

ICFHR 2018 Competition on Document Image Analysis Tasks for Southeast Asian Palm Leaf Manuscripts

Made Windu Antara Kesiman^{1,2}, Dona Valy^{3,4}, Jean-Christophe Burie¹, Erick Paulus⁵, Mira Suryani⁵, Setiawan Hadi⁵, Michel Verleysen³, Sophea Chhun⁴, and Jean-Marc Ogier¹

¹LSi ULR France, ²LCI Undiksha Indonesia, ³ICTEAM UCL Belgium, ⁴TTC Cambodia, ⁵Unpad Indonesia



Background and Goal

- ✓ The existence of ancient palm leaf manuscripts in Southeast Asia is very important both in term of quantity and variety of historical contents.
- ✓ An effort to explore Document Image Analysis (DIA) research for palm leaf manuscripts collection as the heritage documents from Southeast Asia.
- ✓ A new challenge for DIA researchers because it uses palm leaf as writing media and also with a language and script that have never been analyzed before.

Source of the Collections

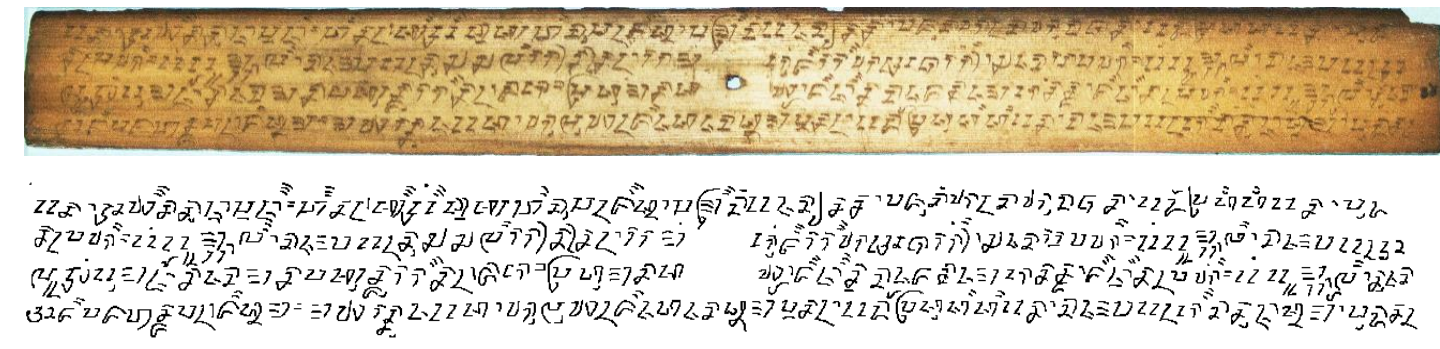
- ☐ Balinese Palm Leaf Manuscripts – Bali – Indonesia:
 - ✓ 23 different collections from 5 different locations (regions)
- ☐ Khmer Palm Leaf Manuscripts – Cambodia:
 - ✓ from Buddhist temples in different locations throughout Cambodia
- ☐ Sundanese Palm Leaf Manuscripts – West Java – Indonesia:
 - ✓ from Situs Kabuyutan Ciburu, Garut, West Java

Challenges: Task, Dataset, Track, Protocol, Evaluation

Challenge A. Binarization (1 Track Mixed)

PALM LEAF MANUSCRIPT DATASETS FOR BINARIZATION TASK

Manuscripts	Train	Test	Dataset
Balinese	50 pages	50 pages	Extracted from AMADI_LontarSet ¹ [3], [9]
Khmer	23 pages	23 pages	Extracted from EFEO ² [13]
Sundanese	31 pages	30 pages	Extracted from Sunda Dataset ICDAR2017 [11]



Evaluation [1]:

- ✓ F-Measure (FM)
- ✓ Peak Signal Noise Ratio (PSNR)
- ✓ Negative Rate Metric (NRM)

