

# Security and Privacy for the Internet of Things: an overview of the project

<sup>+</sup>S. Aroua<sup>1</sup>, R. Champagnat<sup>1</sup>, M. Coustaty<sup>1\*</sup>, G. Falquet<sup>2</sup>, S. Ghadfi<sup>2</sup>, Y. Ghamri-Doudane<sup>1</sup>, P. Gomez-Krämer<sup>1</sup>, G. Howells<sup>2</sup>, K. D. McDonald-Maier<sup>4</sup>, J. Murphy<sup>3</sup>, M. Rabah<sup>1</sup>, K. Rouis<sup>1</sup>, N. Sidère<sup>1</sup> and N. Tamani<sup>1</sup>

**Abstract**—As the adoption of digital technologies expands, it becomes vital to build trust and confidence in the integrity of such technology. The SPIRIT project investigates the proof of concept of employing novel secure and privacy-ensuring techniques in services set-up in the Internet of Things (IoT) environment, aiming to increase the trust of users in IoT-based systems. The proposed system integrates three highly novel technology concepts developed by the consortium partners. Specifically, a technology, termed ICMetrics, for deriving encryption keys directly from the operating characteristics of digital devices; secondly, a technology based on a content-based signature of user data in order to ensure the integrity of sent data upon arrival; a third technology, termed semantic firewall, which is able to allow or deny the transmission of data derived from an IoT device according to the information contained within the data and the information gathered about the requester.

## I. INTRODUCTION

As the adoption of digital technologies expands, it becomes vital to build trust and confidence in the integrity of such technology. This paper presents one solution to this need developed in the framework of the SPIRIT project. The aim is to investigate a secure and trustworthy platform to enable trust and privacy for IoT devices, mainly documents and personal data they produce and share. Figure 1 illustrates the overall system and shows the different steps the data flows through from their generation (IoT devices), their aggregation (data owner) to their reception (data consumers).

With the development of connected devices, many new hardware appear and the primary issue to address in a secure IoT oriented system is to undertake technology enhancement in order to overcome the lack of user confidence which inhibits the use of IoT technology. Within this framework, we research and establish the feasibility of adapting the developed ICMetric technology of deriving encryption keys from traditional computing systems and smartphones to

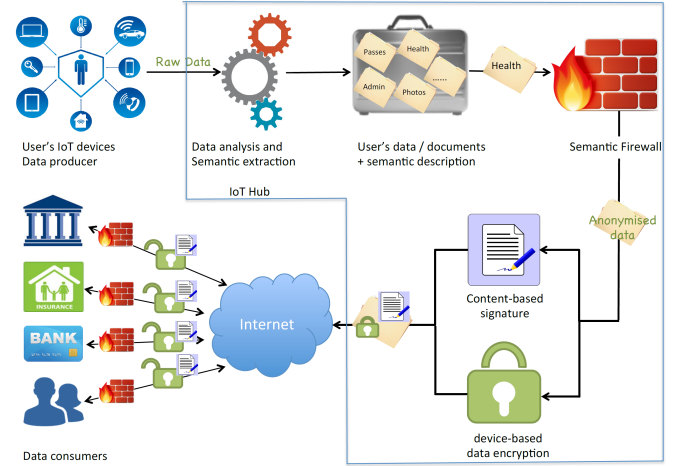


Fig. 1. Global overview of the SPIRIT project

enhance and enable security on IoT devices (nodes) and hubs (networked services), with a focus on features derived from the device software, hardware associated with the devices and employ the data so generated by the devices in a further authentication process based on content data analysis. This is a highly risky endeavour, as IoT devices are inherently simpler than the previously investigated systems and offers the challenge of determining characteristics of IoT devices which uniquely distinguish them.

In addition to device authentication technology, the second main objective in designing a privacy-by-design system is to ensure the authenticity of shared content by improving upon the research of creating a content-based signature of IoT generated user data. The signature provides a novel tool for authenticating the entire content of an IoT device generated document image in order to discover falsification. This signature is based on the documents content (text and graphics) and structure (spatial relationships). Thanks to a hashing of the documents information during the signature computation, no information from the original document will be deduced from its signature alone.

