

Indiscapes: Instance Segmentation Networks for Layout Parsing of Historical Indic Manuscripts

Abhishek Prusty, Sowmya Aitha, Abhishek Trivedi, Ravi Kiran Sarvadevabhatla

Centre for Visual Information Technology (CVIT)

International Institute of Information Technology, Hyderabad (IIIT-H)

Gachibowli, Hyderabad 500032, INDIA.

{abhishek.prusty@students., sowmya.aitha@research., abhishek.trivedi@research., ravi.kiran@}iiit.ac.in

Abstract—Historical palm-leaf manuscript and early paper documents from Indian subcontinent form an important part of the world’s literary and cultural heritage. Despite their importance, large-scale annotated Indic manuscript image datasets do not exist. To address this deficiency, we introduce Indiscapes, the first ever dataset with multi-regional layout annotations for historical Indic manuscripts. To address the challenge of large diversity in scripts and presence of dense, irregular layout elements (e.g. text lines, pictures, multiple documents per image), we adapt a Fully Convolutional Deep Neural Network architecture for fully automatic, instance-level spatial layout parsing of manuscript images. We demonstrate the effectiveness of proposed architecture on images from the Indiscapes dataset. For annotation flexibility and keeping the non-technical nature of domain experts in mind, we also contribute a custom, web-based GUI annotation tool and a dashboard-style analytics portal. Overall, our contributions set the stage for enabling downstream applications such as OCR and word-spotting in historical Indic manuscripts at scale.

Keywords—Document Layout Parsing; Palm-leaf manuscripts ; Semantic Instance Segmentation ; Deep Neural Networks, Indic

I. INTRODUCTION

The collection and analysis of historical document images is a key component in the preservation of culture and heritage. Given its importance, a number of active research efforts exist across the world [1]–[6]. In this paper, we focus on palm-leaf and early paper documents from the Indian subcontinent. In contrast with modern or recent era documents, such manuscripts are considerably more fragile, prone to degradation from elements of nature and tend to have a short shelf life [7]–[9]. More worryingly, the domain experts who can decipher such content are small in number and dwindling. Therefore, it is essential to access the content within these documents before it is lost forever.

Surprisingly, no large-scale annotated Indic manuscript image datasets exist for the benefit of researchers in the community. In this paper, we take a significant step to address this gap by creating such a dataset. Given the large diversity in language, script and non-textual regional elements in these manuscripts, spatial layout parsing is crucial in enabling downstream applications such as OCR, word-spotting, style-and-content based retrieval and clustering. For this reason, we first tackle the problem of creating a diverse, annotated *spatial layout* dataset. This has the immediate advantage of bypassing

the hurdle of language and script familiarity for annotators since layout annotation does not require any special expertise unlike text annotation.

In general, manuscripts from Indian subcontinent pose many unique challenges (Figure 1). To begin with, the documents exhibit a large multiplicity of languages. This is further magnified by variations in intra-language script systems. Along with text, manuscripts may contain pictures, tables, non-pictorial decorative elements in non-standard layouts. A unique aspect of Indic and South-East Asian manuscripts is the frequent presence of holes punched in the document for the purpose of binding [7], [9], [10]. These holes cause unnatural gaps within text lines. The physical dimensions of the manuscripts are typically smaller compared to other historical documents, resulting in a dense content layout. Sometimes, multiple manuscript pages are present in a single image. Moreover, imaging-related factors such as varying scan quality play a role as well. Given all of these challenges, it is important to develop robust and scalable approaches for the problem of layout parsing. In addition, given the typical non-technical nature of domain experts who study manuscripts, it is also important to develop easy-to-use graphical interfaces for annotation, post-annotation visualization and analytics.

We make the following contributions:

- We introduce Indiscapes, the first ever historical Indic manuscript dataset with detailed spatial layout annotations (Section III).
- We adapt a deep neural network architecture for instance-level spatial layout parsing of historical manuscript images (Section IV-A).
- We also introduce a lightweight web-based GUI for annotation and dashboard-style analytics keeping in mind the non-technical domain experts and the unique layout-level challenges of Indic manuscripts (Section III-B).

II. RELATED WORK

A number of research groups have invested significant efforts in the creation and maintenance of annotated, publicly available historical manuscript image datasets [1]–[4], [11]–[13]. Other collections contain character-level and word-level spatial annotations for South-East Asian palm-leaf manuscripts [5], [10], [14]. In these latter set of works,



Fig. 1: The five images on the left, enclosed by pink dotted line, are from the BHOOMI palm leaf manuscript collection while the remaining images (enclosed by blue dotted line) are from the 'Penn-in-Hand' collection (refer to Section III). Note the inter-collection differences, closely spaced and unevenly written text lines, presence of various non-textual layout regions (pictures, holes, library stamps), physical degradation and presence of multiple manuscripts per image. All of these factors pose great challenges for annotation and machine-based parsing.

annotations for lines are obtained by considering the polygonal region formed by union of character bounding boxes as a line. While studies on Indic palm-leaf and paper-based manuscripts exist, these are typically conducted on small and often, private collections of documents [15]–[21]. No publicly available large-scale, annotated dataset of historical Indic manuscripts exists to the best of our knowledge. In contrast with existing collections, our proposed dataset contains a much larger diversity in terms of document type (palm-leaf and early paper), scripts and annotated layout elements (see Tables I,III). An additional level of complexity arises from the presence of multiple manuscript pages within a single image (see Fig. 1).