Challenge B. Text Line Segmentation (1 Track Mixed)

PALM LEAF MANUSCRIPT DATASETS FOR TEXT LINE SEGMENTATION TASK

Manuscripts	Train	Test	Dataset
Balinese	47 pages	49 pages	Extracted from AMADI_LontarSet [9]
Khmer	50 pages	200 pages	Extracted from SleukRith Set [10]
Sundanese	31 pages	30 pages	Extracted from Sunda Dataset [11]

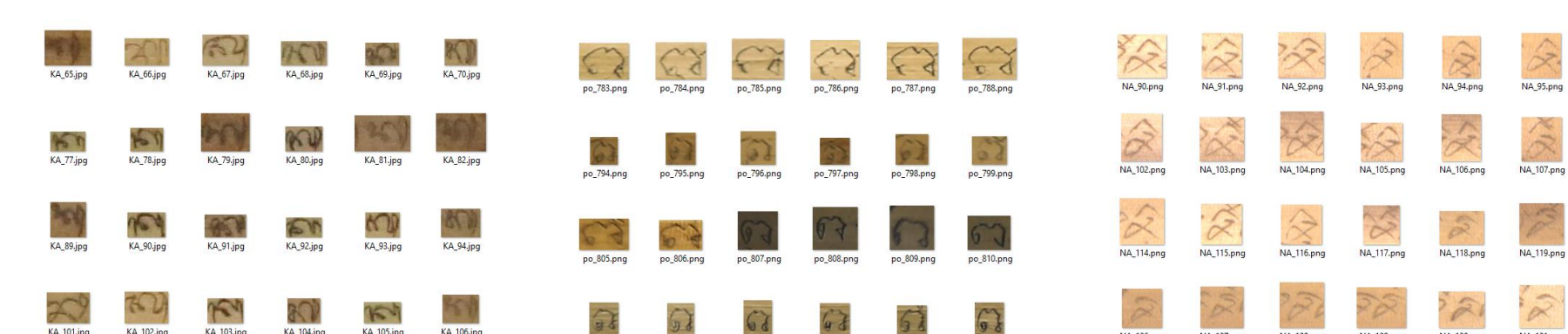
Evaluation [2]:

- ✓ One-to-one match score (o2o)
- ✓ Detection Rate (DR)
- ✓ Recognition Accuracy (RA)
- ✓ Performance Metric (FM)

Challenge C. Isolated Character/Glyph Recognition (3 Single Tracks)

PALM LEAF MANUSCRIPT DATASETS FOR ISOLATED CHARACTER/GLYPH RECOGNITION TASK

Manuscripts	Classes	Train (images)	Test (images)	Dataset
Balinese	133	11,710	7,673	AMADI_LontarSet [3], [6], [9]
Khmer	111	113,206	90,669	SleukRith Set [10]
Sundanese	60	4,555	2,816	Sunda Dataset [11]



Evaluation [1]:

- ✓ Recognition Rate

Challenge D. Word Transliteration (3 Single Tracks and 1 Track Mixed)

PALM LEAF MANUSCRIPT DATASETS FOR WORD TRANSLITERATION TASK

Manuscripts	Train (images)	Test (images)	Text	Dataset
Balinese	15,022 from 130 pages	10,475 from 100 pages	Latin	AMADI_LontarSet [3], [9]
Khmer	16,333 (part of 657 pages)	7,791 (part of 657 pages)	Latin and Khmer	SleukRith Set [10]
Sundanese	1,427 from 20 pages	318 from 10 pages	Latin	Sunda Dataset [11]



- ✓ Evaluation [3]: Character Error Rate (CER)

The Participants and Evaluation Results

❖ 22 research groups registered, 8 research groups submitted their results: 4 groups for Challenge A, 1 group for Challenge B, 2 groups for Challenge C, and 2 groups for Challenge D.

