

CAMIO: A Corpus for OCR in Multiple Languages

Michael Arrigo¹, Stephanie Strassel¹, Nolan King², Thao Tran², Lisa Mason² Linguistic Data Consortium (1) U.S. Dept. of Defense (2)



Introduction

- CAMIO: Corpus of Annotated Multilingual Images for OCR
- New LDC resource to support the development and evaluation of optical character recognition-related technologies like script ID, language ID, text localization, OCR decoding, keyword search and end-to-end OCR
- Systematically constructed to address gaps in existing corpora*
- CAMIO includes machine-printed text, covering
 - 35 languages across 24 unique scripts, up to 2500 docs/language
 - Scanned docs and images in a broad set of document domains
 - Variety of genres, attributes and scanning/capture artifacts
 - Most documents subject to text localization and reading order
 - Subset of documents in 13 languages also transcribed
- Data to be published in LDC Catalog and used in future evaluations

*Huang, Z., Chen, K., He, J., Bai, X., Karatzas, D., Lu, S., Jawahar, C.V. (2019). ICDAR2019 Competition on Scanned Receipt OCR and Information Extraction. 2019 Internatio Conference on Document Analysis and Recognition (ICDAR), pp. 1516-20.



CAMIO Corpus Pipeline

	LDC Annotator	Crowd Worker	Native Speaker Required
1 Collection	\checkmark	\checkmark	
2 Auditing	\checkmark		if collection by non-native
3 Text Localization	\checkmark		
4 Reading Order	\checkmark		\checkmark
5 Transcription	\checkmark		✓

Each stage is a check on prior stage's quality, with feedback
Additional quality review at end of pipeline



Data Collection and Auditing



CAMIO Languages

- 35 language-script pairs
- 24 unique scripts reflected
- 13 languages selected for transcription, covering 10 scripts

Amharic (Ge'ez)	Hebrew (Hebrew)	Swahili (Latin)
Arabic (Arabic)	Hindi (Devanagari)	Tagalog (Latin)
Armenian (Armenian)	Hungarian (Latin)	Tamil (Tamil)
Bengali (Eastern Nagari)	Japanese (Japanese)	Telugu (Telugu)
Burmese (Burmese)	Kannada (Kannada)	Thai (Thai)
Cambodian (Khmer)	Korean (Hangul)	Tibetan (Tibetan)
Chinese (Simplified)	Malayalam (Malayalam)	Tigrinya (Ge'ez)
Dari (Arabic)	Maldivian (Thanna)	Ukrainian (Cyrillic)
English (Latin)	Oriya (Oriya)	Urdu (Arabic)
Farsi (Arabic)	Pashto (Arabic)	Uyghur (Arabic)
Georgian (Georgian-Mkhedruli)	Russian (Cyrillic)	Vietnamese (Latin)
Greek (Greek)	Sinhalese (Sinhala)	



- Originally planned for 2500 images of machine-printed text collected per language
- Revised to target a variable number of images per language
 - Reflecting data and annotator availability
 - And addition of new annotation tasks (e.g., reading order)
 - As well as prevalence of text-heavy images for some languages

Images	Languages	
2500	Arabic, Chinese, English, Farsi, Hindi, Japanese, Kannada, Korean, Russian, Tamil, Thai, Urdu, Vietnamese	
2000	Amharic, Armenian, Burmese, Dari, Greek, Hungarian, Malayalam, Odia, Pashto, Swahili, Tagalog, Telugu, Ukrainian	
1500	Bengali, Cambodian, Georgian	
1000	Hebrew, Tibetan, Tigrinya	
500	Maldivian, Sinhalese, Uyghur	ſ



- Range of desired document features in collection, plus broad topical variety including both formal and informal content
- Attempted to achieve representation from every feature and plausible feature combination for all languages
 - Not intended to be a fully balanced corpus (cost prohibitive)
- Data scouting assignments specified desired feature combos

Document		Document	
Genre	Document Domain	Attributes	Scanning/Capture Artifacts
Book	Text-heavy documents	Tables	Varied DPI/resolutions
Card/Slide	Unconstrained text	Multi-column	Color/grayscale/black & white
Periodical	Overlaid text	Fielded text	Warping
Record	Diagrams with text	Multi-script	Text runoff
Scene text	Varied content	Multilingual	Occluded text
Webpage		Handwriting	Distance from camera
		Text with images	Perspective
		Other complex layout	Lighting
			Skew, slant, rotation
			Noise

Document Examples







বি, বস্থ, এণ্ড কোম্পানীর

হাতীমার্কা সালসা।

ইহা ঠিক সালসা নহে, তবে সালসা নাম না দিলে. ইহার গুণাবলীর বিষয় কিছুই জনয়জম করিতে সমথ হইবেন না, সেই জন্ত সালসা নাম দিতে হইল। কামরা ইংরাজা-ভাবাপন্ন হইয়া পড়িতেছি, এই আয়ুর্বেষীয় ঔষধেয় নাম তাই বিজাতীয় ভাবায় করিতে বাধ্য হই-লাম—লচে২ উপায় নাই। বলুন হেখি, সোমরস নাম দিলে সাধারণে কি বুকিবেন ?