A number of contributions can also be found for the task of historical document layout parsing [22]–[25]. Wei et al. [23] explore the effect of using a hybrid feature selection method while using autoencoders for semantic segmentation in five historical English and Medieval European manuscript datasets. Chen et al. [25] explore the use of Fully Convolutional Networks (FCN) for the same datasets. Barakat et al. [26] propose a FCN for segmenting closely spaced, arbitrarily oriented text lines from an Arabic manuscript dataset. The mentioned approaches, coupled with efforts to conduct competitions on various aspects of historical document layout analysis have aided progress in this area [27]–[29]. A variety of layout parsing approaches, including those employing the modern paradigm of deep learning, have been proposed for Indic [18], [20], [21], [30] and South-East Asian [14], [24], [31]–[33] palm-leaf and paper manuscript images. However, existing approaches typically employ brittle hand-crafted features or demonstrate performance on datasets which are limited in terms of layout diversity. Similar to many recent works, we employ Fully Convolutional Networks in our approach. However, a crucial distinction lies in our formulation of layout parsing as an *instance* segmentation problem, rather than just

a *semantic* segmentation problem. This avoids the problem of closely spaced layout regions (e.g. lines) being perceived as contiguous blobs.

The ready availability of annotation and analysis tools has facilitated progress in creation and analysis of historical document manuscripts [34]–[36]. The tool we propose in the paper contains many of the features found in existing annotation systems. However, some of these systems are primarily oriented towards single-user, offline annotation and do not enable a unified management of annotation process and monitoring of annotator performance. In contrast, our web-based system addresses these aspects and provides additional capabilities. Many of the additional features in our system are tailored for annotation and examining annotation analytics for documents with dense and irregular layout elements, especially those found in Indic manuscripts. In this respect, our annotation system is closer to the recent trend of collaborative, cloud/web-based annotation systems and services [37]–[39].

III. INDISCAPES: THE INDIC MANUSCRIPT DATASET

The Indic manuscript document images in our dataset are obtained from two sources. The first source is the publicly available Indic manuscript collection from University of Pennsylvania's Rare Book and Manuscript Library [40], also referred to as Penn-in-Hand (PIH). From the 2,880 Indic manuscript book-sets¹, we carefully curated 193 manuscript images for annotation. Our curated selection aims to maximize the diversity of the dataset in terms of various attributes such as the extent of document degradation, script language, presence of non-textual elements (e.g. pictures, tables) and number of lines. Some images contain multiple manuscript pages stacked vertically or horizontally (see bottom-left image in Figure 1). The second source for manuscript images in our dataset is

¹A book-set is a sequence of manuscript images.

	Character Line Segment (CLS)	Character Component (CC)	Hole (H)	Page Boundary (PB)	Library Marker (LM)	Decorator (D)	Picture (P)	Physical Degradation (PD)	Boundary Line (BL)
PIH	2401	494	—	256	32	59	94	34	395
BHOOMI	2440	210	565	316	133	—	—	2078	—
Combined	4841	704	565	572	165	59	94	2112	395

TABLE I: Counts for various annotated region types in INDISCAPES dataset. The abbreviations used for region types are given below each region type.

	Train	Validation	Test	Total
PIH	116	28	49	193
BHOOMI	236	59	20	315
Total	352	87	69	508

TABLE II: Dataset splits used for learning and inference.

Script	Source	Document Count
Devanagari	PIH	193
Nandinagari	BHOOMI	2
Telugu	BHOOMI	75
Grantha	BHOOMI	238

TABLE III: Scripts in the INDISCAPES dataset.

BHOOMI, an assorted collection of 315 images sourced from multiple Oriental Research Institutes and libraries across India. As with the first collection, we chose a subset intended to maximize the overall diversity of the dataset. However, this latter set of images are characterized by a relatively inferior document quality, presence of multiple languages and from a layout point of view, predominantly contain long, closely and irregularly spaced text lines, binding holes and degradations (Figure 1). Though some document images contain multiple manuscripts, we do not attempt to split the image into multiple pages. While this poses a challenge for annotation and automatic image parsing, retaining such images in the dataset eliminates manual/semi-automatic intervention. As our results show, our approach can successfully handle such multi-page documents, thereby making it truly an end-to-end system.

Overall, our dataset contains 508 annotated Indic manuscripts. Some salient aspects of the dataset can be viewed in Table I and a pictorial illustration of layout regions can be viewed in Figure 4. Note that multiple regions can overlap, unlike existing historical document datasets which typically contain disjoint region annotations.

For the rest of the section, we discuss the challenges associated with annotating Indic manuscripts (Section III-A) and our web-based annotation tool (Section III-B).

A. Annotation Challenges

A variety of unique challenges exist in the context of annotating Indic manuscript layouts. The challenges arise from three major sources.

Content: The documents are written in a large variety of Indic languages. Some languages even exhibit intra-language script variations. A large pool of annotators familiar with the languages and scripts present in the corpus is required to ensure proper annotation of lines and character components.

