UNIVERSIDADE FEDERAL DE PERNAMBUCO
CENTRO DE INFORMÁTICA
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA DA COMPUTAÇÃO

DANIELA RAPOSO NUNES DE MELLO

**Which is the most suitable scanner resolution for documents?**

Recife

2025

DANIELA RAPOSO NUNES DE MELLO

**Which is the most suitable scanner resolution for documents?**

Dissertação de Mestrado apresentada ao Programa de Pós-graduação em Ciência da Computação do Centro de Informática da Universidade Federal de Pernambuco, como requisito parcial para obtenção do grau de Mestre em Ciência da Computação.

**Área de Concentração**: Processamento e Síntese de Mídias e Interação

**Orientador**: Prof. Dr. Rafael Dueire Lins

**Coorientador**: Prof. Dr. Raimundo Oliveira

Recife

2025

.Catalogação de Publicação na Fonte. UFPE - Biblioteca Central

**Daniela Raposo Nunes de Mello**

**"Which is the most suitable scanner resolution for documents?**

> Dissertação de mestrado apresentada ao Programa de Pós-Graduação em Ciência da Computação da Universidade Federal de Pernambuco, como requisito parcial para a obtenção do título de Mestre em Ciência da Computação. Área de Concentração: Mídia e Interação.

Aprovado em: 27/01/2025.

## BANCA EXAMINADORA

_____
Prof. Dr. Carlos Alexandre Barros de Mello
Centro de Informática / UFPE

_____
Prof. Dr. Gabriel de França Pereira e Silva
Unidade Acadêmica do Cabo de Santo Agostinho/UFRPE

_____
Prof. Dr. Rafael Dueire Lins
Centro de Informática / UFPE
(**orientador**)

# RESUMO

Definir a resolução correta de imagem para digitalização de documentos é essencial para garantir a preservação precisa de todas as informações necessárias e, ao mesmo tempo, otimizar a eficiência de aquisição, o tempo de processamento, a demanda de armazenamento e a largura de banda de transmissão. É fundamental que a resolução de digitalização garanta que o documento digitalizado, quando impresso em alta resolução nas dimensões originais, mantenha integralmente as informações do documento original, de forma a propiciar a preservação e manutenção da informação. Apesar de sua importância para a engenharia e preservação de documentos, os atuais padrões e recomendações de digitalização baseiam-se predominantemente em evidências experimentais e conhecimento empírico. Considerando esses parâmetros, a questão crucial "Qual é a resolução de digitalização mais adequada para documentos?" foi recentemente levantada pelo Professor George Nagy, um dos pesquisadores pioneiros na Engenharia de Documentos, destacando uma preocupação de longa data que não foi devidamente esclarecida pela comunidade de pesquisa. Este trabalho propõe uma resposta fundamentada nos princípios do teorema de Nyquist-Shannon da Teoria da Informação, esta abordagem tem por objetivo fornecer uma fundamentação teórica para a definição adequada da resolução na digitalização de documentos, de forma a assegurar a preservação das informações.

**Palavras-chaves**: Resolução de imagem. Digitalização de documentos. Teorema de Nyquist-Shannon

## ABSTRACT

Defining the correct image resolution for document digitization is essential to ensure the accurate preservation of all necessary information while optimizing acquisition efficiency, processing time, storage demand, and the transmission bandwidth. It is crucial that the digitization resolution ensures that the digitized document, when printed in a high resolution at its original dimensions, fully retains the information from the original document, thus enabling the preservation and maintenance of the information. Despite its importance for document engineering and preservation, current digitization standards and recommendations are predominantly based on experimental evidence and empirical knowledge. Considering these parameters, the crucial question, "Which is the most suitable resolution for document digitization?" was recently raised by Professor George Nagy, one of the pioneers in Document Engineering, highlighting a long-standing concern that has not been adequately addressed by the research community. This work proposes an answer based on the principles of the Nyquist-Shannon theorem from Information Theory, this approach aims to provide a theoretical foundation for determining the appropriate resolution for document digitization, ensuring the preservation of information.

**Keywords**: Image resolution. Document digitization. Nyquist-Shannon Theorem.

# LIST OF FIGURES

# LIST OF CODES

# LIST OF BOARDS

# LIST OF TABLES

# LIST OF ACRONYMS

**BMP**            Windows Bitmap

**bpp**            bits per pixel

**cm**             centimeters

**DocEng'24**      ACM Symposium on Document Engineering 2024

**dpi**            dots per inch

**HF**             Highest frequency

**ICDAR**          International Conference on Document Analysis and Recognition

**in**             inches

**LD**             Levenshtein Distance

**lp/mm**          line pairs per millimeter

**MB**             megabytes

**mm**             millimeters

**OCR**            Optical Character Recognition

**PNG**            Portable Networks Graphics

**PSM**            Page Segmentation Mode

**pt**             points

**RGB**            Red Green Blue

**TIFF**           Tagged Image Format

# LIST OF SYMBOLS

$\Delta$           Delta

$\sigma$           Sigma

$\mu$           Mi

$\pi$           Pi

# CONTENTS

# 1 INTRODUCTION

The correct determination of scanning resolution is essential to ensure that all information present in a document, especially in historical documents, is accurately captured. According to The Library of Congress (2014), the choice of spatial and tonal resolution directly impacts the quality of document preservation in digital images. If the scanning resolution is set too low, critical elements of the document may be lost, conversely, excessively high resolution drastically increases scanner acquisition time yielding too large images that are slow to process and claim unnecessary storage space and computer bandwidth for network transmission.

Back in 1991, Rafael Dueire Lins was challenged by Graziela Peregrino, head of the Documentation Department of the Joaquim Nabuco Foundation in Recife, Brazil, to find a way to preserve all the information on the rich file of letters of Joaquim Nabuco[1] to future generations Lins (2011). Using the only scanner available at the Universidade Federal de Pernambuco in Brazil then, Lins digitized several letters of Joaquim Nabuco in 75, 150, 200, 300, 500, and 700 dpi in true-color and printed the images using a high-quality color laser printer. The printed documents were handed to the documentation experts of the Joaquim Nabuco Foundation, who concluded that the 150 dpi images had good enough quality to preserve all the information of the original document. As there was a huge international interest in the letters of the Nabuco bequest by social science researchers from all over the world, the 200 dpi scanner resolution was adopted by then, as it was compatible with FAX-machine standards, the technology available to send documents via computer/phone networks.

Over the years, prominent institutions dedicated to document digitization and preservation, including The Digital Library Federation Benchmark Working Group (2002), The Library of Congress (2014) and Forschungsgemeinschaft (2016), have conducted empirical studies to identify and recommend optimal scanning resolutions for different document types. The technical guidelines provided by those institutions are primarily based on visual inspection of the images generated in which documents were scanned at different resolutions, with the optimal resolution selected based on the visual "quality" of the resulting images. Among the findings, it is possible to observe a consensus of approximately 300 dpi as the ideal minimum resolution for grayscale and RGB-mode documents. According to Puglia et al. (2004), this resolution can also serve as a baseline for other document types, such as illustrations, maps, plans, ma-

---

[1] Brazilian statesman, writer, and diplomat, the first Brazilian ambassador to the U.S.A, one of the key figures in the campaign for freeing black slaves in Brazil (b.1861-d.1910).

nuscripts, and large-format documents. Despite the agreement of using 300 dpi among those important institutions, there is no technical justification beyond empirical data analysis to support such recommendations.

During the welcome reception of the International Conference on Document Analysis and Recognition (ICDAR) 2023, Professor George Nagy asked Daniel Lopesti and Rafael Dueire Lins, the fundamental question: "Which is the most suitable scanner resolution for documents?". Rafael Dueire Lins promptly answered that the proper way to determine the resolution, considering the maintenance and preservation of information, would be through the Nyquist-Shannon theorem. Such answer was immediately acknowledged as providing the correct solution to the problem.

Whenever converting an analog signal into a digital representation two phases are involved. The first one is determining the number of samples one must take. As explained by Lathi (2005), the Nyquist-Shannon Theorem (1938) teaches that any (continuous) signal sampled with twice its maximum frequency may be re-constructed without losses. The second step in the digitization process is the quantization error, the maximum error tolerated in the representation of each sample.

Although the response given was immediately recognized as providing the solution to the question, as the theorem provides a theoretical basis for defining the minimum resolution required to capture all document details without information loss, ensuring that all the details of the document are retained whenever viewed or printed on the original size the validity of the solution had still to be proved. The suggested solution needed to be detailed and proved correct, which motivated the research developed in this dissertation.

The initial results of this dissertation Lins et al. (2024) was originally published in a short paper in the proceedings of the ACM Symposium on Document Engineering 2024 (DocEng'24), which was sent to Professors George Nagy, Daniel Lopesti, and Elisa Barney-Smith. Professor Nagy main feedback was:

> "... your experiments provide reassurance about the adequacy of generally accepted digitization parameters. Testing both bi-level and color scans, and using the spectral density for estimating the range of image frequencies, seem sound."

Professor Elisa Barney-Smith informed us that:

"It is an interesting paper. It was the topic of a WG at the DAS 2002 workshop, but no detailed experiments were run, just a discussion and WG report paper (published in IJDAR in 2004)."

Besides that, Professor Barney-Smith listed some important papers of hers on image scanning Barney-Smith (1998), Barney-Smith and Andersen (2005) showing that the quality of scanner output has evolved drastically:

"Some information about the scanner itself can be measured through the use of various test targets, including a simple knife edge (to return the step response of the optics). See several of my early papers where I spent a lot of time exploring this. Scanners today are much higher quality than scanners 30 years ago, we have the disk space and the processing power to process over-sampled images, so it isn't as critical today, but people still explore the scanner optics."

This point raised by Professor Barney-Smith is very important because the Shannon-Nyquist theorem does not address *quantization errors*. But it is also relevant to notice that such errors are device-dependent and that their effect would be more noticeable in higher-resolution images. Another concern raised by Professor Barney-Smith was:

"Another point you should mention in future papers would be a discussion of WHAT information you want to extract from a scan of a document. You touched on this when you mentioned historical documents vs maps... If it is just the text content, it has a much lower frequency content than if one were looking at the paper fibers of the same document."

It is important to stress here that this paper aims to determine the "correct" scanner resolution to keep all the document information making a high-definition same-size print of the scanned document indistinguishable from the original. However, on two occasions scanning documents in 300 dpi was enough to distinguish the age of documents as described by Barboza, Lins and Jesus (2013) and even to detect frauds with latter-added strokes in hand-written documents as detailed by Barbosa et al. (2014).

This dissertation also widens the scope of Lins et al. (2024) in two ways. The first one was the inclusion of more complex and fine-detailed images that employ tiny fonts that would generate very high-frequency components. The second path uses Optical Character Recognition (OCR) to analyze how the resolution variation affects the recognition rate.

## 1.1 OBJECTIVES

The general objective in this dissertation is to bring a theoretical basis and show evidences that the most suitable resolution for document digitization are based on the principles of the Nyquist-Shannon theorem from Information Theory and ensure the accurate preservation of all necessary document information. Thus, this dissertation expands the scope of the study published in Lins et al. (2024), incorporating new experiments, additional images and detailed analyses.

Such target is reached by:

- Enlarging and evaluating the datasets of document used considering a diversity of detailed images, including those with tiny fonts and high-frequency components, at multiple resolutions to represent a much wider range of document types in the experiments performed.

- Investigating the effects of the binarization process in enhancing the recognition of text and fine details at different resolutions.

- Using OCR tools to evaluate the impact of resolution on recognition rates, using the Levenshtein Distance to analyze the effects of scanner resolution in the accuracy of the transcription of the texts.

- Experimentally validating the Nyquist-Shannon theorem as a basis for determining optimal scanner resolution, considering several practical implications.

## 2 THEORETICAL BACKGROUND

According to Gonzalez and Woods (2018), in the context of digital image processing, an image can be described as a two-dimensional function, *f(x,y)*, where *x* and *y* represent spatial coordinates of a plane and the amplitude of *f* is referred to as the intensity or gray level of that specific point. If *x, y* and the intensity values of f are all finite and discrete, the image is considered a digital image. The value of *f* at spatial coordinates *(x, y)* is regarded as a scalar quantity physically determined by the image source, with values being proportional to the energy of the image. Most images are acquired through the interaction between the lighting source falling on the scene and the reflection or absorption of energy carried out by the elements within the photographed scene. Thus, the function *f(x,y)* can be characterized as the outcome of the combination of the illumination and reflectance elements, defined respectively by *i(x,y)* and *r(x,y)*:

$$f(x,y) = i(x,y)r(x,y) \tag{2.1}$$

where

$$0 \leq i(x,y) < \infty \tag{2.2}$$

and

$$0 \leq r(x,y) \leq 1 \tag{2.3}$$

Gonzalez and Woods (2018) asserts that reflectance is determined by the properties of the imaged objects and constrained between 0 (indicating total absorption) and 1 (indicating total reflectance). In instances of images formed through transmission, like a chest X-ray, the function would rely on the transmissivity function rather than the reflectance function. However, the boundaries remain consistent, and the image function remains represented as the product in equation 2.1.

Reference Sonka, Hlavac and Boyle (2015) states that computerized image processing relies on digital image functions. Therefore, an image must be represented using suitable discrete data, such as a matrix. In Figure 1, there is a representation of a digital image acquisition process using array sensors. In such an example, the imaging system collects the incoming energy and focuses it onto an image plane. Since the array sensor coincides with the focal plane, the resulting outputs will be proportional to each sensor's integral of the light received.

These outputs are then converted to an analog signal by digital and analog circuitry, which is then digitized, resulting in a digital image.

Figure 1 – Representation of an digital image acquisition.



**Source:** GONZALEZ; WOODS (2018)

Considering the diversity of potential applications for digital image processing, Gonzalez and Woods (2018) highlights that images can be categorized by source, such as X-ray, visual, and infrared. The most prevalent energy source utilized today is the electromagnetic spectrum. This spectrum is grounded in the revelation that the resultant light beam is not white but comprises a continuous spectrum of colors extending from violet to red. Bushong (2013) describes that when sunlight passes through a prism, it becomes clear that white light consists of photons spanning a spectrum of wavelengths, and the prism functions to disperse and organize the emitted light into colors by refracting various wavelengths through distinct angles. Visible light comprises only a portion of the electromagnetic spectrum, extending over 25 orders of magnitude. Figure 2 illustrates the values of energy, frequency, and wavelength and identifies the three imaging windows of X-ray, visual, and MR imaging according to their associated values.

Figure 2 – Representation of different dimensions of pairs of lines.

## 2.1  IMAGE RESOLUTION

Regarding the quality of a digital image, Sonka, Hlavac and Boyle (2015) emphasizes that it is directly proportional to its resolution. The author defines that the resolution of an image can be categorized into four main types: spatial resolution, spectral resolution, radiometric resolution and temporal resolution. This work will focus more specifically on spatial resolution and its impact on document scanning.

### 2.1.1  Spatial resolution

Bushberg (2012) states that spatial resolution refers to the ability of an imaging system to identify small details, that is, to distinguish small nearby elements that have different contrasts in an image. Sonka, Hlavac and Boyle (2015), highlights that the distance of the samples in the image plane determines this aspect. Therefore, the higher the spatial resolution, the better the visualization of details in an image.

Gonzalez and Woods (2018) add that, when considering the measurement of the smallest

detail discernible in the image, the resolution in space is usually quantified in terms of line pairs per unit distance or dots (pixels) per unit distance. According to Bushong (2013), a line pair is defined by a high-contrast line separated by a space of width equal to that of the line. For radiological images, for example, it is common to use the unit line pairs per millimeter (lp/mm), referring to the lp/mm present in the image composition. According to Gonzalez and Woods (2018), one point per unit of distance refers to the size of the pixel present in the image composition. This measurement is generally used for printing files and is usually represented by the number of pixels in an inch, using the unit of dots per inch (dpi).

In Figures 3 and 4, respectively, representations of dimensions of pairs of lines and points (pixels) that can be applied concerning distance to determine spatial resolution are illustrated.

Figure 3 – Representation of different dimensions of pairs of lines.



**Source:** BUSHONG (2013)

Figure 4 – Representation of different sizes of objects formed by points (pixels).



**Source:** BUSHONG (2013)

In figure 5 it is possible to see an example of the impact generated in an image at different spatial resolutions, considering the image at resolutions of 930 dpi, 300 dpi, 150 dpi, and 72 dpi. Gonzalez and Woods (2018), highlights that the difference in spatial resolution also affects the size of the image. However, in terms of comparing the impact of spatial resolution, the images were zoomed in for better visualization.

Figure 5 – Variation of spatial resolution in an image.



(a) 930 dpi  (b) 300 dpi

(c) 150 dpi  (d) 72 dpi

**Source:** GONZALEZ; WOODS (2018)

## 2.2 HISTORICAL DOCUMENTS IMAGE QUALITY

According to Kenney and Chapman (2001), determining image quality requirements depends on the characteristics of the type of document to be scanned. Capturing a finely detailed

historical document written in pen and ink, for example, requires a higher resolution to correctly identify its elements. The author defines that, in general, documents can be classified into four identification categories as illustrated in table 1.

