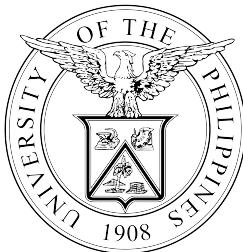


Text Analytics of the Qur'ān

by

Al-Ahmadgaid Bahauddin Asaad

Submitted to the *Institute of Islamic Studies*
in partial fulfillment of the requirements for the degree of
Master of Arts in Islamic Studies



UNIVERSITY OF THE PHILIPPINES DILIMAN
Diliman, Quezon City

May 2025

Abstract

The interest of the paper is to provide a comprehensive text analytics of the Qur'an. This is by utilizing Statistical and Machine Learning methods to computationally analyze the said scripture. Specifically, the following procedures have been done: descriptive analyses of the structure of the Qur'an, morphological analyses of the Qur'an. Further, the data used for the Qur'anic Arabic corpus is the one in QuranTree.jl by Asaad (2022). The computational software used is Julia programming language.

Chapter 1

Introduction

The paper studies the Qur’ān through the use of scientific computing, and the new perspective and methodology it brings will therefore be of interest, respectively, to researchers from Islamic studies, and researchers from the field of Statistics, Machine Learning, and Data Science. Thus, it is necessary to provide the essential background on Quranic studies to give more context on the rationales of the paper.

1.1 Background

The Qur’ān or *al-qur’ān* ﴿الْقُرْآن﴾ meaning *the recitation*, the holy book of Islam, is revered by 1.9 billion (according to 2020 projection of Cooperman et al. (2011, p. 13)) Muslims across the globe as the literal words of God. Muslims believed that the Qur’ān was gradually revealed (Qur’ān 25:32) to Prophet Muhammad ﷺ through angel *gibrīl* چَبِرْلِيلَ or Gabriel (Qur’ān 2:97). The Qur’ān contains 77,429 Arabic words in total, which covers only 56 percent of the Greek New Testament which has 138,020 words in total (Sinai, 2017, p. 11).

The Qur’ān is divided into *sūrahs* سورٌ which are the equivalent of chapters, each containing *ayāt* آيات (meaning *signs*), which are the equivalent of verses. The *sūrahs* سورٌ are not arranged in chronological order as in the Bible’s books and chapters, but rather arranged in monotonically decreasing length of number of verses after the first *sūrah* سُورَة (see Figure 1.1). The *sūrah* سُورَة of the Qur’ān can be categorized into two types: the *makkiyya* مَكْيَّة (Meccan) and *madaniyya* مَدْنَيَّة (Medinan). The categories refer to the geographical location of where the *sūrah* سُورَة was revealed. Figure 1.1 shows the grouping of the *sūrahs* سورٌ. Note that some of the *sūrahs* سورٌ have mixed geographical locations¹, that is, a few of the *ayāt* آيات in it were revealed in other geographical location apart from the geographical location of the rest of the *ayāt* آيات. Therefore, the categorization in Figure 1.1 highlights the geographical

¹see list of the location in https://tanzil.net/docs/revelation_order

location of the majority of the *ayāt* آيات in the *sūrah* سُورَة.

The Qur'ān (meaning *the recitation*) was revealed *orally* by angel ḡibrīl جَبْرِيل to Prophet Muhammad ﷺ and passed onto other believers through oral tradition (reciting the Qur'ān to students repeatedly so as to memorize it, instead of writing it down and let the believers read it and memorize it). Memorizing 77,429 Arabic words of the Qur'ān through oral transmission can be a difficult task, but what aids this memorization is the rhythm feature of the Qur'ān. According to one Orientalist, Sinai (2017), "rhyme, however, or rather a periodically recurrent word-final assonance, is a feature of the Qur'ān throughout, and it naturally partitions the *sūrah* سُورَة." Indeed, because of this feature, it makes it easy to memorize the entire Qur'ān, and the one who do so is called *hafiz* حَفِظ meaning *one who remembers* or