Layout: Unlike some of the existing datasets, Indic manuscripts contain non-textual elements such as color pictures, tables and document decorations. These elements are frequently interspersed with text in non-standard layouts. In many cases, the manuscripts contain one or more physical holes, designed for a thread-like material to pass through and bind the leaves together as a book. Such holes vary in terms of spatial location, count and hole diameter. When the holes are present in the middle of the document, they cause a break in the contiguity of lines. In some documents, the line contiguity is broken by a ‘virtual’ hole-like gap, possibly intended for creation of the punched hole at a future time. In many cases, the separation between lines is extremely small. The handwritten nature of these documents and the surface material result in extremely uneven lines, necessitating meticulous and slow annotation. If multiple manuscript pages are present, the stacking order could be horizontal or vertical. Overall, the sheer variety in layout elements poses a significant challenge, not only for annotation, but also for automated layout parsing.

Degradations: Historical Indic manuscripts tend to be inherently fragile and prone to damage due to various sources – wood-and-leaf-boring insects, humidity seepage, improper storage and handling etc. While some degradations cause the edges of the document to become frayed, others manifest as irregularly shaped perforations in the document interior. It may be important to identify such degradations before attempting lexically-focused tasks such as OCR or word-spotting.

B. Annotation Tool

Keeping the aforementioned challenges in mind, we introduce a new browser-based annotation tool (see Figure 2). The tool is designed to operate both stand-alone and as a web-service. The web-service mode enables features such as distributed parallel sessions by registered annotators, dashboard-based live session monitoring and a wide variety of annotation-related analytics. On the front-end, a freehand region option is provided alongside the usual rectangle and polygon to enable maximum annotation flexibility. The web-service version also features a ‘Correction-mode’ which enables annotators to correct existing annotations from previous annotators. Additionally, the tool has been designed to enable

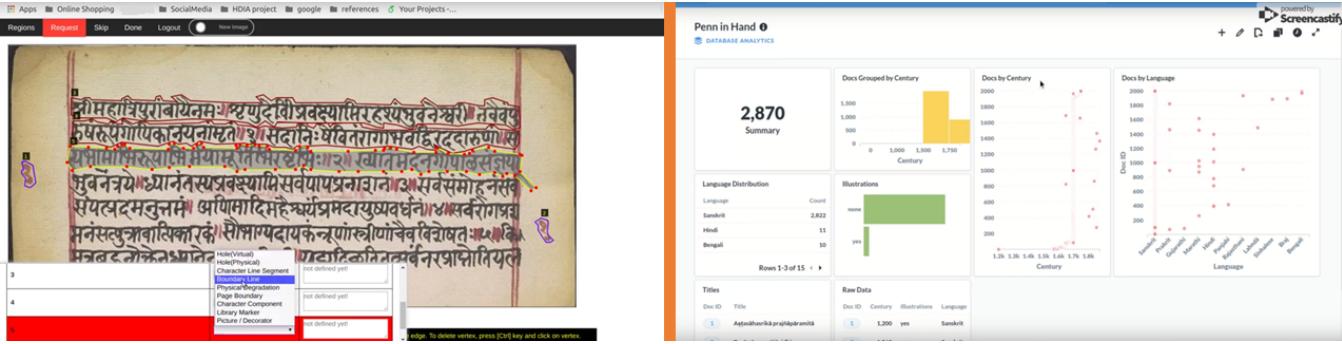


Fig. 2: Screenshots of our web-based annotator (left) and analytics dashboard (right).

lexical (text) annotations in future.

IV. INDIC MANUSCRIPT LAYOUT PARSING

To succeed at layout parsing of manuscripts, we require a system which can accurately localize various types of regions (e.g. text lines, isolated character components, physical degradation, pictures, holes). More importantly, we require a system which can isolate individual *instances* of each region (e.g. multiple text lines) in the manuscript image. Also, in our case, the annotation regions for manuscripts are not disjoint and can overlap (e.g. The annotation region for a text line can overlap with the annotation region of a hole (see Figure 4)). Therefore, we require a system which can accommodate such overlaps. To meet all of these requirements, we model our problem as one of semantic *instance-level* segmentation and employ the Mask R-CNN [41] architecture which has proven to be very effective at the task of object-instance segmentation in photos. Next, we briefly describe the Mask R-CNN architecture and our modifications of the same. Subsequently, we provide details related to implementation (Section IV-B), model training (Section IV-B1) and inference (Section IV-B2).

A. Network Architecture

The Mask-RCNN architecture contains three stages as described below (see Figure 3).

Backbone: The first stage, referred to as the backbone, is used to extract features from the input image. It consists of a convolutional network combined with a feature-pyramid network [42], thereby enabling multi-scale features to be extracted. We use the first four blocks of ResNet-50 [43] as the convolutional network.

Region Proposal Network (RPN): This is a convolutional network which scans the pyramid feature map generated by the backbone network and generates rectangular regions commonly called ‘object proposals’ which are likely to contain objects of interest. For each level of the feature pyramid and for each spatial location at a given level, a set of level-specific bounding boxes called anchors are generated. The anchors typically span a range of aspect ratios (e.g. 1 : 2, 1 : 1, 2 : 1) for flexibility in detection. For each anchor, the RPN network predicts (i) the probability of an object being present (‘objectness score’) (ii) offset coordinates of a bounding box relative to

location of the anchor. The generated bounding boxes are first filtered according to the ‘objectness score’. From boxes which survive the filtering, those that overlap with the underlying object above a certain threshold are chosen. After applying non-maximal suppression to remove overlapping boxes with relatively smaller objectness scores, the final set of boxes which remain are termed ‘object proposals’ or Regions-of-Interest (RoI).

