

# Layout-Aware Zero-Shot Learning for Visual Document Matching

## Qualificação de Mestrado

Lucas de Almeida Bandeira Macedo

Universidade de Brasília  
Departamento de Ciência da Computação

Orientador: Prof. Dr. Pedro Garcia Freitas  
Coorientador: Prof. Dr. Bruno Luiggi Macchiavello Espinoza

Outubro de 2025

1 Introdução

2 Metodologia

3 Resultados

4 Conclusão

# Introdução

# Contexto - Documentos e Compliance

- Ambiente Bancário
- Documentos físicos
- Imagens de documentos
- Exemplos:

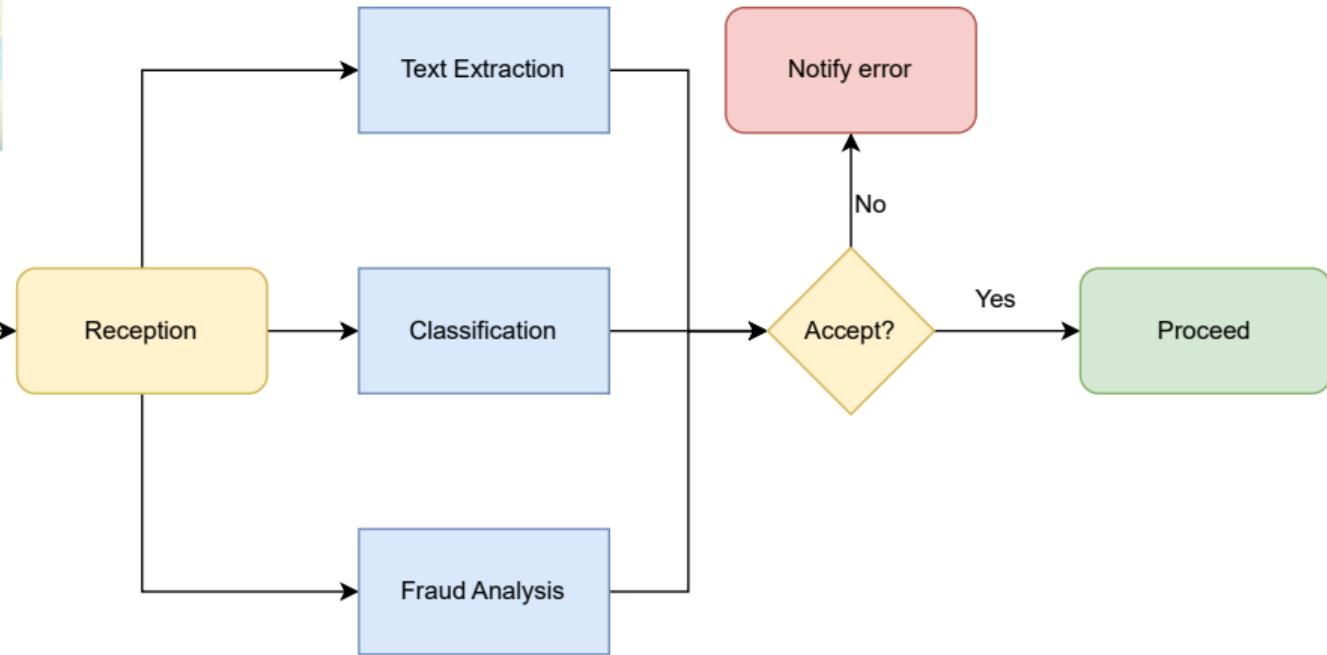


RECEBIMENTO DE ENTREGA DA DECLARAÇÃO DE AJUSTE ANUAL - OPÇÃO PELO DESCONTO SIMPLIFICADO DECLARAÇÃO ORIGINAL			
MINIST&P;RIO DA FAZENDA SECRETARIA DA RECEITA FEDERAL DO BRAS&P;L EXERCIC&P;O 2017 ANO-CALENDÁRIO 2016			
IDENTIFICAÇÃO DO DECLARANTE			
CPF do declarante 352.651.378-08	Nome do declarante ERIKA TOMAZELLA	Telefone (11) 43791221	
Endereço RUA RUA JUSTINO ALVES BATISTA		Número 99	Complemento AP-64 BL P
Bairro/Localidade VILA YOLANDA	CEP 06126-120	Município DSASCO	UF/SP SP
(Valores em Reais)			
TOTAL RENDIMENTOS TRIBUTÁVEIS		62.200,42	
IMPOSTO DEVIDO		3.562,59	
IMPOSTO A RESTITUIR		781,16	

# Contexto - Fluxo de Compliance



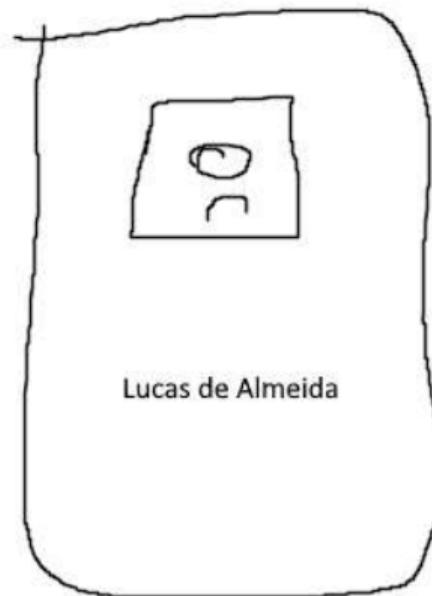
Actor



- Assegurar que o documento está correto
- Documentos não-digitais
- Evita fraudes

# Contexto - Classificação de Imagem

- Assegurar que o documento está correto
- Documentos não-digitais
- Evita fraudes



## Classificação Tradicional:

- Categorização em classes predefinidas
- Cross-Entropy Loss

## Desempenho Atual:

- Bakkali et al. (2021): 97.70% de acurácia no RVL-CDIP

## Classificação Tradicional:

- Categorização em classes predefinidas
- Cross-Entropy Loss

## Desempenho Atual:

- Bakkali et al. (2021): 97.70% de acurácia no RVL-CDIP

## O Problema:

- Novos layouts de documentos
- Classes completamente novas
- Necessidade de retreinamento
- Semanas/meses de engenharia de dados e treinamento

# Zero-Shot Learning

Permite que o modelo reconheça elementos de classes nunca vistas no treinamento

## Desafios

- Falta de dataset especializado
  - Imagens de Documento
  - Generalização
  - Divisão treino e teste zero-shot
- Ausência de metodologia estado-da-arte
  - Paradigma ZSL
  - Capacidade de classificar

