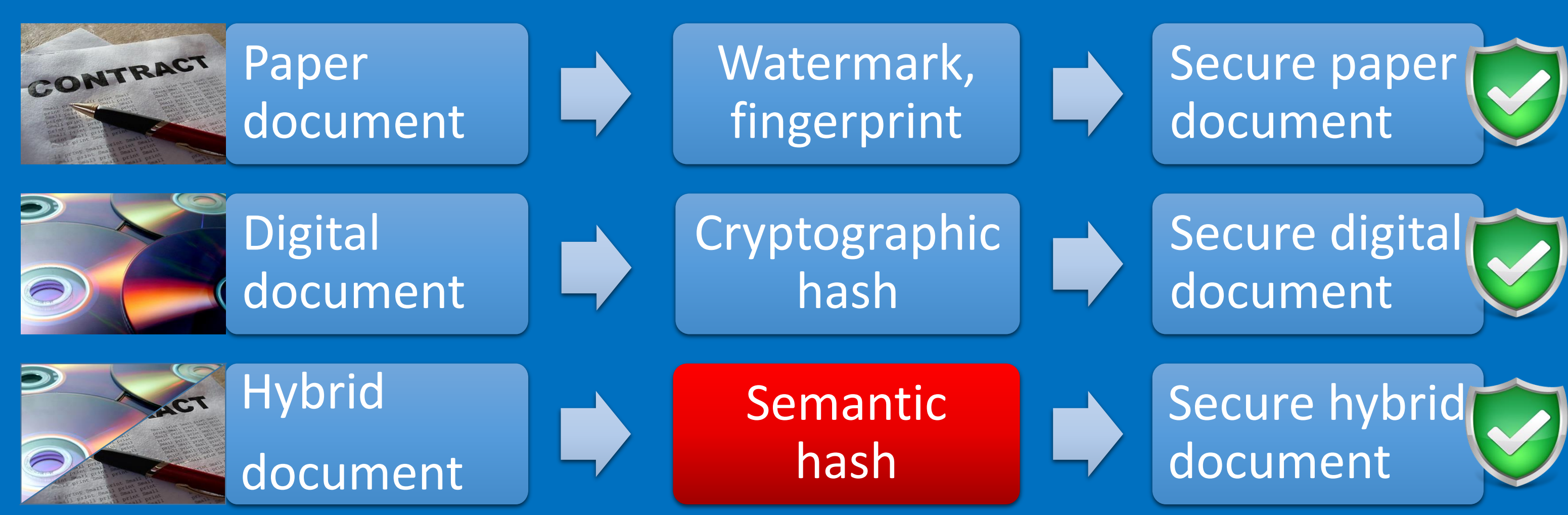
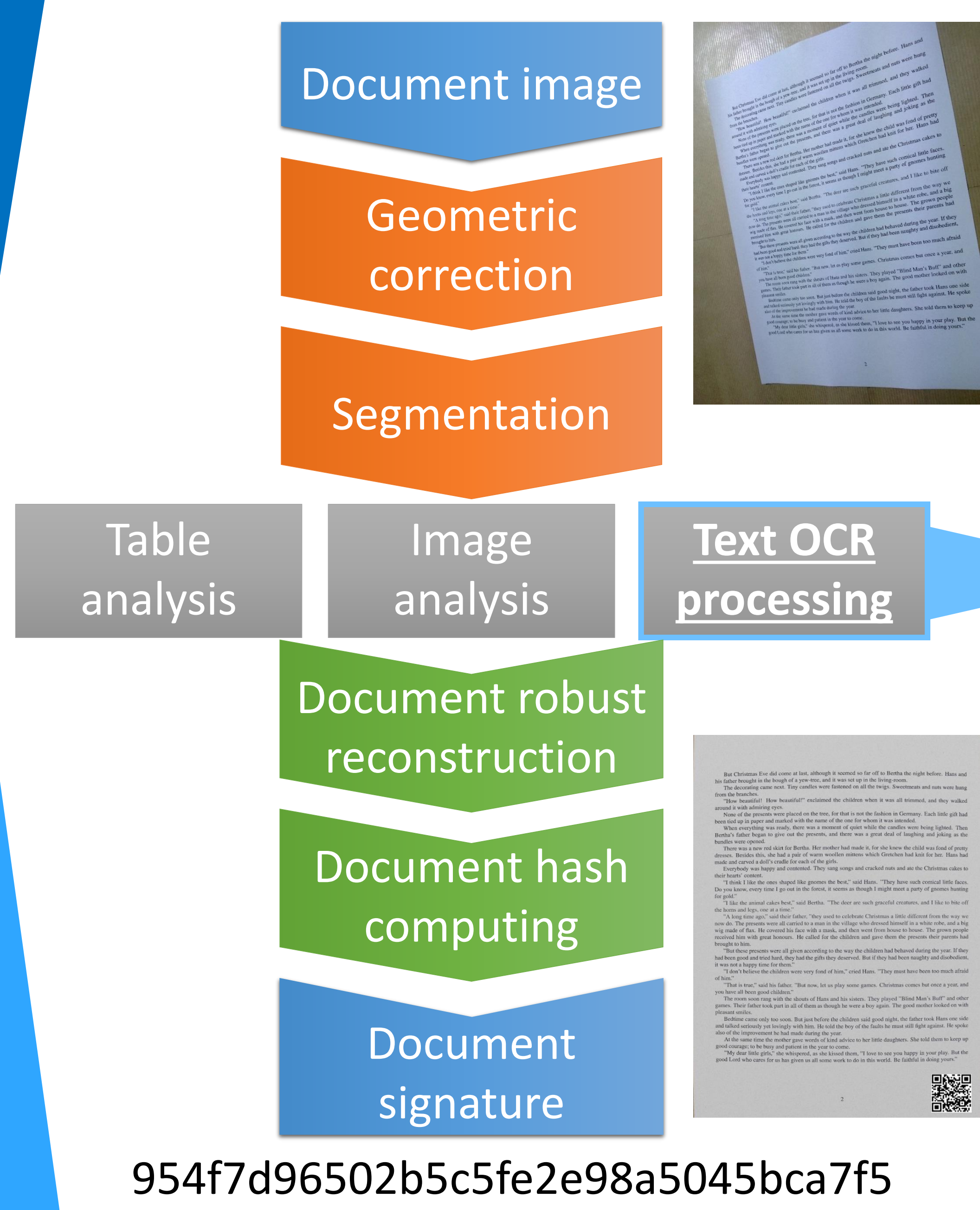