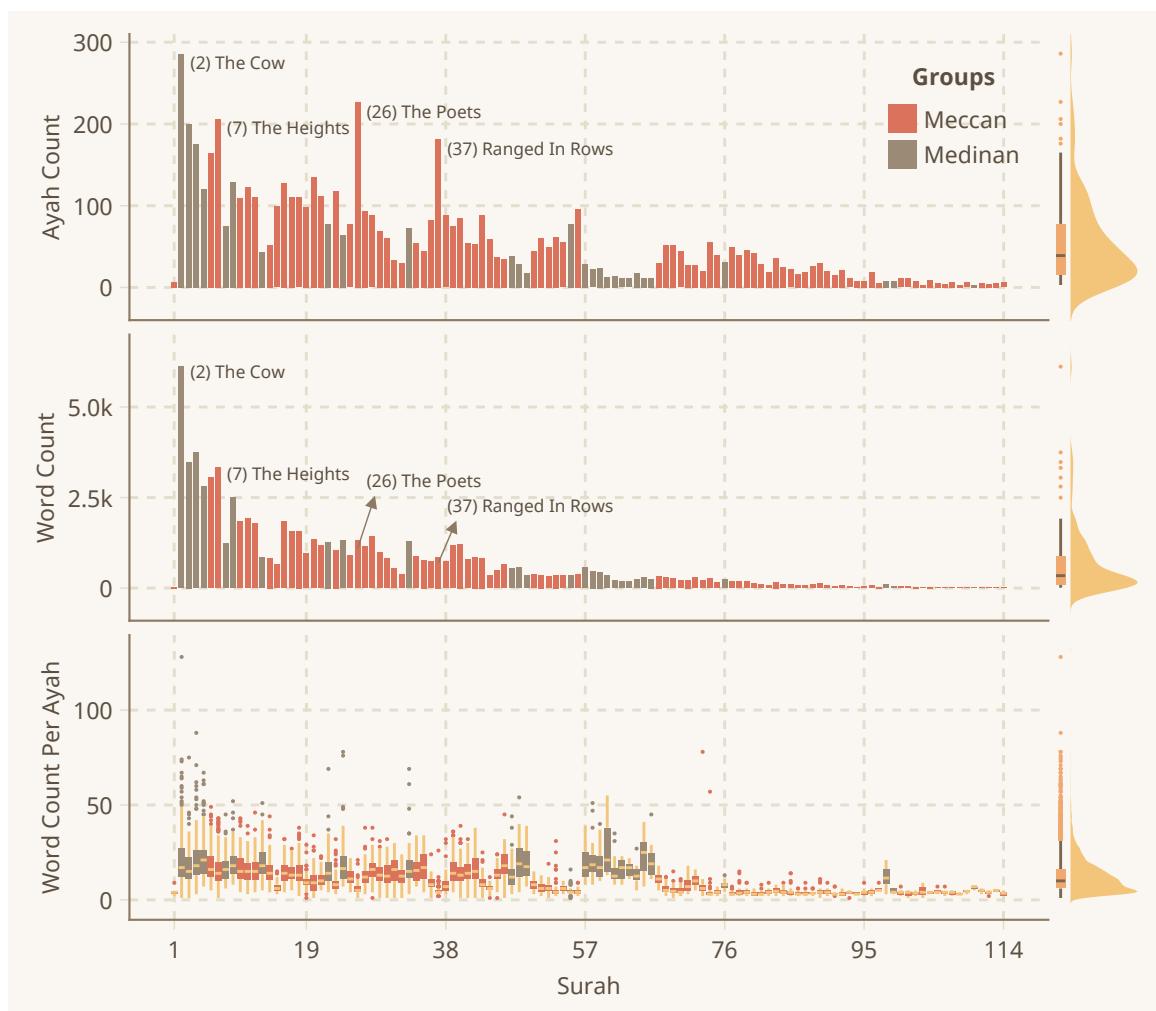


Figure 1.1: Statistics of the words and *ayāt* آيات (verses) of the Qur'ān

keeper. Qur'ān memorization contest is a common event in Muslim countries, the Philippine embassy has hosted one in 2022 in Saudi Arabia (Manila Bulletin, 2022).