## Contribuições

### ① Novo dataset LA-CDIP

- Classificação ZSL
- Derivado do RVL-CDIP

### ② Abordagem de Visual Document Matching (VDM)

- Similaridade de documentos
- Metric Learning
- Generalização Zero-Shot

### ③ Avaliação sistemática

- Benchmark extensivo
- Comparação com LLM

# Metodologia

# Dataset - Motivação

## Datasets Disponíveis

- PubLayNet, DocLayNet
- DocVQA
- CORD, SROIE
- RVL-CDIP

## Towards Zero-Shot Document Image Classification

Lucas Macedo, João Paulo Costa, João Pedro Felix de Almeida, Pedro Freitas, and Li Weigang  
*Department of Computer Science, University of Brasília, Brasília, Brazil*

**Abstract**—Classification is a fundamental tool to automate the process of categorizing documents in many real-world applications, such as information management, financial document fraud detection, regulatory compliance, and many others. Because of this broad spectrum of applications, document classification is one of the most active research areas in computer vision. However, documents often change in terms of format and their visual patterns, which may impact a simple classification model. Moreover, most companies are retraining often domain-specific classifiers, which is time-consuming and costly. This work is demanding new data. Therefore, techniques capable of classifying documents by simply observing new data, without necessarily requiring annotated examples, are highly desired. This is the case for a wide variety of applications. In this context, Zero-Shot Learning (ZSL) is especially suitable for document classification because it can handle new categories without retraining. In this work, we tackle the gap involving Zero-Shot Document Image Classification (ZS-DIC), where we classify documents that have not been seen by the model during training. To address this, we built LayDoc, a Zero-Shot Document Image Classification Processing (LA-CDIP), a dataset tailored for this problem. LA-CDIP prioritizes structural consistency, allowing models to classify documents from different domains. To evaluate the performance of this dataset, we developed a series of Siamese Neural Networks (SNNs) based on a variety of common neural architectures, such as ResNet, EfficientNet, VGG, and DenseNet. As a result, our proposed ZSL-based method achieves Equal Error Rates (EER) no higher than 5%. The code of the proposed method is available at <https://github.com/AlMHD/la-CDIP>.

To overcome the limitations imposed by the need for labeled training data, zero-shot learning paradigm

This paradigm defines a task where the classes in the test set are disjoint from those in the training set, thus forcing the model to generalize effectively to novel categories [3], [6]. Despite being a considerably under-researched area within document analysis, ZS-DIC provides valuable contributions.

Although RVL-CDIP is the most widely used dataset for document classification, it is not suitable for ZS-DIC.

DocVQA [7] focuses on Visual Question Answering (VQA).

DocLayNet [8] is specialized in layout segmentation, and

CORD [9], SROIE [10], FUNSD [11], and XFLUND [12] were created for information extraction. While these datasets advance their respective fields, their label characteristics prevent them from being used for ZS-DIC.

To address these challenges, we introduce LA-CDIP, a dataset designed for the ZS-DIC challenge. LA-CDIP is a document-image dataset that re-purposes the images from RVL-CDIP under a new classification scheme. Unlike RVL-CDIP, which classifies documents by their purpose (e.g., emails, letters), our dataset arranges them by their visual patterns, grouping documents with similar visual patterns together. This reorganization enables ZS-DIC by ensuring that each class exhibits high consistency in both visual and textual patterns.

We propose a Visual Document Matching (VDM) framework for this dataset. This framework leverages Contrastive Learning (CL) to generate a general representation for two documents. A pre-defined similarity threshold is used to determine if documents belong to the same class. In a practical application, by providing a reference for each class, the framework can be deployed in two modes: as a binary classifier for verification (comparing each incoming document with a single class reference) or as a multi-class classifier for identification (comparing each incoming document with all available references to select the most similar ones).

The remainder of this paper is structured as follows. Section II details our LA-CDIP dataset. Section III describes the proposed Siamese network modeling. Sections IV and Section V present the experimental results and conclusions, respectively.

978-1-6654-8951-6/23/\$31.00 ©2023 IEEE

# Dataset - Motivação

## Datasets Disponíveis

- PubLayNet, DocLayNet
- DocVQA
- CORD, SROIE
- RVL-CDIP

### Towards Zero-Shot Document Image Classification

Lucas Macedo, João Paulo Costa, João Pedro Felix de Almeida, Pedro Freitas, and Li Weigang  
*Department of Computer Science, University of Braga, Braga, Brazil*

**Abstract**—Classification is a fundamental tool to automate the process of categorizing documents in many real-world applications, such as information management, financial document processing, document retrieval, document classification, fraud detection, regulatory compliance, and many others. Because of this broad spectrum of applications, document classification has been a subject of extensive research. However, documents often change in terms of format and their visual pattern, which may impair a simple classification model. Moreover, model customization and retraining often demands large efforts, which is impractical for many companies dealing with new data. Therefore, techniques capable of classifying documents by simply observing new data, without necessarily re-training the model, are highly demanded. In this context, Zero-Shot Learning (ZSL) is especially suitable for document classification because it can handle new classes without training on them. In this work, we tackle the gap involving Zero-Shot Document Image Classification (ZS-DIC), where we classify documents that have not been seen by the model during training. To achieve this, we built a novel dataset, called Document Image Processing (LA-CDIP), a dataset tailored for this problem. LA-CDIP prioritizes structural consistency, allowing models to classify unseen documents based on their visual features. To build this dataset, we developed a series of Siamese Neural Networks (SNNs) based on a variety of computer vision neural architectures, namely, ResNet, DenseNet, and EfficientNet. The proposed ZSL-based method achieves Equal Error Rates (EER) under 5%. The code of the proposed method is available at <https://github.com/LMBM/la-cdip>.

**dataset [2].** This dataset comprises 300,000 documents divided into 16 classes based on their purpose, such as letters and emails. RVL-CDIP has been a subject of extensive research for traditional classification, with some works achieving over 97% accuracy [3], [4] in specific tasks. However, existing datasets do not consider changes in format and appearance, and new document types are excessively generated. In such dynamic environments, new data must be acquired to retrain the classification model.

To overcome the limitations imposed by the need for labeled training data, ZSL emerges as the most appropriate paradigm. This approach allows learning a model that can classify new data that are disjoint from those in the training set, thus forcing the model to generalize efficiently to novel categories [5], [6]. Despite being a considerably under-researched area within document analysis, ZS-DIC provides valuable contributions.

Although RVL-CDIP is the most widely used dataset for document classification, other datasets such as DocVQA [7] focused on Visual Question Answering (VQA), DocLayNet [8], SROIE [10], FUNSD [11], and XFUND [12] were created for information extraction. While these datasets advance their respective fields, their label characteristics prevent their application to ZS-DIC.

To address these shortcomings, we introduce LA-CDIP, a dataset designed for the ZS-DIC challenge. LA-CDIP is a document-image dataset that re-purposes the images from RVL-CDIP under a new classification scheme. Unlike RVL-CDIP, which classifies documents by their purpose (e.g., emails, letters), our dataset arranges them by their visual features, grouping together documents from a same class. The resulting structure enables ZS-DIC by ensuring that each class exhibits high consistency in both visual and textual patterns.

