

Master's thesis project in progress

Jawad Zaheer, Tuomas Aura
Secure Systems Group, Aalto University, Espoo, Finland
jawad.zaheer@aalto.fi
tuomas.aura@aalto.fi

1. Objective

Watermarking PDF documents for tracking ownership and information leaks online.

Exercise Sheet 10: Cloud Monitoring and Management

Questions and Answers

1. **Monitoring Services:** Compare cloud monitoring services like CloudWatch, Azure Monitor, and Stackdriver.

Answer:

- **CloudWatch (AWS):** Monitors resources and applications, collects and tracks metrics.
- **Azure Monitor:** Collects data from applications, infrastructure, and networks, providing full-stack monitoring.
- **Stackdriver (Google Cloud):** Offers monitoring, logging, and diagnostics for applications on Google Cloud and AWS.

2. **Performance Metrics:** Identify key performance metrics to monitor in cloud environments.

Answer: Key metrics include CPU utilization, memory usage, disk I/O, network traffic, response time, and error rates.

3. **Alerting and Automation:** Explain how to set up alerts and automate responses to incidents in the cloud.

Answer: Set up alerts using monitoring tools to trigger notifications or automated actions based on defined thresholds. Automation can include scaling resources, restarting services, or executing predefined scripts.

Exercise Sheet 10: Cloud Monitoring and Management

Questions and Answers

1. **Monitoring Services:** Compare cloud monitoring services like CloudWatch, Azure Monitor, and Stackdriver.

Answer:

- **CloudWatch (AWS):** Monitors resources and applications, collects and tracks metrics.
- **Azure Monitor:** Collects data from applications, infrastructure, and networks, providing full-stack monitoring.
- **Stackdriver (Google Cloud):** Offers monitoring, logging, and diagnostics for applications on Google Cloud and AWS.

2. **Performance Metrics:** Identify key performance metrics to monitor in cloud environments.

Answer: Key metrics include CPU utilization, memory usage, disk I/O, network traffic, response time, and error rates.

3. **Alerting and Automation:** Explain how to set up alerts and automate responses to incidents in the cloud.

Answer: Set up alerts using monitoring tools to trigger notifications or automated actions based on defined thresholds. Automation can include scaling resources, restarting services, or executing predefined scripts.

2. Solution

Encoding bit string

- **PDF normalization:** Convert document to a canonical PDF format.
- **Word space detection:** Locate the spaces in the PDF.
- Generate **personalized pseudorandom bit string** for each copy of the document.
- **Encoding:** Modify word spaces to embed the bit string.

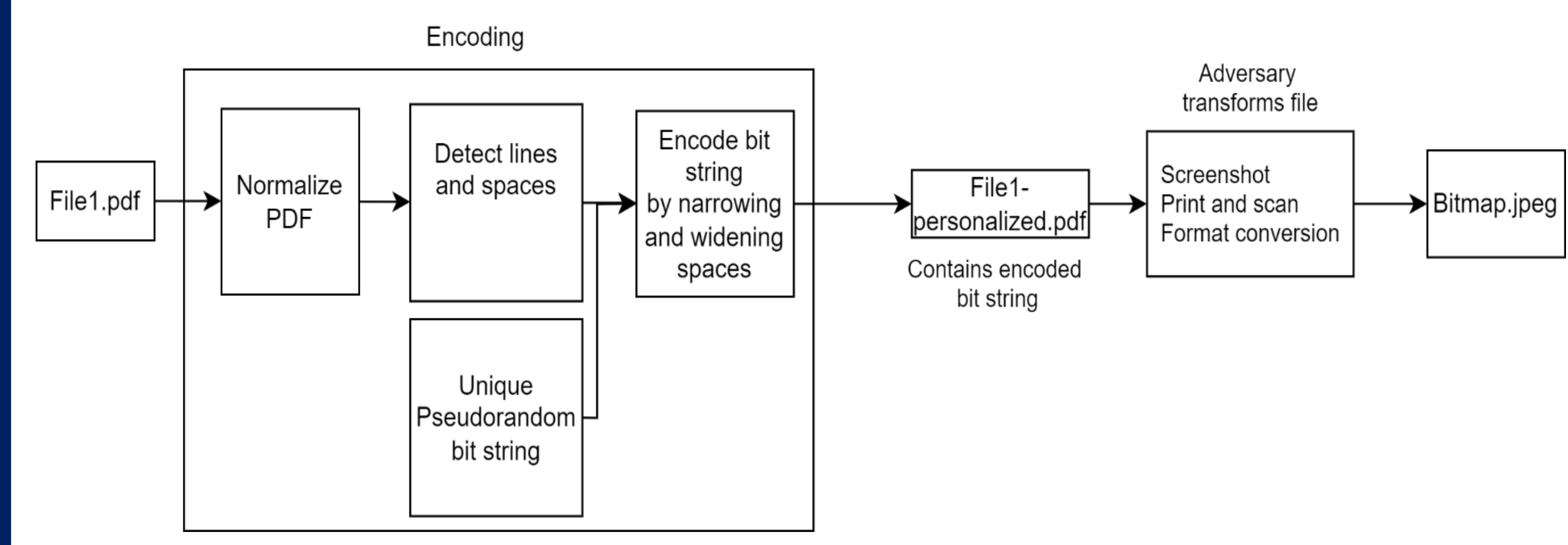


Fig. 1 : Encoding

Decoding bit string

- **Detect lines and word spaces** in the scanned document.
- **Decode bit string** from the image by comparing the word spacing in the original and scanned document.