According to the Muslim tradition, the oral transmission of passing the Qur'ān from a *hafiz* حفظ to new believers was gradually put into writings as requested by the believers themselves. The idea was brought up after the battle of Yamama, where many of the Muslims who died were *qurrā'* قُرَاءٌ (*the one who properly recite the Qur'ān*), and so fearing that their numbers will reduce in other battle fields, Umar ibn al-Khattab عمر بن الخطاب (who became the second caliph) suggested to the first caliph, Abū Bakr 'Abd Allāh ibn 'Abī Quhāfa, أبو بكر عبد الله بن أبي قحافة or short for Abū Bakr أبو بكر, to collect the Qur'ān into writing. Abū Bakr then authorized Zaid ibn Thabit زيد بن ثابت for the task. According to Zaid, he started collecting from the leafless stalks of the date-palm tree and from the pieces of leather and hides and from the stones, and from the chests of men (who had memorized the Qur'ān, i.e. the *hafiz* حفظ)². Long story short, the effort was finally codified by the third caliph, Uthman ibn Affan عثمان بن عفان in the year 645 CE, which was then recopied and distributed to the different regional capitals of the early Islamic empire of that time. The rest of the copies outside this codification were then burned³ down in order to have one standard Qur'ān. The Qur'ān nowadays is therefore assumed to be based on Uthmanic codex because of the story mentioned. That is, if indeed Uthman has ordered to burn other copies of the Qur'ān outside his codification, then what's left should only be based on his codex or archetype, and that should only be the inherited codex of the Muslims today.

The Qur'ān is believed by the Muslims to have been preserved since it was first recited by angel *gibrīl* جبريل or Gabriel to the Prophet ﷺ. Many orientalists had been skeptic about this claim, for example, John Wansbrough theorized that the Qur'ān was collected over a 200-year period (see Wansbrough, 2004, p. 101) after the death of the Prophet ﷺ, instead of within a few years after the death of the Prophet ﷺ. However, recent findings through radiocarbon dating brings forward strong evi-

²see <https://sunnah.com/bukhari:7191>

³see <https://sunnah.com/bukhari:4987>

dence of potential preservation of the whole Qur'ān, which the Muslims believed to be so. For example, the Birmingham Qur'ān manuscript discovered in 2015 is dated to be between 568 and 645 CE with 94.5% accuracy, making it among the oldest Qur'ānic manuscript in the world (see Birmingham University, 2015). Its predicted range of years intersects with the lifetime (570 to 632 CE) of the Prophet ﷺ. What is interesting is that the Birmingham Qur'ān is consistent with the Qur'ān today, word-by-word and letter-by-letter⁴, see Figure 1.2. This is indeed another evidence that the Qur'ān today was codified by Uthman since the discovery of the Birmingham Qur'ān manuscripts have confirm it. In addition to this, the Sana'a Palimpsest is also among the oldest Qur'ān radiocarbon dated to be between 578 CE and 669 CE with 95% accuracy (Sadeghi & Bergmann, 2010), which according to Sinai (2017), "neither does the edited portion of the Sana'a palimpsest offer evidence for additional or missing verses or for a divergent verse order within the *sūrahs* سور." Given these discoveries on the recent Quranic manuscripts, the claim of Wansbrough (2004, p. 101) is now untenable (see Sinai, 2014).

One likely reason as to why the scribes were able to preserve the Qur'ān in the two folios of the Birmingham Qur'ān manuscripts follows from the fact that the Qur'ān is firstly memorized before it was decided to be written. So that, the Topkapi manuscripts that contain 99% (only 23 verses missing)—dated around 701 to 750 CE (Karatay, 1962)⁵, that is, around 69 to 118 years after the death of the Prophet ﷺ—of the Qur'ān today was likely partly written from memory. This is possible since the Qur'ān possesses a rhythmic feature (that aids with memorization) that naturally divides its verses or *ayāt* آيات, and that the Qur'ān memorization competition is still held to this day (as in the example of Manila Bulletin, 2022), apart from the fact that it needs to be recited (any chapter after the first chapter of the Qur'ān according to the choice of the worshipper) every prayer from memory. Bottomline,

⁴The orthography of the Arabic letters in the early days had no diacritics and were written in its basic consonantal skeleton. That is, Arabic orthography and grammar were in their nascent when the Qur'ān was revealed, and therefore has to adjust and catch up, to capture and preserve the proper recitation of the Qur'an.