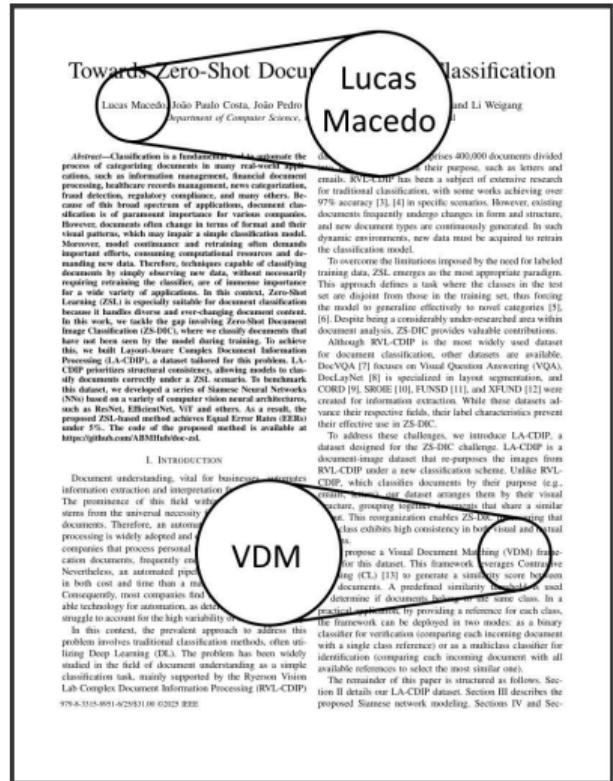
We propose a Visual Document Matching (VDM) framework for this dataset. This framework leverages Contrastive Learning (CL) [13] to generate visual prototypes for each document. A pre-trained similarity threshold is used to determine if documents belong to the same class. It is a practical application, by providing a reference for each class, the framework can be deployed in two modes: as a binary classifier for verification (comparing each incoming document with a single reference selected from a multi-class classifier for document verification) or as a search engine for document retrieval (comparing each incoming document with all available references to select the most similar ones).

The remainder of this paper is structured as follows. Section II details our LA-CDIP dataset. Section III describes the proposed Siamese network modeling. Sections IV and V conclude this paper.

# Dataset - Motivação

## Datasets Disponíveis

- PubLayNet, DocLayNet
- DocVQA
- CORD, SROIE
- RVL-CDIP



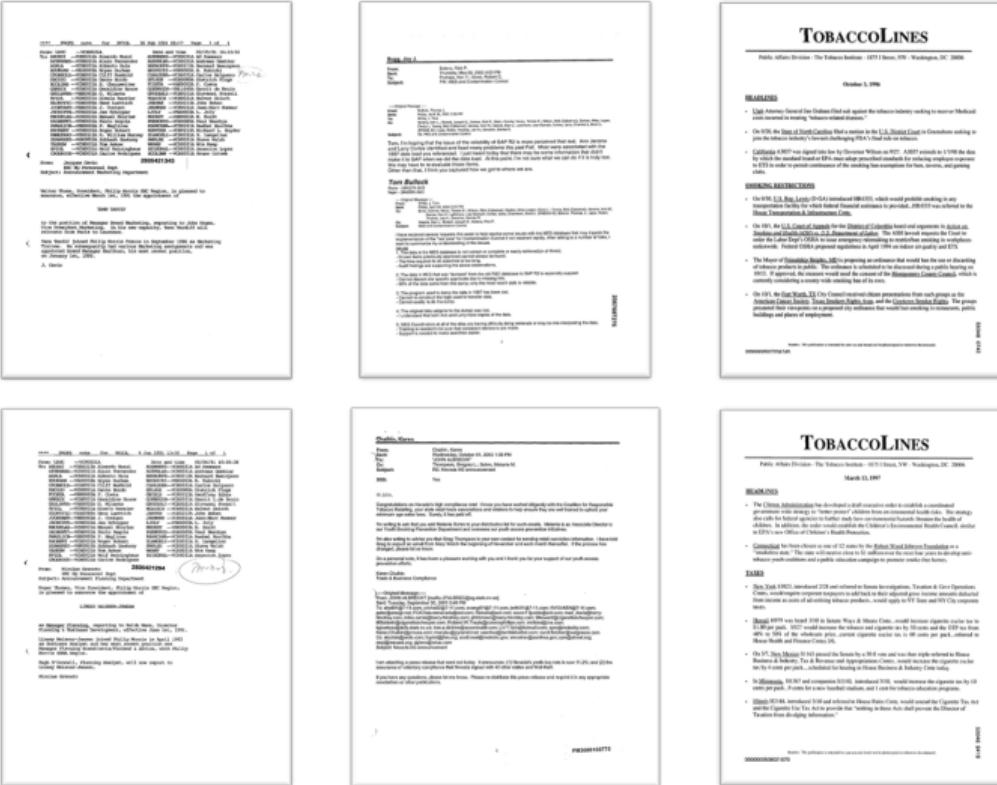
## Datasets Disponíveis

- PubLayNet, DocLayNet
- DocVQA
- CORD, SROIE
- RVL-CDIP

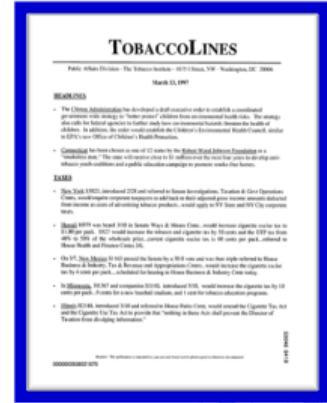
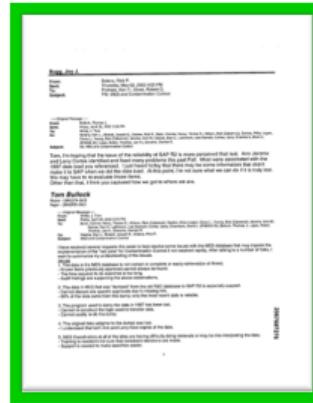
## RVL-CDIP

- 400.000 documentos
- 16 classes
- email, formulário, carta...
- Separado por função

