

# A Novel Approach to Unsupervised Automated Extraction of Standard Cell Library for Reverse Engineering and Hardware Assurance

**Ronald Wilson, Rabin Y. Acharya, Domenic Forte, Navid Asadizanjani and Damon Woodard**  
Florida Institute for Cyber Security, University of Florida, Gainesville  
{ronaldwilson, rabin.acharya}@ufl.edu, {dforte, nasadi, dwoodard}@ece.ufl.edu

## Abstract

Reverse engineering today is supported by several tools, such as ICWorks, that assist in the processing and extraction of logic elements from high definition layer by layer images of integrated circuits. To the best of our knowledge, they all work under the assumption that the standard cell library used in the design process of the integrated circuit is available. However, in situations where reverse engineering is done on commercial off-the-shelf components, this information is not available thereby, rendering the assumption invalid. Until now, this problem has not been addressed. In this paper, we introduce a novel approach for the extraction of standard cell library using the contact layer from these images. The approach is completely automated and does not require any prior knowledge on the construction or layout of the target semiconductor integrated circuit. The performance of the approach is evaluated on two AES designs with 10,000 cells compiled from standard libraries with 32nm and 90nm node technologies having 350 and 340 standard cells respectively. We were able to successfully extract 94% and 60% of the standard cells from the 32nm and 90nm AES designs using the proposed approach. We also perform a case study using a real-world sample extracted from a smartcard. Finally, we also investigate the various challenges involved in the extraction of standard cells from images and the steps involved in resolving them.

## Introduction

Reverse engineering is essentially defined as the in-depth analysis of an end product to reveal its functionality and the techniques involved in its manufacturing. This information is exploited in several ways ranging from re-manufacturing of legacy devices to defending Intellectual Property (IP) rights. In case of semiconductor Integrated Circuits (IC), the reverse engineering (RE) process requires three major steps: identifying the device technology from images of the chip taken layer-by-layer using scanning electron microscope (SEM), extracting its gate-level netlist and inferring its functionality [1]. The first phase of the reverse engineering paradigm involves imaging to recover the node technology employed in its design and obtain a heuristic on the design rules used to produce the IC. The second step involves matching the patterns found in the active, polysilicon, contact and metal 1 layers to a standard cell library and extraction of the netlist. This step also involves examination of the remaining metal layers which show how the cells are

connected to each other in the design. The final step infers the functionality of the circuit or a sub-circuit in the design. There are a plethora of algorithms [2,3,4,5] and automated tools [6,7,8,9] that assist in the identification and extraction of logic elements in the circuit. RE on gate design needs active, contact and metal 1 layers [10]. Inferring cell functionality typically only requires the metal 1 layer [11].

Even though the paradigm is relatively straight forward and simple, there are several limitations that makes RE hard to execute. One such limitation is the requirement of a template standard cell library to be present to extract the gates, filler, and flip flops in the netlist. To the best of our knowledge, all existing algorithms and tools require some heuristic on the standard cell library to work. With the goal of RE being the extraction of information on the product with no prior knowledge, this limitation is counterintuitive. In this paper, we propose a novel method to address this issue.

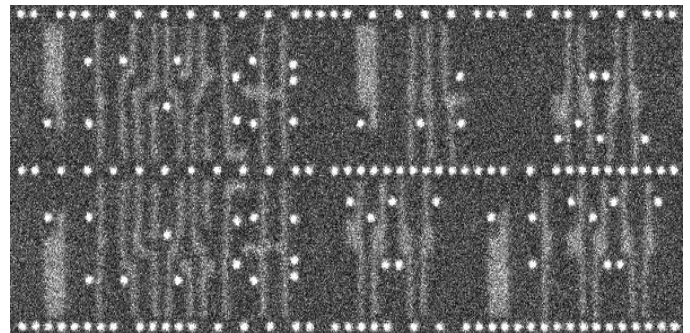


Figure 1: SEM image showing noise in the polysilicon layer

There are several problems associated with the extraction of standard cells from a given device. Some of them are described below in detail:

- *Image acquisition*: In some layers, especially polysilicon, the noise intensity is too high (Fig. 1). This can only be counteracted by increasing the quality of the scan. In a typical imaging modality, like the SEM, a higher quality scan would extend the imaging time frame in the order of several days [12], and even months for advanced technology nodes.
- *Alignment and Stitching*: Due to nature of the imaging modalities used and the scale at which images are acquired (mosaicking), the layers may be misaligned [9].
- *Intra-cell Similarity*: In some standard cell libraries, there are extensive similarities between different cells (Fig. 2). This inherent similarity might cause traditional machine learning algorithms to produce inaccurate results.

- *Scale*: With ever increasing levels of integration in IC, the memory and time complexity involved in processing the amount of data extracted from one IC is nearly intractable. This holds true, especially, when each potential cell pattern has to be compared against all the cells in the library.

