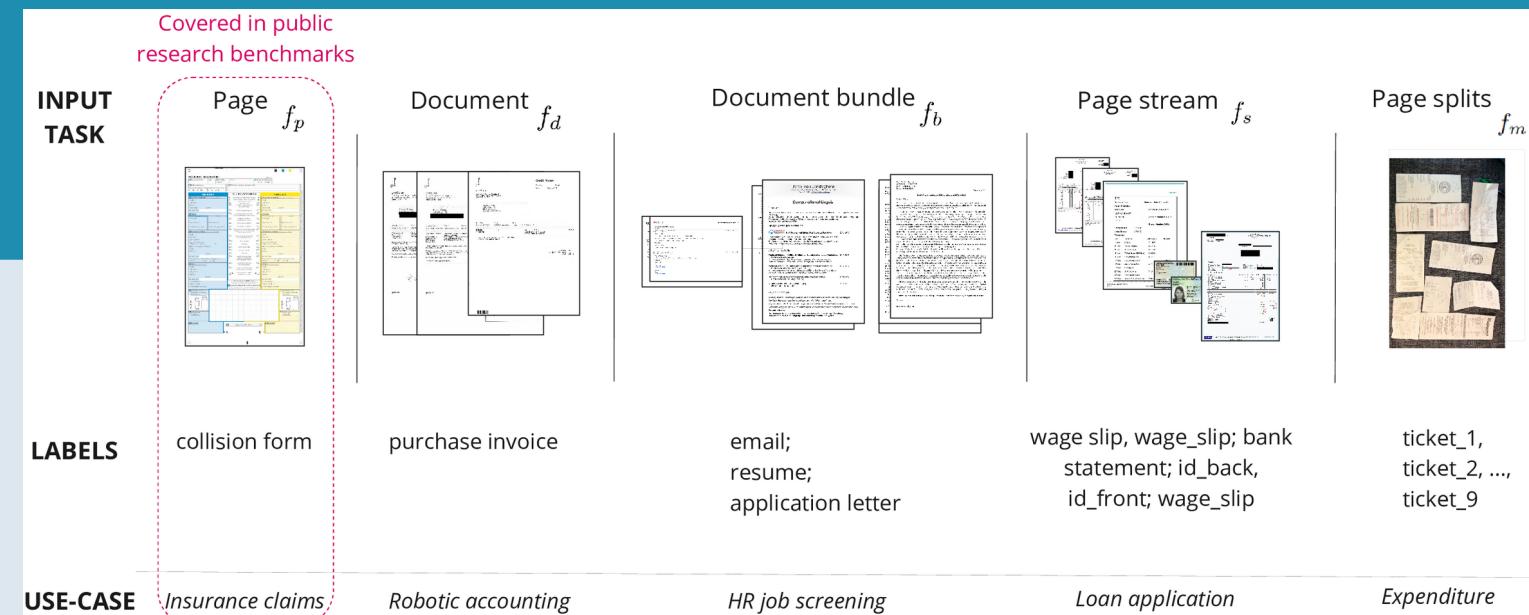


Beyond Document Page Classification: Design, Datasets, Challenges



Jordy Van Landeghem, Sanket Biswas,
Matthew Blaschko, Marie-Francine Moens



Beyond Document Page Classification: *TLDR*;

Move toward evaluation on complete document inputs, as they prevalently occur (**multi-page documents, bundles, page streams, and splits**) across various practical scenarios within real-world DC applications

Covered in public research benchmarks					
INPUT TASK	Page f_p	Document f_d	Document bundle f_b	Page stream f_s	Page splits f_m
LABELS	collision form	purchase invoice	email; resume; application letter	wage slip, wage_slip; bank statement; id_back, id_front; wage_slip	ticket_1, ticket_2, ..., ticket_9
USE-CASE	Insurance claims	Robotic accounting	HR job screening	Loan application	Expenditure

Beyond Document Page Classification: scope

Problem – Observation - Contribution

*The practical task of **long document understanding** is largely underexplored due to several challenges in computation and how to efficiently represent long multimodal input.*

*Document classification should be **evaluated closer to practical applications** and task formulations (documents, bundles & page streams).*

We contribute **novel methodology** and demonstrate that document classification can be considered far from solved.

