# The ASAR 2018 Competition on Physical Layout Analysis of Scanned Arabic Books (PLA-SAB 2018)

Randa Elanwar and Margrit Betke

*Abstract*—**Successful physical layout analysis (PLA) is a key factor in the performance of text recognizers and many other applications. PLA solutions for scanned Arabic documents are few and difficult to compare due to differences in methods, data, and evaluation metrics. To help evaluate the performance of recent Arabic PLA solutions, the ASAR 2018 Competition on Physical Layout Analysis (PLA) was organized. This paper presents the results of this competition. The competition focused on analyzing layouts for Arabic scanned book pages (SAB). PLA-SAB required solutions of two tasks: page-to-block segmentation and block text/non-text classification. In this paper we briefly describe the methods provided by participating teams, present their results for both tasks using the BCE-Arabic benchmarking dataset [1], and make an open call for continuous participation outside the context of ASAR 2018.**

*Index Terms*—**Arabic document, benchmarking dataset, block classification, layout analysis, page segmentation, scanned PDF.**

## 1. Introduction

Physical layout analysis (PLA) of a scanned document is the task of segmenting the layout of the document image and identifying the class to which each image region belongs without using text recognizers or human supervision. Successful physical layout analysis leads to more accurate performance of optical character recognition (OCR) engines and tools for document search and retrieval, word spotting, PDF-to-Word conversion, etc.

Researchers addressing document layout analysis (DLA) problems with regards to Arabic documents have difficulty comparing their system results to other researchers' results. The difficulty comes from the lack of benchmarking datasets that include Arabic documents in general, and datasets that include Arabic layout annotation information specifically. Even limited-size datasets were not publicly available until recently, when BCE-Arabic-v1 [1] was published. Researchers could only report their DLA system results on proprietary datasets containing a few tens to hundreds of images, or customized subsets of data that they chose from public datasets, matching their research needs (e.g., [2], [3]). Consequently, in the PLA literature, researchers have reported a wide variety of methods, datasets, evaluation metrics, and desired targets for system output.

The difficulties in comparing system performance are particularly noticeable when both tasks of PLA are addressed, (1) segmenting the image into a set of nonoverlapping homogeneous blocks and (2) labeling each block according to its content class, for example, most simply, "text" or "non-text." Performance of a newly proposed system on the segmentation task could only be compared to that of basic algorithms. While such a comparison might be acceptable for segmentation evaluation, it is less helpful when classification evaluation is also needed. PLA solutions in the literature are typically not accompanied by published source code and reimplementing them is a nontrivial task, hindered by the fact that implementation details are also often not published. A consequence of all these difficulties is that only a small number of publications address DLA problems, and researchers have shifted to other research areas where data is available.

With the publication of the BCE-Arabic dataset [1], the situation changed. In BCE–Arabic–V1, a dataset of 1,850 document images of Arabic book pages were made available to the community. Creation of a second dataset, BCE–Arabic–V2, taken from 700 different Arabic books, made it possible for us to design a benchmarking competition on physical layout analysis for Arabic documents as part of the 2nd IEEE International Workshop on Arabic and derived Script Analysis and Recognition (ASAR 2018). Our goal was to promote Arabic-DLA research, open channels of cooperation, and enable networking between Arabic-DLA community members from different countries. We hope that the competition encourages the establishment of competitions about other DLA problems and accelerates research towards solutions of these problems.

The rest of this paper is organized as follows: Section 2 describes the competition rules and timeline. Section 3 describes our benchmarking dataset and the format of the ground truth. The performance evaluation metrics and code are described in Section 4. A summary of participating methods and their results are given in Sections 5 and 6, respectively. Conclusions are provided in Section 7.

- *Randa Elanwar is with the Computers and Systems Department, Electronics Research Institute, Egypt. E-mail: randa.elanwar@eri.sci.eg*
- *Margrit Betke is with Department of Computer Science, Boston University, USA. E-mail: betke@bu.edu*

## 2. The Competition

The competition name PLA–SAB 2018 stands for **P**hysical **L**ayout **A**nalysis of **S**canned **A**rabic **B**ook pages. The competition was announced to compare state-of-the-art solutions to a specific layout analysis problem with regards to certain document domain and/or content. This competition should only be a first "edition" or "round," and help researchers understand promising techniques or missed opportunities. The hope is that additional rounds of the competition will be organized that provide further insights. This round of the competition was organized by members of a joint team from Boston University USA, and Electronics Research Institute, Egypt. Competition details can be found at http://www.cs.bu.edu/faculty/betke/ASARLayout-AnalysisCompetition.