ৰি, ৰস্থ এণ্ড কোম্পানীয় হাতীমাৰ্কা সালসা

এক মহাতেজ্ঞখন্ধল। উত্তা চৌলবেশ হইতে আনাত কোন লতা-বিশেষের এমন গুণ বে, এ সালসা সেবনের পাঁচ মিনিট পছেই বেহে এবং মনে মহাক্ষৃতি অন্তত্তুত হইবে। মনে গইবে, পরায়ে বেল কোল বৈহাতিক জিয়া নিম্পদ্ধ হইব। এই মহাশকি- য গণিনী সালসা-ত্বণাগনে মন:গ্রাণ খগীয় হথে বিজেয় হইয়া উঠিবে। এ সালসা সহজ পরীয়েও সেবনীয় ! উত্ত গ্রীয়, বর্ষা, পরু, ব্যক্ত-সের্বাঞ্জে সর্বাঞ্জুতে সেবনীয়

न्ताम	ł	•	

		10	ডা:ৰা	
১নং আরপোয়া শিশি	.۲	1.	1.	
২নং একলোমা শিশি	<i>a</i>	se.	ч.	
ন্ম নেড়লোয়া শিশি		1.	2	
জ্যানুপেৰলে নইলে ' মূল্য	-	আনা	ৰা চারি	আন৷ অধিক

পড়ে। তিন বা চারি শিশি খধবা এক ভলন একর লইলে ডাকরাওল



ર્વોત્ત્વભગદેશનાદ્વેવ શુક્રત્વર્હનાનો ત્વેતુપ્રાના ૧૮ તહેર્લ્ડ્સમાં પુત્રભાષ્ટ્રતાર્વક સંભાઈ વૃતાકેલ્વના પ્રત્ય રેલોંત્ પછા વાણેનાના છેલ ગણના વ્યવવારાહેતા વેત્ત્રે આવું પણ બાહુ સંભાગકે વાણે છે. છે કે કે તે કે તે જો





- <u>Method 1</u>: Trained data scouts search the web for existing images (scanned documents or photos) of text in their language on public websites
 - Primary method for most languages
- Method 2: Data from text repositories
 - Designed to address gaps in Method 1
 - Proved logistically difficult, not utilized
- <u>Method 3</u>: Crowdworkers upload existing images in their language
 - Designed to address gaps in Methods 1/2
 - Most fruitful for Bengali, Cambodian, Dari, Hebrew, Maldivian, Sinhalese and Uyghur





- All collected images audited
- By a native speaker if collection by nonnative speaker
- Verify that image is usable and inlanguage (if applicable)
- Verify or provide document feature judgments
- Simple web form





Annotation



Text Localization Task

- Draw 4-point bounding box around each line of machine printed text
- Specify attributes for each box, including the presence of non-target languages, other scripts, or illegible content
- Custom user interface, used for all stages of annotation





Text Localization GUI





Text Localization Guidelines

Avoid whitespace and non-textual boundaries







Reading Order Task

- Top text localization annotators do reading order annotation
- First verify and correct text localization annotation
- Then specify natural, logical reading order by applying next_id tag to each bounding box
 - Last line in doc has next_id = NONE





- Do follow natural, logical flow of content within section
- Do not break flow of text by crossing between sections





- For each of 13 transcription languages, 1250 images selected for transcription
- Native speaker annotators add orthographic transcript for each line of machine printed, bounding boxed text
 - Native orthography matching the image script
 - Retain capitalization and punctuation
 - Follow reading order
 - Transcribe only what's inside bounding box
 - Don't correct errors transcribe as written
 - Add markup for special content and for orientation



Transcription Guidelines





Transcription Markup

- Additional markup for special text features
- Detailed guidelines with examples in many languages for each category

	Lang	Image	Transcription
Wrong Script	Korean	OCR software 는 이 pixel 들로부터	### ### 는 이 ### 들로부터
Unusual Characters	English	Paradife loft.	Paradile lolt. or [[Paradise]] [[lost]].
Need Context	English	Took photoe of the pageds,	Took photos of the ((pagoda)).
Illegible	Hindi	दिन बहुद्धाका पति चल बसा ।	दिन +++ पति चल बसा ।



Annotator Management and Quality Control

- Over 200 staff and contract annotators on CAMIO, plus crowdworkers for data scouting
 - Dedicated team of project assistants to manage workers
- Pre-annotation QC
 - Crowdworker qualification test
 - Annotator aptitude screening, (in person or virtual) training and testing prior to work being assigned
 - Formal guidelines for each task
- Manual quality checks for all languages, tasks and workers (approx. 5%) plus spot checks to identify pervasive issues
- Pipeline design: each stage of annotation checks prior stage
 - Later stages required more skilled, experienced workers