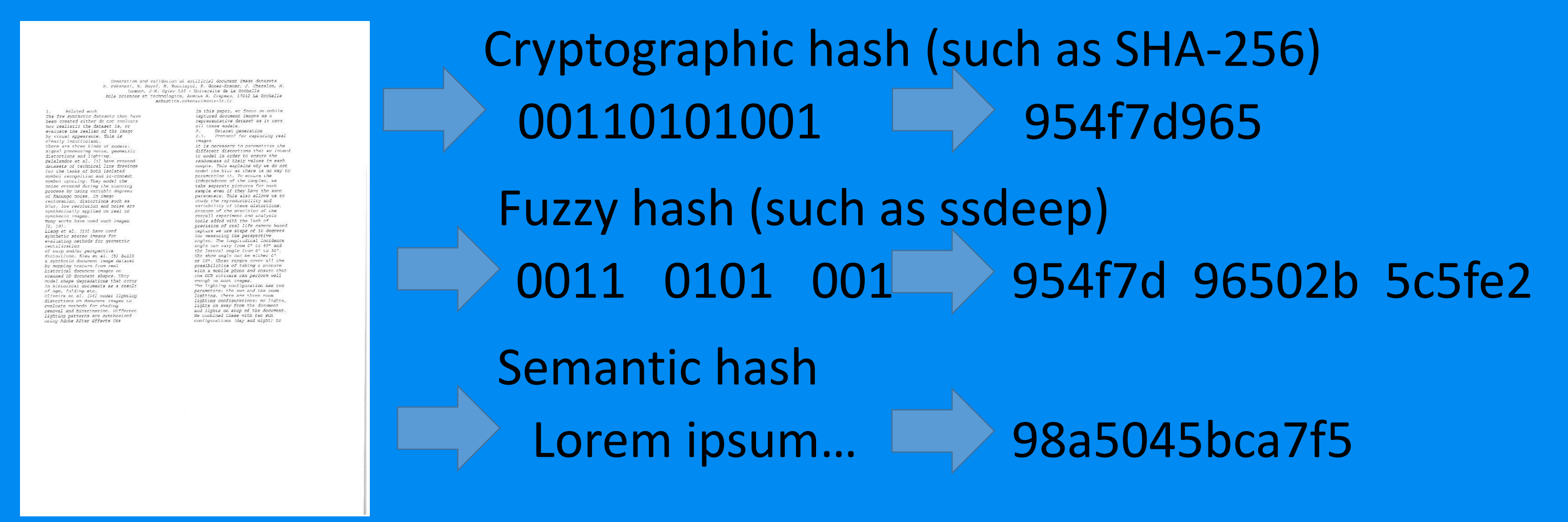
## Context: document security



## Hash computation process

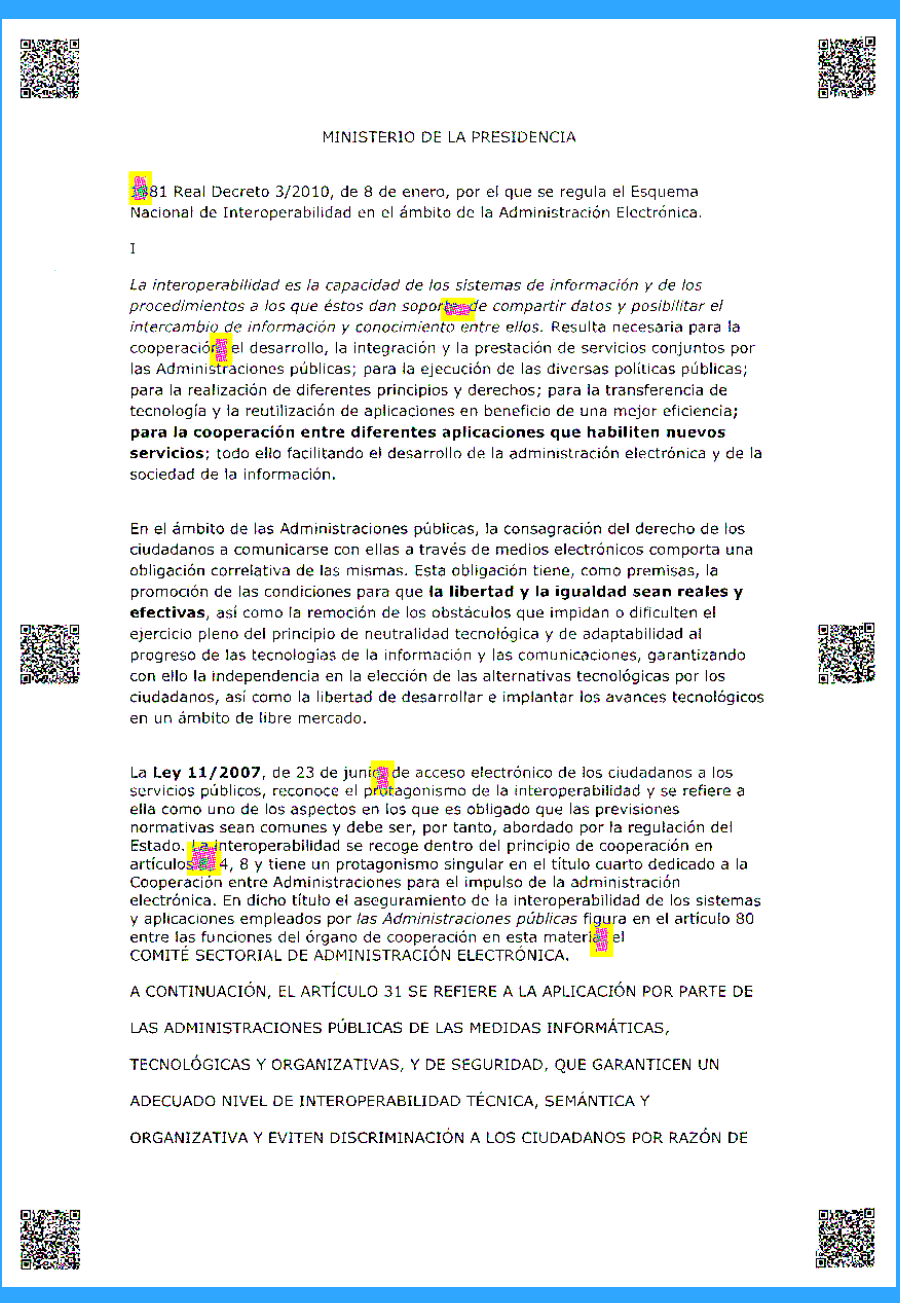


## What is a semantic hash ?



## State of the art and remaining challenges

- ### SIGNED project [1]
- Cut the document image in squares of 64 by 64 pixels
  - Apply a Haar Discrete Wavelet Transform on each square
  - Create a fuzzy hash
  - Other undisclosed pre- and post-processing steps.