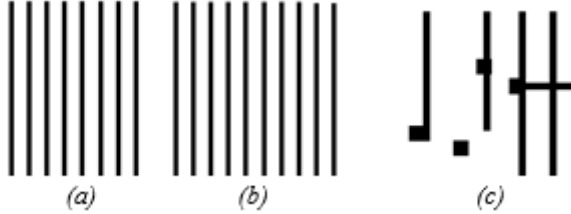


Figure 2: (a, b) Cells showing similarity in polysilicon layer on the 32nm dataset. (c) Complicated shapes in poly-silicon layer on the 90nm dataset.

Our contribution in this paper is two-fold. We investigate the similarities between multiple cells from two standard libraries and discuss the infeasibility of machine learning based methods on unsupervised automated cell extraction. Furthermore, we introduce a novel rule-based method that can successfully extract the standard cell library from just the contact layer in the IC. The results and limitations of our approach are also discussed.

## Dataset

With the strict control on the availability of industry-use standard libraries, we have based our experiments on two standard cell libraries that can be readily reproduced: 32/28nm Educational Design Kit [13] and Synopsys open educational design kit [14]. These libraries contain 350 and 340 standard cells corresponding to 32nm and 90nm node technologies respectively. Using the libraries, an AES design was synthesized into approximately 10,000 standard cells from which we extract the contact layer as a segmented image. A small section of the designs is shown in Fig. 3.

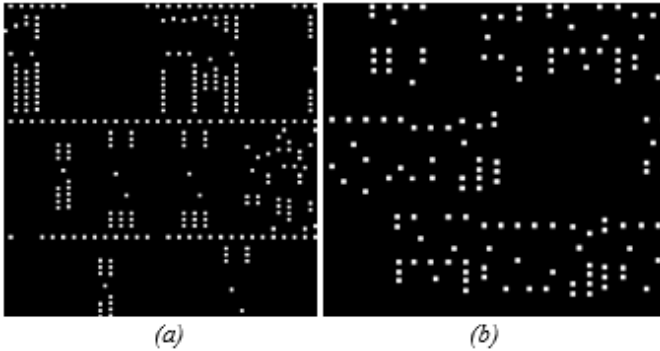


Figure 3: (a) AES design from 32nm dataset (b) AES design from 90nm dataset

In addition, a separability test was performed on all the layers associated with a standard cell for both datasets. The goal of the test is to evaluate the extent of similarity between two different cells from the same standard cell library. For layers other than the contact layer, the images of layers were

convolved against each other and the point of maximum correlation was recorded. For the contact layer, the N-gram length was used. It is defined as the count of consecutive columns of contacts in a cell. For instance, if the n-gram length is 2, every two consecutive columns in the source cell is compared against every two consecutive columns in the target cell and the number of overlapping contacts is recorded. However, only the maximum number of overlapping contacts between two cells for the chosen n-gram length is kept for the final analysis. If the source or target cell does not have enough contact columns to satisfy the chosen n-gram count, the cells are ignored from the analysis. In Fig. 4 and 5, the correlation of patterns for all cells with each other have been plotted.

As expected, the general trend for both plots are monotonically decreasing. This can be attributed to the fact that the number of possible combinations for each n-gram increases with the considered length, making them more unique in the space of all possible combinations. Another interesting point is the major difference in correlation values between the 90nm and 32nm datasets. The reason for this drastic difference can be observed in Fig. 6. The contacts in the 32nm dataset follow a grid-like arrangement with a constant distance between each column. However, the distance between two adjacent columns in the 90nm dataset is not constant. This additional freedom of choosing the distance between two columns make the combination more unique.

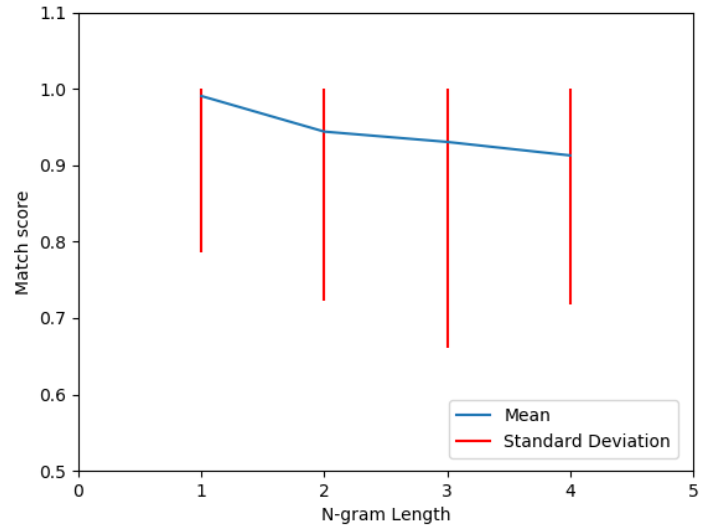


Figure 4: Correlation of patterns in various n-gram length for the 32nm dataset

The raw correlation values for the other layers can be seen in Fig. 7. As expected from the examples shown in Fig. 2, the correlation for the polysilicon layer is much higher than any other layer for the 32nm dataset. The 90nm dataset shows similar correlation values for all layers. The observation obtained from Fig. 7, supports the fact that active (diffusion), contact and metal 1 layers can be used to successfully reverse engineer the IC and infer its functionality if the templates of the standard cells are known. However, if it's not available, the correlation within different layers also makes the features extracted inherently correlated and inseparable for any machine learning algorithm. With the primary goal for the