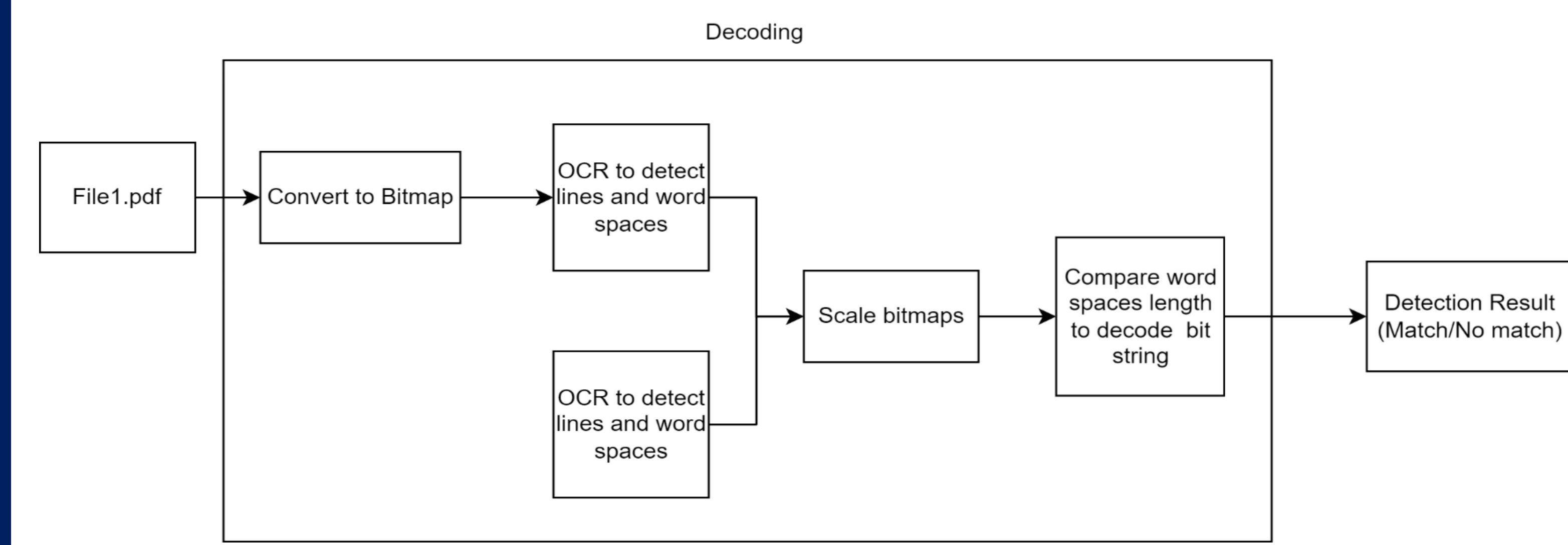


Fig. 1 : Decoding

3. Preliminary Results

Detection rate for sample PDF documents

Document type	Sample size	Bit error rate	Detection rate
Exercise sheets scalable cloud services	10	24.3%	Todo
Mathematics/ Linear Algebra Documents	10	33.8%	Todo
Applied mathematics Documents	10	24.5%	Todo

- A pseudorandom 100-bit string was encoded in all the files with spacing threshold=150.
- Bit error rate was calculated by finding error rate in all lines of a single PDF file and then detecting average error rates for all PDFs.
- Detection rate will be calculated later on by encoding different 100 bit strings in a single PDF and finding out if the required string is accurately detected or not.

4. Discussion

What we achieve

- **Robust watermarking:** Embed information via text spacing, resistant to print-scan attacks [2].
- **Efficient PDF creation:** Quickly generate compressed, watermarked PDFs at scale.
- **Invisible watermarks:** Encode without human detectable changes to PDF structure.

Limitations

- If adversary suspects watermarking, they may detect unusual word spacing.
- If adversary has two watermarked copies of the same document, they can detect and randomize the watermark.
- Document must have a sufficient number of long lines of text.
- Original and modified document should be present for detecting encoded text.
- Detecting word spaces in font changes and in mathematical equations poses significant challenges due to space widths for different font encodings.

5. Summary

This research develops a robust method for embedding and detecting watermarks in PDFs to protect ownership and prevent unauthorized distribution. The main goals are to create a watermarking algorithm that resists common attacks [3], works with LaTeX-converted PDFs, and embeds information invisibly. The process involves detecting line and word spaces in the original PDF, generating a random bit string, and embedding it by modifying these spaces. Both original and watermarked PDFs are then converted to bitmaps, and OCR is used to detect and compare lines and spaces. Differences reveal the embedded bit string, confirming the watermark. This ensures the watermark is imperceptible and robust, preserving the document's integrity and ownership.

6. References

1. Anderson, Ross & Petitcolas, Fabien. On The Limits of Steganography. (1998)
2. Kuribayashi, M., Fukushima, T., Funabiki, N, Data hiding for text document in PDF file.(2018)
3. Li, L., Zhang, H.J., Meng, J.L., Lu, Z.M. Robust pdf watermarking against print-scan attack.(2023).