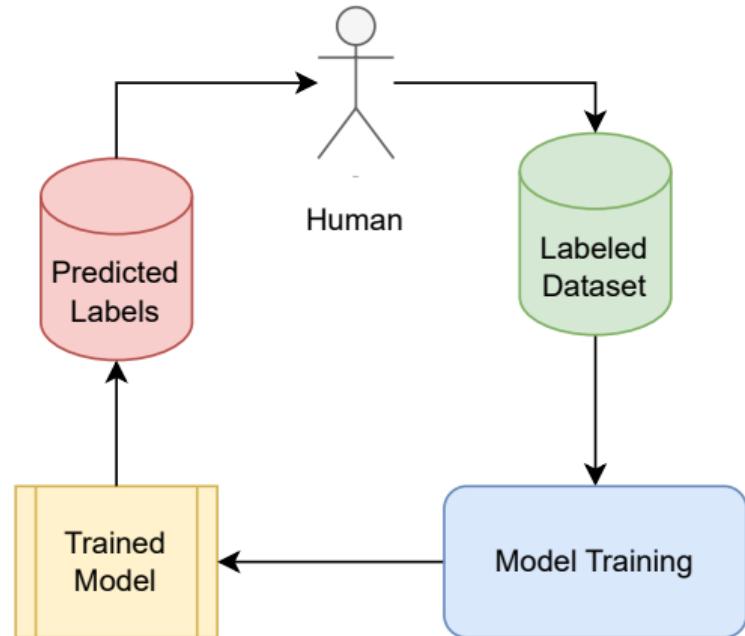
# Dataset - Motivação

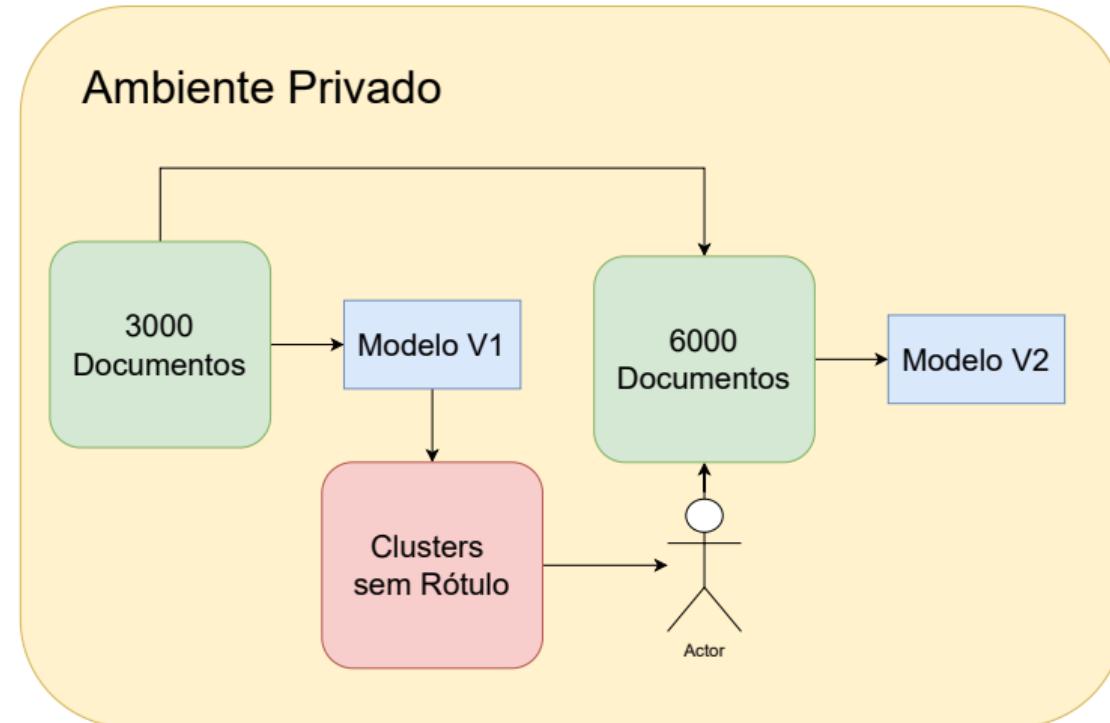


# Dataset - Motivação

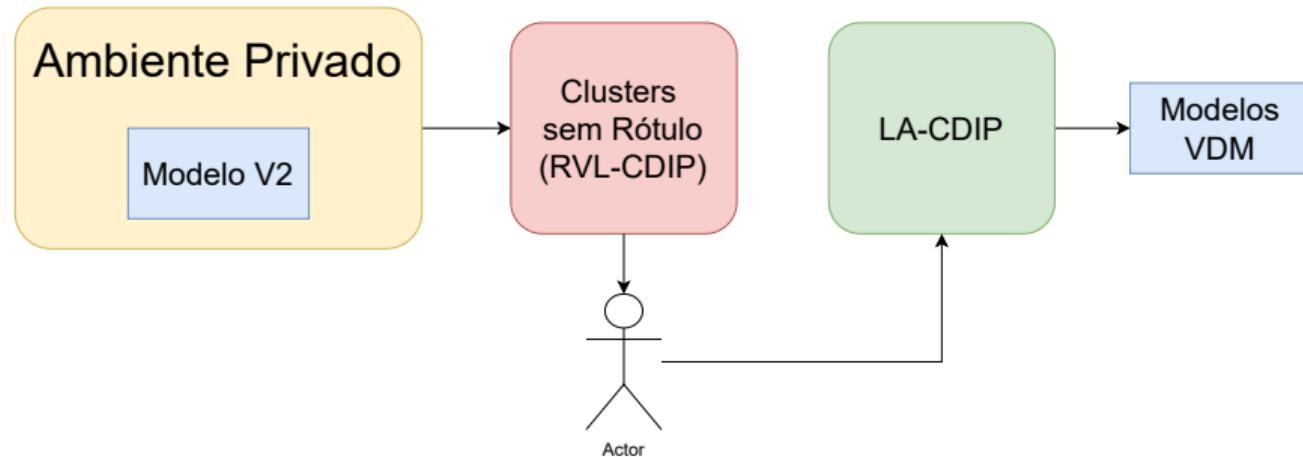


# LA-CDIP Dataset - Active Learning

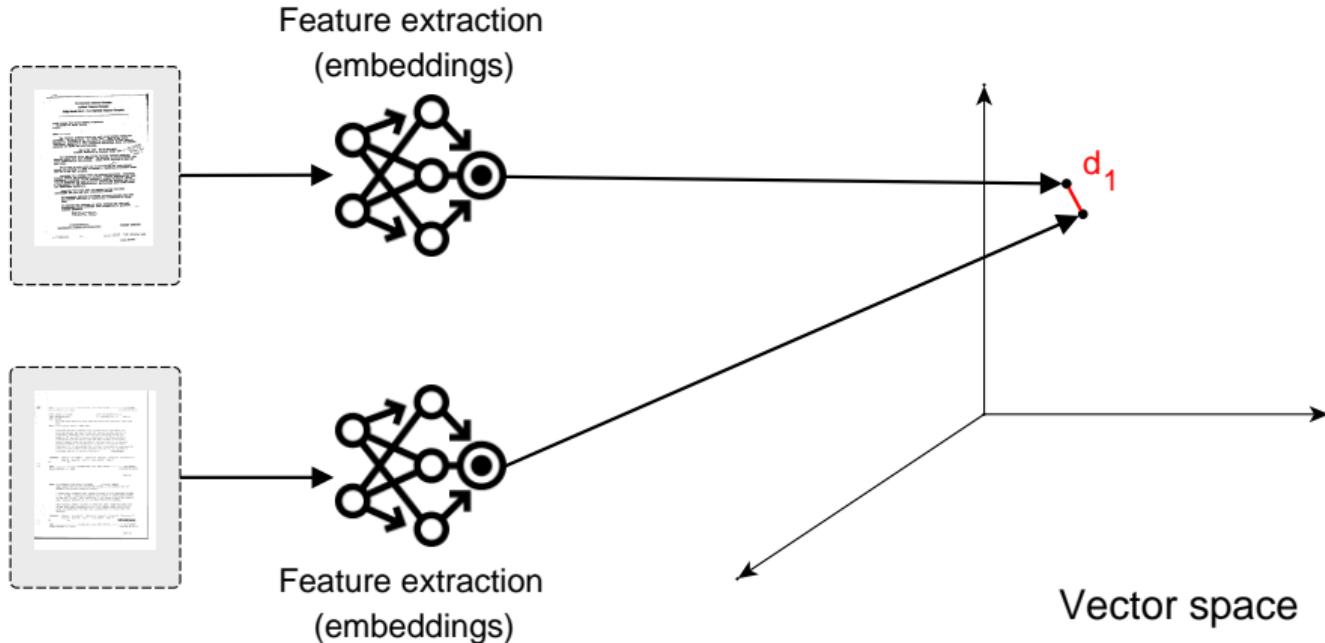




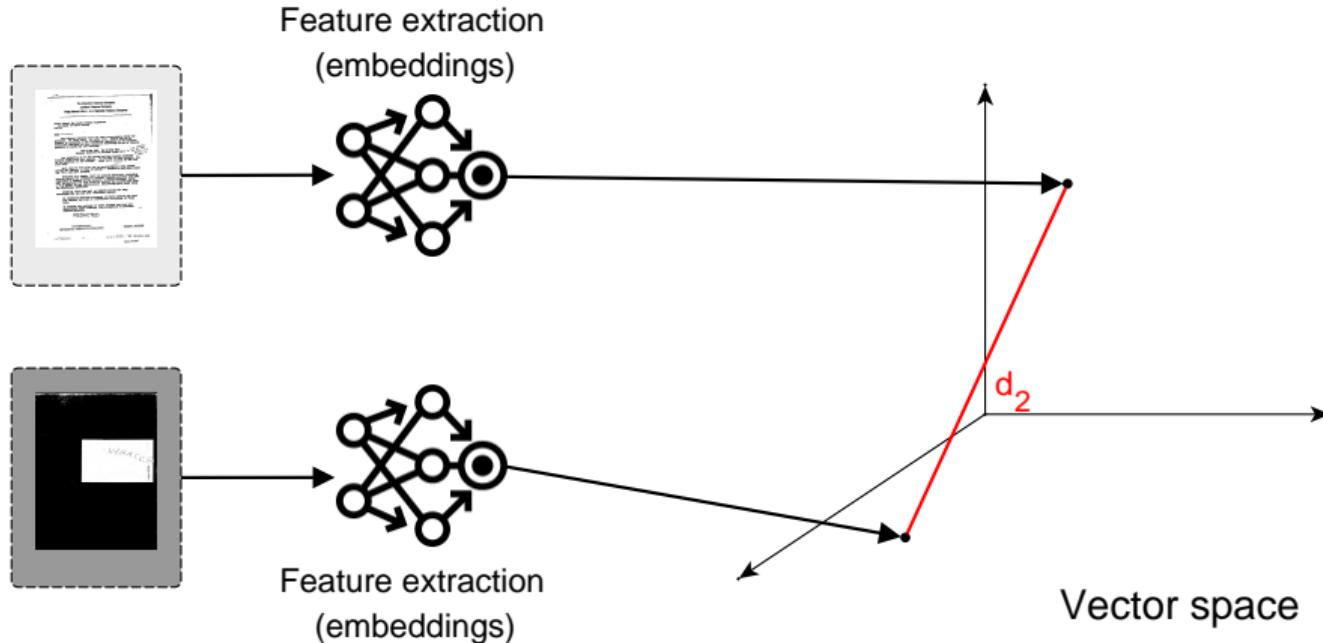
# LA-CDIP Dataset - Rotulação



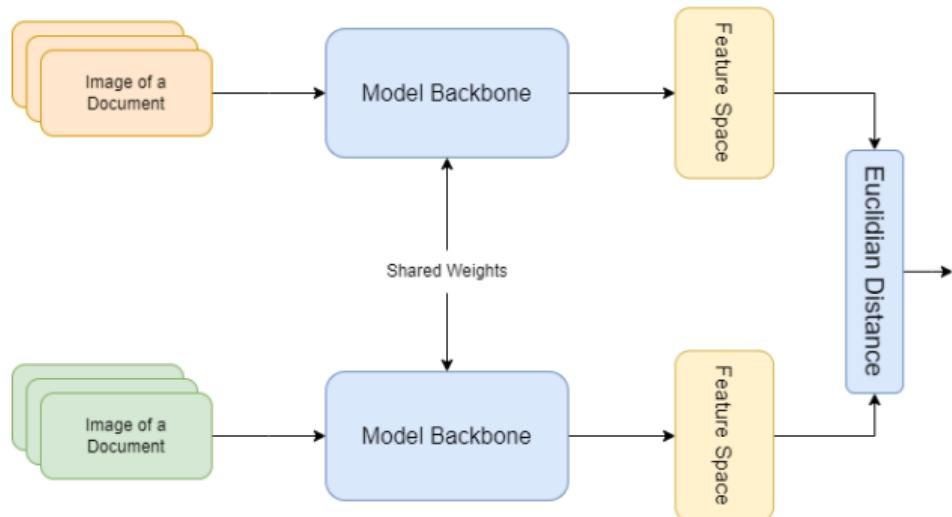
# Visual Document Matching



# Visual Document Matching



# Visual Document Matching - Arquitetura



$$\mathcal{L} = \begin{cases} d(x_1, x_2)^2, & \text{if } y = 1 \\ \max(0, m - d(x_1, x_2))^2, & \text{otherwise.} \end{cases}$$

# Visual Document Matching - Backbone

Model Backbone



- AlexNet
- ResNet
- VGG
- EfficientNet
- MobileNetV3
- Vision Transformer (ViT)

## Modelos Avaliados:

- LLaVA 3.2 Vision
- InternVL 2.5
- Qwen2.5-VL
- GPT-4o (2024-11-20)
- GPT-4o-mini (2024-07-18)

## Avaliação:

- Zero-shot (sem fine-tuning)
- Pontuação de similaridade 0–100
- 5 níveis de categorização

# Benchmarking com LLMs

Reference Image

2020 MARLBORO BAR PROGRAM  
CONTRACT TOP SHEET

GMM/SSM: Emily Durick  
MARKET: Durles  
VENUE NAME: Blackberry's

VENUE ID  
DAL-0121-03

Please check the appropriate box that will identify the type of club and the appropriate contract executed by club owner/manager:  
 EVENT  VISIBILITY  MUSIC  RNB

Please check the appropriate box regarding venue admission policy:  
 AQ18-P  AQ21-P  PAO-P  
 AD18-V  AQ21-V