application of RE being the accurate and complete reconstruction of features inside the IC, the correlation between features must be avoided or dealt with in another way.

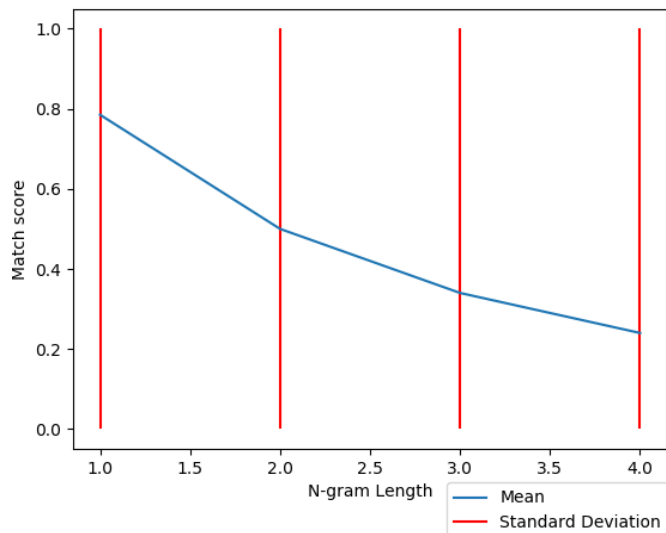


Figure 5: Correlation of patterns in various  $n$ -gram length for the 90nm dataset

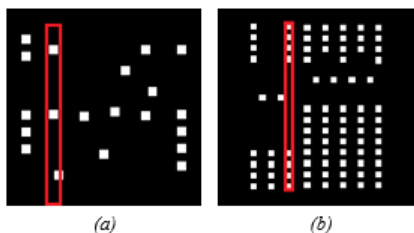


Figure 6: (a) Non grid-like allocation of contacts on the 90nm dataset in cell AND3x1. (b) Grid-like allocation of contacts on the 32nm dataset in cell AND2x4.

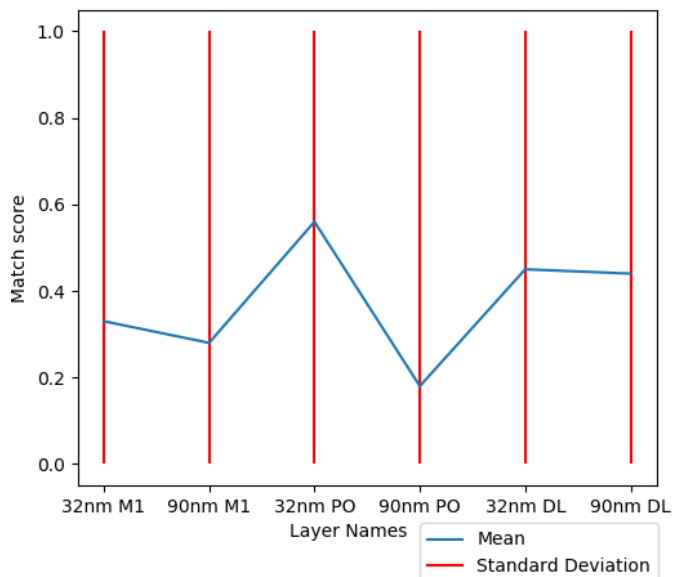


Figure 7: Simple correlation plot for various layers in the standard cell library. (PO: Polysilicon, M1: Metal 1, DL: Diffusion/Active layers)

## Proposed Approach

The first step in the extraction of logic cells from the contact layer is the detection of boundary between two cells. Typically, this can be done by considering the amount of free silicon substrate between the columns of contacts. If the space occupied by the silicon substrate is greater than a certain threshold, the contact can be considered as belonging to two different cells. There are two major flaws to this line of thought. Firstly, some cells do have unusually long space between adjacent contact columns and splitting them would cause over-segmentation in the cell. Secondly, this approach being well known, has prompted the development of several anti-RE approaches. One of them being the insertion of dummy but operational structures into these empty silicon substrate locations [15]. This makes the boundary detection solely based on the separation between contact columns invalid.