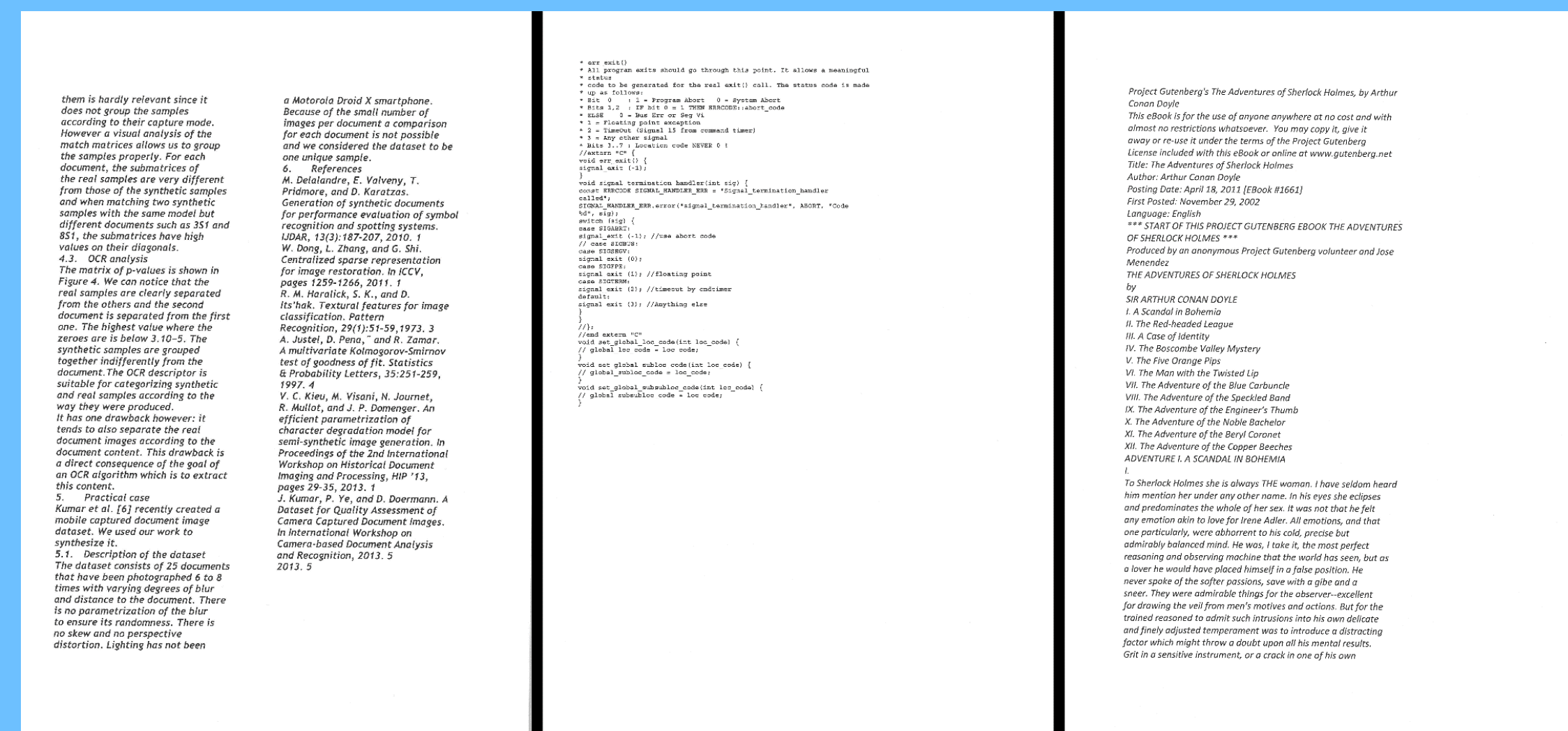
Document secured by the SIGNED project, it takes 6 2D-barcodes to embed the signature

## Text OCR processing [2]

### Test of Tesseract on 28512 document images

#### Dataset:

- 22 texts
- 6 fonts
- 3 font sizes
- 4 font emphasis
- 3 printers
- 3 scanners
- 3 scanning resolutions



**Conclusion: We need to study and improve the robustness of document analysis algorithms.**

### Best case scenario results:

- Accuracy : 99.83%
- Probability of false positive: 53% !!!
- Collision probability : 0.2%
- Other criteria are OK (throughput, hash size...)

### False positives:

- Occurs when two identical documents have different hashes
- Related to the robustness of the OCR algorithm
- Computed as a random draw of two copies of the same document:

$$P_i = \sum_{j=1}^{n_i} \left( \frac{n_{ij}}{\sum_{k=1}^{n_i} n_{ik}} \times \frac{n_{ij} - 1}{\sum_{k=1}^{n_i} n_{ik} - 1} \right)$$

- $n$  documents
- $n_i$  hashes per document
- Each hash is present  $n_{ij}$  times

Strengths	Weaknesses
Probability of false alarm <0.001	Cannot detect the replacement of dots by commas with an error <0.001
Probability of missed detection <0.001 for the replacement of digits	Cannot detect a manipulation smaller than 64x64 pixels at 600dpi (42x42 required)
Collision probability <0.001	Throughput >5s per page
Compatible with current scanners and printers	Hash size >4kB (between 4.8 and 170 kB)

References:  
 [1] A. Malvido Garcia: Secure Imprint Generated for Paper Documents (SIGNED). Technical Report December 2010, Bit Oceans (2013)  
 [2] S. Eskenazi, P. Gomez-Krämer, J.-M. Ogier: When document security brings new challenges to document analysis. International Workshop on Computational Forensics (IWCF), (2014)

For further information



l3i.univ-larochelle.fr