# Text Removal Using Wavelet Transform and Morphological Operations

Author: Diana Banță <D.M.Banta-1@student.tudelft.nl>
Supervisors: Dr. Martin Skrodzki, Dr. Jorge Martinez Castaneda

## Introduction

- **Watermarks** are images embedded in paper used to identify the origins of historic documents (Figure 1) [1].
- **Paper degradation** and **overlapping ink marks** make it difficult to retrieve the shape of the watermark (Figure 2).
- Existing work on watermark retrieval faces limitations due to extent of applicability.
- One algorithm has proven to be highly effective for watermark retrieval, and is the only one that was found to address **text removal** [2]. It still presents limitations with text of certain size and contrast.
- Previous version of a **watermark recognition system** introduced promising line removal method, inspiring text removal concept[3, 4]. .

### Terminology

- **Wavelets**: wave-like oscillations which can decompose an image into multiple scales, allowing analysis over multiple detail levels.
- **Morphological operations**: operations applied to an image to adjust its pixels based on neighboring regions.
- **Contrast enhancement**: adjustment of image features for increasing the image quality and object

## Research Question

*How effective is the joint use of wavelet transform and morphological operations in the removal of text from watermark images, and how does it compare to algorithms using morphological operations and contrast enhancement?*
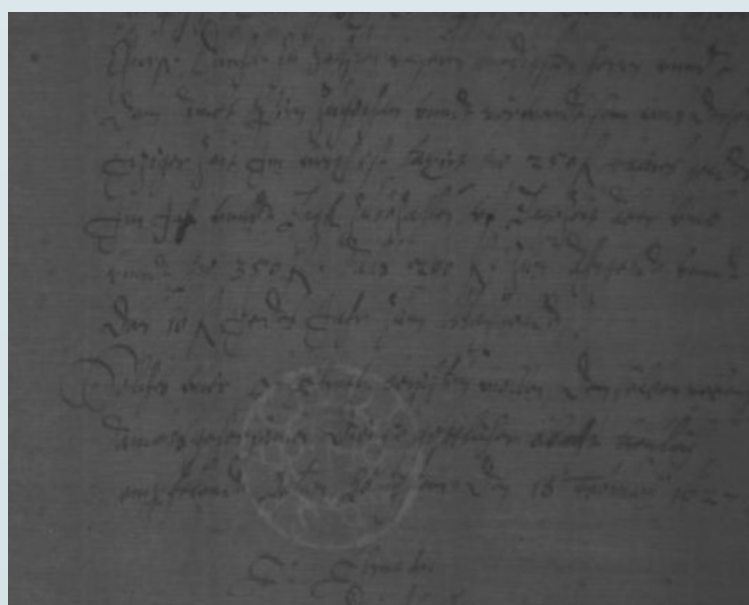

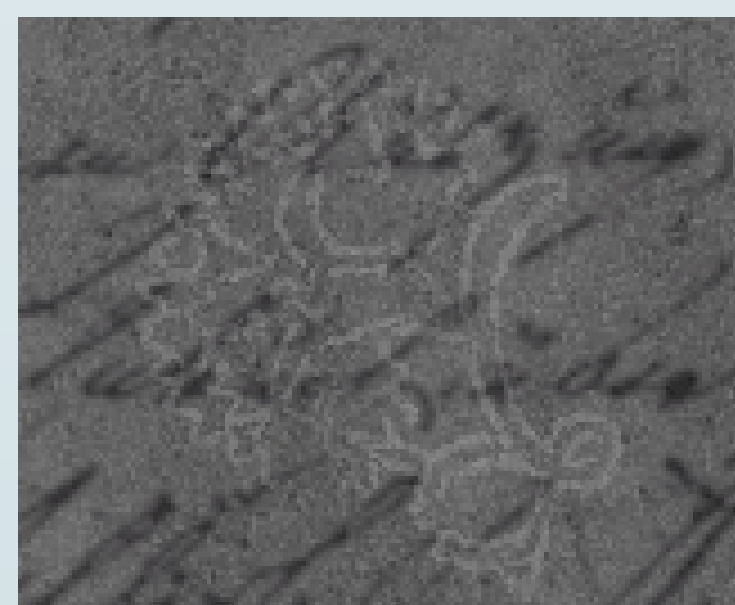Figure 1. Example of synthetic watermark image


Figure 2. Example of synthetic watermark image

## Methodology

### Dataset

Synthetically generated images resembling watermarks overlapped by text, created from three components:

- **Binarized watermark images** from the original dataset.
- **Image backgrounds** with different levels of noise from original dataset.
- **Handwritten text** images from three public databases of old documents [5,6,7].