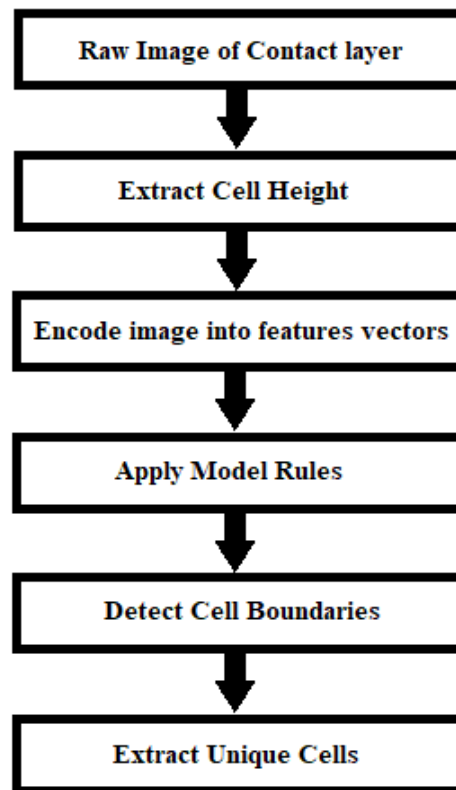


Figure 8: Framework for the proposed algorithm

Our approach accomplishes the extraction of cell boundaries by encoding the contacts into a vector and looking for unique patterns in them. The framework of the proposed algorithm is shown in Fig. 8 and the major steps discussed below.

### Feature Encoding

There are two drawbacks to working on an image. The space/time complexity involved in going through the image and the noise introduced by the imaging modality. The space-time complexity of the algorithm can be considerably reduced by encoding the image into simple numerical vectors.

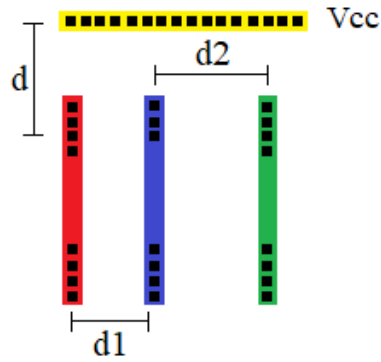


Figure 9: Illustration of features encoded by the algorithm

Each contact column in the AES image, shown in Fig. 3, is encoded into a numerical vector using the distance between each contact and the  $V_{cc}$  line. The distance is indicated as ‘d’ in Fig. 9. For instance, the red contact column in Fig. 9, would have 8 values indicating the distance between the centroid of each of its 8 contacts and the centroid of the  $V_{cc}$  line highlighted in yellow. The  $V_{cc}$  line can be found either on the contact layer or on metal 1. This should alleviate the misalignment to a certain extent since the patterns are encoded using a local point of reference. Following the encoding, every three adjacent contact columns in the image are extracted and saved along with the distances between them. This is depicted in Fig. 9 where the distances between the columns are denoted as  $d_1$  and  $d_2$ . Even though the features are more unique with longer sequences of columns, we have limited the count to 3. This is done to ensure that standard cells having just one contact column can be detected. If the minimum number of columns in a cell in the standard cell library is known beforehand, this parameter can be changed to that value. If the parameter is set to a value larger than the minimum sequence length of contact columns, then all such cells will not be detected.

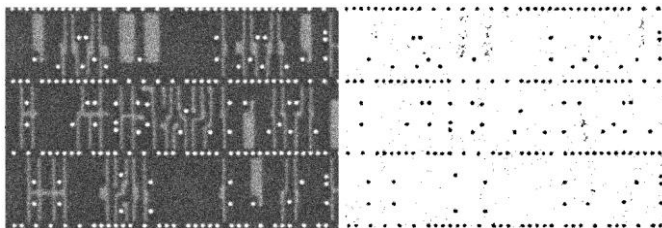


Figure 10: Illustration of contacts extracted from a raw SEM image using simple thresholding

In a typical imaging modality used in RE, like the SEM, the contrast for metallic structures are much higher than non-metals. This helps in the extraction of features even on noisy low-quality images. This is our motivation for considering just the contact layer for the standard library extraction. In Fig. 10, a simple thresholding method was used to extract the contacts from the SEM image. Any pixel having intensity value less than 250 was filtered out.

### Automated Extraction of Cell Height

The  $V_{cc}$  lines were extracted from the AES design using a simple algorithm. It is to be noted that there exists a minimum

distance between two contacts on the  $V_{cc}$  line and distance between two  $V_{cc}$  lines is always constant. We have exploited this information to localize them. Firstly, we measured the distance between contacts in the same row and calculated its median value. Once this is done for every row in the design, the minimum value is taken. Secondly, every row that has median separation equal to the minimum value is marked. The respective plot is shown in Fig. 11. The periodic nature associated with the constant separation between  $V_{cc}$  lines can also be observed in Fig. 11. Since the separation between the lines can be taken as an integral multiple of an arbitrary integer, an arithmetic progression that best describes the periodic nature is fitted to the plot. The common difference of the arithmetic progression defines the cell height in the design. The result of applying this method is shown in Fig. 12.