As for privacy-ensuring aspects, and in order to allow the user to get back control of information he/she is sharing, the last step of our global system relies on a semantic firewall technology. This is based on user privacy regulations to allow or deny the transmission of data derived from an IoT device according to the information contained within the data and the information gathered about the requester, hence ensuring

<sup>+</sup> Authors are order by alphabetical order and contributed in an equal way to the paper

<sup>1</sup>L3i Laboratory, La Rochelle Universit, 23 avenue Albert Einstein, BP 33060 - 17031 La Rochelle - France {firstname.lastname}@univ-lr.fr

\*Corresponding author {mickael.coustaty}@univ-lr.fr

<sup>2</sup>Information Science Institute, University of Geneva, 7 route de Drize, CH-1227 Carouge, Geneva, Switzerland Gilles.Falquet@unige.ch, Sami.Ghadfi@etu.unige.ch

<sup>3</sup>University of Kent, Canterbury, Kent, CT2 7NT {w.g.j.howells, J.Murphy-2060}@kent.ac.uk

<sup>4</sup>School of Computer Science and Electronic Engineering, INW.4.22, Colchester Campus, University of Essex kdm@essex.ac.uk

that access to such data is governed by the access permissions commensurate with the requester.

## II. RELATED WORKS

### A. ICMetrics

ICmetrics [1] represents an exciting new approach for generating unique identifiers for IoT and embedded devices enabling secure encrypted communication between devices potentially significantly reducing both fraudulent activity such as eavesdropping and device cloning. While data encryption techniques are now highly sophisticated and well established, encryption itself cannot necessarily protect against fraudulent data manipulation when the security of encryption keys cannot be absolutely guaranteed. The use of ICMetric authentication represents a novel concept of regulating access to devices and is explicitly aimed at providing protection at the vulnerable points where data access is initiated.

Specifically, ICMetrics possess the significant potential for secure communication from mobile and pervasive computing devices via the direct generation of digital signatures and encryption keys from the internal behavioural characteristics of software and hardware associated with the device. This naturally implies the major advantage that no encryption keys or device characteristic templates are stored. Significantly, characteristics are not limited to non-changing static values but may vary according to determinable patterns. The system will offer the following significant advantages:

- The removal of the need to store any form of template for reference data for validating the device, hence directly addressing the major weakness that the feature templates are accessed and used to circumvent the security afforded by the system.
- The security of the system is as strong as the ICMetric and encryption algorithm employed (there is no back door). The only mechanisms to gain subsequent access are to provide another sample of the ICMetric or to break the cipher employed by the encryption technology.
- The removal of the need for the storage of the private key associated with the encryption system. This is a natural consequence of the system since the key will be uniquely associated with the given ICMetric sample and a further ICMetric sample may be used to regenerate the required private key.

### B. Content-Based Signature

The identification of manipulated documents and the detection of these manipulations are two challenging problems in the document analysis field and in computer vision in general. On natural images, some authors approach these problems by analyzing the integrity of the image at frequency and spatial domains or by looking for specific patterns. Although these techniques proved to work on certain cases, its effectiveness on document images is not guaranteed due to the nature of the forgeries committed on this type of images. Specificities on document image forgery and tampering detection is detailed in [2] and a specific competition [3] has

been organized in 2018 to illustrate the challenges in this field.

The main trend of forgery detection in documents has been focused on the analysis of the content of the document and more precisely, on the appearance and shape of characters. In [4], the authors propose to detect forged text areas by analyzing unusual alignments between paragraphs or text lines. With respect to shape analysis, the authors of [5] present a method that compares intrinsic features from character shapes to identify copy-paste forgery cases. Font type also offers information to detect forgeries. Conditional Random Fields have been used to classify font types and to discriminate if a character is fake or genuine according to font features [6]. Some other works aims at using lower level features to capture inconsistencies at image level. For instance, in [7], authors uses Local Binary Patterns to obtained remarkable results.

Other works focuses on the verification of printer and scanner identification codes to certify the origin of the document and validate its global integrity the global integrity [8], [9]. Hyperspectral imaging techniques have been also used to detect ink color mismatch identified with added text into the original document [10].