Multi-Task Branch Networks: The RoIs obtained from RPN are warped into fixed dimensions and overlaid on feature maps extracted from the backbone to obtain ROI-specific features. These features are fed to three parallel task sub-networks. The first sub-network maps these features to region labels (e.g. Hole,Character-Line-Segment) while the second sub-network maps the ROI features to bounding boxes. The third sub-network is fully convolutional and maps the features to the pixel mask of the underlying region. Note that the ability of the architecture to predict masks independently for each ROI plays a crucial role in obtaining instance segmentations. Another advantage is that it naturally addresses situations where annotations or predictions overlap.

B. Implementation Details

The dataset splits used for training, validation and test phases can be seen in Table II. All manuscript images are adaptively resized to ensure the width does not exceed 1024 pixels. The images are padded with zeros such that the input to the deep network has spatial dimensions of 1024×1024 . The ground truth region masks are initially subjected to a similar resizing procedure. Subsequently, they are downsized to 28×28 in order to match output dimensions of the mask sub-network.

1) Training: The network is initialized with weights obtained from a Mask R-CNN trained on the MS-COCO [44] dataset with a ResNet-50 backbone. We found that this results in faster convergence and stabler training compared to using weights from a Mask-RCNN trained on ImageNet [45] or training from scratch. Within the RPN network, we use custom-designed anchors of 5 different scales and with 3 different aspect ratios. Specifically, we use the following aspect ratios – 1:1, 1:3, 1:10 – keeping in mind the typical spatial extents of the various region classes. We also limit

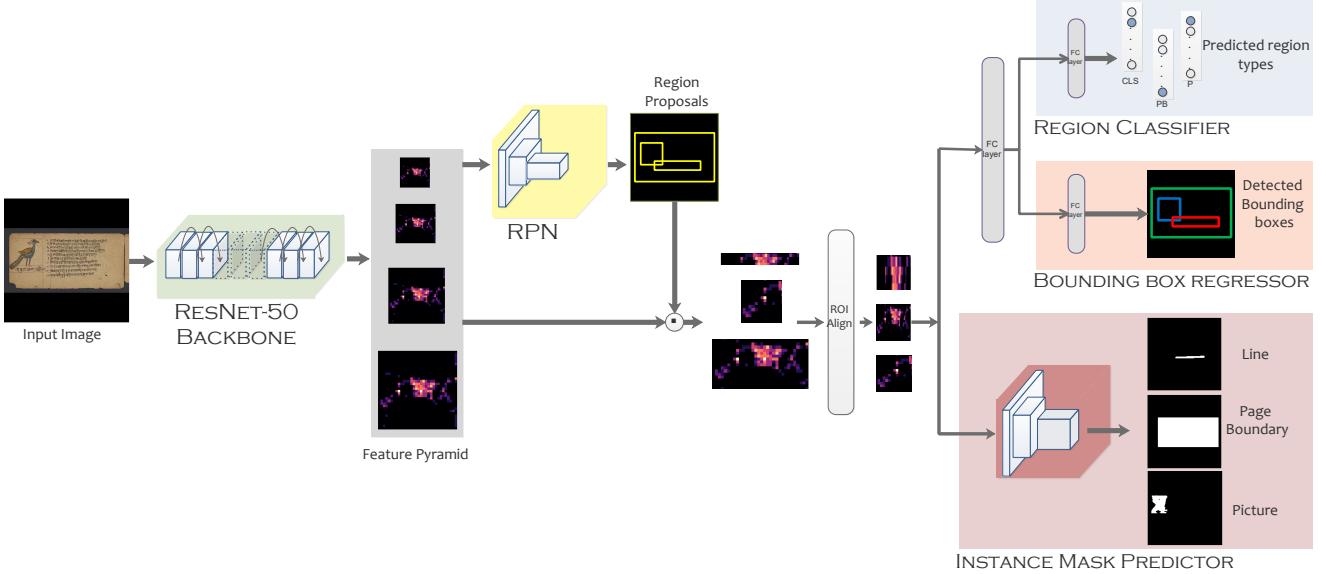


Fig. 3: The architecture adopted for Indic Manuscript Layout Parsing. Refer to Section IV for details.

Dataset ↓	H	Average IoU / Average Per pixel Accuracy								
		CLS	PD	PB	CC	P	D	LM	BL	
PIH	—	74.17/92.57	—	86.90/96.37	52.84/74.85	60.49/82.21	5.23/6.17	50.29/56.97	29.45/43.14	
BHOOMI	79.29/99.95	29.07/43.67	8.72/12.98	91.09/99.22	32.50/47.19	—	—	38.25/49.98	—	
Combined	79.29/99.95	57.77/74.79	8.72/12.98	88.47/97.44	45.87/65.37	60.49/82.21	5.23/6.17	42.93/52.70	29.45/43.14	

TABLE IV: Class-wise average IoUs and per-pixel accuracies on the test set. Refer to Table I for full names of abbreviated region types listed at top of the table.