⁵see also <https://corpuscoranicum.de/en/manuscripts/1977/page/1-410?sura=1&verse=1>

there are many avenues for Qur'ān recitation from memory, and these have helped in its preservation. Moreover, since there are no significant evidence of insertion or malicious intention on addition or revision in all of the extant Qur'ānic manuscripts so far, some Orientalists came up with other theories of insertions on the basis of the literary style of the Qur'ān, see for example Sinai (2017, p. 92), where verse 102 of *sūra l-sāffāt* سُورَةِ الصَّافَاتِ or The Chapter of *Ranged in Rows* is theorized as addition because it is longer compared to other verses in the said chapter, refer to Sinai (2017, p. 92) for his other reasonings. Nonetheless, the Qur'ān is indeed stable based on the extant manuscripts.

Furthermore, the vastness of the early Islamic empire meant that different Muslim regional capitals have covered populations with different Arabic dialects, and so to accommodate these differences, Muslims believed that there were seven variant



Figure 1.2: 20th Century Qur'ān (left) in its fully featured orthographies vs Birmingham Qur'ān dated between 568 and 645 CE (right) in its basic consonantal skeleton. Image from Wikipedia (2015).

readings of the Uthmanic codex. Variant readings are defined as different pronunciations of the same word, in this case seven Uthamnic Qur’ān for seven different pronunciations. The *hadīt* حديث or *narration* comes from Ubayy ibn Ka'b who reported⁶ أبى بن كعب that the Prophet ﷺ was near the tank of Banu Ghifar that گبريل or Gabriel came to him and said: "... Allah has commanded you to recite the Qur’ān to your people in *seven dialects*, and in whichever dialect they would recite, they would be right." Recent work of Sidky (2020) shows that the material evidence on the regional variants is in remarkable agreement with well-attested written variants documented in the traditional Muslim literature.

Muslim and non-Muslim scholars alike have been extremely interested in understanding the unique literary characteristics of the Qur’ān. As mentioned earlier, unlike other books like the Bible (arranged in chronological order), the Qur’ān does not follow any obvious organization. In addition to this, a *sūrah* سورة does not fit the exact definition of a chapter. Indeed, the name attached to a *sūrah* سورة is often decided as the unique entity mentioned in the said *sūrah* سورة, its main purpose is to help early Muslims distinguish which *sūrah* سورة they are talking about, this is contrary to the chapter name where the associated name is obviously the main topic of the chapter. Further, as described by Sinai (2017), "... the compositional unity of the long surahs located at the beginning of the corpus is anything but obvious: at least at first sight, they can appear a flit back and forth between different topics in a largely haphazard manner. This impression is not limited to Western readers: even pre-modern Muslim scholars have often approached their scripture as a quarry of unconnected verses and groups of verses that bear little intrinsic relation to what precedes and follows." It wasn't until Neuwirth (2007), that the compositional unity of the Qur’ān can be observed in tighter literary unities, as Neuwirth (2007) showed that the many of these texts display a tripartite structure and are often constructive around a narrative middle part (Sinai, 2017). Samples of the organizational style of the Qur’ān was shown in Sinai (2017, p. 88).

⁶source: <https://sunnah.com/muslim:821a>

1.2 Rationale of the Study

Attempts at understanding the Qur'ān by Qur'ānic scholars were mostly done through manual process, that is, studying the scriptures by going through its content one-at-a-time manually. However, with the advent of computers some researchers have started using it to aid in their research. The first known to have used computers for studying the Qur'ān was likely Rashad Khalifa in 1968⁷, where he studied the significance of the mysterious initials at the beginning of some *sūrahs* سور. Rashad uploaded the Qur'ān into his computer by transliterating the Arabic letters and other Qur'ānic orthographies into Roman letters and symbols that the computer can easily parse. This approach of using computers to find new insights is more common in the field of science, and it was new for the field of Qur'ānic studies.