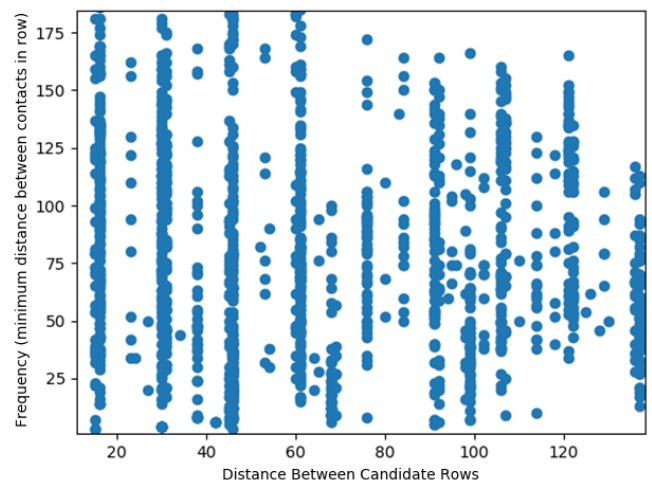


Figure 11: Plot depicting the period nature of  $V_{cc}$  contact lines

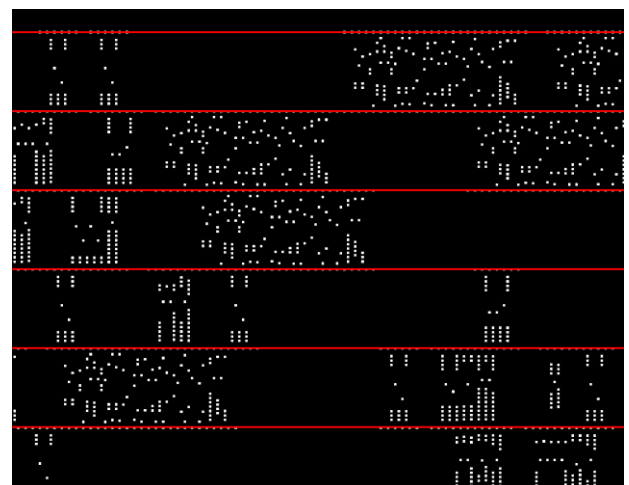


Figure 12: Image showing the  $V_{cc}$  contact lines. Detected lines are shown in red

### Model Rules

The extraction of the standard library requires the application of two rules. They are described below:

- *The inter-contact distance rule:* For any given cell, the distance between its respective contact columns would be the same, i.e., the  $d_1$  and  $d_2$  given in Fig. 9 would be the same if the contact appears inside the cell. However, if the

columns appear on the boundary of the cell, the distance would vary randomly. The information in this rule is captured by the feature encoding.

- *The start-end rule:* The ending boundary of one cell is always followed by the starting boundary of another cell. A contact column that doesn't belong to a cell cannot exist. Any detected boundary that does not satisfy this rule will be ignored.

Since most imaging modalities used in RE captures the IC layers in smaller images and stitches together to form a mosaic, there are concerns for the detection of partial cells in the contact layer. This problem is resolved by the start-end rule since partial cells has either a beginning or an end and not both. Partial cells detected at the edges of the image can be ignored or flagged for further processing later.

### Segmentation

Once the boundary features are detected, the cells can be extracted from the image and added to a candidate library. However, with the application of start-end rule and presence of noise in the image, some under-segmentation of the cell might have occurred. To remedy this, the cells are compared against each other. If a given cell can be expressed as the sum of 'n' cells present in the library in any given order, they are segmented.

## Results

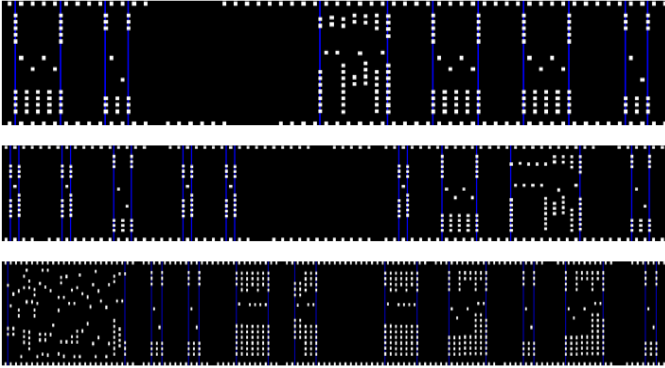


Figure 13: Cells extracted by the algorithm. Blue lines mark the boundary of the cell

Table 1: Summary of the results obtained from applying the algorithm to the AES design

Node Technology	32 nm	90 nm
No. of Unique Cells	64	87
True Match	60 (94%)	27 (31%)
Partial Match	0 (0%)	33 (38%)
Failed Match	4 (6%)	27 (31%)

The image shown in Fig. 13 shows the result of applying the given rules to the segmented image. The boundaries extracted by algorithm are highlighted in blue. Table 1 summarizes the results of the experiment. In Table 1, True Match corresponds to the extraction of the complete cell when compared to ground truth. Partial Matches happen when compound cells

are detected. A compound cell is a combination of multiple ground truth cells. Failed Matches happen when the detected cells are over segmented versions of the ground truth cells.