Only few works have been presented on content-based hashing. One approach is to hash the document's content [11], [12], [13]. The drawback of these methods is that they only can secure the textual content of a document. Its layout or graphics are not secured. An other approach is fuzzy hashing based on the document's signal [14]. The fuzzy hashing breaches the possible confidentiality of the document, but this also allows for the localization of the modification. The main drawback is the big size of the digest.

### C. Document Classification and Knowledge Extraction

Document classification is a task that consists of assigning a category to a document, i.e. recognizing whether a document is an invoice, a purchase order, etc. Each category contains a specific set of knowledge units (e.g. name or address of a Person in an ID card, or Price and Quantity of a product in an Invoice) which are information of interest to be extracted from the document. The process of extracting such information is called **Information Extraction-IE**.

An agreement have been reached between the partners of the project to do IE only from images containing textual content (e.g. invoices, ID cards, insurance cards) so that use-case scenarios involving images can be tested to demonstrate the interaction between the components made by the different project partners to prevent the transmission of data to the wrong requester.

Methods for document image classification can be divided in 4 categories: (1) **Generalized n-gram approaches** which are based on a generalization of n-grams to describe a document in the same way a sentence in Natural Language Processing. Examples of such approaches are available in [15]and [16].(2) **Layout-based approaches** are based on finding documents with similar layout to an input document (some make also use of textual information like [17]) [18],

[19] [17], (3) **Template-based approaches** where the classification is based on finding which templates correspond most to an input document [20] [21], (4) **Text based approaches** make use of keywords to describe documents (e.g. examples of keywords in Invoice documents are "invoice", "qty" or "quantity"). The simplest methods are bag of words models which are based on computing a vector representation for each document and a classifier in a supervised way to learn how to assign a document to the right category.

Once the document is classified, an IE step is needed to extract relevant information for the Semantic Firewall, the goal is either to prevent the transfer of the whole document or just sensitive information in it. Many works have been done for IE from Image documents like in InforMys [22], Intellix [23], InDUS [24], CloudScan [25] or CUTIE [26]. The reader could find more details in those papers.

#### D. User Privacy Protection and Semantic Firewall

User Privacy Protection is a step further in user data protection mechanisms, which aims at enabling users to take a full control of their own personal data. It means to have a right to consent/deny sharing data with any external entity, to be informed about all the processes the data are subject to and the conclusions drawn from them, to be able to retract the consent for any reason, and to delete for good the data collected by service providers. Different privacy preserving mechanisms have been developed [27]. In addition to cryptography, we can consider the following privacy preserving methods:

1) *Anonymization techniques*: These techniques avoid the association between the user identity and his/her private data.  $K$ -anonymity constructs a new  $k$ -anonymous table with  $k$  records having the same Quasi-Identifier-Attribute. Every record occurs at least  $k$  times. This technique is not suitable for high dimensional data [28].

2) *Differential privacy*: It is used in statistical database. The basic idea is to perturb the raw records of users randomly to reduce the risk of privacy disclosure. However, [29] proves that differential privacy cannot provide guarantees in terms of average or maximum information leakage.

3) *Obfuscation/Perturbation*: The basic idea is to add noise or transform the user's private data to protect the association between a user and his/her sensitive data. These techniques incur an important data degradation and information loss.

The European regulation "General Data Protection Regulation" (GDPR - <https://gdpr-info.eu>) imposes on public institutions and service providers to provide a trustworthy environment for data sharing and usage. As Privacy concept has a gradual nature, then it should not be perceived as an all/nothing manner [30]. This why using semantic firewall is an appropriate solution. Semantics means using context information surrounding the data and the requester to decide whereas the sharing is legitimate or not. The knowledge can be based on logical rules such as in [31], [30], or learned from machine learning methods [32].

### III. PROPOSED SYSTEMS

Based on the literature review presented before, we now introduce our proof of concept on privacy management in content sharing based on the combination of all those steps.

#### A. ICMetrics

The first step proposed is the ICMetric system for securing the channel between devices. This is a two phase system with each phase generally operating as follows:-

Calibration phase (applied once only per application domain)

- For each sample: measure the desired feature values.
- Generate feature distributions for each feature illustrating the frequency of each occurrence of each discrete value for each sample circuit.
- Normalize the feature distributions generating normalisation maps for each feature.