Beyond Document Page Classification: *goals*

Position paper with main objectives:

- I. Benchmark closer to applied document classification scenarios
- II. Experimental study on multi-page inference methods
- III. Reflect on evaluation practices & moving beyond *iid* test sets
- IV. Propose guidelines to foster document dataset construction efforts

Contributions

1. Redesign and formalization of multi-page DC scenarios
2. Construction of two novel datasets **RVL-CDIP_MP** and **RVL-CDIP-N_MP**
3. Comprehensive experimental analysis of the novel datasets
4. Survey of challenges stalling document classification progress

à huggingface.co/datasets/bdpc/rvl_cdip_mp
à huggingface.co/datasets/bdpc/rvl_cdip_n_mp



I. Problem formulation

Classification tasks

Page classification



$$f_p : \mathcal{X}_p \rightarrow \mathcal{Y},$$

where $\mathcal{Y} = [C]$ for C mutually exclusive categories.

Document classification



$$f_d : \mathcal{X}_d \rightarrow \mathcal{Y},$$

where $\mathcal{Y} = [K]$ for K mutually exclusive categories.

Document bundle classification



$$f_b : \mathcal{X}_b \rightarrow \mathcal{Y}, \text{ where } \mathcal{Y} \text{ is a product space of } B \text{ documents,}$$

$$\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_B, \text{ with } \{\mathcal{Y}_j = [K] : j \in [B]\}.$$

Page stream classification



$$f_s : \mathcal{X}_d \rightarrow \mathcal{Y}, \text{ where } \mathcal{Y} \text{ is a product space of } L \text{ pages,}$$

$$\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_L, \text{ with } \{\mathcal{Y}_j = [C] : j \in [L]\}.$$

Page splitting



$$f_m : \mathcal{X}_p \rightarrow \mathcal{Y}, \text{ where } \mathcal{Y} = \mathbb{Z}^C.$$



II. Balancing research & applications

DU datasets and document sources

Dataset	Size	Data Source	Domain	Task	OCR	Layout
IIT-CDIP [28]	35.5M	UCSF-IDL	Industry	Pretrain	✗	✗
RVL-CDIP [17]	400K	UCSF-IDL	Industry	DC	✗	✗
RVL-CDIP-N [25]	1K	Document Cloud	Industry	DC	✗	✗
TAB [36]	44.8K	UCSF-IDL	Industry	DC	✗	✗
FUNSD [21]	199	UCSF-IDL	Industry	KIE	✓	✗
SP-DocVQA [35]	12K	UCSF-IDL	Industry	QA	✓	✗
OCR-IDL [6]	26M	UCSF-IDL	Industry	Pretrain	✓	✗
FinTabNet [58]	89.7K	Annual Reports S&P	Finance	TSR	✗	✓
Kleister-NDA [47]	3.2K	EDGAR	US NDAs	KIE	✓	✗
Kleister-Charity [47]	61.6K	UK Charity Commission	Legal	KIE	✓	✗
DeepForm [48]	20K	FCC Inspection	Forms broadcast	KIE	✓	✗
TAT-QA [61]	2.8K	Open WorldBank	Finance	QA	✓	✗
PubLayNet [60]	360K	PubMed Central	Scientific	DLA	✗	✓
DocBank [30]	500K	arxiv	Scientific	DLA	✓	✓
PubTabNet [59]	568K	PubMed Central	Scientific	TSR	✗	✓
DUDE [53]	40K	Mixed	Multi-domain	QA	✓	✗
Docile [45]	106K	EDGAR & synthetic	Industry	KIE	✓	✗
CC-PDF [51]	1.1M	Common-Crawl (2010-22)	Multi-domain	Pretrain	✗	✗