Table 1 – Documents Category

| Document Category | Details |
|---|---|
| TEXTS/DASH | Can be produced manually, mechanographed or by machine. Usually in black and white. Includes books, manuscripts, newspapers, reports, typed or laser-printed documents, architectural plans, maps, line drawings, etchings, lithographs, and musical scores. |
| HALF TONE | In color or black and white. Reproductions, usually created from a photograph, made up of small dots or squares or lines, which are used to represent continuous tones. Most photographs in publications are halftones. |
| CONTINUOUS TONE | Color or black and white. Includes graphic illustrations in which all gray and color values can be reproduced: photographs, crayons, chalk and some pencil drawings, acrylics, watercolors and photographically reproduced facsimiles. |
| MIXED | Color or black and white. Refers to items containing halftone text and images or continuous tone text and images, such as newspapers, magazines, picture books, programs, and sheet music covers. Does not include text and line drawings together. |

**Source:** Prepared by the author based on (KENNEY; CHAPMAN, 2001), 2024

Kenney and Chapman (2001) emphasizes, that in addition to document categories, other important factors to be considered in the scanning process are: level of detail, type of paper, level of contrast between ink and paper, sharpness, production process (machine-made or hand-made) and quality of the document, considering if it is damaged, stained or incomplete. The quality of the digital image is inherently dependent on the quality of the source. For example, scanning a blurred image will result in a blurred digital copy, regardless of the resolution. In addition, the choice of hardware and software used in the digitization process can also influence the final output.

### 2.2.1 Historical documents resolution

According to Nagy and Seth (1984), the resolution applied by the scanner in the document scanning process is defined by the level of detail required in the image. When considering a

two-level image (black and white), for example, 400 lines per inch would be considered an adequate value for spatial resolution.

The resolution for scanning documents for printing, according to Gonzalez and Woods (2018), is defined in terms of dpi. Table 2 contains different categories of documents followed by their spatial resolution information for printing provided by Gonzalez and Woods (2018) and Kenney and Chapman (2001).

Table 2 – Documents and their spatial resolutions. (**S.R.1** - Spatial Resolution provided by Gonzalez and Woods (2018), **S.R.2** - Spatial Resolution provided by Kenney and Chapman (2001))

| Document | S. R. 1 | S. R. 2 |
|---|---|---|
| **Newspapers** | $\approx 75$ | $\approx 120$ or $300$ |
| **Magazines** | $\approx 133$ | $-$ |
| **Flyers** | $\approx 175$ | $\approx 600$ |
| **Books** | $\approx 2400$ | $\approx 120$ or $300$ |
| **Monographs** | $-$ | $\approx 600$ |
| **Periodicals** | $-$ | $\approx 600$ |
| **Office documents** | $-$ | $\approx 300$ |
| **Tecnical draws** | $-$ | $\approx 600$ |
| **Maps** | $-$ | $\approx 600$ |
| **Handwritten materials** | $-$ | $\approx 300$ |

**Source:** Prepared by the author based on Gonzalez and Woods (2018) and Kenney and Chapman (2001), 2024

## 2.3   CONTINUOUS-TIME AND DISCRETE-TIME SIGNALS

Gonzalez and Woods (2018) explains that a continuous signal must be converted into a sequence of discrete values to be processed by a computer, a process that involves both sampling and quantization. Lathi (2005) emphasizes that continuous-time signals can be effectively analyzed by processing their discrete samples using discrete-time systems. To ensure accurate reconstruction of the original signal from its samples, the sampling rate must be sufficiently high, minimizing reconstruction errors. The sampling theorem provides the theoretical foundation for this process, defining the conditions under which continuous-time signals can be reliably sampled and reconstructed from their discrete sampled values.

Lathi (2005) further describes sampling theory as the connection between continuous-time and discrete-time representations. A continuous-time signal is expressed as a sequence of pulses and is specified in a continuum of time values $t$, while a discrete-time signal represents the same

information numerically, as discrete values of *t*. Thus, both formats represent the same data, ensuring that each sample's information is equivalent, regardless of the representation. Some examples of continuous-time signals include outputs from telephones and analog cameras, while examples of discrete-time signals include monthly sales figures of a corporation and daily stock market averages. Figure 6(a) illustrates a continuous function $f(t)$ sampled at regular intervals $\Delta T$, where each sample corresponds to a specific point on the continuous function. In contrast, Figure 6(b) represents these samples in the discrete $x$-domain, where the samples are indexed sequentially by integers (0, 1, 2, and 3) instead of their exact positions in time. This highlights that the sampled values at these points remain the same, ensuring information equivalence between the two domains.

Figure 6 – (a) *Left figure:* Continuous function $f(t)$ sampled $\Delta T$ units apart.
(b) *Right figure:* Samples in the $x$-domain, where $x$ is discrete.



**Source:** GONZALEZ; WOODS (2018)

## 2.4 THE NYQUIST-SHANNON THEOREM

According to Lathi (2005), the basis of the Nyquist-Shannon Theorem is the Sampling Theorem, originally known as the "Interpolation Formula". Siebert (1986) describes that the general formulation of the "Sampling Theorem" was credited by H. S. Black to Cauchy in 1841 and around the 1920s, Nyquist, Carson, and Hartley revisited and refined the established concepts, laying the foundation for modern communication theory.

In 1948, the work of C. E. Shannon further formalized and popularized the theorem, Shannon (1949) even points out in his work that the theorem had already been portrayed previously by mathematicians in other forms, but had not appeared explicitly in the communication theory literature until Nyquist returned to work on the essential idea of the theorem and stated

that, considering a $f(t)$ limited in bandwidth to a maximum frequency $W$ and limited to the time interval $T$, with samples taken at intervals of $\frac{1}{2W}$ seconds apart, a total of $2TW$ samples will fully capture the information within the interval. Therefore, approximately $2TW$ numbers are enough to specify the function, based on a Fourier series expansion over the time interval $T$.

A band-limited function is defined as a function $f(t)$ whose Fourier transform is zero for frequency values outside a finite interval $[-\mu_{max}, \mu_{max}]$, as described by Gonzalez and Woods (2018). This restriction implies that, for a function to be confined within a specific time interval $T$, all frequency values outside this range must be zero, as stated by Shannon (1949).

Shannon articulated the theorem as follows:

> If a function of time $f(t)$ is limited to the band from $0$ to $W$ cycles per second, it is completely determined by giving its ordinates at a series of discrete points spaced $\frac{1}{2W}$ seconds apart (SHANNON, 1948, p. 34).

According to Lathi (2005), the samples can be taken arbitrarily at any instants, resulting in different intervals, as long as the sampling instants are recorded and on average of $2B$ samples per second. Gonzalez and Woods (2018) states that the term *Nyquist rate* defines a sampling rate equal to twice the highest frequency.

Based on Lee and Varaiya (2011), figure 7 illustrates an example of a signal sampling and reconstruction process based on The Nyquist-Shannon Theory in terms of its basis. Each line of pair graphs in the figure represents a step of the process as follows:

- **Step 1:** The left image represents the continuous-time original signal $x(t)$, and the right image shows its continuous Fourier transform $X(w)$, where the data is in terms of the frequency and band-limited to $[-\pi/T, \pi/T]$.

- **Step 2:** The signal $x(t)$ is sampled at intervals of $T$, resulting in the discrete signal $y(n)$. On the right, the discrete-time Fourier transform (DTFT) $Y(w)$ is shown, where the spectrum is periodically replicated with a period of $2\pi/T$.

- **Step 3:** The discrete-time signal $y(n)$ is converted into a train of impulses $w(t)$ by the Impulse Generator. In the frequency domain, $W(w)$ represents the scaled and periodic spectrum $1/T$.

- **Step 4:** The impulse train $w(t)$ is filtered by the $Sinc_T$ function, which is a scaled version of the sinus cardinalis function, defined as:

$$\text{sinc}_T(t) = \frac{\sin(\pi t/T)}{\pi t/T} \tag{2.4}$$

where:

- $t =$ time

- $T =$ the sampling period

This function acts as an interpolation kernel, ensuring that each sampled point contributes correctly to the reconstruction of the original signal. In frequency domain terms, $Sinc_T$, performs low-pass filtering, removing spectral replicas and preserving the original spectrum. Then, the inverse Fourier transform reconstructs the continuous-time signal $z(t)$, shown on the left, which matches the original signal $x(t)$.

Figure 7 – Steps in the Nyquist-Shannon sampling theorem.



**Source:** LEE; VARAIYA (2011)

### 2.4.1  The Nyquist-Shannon Theorem applied to Images

In terms of a digital image, Sonka, Hlavac and Boyle (2015) describes that a continuous function $f(x, y)$ can be sampled by applying a discrete grid of sampling points in the plane, or the image function can also be expanded by applying an orthonormal function as a basis, such as the Fourier transform, resulting in coefficients of the expansion that represent the digitized image.

Colarusso et al. (1999) points out that the Nyquist-Shannon theorem is essential to direct imaging applications, as the signal is sampled through discrete pixel elements in an array. Sonka, Hlavac and Boyle (2015) also emphasizes that, there is a connection between the digital sampling density and the details that an image will contain. The Nyquist-Shannon principle implies that the sampling rate should be at least twice the highest spatial frequency to ensure accuracy. For example, Colarusso et al. (1999) describes that considering a diffraction-limited arrangement, the product of the lateral spatial resolution and the total magnification must be twice the pixel size of the image, as defined by Equation 2.5:

$$M_{tot} R_L \geq 2p \qquad (2.5)$$

where:

- $M_{tot}$:total magnification

- $R_L$: lateral spatial resolution

- $p$: pixel size

Adhering to the Nyquist rate helps maintain the spatial fidelity of images and prevents sampling artifacts. In this context, oversampling does not provide additional information, resulting in a scenario commonly referred to as "empty magnification."Conversely, undersampling below the Nyquist-Shannon limit leads to distortions or artifacts, such as Moiré patterns, a phenomenon known as spatial aliasing Gonzalez and Woods (2018). Sonka, Hlavac and Boyle (2015) describes that aliasing happens when there is an overlapping of the periodically repeated results of the Fourier transform $f(u, v)$ of an image with band-limited spectrum, a result that can be prevented by choosing a sampling interval that it is less than half of the smallest interesting detail in the image.

Colarusso et al. (1999) explains that the resolution of an imaging device is primarily determined by its maximum capability. In terms of document scanning, Barney-Smith (1998) points out that the bitmap resulting from a two-level scan depends greatly on the parameters of the scanner used. The most common and simplest way to digitize a document in companies, or even at home, is through a flatbed scanner. Federal Agencies Digital Guidelines Initiative (2016) states that these models are easy to use, versatile and can achieve different resolution ranges, with values from 75 to 1200 dpi depending on the model used. Therefore, selecting the appropriate scanner resolution is essential to capture all document details effectively while avoiding the drawbacks of prolonged image acquisition and processing times, excessive storage requirements, and unnecessary network bandwidth usage.

## 2.5 EVALUATION MEASURE

### 2.5.1 Levenshtein Distance

Defined by Levenshtein (1966) in "*Binary codes capable of correcting deletions, insertions and reversals*", the Levenshtein Distance (LD) refers to the minimal number of insertions, deletions and replacements needed to to make two strings identical. The metric was proposed in the context of the code theory to solve problems in detecting and correcting transmission data. According to Navarro (2001) the Levenshtein distance is also known as *Edit distance* and, in a simplified definition, all the operations involved to guarantee the string matching are assigned as a cost of 1. The LD is widely used in applications such as evaluating spell-checkers, OCR systems, and DNA sequence analysis.

Shuwandy et al. (2020) defines that mathematically, the Levenshtein distance between two strings $a, b$ with respectively length of $|a|$ and $|b|$ is given by $lev_{a,b}$ and detailed in Equation 2.6. The $lev_{a,b}(i, j)$ is the distance between the first $i$ characters of $a$ and the first $j$ characters of $b$.

$$lev_{a,b}(i,j) = \begin{cases} \max(i,j) \text{ if } \min(i,j) = 0, \\ \\ \min \begin{cases} lev_{a,b}(i-1,j) + 1, & \text{(deletion)}, \\ \\ lev_{a,b}(i,j-1) + 1, & \text{(insertion)}, \\ \\ lev_{a,b}(i-1,j-1) + 1_{(a_i \neq b_j)}, & \text{(match or mismatch)}. \end{cases} \end{cases} \quad (2.6)$$

According to Shuwandy et al. (2020), the first element in the minimum refers to deletion, the second to insertion and the third refers to the match or mismatch of symbols, depending on their similarity. The $1_{(a_i \neq b_j)}$ is the indicator function in the match or mismatch elements, where:

$$\begin{cases} 1_{(a_i \neq b_j)} = 0 \text{ if } a_i = b_j \\ 1_{(a_i \neq b_j)} = 1 \text{ if } a_i \neq b_j \end{cases} \quad (2.7)$$

Shuwandy et al. (2020) reports an LD example between the words 'GEEI' and 'GEELY' to evaluate their similarity, considering 'GEELY' as the password of a system. Board 1 shows that the LD is 2, because of the two last characters.

Board 1 – Example of Levenshtein distance metric

| G | E | E | I | |
|---|---|---|---|---|
| G | E | E | L | Y |

**Source:** SHUWANDY et al. (2020)

# 3 MATERIALS AND METHODS

A distinct set of diversified files, representing the scope of document engineering, was selected to establish the optimal scanning resolution that balances image quality, processing time and storage space in document digitization. This selection was referred to as the "Test dataset"and is broadly divided into two main categories: text-only documents ("Text") and documents containing both text and images ("General"). The test dataset includes twenty-eight distinct documents scanned at resolutions: 75, 150, 200, 300, 500, 700, and 1200 dpi. These documents are book pages and single paper files in different papers, font styles, sizes, formats and colors. The content includes:

- Typed Texts, in English, Portuguese and Spanish;

- Handwritten Texts in Portuguese;

- Maps;

- Architectural Plan,

- Images and Illustrations.

In figure 8 it is possible to observe samples of images from the text-only dataset followed by samples of images from the general dataset in figure 9.

Figure 8 – Samples of images from the Text-only Dataset



**Source:** The author (2024)

Figure 9 – Samples of images from the General Dataset.



**Source:** The author (2024)

The eighteen images of English texts were classified into the "Text Dataset", while the remaining ten images were organized into a separate group referred to as the "General Dataset". The flowchart in Figure 10 presents the process carried out to prepare each data set for the Fourier Transform application.

Figure 10 – Datasets Workflow



**Source:** The author (2024)

As mentioned in the previous section, the Nyquist-Shannon theorem applies to continuous signals that are bandwidth-limited. It is a principle in Fourier analysis that sharper signals (those with higher derivatives in the time or spatial domain) exhibit higher frequency components. When signals display "discontinuities" in time or space, a phenomenon known as "frequency overshooting" occurs, resulting in intense oscillations in the frequency domain. By closely examining a color document image, one observes that the pixel colors at boundaries transition smoothly, whereas binary images have sharp edges, which can introduce perceived

signal discontinuities. The selected documents were examined in true-color and binary formats, to facilitate a more thorough frequency domain analysis.

To determine the optimal scanning resolution, the approach here involved analyzing the spectral density of various resolutions of the same image to locate the "maximum frequency". The smallest resolution to yield such a frequency is considered the most efficient choice, balancing scanning time, storage requirements, network transmission bandwidth, and processing time, as processing time typically scales with file size. This study also evaluates the effect of different scanning resolutions on the transcription quality of text-only documents by comparing the Levenshtein distances between transcribed texts and their ground-truth counterparts.

## 3.1 SCANNING PROCESS

In this study, the first step in determining the optimal resolution for scanned documents was to digitize each document at 75, 150, 200, 300, 500, 700, and 1200 dpi. The files were scanned in color mode and saved in Portable Networks Graphics (PNG) format, a widely used raster format known for its efficient lossless compression and high bit-depth compatibility, as described by Puglia et al. (2004).

The scanning process was performed using three scanners: Epson EcoTank L4160, Epson EcoTank L8180, and HP LaserJet Pro Color M479FDW. The technical specifications of each device, as provided in their respective manuals, are summarized in Table 3. Based on these specifications, the first two scanners were capable of digitizing documents at all seven predefined resolutions. However, the HP LaserJet Pro Color M479FDW scanner was limited to the following resolutions: 75, 200, 300, and 1200 dpi.

Table 3 – Scanners specifications.

| Specification | Epson L4160 | Epson L8180 | HP M479FDW |
|---|---|---|---|
| Scanner type | Flatbed | Flatbed | Flatbed |
| Photoelectric device | CIS | CIS | CIS |
| Effective pixels | 10200 × 14040 pixels (1200 dpi) | - | - |
| Maximum document size | US letter or A4 (8.5 × 11.7 inches (in) or 216 × 297 millimeters (mm)) | Legal (8.5 × 14 in or 216 × 356 mm) | Legal (8.5 × 14 in or 216 × 356 mm) |
| Optical resolution | 1200 x 2.400 dpi | 1.200 x 4.800 dpi | 1200 x 1200 dpi |
| Output resolution | 50 to 9600 dpi in 1 dpi increments | - | - |
| Image data | **Color:** 48 bits per pixel (bpp) (internal), 24 bpp (external). **Grayscale:** 16 bpp (int.), 8 bpp (ext.). **Black and White:** 16 bpp (int.), 1 bpp (ext.) | **Color:** 48 bpp (internal), 24 bpp (external). **Grayscale:** 16 bpp (int.), 8 bpp (ext.). **Black and White:** 16 bpp (int.), 1 bpp (ext.) | - |
| Light source | LED | LED | LED |

**Source:** Prepared by the author based on EPSON (2019), EPSON (2021) and HP (2019), 2025.

## 3.2 CROPPING PROCESS

The Test Dataset was processed using a cropping algorithm, presented in Appendix A, to remove visible scanner edges and ensure a clear focus on the document page.

The Source Code 1 begins by selecting the image to be processed and defining its path. After verifying the existence of the image according to the specified path, the function

"cv2.imread(path)" loads the found images and displays an error message if the image is not located. The cropping parameters are manually and individually adjusted for each image, considering its specific characteristics. The parameters 'top', 'bottom', 'left', and 'right' define the cropping percentages for the top, bottom, left, and right edges of the image, respectively. These values are stored in the dictionary "cut_values", ensuring a clear mapping between cropping directions and their respective values.