## Discussion

As observed from Table 1, there is a significant deviation between the cell extraction between 32nm and 90nm node technologies. The reason for this discrepancy can be seen in Fig. 6. The segmentation and extraction of cells is solely depended on the encoded features. Since the AES design based on the 90 nm node technology does not follow a grid-like arrangement of contacts, it is difficult to encode them into a discretized feature space. Depending on the noise in the image, the floating contact in Fig. 6 can be either included in the current column or one of the adjacent columns. This uncertainty in its absolute position causes the encoding to change at times. Although inconvenient, the encoding is necessary to reduce the complexity of the algorithm and enable it to scale to present day billion transistor ICs.

The partial and failed matches are caused either due to over-segmentation or under-segmentation of standard cells. In the AES design, there were situations where two standard cells always occurred together. With the lack of variance in distance measures for the encoded features for that specific standard cell, they were grouped together causing a partial match. Segmentation of cell only happen when evidence of variation in distance in feature encoding is detected, i.e., d1 or d2 must be different.

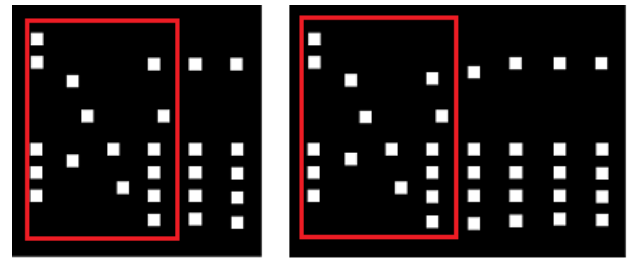


Figure 14: Image showing subset of similar sequences (highlighted in red) between two different cells (a) AND2x2 (b) AND2x4

Failed Matches happen when a subset of contact columns repeat across multiple standard cells or when variation in encoded features occur due to floating contacts as shown in Fig. 6(a). This can be easily remedied in design that follow a strict grid-like allocation of contacts, thus, explaining the higher rate of failed matches in the 90nm AES design. Fig. 14 shows an instance where the algorithm caused a failed match due to over-segmentation.

## Case Study

A case study was performed on a real-world sample extracted from a Smartcard using 190nm node technology. The image was acquired using the TESCAN FERA3 dual-beam FIB-SEM with imaging potential set at 12 kV for higher



penetration. Delayering was performed from the backside using the approach described in [18]. Imaging at a higher potential also helped reduce the influence of uneven delayering. The acquired image is shown in Fig. 15. The other imaging parameters such as magnification, dwelling time and resolution were set to 100 $\mu$ m, 32 $\mu$ s/pixel and 1024x1024 pixels respectively.