Operation phase (applied each time an encryption key is desired for a given circuit)

- Measure features for the given ICMetric for which an encryption key is desired.
- Apply the normalisation maps to generate values suitable for key generation.
- Apply the key generation algorithm.

This SPiRiT work has developed practical approaches for analyzing the feature distributions associated with the behaviour of IoT devices[33]. In the context of ICMetrics, the proposed analysis presents some novel challenges as compared to many traditional pattern recognition tasks, because the distribution of values exhibited by the features or characteristics being investigated is more diverse than that found in traditional pattern recognition tasks, often a consequence of the software associated with the device operating in a number of distinct states. The problem of incorporating pattern features with unusual distributions is well known within pattern recognition problems, even if not always easily addressed. The problem is, however, more acute when features are derived with the additional requirement that the necessary information should not be easily available to a system attempting to clone the desired software. This is further exacerbated with features derived from characteristics of device usage, which nevertheless offer additional useful device defining features. The work has developed some novel analysis techniques to model pattern features with highly non-standard distributions derived from low-level behavioural characteristics, which address this issue to ensure that a digital signature is not easily faked.

#### B. Content-Based Signature

The content-based signature [34] step is the second one where we want to guarantee that the content read by the consumer has not been forged or modified from the original content sent. In order to create an automatic document security system one needs to secure the textual content but also the graphical content of the document. We then propose to use a hashing algorithm capable of securing

the graphical parts of paper and digital documents with unprecedented performance and a very small digest. The main challenge for such an algorithm is that of stability, in particular with respect to print and scan noise. We define the generic notion of stability and how to evaluate it. To achieve such performance we use both dense local information and global descriptors and then defined the ASYCHA method. We have tested our method on two datasets totaling nearly 45000 images and compared our method to two state-of-the-art methods. We finally obtained a significant improvement except on the False Negative Rate measure. We are actually working on this point to enhance this step. More details can be found in [35].

### C. Document Classification and Knowledge Extraction

The Document Classification and Knowledge Extraction modules aims at extracting knowledge from documents. The extracted knowledge consists of RDF triples which will be used by the Semantic Firewall to allow or deny access to data. If the access-rules used by the Semantic Firewall prohibit the requester from accessing sensitive information of the data owner then a typical action is to obfuscate such information.

The figure 2 gives an overview of the processes involved in the Document Classification and Knowledge Extraction modules.

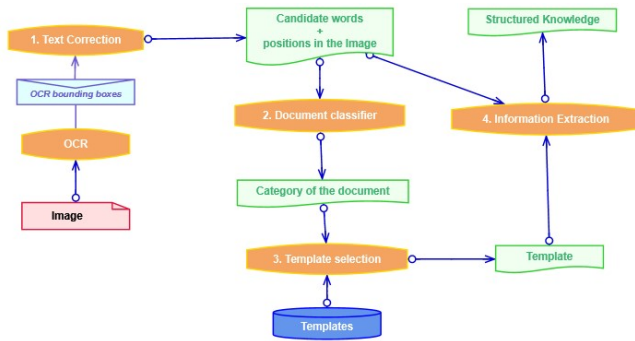


Fig. 2. Diagram showing the processes involved in Document classification and Information Extraction

We use the text output by an OCR engine to classify documents. Since such text is not perfect, there's a need for correcting this output (for brevity, we call it text correction). We designed two approaches for text correction, both approaches consider the word to correct as a query then the goal is to find the most relevant word in the vocabulary, where a word is considered relevant if it's character-wise structure is similar to the query (e.g. the word "designed" is more relevant to the query "desliigned" than the word "aligned", a reference measure of similarity between tokens is the Levenshtein distance): (1) the first approach is based on learning the maximal patterns that are common to both the query and words in the vocabulary, this helps to find fast a cluster of candidate words that are similar to the query, these candidates are then pruned using the Levenshtein distance,

(2) the second approach is based on learning character-wise embeddings (we used fastText [36] to compute these embeddings) of the vocabulary words, the learned vector representations are then stored in a KD-tree like structure and the search is done by an Approximate Nearest Neighbor search by transforming the query into a feature vector (following [36], the embedding of a query is the sum of the embeddings of its n-grams), at last pruning is done on the retrieved words using the Levenshtein distance.