GMM/SSM SIGNATURE:        DATE: 2-3-00

SELL-IN APPROVAL: X DATE:             

Image to Compare

2020 MARLBORO BAR PROGRAM  
CONTRACT TOP SHEET

GMM/SSM: Lindy Jenkins  
MARKET: Charlotte  
VENUE NAME: Coch's Sports Bar & Grill

VENUE ID  
LHA-0080-01

Please check the appropriate box that will identify the type of club and the appropriate contract executed by club owner/manager:  
 EVENT  VISIBILITY  MUSIC  RNB

Please check the appropriate box regarding venue admission policy:  
 AQ18-P  AQ21-P  PAO-P  
 AD18-V  AQ21-V

GMM/SSM SIGNATURE: Lindy Jenkins DATE: 2-20-00

SELL-IN APPROVAL:        DATE: 3-13-2001

Similarity Score: 98  
Category: Nearly Identical

Reference Image

2020 MARLBORO BAR PROGRAM  
CONTRACT TOP SHEET

GMM/SSM: Emily Durick  
MARKET: Durles  
VENUE NAME: Blackberry's

VENUE ID  
DAL-0121-03

Please check the appropriate box that will identify the type of club and the appropriate contract executed by club owner/manager:  
 EVENT  VISIBILITY  MUSIC  RNB

Please check the appropriate box regarding venue admission policy:  
 AQ18-P  AQ21-P  PAO-P  
 AD18-V  AQ21-V

GMM/SSM SIGNATURE:        DATE: 2-3-00

SELL-IN APPROVAL: X DATE:             

Image to Compare

Action T:N Request

Reference Action Network, 1875 Eye Street, N.W., Washington, D.C. 20006 800-424-9870

March 9, 1981  
W-P-Please give O.K. to Kelly office  
C-Show to A.D.S.  
TO: TAN Corporate Coordinator  
Mr. Charles T. McCarty Mr. Shepard P. Pollio  
Mr. V.B. Dey, Jr. Mr. Edward A. Morrison, Jr.  
Mr. Curtis H. Judge Mr. Manuel Leitao, Jr.

FROM: Jack Kelly                
RE: Maine Legislation - L.D. 195, L.D. 599

BACKGROUND  
Two bills currently remain under consideration in the State of Maine: L.D. 399 and L.D. 599. L.D. 246 was withdrawn by its sponsor, Senator Barbara Gill, on March 5, 1981. A third piece of legislation may be added in the next month.

On February 17, 1981 a public hearing was held by the Joint Health & Institutional Services Committee on the remaining bills. After a lengthy debate and action by the committee on March 5, 1981 in a work session, the committee voted 4 to 3 against L.D. 399. However, a minority report written by the committee chairman, Senator Barbara Gill, will be presented to the legislature.

You have previously approved action by TAN members in the state to impact on this legislation. This second request is designed to seek an additional action by our efforts to defeat the pending legislation in the State of Maine.

ACTION REQUESTED  
At this time, we request that the State Director be given permission to contact your organization, TAN members in the state and request that they take the following action at their own time:

1. Personally solicit a minimum of 60 signatures on a public smoking petition (See Attachment A) and return petition to the State Director by March 31, 1981.

Similarity Score: 15  
Category: Completely Different

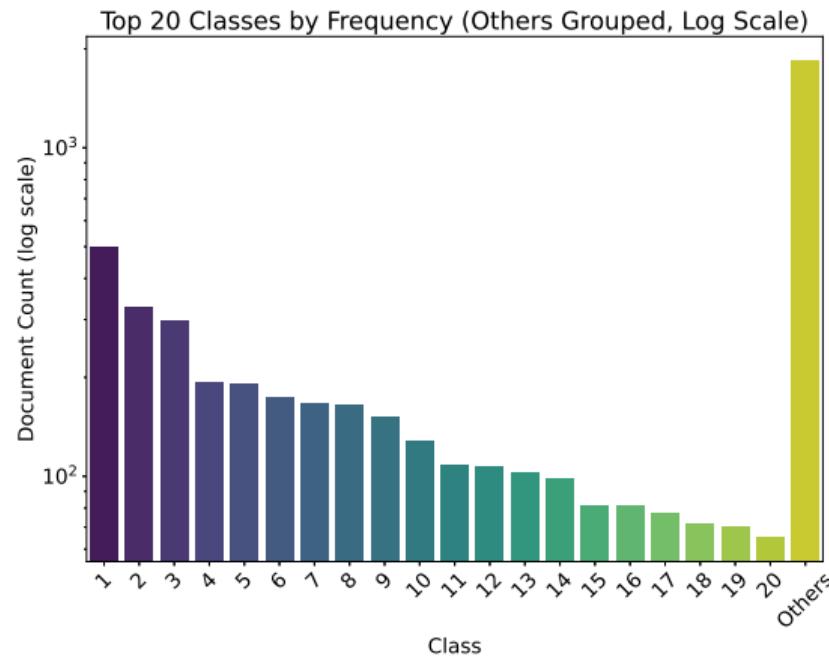
# Resultados

## Composição:

- 4.993 documentos
- 144 classes diferentes
- Min: 2 documentos/classe
- Max: 497 documentos/classe
- Mediana: 13 documentos/classe

## Splits:

- ZSL: separação completa treino/teste
- GZSL: 50% overlap de classes
- 5-fold cross-validation



## Equal Error Rate (EER)

- Ponto onde FAR = FRR
- FAR: False Acceptance Rate
- FRR: False Rejection Rate

$$FAR(\tau) = \frac{\text{False Acceptances}}{\text{Total Negatives}}$$

$$FRR(\tau) = \frac{\text{False Rejections}}{\text{Total Positives}}$$

### Protocolo de Teste:

Para cada documento: 1 par similar + 1 par dissimilar

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	<u>2.10</u>	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	<u>5.73</u>	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	<u>2.75</u>	1.33

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	<u>2.10</u>	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	<u>5.73</u>	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	2.75	1.33

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	2.10	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	2.75	1.33

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	2.10	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	2.75	1.33

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	1.54	4.98	1.51
	34	21M	4.32	2.10	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	0.95
	1	6M	3.93	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	2.75	1.33

# Resultados - Visão Geral

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
AlexNet		57M	8.92	5.45	17.33	6.31
VGG	11	129M	7.47	5.01	14.24	3.95
	13	129M	7.03	4.79	9.30	3.95
	16	134M	8.29	5.23	14.74	4.82
	19	139M	7.30	4.57	17.08	3.90
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	<u>2.10</u>	4.13	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
MobileNetV3	Small	1M	7.98	5.06	12.74	5.26
	Large	4M	8.16	4.27	8.45	4.43
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26
Llama	3.2	11B	—	—	13.95	21.90
InternVL	2.5	8B	—	—	8.58	10.40
Qwen-VL	2.5	7B	—	—	6.61	4.20
GPT-4o mini	2024-07-18	*	—	—	4.70	4.07
GPT-4o	2024-11-20	*	—	—	2.75	1.33