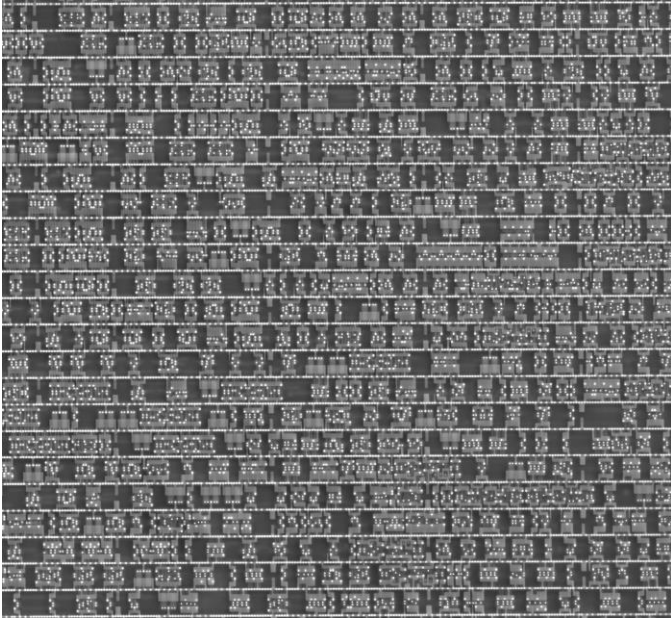


Figure 15: A Raw SEM image of the contact layer

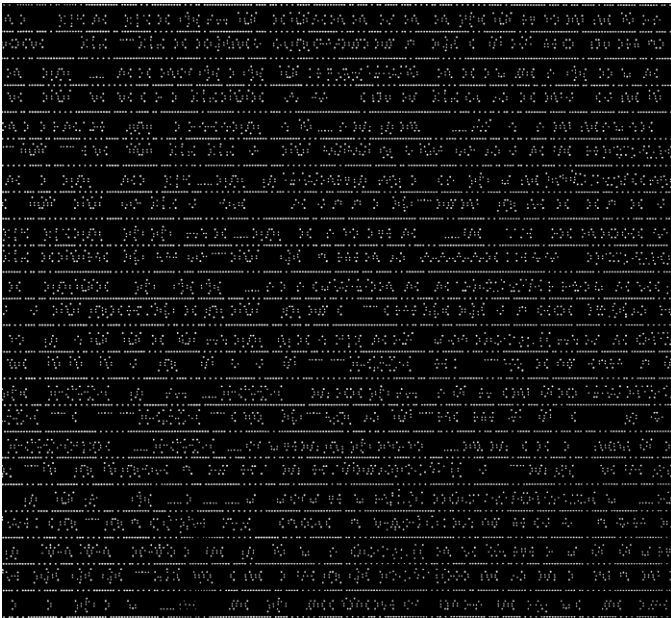


Figure 16: Binarized image of the contact layer with threshold set at 225. All structures except the contacts are removed after thresholding

Following acquisition, the image was binarized. Since the contacts in the image are much brighter than the surrounding structures, a simple thresholding was used to separate out the contacts. No specific steps were taken to account for the variation in size of the contacts due to depth. Any pixel in the

original image with value over 225 is taken as a pixel belonging to a contact. The rest of the pixels are set to zero. The resultant binarized image is shown in Fig. 16. Following the sequence of steps discussed in this paper, the image was segmented into strips after detecting the  $V_{cc}$  power supply lines. Some of the strips, after the extraction, is shown in Fig. 17. Finally, the rules were applied to the strip to extract the standard cells. Some of the standard cells extracted are shown in Fig. 18. The functionality of the cells was identified with the help of a Subject Matter Expert.



Figure 17: Image showing strips of contact extracted from the binarized image. The  $V_{cc}$  supply lines in the image are removed after detection

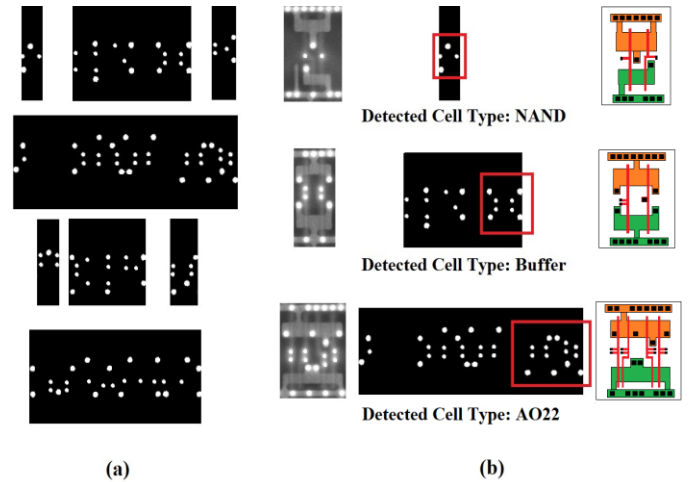


Figure 18: (a) Some of the cells detected by the algorithm. The conjoined cells are the result of limited data provided to the algorithm in conjunction with the non-grid like arrangement of contacts (b) Cells detected by the algorithm with their functionality decoded by a Subject Matter Expert (SME). True individual cells are highlighted in red bounding boxes.

## Conclusion

In this paper, we were able to resolve a problem that has not been addressed in RE and hardware assurance literatures- the extraction of standard cell libraries from images of IC layers. With most of the existing algorithms concentrating on matching extracted features to templates in libraries, the extent to which effective RE can be performed was limited. We have extended the scope of RE to ICs of completely unknown architectures and provided a new avenue for RE research.

In conclusion, we were able to achieve our goal successfully along with the following advantages:

- Extraction of the library was done only using information from the contact layer. The features in this layer are considerably more resistant to noise than other layers [16].
- Using local landmarks for extracting features, we were able to handle alignment problems in image mosaics up to a certain extent.
- With the limited possible combination of contacts in the layer and the encoding of images into feature space, the number of features extracted are exponentially lower than the number of gates in the IC thereby reducing memory and time complexity.
- The approach is completely automated and does not require any prior knowledge on the design rules of the target IC. The only human interaction required is to confirm the validity of the extracted cells. By automating the process, we have also taken out biases and incorrect annotation introduced by human factors.

In a future work, we will also be looking into improving the cell extraction by incorporating information from other layers. This information will also be used to aid in automatically identifying the cell's functionality.