	AP_{50}	AP_{75}	AP
PIH	79.78	60.11	49.64
Bhoomi	36.88	14.95	18.00
Combined	64.76	44.30	38.57

TABLE V: AP at IoU thresholds 50, 75 and overall AP averaged over IoU range for test set.

the number of RoIs ('object proposals') to 512. We use categorical cross entropy loss \mathcal{L}_{RPN} for RPN classification network. Within the task branches, we use categorical cross entropy loss \mathcal{L}_r for region classification branch, smooth L1 loss [46] (\mathcal{L}_{bb}) for final bounding box prediction and per-pixel binary cross entropy loss \mathcal{L}_{mask} for mask prediction. The total loss is a convex combination of these losses, i.e. $\mathcal{L} = \lambda_{RPN}\mathcal{L}_{RPN} + \lambda_r\mathcal{L}_r + \lambda_{bb}\mathcal{L}_{bb} + \lambda_{mask}\mathcal{L}_{mask}$. The weighting factors (λ s) are set to 1. However, to ensure priority for our task of interest namely mask prediction, we set $\lambda_{mask} = 2$. For optimization, we use Stochastic Gradient Descent (SGD) optimizer with a gradient norm clipping value of 0.5. The batch size, momentum and weight decay are set to 1, 0.9 and 10^{-3} respectively. Given the relatively smaller size of our manuscript dataset compared to the photo dataset

(MS-COCO) used to originally train the base Mask R-CNN, we adopt a multi-stage training strategy. For the first stage (30 epochs), we train only the task branch sub-networks using a learning rate of 10^{-3} while freezing weights in the rest of the overall network. This ensures that the task branches are fine-tuned for the types of regions contained in manuscript images. For the second stage (20 epochs), we additionally train stage-4 and up of the backbone ResNet-50. This enables extraction of appropriate semantic features from manuscript images. The omission of the initial 3 stages in the backbone for training is due to the fact that they provide generic, re-usable low-level features. To ensure priority coverage of hard-to-localize regions, we use focal loss [47] for mask generation. For the final stage (15 epochs), we train the entire network using a learning rate of 10^{-4} .

2) *Inference*: During inference, the images are rescaled and processed using the procedure described at the beginning of the subsection. The number of RoIs retained after non-maximal suppression (NMS) from the RPN is set to 1000. From these, we choose the top 100 region detections with objectness score exceeding 0.5 and feed the corresponding RoIs to the mask branch sub-network for mask generation. It is important to note that this strategy is different from the parallel generation of outputs and use of the task sub-networks during training. The generated masks are then binarized using an empirically chosen threshold of 0.4 and rescaled to their