The competition provided a the following:
(1) A **benchmarking dataset** for testing physical layout analysis solutions, which contains an annotated test set of scanned Arabic book page samples with a wide variety of content and appearance, and a
(2) A **full evaluation scheme** by offering code to compute a set of evaluation metrics for quantitative evaluation of the competition tasks.

The competition involved two tasks (see Fig. 1):

**Task 1**: Provide a block-level segmentation to the benchmarking data, i.e., bounding box coordinates information of each block, and

**Task 2**: Identify the block type as text or non-text.

The challenge was open (not "blind") in the sense that participants could download the data, run their own algorithms, and provide the organizers with their results directly instead of sending their executable programs for the organizers to run.

The competition was announced on the ASAR website on January 1, 2018. A two-week registration period followed. During this period, teams could either register new systems, which had not been previously published, or an existing system, which had been published. On the registration form, participating teams committed not to modify their system or manipulate the results during the testing period, and we trusted their academic integrity when reporting final results. The benchmarking dataset and corresponding ground truth were made available on January 16, 2018, giving participants a three-week testing period to run their system on test images, obtain a ground-truth comparison, and submit a paper on a new system or a report on a previously published system with the challenge results included.

A total of 15 teams from six countries answered the registration call, out of which 12 registered as competitors and three were interested in obtaining the competition dataset. Three submissions for the challenge were eventually received.



Figure 1. Visual representation of physical layout analysis tasks: Page-to-block segmentation and classifying each block according to its content (text or non-text).

## 3. The Competition Dataset

The competition dataset was chosen from the dataset BCE-Arabic–V2, which was provided by the second phase of BCE-Arabic data collection project and can be found at *http://www.cs.bu.edu/faculty/betke/BCE*. BCE-Arabic–V2 consists of book pages of various contents and have layout problems that involve dealing with decorative backgrounds, pages produced by different printing technologies, low quality paper effects, etc. (see Figure 2).



Figure 2. Examples of challenging layouts in BCE-Arabic–v2.

The competition dataset contains 90 images provided in three equal-size sets A, B, and C according to layout content. The 90 images contain a total of 1,112 blocks, in particular, 927 text blocks, 149 image blocks, and 36 graphics or decorations blocks:

- **Set A:** 30 images of single-column layouts with plain backgrounds and rectangular block shapes (297

text blocks, 53 image blocks)

- **Set B:** 30 images of double-column "simple" layouts (379 text blocks, 54 image blocks, 19 graphics blocks)
- **Set C:** 30 images of "complex" layouts (251 text blocks, 42 image blocks, 17 graphics blocks)

The categorization into "simple" and "complex" layouts was performed by the competition organizers and is somewhat subjective.

Ground-truth annotations were performed by the competition organizers using the Alethia tool [5]. They were provided to the competition participants in PAGE XML format [6]. The selection of the Alethia tool for ground-truth annotation was informed by a careful study by Saad et al. [1], which compared several document annotation tools and found Alethia most valuable for annotating the layout of Arabic book pages.

## 4. Methodology of Performance Evaluation

Physical layout analysis involves two main tasks: page segmentation, breaking down the page image into its main components, here called "blocks," and classification of the blocks into to domainbased classes (usually 'text' and 'non-text'). Accordingly, the evaluation of the competition systems assessed both the segmentation task and the classification task.

If the analysis primitive of a system is a pixel or a connected component, the system typically addresses the classification task first, and the ground truth available is also labeled with regards to pixels or connected components. The system then performs segmentation to compute blocks, and block classification comes next. If the analysis primitive of a system is a block, it typically performs segmentation first and then classification, and the ground truth is provided per block. The competition did not restrict systems to a certain analysis primitive. For example, performing PLA by skipping the page-to-block segmentation step and working on pixels or CCs directly avoids propagating any potential segmentation errors to the classification step and may thus result in superior classification results.