The cropping function cut_and_save calculates the number of pixels to be cropped based on the defined percentages and performs the cropping operation using image[top:bottom, left:right]. This function takes the following parameters:

- image: the image to be cropped;

- output_path: the path to save the cropped image;

- resolution: the resolution of the image being processed;

- cut_percent: the cropping percentages.

It is important to highlight that the cropping percentages were applied across all different image resolutions, in line 43 of the Source Code 1, it is established as follows:

```
top, bottom, left, right = [int(p * max(image.shape) / 100) for p in
    cut_percent]
```

The line defines that, the cropping values are calculated as a percentage of the maximum image dimension (max(image.shape)). Thus, the cropping process is proportional to the image size, adapting to different resolutions. This approach helps maintain relative consistency in the visual content regardless of varying resolutions.

At the end of the process, the new cropped images are saved in PNG format and the applied cropping parameters are saved in an Excel file. An example of the cropping result for Image 01 across different resolutions is shown in Figure 11.

Figure 11 – Example of the different resolutions of Image 01 in the original scan size and after the cropping process.



**Source:** The author (2025)

## 3.3 BINARIZATION PROCESS

After the cropping process, the images were subjected to binarization. According to Lins et al. (2021), the binarization of document images remains a significant research area, particularly due to its importance in processing stages of document analysis and optical character recognition (OCR). In this study, the Otsu binarization algorithm Otsu (1979) was applied to the test images due to its accessibility, ease of implementation, and consistent image quality. The implementation of this method can be found in Appendix B.

Initially, the Source Code 2 defines the specific image to be processed and sets its corresponding path. Subsequently, a validation step is performed to verify if the image exists. If the image is successfully found, it is loaded using the function `cv2.imread(path)` from the OpenCV library. In the event of an invalid path or missing image, an error message is displayed.

The main function responsible for the binarization process is `binarize_and_save`. This function begins by converting the original image to grayscale using the OpenCV function `cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)`. This conversion is essential because Otsu's thresholding technique operates more effectively on grayscale images, where pixel intensities range from 0 (black) to 255 (white). According to the OpenCV documentation Doxygen

(2024), the use of the conversion code cv2.COLOR_BGR2GRAY ensures a reduction in compu-
tational complexity, facilitating the accurate application of the binarization technique.
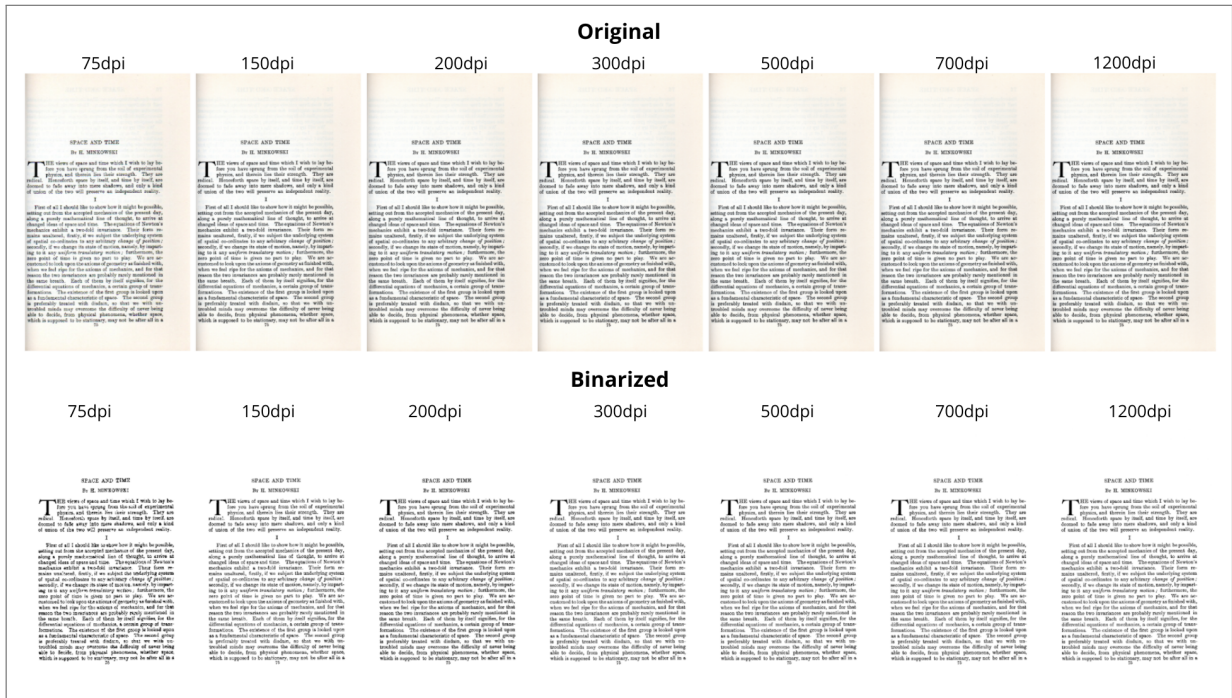
As defines by Doxygen (2024), the subsequent step in line 30 of the Source Code 2, applies
the Otsu thresholding method using the following OpenCV flags:

- cv2.THRESH_BINARY: Performs binary thresholding, where pixel values below the th-
  reshold are set to 0 (black) and values equal to or above the threshold are set to 255
  (white);

- cv2.THRESH_OTSU: Automatically calculates the optimal global threshold for the image,
  minimizing the intra-class variance between the foreground and background pixel inten-
  sities.

Otsu's method is effective with images that exhibit a bimodal histogram, that presents two
distinct peaks corresponding to the foreground and background, by automatically selecting the
ideal threshold. It eliminates the subjectivity associated with manually determining a suitable
threshold value, thus reducing the likelihood of misclassification.

Finally, the binarized image is saved to the specified path using the cv2.imwrite(output_path,
binarized_image) function. This approach ensures that the resulting binarized images can
be utilized for further analytical steps. Figure 12 presents an example of the binarized result
for Image 01 in its different resolutions.

Figure 12 – Example of the different resolutions of Image 01 in the original scan color and after the binarization process.



**Source:** The author (2025)

## 3.4 SPECTRAL ANALYSIS

The images were prepared as described previously and then the Fourier Transform was applied to each one using the NumPy library in Python, a tool widely used for scientific computing. The source code of the Fourier transform algorithm applied in the dataset is in Appendix C

The Source Code 3 starts locating the defined images, followed by checking for their existence. The images are loaded in grayscale mode using "cv2.IMREAD_GRAYSCALE", considering that Red Green Blue (RGB) images require separate transforms for each channel, the use of grayscale images reduces complexity, improving computational efficiency while still preserving meaningful frequency data based on the different intensities.

For each image, the Fourier Transform was applied considering the following code, located in lines 75-77 of the Source Code 3

```
transform_image = np.fft.fft2(image)
center_transform_image = np.fft.fftshift(transform_image)
magnitude_spectrum = np.log(np.abs(center_transform_image) + 1)
```

where:

- `np.fft.fft2(image)`: computes the 2D Fourier Transform, performing the calculation of the two-dimensional discrete Fourier transform across specified axes within an M-dimensional matrix;

- `np.fft.fftshift()`: moves to center the zero-frequency component within the spectrum, providing better visualization;

- `np.log(np.abs(center_transform_image) + 1)`: converts the frequency domain to a log scale for improved contrast, calculating the logarithm of the magnitude of the spectrum.

The function `calculate_max_frequency` finds the maximum frequency component in the transformed image, corresponding to the dominant structure in the image. The highest frequency was identified with the `numpy.argmax`, which finds the indices of maximum values and was used here to pinpoint the peak frequency within the magnitude spectrum as defined by Numpy (2024). With the magnitude spectrum results, it was possible to plot a frequency map and frequency range for each image resolution.

Following the same Fourier-based spectral analysis on the Source Code 3, the Source Code 4 shows the results in a Spectral Density Plot. The Spectral Density function $S(f)$ represents the distribution of energy across different frequency components and, as described by Oppenheim, Willsky and Nawab (1994), is defined as:

$$S(f) = |X(f)|^2 \tag{3.1}$$

where

- X(f) = Fourier transform of the signal x(t)

Considering that the signal has finite energy:

$$E = \int_{-\infty}^{\infty} |x(t)|^2 \, dt < \infty \tag{3.2}$$

To estimate the Spectral Density, the magnitude spectrum was calculated using the Fourier Transform, and then a normalized histogram of the frequency magnitudes was created. This histogram represents the relative occurrence of different frequencies in the image and the spectral analysis functions were defined as:

```python
def calculate_spectral_energy(magnitude_spectrum):
    return np.sum(magnitude_spectrum**2)
```

- `calculate_spectral_energy`: calculates the total spectral energy by summing the squares of the magnitude of the spectrum. This value represents the total energy of the signal in the analyzed frequency components;

```python
def calculate_max_frequency(magnitude_spectrum):
    row_center, col_center = np.unravel_index(np.argmax(magnitude_spectrum),
        magnitude_spectrum.shape)
    return magnitude_spectrum[row_center, col_center]
```

- `calculate_max_frequency`: finds the pixel with the highest frequency magnitude in the spectrum;
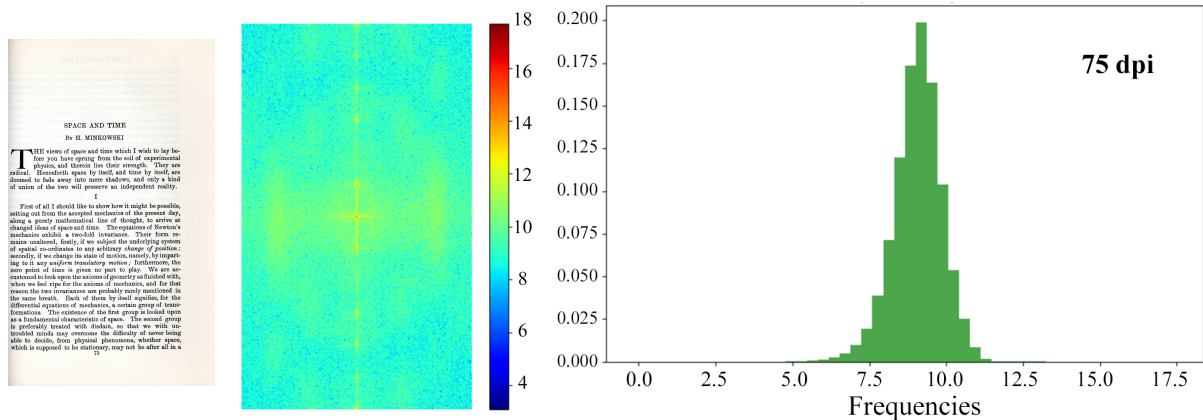
```python
def calculate_frequency_distribution(magnitude_spectrum):
    hist, bins = np.histogram(magnitude_spectrum.flatten(), bins=50, range
        =(0, magnitude_spectrum.max()))
    hist = hist / hist.sum()
    predominant_frequencies = bins[:-1][hist > 0.1]
    return hist, bins, predominant_frequencies
```

- `calculate_frequency_distribution`: creates a histogram of frequency magnitudes to identify the frequency distribution of the image. The relative frequency threshold above 0.1 was used, as defined by "bins[:-1][hist > 0.1]", which facilitated the generation of a frequency recurrence graph for the spectrum, normalizing the values so that the sum is 1 and considering values above 10% of occurrence according to Numpy (2024);

The frequencies in the resulting spectrum are based on the image resolutions and correspond to cycles per pixel. The Fourier Transform allowed the extraction of a frequency map, the identification of the frequency range, the determination of the highest frequency and the creation of a Spectral Density plot for each image, all using the Numpy library based on NumPy (2024) and subsequently visualized with the Matplotlib library as stated in the Matplotlib (2024) documentation.

Figures 13 to 19 show the frequency maps and ranges for Image 01 across different resolutions, accompanied by the Spectral Density graph. The maps are displayed using Matplotlib's jet colormap model described in Matplotlib (2024).

Figure 13 – Image 01, the frequency map and the Spectral Density of 75 dpi resolution.



**Source:** The author (2024)

Figure 14 – Image 01, the frequency map and the Spectral Density of 150 dpi resolution.



**Source:** The author (2024)

Figure 15 – Image 01, the frequency map and the Spectral Density of 200 dpi resolution.



**Source:** The author (2024)

Figure 16 – Image 01, the frequency map and the Spectral Density of 300 dpi resolution.



**Source:** The author (2024)

Figure 17 – Image 01, the frequency map and the Spectral Density of 500 dpi resolution.



**Source:** The author (2024)

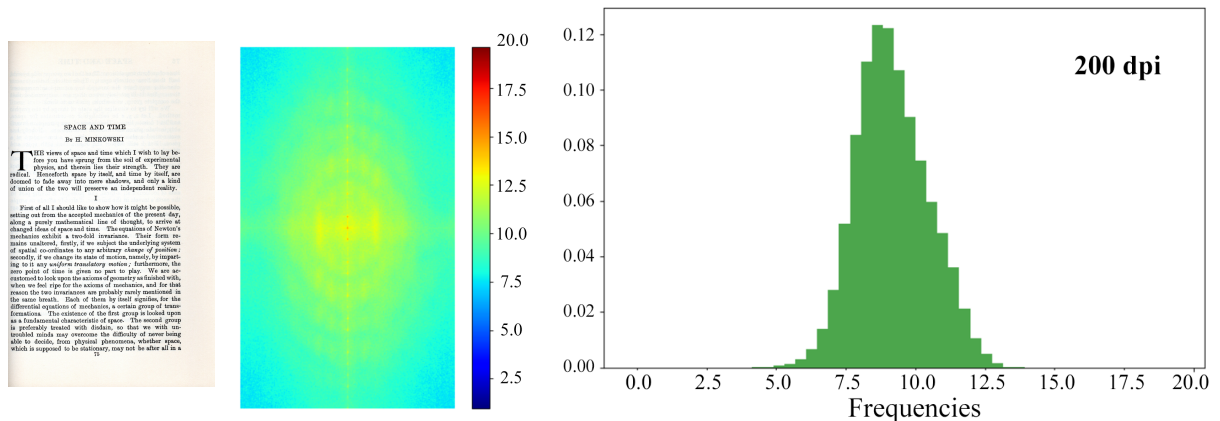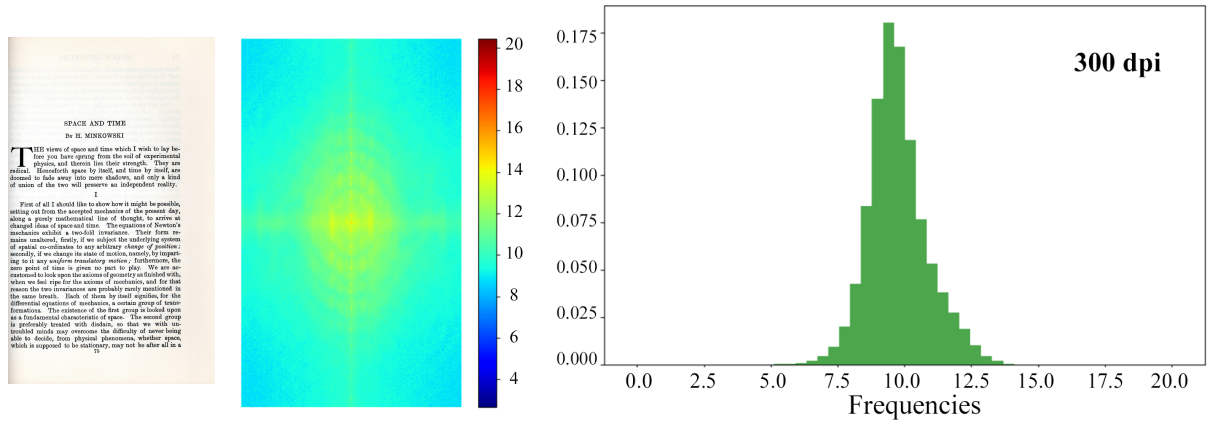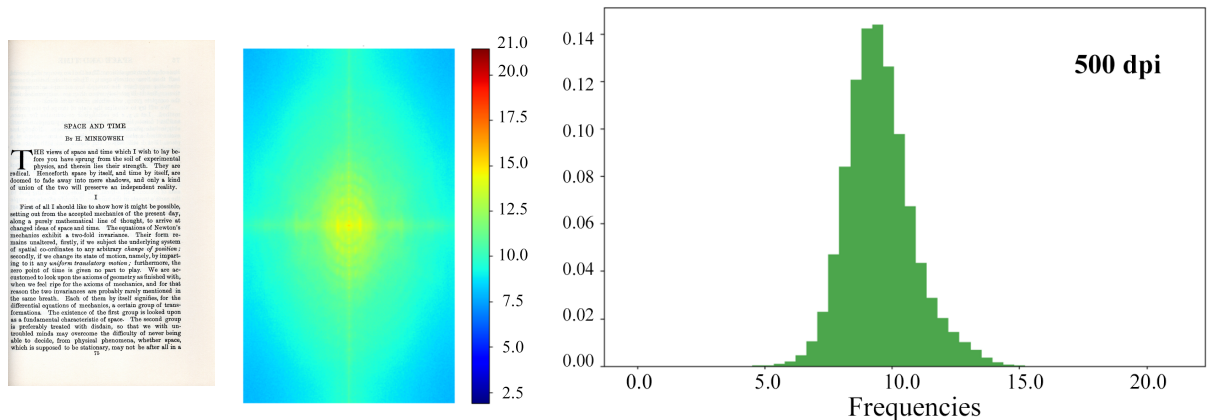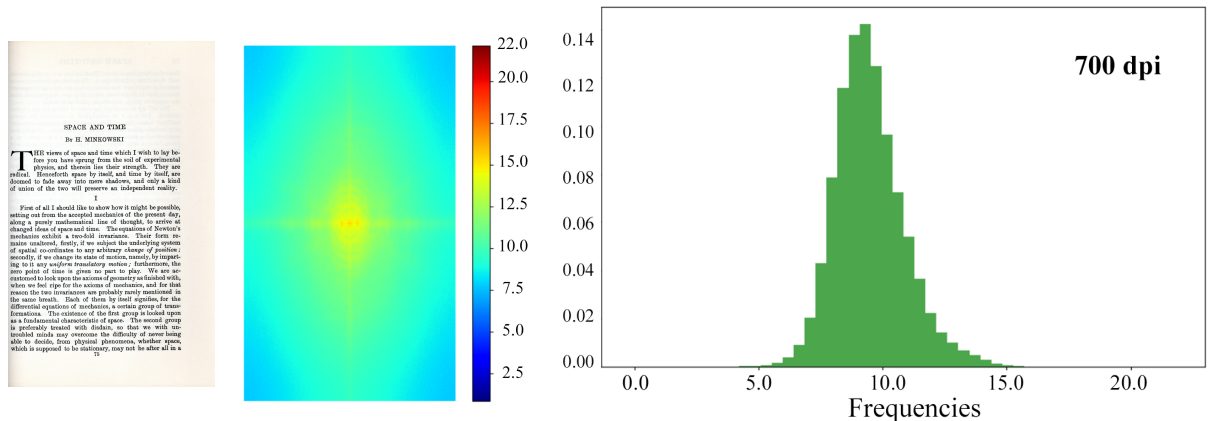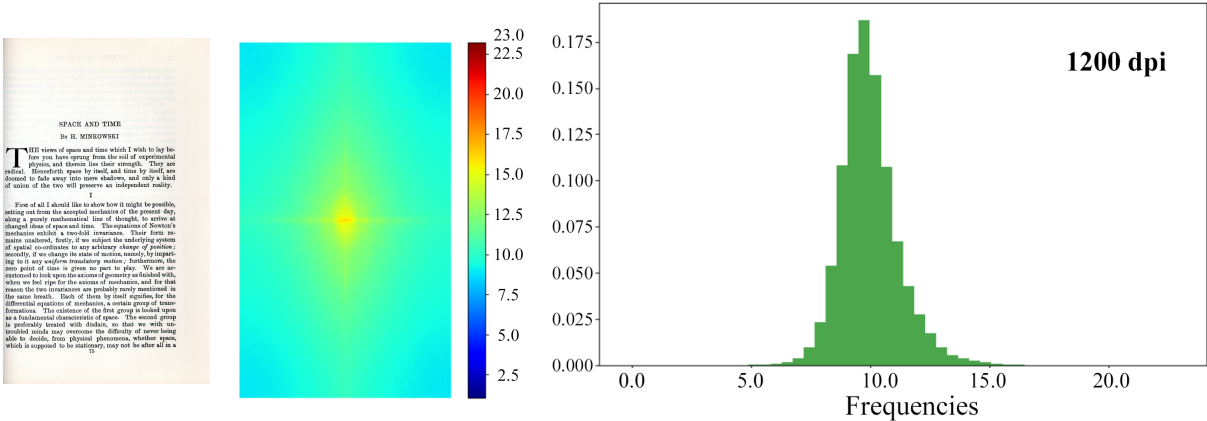Figure 18 – Image 01, the frequency map and the Spectral Density of 700 dpi resolution.