# Resultados - ResNet e EfficientNet

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
ResNet	18	11M	5.03	<b><u>1.54</u></b>	4.98	1.51
	34	21M	4.32	2.10	<b><u>4.13</u></b>	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
EfficientNet	0	4M	4.41	2.27	6.02	<b><u>0.95</u></b>
	1	6M	<b><u>3.93</u></b>	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34

# Resultados - ResNet, EfficientNet e ViT

Architecture	Edition	Params	ZSL	GZSL	Test ZSL	Test GZSL
ResNet	18	11M	5.03	<u>1.54</u>	4.98	1.51
	34	21M	4.32	2.10	<u>4.13</u>	1.53
	50	23M	6.90	3.39	10.34	2.21
	101	42M	8.20	2.72	11.31	1.98
	152	58M	9.44	3.38	12.70	2.39
EfficientNet	0	4M	4.41	2.27	6.02	<u>0.95</u>
	1	6M	<u>3.93</u>	3.54	8.88	2.70
	2	7M	5.73	2.61	7.29	2.14
	3	10M	5.65	3.64	7.37	2.34
ViT	Base	87M	12.43	7.97	19.72	5.19
	Large	305M	13.16	7.57	19.88	5.26

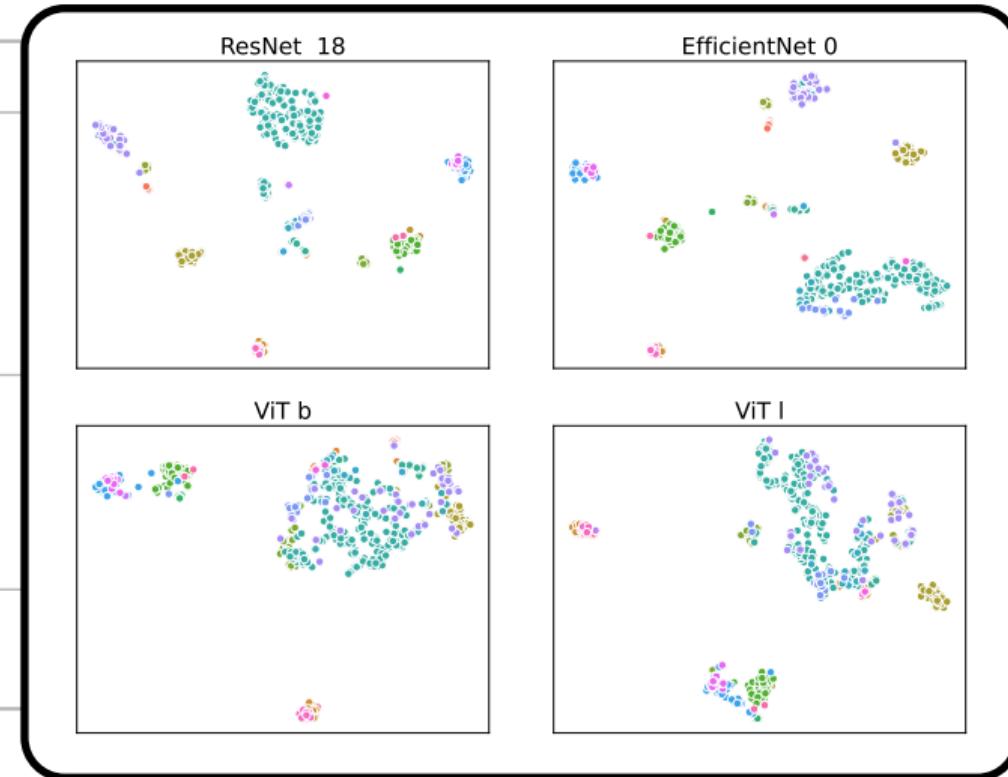
# Resultados - Visualização TSNE

## Architecture

ResNet

EfficientNet

ViT



## Test GZSL

1.51  
1.53  
2.21  
1.98  
2.39

**0.95**  
2.70  
2.14  
2.34

5.19  
5.26

# Resultados - Large Language Models

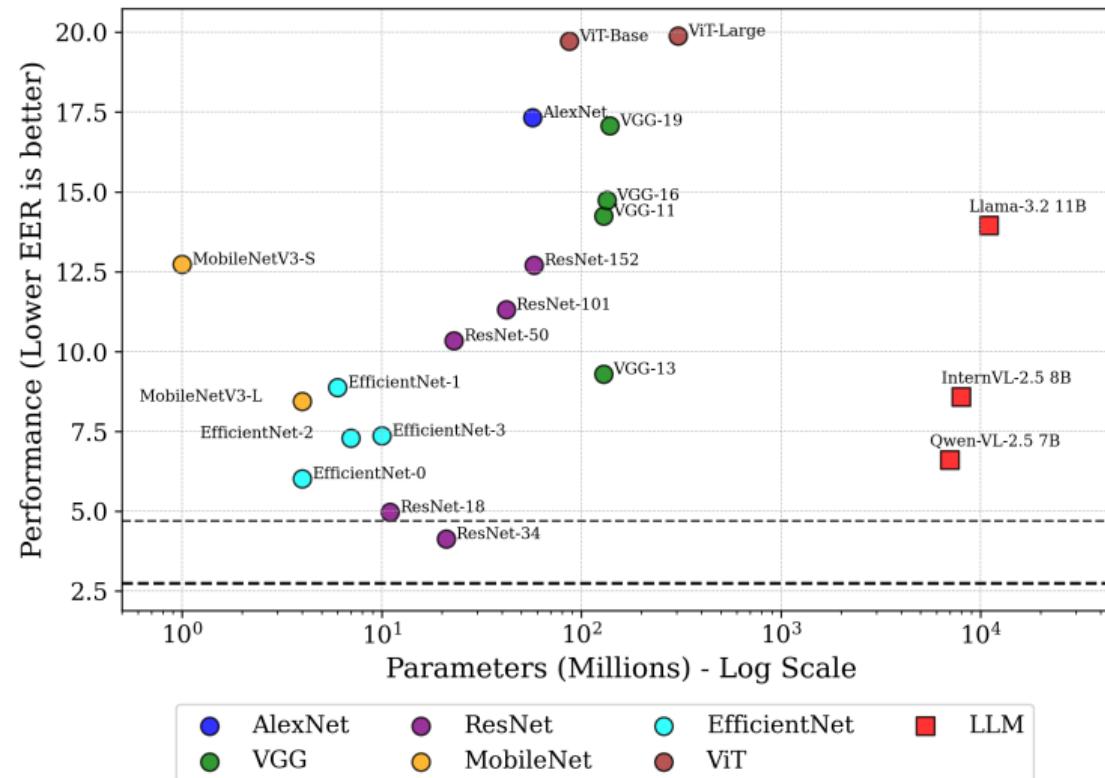
Model	Version	Params	Test ZSL	Test GZSL
Llama	3.2	11B	13.95	21.90
InternVL	2.5	8B	8.58	10.40
Qwen-VL	2.5	7B	6.61	4.20
GPT-4o mini	2024-07-18	*	4.70	4.07
GPT-4o	2024-11-20	*	<u>2.75</u>	<u>1.33</u>