We also created a text-based classifier which uses n-gram feature vectors (Frequent Sequential Patterns) that are representative of a given class of documents. The ranking of patterns is done based on how much a pattern is frequent in a given class compared to the other classes, the number of top ranked patterns to select is a hyper-parameter to set by cross-validation, these patterns are learned automatically in a pre-processing step, so each document is represented by a feature vector which is used as an input for a multi-class classifier. We experimented with different kinds of models for classification (multinomial logistic regression, SVM, feed forward neural networks with hidden layers), the one that performed well (at least 92% of macro F-Score) within reasonable runtime is the multinomial logistic regression model. The experiments we conducted for classification were done on a corpus of wikipedia articles containing 21 categories<sup>1</sup>.

Finally, we made a method to do Information Extraction from images which is robust to scale changes, translation (when an image is shifted), and cases where a part of the image is missing.

### D. User Privacy Protection with Semantic Firewall

The concept of "semantic firewall" provides an environment where data owners and data consumers can agree upon a set of data sharing rules, which are compliant with GDPR. For instance, it becomes possible to grant/deny access to data, as performed by any Access Control rule, but the owner can also specify privacy functions that can be used to allow data to be shared in a distorted form such as aggregation, obfuscation, anonymization, etc. This work is an extension of a previously published work [31].

1) *Logical-based access control*: A core of SPIRIT ontology is developed to gather the main concepts and the relationships among them. The following Fig. 4 sums up the logical data model describing IoT data and privacy concepts and functions to use to protect personal data.

2) *Privacy-aware query answering*: We express the problem of privacy preservation as a query rewriting process in a knowledge base, such that a given query is extended with the information about the requester and its privacy restrictions defined in the ontology. The following Fig. 3 displays the main steps of our privacy-aware query answering process.

The privacy algorithm starts by extracting the context of a query, based on requester identification information

<sup>1</sup>astronomy, biology, economy, food, football, genetics, geography, health, heritage, informatics, literature, mathematics, medicine, movie, music, politics, religion, rugby, sculpture, skating, tennis. Still

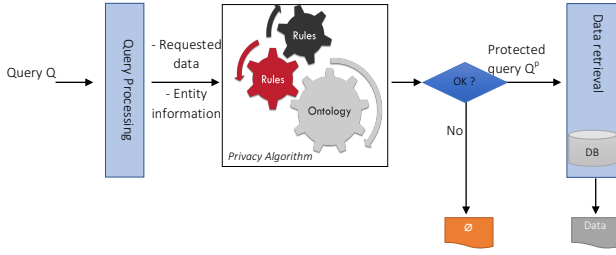


Fig. 3. Privacy-aware query answering main steps.

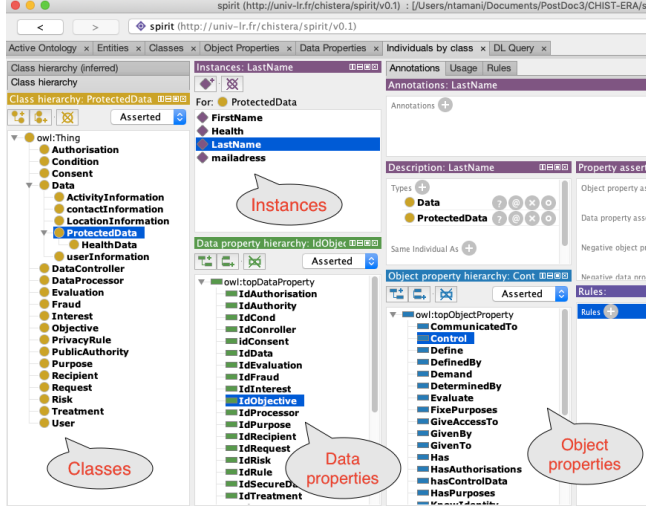


Fig. 4. Fragment of SPIRIT Ontology.

included in the query. It uses the ontology to check the access rights defined by the data owner as privacy rules. If the requester has not any access right to the data, then the query is rejected. If the requester has no access limit to the data then the privacy-aware query is similar to the original query ( $Q^P = Q$ ). If the requester has access to data under some privacy rules, then a privacy-aware query, denoted by  $Q^P$ , is generated based on privacy functions to apply on the requested data. Once query  $Q^P$  is obtained, it is processed by the data retrieval module to fetch the protected data to the requester.