Components were randomly selected and overlapped with random transparency, size, and position values. This resembles non-synthetic images while allowing control over text and watermark variables.

### Compared Algorithms

- Algorithm using contrast enhancement and morphological operations for estimating and removing foreground and text interference from an image.
- Proposed algorithm using wavelet transform and morphological operations for text localization and removal. Wavelet domain coefficients were used for creating a mask over the pixels containing text in the original image. Thresholding and morphological operations were then used to estimate the background intensity and replace text pixels.

### Experiments

Four datasets with varying text widths relative to watermark contour thickness, from very thin (Fig. 1-3) to very thick (Fig. 4-6), were used. Algorithms were evaluated on each dataset based on: **Original Watermark Conservation**, **Text Removed Successfully**, and **Processing Time**

The metrics used for comparing these results were **SSIM**, **MSE**, and **PSNR [7]**.

## Results

Table 1: Comparison of Evaluation Metrics of Proposed Approaches. _wav denotes the values of metrics computed for the proposed wavelet algorithm, while _ip denotes the metrics computed for the algorithm from literature. The values in bold correspond to the better score for each metric in each category.

| Type of Dataset | Evaluation Criteria | SSIM_wav | SSIM_ip | MSE_wav | MSE_ip | PSNR_wav | PSNR_ip |
|---|---|---|---|---|---|---|---|
| Thin Text | Watermark Conservation | **0.9209** | 0.8704 | **9.0006** | 9.6138 | 39.2003 | **39.1286** |
| | Text Removed | **0.9354** | 0.8973 | 7.4648 | **6.3803** | 39.8297 | **40.6988** |
| | Image Preservation | **0.9862** | 0.9728 | **1.6130** | 1.7505 | 47.0520 | **47.2130** |
| Very Thin Text | Watermark Conservation | 0.8727 | **0.8746** | 17.7264 | **17.8141** | 36.7704 | 36.1226 |
| | Text Removed | 0.8846 | **0.9197** | 13.5218 | **10.5727** | 37.5582 | **38.4590** |
| | Image Preservation | **0.9673** | 0.9664 | **3.7932** | 4.4782 | 43.4034 | 42.7634 |
| Thick Text | Watermark Conservation | **0.9261** | 0.7964 | **9.1501** | 14.4470 | 39.4024 | 36.9138 |
| | Text Removed | **0.9394** | 0.8370 | **7.6228** | 10.0531 | 39.8929 | 38.5258 |
| | Image Preservation | **0.9832** | 0.9499 | **2.2164** | 3.3884 | 45.8250 | 43.6180 |
| Very Thick Text | Watermark Conservation | **0.8986** | 0.7609 | **10.1473** | 13.8461 | 38.4659 | 37.3121 |
| | Text Removed | **0.9153** | 0.8177 | **8.4766** | 9.0970 | 39.2919 | 39.0829 |
| | Image Preservation | **0.9792** | 0.9498 | **2.1388** | 2.6673 | 45.4478 | 44.6648 |

Table 2: Comparison of Evaluation Metrics of Proposed Approaches. _wav denotes the values of metrics computed for the proposed wavelet algorithm, while _ip denotes the metrics computed for the algorithm from literature. The values in bold correspond to the better score for each metric in each category.