Additional automatic, manual sanity checks by external tee



Challenges and Solutions



Team Management Challenges

- Very large annotation team
 - Increased need for quality control
 - Difficult to maintain consistent communication with contractors
 - Solution: Rely on more project assistants than initially planned to increase communication and boost productivity
- Scarceness of available native speakers and data for some languages
 - Recruitment took longer than expected
 - Solution: Hire additional Penn students who are native speakers



Data Challenges

- Text-heavy images
 - Decreased annotation rate for Text Localization and Reading Order
 - More annotations meant less data
 - Solutions
 - Improvements to annotation tool functionality and userfriendliness
 - Sequestering of very textheavy documents
 - Additional collection to get appropriately sized images into pipeline





Results



- Produced baseline performance results for text localization and OCR decoding
- Using open-source Tesseract OCR engine for a subset of transcribed images
- 50/10/40 train/validation/test split
 - Nominally 625 train, 125 validation, and 500 test images per language

*Smith, R., Antonova, D., Lee, D. (Jul 25, 2009) Adapting the Tesseract open source OCR engine for multilingual OCR. Proceedings of the International Workshop on multilingual ocr, pp. 1-8



- Tesserocr python wrapper for Tesseract 4.0.0
- Comparing Tesserocr SINGLE_BLOCK output with ground truth
- Intersection over union (IoU) measures overlap between bounding boxes
- Overlap between Tesseract text boxes and ground truth annotation

CAMIO Test Set	F1 Score	Precision	Recall
Kannada	36.90%	30.40%	47.20%
Tamil	30.40%	26.20%	36.30%
Arabic	27.50%	22.60%	35.30%
Urdu	27.50%	22.90%	34.50%
Korean	27.00%	21.90%	35.40%
Hindi	25.20%	22.10%	29.40%
Farsi	23.80%	18.50%	33.40%
Russian	19.90%	18.20%	22.00%
Vietnamese	18.30%	14.00%	26.50%
English	17.30%	15.20%	20.10%
Thai	16.00%	15.00%	17.10%
Chinese (Simplified)	15.60%	13.00%	19.30%
Japanese	13.50%	11.30%	16.90%

Precision: ground truth boxes with IoU of at least 0.5 divided by total number of Tesseract boxes found

Recall: Tesseract boxes with IoU of at least 0.5 divided by number of ground truth boxes

F1: harmonic mean of precision and recall 2* *Precision* * *Recall /* (*Precision* + *Recall*)



OCR Decoding Results

- Character Error Rate (CER) based on Levenshtein edit distance
- Minimum number of edits per line to modify system output to match the ground truth, divided by the number of characters in the ground truth
- Significant performance degradation possible with CER > 5%*

Normalization Procedure	CAMIO Test Set	Character Error Rate (CER)	Likely due to		
1. Conversion of characters to	Urdu	68.20%	fonts like		
lower-case	Japanese	39.80% 🗨	Nastaliq		
2. Removal of extra spaces	Chinese (Simplified)	34.90%			
(down to single spaces)	Tamil	24.50%			
2. Demoval of numerican	Arabic	24.00%			
3. Removal of punctuation	Farsi	23.00%			
4. Apply Unicode compatibility	Korean	18.90%			
decomposition, followed by	Thai	18.10%	LIKEIY due to		
canonical composition (NFKC)	Vietnamese	17.00%	vertical text		
	Hindi	16.40%			
*Bazzo, G.T., Lorentz, G.A., Suarez Vargas, D., Moreira, V.P. (2020). Assessing the Impact of	Russian	15.30%			
OCR Errors in Information Retrieval. In Advances in Information Retrieval, ECIR 2020	Kannada	12.40%			
Lecture Notes in Computer Science, vol 12036.	English	12.10%			
CAMIO:A Corpus for OCR in Multiple Languages – LREC 2022, Marseille, France – June 21-23, 2022					



Corpus Release

- Image data in two versions: original (any format) and converted (PNG)
- Annotations in unified XML format
- Guidelines and documentation





Conclusion

- Corpus of Annotated Multilingual Images for OCR (CAMIO) is a new resource for computer vision and document analysis
- Poor baseline performance in text localization and OCR reflects challenging data that will stimulate future research
- CAMIO is expected to appear in LDC's catalog starting later this year
 - Part 1: Transcribed languages
 - Part 2: Untranscribed languages
 - Some data will remain unpublished for use in future evaluation campaigns
- More annotation planned including doc zoning, translation, transcription

	Collection/	
CAMIO Corpus	Annotation	Transcription
Languages	35	13
Total images	69,440	15,724
Boxed images	59,982	15,724
Boxes	2,339,358	311,619
Average boxes per image	~39	~20
Tokens	2,352,411	2,352,411