3) *Proof of concept*: We developed a small proof of concept as a first step before integration of the semantic firewall with the other modules developed by the other partners of the project.

The use-case we implemented shows the different rules for data sharing that we instantiated as follows:

My employer is requesting my sport activity data: for my employer, the rule implemented is: “do not share the time and location of my activities, aggregate the rest of the data”. Figure 5 displays the results computed for the my employer, university of La Rochelle (ULR).

My doctor practitioner is requesting the same data: the rule is “share all the data with no limit”. Figure 6 displays the results computed for the my doctor, doctor AZAR.

4) *Remaining functions to implement*: For evolving the prototype, the next step is implementing the privacy toolkit

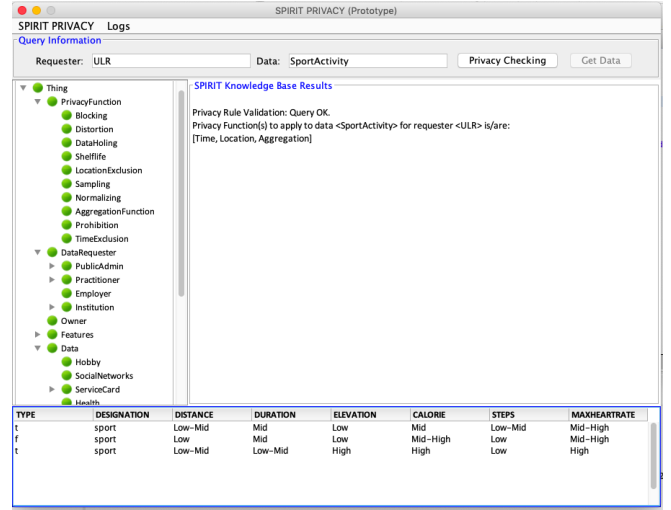


Fig. 5. Aggregated data sent to my employer.

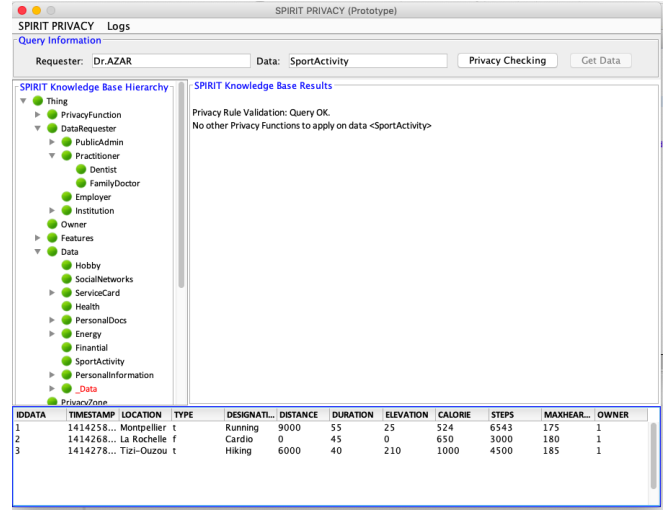


Fig. 6. Data sent to my doctor practitioner.

defining the privacy functions mentioned within the SPIRIT ontology such as anonymization, obfuscation, data shelf-life, etc. Within the SPIRIT integration chain, the semantic firewall is located in between the knowledge extraction module and the encryption module. Therefore, we need to implement the interface with the knowledge extraction module to feed the ontology with required data according to the use cases considered, and the format of the data to send to the encryption module.

## IV. CONCLUSIONS

The SPIRIT project investigates the proof of concept of employing novel secure and privacy-ensuring techniques in services set-up in the Internet of Things environment, aiming to increase the trust of users in IoT-based systems. We then presented our system which integrates three highly novel technology concepts developed by the consortium partners. Specifically, a technology, termed ICMetrics, for deriving encryption keys directly from the operating characteristics of



digital devices; secondly, a technology based on a content-based signature of user data in order to ensure the integrity of sent data upon arrival; thirdly, a Semantic firewall which is able to allow or deny the transmission of data derived from an IoT device according to the information contained within the data and the information gathered about the requester. Even if no global evaluation of our system can be proposed in its current state, each part of the system is currently evaluated and first results are really encouraging.