**Source:** The author (2024)

Figure 19 – Image 01, the frequency map and the Spectral Density of 1200 dpi resolution.



**Source:** The author (2024)

## 3.5 OPTICAL CHARACTER RECOGNITION (OCR)

According to Puglia et al. (2004), digital images must meet quality standards to ensure that OCR conversion achieves a certain level of accuracy. Barney-Smith and Andersen (2005) highlights that specific combinations of parameters governing scanning degradations can cause characters to appear either similar or significantly different when degraded. Beyond the influence of the optical resolution of the scanner on the document quality, factors such as contamination, manual underlining, or annotations can also pose challenges and interfere with the OCR process, as detailed by Forschungsgemeinschaft (2016).

Given the impact of document resolution on transcription accuracy, text-based documents were transcribed using optical character recognition (OCR), and the Levenshtein Distance was calculated for each resulting file. The implemented algorithm is presented in Appendix D.

The Source Code 5 begins by selecting the image to be processed and defining its path. After verifying the existence of the image according to the specified path, the function "cv2.imread(path)" loads the images and displays an error message if the image is not located. For the OCR process, Python-tesseract, a Python wrapper for the Google Tesseract-OCR Engine, was used, as referenced in PyPi (2024). Specifically, Tesseract OCR version 4.1.1 was applied and accessed by the pytesseract library (version 0.3.13) within the Python environment. The algorithm was executed without using the language dictionary or loading frequent words, this approach eliminates the influence of predefined words on OCRbased word recognition, prioritizing character-level identification for evaluating text element recognition. Based on doxygen (2015) and tesseract-ocr (2024), these settings were defined in line 60 of the Source Code 5, using the commands:

```
tesseract_config = --psm 6 load_system_dawg=0 load_freq_dawg=0]
```

where:

- – -psm 6: sets the Page Segmentation Mode (PSM) to assume a single uniform block of text;

- load_system_dawḡ0 load_freq_dawḡ0: disables the dictionary-based correction to allow a raw text extraction without auto-correction.

The OCR transcription of each document was saved in the correspondent TXT file format and the Python-Levenshtein library was employed by applying the code "`distance = levenshtein(reference_text, ocr_text)`" to calculate the Levenshtein distance between each transcribed text and the corresponding ground-truth text as described in PyPi (2024), the LD results were then saved in an Excel file and analyzed considering the OCR accuracy at the different image resolutions, based on measuring how many character changes (insertions, deletions, substitutions) were needed.

# 4  RESULTS

The data collected using the scanner Epson L4160 for the color and binary versions of the documents are summarized in four tables. Tables 4 and 5 contain data from the Text Dataset, while Tables 7 and 8 contain information from the General Dataset.

## 4.1   TEXT DATASET

Table 4 contains detailed information for the 18 images of the text dataset considering the colored images scanned with the Epson L4160 scanner. The information includes image size in centimeters (cm), text size in points (pt), resolution in dpi, file size for both PNG and Windows Bitmap (BMP) formats in megabytes (MB), the Highest frequency (HF) and Levenshtein distance measurements for each document across the seven previously defined resolution levels.

Table 4 – Text Dataset - Original images results (**Image size** in centimeters (cm), **Font size** in points (pt), **dpi** - image resolution in dots per inch, **png/bmp** - file size of the png/bmp image in megabytes (MB), **HF** - Highest frequency, **LD** - Levenshtein distance.)
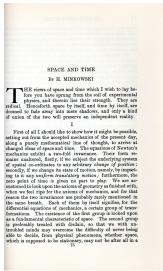
| Image | dpi | png | bmp | HF | LD | Image | png | bmp | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|
| | 75 | 0.35 | 0.59 | 18 | 97 | | 0.25 | 0.6 | 18 | 101 |
| | 150 | 1.28 | 2.37 | 19 | 12 | | 0.75 | 2.4 | 19 | 11 |
| | 200 | 2.28 | 4.21 | 20 | 22 | | 1.21 | 4.28 | 20 | 11 |
| | 300 | 5.12 | 9.47 | 20 | 21 | | 2.39 | 9.61 | 21 | 9 |
| | 500 | 13.74 | 26.27 | 21 | 22 | | 5.74 | 26.69 | 22 | 12 |
| | 700 | 25.82 | 51.49 | 22 | 24 | | 9.91 | 52.27 | 22 | 11 |
| 11.89 x 18.95 Font size: 10 | 1200 | 69.81 | 151.27 | 23 | 28 | 11.45 x 21.00 Font size: 10 | 25.18 | 153.59 | 23 | 12 |
| | 75 | 0.39 | 0.75 | 18 | 490 | | 0.32 | 0.48 | 17 | 945 |
| | 150 | 1.29 | 2.98 | 19 | 164 | | 1.18 | 1.92 | 19 | 912 |
| | 200 | 2.25 | 5.29 | 20 | 164 | | 1.99 | 3.41 | 19 | 909 |
| | 300 | 4.51 | 11.9 | 21 | 170 | | 4.76 | 7.66 | 20 | 912 |
| | 500 | 12.12 | 33.03 | 22 | 172 | | 12.2 | 21.26 | 21 | 914 |
| | 700 | 21.71 | 64.75 | 22 | 178 | 11.62 x 16.56 Font size: 12 | 23.02 | 41.67 | 22 | 916 |
| 14.19 x 21.00 Font size: 10 | 1200 | 54.9 | 190.23 | 23 | 180 | | 59.57 | 122.44 | 23 | 913 |

| Image | dpi | png | bmp | HF | LD | Image | png | bmp | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|
| 15.68 x 22.28 Font size: 12 | 75 | 0.42 | 0.87 | 18 | 28 | 12.94 x 21.00 Font size: 11 | 0.27 | 0.68 | 18 | 196 |
| | 150 | 1.44 | 3.49 | 19 | 14 | | 0.87 | 2.71 | 19 | 6 |
| | 200 | 2.54 | 6.2 | 20 | 14 | | 1.44 | 4.82 | 20 | 6 |
| | 300 | 5.47 | 13.94 | 21 | 18 | | 2.84 | 10.85 | 21 | 9 |
| | 500 | 14.39 | 38.69 | 22 | 19 | | 6.94 | 30.1 | 22 | 8 |
| | 700 | 26.06 | 75.8 | 23 | 20 | | 11.95 | 58.98 | 22 | 9 |
| | 1200 | 66.77 | 222.75 | 24 | 19 | | 30.5 | 173.31 | 23 | 12 |
| 11.04 x 21.00 Font size: 12 | 75 | 0.35 | 0.58 | 18 | 36 | 15.10 x 22.28 Font size: 10 | 0.56 | 0.84 | 18 | 3071 |
| | 150 | 1.24 | 2.31 | 19 | 7 | | 1.96 | 3.35 | 19 | 3110 |
| | 200 | 2.14 | 4.11 | 20 | 7 | | 3.29 | 5.97 | 20 | 3104 |
| | 300 | 4.49 | 9.26 | 20 | 7 | | 6.92 | 13.4 | 21 | 3205 |
| | 500 | 11.81 | 25.7 | 21 | 7 | | 17.68 | 37.24 | 22 | 3211 |
| | 700 | 21.68 | 50.35 | 22 | 7 | | 31.94 | 72.95 | 22 | 3201 |
| | 1200 | 56.13 | 147.98 | 23 | 12 | | 81.99 | 214.32 | 24 | 3205 |
| 16.29 x 23.47 Font size: 12 | 75 | 0.34 | 0.95 | 18 | 66 | 11.89 x 19.95 Font size: 12 | 0.42 | 0.59 | 18 | 59 |
| | 150 | 1.06 | 3.81 | 20 | 23 | | 1.43 | 2.37 | 19 | 1 |
| | 200 | 1.86 | 6.77 | 20 | 23 | | 2.46 | 4.21 | 20 | 1 |
| | 300 | 3.9 | 15.23 | 21 | 25 | | 5.36 | 9.47 | 20 | 1 |
| | 500 | 10.07 | 42.31 | 22 | 24 | | 14.11 | 26.27 | 21 | 1 |
| | 700 | 18.12 | 82.89 | 23 | 23 | | 25.97 | 51.49 | 22 | 2 |
| | 1200 | 46.59 | 243.54 | 24 | 27 | | 68.84 | 151.27 | 23 | 4 |
| 8.57 x 11.89 Font size: 7 | 75 | 0.15 | 0.25 | 17 | 1638 | 7.96 x 8.94 Font size: 7 | 0.05 | 0.18 | 17 | 442 |
| | 150 | 0.53 | 1.01 | 18 | 55 | | 0.18 | 0.71 | 18 | 48 |
| | 200 | 0.87 | 1.8 | 19 | 46 | | 0.33 | 1.26 | 19 | 47 |
| | 300 | 1.87 | 4.05 | 20 | 44 | | 0.74 | 2.83 | 19 | 39 |
| | 500 | 4.62 | 11.25 | 21 | 45 | | 2 | 7.85 | 20 | 42 |
| | 700 | 8.41 | 22.04 | 21 | 43 | | 3.85 | 15.39 | 21 | 39 |
| | 1200 | 22.18 | 64.78 | 22 | 45 | | 10.33 | 45.2 | 22 | 44 |
| 4.71 x 4.47 Font size: 5 | 75 | 0.03 | 0.05 | 15 | 210 | 5.01 x 5.35 Font size: 5 | 0.03 | 0.07 | 15 | 279 |
| | 150 | 0.11 | 0.21 | 16 | 56 | | 0.12 | 0.27 | 17 | 71 |
| | 200 | 0.18 | 0.37 | 17 | 52 | | 0.22 | 0.47 | 17 | 62 |
| | 300 | 0.43 | 0.84 | 18 | 39 | | 0.5 | 1.07 | 18 | 32 |
| | 500 | 1.12 | 2.31 | 19 | 20 | | 1.3 | 2.95 | 19 | 25 |
| | 700 | 2.13 | 4.53 | 20 | 29 | | 2.5 | 5.79 | 20 | 37 |
| | 1200 | 5.69 | 13.32 | 21 | 49 | | 6.66 | 17 | 21 | 54 |

| Image | dpi | png | bmp | HF | LD | Image | png | bmp | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|
| 4.71 x 5.38 Font size: 5 | 75 | 0.03 | 0.06 | 15 | 279 | 5.59 x 4.78 Font size: 4 | 0.03 | 0.07 | 16 | 786 |
|  | 150 | 0.12 | 0.25 | 17 | 71 |  | 0.12 | 0.27 | 17 | 282 |
|  | 200 | 0.21 | 0.45 | 17 | 62 |  | 0.2 | 0.47 | 17 | 41 |
|  | 300 | 0.5 | 1 | 18 | 32 |  | 0.44 | 1.06 | 18 | 24 |
|  | 500 | 1.3 | 2.78 | 19 | 25 |  | 1.11 | 2.94 | 19 | 20 |
|  | 700 | 2.45 | 5.44 | 20 | 37 |  | 2.04 | 5.76 | 20 | 21 |
|  | 1200 | 6.51 | 15.99 | 21 | 54 |  | 5.33 | 16.91 | 21 | 22 |
| 8.57 x 8.64 Font size: 5 | 75 | 0.1 | 0.18 | 17 | 2586 | 18.05 x 25.26 Font size: 8 | 0.87 | 1.14 | 18 | 4049 |
|  | 150 | 0.4 | 0.74 | 18 | 1906 |  | 3.04 | 4.54 | 20 | 4152 |
|  | 200 | 0.66 | 1.31 | 19 | 1960 |  | 5.12 | 8.09 | 20 | 4158 |
|  | 300 | 1.44 | 2.94 | 19 | 1978 |  | 10.94 | 18.18 | 21 | 4153 |
|  | 500 | 3.57 | 8.16 | 20 | 1982 |  | 27.47 | 50.5 | 22 | 4155 |
|  | 700 | 6.43 | 15.99 | 21 | 1976 |  | 50.1 | 98.93 | 23 | 4158 |
|  | 1200 | 16.94 | 46.97 | 22 | 1978 |  | 130.4 | 290.75 | 24 | 4169 |

**Source:** The author (2024)

Table 5 shows data for the 18 images in the Text Dataset considering the binarized images scanned with the Epson L4160 scanner. The information includes image size in cm, text size in pt, resolution in dpi, file size for PNG, BMP and Tagged Image Format (TIFF) formats in MB, the highest frequency and Levenshtein distance measurements for each document across the seven previously defined resolution levels.