EVALUATION RESULTS OF CHALLENGE A

Group	FM	NRM	PSNR
G12	49.60	0.16	23.43
G17	58.87	0.17	28.71
G18	37.87	0.31	22.87
G22	31.17	0.35	19.10

EVALUATION RESULTS OF CHALLENGE A PER COLLECTION

Group	Bali			Khmer			Sunda		
	FM	NRM	PSNR	FM	NRM	PSNR	FM	NRM	PSNR
G12	52.90	0.16	26.53	44.38	0.15	20.73	48.09	0.17	20.32
G17	56.45	0.19	29.53	66.94	0.11	30.70	56.72	0.20	25.82
G18	47.54	0.22	26.16	1.03	0.59	13.99	49.99	0.24	24.20
G22	40.87	0.30	25.95	3.82	0.62	9.88	35.96	0.22	14.74
G2-2016	68.7%	0.13	33.39	-	-	-	-	-	-

EVALUATION RESULTS OF CHALLENGE C (% RECOG. RATE)

Group	Track 1 Balinese	Track 2 Khmer	Track 3 Sundanese
G20	92.17	97.21	86.54
G21	No Result	No Result	84.98
G1-2016 [3]	88.39	-	-
CNN [24]	85.39	93.96	79.05
HF [24]	85.63	92.44	79.33

EVALUATION RESULTS OF CHALLENGE D (% CER)

Group	Track 1 Balinese	Track 2 Khmer	Track 3 Sundanese	Track 4 Mixed
G3	11.59	4.51	9.68	7.16
G13	9.54	3.38	8.81	5.62
LSTM [24]	39.70	-	-	-

EVALUATION RESULTS OF CHALLENGE B

Group	N	M	o2o	DR	RA	FM
G17	1,271	1,577	962	75.68%	61.00%	67.55%

Evaluation Results of Challenge B per collection

Group	N	M	o2o	DR	RA	FM
G17-Bali	182	264	109	59.89%	41.28%	48.87%
APP-Bali [24]	182	191	164	90.10%	85.86%	87.93%
G17-Khmer	971	1153	778	80.12%	67.47%	73.25%
APP-Khmer [24]	971	990	910	93.71%	91.91%	92.80%
G17-Sunda	118	160	75	63.55%	46.87%	53.95%
APP-Sunda	-	-	-	-	-	-

- ❖ A: G17 uses difference of Gaussian and non-linear enhancement
- ❖ C: G20 uses a very deep convolutional neural network (100 layers) with dense connection
- ❖ D: G13 uses a convolutional neural network encoder and a recurrent neural network decoder equipped with an attention mechanism

The Winners

- ✓ Challenge A: Deepak Kumar, from Department of Electronics & Communication Engineering, Dayananda Sagar Academy of Technology and Management (DSATM), Bengaluru, India.
- ✓ Challenge B: No winner.
- ✓ Challenge C: Zi-Rui Wang, Jun Du and Wen-Chao Wang, from University of Science and Technology of China, China.
- ✓ Challenge D: Jianshu Zhang, Jun Du, and Lirong Dai, from National Engineering Laboratory for Speech and Language Information Processing, University of Science and Technology of China, China.

[1] J.-C. Burie, M. Coustaty, S. Hadi, M.W.A. Kesiman, J.-M. Ogier, E. Paulus, K. Sok, I.M.G. Sunarya, D. Valy, ICFHR 2016 Competition on the Analysis of Handwritten Text in Images of Balinese Palm Leaf Manuscripts, in: 15th Int. Conf. Front. Handwrit. Recognit. 2016, Shenzhen, China, 2016; pp. 596–601. doi:10.1109/ICFHR.2016.107.
 [2] N. Stamatopoulos, B. Gatos, G. Louloudis, U. Pal, A. Alaei, ICDAR 2013 Handwriting Segmentation Contest, in: IEEE, 2013; pp. 1402–1406. doi:10.1109/ICDAR.2013.283.
 [3] https://github.com/tmbdev/ocropy/blob/master/ocropus-errs