Indeed, to proceed with the use of scientific computing, the Qur'ān will be treated as the data that needs to be analyzed using what is called Natural Language Processing (NLP), a branch of Machine Learning (ML) that aims to understand natural language, such as Arabic. To instruct the computer to do Statistical analyses, ML, or NLP, one needs to use a *software application* or a formal language called *programming language*. There are several programming languages that the computer can understand. The popular one for researchers in the field of science are Python, R, and sometimes Julia (Bezanson et al., 2017). These programming languages will be used to construct instructions for computer. Therefore, if the data is the Qur'ān, then obviously there should be a way to interface with the Qur'ān using any of these programming languages, or to upload it and encode the Arabic letters into something that can be easily parsed by the computer, like what Rashad Khalifa did. There are, however, some programming languages with libraries or packages for interfacing with the Qur'an, and this is true for Python, R, and Julia. For this study, Julia programming language is used since its library for interfacing the

⁷see https://www.masjidtucson.org/quran/miracle/a_profound_miracle_sura68nun133.html

Qurān, QuranTree.jl⁸, is based on Tanzil⁹ for the Qurānic Arabic texts, and Dukes and Habash (2010) for morphological annotation, which both R and Python's libraries for Qurān does not have morphological annotation from Dukes and Habash (2010).

Machine Learning is a branch of Artificial Intelligence that aims to characterize data (in this case Arabic texts) by learning its features. These learned features are captured by an artifact called *model*. Basically, it is a model of some real data, meaning if there is a good representation of this real data, it is the model that has captured/learned/characterized the characteristics of the said real data. Therefore, a model is not exactly the real data, but an estimate of the characteristics of that data. To help understand this, and relate it to the fashion industry, which the author assumes most readers will be familiar with, a model in the fashion industry aims to capture the characteristics of the target customers, so for a clothing company, they hire a model which has the best features of their target customers. Hence, they hire Asian models to target Asian customers. So that, when these models wore the clothes sold by the said company, the potential customer will more or less be able to relate to the model, and be able to imagine themselves wearing that same clothes as well, which help them incline to buying the said clothing. The model therefore does not necessarily have the looks of every target Asian customer, but at least in terms of height, skin tone, hair, and other common Asian features, the model will likely have it. The question now is, what are the benefits that this model can bring to the clothing company? Well, the clothing company will be able to create products that are tailored to their Asian customers using the said model, since the company will have the right baseline measurements needed. Relating this analogy to the technical concept discussed earlier, you can think of the target customer as the real data, and the model as the same technical term use in Machine Learning and Statistics, but this time this technical model is expected to capture the characteristics of the real data analogous to fashion model that is expected to capture the characteristics

⁸see <https://alstat.github.io/QuranTree.jl/stable/>

⁹<https://tanzil.net/download/>

of the target customer. This Statistical or Machine Learning model brings the following benefit: researchers will be able to study the real data by simply using the model to answer questions that are not available in the sampled real data.

1.3 Objectives

The objective of the paper is to answer the following research questions:

1. What are the thematic themes of the surah that can be extracted by a Statistical and Artificial Intelligence (AI) models?
2. What other insights can be extracted by AI and how it compares or supplements the related findings from Muslim's traditional sources?
3. How this combination of AI and Muslim's traditional literatures help in understanding the Qurān?

1.4 Significance of the Study

While the Qur'ān has been extensively studied by Muslims and non-Muslims scholars alike, especially in the topic of Meccan and Medinan surahs, there is still a lot to uncover from the perspective of Computational Statistics. Hence, the significance of this study is that it brings forward new ways of extracting insights from the Qur'ān by leveraging Computations, Statistics, Machine Learning, and AI, that is still in its early stage in the field of Qur'ānic Studies. Therefore, this new perspective or process of studying the scripture not only aids the scholars of the Islamic Studies, but may also contribute indirectly to community development and policy makers who use Qur'ān as part of their decision making.

1.5 Scope of the Study

The paper will cover all chapters of the Qur'ān in the analyses. It will also not delve too much into the tafsir of each of the verses in the analyses, but will do on a few that may need further context.

1.6 Thesis Organization

The paper is organized as follows: Chapter 2 will discuss the related literatures, Chapter 3 will discuss the methodology, Chapter 4 will present the results and discussions, and finally Chapter 5 will contain the conclusion and recommendation. The references and appendices are placed after the Chapter 5.