DC datasets - lacking variety, scale and multi-page

Proposed multi-page datasets

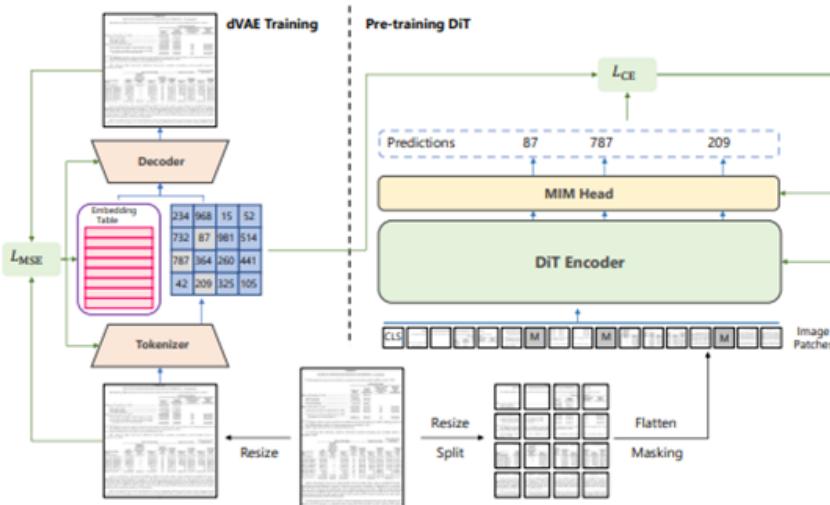
Dataset	Purpose	#d	#p	$ \mathcal{Y} $	Language	Color depth
NIST [8]	f_s		5590	20	English	Grayscale
MARG [32]	f_s		1553	2	English	RGB
Tobacco-800 [62]	f_s		800	2	English	Grayscale
TAB [36]	f_s		44.8K	2	English	Grayscale
Tobacco-3482 [23]	f_p		3482	10	English	Grayscale
RVL-CDIP [17]	pre-training, f_p		400K	16	English	Grayscale
RVL-CDIP-N [25]	f_p , OOD		1002	16	English	RGB
RVL-CDIP-O [25]	f_p , OOD		3415	1	English/Mixed	RGB
RVL-CDIP_MP	f_d	$\pm 400K$	$\mathbb{E}[L] = 5$	16	English	Grayscale
RVL-CDIP-N_MP	f_d , OOD	1002	$\mathbb{E}[L] = 10$	16	English	RGB



III. Multipage inference experiments

Tease some issues and strategies when naively scaling beyond page-level DC.

Model: Document Image Transformer

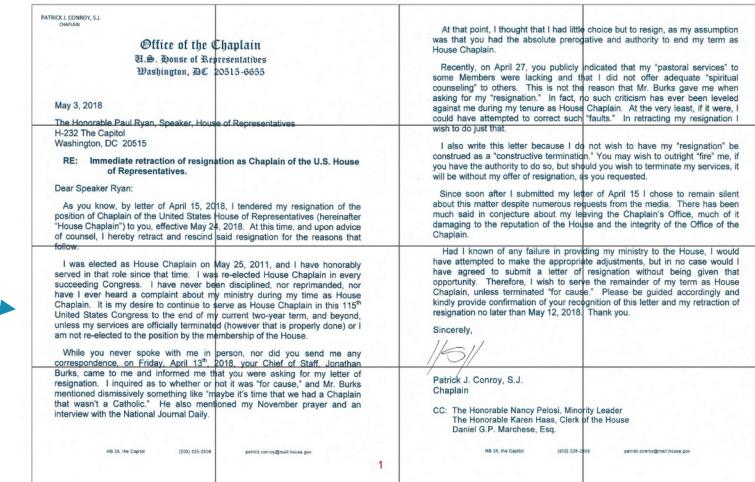


Inference	Strategy	Scope
sample	first	page
	second	page
	last	page
	max confidence	page
	soft voting	page
	hard voting	page
sequence	grid	document
	document	document
grid document	grid	document
	(not tested)	document

$$\text{MaxConf}(x, y) = \arg \max_{\substack{l \in [L] \\ k \in [K]}} [\tilde{f}_p(x, y)]_k^l$$

$$\text{SoftConf}(x, y) = \arg \max_{k \in [K]} \sum_{l=1}^L [\tilde{f}_p(x, y)]^l$$

$$\text{HardVote}(x, y) = \arg \max_{k \in [K]} \sum_{l=1}^L e_{\hat{y}^l},$$





III. Results

Strategy	Acc↑	F1↑	$F1_M \uparrow$	ECE↓	AURC↓
$f_p\$$ [29]	93.345	93.351	93.335	0.075	0.010
first	91.291	91.286	91.271	0.073	0.014
second	87.295	87.305	87.277	0.070	0.029
last	85.091	85.060	85.028	0.072	0.038
MaxConf	91.407	91.453	91.344	0.124	0.006
SoftVote	91.220	91.185	91.236	0.134	0.004
HardVote	85.995	86.182	85.781	0.085	0.018
grid	72.642	72.045	73.266	0.109	0.042