## ACKNOWLEDGMENT

The research leading to these results has received funding from CHIST-ERA which is a transnational R&D program jointly funded by the national funding organizations within the Framework Program 7, and by the Région Nouvelle-Aquitaine and the ERDF funds from the European Commission under the CPER-FEDER program.

## REFERENCES

- [1] Y. Kovalchuk, K. D. McDonald-Maier, and G. Howells, "Overview of icmetrics technology security infrastructure for autonomous and intelligent healthcare system," *International Journal of u- and e-Service, Science and Technology*, vol. 4 (3), 2011.
- [2] F. Cruz, N. Sidère, M. Coustaty, V. P. D'Andecy, and J.-M. Ogier, "Categorization of document image tampering techniques and how to identify them," in *CVAUI/IWCF/MIPPSNA@ICPR*, ser. Lecture Notes in Computer Science, vol. 11188, 2018, pp. 117–124.
- [3] C. Artaud, N. Sidère, A. Doucet, J. Ogier, and V. P. D'Andecy, "Find it! fraud detection contest report," in *International Conference on Pattern Recognition (ICPR)*, 2018, pp. 13–18.
- [4] J. van Beusekom, F. Shafait, and T. M. Breuel, "Text-line examination for document forgery detection," *International Journal on Document Analysis and Recognition (IJDR)*, vol. 16, no. 2, pp. 189–207, 2013.
- [5] R. Bertrand, P. Gomez-Krämer, O. R. Terrades, P. Franco, and J.-M. Ogier, "A System Based On Intrinsic Features For Fraudulent Document Detection," in *International Conference on Document Analysis and Recognition (ICDAR)*, vol. 12, Washington, DC, United States, 2013, pp. 106–110.
- [6] R. Bertrand, O. R. Terrades, P. Gomez-Krämer, P. Franco, and J.-M. Ogier, "A conditional random field model for font forgery detection," in *International Conference on Document Analysis and Recognition (ICDAR)*, 2015, pp. 576–580.
- [7] F. Cruz, N. Sidère, M. Coustaty, V. P. D'Andecy, and J.-M. Ogier, "Local binary patterns for document forgery detection," in *International Conference on Document Analysis and Recognition (ICDAR)*, 2017, pp. 1223–1228.
- [8] J. van Beusekom, F. Shafait, and T. M. Breuel, "Automatic authentication of color laser print-outs using machine identification codes," *Pattern Analysis and Applications*, vol. 16, no. 4, pp. 663–678, 2013.
- [9] S. Shang, N. Memon, and X. Kong, "Detecting documents forged by printing and copying," *EURASIP Journal on Advances in Signal Processing*, vol. 2014, no. 1, p. 140, 2014.
- [10] Z. Khan, F. Shafait, and A. Mian, "Hyperspectral imaging for ink mismatch detection," in *2International Conference on Document Analysis and Recognition (ICDAR)*, 2013, pp. 877–881.
- [11] D. M. Shimizu and H. Y. Kim, "Perceptual hashing for hardcopy document authentication using morphological segmentation," in *International Symposium on Mathematical Morphology*, 2007, pp. 77–78.
- [12] R. Villán, S. Voloshynovskiy, O. Koval, F. Deguillaume, and T. Pun, "Tamper-proofing of electronic and printed text documents via robust hashing and data-hiding," in *Security, Steganography, and Watermarking of Multimedia Contents*, 2007, p. 65051T.
- [13] L. Tan, X. Sun, Z. Zhou, and W. Zhang, "Perceptual text image hashing based on shape recognition," *Advances in Information Sciences and Service Sciences (AISSS)*, vol. 3, no. 8, pp. 1–7, 2011.
- [14] A. Malvido García, "Secure Imprint Generated for Paper Documents (SIGNED)," Bit Oceans, Tech. Rep. December 2010, 2013.
- [15] R. Brugger, A. Zramdini, and R. Ingold, "Modeling documents for structure recognition using generalized n-grams," in *Proceedings of the Fourth International Conference on Document Analysis and Recognition*, 1997.
- [16] A. Soffer, "Image categorization using texture features," in *Proceedings of the Fourth International Conference on Document Analysis and Recognition*, 1997.