Chapter 2

Literature Review

As mentioned in the previous chapter, the earliest work in Qur'ānic studies using computer was likely the work of Rashad Khalifa in 1968¹, which led to one of his book entitled 'The Computer Speaks: God's Message to the World' (*see Khalifa (1981)*). While the work of Rashad started at studying the mystery letters in the beginning of some *sūrahs* سور (for example Qur'ān 2:1, 3:1, 7:1, etc.), it quickly went on to cover what he calls other *mathematical miracles*, all of which are covered in Khalifa (1981). Rashad went on to claim that these findings meant that God revealed His words through this mathematical patterns through out the Qur'ān, and that those verses that were off and did not conform to this discovered mathematical patterns led him to extensive investigation of the said verses, and concluded that those could be or surely be an insertion that should not have been in the Qur'ān in the first place. There are two verses that were off, and Rashad called these verses as *false verses*², these are the last two *ayāt* أیات of *sūra l-tāwba* سورة التوبة or The Chapter of Repentance. These two verses were removed in Rashad's Qur'ān translation³. Rashad believed so much on his findings that he claimed to be a messenger⁴ with this new findings and that the Qur'ān nowadays should conform to his found mathematical patterns. This self-proclamation led to his assassination.

Fast forward to 20th century, Thabet (2004) started developing a stemmer system for the Qur'ān. A stemmer system is a system for trimming inflected words into its basic form, which is the root. For example, in English the root word for *computational*, *computer*, *computation*, and *computerize* is *compute*. Therefore, from the root forms different stems representing the different words. Hence, the idea of stemming is to trim these words into its basic form. The use case of stemming is

¹https://www.masjidtucson.org/quran/miracle/a_profound_miracle_sura68nun133.html

²<https://submission.org/App24.html>

³<https://www.masjidtucson.org/quran/frames/>

⁴https://www.masjidtucson.org/submission/faq/rashad_khalifa_summary.html

finding groups of words, which by using the root of the word makes it easy to find relations or similarity. This is the work done by Thabet (2004), which according to the author it is even more a challenge for Arabic words since it is highly inflected, and more so for the Classical Arabic texts like the Qur'ān.

Building on the work of Thabet (2004), Thabet (2005) used a statistical methodology for clustering the chapters of the Qur'ān, in particular using the agglomerative hierarchical clustering. The data processing makes use of the stemming methodology in Thabet (2004) to remove the different inflections on the Qur'ānic words. Moving on, the work of Noordin and Othman (2006) focused on information retrieval of Qur'ānic content by surveying 125 websites and investigating how Qur'ānic informations are presented, their aim is to propose a system for retrieving these information.

The work of Sharaf and Atwell (2009) studied knowledge representation of the Qur'ānic verb valences using FrameNet frames, the output of which is a lexical database of the corpus of Qur'ān verbs. Further, the work of Sharaf and Atwell (2012a) came up with corpus for the annotations of the Qur'ānic pronouns, the authors named it as QurAna. Building on this work, Sharaf and Atwell (2012b) came up with a corpus for studying Qur'ānic relatedness based on the commentary of Ibn Kathir ﴿ابن كثیر﴾, the authors named this corpus as QurSim.

Moving on, an unpublished work by Nassourou (2011) used Machine Learning to study.

Chapter 3

Methodology

Chapter 4

Results and Discussions

This chapter will discuss the results of the study. The chapter is organized as follows: Section 4.1 discusses

4.1 Descriptive Statistics

In this section, the count statistics for the verses and words of the Qur'an. Figure 1.1 shows the count of *ayāt* آيات (verses) and word of the Qur'an.