Table 4. Base classification accuracy of DiT-base [29] (finetuned on RVL-CDIP) evaluated on the test set of RVL-CDIP_MP per baseline f_d strategy. Best results per metric are boldfaced. \$ refers to our reproduction of results.

Hypothesis: Summary-detail document construction

Inefficient (L pages) and dependent on calibration of f_p

Dataset	Strategy	Acc↑	Δ
RVL-CDIP_MP	first+second ^(*)	93.795	2.504
	first+last ^(*)	93.675	2.384
	second+last ^(*)	89.709	-1.583
	first+second/last ^(*)	94.454	3.163
RVL-CDIP_N_MP	first+second ^(*)	83.638	4.878
	first+last ^(*)	83.130	4.370
	second+last ^(*)	71.545	-7.215
	first+second/last ^(*)	84.553	5.793

Table 6. Best-case classification accuracy indicated with ^(*) when combining 'knowledge' over different pages. Δ refers to the absolute difference with the first page only.

Proof of concept

Multi-page document representations are promising for improving document classification



IV. Challenges and Guidelines

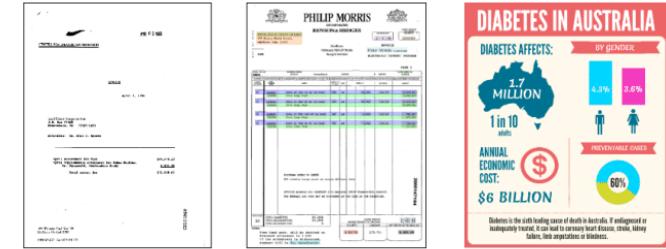
(1) Divergence of tasks

- Only page classification (f_p) research benchmarking
- DC task formulations can help to formulate f in practice

(2) Divergence of label space

- Too simplified label sets in benchmarking
 - $K=16 \leftrightarrow$ industry reqs: $K \sim 50-400$
- Label noise
 - Relabeling campaigns
- MECE principle for construction of doc-level label sets

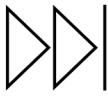
(3) Divergence of input data



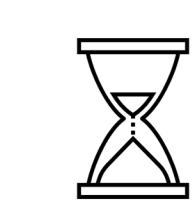
Public document collections
Data synthesis
Anonymization

(4) Maturity of evaluation methodology

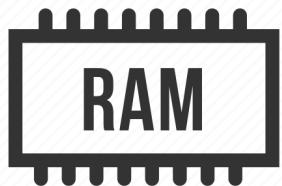
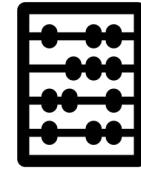
- iid assumption is unmaintainable in production



IV. (4) Evaluation beyond accuracy and *iid* setting



calibration



Covariate shifts

CALIFORNIA MINI GUIDE
Your essential guide to travelling in California from The World Was Here First



WHAT TO KNOW BEFORE VISITING CALIFORNIA

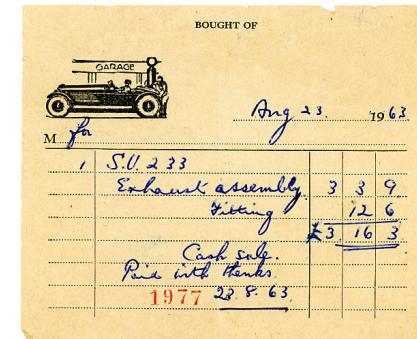
California is one of the most visited states in the US and it packs so much into it that it's hard to see where to see what. With a population larger than those of Canada and Australia and with the world's fifth-largest economy, California can feel like its own country. With a landscape as diverse as its people and cuisine, the perfect California trip tops every traveller's wish list.