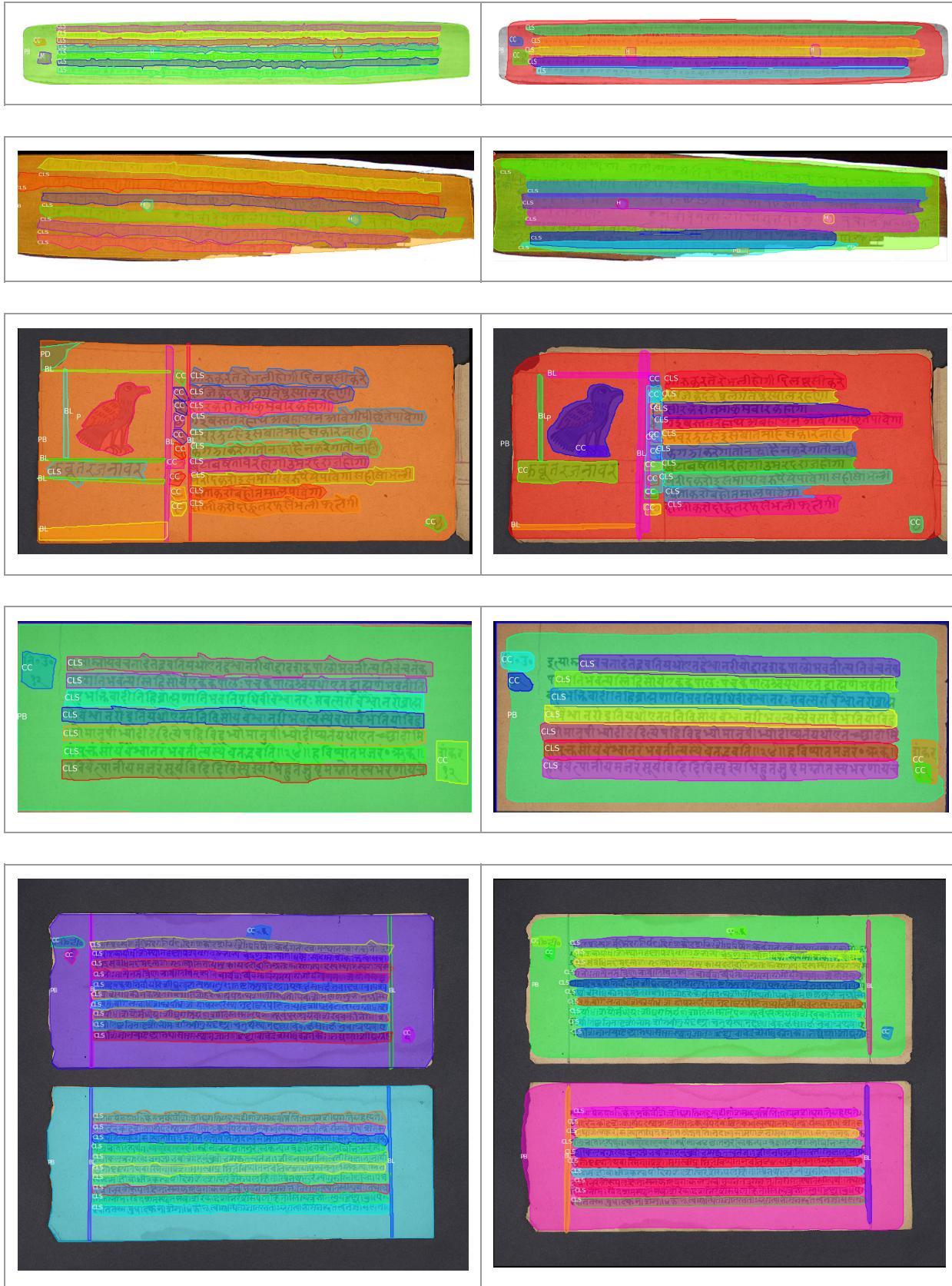


Fig. 4: Ground truth annotations (left) and predicted instance segmentations (right) for test set images. Note that we use colored shading only to visualize individual region instances and not to color-code region types. The region label abbreviations are shown alongside the regions. CLS : Character Line Segment, PB : Page Boundary, H : Hole, BL : Boundary Line, CC : Character Component, PD : Physical Degradation.

original size using bilinear interpolation. On these generated masks, NMS with a threshold value of 0.5 is applied to obtain the final set of predicted masks.

C. Evaluation

For quantitative evaluation, we compute Average Precision (AP) for a particular IoU threshold, a measure widely reported in instance segmentation literature [44], [48]. We specifically report AP_{50} and AP_{75} , corresponding to AP at IoU thresholds 50 and 75 respectively [41]. In addition, we report an overall score by averaging AP at different IoU thresholds ranging from 0.5 to 0.95 in steps of 0.05.

The AP measure characterizes performance at document level. To characterize performance for each region type, we report two additional measures [25] – average class-wise IoU (cwIoU) and average class-wise per-pixel accuracy (cwAcc). Consider a fixed test document k . Suppose there are r_i regions of class i and let IoU_r denote the IoU score for one such region r , i.e. $1 \leq r \leq r_i$. The per-class IoU score for class i and document k is computed as $cwIoU_i^d = \frac{\sum_r IoU_r}{r_i}$. Suppose there are N_i documents containing at least a single region of class i in ground-truth. The overall per-class IoU score for class i is computed as $cwIoU_i = \frac{\sum_d cwIoU_i^d}{N_i}$. In a similar manner, we define class-wise pixel accuracy $pwAcc_i^d$ at document level and average it across all the documents containing class i , i.e. $cwAcc_i = \frac{\sum_d pwAcc_i^d}{N_i}$. Note that our approach for computing class-wise scores prevents documents with a relatively larger number of class instances from dominating the score and in this sense, differs from existing approaches [25]

V. RESULTS

We report quantitative results using the measures described in Section IV-C. Table IV reports Average Precision and Table V reports class-wise average IOUs and per-pixel accuracies. Qualitative results can be viewed in Figure 4. Despite the challenges posed by manuscripts, our model performs reasonably well across a variety of classes. As the qualitative results indicate, the model predicts accurate masks for almost all the regions. The results also indicate that our model handles overlap between *Holes* and *Character line segments* well. From ablative experiments, we found that our choice of focal loss was crucial in obtaining accurate mask boundaries. Unlike traditional semantic segmentation which would have produced a single blob-like region for line segments, our instance-based approach isolates each text line separately. Additionally, the clear demarcation between *Page-Boundary* and background indicates that our system identifies semantically relevant regions for downstream analysis. As the result at the bottom of Figure 4 shows, our system can even handle images with multiple pages, thus removing the need for any pre-processing related to isolation of individual pages.

From quantitative results, we observe that *Holes*, *Character line segments*, *Page boundary* and *Pictures* are parsed the best while *Physical degradations* are difficult to parse due to the relatively small footprint and inconsistent patterns in degradations. The results show that performance for Penn

in Hand (PIH) documents is better compared to BHOOMI manuscripts. We conjecture that the presence of closely spaced and unevenly written lines in latter is the cause. In our approach, two (or more) objects may share the same bounding box in terms of overlap and it is not possible to determine which box to choose during mask prediction. Consequently, an underlying line's boundary may either end up not being detected or the predicted mask might be poorly localized. However, this is not a systemic problem since our model achieves good performance even for very dense BHOOMI document line layouts.

VI. CONCLUSION

Via this paper, we propose Indiscapes, the first dataset with layout annotations for historical Indic manuscripts. We believe that the availability of layout annotations will play a crucial role in reducing the overall complexity for OCR and other tasks such as word-spotting, style-and-content based retrieval. In the long-term, we intend to expand the dataset, not only numerically but also in terms of layout, script and language diversity. As a significant contribution, we have also adapted a deep-network based instance segmentation framework custom modified for fully automatic layout parsing. Given the general nature of our framework, advances in instance segmentation approaches can be leveraged thereby improving performance over time. Our proposed web-based annotator system, although designed for Indic manuscripts, is flexible, and could be reused for similar manuscripts from Asian subcontinent. We intend to expand the capabilities of our annotator system in many useful ways. For instance, the layout estimated by our deep-network could be provided to annotators for correction, thus reducing annotation efforts. Finally, we plan to have our dataset, instance segmentation system and annotator system publicly available. This would enable large-scale data collection and automated analysis efforts for Indic as well as other historical Asian manuscripts. The repositories related to the systems presented in this paper and the Indiscapes dataset can be accessed at <https://ihdia.iit.ac.in>.