| Type of Dataset | Evaluation Criteria | SSIM_wav | SSIM_ip | MSE_wav | MSE_ip | PSNR_wav | PSNR_ip |
|---|---|---|---|---|---|---|---|
| Combined Small | Watermark Conservation | **0.8441** | 0.7736 | **18.1949** | 22.6296 | 36.2197 | 34.9482 |
| | Text Removed | **0.8652** | 0.8435 | 14.6018 | **14.3810** | 36.9065 | **36.9670** |
| | Original Image Preservation | **0.9567** | 0.9344 | **4.6923** | 6.6957 | 42.1373 | 40.7173 |
| Combined Large | Watermark Conservation | **0.9123** | 0.7907 | **9.8453** | 14.7049 | 39.1842 | 37.1416 |
| | Text Removed | **0.9284** | 0.8426 | **7.5011** | 9.0229 | 39.9298 | 39.2211 |
| | Original Image Preservation | **0.9811** | 0.9537 | **2.0425** | 2.9229 | 46.7739 | 45.2896 |


Figure 7. Result of baseline algorithm on Figure 4.


Figure 5. Result of baseline algorithm on Figure 4.


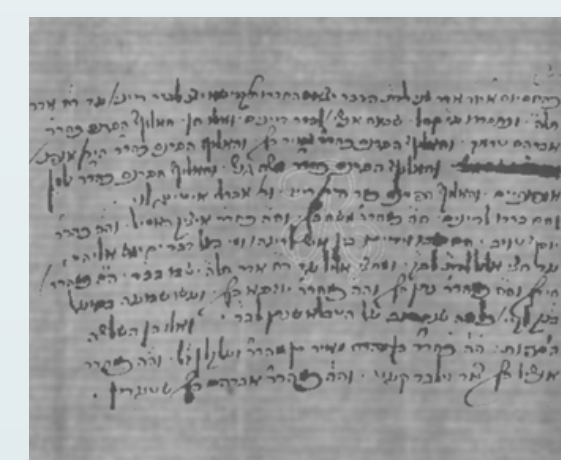Figure 5. Result of baseline algorithm on Figure 4.


Figure 4. Example of image in the 'Very Thick Text' dataset


Figure 5. Result of baseline algorithm on Figure 4.


Figure 6. Result of wavelet algorithm on Figure 4.
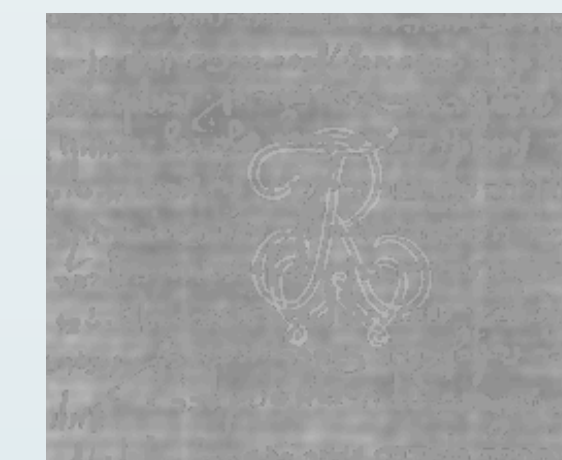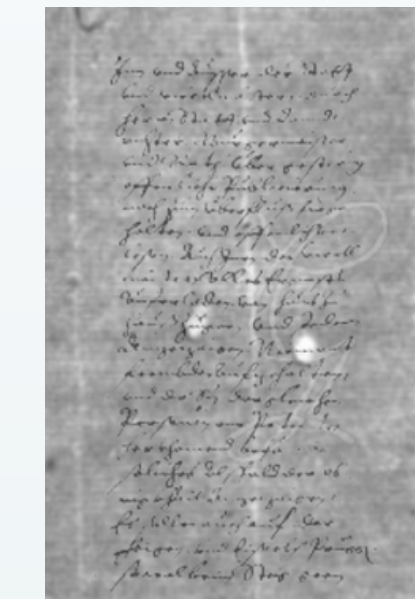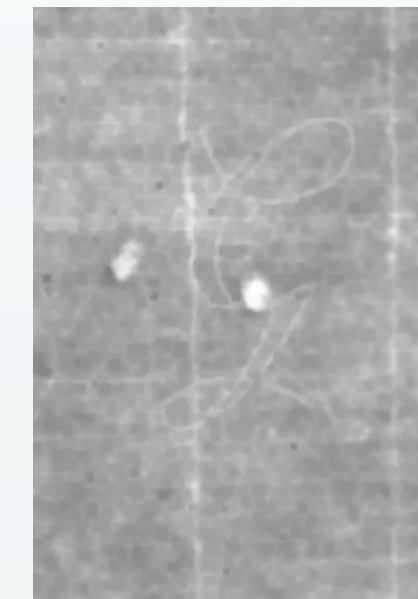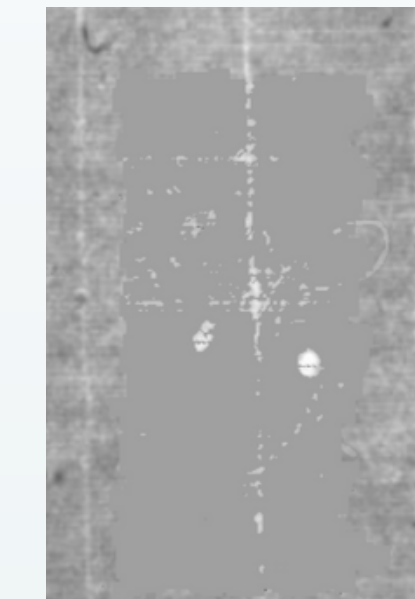
- Results are shown in Table 1 and Table 2.
- The proposed wavelet algorithm outperforms the baseline algorithm for most datasets.
- The worst performance is achieved for the case where text width is thinner than the watermark contour.
- Both algorithms had lowest personal scores when text was significantly thinner than the watermark
- Larger scores are obtained for when images are larger in size.
- The wavelet trasnform algorithm obtains lowest overall values for the metrics regarding text removed.