As was mentioned in Section 2, the evaluation procedure of the competition was not "blind." The benchmarking dataset was sent to participants after the registration period. During the testing phase, the evaluation code was used by participants to evaluate their system performance. The segmentation results of a system must be reported as contours (mainly bounding boxes of blocks) indicated by multiple vertices, and the classification results as the text or non-text label of each block.

The evaluation code for the competition was written in C++ with Visual Studio 2013 and OpenCV 3.1. It was provided to the competition participants as an executable file. The evaluation code takes as input the original image,

the results of a participating system, and ground truth in PAGE format. The output of the evaluation code is an XML file that contains the evaluation statistics based on the metrics below.

The block segmentation evaluation was performed based on metrics inspired by the work of Shafait et al. [4], who compared the performance of six classical page segmentation algorithms.

- The over-segmentation error (OSE) counts the number of oversegmented blocks per image and divides it by the number of images in the dataset.
- The under-segmentation error (USE) counts the number of under-segmented blocks per image and divides it by the number of images in the dataset.
- The correct-segmentation (CS) metric counts the number of correctly segmented blocks per image and divides it by the number of images in the dataset.
- The missed-segmentation error (MSE) compares the number of missed blocks per image and divides it by the number of images in the dataset.
- The false alarm error (FA) counts the number of false alarms per image and divides it by the number of images in the dataset. (False alarms occur when border noise is present.)
- The overall block error rate combines the OSE, MSE, and USE metrics and compares the sum to the total number of ground truth blocks $\rho$:

$$\rho = \frac{\text{OSE} + \text{USE} + \text{MSE}}{\text{Total No. Blocks}} \qquad (1)$$

The block error rate $\rho$ can be reported per image or averaged over the images in a dataset.

The average black pixel rate (AvgBPR) is also computed. It represents the number of black pixels contained in the segmented blocks compared to the corresponding blocks in the ground truth image. This rate is computed per image and then averaged over the dataset.

Classification evaluation was performed based on the standard definitions for precision (Pr), recall (Rec), F1-measure (F1), and average class accuracy (Acc) on both pixel and block levels for text and non-text classes.

## 5. Participating Methods

We asked each participating team to provide us with a summary description of their submitted systems together with their results. For the completeness of this paper we present these summaries here. If the proposed system was new, the participating team had the option to submit a research paper to ASAR 2018, which, if accepted in peer review, appears in the same proceedings as this paper.

## 5.1. The RFAAD Method

Competition results for the "RFAAD Method" was submitted by a team from the Electronics Research Institute in Egypt, which included R. Saad, R. Elanwar, N. Kader, S. Mashali, and Boston University collaborator M. Betke. RFAAD performs layout analysis using a Random Forests classifier based on structural features from the connected components (CCs) of a document image [8].

Initially, the input gray-scale image is binarized using Otsu's method. Then a median filter is applied to the binarized image for noise removal. A sequence of morphological operations is performed to help merge broken small-size CCs with larger CCs and help distinguish non-text components from text components, as the former will appear much larger.

After preprocessing, structural features are extracted from the image CCs inspired by a method by Le et al. [7]. These are normalized CCs centers, normalized width and height of bounding boxes of CCs, elongation, solidity, log-normal distribution of height, Hu moments, mean and standard deviation of CCs, stroke width, and nearest neighbor features. The extracted features are further normalized to zero mean and unity standard deviation before training a Random Forests classifier.

RFAAD was trained and tested using BCE-Arabic-v1 dataset [1]. It was implemented using the Java-based WEKA-library and run on a core i7 PC. System evaluation results are reported on the CC level (average CC classification accuracy). Block building was also performed by grouping adjacent same-class CCs, and segmentation evaluation was reported together with pixel-level and block-level classification evaluation on the competition datasets.

## 5.2. The FCN-based Method

The "FCN-based Method" was submitted by a team from Ben-Gurion University of the Negev, Israel, which included A. Droby, B. Barakat, and J. El-Sana, and will be presented at ASAR 2018 [9]. FCN stands for "fully convolutional network." It is a deep learning system that used a FCN architecture for object segmentation, in particular, a 16-layer VGG network. The system has three output channels corresponding to the number of classes: text, non-text, and background.

