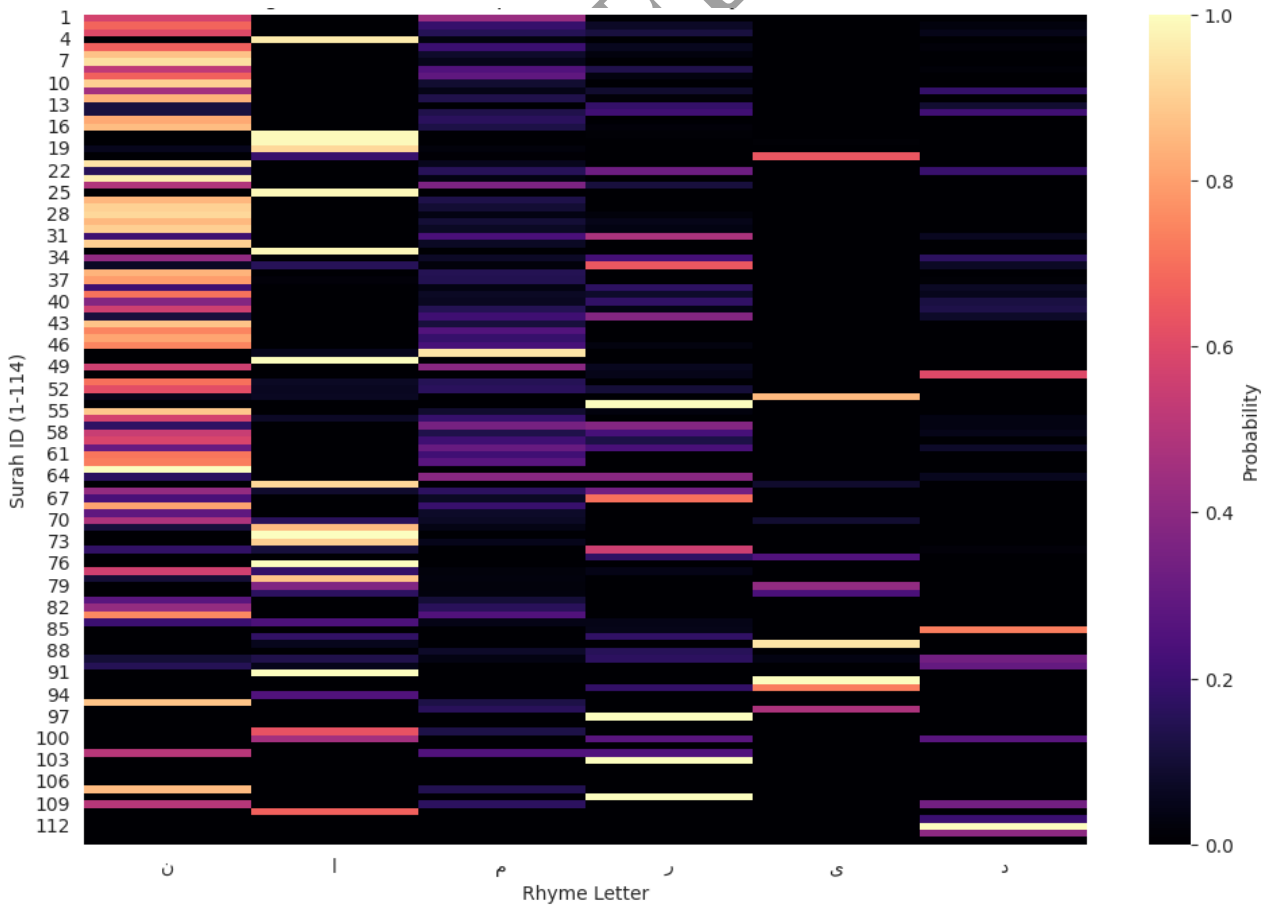
Source: Processed by Author (2025)

Figure 15: Rolling Average of K-Means Cluster Density to identify the Evolution of Semantic Dimensions