- [17] A. Z. A. C. Lucia Noce, Ignazio Gallo, "Embedded textual content for document image classification with convolutional neural networks," in *Symposium on Document Engineering (DocEng)*, 2016.
- [18] E. Appiani, F. Cesarini, A. Colla, M. Diligenti, M. Gori, S. Marinai, and G. Soda, "Automatic document classification and indexing in high-volume applications," *International Journal on Document Analysis and Recognition*, vol. 4, no. 2, pp. 69–83, 2001.
- [19] M. G. M. Diligenti, P. Frasconi, "Hidden tree markov models for document image classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, pp. 519–523, 2003.
- [20] P. Sarkar, "Learning image anchor templates for document classification and data extraction," in *International Conference on Pattern Recognition (ICPR)*, 2010.
- [21] H. Peng, F. Long, Z. Chi, and W.-C. Siu, "Document image template matching based on component block list," *Pattern Recognition Letters*, vol. 22, pp. 1033–1042, 2001.
- [22] F. Cesarini, M. Gori, S. Marinai, and G. Soda, "Informys: a flexible invoice-like form-reader system," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 730–745, 1998.
- [23] D. E. A. S. M. B. C. W. Daniel Schuster, Klemens Muthmann, "Intellix – end-user trained information extraction for document archiving," in *International Conference on Document Analysis and Recognition (ICDAR)*, 2013.
- [24] V. P. D'Andecy, A. Joseph, and J.-M. Ogier, "Indus: Incremental document understanding system focus on document classification," in *International Workshop on Document Analysis Systems (DAS)*, 2018.
- [25] F. L. Rasmus Berg Palm, Ole Winther, "Cloudscan - a configuration-free invoice analysis system using recurrent neural networks," in *International Conference on Document Analysis and Recognition (ICDAR)*, 2017.
- [26] X. W. Xiaohui Zhao, Zhuo Wu, "Cutie: Learning to understand documents with convolutional universal text information extractor," 2019. [Online]. Available: <https://arxiv.org/abs/1903.12363>
- [27] I. J. Vergara-Laurens, L. G. Jaimes, and M. A. Labrador, "Privacy-preserving mechanisms for crowdsensing: Survey and research challenges," *IEEE Internet of Things Journal*, vol. 4, pp. 855–869, 2017.
- [28] I. Wagner and D. Eckhoff, "Technical privacy metrics: a systematic survey," *ACM Computing Surveys*, vol. 51, 2018.
- [29] F. Du Pin Calmon and N. Fawaz, "Privacy against statistical inference," in *50th Annual Allerton Conference on Communication, Control, and Computing*, pp. 1401–1408, 2012.
- [30] S. Alboaie, L. Alboaie, and A. Panu, "Levels of privacy for ehealth systems in the cloud era," in *24th International Conference on Information Systems Development (ISD2015 HARBIN)*, 2015.
- [31] N. Tamani and Y. Ghamri-Doudane, "Towards a user privacy preservation system for iot environments: a habit-based approach," in *International Conference on Fuzzy Systems (FUZZ-IEEE)*, 2016, pp. 2425–2432.
- [32] D. Hu, X. Hu, W. Jiang, S. Zheng, and Z. qiu Zhao, "Intelligent digital image firewall system for filtering privacy or sensitive images," *Cognitive Systems Research*, vol. 53, pp. 85 – 97, 2019, advanced Intelligent Computing.
- [33] X. Zhai, K. Appiah, S. Ehsan, G. Howells, H. Hu, D. Gu, and K. D. McDonald-Maier, "A method for detecting abnormal program behavior on embedded devices," *IEEE Transactions on Information Forensics and Security*, vol. 10 (8), 2015.
- [34] S. Eskenazi, P. Gomez-Krämer, and J.-M. Ogier, "When document security brings new challenges to document analysis," in *International Workshop on Computational Forensics (IWCF)*, 2015, pp. 104–116.
- [35] S. Eskenazi, B. Bodin, P. Gomez-Krämer, and J.-M. Ogier, "A perceptual image hashing algorithm for hybrid document security," in *International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 741–746.
- [36] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2016.