The system was trained using the BCE-Arabic V1 dataset. Training and validation images were binarized. Then a random generation was performed of 100,000 and 15,000 patches of size $320 \times 320$ for training and validation sets, respectively. Training was performed with Stochastic Gradient Descent (SGD) with a momentum of 0.9 and a learning rate of 0.001. The VGG network was initialized with its publicly available pre-trained weights until least validation loss. Experiments were conducted on Keras and run on a single Nvidia 1080 GTX.

The test images of the competition dataset were processed as follows: The method first binarizes the test images and trims their margins by 3% of the rows and 10% of the columns. Horizontal lines are removed using morphological operations. During prediction of test patches, marginal regions that are less than patch size are filled with background pixels. A post-processing method is employed to denoise the results and extract well-defined classified zones. Each class (e.g. text and non-text) is considered separately:
- First, using morphological operations, small classified zones are removed and ragged edges are smoothed.
- Second, the contours of each connected component (i.e. classified zones) are extracted. Then, for each extracted contour, the following values are considered:

- AP, the total number of pixels inside the contours,
- CP, the number of classified pixels inside the contours.

Only the connected components that satisfy $AP \leq M$ and $CP/AP \leq \alpha$ are considered, where M and $\alpha$ are constants (M = 100 and $\alpha$ = 0.5).
To avoid over–segmentation, if two contours for different connected components intersect, the system discards the connected component with the smaller area.
- Third, for connected components classified as text, bounding boxes are defined; for connected components classified as non-text, simplified contours are computed with the "Ramer Douglas Peucker algorithm."

## 5.3. The Adaptive thresholding-based Method

The "Adaptive thresholding-based Method" was submitted by a team of collaborators from the Prince Mohammad Bin Fahd University, Saudi Arabia, and the University of Malaysia Sarawak, Malaysia, Prince Mohammad Bin Fahd University, Saudi Arabia, which included M.A. Al-Dobais, F.A.G. Alrasheed, G. Latif, and L. Alzubaidi [10]. The method is a heuristic rule-based algorithm for analyzing a document image. The input images are resized to 1260 × 920 pixels. Unlike known binarizartion methods such as Outsu's, where a fixed threshold value is selected by maximizing the variance of the image, in linear adaptive thresholding, the threshold value for every pixel of the image is selected based on its surrounding pixel values. This helps overcoming irregular illumination of the input image and providing better binarization quality.

Five percent of the binarized images were cropped before erosion and dilation operations were performed to help define the merged structure of the text and image regions and ignore the outliers.

Geometric features of regions outlined by bounding boxes in the binary image were used to classify these regions as text, non-text, or noise by the following heuristic rules:

a. Text region: 12 pixels < bounding box height ≤ 35 pixels AND box width > 85 pixels.

b. Figure region: Bounding box height > 35 pixels AND box width > 80 pixels.

c. A segmented text region with a bounding box that overlaps the bounding box of a non-text region is discarded, which means it is considered to be part of the non-text region.

d. Segmented regions with shape properties different from those mentioned in a.–c. are considered regions with noise.

## 6. Results

The evaluation metrics for both text and non-text classes were calculated for every document and then averaged for dataset A, B and C, respectively. The segmentation results for the three participating systems are given in Tables 1, 2 and 3 for sets A, B and C, respectively, while the classification results are shown in Tables 4, 5 and 6.



Figure 3. Examples of challenging layouts in set A.



Figure 4. Examples of challenging layouts in set B.

The data in the Tables 1–3 reveal that the segmentation results of the RFAAD system are the most accurate. It has a higher average black pixel segmentation rate (AvgBPR)



Figure 5. Examples of challenging layouts in set C.