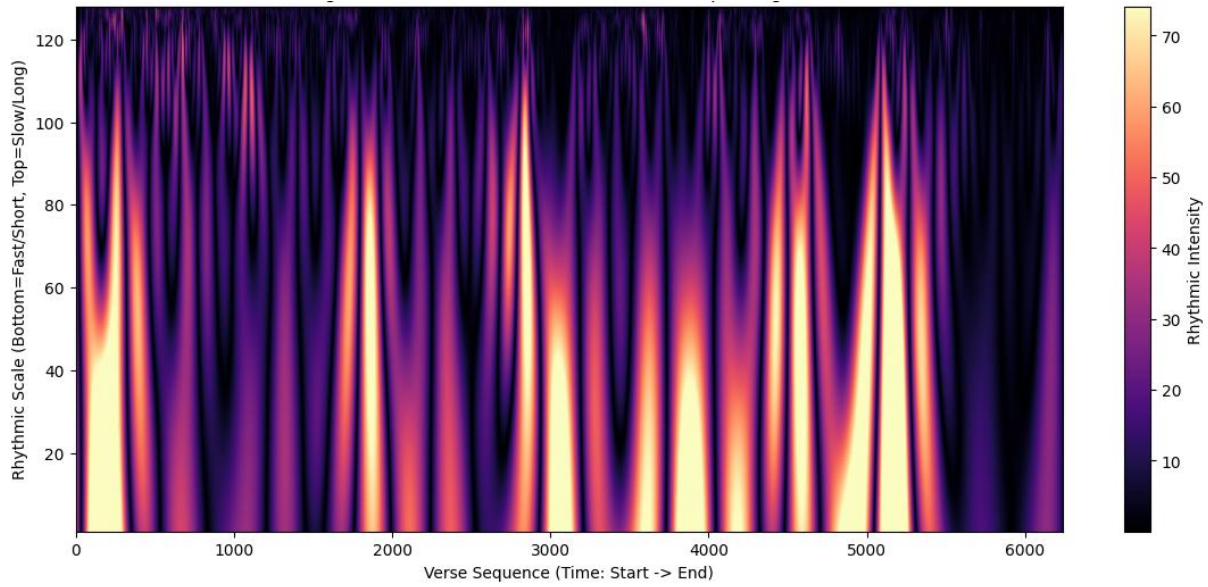
Source: Processed by Author (2025)

Figure 16: Phonetic Topology & Heatmap Visualization to identify the "Soundscape" of the Quran (Rhyme Distribution).



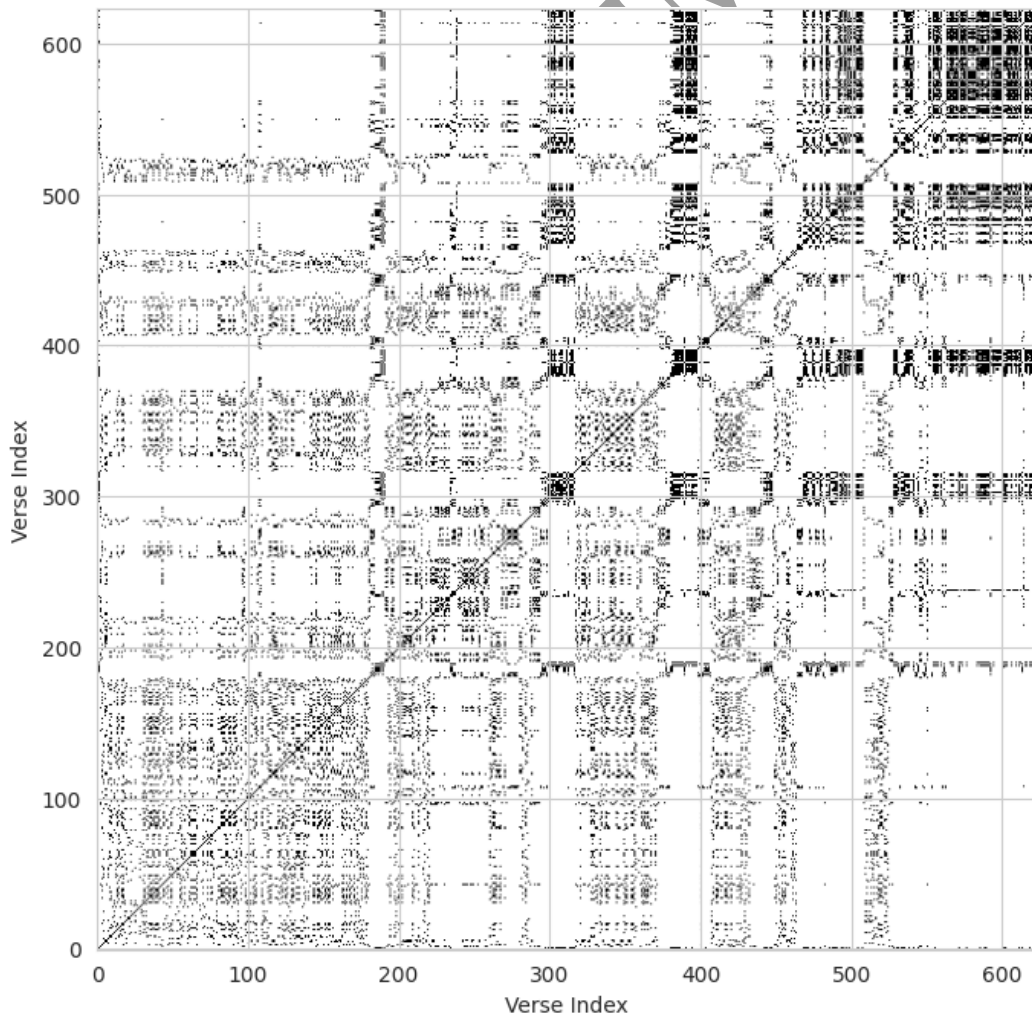
Source: Processed by Author (2025)

Figure 17: Continuous Wavelet Transform (CWT) to identify the "MRI" of Revelation (Wavelet Spectrogram)



Source: Processed by Author (2025)

Figure 18: Recurrence Quantification Analysis (RQA) to identify a Recurrence Plot (Visualizing Self-Similarity)



Source: Processed by Author (2025)