Table 5 – Text Dataset - Binarized images results (**Image size** in centimeters (cm), **Font size** in points (pt), **dpi** - image resolution in dots per inch, **png/bmp/tiff** - file size of the png/bmp/tiff image in megabytes (MB), **HF** - Highest frequency, **LD** - Levenshtein distance).

| Image | dpi | png | bmp | tiff | HF | LD | Image | png | bmp | tiff | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 11.89 x 18.95 Font size: 10 | 75 | 0.01 | 0.2 | 0.01 | 18 | 229 | 11.45 x 21.00 Font size: 10 | 0.01 | 0.2 | 0.01 | 18 | 264 |
|  | 150 | 0.04 | 0.79 | 0.02 | 19 | 16 |  | 0.04 | 0.8 | 0.02 | 19 | 12 |
|  | 200 | 0.06 | 1.4 | 0.02 | 20 | 12 |  | 0.06 | 1.43 | 0.02 | 20 | 14 |
|  | 300 | 0.12 | 3.16 | 0.03 | 20 | 10 |  | 0.11 | 3.21 | 0.03 | 21 | 13 |
|  | 500 | 0.25 | 8.77 | 0.06 | 21 | 12 |  | 0.24 | 8.9 | 0.05 | 22 | 12 |
|  | 700 | 0.38 | 17.16 | 0.08 | 22 | 14 |  | 0.36 | 17.42 | 0.08 | 22 | 11 |
|  | 1200 | 0.73 | 50.42 | 0.15 | 23 | 15 |  | 0.68 | 51.22 | 0.14 | 23 | 12 |

| Image | dpi | png | bmp | tiff | HF | LD | Image | png | bmp | tiff | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 14.19 x 21.00 Font size: 10 | 75 | 0.02 | 0.25 | 0.01 | 18 | 839 | 11.62 x 16.56 Font size: 12 | 0.01 | 0.16 | 0.01 | 17 | 983 |
| | 150 | 0.05 | 0.99 | 0.03 | 19 | 172 | | 0.03 | 0.64 | 0.01 | 19 | 912 |
| | 200 | 0.09 | 1.77 | 0.03 | 20 | 164 | | 0.04 | 1.14 | 0.01 | 19 | 914 |
| | 300 | 0.17 | 3.97 | 0.05 | 21 | 166 | | 0.08 | 2.55 | 0.02 | 20 | 911 |
| | 500 | 0.36 | 11.01 | 0.08 | 22 | 169 | | 0.16 | 7.09 | 0.04 | 21 | 912 |
| | 700 | 0.54 | 21.6 | 0.12 | 22 | 179 | | 0.25 | 13.89 | 0.05 | 22 | 917 |
| | 1200 | 1.02 | 63.41 | 0.22 | 23 | 178 | | 0.48 | 40.81 | 0.1 | 23 | 915 |
| 15.68 x 22.28 Font size: 12 | 75 | 0.02 | 0.29 | 0.01 | 18 | 66 | 12.94 x 21.00 Font size: 11 | 0.01 | 0.23 | 0.01 | 18 | 395 |
| | 150 | 0.05 | 1.17 | 0.02 | 19 | 14 | | 0.04 | 0.9 | 0.02 | 19 | 10 |
| | 200 | 0.08 | 2.07 | 0.03 | 20 | 14 | | 0.06 | 1.61 | 0.02 | 20 | 9 |
| | 300 | 0.16 | 4.65 | 0.04 | 21 | 21 | | 0.12 | 3.62 | 0.03 | 21 | 7 |
| | 500 | 0.31 | 12.9 | 0.07 | 22 | 19 | | 0.24 | 10.03 | 0.06 | 22 | 9 |
| | 700 | 0.47 | 25.28 | 0.1 | 23 | 20 | | 0.37 | 19.68 | 0.09 | 22 | 10 |
| | 1200 | 0.89 | 74.25 | 0.19 | 24 | 20 | | 0.71 | 57.77 | 0.16 | 23 | 14 |
| 11.04 x 21.00 Font size: 12 | 75 | 0.02 | 0.19 | 0.01 | 18 | 107 | 15.10 x 22.28 Font size: 10 | 0.03 | 0.28 | 0.02 | 18 | 3090 |
| | 150 | 0.05 | 0.77 | 0.02 | 19 | 12 | | 0.08 | 1.12 | 0.04 | 19 | 3187 |
| | 200 | 0.08 | 1.38 | 0.03 | 20 | 7 | | 0.13 | 1.99 | 0.05 | 20 | 3163 |
| | 300 | 0.15 | 3.09 | 0.04 | 20 | 7 | | 0.24 | 4.47 | 0.07 | 21 | 3204 |
| | 500 | 0.3 | 8.57 | 0.07 | 21 | 7 | | 0.52 | 12.41 | 0.12 | 22 | 3208 |
| | 700 | 0.45 | 16.78 | 0.1 | 22 | 7 | | 0.78 | 24.32 | 0.18 | 22 | 3205 |
| | 1200 | 0.85 | 49.33 | 0.18 | 23 | 12 | | 1.47 | 71.44 | 0.32 | 24 | 3210 |
| 16.29 x 23.47 Font size: 12 | 75 | 0.02 | 0.32 | 0.01 | 18 | 79 | 11.89 x 19.95 Font size: 12 | 0.02 | 0.2 | 0.01 | 18 | 159 |
| | 150 | 0.05 | 1.27 | 0.02 | 20 | 29 | | 0.05 | 0.79 | 0.02 | 19 | 2 |
| | 200 | 0.07 | 2.26 | 0.02 | 20 | 20 | | 0.09 | 1.4 | 0.03 | 20 | 2 |
| | 300 | 0.12 | 5.09 | 0.03 | 21 | 25 | | 0.17 | 3.16 | 0.04 | 20 | 1 |
| | 500 | 0.25 | 14.11 | 0.05 | 22 | 26 | | 0.33 | 8.77 | 0.07 | 21 | 1 |
| | 700 | 0.38 | 27.65 | 0.07 | 23 | 20 | | 0.5 | 17.16 | 0.11 | 22 | 2 |
| | 1200 | 0.71 | 81.21 | 0.14 | 24 | 22 | | 0.94 | 50.42 | 0.19 | 23 | 3 |
| 8.57 x 11.89 Font size: 7 | 75 | 0.007 | 0.09 | 0.006 | 17 | 1602 | 7.96 x 8.94 Font size: 7 | 0.002 | 0.06 | 0.002 | 17 | 442 |
| | 150 | 0.02 | 0.34 | 0.01 | 18 | 107 | | 0.006 | 0.24 | 0.003 | 18 | 49 |
| | 200 | 0.03 | 0.6 | 0.02 | 19 | 63 | | 0.009 | 0.42 | 0.004 | 19 | 46 |
| | 300 | 0.06 | 1.36 | 0.02 | 20 | 47 | | 0.02 | 0.94 | 0.005 | 19 | 41 |
| | 500 | 0.14 | 3.76 | 0.04 | 21 | 46 | | 0.04 | 2.62 | 0.009 | 20 | 43 |
| | 700 | 0.23 | 7.36 | 0.05 | 21 | 43 | | 0.06 | 5.13 | 0.01 | 21 | 45 |
| | 1200 | 0.45 | 21.6 | 0.1 | 22 | 47 | | 0.14 | 15.08 | 0.04 | 22 | 58 |

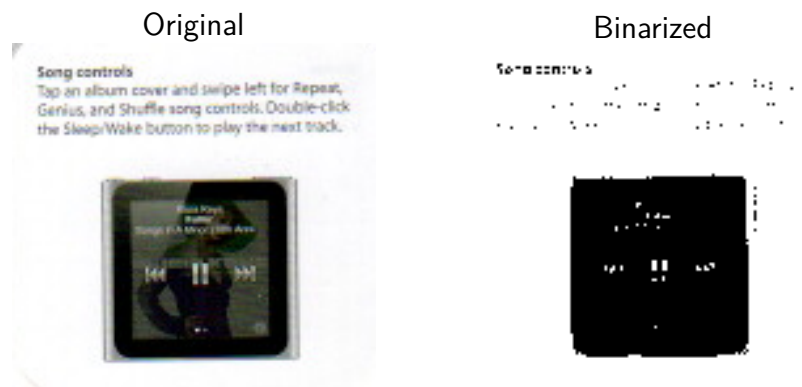| Image | dpi | png | bmp | tiff | HF | LD | Image | png | bmp | tiff | HF | LD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Song controls — Tap an album cover and swipe left for Repeat, Genius, and Shuffle song controls. Double-click the Sleep/Wake button to play the next track. 4.71 x 4.47 Font size: 5 | 75 | 0.001 | 0.02 | 0.0004 | 15 | 210 | Controles de música — Toque na capa de um álbum e passe o dedo para a esquerda para que apareçam os controles musicais Repetir, Genius e Aleatório. Clique duas vezes no botão Repouso/Despertar para reproduzir a próxima faixa. 5.01 x 5.35 Font size: 5 | 0.0008 | 0.02 | 0.0004 | 15 | 279 |
| | 150 | 0.003 | 0.07 | 0.001 | 16 | 85 | | 0.004 | 0.09 | 0.002 | 17 | 104 |
| | 200 | 0.005 | 0.13 | 0.002 | 17 | 53 | | 0.005 | 0.16 | 0.002 | 17 | 65 |
| | 300 | 0.009 | 0.28 | 0.004 | 18 | 43 | | 0.01 | 0.36 | 0.004 | 18 | 53 |
| | 500 | 0.02 | 0.77 | 0.008 | 19 | 25 | | 0.02 | 0.99 | 0.009 | 19 | 39 |
| | 700 | 0.04 | 1.51 | 0.01 | 19 | 24 | | 0.04 | 1.93 | 0.01 | 20 | 45 |
| | 1200 | 0.08 | 4.44 | 0.02 | 21 | 49 | | 0.09 | 5.67 | 0.03 | 21 | 48 |
| Controles de canción — Pulse la portada de un álbum y deslícese hacia la izquierda para que aparezcan los controles de canción Repetición, Genius y Aleatorio. Haga doble clic en el botón de reposo/activación para reproducir la pista siguiente. 4.71 x 5.38 Font size: 5 | 75 | 0.0008 | 0.02 | 0.0005 | 15 | 294 | Informações Importantes Sobre o Uso do iPod para Consumidores Brasileiros 5.59 x 4.78 Font size: 4 | 0.002 | 0.02 | 0.002 | 15 | 783 |
| | 150 | 0.004 | 0.09 | 0.002 | 17 | 116 | | 0.005 | 0.09 | 0.004 | 17 | 360 |
| | 200 | 0.006 | 0.15 | 0.002 | 17 | 69 | | 0.008 | 0.16 | 0.005 | 17 | 65 |
| | 300 | 0.01 | 0.34 | 0.005 | 18 | 47 | | 0.02 | 0.35 | 0.007 | 18 | 32 |
| | 500 | 0.03 | 0.93 | 0.009 | 19 | 52 | | 0.04 | 0.98 | 0.01 | 19 | 22 |
| | 700 | 0.04 | 1.82 | 0.01 | 20 | 48 | | 0.06 | 1.92 | 0.02 | 20 | 25 |
| | 1200 | 0.09 | 5.33 | 0.03 | 21 | 42 | | 0.13 | 5.64 | 0.03 | 21 | 20 |
| Safety and handling 8.57 x 8.64 Font size: 5 | 75 | 0.006 | 0.06 | 0.006 | 16 | 2586 | Benalet 18.05 x 25.26 Font size: 8 | 0.04 | 0.38 | 0.03 | 18 | 4066 |
| | 150 | 0.02 | 0.25 | 0.01 | 18 | 2325 | | 0.11 | 1.52 | 0.05 | 20 | 4152 |
| | 200 | 0.03 | 0.44 | 0.02 | 18 | 1976 | | 0.18 | 2.7 | 0.07 | 20 | 4153 |
| | 300 | 0.06 | 0.98 | 0.03 | 19 | 1981 | | 0.35 | 6.06 | 0.11 | 21 | 4148 |
| | 500 | 0.13 | 2.72 | 0.04 | 20 | 1983 | | 0.75 | 16.83 | 0.18 | 22 | 4156 |
| | 700 | 0.21 | 5.34 | 0.06 | 21 | 1979 | | 1.13 | 32.98 | 0.25 | 23 | 4150 |
| | 1200 | 0.45 | 15.66 | 0.1 | 22 | 1982 | | 2.15 | 96.93 | 0.46 | 24 | 4162 |

**Source:** The author (2024)

The analysis of the images in the Text Dataset shows that applying a low scan resolution, such as 75 dpi, leads to very poor OCR performance, with minimal document details preserved. In these cases, the Levenshtein distances were significantly higher than those at 150 dpi, highlighting that low resolutions hinder accurate recognition of document features, including finer character details. For binarized images, low resolution further exacerbates this issue, especially in images with small details. For instance, a small image from the dataset is illustrated in Figure 20, measuring 4.71 x 4.47 centimeters and with a font size of 5 pt. The figure shows that during binarization all text content was lost, resulting in an illegible image.

At higher resolutions, when analyzing the density of frequencies exceeding the maximum frequency obtained at 300 dpi, the difference in spectral densities observed approaches the level of scanning noise, with values typically below 0.001. This suggests that increasing the resolution beyond

Figure 20 – Small image, measuring 4.71 x 4.47 centimeters, in the resolution of 75 dpi.



**Source:** The author (2024)

300 dpi adds minimal useful detail, as the spectral density for such high frequencies does not significantly contribute to document clarity.

The maximum frequency observed in the 200 to 300 dpi range tends to be consistent, generally falling between 19 and 21. This stability in maximum frequency across these resolutions indicates that further increases in resolution yield diminishing returns in terms of detail capture. Consequently, the maximum frequency difference between 300 dpi and even higher resolutions is minor, with a maximum discrepancy of about three units. Such findings imply that 300 dpi may be an optimal balance point, as higher resolutions lead to larger file sizes without substantial improvements in detail discernment or spectral information.

Table 6 presents the percentage of frequency density exceeding the maximum frequencies obtained at 200 dpi and 300 dpi for the original text images. The results show that the use of different scanners did not result in significant changes in the frequency results, and the observed difference in density values can be considered a scanning error of the device.

Table 6 – Proportion of highest frequency density above the maximum frequency of 200 and 300 dpi in the General dataset - Original images
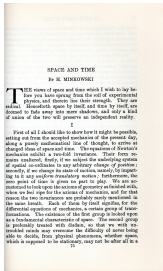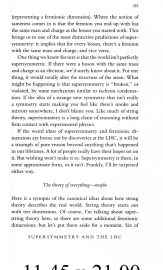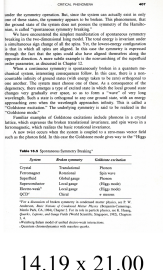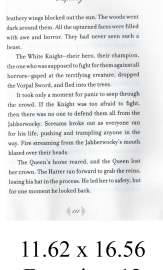
| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
| 11.89 x 18.95 Font size: 10 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1761% | 0.0000% | 0.1832% | 0.0000% | 0.3225% | 0.0000% |
| | 500 | 0.5880% | 0.1843% | 0.5502% | 0.1455% | - | - |
| | 700 | 0.9939% | 0.4210% | 0.8421% | 0.3093% | - | - |
| | 1200 | 1.2763% | 0.7350% | 0.9576% | 0.3971% | 1.7536% | 0.8336% |
| 11.45 x 21.00 Font size: 10 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.2400% | 0.0000% | 0.1377% | 0.0000% | 0.1707% | 0.0000% |
| | 500 | 1.2415% | 0.4823% | 0.6000% | 0.2044% | - | - |
| | 700 | 1.7568% | 0.8456% | 0.8859% | 0.3397% | - | - |
| | 1200 | 1.8204% | 1.0146% | 1.2866% | 0.6932% | 1.7131% | 0.8698% |
| 14.19 x 21.00 Font size: 10 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.1362% | 0.0000% | 0.4593% | 0.0000% |
| | 500 | 0.4851% | 0.1886% | 0.5081% | 0.1905% | - | - |
| | 700 | 0.8319% | 0.3965% | 0.7758% | 0.3259% | - | - |
| | 1200 | 1.0326% | 0.7835% | 1.1232% | 0.6761% | 1.8980% | 0.4013% |
| 11.62 x 16.56 Font size: 12 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1707% | 0.0000% | 0.1447% | 0.0000% | 0.1925% | 0.0000% |
| | 500 | 0.7210% | 0.1802% | 0.4661% | 0.1687% | - | - |
| | 700 | 0.9952% | 0.4215% | 0.7909% | 0.3424% | - | - |
| | 1200 | 1.1225% | 0.6127% | 0.9789% | 0.5623% | 1.0115% | 0.5206% |
| 15.68 x 22.28 Font size: 12 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1576% | 0.0000% | 0.1611% | 0.0000% | 0.1615% | 0.0000% |
| | 500 | 0.6866% | 0.2405% | 0.7174% | 0.1990% | - | - |
| | 700 | 0.8531% | 0.3534% | 0.9929% | 0.4385% | - | - |
| | 1200 | 1.0343% | 0.6024% | 1.2005% | 0.7363% | 1.3622% | 0.7875% |

| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
| 12.94 x 21.00 Font size: 11 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.4237% | 0.0000% | 0.1810% | 0.0000% |
| | 500 | 0.2566% | 0.1633% | 1.1932% | 0.1792% | - | - |
| | 700 | 0.3839% | 0.3000% | 1.5859% | 0.3264% | - | - |
| | 1200 | 0.6848% | 0.3939% | 1.6612% | 0.5077% | 1.5448% | 0.8194% |
| 11.04 x 21.00 Font size: 12 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 3.8212% |
| | 300 | 0.0000% | 0.0000% | 0.2335% | 0.0000% | 0.1679% | 8.0493% |
| | 500 | 0.1417% | 0.1417% | 0.7729% | 0.2639% | - | - |
| | 700 | 0.1387% | 0.1387% | 1.0906% | 0.5209% | - | - |
| | 1200 | 0.0000% | 0.0000% | 1.2987% | 0.8255% | 0.8636% | 14.3764% |
| 15.10 x 22.28 Font size: 10 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1622% | 0.0000% | 0.2316% | 0.0000% | 0.1615% | 0.0000% |
| | 500 | 0.5178% | 0.1711% | 0.7481% | 0.2823% | - | - |
| | 700 | 1.0329% | 0.4078% | 1.1366% | 0.4747% | - | - |
| | 1200 | 1.1921% | 0.6618% | 1.3727% | 0.7650% | 0.9998% | 0.4669% |
| 16.29 x 23.47 Font size: 12 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.2068% | 0.0000% | 0.4196% | 0.0000% | 0.1668% | 0.0000% |
| | 500 | 1.1964% | 0.4877% | 1.2297% | 0.1695% | - | - |
| | 700 | 1.5840% | 0.8494% | 1.7359% | 0.3271% | - | - |
| | 1200 | 1.3543% | 0.9224% | 1.4636% | 0.5573% | 1.1400% | 0.7245% |
| 11.89 x 19.95 Font size: 12 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.4489% | 0.0000% | 0.2411% | 0.0000% | 0.1794% | 0.0000% |
| | 500 | 0.9904% | 0.1568% | 0.9393% | 0.4603% | - | - |
| | 700 | 1.3210% | 0.2917% | 1.2009% | 0.6165% | - | - |
| | 1200 | 1.4147% | 0.4168% | 1.3511% | 0.7556% | 1.3881% | 0.7763% |

| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
| SANTO EZEQUIEL MORENO<br><br>8.57 x 11.89<br>Font size: 7 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.8463% | 0.0000% | 0.1748% | 0.0000% | 0.3930% | 0.0000% |
| | 500 | 2.5246% | 0.2690% | 1.0470% | 0.2753% | - | - |
| | 700 | 3.4287% | 0.6810% | 1.7451% | 0.6810% | - | - |
| | 1200 | 3.7217% | 1.1440% | 2.3127% | 1.2318% | 3.8615% | 1.6708% |
| 7.96 x 8.94<br>Font size: 7 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.0099% | 0.0000% | 0.0095% | 0.0000% | 0.1924% | 0.0000% |
| | 500 | 0.5667% | 0.5317% | 0.5033% | 0.4721% | - | - |
| | 700 | 1.2616% | 1.1879% | 1.0458% | 0.9786% | - | - |
| | 1200 | 1.8711% | 1.8147% | 1.5301% | 1.4701% | 2.0657% | 0.9401% |
| Song controls<br>4.71 x 4.47<br>Font size: 5 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0872% | 0.0000% | 0.0795% | - | - |
| | 200 | 0.0000% | 0.1244% | 0.0000% | 0.1796% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0442% | 0.0000% |
| | 500 | 0.0000% | 0.0636% | 0.0000% | 0.0629% | - | - |
| | 700 | 0.0000% | 0.1703% | 0.0000% | 0.1939% | - | - |
| | 1200 | 0.1592% | 0.5337% | 0.1721% | 0.6931% | 1.9900% | 1.7066% |
| Controles de música<br>5.01 x 5.35<br>Font size: 5 | 75 | 0.0000% | 0.0113% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.1026% | 0.0000% | 0.0939% | - | - |
| | 200 | 0.0000% | 0.0807% | 0.0000% | 0.1832% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.2167% | 0.0000% |
| | 500 | 0.0941% | 0.1968% | 0.0000% | 0.0617% | - | - |
| | 700 | 0.0605% | 0.1967% | 0.0000% | 0.1907% | - | - |
| | 1200 | 0.3266% | 0.7406% | 0.1496% | 0.6592% | 2.2168% | 0.7389% |
| Controles de canción<br>4.71 x 5.38<br>Font size: 5 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.1087% | 0.0000% | 0.2236% | - | - |
| | 200 | 0.0000% | 0.0815% | 0.0000% | 0.1731% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0543% | 0.0000% |
| | 500 | 0.0872% | 0.2101% | 0.0000% | 0.0710% | - | - |
| | 700 | 0.1865% | 0.3763% | 0.0000% | 0.2615% | - | - |
| | 1200 | 0.4501% | 0.9327% | 0.1572% | 0.8157% | 1.0271% | 0.8403% |

| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
| **5.59 x 4.78** **Font size: 4** | 75 | 0.0000% | 0.0000% | 0.1722% | 0.3560% | 0.0000% | 0.0000% |
| | 150 | 0.0107% | 0.0000% | 0.0000% | 0.1598% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.1164% | 0.0000% | 0.0000% |
| | 300 | 0.2360% | 0.0000% | 0.0000% | 0.0000% | 0.2684% | 0.0000% |
| | 500 | 1.2629% | 0.3359% | 0.0000% | 0.0000% | - | - |
| | 700 | 2.7519% | 0.8205% | 0.0000% | 0.0000% | - | - |
| | 1200 | 3.9866% | 2.0043% | 0.0000% | 0.0000% | 4.7040% | 2.0457% |
| **8.57 x 8.64** **Font size: 5** | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1662% | 0.0000% | 0.1666% | 0.0000% | 0.3770% | 0.0000% |
| | 500 | 0.7108% | 0.1727% | 0.7647% | 0.1858% | - | - |
| | 700 | 1.8023% | 0.5432% | 2.0381% | 0.6429% | - | - |
| | 1200 | 3.4152% | 1.7602% | 3.4642% | 1.7633% | 6.0758% | 2.7101% |
| **18.05 x 25.26** **Font size: 8** | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1725% | 0.0000% | 0.1515% | 0.0000% | 0.1584% | 0.0000% |
| | 500 | 0.6136% | 0.1534% | 0.5106% | 0.1687% | - | - |
| | 700 | 0.9497% | 0.3855% | 0.8566% | 0.3639% | - | - |
| | 1200 | 1.2721% | 0.6811% | 1.2115% | 0.6637% | 0.8021% | 0.3238% |

**Source:** The author (2024)

The results shown in Table 6 indicate that increasing the resolution does not lead to a significant rise in higher frequencies. The slight increase in higher frequency density observed is likely due to digitization noise.

Figure 21 depicts how often each resolution achieved the minimum Levenshtein Distance value, indicating the optimal Levenshtein Distance performance across various resolutions for each image. Notably, the 1200 dpi resolution did not achieve the best results for any of the 18 images assessed. In contrast, the 150 dpi resolution achieved the best outcomes in terms of the minimum Levenshtein Distance for 8 of the 18 images examined.

In some cases, as demonstrated in Figure 22, the identical Levenshtein distance outcomes were achieved across a range of five different resolutions. Specifically, the resolutions 150, 200, 300, 500, and 700 dpi all yielded the same text recognition results.

Therefore, in terms of choosing the best resolution equalizing all the digitizing aspects, it could be said that 200 dpi was the optimal resolution in this case, considering the identification of image

Figure 21 – Quantitative of image resolutions with minimal Levenshtein distance - 18 images ($\mu$: mean, $\sigma^2$: variance).



**μ: 25.8065%**
**σ²: 0.1915**

**μ: 22.5806%**
**σ²: 0.1748**

**μ: 19.3548%**
**σ²: 0.1561**

**μ: 12.9032%**
**σ²: 0.1124**

**μ: 12.9032%**
**σ²: 0.1124**

**μ: 6.4516%**
**σ²: 0.0604**

**μ: 0.0000%**
**σ²: 0.0000**

**Source:** The author (2024)

Figure 22 – Image illustrating the same Levenshtein Distance results to different resolutions.

| Image | Resolution (dpi) | Image File Size (MB) | | Highest Frequency | Levenshtein Distance |
|---|---|---|---|---|---|
| | | **PNG** | **BMP** | | |
| | 75 | 0.35 | 0.58 | 18 | 36 |
| | 150 | 1.24 | 2.31 | 19 | 7 |
| | 200 | 2.14 | 4.11 | 20 | 7 |
| | 300 | 4.49 | 9.26 | 20 | 7 |
| | 500 | 11.81 | 25.7 | 21 | 7 |
| | 700 | 21.68 | 50.35 | 22 | 7 |
| | 1200 | 56.13 | 147.98 | 23 | 12 |

**Source:** The author (2024)

details, file size and processing time linked.

The complete text dataset contains images with varying font sizes, ranging from 4 to 12 pt. When analyzing only the first ten images with a standard font size of 10 to 12 pt, Figure 23 shows

that the best Levenshtein distance was generally achieved at 150 or 200 dpi. This suggests that these resolutions are likely aligned with the OCR training data. Notably, even for images with smaller text sizes, 150 and 200 dpi remained the best resolutions. Moreover, these results remain optimal even when evaluating all eighteen images, as illustrated in Figure 23, where the performance is more balanced across all resolutions.

However, the results may vary depending on the use of a dictionary in the OCR process and the structural characteristics of the pages, such as whether they contain single or double columns. In general, the best Levenshtein distance was consistently found at 150 or 200 dpi, demonstrating that these resolutions already provide good results in terms of information retention and preservation.

Figure 23 – Quantitative of image resolutions with minimal Levenshtein distance - 10 images ($\mu$: mean, $\sigma^2$: variance).



**Source:** The author (2024)

## 4.2 GENERAL DATASET

Table 7 display the results obtained using the scanner Epson L4160 for the colored images in the General Dataset. The information includes image size in cm, resolution in dpi, file size for both PNG

and BMP formats in MB and the highest frequency for each image across the seven resolution levels.

Table 7 – General Dataset - Original images results. (**Image size** in centimeters (cm), **Font size** in points (pt), **dpi** - image resolution in dots per inch, **png/bmp** - file size of the png/bmp image in megabytes (MB), **HF** - Highest frequency).

| Image | dpi | png | bmp | HF | Image | png | bmp | HF |
|---|---|---|---|---|---|---|---|---|
| | 75 | 0.47 | 0.72 | 17 | | 0.36 | 0.67 | 18 |
| | 150 | 1.75 | 2.88 | 19 | | 1.18 | 2.66 | 19 |
| | 200 | 3.04 | 5.13 | 19 | | 1.97 | 4.74 | 20 |
| | 300 | 7.23 | 11.54 | 20 | | 4.01 | 10.66 | 21 |
| | 500 | 19.3 | 32.05 | 21 | | 10.03 | 29.62 | 22 |
| | 700 | 36.61 | 62.8 | 22 | | 17.38 | 58.03 | 22 |
| 13.75 x 21.00 | 1200 | 99.61 | 184.55 | 23 | 12.70 x 21.00 | 44.3 | 170.47 | 23 |
| | 75 | 0.35 | 0.59 | 18 | | 0.69 | 0.98 | 18 |
| | 150 | 1.25 | 2.37 | 19 | | 2.54 | 3.9 | 19 |
| | 200 | 2.22 | 4.22 | 20 | | 4.4 | 6.93 | 20 |
| | 300 | 4.91 | 9.48 | 20 | | 11.78 | 15.59 | 21 |
| | 500 | 13.27 | 26.34 | 21 | | 31.04 | 43.3 | 22 |
| | 700 | 24.71 | 51.64 | 22 | | 60.25 | 84.83 | 23 |
| 12.29 x 19.34 | 1200 | 65.27 | 151.69 | 23 | 16.90 x 23.16 | 162.8 | 249.23 | 24 |
| | 75 | 0.62 | 0.95 | 18 | | 0.62 | 1 | 18 |
| | 150 | 2.23 | 3.78 | 20 | | 2.48 | 3.99 | 20 |
| | 200 | 3.76 | 6.73 | 20 | | 4.29 | 7.1 | 20 |
| | 300 | 8.5 | 15.12 | 21 | | 10.4 | 15.97 | 21 |
| | 500 | 22.11 | 42.02 | 22 | | 27.18 | 44.36 | 22 |
| | 700 | 41.49 | 82.31 | 23 | 20.49 x 19.54 | 52.34 | 86.92 | 23 |
| 15.41 x 24.65 | 1200 | 111.64 | 241.84 | 24 | | 137.95 | 255.43 | 24 |
| | 75 | 0.79 | 1.55 | 19 | | 1.17 | 1.55 | 19 |
| | 150 | 2.89 | 6.22 | 20 | | 4.28 | 6.22 | 20 |
| | 200 | 5.21 | 11.07 | 21 | | 7.28 | 11.07 | 20 |
| | 300 | 12.22 | 24.89 | 21 | | 16.61 | 24.89 | 21 |
| | 500 | 32.88 | 69.17 | 22 | | 43.1 | 69.17 | 22 |
| | 700 | 62.73 | 135.52 | 23 | | 79.64 | 135.52 | 23 |
| 21.00 x 29.67 | 1200 | 175.47 | 398.27 | 24 | 21.00 x 29.67 | 214.77 | 398.27 | 24 |

| Image | dpi | png | bmp | HF | Image | png | bmp | HF |
|---|---|---|---|---|---|---|---|---|
| | **75** | 0.32 | 0.5 | 18 | | 0.78 | 1.28 | 18 |
| | **150** | 1.12 | 2 | 19 | | 2.73 | 5.14 | 20 |
| | **200** | 1.91 | 3.56 | 19 | | 4.91 | 9.13 | 20 |
| | **300** | 4.16 | 8.01 | 20 | | 12.67 | 20.51 | 21 |
| | **500** | 10.81 | 22.25 | 21 | | 33.31 | 56.97 | 22 |
| | **700** | 19.92 | 43.61 | 22 | 26.75 x 19.24 | 61.64 | 111.64 | 23 |
| 10.60 x 18.90 | **1200** | 52.31 | 128.08 | 23 | | 166.26 | 328.13 | 24 |

**Source:** The author (2024)

Table 8 show the results obtained using the scanner Epson L4160 to digitize the binarized images in the General Dataset, presenting the same information present in Table 7 for each image across the seven resolution levels but with the tiff file size beyond the PNG and BMP formats.

Table 8 – General Dataset - Binarized images results (**Image size** in centimeters (cm), **Font size** in points (pt), **dpi** - image resolution in dots per inch, **png/bmp/tiff** - file size of the png/bmp/tiff image in megabytes (MB), **HF** - Highest frequency).
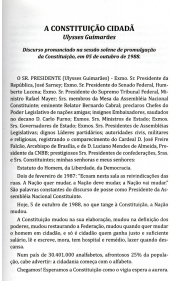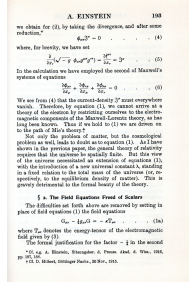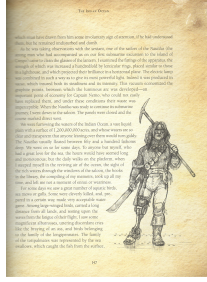
| Image | dpi | png | bmp | tiff | HF | Image | png | bmp | tiff | HF |
|---|---|---|---|---|---|---|---|---|---|---|
| | **75** | 0.01 | 0.24 | 0.01 | 17 | | 0.02 | 0.22 | 0.01 | 18 |
| | **150** | 0.04 | 0.96 | 0.02 | 19 | | 0.05 | 0.89 | 0.02 | 19 |
| | **200** | 0.06 | 1.71 | 0.03 | 19 | | 0.08 | 1.58 | 0.03 | 20 |
| | **300** | 0.15 | 3.85 | 0.08 | 20 | | 0.15 | 3.56 | 0.04 | 21 |
| | **500** | 0.34 | 10.7 | 0.19 | 21 | | 0.29 | 9.88 | 0.06 | 22 |
| | **700** | 0.61 | 20.94 | 0.31 | 22 | | 0.44 | 19.35 | 0.09 | 22 |
| 13.75 x 21.00 | **1200** | 1.46 | 61.52 | 0.65 | 23 | 12.70 x 21.00 | 0.81 | 56.82 | 0.16 | 23 |
| | **75** | 0.01 | 0.2 | 0.01 | 18 | | 0.03 | 0.33 | 0.02 | 18 |
| | **150** | 0.04 | 0.79 | 0.02 | 19 | | 0.08 | 1.3 | 0.04 | 19 |
| | **200** | 0.06 | 1.41 | 0.02 | 20 | | 0.13 | 2.31 | 0.06 | 20 |
| | **300** | 0.12 | 3.16 | 0.03 | 20 | | 0.3 | 5.2 | 0.19 | 21 |
| | **500** | 0.24 | 8.78 | 0.05 | 21 | | 0.69 | 14.44 | 0.38 | 22 |
| | **700** | 0.37 | 17.21 | 0.08 | 22 | | 1.3 | 28.29 | 0.7 | 23 |
| 12.29 x 19.34 | **1200** | 0.7 | 50.58 | 0.14 | 23 | 16.90 x 23.16 | 3.34 | 83.11 | 1.52 | 24 |

| Image | dpi | png | bmp | tiff | HF | Image | png | bmp | tiff | HF |
|---|---|---|---|---|---|---|---|---|---|---|
| | 75 | 0.02 | 0.32 | 0.02 | 18 | | 0.01 | 0.23 | 0.01 | 18 |
| | 150 | 0.07 | 1.26 | 0.03 | 20 | | 0.02 | 0.34 | 0.01 | 18 |
| | 200 | 0.11 | 2.24 | 0.04 | 20 | | 0.06 | 1.33 | 0.02 | 20 |
| | 300 | 0.2 | 5.04 | 0.06 | 21 | | 0.08 | 2.37 | 0.03 | 20 |
| | 500 | 0.42 | 14.01 | 0.09 | 22 | | 0.35 | 14.79 | 0.14 | 22 |
| | 700 | 0.64 | 27.44 | 0.14 | 23 | 20.49 x 19.54 | 0.59 | 28.98 | 0.22 | 23 |
| 15.41 x 24.65 | 1200 | 1.2 | 80.61 | 0.25 | 24 | | 1.32 | 85.14 | 0.42 | 24 |
| | 75 | 0.02 | 0.52 | 0.01 | 19 | | 0.05 | 0.52 | 0.03 | 19 |
| | 150 | 0.06 | 2.07 | 0.02 | 20 | | 0.14 | 2.07 | 0.05 | 20 |
| | 200 | 0.09 | 3.69 | 0.03 | 21 | | 0.21 | 3.69 | 0.07 | 21 |
| | 300 | 0.18 | 8.31 | 0.06 | 21 | | 0.37 | 8.31 | 0.12 | 21 |
| | 500 | 0.36 | 23.06 | 0.1 | 22 | | 0.75 | 23.06 | 0.21 | 22 |
| | 700 | 0.56 | 45.2 | 0.16 | 23 | | 1.15 | 45.2 | 0.31 | 23 |
| 21.00 x 29.67 | 1200 | 1.19 | 132.76 | 0.36 | 24 | 21.00 x 29.67 | 2.28 | 132.76 | 0.6 | 24 |
| | 75 | 0.01 | 0.17 | 0.01 | 18 | | 0.02 | 0.43 | 0.01 | 18 |
| | 150 | 0.03 | 0.67 | 0.02 | 19 | | 0.07 | 1.71 | 0.04 | 20 |
| | 200 | 0.05 | 1.19 | 0.02 | 20 | | 0.11 | 3.04 | 0.06 | 20 |
| | 300 | 0.1 | 2.68 | 0.03 | 20 | | 0.31 | 6.84 | 0.22 | 21 |
| | 500 | 0.22 | 7.43 | 0.05 | 21 | | 0.72 | 19 | 0.37 | 22 |
| | 700 | 0.33 | 14.55 | 0.07 | 22 | 26.75 x 19.24 | 1.27 | 37.22 | 0.53 | 23 |
| 10.60 x 18.90 | 1200 | 0.62 | 42.69 | 0.13 | 23 | | 2.98 | 109.4 | 0.98 | 24 |

**Source:** The author (2024)

When analyzing the General Dataset, similar trends were observed regarding maximum frequency variation and spectral density as in the Text Dataset. For example, Figures 4 to 10 in Table 7 illustrate the spectral density across frequencies, showing that the density beyond the 200 or 300 dpi cutoff is minimal, a pattern consistent across all images tested in the dataset. Table 9 also illustrates the percentage of frequencies higher than the maximum frequencies obtained in 200dpi and 300dpi, demonstrating that higher resolutions do not result in a significant amount of higher frequencies.