TABLE 1. SEGMENTATION PERFORMANCE OF THE 3 SYSTEMS ON THE ASAR 2018 BENCHMARKING DATASET A

|  | AvgBPR | CS | OSE | USE | MSE | FA | $\rho$ |
|---|---|---|---|---|---|---|---|
| **RFAAD** | 86.50 | **10.10** | **1.37** | **1.00** | 0.43 | 2.50 | **2.80** |
| **FCN-based** | 77.00 | 9.06 | 3.67 | 1.94 | 2.43 | **1.83** | 8.04 |
| **Adap. Th.** | **90.68** | 6.07 | 8.50 | 3.20 | 1.40 | 6.10 | 13.10 |

TABLE 2. SEGMENTATION PERFORMANCE OF THE 3 SYSTEMS ON THE ASAR 2018 BENCHMARKING DATASET B

|  | AvgBPR | CS | OSE | USE | MSE | FA | $\rho$ |
|---|---|---|---|---|---|---|---|
| **RFAAD** | **90.00** | **13.25** | **1.36** | **1.68** | 1.07 | 3.50 | **4.10** |
| **FCN-based** | 76.50 | 9.13 | 3.07 | 4.13 | 2.83 | **1.04** | 10.03 |
| **Adap. Th.** | 83.99 | 6.5 | 15.47 | 4.57 | 1.23 | 6.5 | 21.27 |

TABLE 3. SEGMENTATION PERFORMANCE OF THE 3 SYSTEMS ON ASAR 2018 BENCHMARKING DATASET C

|  | AvgBPR | CS | OSE | USE | MSE | FA | $\rho$ |
|---|---|---|---|---|---|---|---|
| **RFAAD** | **78.66** | **9.40** | **0.66** | **0.59** | 1.45 | 3.60 | **2.70** |
| **FCN-based** | 70.00 | 6.90 | 3.79 | 0.86 | 1.83 | **1.73** | 6.48 |
| **Adap. Th.** | 77.71 | 6.53 | 5.93 | 3.07 | **0.93** | 2.97 | 9.93 |

for datasets B and C and a higher correct segmentation rate (CS) than the other two methods for the three datasets A–C. Fewer over–segmentation, under–segmentation, and missed segmentation errors were reported than for the other methods for datasets A–C. This means that the overall error rate $\rho$ of the RFAAD system is lower for all three datasets. However, the RFAAD system has a higher false alarm rate than the FCN-based system.

The classification results shown in Tables 4–6 do not show a clear winner among the three methods or datasets. The Adaptive thresholding-based method performed well on the classification task for set A. The FCN-based method worked well on the classification task of sets B and C, particularly with respect to the block-level performance measures.

TABLE 4. CLASSIFICATION PERFORMANCE OF THE 3 SYSTEMS ON THE ASAR 2018 BENCHMARKING DATASET A

| | Pixels (%) | | | | Blocks (%) | | | |
|---|---|---|---|---|---|---|---|---|
| | Pr | Rec | F1 | Acc | Pr | Rec | F1 | Acc |
| RFAAD | 62 | **98** | 66 | 69 | 81 | 86 | 82 | 75 |
| FCN-based | 82 | 94 | 87 | 80 | **99** | 83 | 88 | **97** |
| Adap. Th. | **89** | 90 | **90** | **89** | 92 | **95** | **93** | 90 |

TABLE 5. CLASSIFICATION PERFORMANCE OF THE 3 SYSTEMS ON THE ASAR 2018 BENCHMARKING DATASET B

| | Pixels (%) | | | | Blocks (%) | | | |
|---|---|---|---|---|---|---|---|---|
| | Pr | Rec | F1 | Acc | Pr | Rec | F1 | Acc |
| RFAAD | 50 | **97** | 56 | 59 | 79 | 81 | 79 | 71 |
| FCN-based | **88** | 94 | **90** | **87** | **98** | **97** | **97** | **99** |
| Adap. Th. | 86 | 90 | 87 | 86 | 81 | 83 | 82 | 87 |

TABLE 6. CLASSIFICATION PERFORMANCE OF THE 3 SYSTEMS ON THE ASAR 2018 BENCHMARKING DATASET C

| | Pixels (%) | | | | Blocks (%) | | | |
|---|---|---|---|---|---|---|---|---|
| | Pr | Rec | F1 | Acc | Pr | Rec | F1 | Acc |
| RFAAD | 64 | **96** | 70 | 71 | 75 | 83 | 77 | 69 |
| FCN-based | 71 | 95 | **81** | 75 | **99** | **85** | **90** | **93** |
| Adap. Th. | 71 | 84 | 76 | **82** | 74 | 78 | 76 | 82 |