A trip to California is an experience that every traveller should have at least once in their lifetime. If you're planning on visiting the Golden State, follow this advice to get you started.

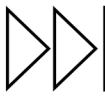
Near OOD



Subclass shift



Concept drift

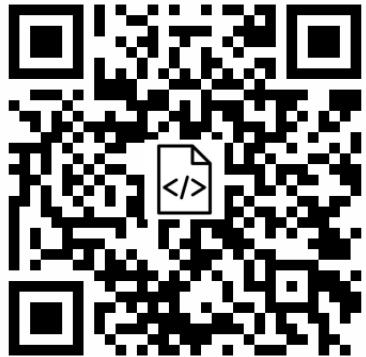


Closing remarks

- Pivotal step forward in establishing methodology for multipage DC
- Covered challenges and limitations hindering progress
- Experimental study shows promise from advancing multi-page document representations and inference
- Recommendations for future DC dataset constructions efforts
 - Type and nature of document data
 - Variety and quality of the classification label set
 - Focus on DC scenarios closer to applications

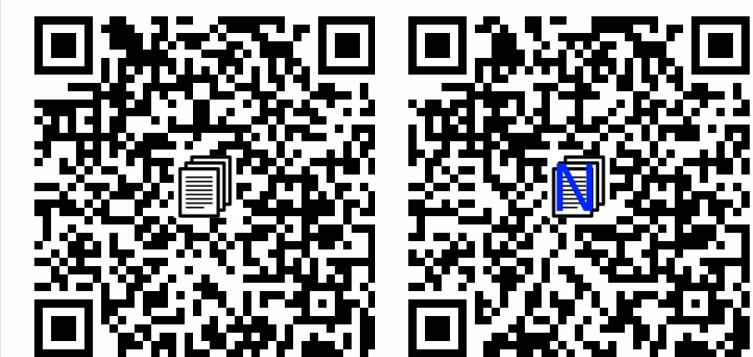
Questions? Future collaborations?

[CODE](#)



huggingface.co/bdpc/src

[DATASETS](#)



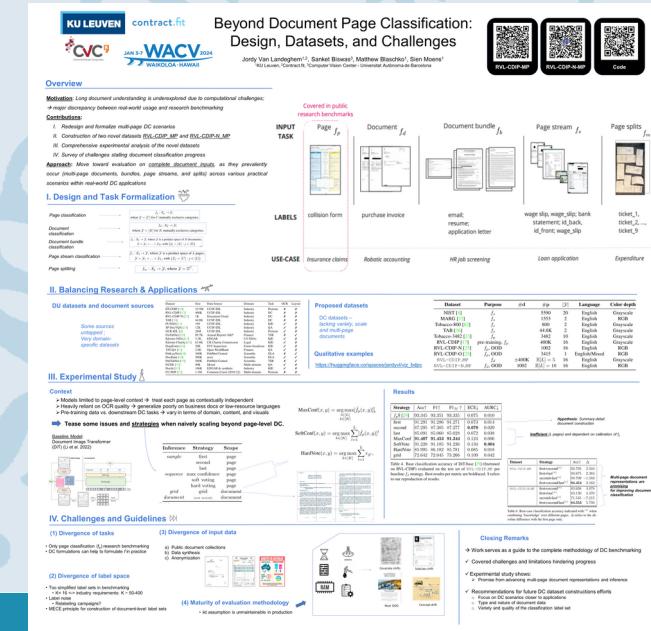
huggingface.co/datasets/bdpc/rvl_cdip_mp
huggingface.co/datasets/bdpc/rvl_cdip_n_mp



See you at WACV 2024!



jordy@contract.fit



The document is a comprehensive report on document page classification. It starts with an overview of the field, highlighting the need for better real-world understanding due to computational challenges. It details the design and task formalization, including input tasks like page classification, document identification, and document bundle identification; labels such as collision form, purchase invoice, email, resume, application letter, wage slip, wage_slip_bank statement, etc.; and use-cases like insurance claims, robotic accounting, HR job screening, and loan applications. A significant part of the document is dedicated to datasets, listing various datasets with their sizes, languages, and characteristics. It also provides qualitative examples and experimental study results, including tables for document counts and language distribution. The report concludes with challenges and guidelines for future research.