Based on the observations, non-historically significant files can be digitized at 200 dpi, which results in raster file sizes 2.25 times smaller than those at 300 dpi and offers over twice the file size reduction in PNG format.

Table 9 – Proportion of highest frequency density above the maximum frequency of 200 and 300 dpi in the Text dataset - Original images

| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
|  13.75 x 21.00 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.1829% | 0.0000% |
| | 500 | 0.1739% | 0.1739% | 0.1991% | 0.1991% | - | - |
| | 700 | 0.3401% | 0.3401% | 0.3627% | 0.3627% | - | - |
| | 1200 | 0.4561% | 0.4561% | 0.7504% | 0.7504% | 1.4719% | 0.6643% |
|  12.70 x 21.00 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1567% | 0.0000% | 0.1400% | 0.0000% | 0.1653% | 0.0000% |
| | 500 | 0.5437% | 0.1614% | 0.5848% | 0.1904% | - | - |
| | 700 | 0.8412% | 0.3623% | 0.8718% | 0.3346% | - | - |
| | 1200 | 1.0642% | 0.5784% | 0.9640% | 0.5272% | 0.6736% | 0.3319% |
|  12.29 x 19.34 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1658% | 0.0000% | 0.1392% | 0.0000% | 0.2803% | 0.0000% |
| | 500 | 0.6387% | 0.2244% | 0.5532% | 0.2114% | - | - |
| | 700 | 0.9439% | 0.4204% | 0.8171% | 0.3416% | - | - |
| | 1200 | 1.2231% | 0.7238% | 1.1250% | 0.6915% | 1.6518% | 0.8943% |
|  16.90 x 23.16 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 1.2229% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1750% | 0.0000% | 0.2145% | 0.0000% | 4.5415% | 0.0000% |
| | 500 | 0.7548% | 0.2599% | 0.8297% | 0.2478% | - | - |
| | 700 | 1.4052% | 0.5967% | 1.5002% | 0.5903% | - | - |
| | 1200 | 1.6978% | 0.8975% | 0.1737% | 0.0000% | 15.6025% | 0.9139% |
|  15.41 x 24.65 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.2299% | 0.0000% | 0.2531% | 0.0000% | 0.2568% | 0.0000% |
| | 500 | 0.9183% | 0.3665% | 1.0455% | 0.4420% | - | - |
| | 700 | 1.3367% | 0.6662% | 1.5035% | 0.7234% | - | - |
| | 1200 | 1.2570% | 0.6966% | 1.3554% | 0.7342% | 1.7939% | 0.9430% |

| Image | dpi | Epson L4160 | | Epson L8180 | | HP M479FDW | |
|---|---|---|---|---|---|---|---|
| | | 200dpi | 300dpi | 200dpi | 300dpi | 200dpi | 300dpi |
|  20.49 x 19.54 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1870% | 0.0000% | 0.0135% | 0.0000% | 0.1891% | 0.0000% |
| | 500 | 0.3562% | 0.1450% | 0.2116% | 0.1864% | - | - |
| | 700 | 0.3787% | 0.0958% | 0.4905% | 0.4638% | - | - |
| | 1200 | 0.7148% | 0.3492% | 0.7261% | 0.6995% | 0.8922% | 0.4309% |
|  21.00 x 29.67 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1520% | 0.0000% | 0.1537% | 0.0000% | 0.1560% | 0.0000% |
| | 500 | 0.5009% | 0.1817% | 0.5374% | 0.1549% | - | - |
| | 700 | 0.6664% | 0.3228% | 0.7438% | 0.3294% | - | - |
| | 1200 | 0.8160% | 0.5120% | 0.1691% | 0.0000% | 0.6554% | 0.3202% |
|  21.00 x 29.67 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1652% | 0.0000% | 0.1521% | 0.0000% | 0.1327% | 0.0000% |
| | 500 | 0.3020% | 0.1656% | 0.2983% | 0.1591% | - | - |
| | 700 | 0.3620% | 0.1034% | 0.3502% | 0.0973% | - | - |
| | 1200 | 0.4672% | 0.2773% | 0.4971% | 0.2937% | 0.6587% | 0.3190% |
|  10.60 x 18.90 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.2295% | 0.0000% | 0.1542% | 0.0000% | 0.1591% | 0.0000% |
| | 500 | 1.1176% | 0.4638% | 0.7368% | 0.2220% | - | - |
| | 700 | 1.4892% | 0.7468% | 1.0632% | 0.4667% | - | - |
| | 1200 | 1.5833% | 0.8906% | 1.2826% | 0.7881% | 1.1653% | 0.5329% |
|  26.75 x 19.24 | 75 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 150 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | - | - |
| | 200 | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% | 0.0000% |
| | 300 | 0.1722% | 0.0000% | 0.1572% | 0.0000% | 2.0679% | 0.0000% |
| | 500 | 0.1495% | 0.0093% | 0.1572% | 0.0196% | - | - |
| | 700 | 0.3703% | 0.0994% | 0.3599% | 0.1028% | - | - |
| | 1200 | 0.7404% | 0.3279% | 0.7280% | 0.3697% | 7.7849% | 0.5799% |

**Source:** The author (2024)

# 5  CONCLUSIONS

Professor George Nagy's inquiry regarding the most suitable scanner resolution for capturing all details and information from a document can be addressed by considering the optimal balance between quality and efficiency, particularly through the lens of the Nyquist-Shannon theorem. The analysis of spectral density and frequency maps of images is essential in determining the ideal digitization resolution, as it highlights the primary frequencies present in each image and reinforces the importance of operating within the limits established by this theorem.

As demonstrated in this study, scanning at 150 or 200 dpi generally captures all the essential details required for accurate OCR transcription of standard documents. These resolutions align well with the principles outlined by the Nyquist theorem, ensuring that necessary frequencies are adequately represented while maintaining manageable file sizes. This efficiency makes 150 dpi and 200 dpi the optimal minimum resolutions for storage and transmission, allowing for excellent lossless compression using formats like PNG for color images and TIFF-G4 for binary files.

Scanning at 300 dpi provides an additional "safety margin" by capturing finer details, which may be relevant for historical documents or those with very small text. While modern storage capacities and processing power can accommodate this higher resolution, it comes with trade-offs, including longer scanning times, increased storage requirements, greater network bandwidth usage, and extended processing durations. Notably, even with this increase in resolution, the primary frequency ranges of each image remain unchanged, as observed in the frequency maps and spectral density analyses. A comparative analysis with different scanners also confirms that increasing the resolution does not lead to a significant rise in higher frequency density. The slight increase observed is likely due to digitization noise rather than additional meaningful details.

Furthermore, this study has shown that binarizing document images does not significantly affect the maximum frequency captured, suggesting that increasing the digitization resolution beyond 300 dpi is generally unnecessary. Higher resolutions do not necessarily translate into proportional improvements in image quality and can lead to inefficiencies, such as increased scanning time, excessive storage and bandwidth consumption, and longer processing durations due to larger raster file sizes.

Thus, the findings suggest that resolutions between 150 and 300 dpi provide the best balance between quality and efficiency for most digitization needs, ensuring accurate OCR transcription and effective digital storage without unnecessary computational overhead.

# REFERENCES

BARBOSA, R. d. S.; LINS, R. D.; LIRA, E. D. F. D.; CAMARA, A. C. A. Later added strokes or text - fraud detection in documents written with ballpoint pens. In: *2014 14th International Conference on Frontiers in Handwriting Recognition*. [S.l.: s.n.], 2014. p. 517–522.

BARBOZA, R. d. S.; LINS, R. D.; JESUS, D. M. d. A color-based model to determine the age of documents for forensic purposes. In: *2013 12th International Conference on Document Analysis and Recognition*. [S.l.: s.n.], 2013. p. 1350–1354.

BARNEY-SMITH, E. H. Characterization of image degradation caused by scanning. *Pattern Recognition Letters*, 1998.

BARNEY-SMITH, E. H.; ANDERSEN, T. Text degradations and ocr training. *Eighth International Conference on Document Analysis and Recognition, 2005. Proceedings*, 2005.

BUSHBERG, J. T. *The essential physics of medical imaging*. [S.l.]: Bushberg, Jerrold T., 2012.

BUSHONG, S. C. *Radiologic science for technologists: Physics, biology and protection*. [S.l.]: Elsevier, 2013.

COLARUSSO, P. et al. Raman and infrared microspectroscopy. *Encyclopedia of Spectroscopy and Spectrometry*, Elsevier, p. 1945–1954, 1999.

doxygen. *tesseract::Dict Class Reference*. 2015. Last accessed 29 Jan 2025. Available at: <https://tesseract-ocr.github.io/tessapi/3.x/a00362.html#aa94c8519d79413969cbf94c9a9efcbd1>.

DOXYGEN. *OpenCV: Open Source Computer Vision. In: Image Processing in OpenCV: Image Thresholding*. 2024. Last accessed Jan 09, 2024. Available at: <https://docs.opencv.org/3.4/d7/d4d/tutorial_py_thresholding.html>.

EPSON. *L4160 User's Guide*. [S.l.], 2019.

EPSON. *L8180 User's Guide*. [S.l.], 2021.

Federal Agencies Digital Guidelines Initiative. *Technical Guidelines for Digitizing Cultural Heritage Materials - Creation of Raster Image Files*. [S.l.], 2016.

FORSCHUNGSGEMEINSCHAFT, D. *DFG-Praxisregeln, „Digitalisierung"*. [S.l.], 2016.

GONZALEZ, R. C.; WOODS, R. E. *Digital image processing*. fourth. [S.l.]: Pearson, 2018.

HP. *MFP HP Color LaserJet Pro série M478-M479*. [S.l.], 2019.

KENNEY, A. R.; CHAPMAN, S. *Requisitos de resolucao digital para textos : metodos para o estabelecimento de criterios de qualidade de imagem*. Rio de Janeiro: Projeto Conservação Preventiva em Bibliotecas e Arquivos, 2001. Caderno temático produzido pelo Projeto Conservação Preventiva em Bibliotecas e Arquivos.

LATHI, B. P. *Linear Systems and Signals*. Second. [S.l.]: Oxford University Press, 2005.

LEE, E. A.; VARAIYA, P. *Structure and Interpretation of Signals and Systems*. Second. [S.l.]: Pearson Education, 2011.

LEVENSHTEIN, V. I. Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics - Doklady*, 1966.

LINS, R. D. Two decades of document processing in latin america. *Journal of Universal Computer Science*, v. 17, n. 1, p. 151–161, 2011.

LINS, R. D.; MELLO; DE, D. R. N.; OLIVEIRA, R. C. d. Which is the most suitable scanner resolution for documents? Detailing the answer given to the question raised by Professor George Nagy. *DocEng '24: Proceedings of the ACM Symposium on Document Engineering 2024*, 2024.

LINS, R. D. et al. ICDAR 2021 Competition on Time-Quality Document Image Binarization. In: *ICDAR 2021*. [S.l.: s.n.], 2021. p. 1539–1546.

MATPLOTLIB. *Matplotlib 3.5.3 documentation*. 2024. Last accessed May 15, 2024. Available at: <https://matplotlib.org/3.5.3/index.html>.

NAGY, G.; SETH, S. Hierarchical representation of optically scanned documents. In: SEVENTH INTERNATIONAL CONFERENCE ON PATTERN RECOGNITION. *Proceedings Volume 1*. Montreal, Canada: Computer Society Press, 1984. p. 347–349. ISBN 0-81860545-6.

NAVARRO, G. A guided tour to approximate string matching. *ACM Computing Surveys (CSUR)*, 2001.

Numpy. *Sorting, searching, and counting*. 2024. Last accessed Feb 3, 2024. Available at: <https://numpy.org/doc/stable/reference/routines.sort.html>.

NUMPY. *Statistics*. 2024. Last accessed May 15, 2024. Available at: <https://numpy.org/doc/stable/reference/routines.statistics.html>.

NUMPY. *What is NumPy?* 2024. Last accessed Jan 27, 2024. Available at: <https://numpy.org/doc/stable/user/whatisnumpy.html>.

OPPENHEIM, A. V.; WILLSKY, A. S.; NAWAB, S. H. *Signals Systems*. [S.l.]: Prentice Hall, 1994.

OTSU, N. A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, v. 1, p. 62–66, September 1979.

PUGLIA, S. et al. *Technical Guidelines for Digitizing Archival Materials for Electronic Access: Creation of Production Master Files – Raster Images, For the Following Record Types- Textual, Graphic Illustrations/Artwork/Originals, Maps, Plans, Oversized, Photographs, Aerial Photographs, and Objects/Artifacts*. [S.l.], 2004.

PYPI. *pytesseract 0.3.10*. 2024. Last accessed March 20, 2024. Available at: <https://pypi.org/project/pytesseract/>.

PyPi. *Python-Levenshtein 0.23.0*. 2024. Last accessed April 24, 2024. Available at: <https://pypi.org/project/Python-Levenshtein/#description>.

SHANNON, C. E. A mathematical theory of communication. *Reprinted with corrections from The Bell System Technical Journal*, 1948.

SHANNON, C. E. Communication in the presence of noise. *Proceedings of the IEEE*, 1949.

SHUWANDY, M. L.; ZAIDAN, B.; ZAIDAN, A.; ALBAHRI, A.; ALAMOODI, A.; ALBAHRI, O.; ALAZAB, M. Health authentication approach based 3d touchscreen and microphone sensors for real-time remote healthcare monitoring system: Comprehensive review, open issues and methodological aspects. *Computer Science Review*, 2020.

SIEBERT, W. M. *Circuits, Signals, and Systems*. [S.l.]: MIT Press, 1986.

SONKA, M.; HLAVAC, V.; BOYLE, R. *Image Processing, Analysis, and Machine Vision*. United States of America: Cengage, 2015.

tesseract-ocr. *Tesseract documentation: A list of useful control parameters and config files*. 2024. Last accessed 29 Jan 2025. Available at: <https://tesseract-ocr.github.io/tessapi/3.x/a00362.html#aa94c8519d79413969cbf94c9a9efcbd1>.

The Digital Library Federation Benchmark Working Group. *Benchmark for Faithful Digital Reproductions of Monographs and Serials*. [S.l.], 2002.

The Library of Congress. *Guidelines for Electronic Preservation of Visual Materials*. 2014. Last accessed June 5, 2024. Available at: <https://www.loc.gov/preservation/resources/rt/guide/guid_dig.html#dig9>.

**Appendix**

# Appendix  A  –  CROPPING ALGORITHM

Source Code 1 – Cropping Algorithm applied in the dataset.

```
1  # Selecting the image
   image_file = "Home_original" #@param ["Home_original","Callidus_original", "
       Saltu_original"]
3  selected_image = "Image 01" #@param ["Image 01", "Image 02", "Image 03", "Image
       04", "Image 05", "Image 06", "Image 07", "Image 08", "Image 09", "Image 10",
       "Image 11", "Image 12", "Image 13", "Image 14", "Image 15", "Image 16", "
       Image 17", "Image 18", "Image 19", "Image 20", "Image 21", "Image 22", "Image
        23", "Image 24", "Image 25", "Image 26", "Image 27", "Image 28"]

5  # Setting destination folder based on image_file
   if image_file == "Callidus_original":
7      cropped_folder = "Callidus_cropped"
   elif image_file == "Saltu_original":
9      cropped_folder = "Saltu_cropped"
   elif image_file == "Home_original":
11     cropped_folder = "Home_cropped"
   else:
13     raise ValueError("Invalid image_file. Choose 'Home_original' or '
           Callidus_original'.")

15 image_resolutions = ["75ppp", "150ppp", "200ppp", "300ppp", "500ppp", "700ppp", "
       1200ppp"]

17 # Creating a list of paths to the original images
   images_paths = [f"{images_path}/{selected_image}/{image_file}/{res}.png" for res
       in image_resolutions]
19
   # Initializing the 'images' list before the loop
21 images = []

23 # Checking if images exist
   timages = []
25 for path in images_paths:
       img = cv2.imread(path)
27     if img is None:
           print(f"Error: Image not found at {path}")
29     else:
           print(f"Image {path} loaded")
31         images.append(img)

33 # Defining the cutting parameters
   top = 0   # @param {type:"slider", min:0, max:85, step:1}
35 bottom = 5  # @param {type:"slider", min:0, max:90, step:1}
```

```python
    left = 2  # @param {type:"slider", min:0, max:75, step:1}
37  right = 12  # @param {type:"slider", min:0, max:75, step:1}
    cut_values = {"Top": top, "Bottom": bottom, "Left": left, "Right": right}

39
    # Function to crop and save the image
41  def cut_and_save(image, output_path, resolution, cut_percent=(top, bottom, left,
        right)):
        if image is not None:
43          top, bottom, left, right = [int(p * max(image.shape) / 100) for p in
                cut_percent]
            cut_image = image[top:image.shape[0] - bottom, left:image.shape[1] -
                right]
45          cv2.imwrite(output_path, cut_image)
            return cut_image
47      else:
            print(f"The image ({resolution}) is empty.")
49          return None


51  # Function to save the cutoff values in an Excel file
    def save_cut_values(cut_values, excel_path):
53      df_corte = pd.DataFrame([cut_values], columns=['Top', 'Bottom', 'Left', '
            Right'])
        df_corte.to_excel(excel_path, index=False)
55      print(f"Cutoff values saved in: {excel_path}")


57  # Creating output path list of cropped images
    output_paths = [f"{images_path}/{selected_image}/{cropped_folder}/{res}.png" for
        res in image_resolutions]

59
    # Cropping and saving images
61  cut_images = []
    for i, image in enumerate(images):
63      resolution = image_resolutions[i]
        output_path = output_paths[i]
65      cut_image = cut_and_save(image, output_path, resolution)
        cut_images.append(cut_image)

67
    # Saving Cutoff Values to an Excel File
69  excel_path = f"{images_path}/{selected_image}/{cropped_folder}/excel_data/
        cut_values_{selected_image}.xlsx"


71  # Creating the directory if it doesn't exist
    os.makedirs(os.path.dirname(excel_path), exist_ok=True)
73  save_cut_values(cut_values, excel_path)


75  # Showing images
    titles_originais = [f'{selected_image} - Original ({res})' for res in
```

```
        image_resolutions]
77  titles_cut_images = [f'{selected_image} - Coropped ({res})' for res in
        image_resolutions]
    exibir_imagens(images + cut_images, titles_originais + titles_cut_images, 2, 7)
```