## 7. Conclusion

There is a crucial need to find a robust solution to the physical layout analysis problem of scanned Arabic documents of all types. An overview of the organization and the results of the ASAR 2018 Layout Analysis Challenge on Scanned Arabic Book Pages was given in this paper. Three new systems from three different Arabic-speaking countries were submitted to the competition. Each of the proposed methods used one of the following known approaches to PLA: rule based, traditional learning based, and deep learning based. Despite the relatively challenging page layout of the book pages in the benchmarking dataset, good results were achieved by the submitted methods. The competition had one winner with top results for the segmentation task (RFAAD), one winner with top results for the classification task for set A (Adap. Th.), and one winner with top results for the classification task for sets B and C (FCN-based). The running times of the competition methods were not compared since the participating teams most likely used different hardware for testing. The detailed quantitative and qualitative results of the competition are published on the challenge website http://www.cs.bu.edu/faculty/betke/ASAR-LayoutAnalysisCompetition including visualizations.

Page analysis for complex layouts of Arabic book pages is still an unsolved problem, and further research is needed to reach high success rates. We hope the PLA-SAB 2018 Competition will trigger additional research on this exciting topic. A new competition round is planned for 2019. The data will be more challenging, and more tools of comparison will be provided. We encourage new submissions to this challenge at the challenge website.

## References

[1] R. S. M. Saad, R. I. Elanwar, N. S. Abdel Kader, S. Mashali, and M. Betke. BCE-Arabic-v1 dataset: A step towards interpreting Arabic document images for people with visual impairments. In ACM 9th Annual International Conference on Pervasive Technologies Related to Assistive Environments (PETRA'16), pp. 25–32, Corfu, Greece, 2016.

[2] A. Alshameri, S. Abdou, and K. Mostafa. A combined algorithm for layout analysis of Arabic document images and text lines extracon. Internaonal Journal of Computer Applicaons, 49(23), pp. 30-37, 2012.

[3] S. S. Bukhari, T. M. Breuel, A. Asi, and J. ElSana. Layout analysis for Arabic historical document images using machine learning. International Conference on Frontiers in Handwriting Recognition (ICFHR), pp. 639-644, 2012.

[4] F. Shafait, D. Keysers, and T.M. Breuel. Performance evaluation and benchmarking of six-page segmentation algorithms. IEEE Transactions on Pattern Analysis and Machine Intelligence, 30(6), pp. 941-954, 2008.

[5] C. Clausner, S. Pletschacher, and A. Antonacopoulos. Aletheia - an advanced document layout and text ground-truthing system for production environments. IEEE International Conference on Document Analysis and Recognition (ICDAR), pp. 48–52, 2011.

[6] S. Pletschacher and A. Antonacopoulos. The PAGE (Page Analysis and Ground-Truth Elements) format framework. In 20th International Conference on Pattern Recognition (ICPR), pp. 257–260, 2010.

[7] V. P. Le, M. Nayef, M. Visani, J. M. Ogier, and C. D. Tran. Text and Non-text Segmentation using Connected Component-based Features. In Proceedings of the IEEE 13th International Conference on Document Analysis and Recognition (ICDAR), pp. 1096–1100, 2015.

[8] R. S. M. Saad, R. Elanwar, N. S. Abdel Kader, S. Mashali, and M. Betke. ASAR 2018 Layout Analysis Challenge: Using Random Forests to Analyze Scanned Arabic Books. 2nd IEEE International Workshop on Arabic and derived Script Analysis and Recognition (ASAR 2018), London, March 2018. 6 pages.

[9] A. Droby, B. Barakat, and J. El-Sana. Binarization Free Layout Analysis for Arabic Historical Documents Using Fully Convolutional Networks. 2nd IEEE International Workshop on Arabic and derived Script Analysis and Recognition (ASAR 2018), London, March 2018.

[10] M.A. Al-Dobais, F.A.G. Alrasheed, G. Latif, and L. Alzubaidi. Adaptive Thresholding and Geometric Features based Physical Layout Analysis of Scanned Arabic Books. 2nd IEEE International Workshop on Arabic and derived Script Analysis and Recognition (ASAR 2018), London, March 2018.