4.1.1 Verses

4.2 Morphological Analysis

4.3 Structure Analysis

4.3.1 Concentric Structure

4.3.2 Mathematical Structure

Chapter 5

Conclusion and Recommendation

References

- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. B. (2017). Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1), 65–98. Retrieved from <https://pubs.siam.org/doi/10.1137/141000671> doi: 10.1137/141000671
- Birmingham University. (2015). *Birmingham qur'an manuscript dated among the oldest in the world*. Birmingham University. (Available at: <https://www.birmingham.ac.uk/news-archive/2015/birmingham-quran-manuscript-dated-among-the-oldest-in-the-world> (Accessed: July 8th, 2023))
- Cooperman, A., O'Connell, E., & Stencel, S. (2011). *the future of the global muslim population* (Tech. Rep.). Pew Research Center.
- Dukes, K., & Habash, N. (2010, May). Morphological annotation of Quranic Arabic. In N. Calzolari et al. (Eds.), *Proceedings of the seventh international conference on language resources and evaluation (LREC'10)*. Valletta, Malta: European Language Resources Association (ELRA). Retrieved from http://www.lrec-conf.org/proceedings/lrec2010/pdf/276_Paper.pdf
- Karatay, F. E. (1962). *Topkapı sarayı müzesi kütüphanesi arapça yazmalar katalogu. kur'an, kur'an ilimleri, tefsirler no. 1 - 2171*. Küçükaydın Matbaası, İstanbul.
- Khalifa, R. (1981). *The computer speaks: God's message to the world*. Renaissance Production.
- Manila Bulletin. (2022). *PH embassy in riyadh hosts first asian qur'an memorization contest*. Manila Bulletin. (Available at: <https://mb.com.ph/2022/04/30/ph-embassy-in-riyadh-hosts-first-asian-quran-memorization-contest/> (Accessed: July 8th, 2023))
- Nassourou, M. (2011). *Using machine learning algorithms for categorizing quranic chapters by major phases of prophet mohammad's messengership*.
- Neuwirth, A. (2007). *Studien zur komposition der mekkanischen suren*. Berlin, Boston: De Gruyter. Retrieved 2023-07-09, from <https://doi.org/10.1515/9783110920383> doi: doi:10.1515/9783110920383

- Noordin, M., & Othman, R. (2006). An information retrieval system for quranic texts: A proposed system design. In *2006 2nd international conference on information and communication technologies* (Vol. 1, p. 1704-1709).
- Sadeghi, B., & Bergmann, U. (2010). The codex of a companion of the prophet and the qurān of the prophet. *Arabica*, 57(4), 343 - 436. doi: \url{https://doi.org/10.1163/157005810X504518}
- Sharaf, A., & Atwell, E. (2009). Knowledge representation of the quran through frame semantics: a corpus-based approach. In *Proceedings of the fifth corpus linguistics conference*. The Fifth Corpus Linguistics Conference. Retrieved from <https://api.semanticscholar.org/CorpusID:18278736>
- Sharaf, A., & Atwell, E. (2012a, May). QurAna: Corpus of the Quran annotated with pronominal anaphora. In N. Calzolari et al. (Eds.), *Proceedings of the eighth international conference on language resources and evaluation (LREC'12)* (pp. 130–137). Istanbul, Turkey: European Language Resources Association (ELRA). Retrieved from http://www.lrec-conf.org/proceedings/lrec2012/pdf/123_Paper.pdf
- Sharaf, A., & Atwell, E. (2012b, May). QurSim: A corpus for evaluation of relatedness in short texts. In N. Calzolari et al. (Eds.), *Proceedings of the eighth international conference on language resources and evaluation (LREC'12)* (pp. 2295–2302). Istanbul, Turkey: European Language Resources Association (ELRA). Retrieved from http://www.lrec-conf.org/proceedings/lrec2012/pdf/190_Paper.pdf
- Sidky, H. (2020). On the regionality of qur'anic codices. *Journal of International Quranic Studies Association*. doi: <http://dx.doi.org/10.5913/jiqsa.5.2020.a005>
- Sinai, N. (2014). When did the consonantal skeleton of the quran reach closure? part ii. *Bulletin of the School of Oriental and African Studies*, 77(3), 509–521. doi: 10.1017/S0041977X14000111
- Sinai, N. (2017). *The qur'an: A historical-critical introduction*. Edinburgh University Press Ltd.

- Thabet, N. (2004). Stemming the qur'an. In *Proceedings of the workshop on computational approaches to arabic script-based languages* (p. 85-88). USA: Association for Computational Linguistics.
- Thabet, N. (2005). Understanding the thematic structure of the qur'an: an exploratory multivariate approach. In *Proceedings of the acl student research workshop* (p. 7-12). USA: Association for Computational Linguistics.
- Wansbrough, J. (2004). *Quranic studies: Sources and methods of scriptural interpretation*. Prometheus Books.
- Wikipedia. (2015). *Comparison of a 20th-century edition of the quran (left) and the birmingham quran manuscript (right)*. Wikipedia. (Available at: https://en.wikipedia.org/wiki/Birmingham_Quran_manuscript (Accessed: July 9th, 2023))