**Source:** The author (2024)

# Appendix B − OTSU BINARIZATION ALGORITHM

Source Code 2 − Otsu binarization Algorithm applied in the dataset.

```python
# Selecting the image
image_file = "Home_cropped" #@param ["Home_cropped","Callidus_cropped", "
    Saltu_cropped"]
selected_image = "Image 01" #@param ["Image 01", "Image 02", "Image 03", "Image
     04", "Image 05", "Image 06", "Image 07", "Image 08", "Image 09", "Image 10",
     "Image 11", "Image 12", "Image 13", "Image 14", "Image 15", "Image 16", "
    Image 17", "Image 18", "Image 19", "Image 20", "Image 21", "Image 22", "Image
     23", "Image 24", "Image 25", "Image 26", "Image 27", "Image 28"]

image_resolutions = ["75ppp", "150ppp", "200ppp", "300ppp", "500ppp", "700ppp", "
    1200ppp"]

# Creating a list of paths to the original images
images_paths = [f"{images_path}/{selected_image}/{image_file}/{res}.png" for res
    in image_resolutions]

# Initializing the 'images' list before the loop
images = []

# Checking if images exist
timages = []
for path in images_paths:
    img = cv2.imread(path)
    if img is None:
        print(f"Error: Image not found at {path}")
    else:
        print(f"Image {path} loaded")
        images.append(img)

# Function to binarize and save the image
def binarize_and_save(image, output_path, resolution):
    if image is not None:
        # Converting the image to grayscale
        gray_image = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)

        # Applying Otsu thresholding to grayscale image
        _, binarized_image = cv2.threshold(gray_image, 0, 255, cv2.THRESH_BINARY
            + cv2.THRESH_OTSU)

        cv2.imwrite(output_path, binarized_image)
        return binarized_image
    else:
        print(f"The image ({resolution}) is empty.")
```

```
36          return None

38 # Binarizing the images and saving them to the binarized_images list
   binarized_images = []
40 for i, image in enumerate(images):
       resolution = images_paths[i].split("/")[-1].split(".")[0]
42     output_path = f"{images_path}/{selected_image}/{image_file}/binarized_{
           resolution}.png"

44     binarized_image = binarize_and_save(image, output_path, resolution)

46     binarized_images.append(binarized_image)

48 # Showing the original and binarized images
   titles_originais = [f'{selected_image} - Original ({images_paths[i].split("/")
       [-1].split(".")[0]})' for i in range(len(images))]
50 titles_binarized = [f'{selected_image} - Binarized ({images_paths[i].split("/")
       [-1].split(".")[0]})' for i in range(len(binarized_images))]
   exibir_imagens(images + binarized_images, titles_originais + titles_binarized, 2,
        7)
```

**Source:** The author (2024)

# Appendix C – SPECTRAL ANALYSIS ALGORITHMS

Source Code 3 – Fourier Transform Algorithm applied in the dataset.

```python
# Selecting parameters
selected_image = "Image 01"  #@param ["Image 01", "Image 02", "Image 03", "Image
    04", "Image 05", "Image 06", "Image 07", "Image 08", "Image 09", "Image 10",
    "Image 11", "Image 12", "Image 13", "Image 14", "Image 15", "Image 16", "
    Image 17", "Image 18", "Image 19", "Image 20", "Image 21", "Image 22", "Image
     23", "Image 24", "Image 25", "Image 26", "Image 27", "Image 28"]
image_file = "Home_cropped"  #@param ["Home_cropped", "Callidus_cropped", "
    Saltu_cropped"]
color_image = "Color Image"  #@param ["Color Image", "Binarized Image"]

# Image resolutions for analysis
image_resolutions = ["75ppp", "150ppp", "200ppp", "300ppp", "500ppp", "700ppp", "
    1200ppp"]

# Image paths
images_path = "/content/drive/MyDrive/ScanImages"
if color_image == "Color Image":
    images_paths = [f"{images_path}/{selected_image}/{image_file}/{res}.png" for
        res in image_resolutions]
else:
    images_paths = [f"{images_path}/{selected_image}/{image_file}/binarized_{res
        }.png" for res in image_resolutions]

# Loading images
images = []
for path in images_paths:
    img = cv2.imread(path, cv2.IMREAD_GRAYSCALE)
    if img is not None:
        images.append(img)
    else:
        print(f"Error: Image not found at {path}")

def plot_images(images, titles, rows, cols, variables, frequencies):
    plt.figure(figsize=(18, 6 * rows))
    for i in range(len(images)):
        plt.subplot(rows, cols, i + 1)
        if i % cols == 1:  # Frequency Map
            im = plt.imshow(images[i], cmap='jet')
            plt.colorbar(im, ax=plt.gca(), fraction=0.05, pad=0.04)
        else:  # Original and Reconstructed Image
            plt.imshow(images[i], cmap='gray')
        plt.title(titles[i])
        plt.axis('off')
```

```python
            if variables and frequencies and i % cols == 0 and i // cols < len(
                variables) and i // cols < len(frequencies):
                plt.text(0.5, 1.1, f"Variable: {variables[i // cols]}\nMax Frequency:
                    {int(frequencies[i // cols])}",
                        horizontalalignment='center', verticalalignment='center',
                        transform=plt.gca().transAxes)
    plt.tight_layout()
    plt.show()


def save_to_excel(variables, frequencies, excel_path):
    os.makedirs(os.path.dirname(excel_path), exist_ok=True)  # Ensure directory
        exists
    data = {'Variable': variables, 'Max Frequency': frequencies}
    df = pd.DataFrame(data)
    df.to_excel(excel_path, index=False)
    print(f'Data saved to {excel_path}')

def save_frequency_maps(frequency_maps, frequency_map_dir, image_variables):
    os.makedirs(frequency_map_dir, exist_ok=True)
    for i, freq_map in enumerate(frequency_maps):
        freq_map_path = os.path.join(frequency_map_dir, f"{image_variables[i]}
            _frequency_map.png")
        plt.imsave(freq_map_path, freq_map, cmap='jet')
    print(f'Frequency maps saved to {frequency_map_dir}')


# Fourier Transform and Frequency Analysis
def calculate_max_frequency(magnitude_spectrum):
    row_center, col_center = np.unravel_index(np.argmax(magnitude_spectrum),
        magnitude_spectrum.shape)
    return magnitude_spectrum[row_center, col_center]


# Configuring layout for data visualization
rows = len(images)
cols = 3  # Original Image, Frequency Map, Reconstructed Image
all_images = []
image_variables = []
frequencies = []
frequency_maps = []


for image_path in images_paths:
    image_variables.append(image_path.split("/")[-1].split(".")[0])


# Processing each resolution
for image in images:
    transform_image = np.fft.fft2(image)
    center_transform_image = np.fft.fftshift(transform_image)
```

```
77      magnitude_spectrum = np.log(np.abs(center_transform_image) + 1)


79      # Image Reconstruction and Frequency Map Saving
        reconstructed_image = np.abs(np.fft.ifft2(np.fft.ifftshift(
            center_transform_image))).astype(np.uint8)
81
        max_frequency = calculate_max_frequency(magnitude_spectrum)
83      frequencies.append(max_frequency)
        frequency_maps.append(magnitude_spectrum)
85
        all_images.extend([image, magnitude_spectrum, reconstructed_image])
87
    titles = ["Original Image", "Frequency Map", "Reconstructed Image"] * len(images)
89
    # Display images and max frequency
91  plot_images(all_images, titles, rows, cols, image_variables, frequencies)


93  # Save results to Excel
    output_excel_path = os.path.join(images_path, selected_image, image_file, "
        excel_data", f"{selected_image}_frequencies.xlsx" if color_image == "Color
        Image" else f"binarized_{selected_image}_frequencies.xlsx")
95  save_to_excel(image_variables, frequencies, output_excel_path)


97  # Save frequency maps
    frequency_map_dir = os.path.join(images_path, selected_image, image_file, "
        frequency_maps")
99  save_frequency_maps(frequency_maps, frequency_map_dir, image_variables)
```

**Source:** The author (2025)

Source Code 4 – Spectral Density Plot Algorithm applied in the dataset.

```python
# Selecting parameters
selected_image = "Image 01"  #@param ["Image 01", "Image 02", "Image 03", "Image
    04", "Image 05", "Image 06", "Image 07", "Image 08", "Image 09", "Image 10",
    "Image 11", "Image 12", "Image 13", "Image 14", "Image 15", "Image 16", "
    Image 17", "Image 18", "Image 19", "Image 20", "Image 21", "Image 22", "Image
     23", "Image 24", "Image 25", "Image 26", "Image 27", "Image 28"]
image_file = "Home_cropped"  #@param ["Home_cropped", "Callidus_cropped", "
    Saltu_cropped"]
color_image = "Color Image"  #@param ["Color Image", "Binarized Image"]

# Image resolutions for analysis
image_resolutions = ["75ppp", "150ppp", "200ppp", "300ppp", "500ppp", "700ppp", "
    1200ppp"]
resolutions_dpi = ["75dpi", "150dpi", "200dpi", "300dpi", "500dpi", "700dpi", "
    1200dpi"]

# Create image paths
if color_image == "Color Image":
    images_paths = [f"{images_path}/{selected_image}/{image_size}/{res}.png" for
        res in resolutions]
else:
    images_paths = [f"{images_path}/{selected_image}/{image_size}/binarized_{res
        }.png" for res in resolutions]

# Load images, ensuring all are correctly loaded
images = [cv2.imread(path, cv2.IMREAD_GRAYSCALE) for path in images_paths]

# Define the default output format for images
pio.kaleido.scope.default_format = "png"

# Function to calculate spectral energy
def calculate_spectral_energy(magnitude_spectrum):
    return np.sum(magnitude_spectrum**2)

# Function to compute the maximum frequency
def calculate_max_frequency(magnitude_spectrum):
    row_center, col_center = np.unravel_index(np.argmax(magnitude_spectrum),
        magnitude_spectrum.shape)
    return magnitude_spectrum[row_center, col_center]

# Function to calculate the frequency distribution
def calculate_frequency_distribution(magnitude_spectrum):
    hist, bins = np.histogram(magnitude_spectrum.flatten(), bins=50, range=(0,
        magnitude_spectrum.max()))
    hist = hist / hist.sum()
    predominant_frequencies = bins[:-1][hist > 0.1]
```

```python
        return hist, bins, predominant_frequencies
37
    # Function to display and save frequency recurrence using Plotly
39  def plot_frequency_recurrence(hist, bins, title, xlabel, ylabel, save_path):
        fig = go.Figure()
41      bin_width = bins[1] - bins[0]

43      fig.add_trace(go.Bar(
            x=bins[:-1],
45          y=hist,
            marker=dict(color='#008000', line=dict(width=0)),
47          width=bin_width,
            opacity=1.0
49      ))

51      fig.update_layout(
            title=title,
53          xaxis_title=xlabel,
            yaxis_title=ylabel,
55          plot_bgcolor="#E5E5E5",
            xaxis=dict(showgrid=False),
57          yaxis=dict(showgrid=True, gridcolor="white"),
            bargap=0,
59          bargroupgap=0,
        )
61
        fig.write_image(save_path, engine="kaleido")
63      fig.show()

65  # Directory to save the graphs
    save_dir = os.path.join(images_path, "espectral_density", selected_image,
        color_image, image_size, "original")
67  os.makedirs(save_dir, exist_ok=True)

69  # Apply processing steps for each resolution
    frequencies = []
71  energies = []
    titles = [f"{selected_image} - {res}" for res in resolutions_dpi]
73
    for i, (image, res, res_dpi) in enumerate(zip(images, resolutions,
        resolutions_dpi)):
75      if image is None:
            print(f"Image {res} was not loaded correctly, skipping...")
77          continue

79      transform_image = np.fft.fft2(image)
        center_transform_image = np.fft.fftshift(transform_image)
```

```python
81      magnitude_spectrum = np.log(np.abs(center_transform_image) + 1)

83      max_frequency = calculate_max_frequency(magnitude_spectrum)
        frequencies.append(max_frequency)
85
        energy = calculate_spectral_energy(magnitude_spectrum)
87      energies.append(energy)

89      hist, bins, _ = calculate_frequency_distribution(magnitude_spectrum)

91      save_path = os.path.join(save_dir, f"spectral_density_{res_dpi}.png")
        print(f"Saving graph at: {save_path}")
93      plot_frequency_recurrence(hist, bins, f"{image_size} | Spectral Density - {
            titles[i]}: {color_image}", "Frequency", "Recurrence", save_path)


95  # Create DataFrame with results
    data = {
97      'Resolution': resolutions_dpi,
        'Max Frequency': frequencies,
99      'Energy': energies
    }
101 df = pd.DataFrame(data)


103 # Path to save the Excel file
    output_excel_path = os.path.join(images_path, selected_image, image_size, "
        excel_data", f"{selected_image}_spectral_analysis.xlsx")
105 os.makedirs(os.path.dirname(output_excel_path), exist_ok=True)
    df.to_excel(output_excel_path, index=False)
```

**Source:** The author (2025)

## Appendix D – OCR AND LEVENSHTEIN DISTANCE ALGORITHM

Source Code 5 – OCR and the Levenshtein Distance Algorithm applied in the dataset.

```python
# Selecting parameters
selected_image = "Image 18"  #@param ["Image 01", "Image 02", "Image 03", "Image
    04", "Image 05", "Image 06", "Image 07", "Image 08", "Image 09", "Image 10",
    "Image 11", "Image 12", "Image 13", "Image 14", "Image 15", "Image 16", "
    Image 17", "Image 18", "Image 19", "Image 20", "Image 21", "Image 22", "Image
     23", "Image 24", "Image 25", "Image 26", "Image 27", "Image 28"]
image_file = "Home_cropped"  #@param ["Home_cropped", "Callidus_cropped", "
    Saltu_cropped"]
color_image = "Color Image"  #@param ["Color Image", "Binarized Image"]


# Image resolutions for OCR
image_resolutions = ["75ppp", "150ppp", "200ppp", "300ppp", "500ppp", "700ppp", "
    1200ppp"]

# Image paths
images_path = "/content/drive/MyDrive/ScanImages"
if color_image == "Color Image":
    images_paths = [f"{images_path}/{selected_image}/{image_file}/{res}.png" for
        res in image_resolutions]
else:
    images_paths = [f"{images_path}/{selected_image}/{image_file}/binarized_{res
        }.png" for res in image_resolutions]

# Loading images
images = []
for path in images_paths:
    img = cv2.imread(path, cv2.IMREAD_GRAYSCALE)
    if img is not None:
        images.append(img)
    else:
        print(f"Error: Image not found at {path}")

# Directory to save OCR results
output_dir = f"/content/drive/MyDrive/ScanImages/OCR_Results/{selected_image}/{
    image_file}/{color_image}"
os.makedirs(output_dir, exist_ok=True)

# Tesseract configuration to disable dictionary
tesseract_config = "--psm 6 load_system_dawg=0 load_freq_dawg=0"
ocr_texts = []

# Performing OCR and saving results
for i, image in enumerate(images):
```

```python
        if image is None:
            print(f"Error processing image at {images_paths[i]}")
            continue

        # Performing OCR
        text = pytesseract.image_to_string(image, config=tesseract_config)

        # Saving OCR result
        output_file = os.path.join(output_dir, f"{image_resolutions[i]}_nd_ocr_result
            .txt")
        with open(output_file, 'w') as file:
            file.write(text)
        ocr_texts.append(text)
        print(f"Extracted text for resolution {image_resolutions[i]} saved to {
            output_file}")

# Reading the reference text
def read_file(file_path):
    with open(file_path, 'r', encoding='utf-8') as file:
        return file.read()

reference_path = f"/content/drive/MyDrive/ScanImages/References/{selected_image}
    _reference.txt"
reference_text = read_file(reference_path)

# Computing Levenshtein distances
distances = []
for i, ocr_text in enumerate(ocr_texts):
    distance = levenshtein(reference_text, ocr_text)
    distances.append({
        'Resolution': image_resolutions[i],
        'Levenshtein Distance': distance
    })

# Saving results to Excel
output_lev_dir = f"/content/drive/MyDrive/ScanImages/Levenshtein_Dist/{
    selected_image}/{image_file}/{color_image}"
os.makedirs(output_lev_dir, exist_ok=True)
output_file_path = os.path.join(output_lev_dir, f"{selected_image}_nd_levenshtein
    .xlsx" if color_image == "Color Image" else f"binarized_{selected_image}
    _levenshtein.xlsx")
df = pd.DataFrame(distances)
df.to_excel(output_file_path, index=False)

print(f"Levenshtein distances saved to: {output_file_path}")
```

**Source:** The author (2025)