ACKNOWLEDGMENT

We would like to thank Dr. Sai Susarla for enabling access to the Bhoomi document collection. We also thank Poreddy Mourya Kumar Reddy, Gollapudi Sai Vamsi Krishna for their contributions related to dashboard and various annotators for their labelling efforts.

REFERENCES

- [1] C. Reul, M. Dittrich, and M. Gruner, “Case study of a highly automated layout analysis and ocr of an incunabulum:’der heiligen leben’(1488),” in *Proc. 2nd Intl. Conf. on Digital Access to Textual Cultural Heritage*. ACM, 2017, pp. 155–160. [1](#)
- [2] U. Springmann and A. Luedeling, “Ocr of historical printings with an application to building diachronic corpora: A case study using the ridges herbal corpus,” *Digital Humanities Quarterly*, no. 2, 2017. [1](#)
- [3] F. Simistira, M. Seuret, N. Eichenberger, A. Garz, M. Liwicki, and R. Ingold, “Diva-hisdb: A precisely annotated large dataset of challenging medieval manuscripts,” in *ICFHR*. IEEE, 2016, pp. 471–476. [1](#)

- [4] A. Pappo-Toledano, F. Chen, G. Latif, and L. Alzubaidi, "Adaptive thresholding and geometric features based physical layout analysis of scanned arabic books," *2018 IEEE 2nd Intl. Workshop on Arabic and Derived Script Analysis and Recognition (ASAR)*, pp. 171–176, 2018. 1
- [5] M. W. A. Kesiman, J.-C. Burie, G. N. M. A. Wibawantara, I. M. G. Sunarya, and J.-M. Ogier, "Amadi_lonterset: The first handwritten balinese palm leaf manuscripts dataset," in *ICFHR*. IEEE, 2016, pp. 168–173. 1
- [6] K. Chen, M. Seuret, M. Liwicki, J. Hennebert, and R. Ingold, "Page segmentation of historical document images with convolutional autoencoders," in *ICDAR*. IEEE, 2015, pp. 1011–1015. 1
- [7] J. Sahoo, "A selective review of scholarly communications on palm leaf manuscripts," *Library Philosophy and Practice (e-journal)*, 2016. 1
- [8] Y. B. Rachman, "Palm leaf manuscripts from royal surakarta, indonesia: Deterioration phenomena and care practices," *Intl. Journal for the Preservation of Library and Archival Material*, vol. 39, no. 4, pp. 235–247, 2018. 1
- [9] D. U. Kumar, G. Sreekumar, and U. Athvankar, "Traditional writing system in southern indiapalm leaf manuscripts," *Design Thoughts*, vol. 9, 2009. 1
- [10] D. Valy, M. Verleysen, S. Chhun, and J.-C. Burie, "A new khmer palm leaf manuscript dataset for document analysis and recognition: Sleukrith set," in *Proc. of the 4th Intl. Workshop on Historical Document Imaging and Processing*. ACM, 2017, pp. 1–6. 1
- [11] J. A. Sánchez, V. Bosch, V. Romero, K. Depuydt, and J. De Does, "Handwritten text recognition for historical documents in the transcriptorium project," in *Proc. of the First Intl. Conf. on Digital Access to Textual Cultural Heritage*. ACM, 2014, pp. 111–117. 1
- [12] T. M. Rath and R. Manmatha, "Word spotting for historical documents," *IJDAR*, vol. 9, no. 2-4, pp. 139–152, 2007. 1
- [13] M. Kassis, A. Abdalhaleem, A. Droby, R. Alaasam, and J. El-Sana, "Vm1-hd: The historical arabic documents dataset for recognition systems," in *1st Intl. Workshop on Arabic Script Analysis and Recognition*. IEEE, 2017. 1
- [14] M. Suryani, E. Paulus, S. Hadi, U. A. Darsa, and J.-C. Burie, "The handwritten sundanese palm leaf manuscript dataset from 15th century," in *ICDAR*. IEEE, 2017, pp. 796–800. 1, 2
- [15] C. Clausner, A. Antonacopoulos, T. Derrick, and S. Pletschacher, "Icdar2017 competition on recognition of early indian printed documents-reid2017," in *ICDAR*, vol. 1. IEEE, 2017, pp. 1411–1416. 2
- [16] C. K. Savitha and P. J. Antony, "Machine learning approaches for recognition of offline tulu handwritten scripts," *Journal of Physics: Conference Series*, vol. 1142, p. 012005, nov 2018. 2
- [17] A. Abeysinghe and A. Abeysinghe, "Use of neural networks in archaeology: preservation of assamese manuscripts." International Seminar on Assamese Culture & Heritage, 2018. 2
- [18] P. N. Sastry, T. V. Lakshmi, N. K. Rao, and K. RamaKrishnan, "A 3d approach for palm leaf character recognition using histogram computation and distance profile features," in *Proc. 5th Intl. Conf. on Frontiers in Intelligent Computing: Theory and Applications*. Springer, 2017, pp. 387–395. 2
- [19] N. S. Panyam, V. L. T.R., R. Krishnan, and K. R. N.V., "Modeling of palm leaf character recognition system using transform based techniques," *Pattern Recogn. Lett.*, vol. 84, no. C. Dec. 2016. 2
- [20] Z. Shi, S. Setlur, and V. Govindaraju, "Digital enhancement of palm leaf manuscript images using normalization techniques," in *5th Intl. Conf. On Knowledge Based Computer Systems*, 2004, pp. 19–22. 2