## Conclusions and Future Work

- The proposed algorithm is promising for thick text images.
- Limitations of the approach were found for images where background contrast is high, as well as for images with thin text.
- More work could be done for assessing the performance of the algorithms for non-synthetic data. Currently this is done with few images, only visually.
- Future work could be done for integrating together the two presented algorithms. Additionally, work could be done in adding Fourier transform to enhance the text localization.

## References

[1] L. Muller, "Understanding paper: Structures, watermarks, and a conservator's passion," Harvard Art Museums, May 2021.

[4] B. Munch, P. Trtik, F. Marone, and M. Stampanoni, "Stripe and ring artifact removal with combined wavelet - fourier filtering," Opt. Express, vol. 17, no. 10, pp. 8567–8591, May 2009

[3] D. Banta, S. Kho, A. Lantink, A. Marin, and V. Petkov, "A watermark recognition system: An approach to matching similar watermarks," 2023, last accessed 28 May 2024.

[2] H. A. Hiary, "Paper-based watermark extraction with image processing," Ph.D. dissertation, The University of Leeds, School of Computing, 2008

[5] M. Sonka, V. Hlavac, and R. Boyle, Image Processing, Analysis, and Machine Vision, 3rd ed. Cengage Learning.

[6] M.-S. Song, "Wavelet image compression," in Contemporary Mathematics, D. Han, P. Jorgensen, and D. Larson, Eds. American Mathematical Society, vol. 414, pp. 41–73.

[7] J. A. Sanchez, "Bentham dataset r0." [Online]. Available: https://zenodo.org/record/44519.

[8] J. A. Sanchez, V. Romero, A. H. Toselli, and E. Vidal, "ICFHR2016 competition on handwritten text recogni- tion on the READ dataset," in 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR). IEEE, pp. 630–635. [Online].

[9] I. Rabaev, B. K. Bakarat, and J. El-Sana, "The pinkas dataset."

[10] C. C. Beckner, Jr. and C. L. Matson, "Using mean-squared error to assess visual image quality," vol. 6313, p. 63130E. [Online]. Available: https://ui.adsabs.harvard.edu/abs/2006SPIE.6313E..0EB 10