## References

- [1] Rostami, M., Koushanfar, F. and Karri, R., "A primer on hardware security: Models, methods, and metrics". Proceedings of the IEEE, 102(8), 2014, pp: 1283-1295.
- [2] Mark C. Hansen, Hakan Yalcin, and John P. Hayes. 1999. Unveiling the ISCAS-85 benchmarks: A case study in reverse engineering. IEEE Design and Test of Computers 16, 3, 72–80.
- [3] W. Li, A. Gascon, P. Subramanyan, W. Y. Tan, A. Tiwari, S. Malik, N. Shankar, and S. A. Seshia. 2013. WordRev: Finding word-level structures in a sea of bit-level gates. In Proceedings of the 2013 IEEE ACM Journal on Emerging Technologies in Computing Systems, Vol. 13, No. 1, Article 6, Publication date: April 2016. 6:32
- [4] W. Li, Z. Wasson, and S. A. Seshia. 2012. Reverse engineering circuits using behavioral pattern mining. In Proceedings of the 2012 IEEE International Symposium on Hardware-Oriented Security and Trust (HOST'12). IEEE, Los Alamitos, CA, 83–88.
- [5] P. Subramanyan, N. Tsiskaridze, K. Pasricha, D. Reisman, A. Susnea, and S. Malik. 2013. Reverse engineering digital circuits using functional analysis. In Proceedings of the Conference on Design, Automation, and Test in Europe. 1277–1280.
- [6] Torrance, R., & James, D. (2009, September). The state-of-the-art in IC reverse engineering. In International Workshop on Cryptographic Hardware and Embedded Systems (pp. 363-381). Springer, Berlin, Heidelberg.
- [7] Kenneth, K. Yu, and C. Neil Berglund, U.S. Patent No. 5,086,477
- [8] Ahmed, Haroon., Blythe, Simon and Fraboni, Beatrice, U.S. Patent No. 5,191,213.
- [9] Zavadsky, Vyacheslav L., Val Gont, Edward Keyes, Jason Abt, and Stephen Begg. U.S. Patent 7,580,557.
- [10] Avery, L.R., Crabbe, J. S., Al Sofi, S., Ahmed, H., Cleaver, J. R. A., Weaver, D. J. 2002. Reverse Engineering Complex Application-Specific Integrated Circuits (ASICs), Proceedings of the Diminishing Manufacturing Sources and Material Shortages Conference (Mar. 2002).
- [11] Nohl, K, Evans, D., Plotz, S., Plotz, H. 2008. Reverse-Engineering a Cryptographic RFID Tag, USENIX Security Symposium (Jul. 31 2008).
- [12] N. Vashistha, M T Rahman, H. Shen, D L Woodard, N Asadizanjani and M. Tehranipoor, Detecting Hardware Trojans Inserted by Untrusted Foundry using Physical Inspection and Advanced Image Processing Techniques, Publication pending in Spring Journal of Hardware System and Security (HaSS), Tentatively December 1, 2018.
- [13] Goldman, R., Bartleson, K., Wood, T., Kranen, K., Melikyan, V., & Babayan, E. (2013, December). 32/28nm educational design kit: Capabilities, deployment and future. In 2013 IEEE Asia Pacific Conference on Postgraduate Research in Microelectronics and Electronics (PrimeAsia) (pp. 284-288). IEEE.
- [14] Goldman, R., Bartleson, K., Wood, T., Kranen, K., Cao, C., Melikyan, V., & Markosyan, G. (2009, July). Synopsys' open educational design kit: capabilities, deployment and future. In 2009 IEEE International Conference on Microelectronic Systems Education (pp. 20-24). IEEE.
- [15] Cocchi, R. P., Baukus, J. P., Chow, L. W., & Wang, B. J. (2014, June). Circuit camouflage integration for hardware IP protection. In Proceedings of the 51st Annual Design Automation Conference (pp. 1-5). ACM.
- [16] Blythe, S., Fraboni, B., Lall, S., Ahmed, H. and de Riu, U., 1993. Layout reconstruction of complex silicon chips. IEEE journal of solid state circuits, 28(2), pp.138-145.
- [17] Quadir, S.E., Chen, J., Forte, D., Asadizanjani, N., Shahbazmohamadi, S., Wang, L., Chandy, J. and Tehranipoor, M., 2016. A survey on chip to system reverse engineering. ACM journal on emerging technologies in computing systems (JETC), 13(1), p.6.
- [18] Principe, E.L., Asadizanjani, N., Forte, D., Tehranipoor, M., Chivas, R., DiBattista, M., Silverman, S., Marsh, M., Piche, N. and Mastovich, J., 2017, December. Steps toward automated deprocessing of integrated circuits. In ISTFA 2017: Proceedings from the 43rd International Symposium for Testing and Failure Analysis (p. 285). ASM International.