- [21] D. Sudarsan, P. Vijayakumar, S. Biju, S. Sanu, and S. K. Shivadas, "Digitalization of malayalam palmleaf manuscripts based on contrast-based adaptive binarization and convolutional neural networks," in *Intl. Conf. on Wireless Communications, Signal Processing and Networking (WiSPNET)*, 2018. 2
- [22] C. Wick and F. Puppe, "Fully convolutional neural networks for page segmentation of historical document images," in *DAS*. IEEE, 2018, pp. 287–292. 2
- [23] H. Wei, M. Seuret, K. Chen, A. Fischer, M. Liwicki, and R. Ingold, "Selecting autoencoder features for layout analysis of historical documents," in *Proc. 3rd Intl. Workshop on Historical Document Imaging and Processing*, ser. HIP '15. ACM, 2015, pp. 55–62. 2
- [24] S. S. Bukhari, T. M. Breuel, A. Asi, and J. El-Sana, "Layout analysis for arabic historical document images using machine learning," in *ICFHR 2012*. IEEE, 2012, pp. 639–644. 2
- [25] K. Chen, M. Seuret, J. Hennebert, and R. Ingold, "Convolutional neural networks for page segmentation of historical document images," in *ICDAR*, vol. 1. IEEE, 2017, pp. 965–970. 2, 7
- [26] B. Barakat, A. Droby, M. Kassis, and J. El-Sana, "Text line segmentation for challenging handwritten document images using fully convolutional network," in *ICFHR*. IEEE, 2018, pp. 374–379. 2
- [27] M. W. A. Kesiman, D. Valy, J. Burie, E. Paulus, M. Suryani, S. Hadi, M. Verleysen, S. Chhun, and J. Ogier, "ICFHR 2018 competition on document image analysis tasks for southeast asian palm leaf manuscripts," in *ICFHR*, 2018, pp. 483–488. 2
- [28] *Proc. 4th Intl. Workshop on Historical Document Imaging and Processing, Kyoto, Japan, November 10-11, 2017*. ACM, 2017. 2
- [29] *Proc. 3rd Intl. Wksp on Historical Document Imaging and Processing, HIP@ICDAR 2015*. ACM, 2015. 2
- [30] R. S. Sabeenian, M. E. Paramasivam, P. M. Dinesh, R. Adarsh, and G. R. Kumar, "Classification of handwritten tamil characters in palm leaf manuscripts using svm based smart zoning strategies," in *ICBIP*. ACM, 2017. 2
- [31] M. W. A. Kesiman, D. Valy, J.-C. Burie, E. Paulus, M. Suryani, S. Hadi, M. Verleysen, S. Chhun, and J.-M. Ogier, "Benchmarking of document image analysis tasks for palm leaf manuscripts from southeast asia," *Journal of Imaging*, vol. 4, no. 2, p. 43, 2018. 2
- [32] D. Valy, M. Verleysen, S. Chhun, and J.-C. Burie, "Character and text recognition of khmer historical palm leaf manuscripts," in *ICFHR*, 08 2018, pp. 13–18. 2
- [33] E. Paulus, M. Suryani, and S. Hadi, "Improved line segmentation framework for sundanese old manuscripts," *Journal of Physics: Conference Series*, vol. 978, p. 012001, mar 2018. 2
- [34] D. Doermann, E. Zotkina, and H. Li, "GEDI-a groundtruthing environment for document images," in *Ninth IAPR Intl. Workshop on Document Analysis Systems*, 2010. 2
- [35] A. Garz, M. Seuret, F. Simistira, A. Fischer, and R. Ingold, "Creating ground truth for historical manuscripts with document graphs and scribbling interaction," in *DAS*. IEEE, 2016, pp. 126–131. 2
- [36] C. Clausner, S. Pletschacher, and A. Antonacopoulos, "Aletheia-an advanced document layout and text ground-truthing system for production environments," in *ICDAR*. IEEE, 2011, pp. 48–52. 2
- [37] "Web aletheia." [Online]. Available: <https://github.com/PRIMa-Research-Lab/prima-gwt-lib> 2
- [38] M. Würsch, R. Ingold, and M. Liwicki, "Divaservicesa restful web service for document image analysis methods," *Digital Scholarship in the Humanities*, vol. 32, no. 1, pp. i150–i156, 2016. 2
- [39] B. Gatos, G. Louloudis, T. Causer, K. Grint, V. Romero, J. A. Sánchez, A. H. Toselli, and E. Vidal, "Ground-truth production in the transcriptorium project," in *DAS*. IEEE, 2014, pp. 237–241. 2
- [40] "Penn in hand: Selected manuscripts," http://dla.library.upenn.edu/dla/medren/search.html?fq=collection_facet:IndicManuscripts. 2
- [41] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask r-cnn," *ICCV*, pp. 2980–2988, 2017. 4, 7
- [42] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017, pp. 936–944. 4
- [43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778. 4
- [44] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," *CoRR*, vol. abs/1405.0312, 2014. [Online]. Available: <http://arxiv.org/abs/1405.0312> 4, 7
- [45] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*. IEEE, 2009, pp. 248–255. 4
- [46] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NIPS*, 2015, pp. 91–99. 5
- [47] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017, pp. 2980–2988. 5
- [48] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," *CoRR*, vol. abs/1604.01685, 2016. [Online]. Available: <http://arxiv.org/abs/1604.01685> 7