\* Parameter count not publicly disclosed

# Resultados - Melhores Modelos Visuais vs LLMs

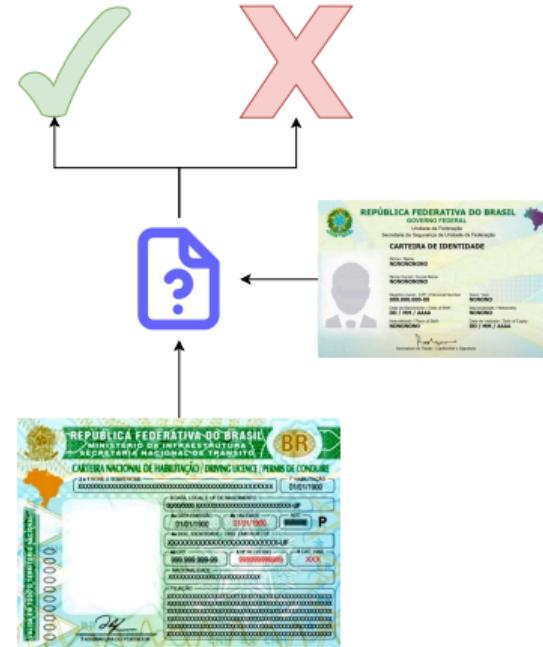
Architecture	Edition	Params	Test ZSL	Test GZSL
ResNet	18	11M	4.98	1.51
	34	21M	4.13	1.53
	50	23M	10.34	2.21
	101	42M	11.31	1.98
	152	58M	12.70	2.39
EfficientNet	0	4M	6.02	<b><u>0.95</u></b>
	1	6M	8.88	2.70
	2	7M	7.29	2.14
	3	10M	7.37	2.34
Llama	3.2	11B	13.95	21.90
InternVL	2.5	8B	8.58	10.40
Qwen-VL	2.5	7B	6.61	4.20
GPT-4o mini	2024-07-18	*	4.70	4.07
GPT-4o	2024-11-20	*	<b><u>2.75</u></b>	1.33

# Comparação: Visual Models vs LLMs



----- GPT-4o Mini 2024-07-18 ----- GPT-4o 2024-11-20

# Comparação: Visual Models vs LLMs



Verificação - Classificação Binária

# Comparação: Visual Models vs LLMs



## Identificação - Classificação Multiclasses

## Conclusão

## ① LA-CDIP Dataset

Categorizado por layout visual para Zero-Shot Learning

## ② Metodologia com VDM

Enforçando comparação de documentos como forma de classificação

## ③ Benchmarking Sistemático

Comparação entre modelos visuais e LLMs

## Trabalhos Futuros:

- Aumentar número de amostras
- Aumentar número de classes
- Incluir fontes adicionais de documentos
- PIBIC

# Publicações

**Visual Document Matching for Zero-Shot Document Classification**

Lucas De Almeida Bandeira Macedo<sup>1,2</sup> (✉) · Joao Pedro Felix De Almeida<sup>1</sup> · Pedro Garcia Furtado<sup>1</sup> · Pedro Freitas<sup>1</sup> · Li Weipeng<sup>2</sup>

University of Brasília, UnB · Brasília, Brazil · Federal District, 70910-900

**Abstract.** Accurate document identification is crucial for ensuring compliance and maintaining consistency across various applications. Consequently, document classification has been extensively studied. However, the dynamic nature of documents often requires single classification instead of multi-class classification. Layout-aware document classification is of paramount importance for various companies. However, documents often change in terms of format and their visual representation, making it challenging to maintain a classifier. Moreover, model continuance and retraining often demands important resources, such as time and computational power, for handling new data. Therefore, techniques capable of classifying documents by simply sharing some data, without necessarily requiring labeled data, are highly desired. In this work, we propose a zero-shot learning framework for visual document matching, addressing document invariance, ensuring incoherence detection, and using losses that encourage a more suitable approach. To address the zero-shot classification problem, we build a Visual Document Matching (VDM) dataset focused on verifying whether two documents share an identical visual layout structure, a particularly effective approach in scenarios where training classes do not align with those encountered during inference. Despite its simplicity, our VDM approach achieves state-of-the-art performance in document layout understanding. To support our study, we introduce Layout-Aware Complex Document Information Processing (LA-CIDP), a dataset containing 100k images from 1000 unique classes. We reorganized the dataset from the RVL-CIDP database to emphasize visual structure over semantic information. Our approach is based on a multi-task learning framework, a common machine learning framework across multiple backbone architectures, including ResNet, EfficientNet, and Vision Transformer (ViT). In another scenario, our method is compared with state-of-the-art document matching methods via verification with cross-validation. Furthermore, our VDM approach outperforms lighter Large Language Models (LLMs) and results GPT-in. These findings demonstrate that our proposed framework can be used as a general-purpose multi-modal module, demonstrating high accuracy with significantly fewer parameters, making our approach more practical for real-world applications.

**Keywords:** Zero-shot learning · Metric Learning · Document Understanding · Document Image Classification · Visual Document Matching

**Towards Zero-Shot Document Image Classification**

Lucas Macedo, João Pedro Costa, João Pedro Félix de Almeida, Pedro Freitas, and Li Weipeng  
Departments of Computer Science, University of Brasília, Brasília, Brazil

**Abstract.** Classification is a fundamental tool to automate the process of categorizing documents in many real-world applications, such as document processing, medical records management, news categorization, fraud detection, regulatory compliance, and many others. Recently, document classification has been extensively studied. However, document classification is of paramount importance for various companies. However, documents often change in terms of format and their visual representation, making it challenging to maintain a classifier. Moreover, model continuance and retraining often demands important resources, such as time and computational power, for handling new data. Therefore, techniques capable of classifying documents by simply sharing some data, without necessarily requiring labeled data, are highly desired. In this work, we propose a zero-shot learning framework for visual document matching, addressing document invariance, ensuring incoherence detection, and using losses that encourage a more suitable approach. To address the zero-shot classification problem, we build a Visual Document Matching (VDM) dataset focused on verifying whether two documents share an identical visual layout structure, a particularly effective approach in scenarios where training classes do not align with those encountered during inference. Despite its simplicity, our VDM approach achieves state-of-the-art performance in document layout understanding. To support our study, we introduce Layout-Aware Complex Document Information Processing (LA-CIDP), a dataset containing 100k images from 1000 unique classes. We reorganized the dataset from the RVL-CIDP database to emphasize visual structure over semantic information. Our approach is based on a multi-task learning framework, a common machine learning framework across multiple backbone architectures, including ResNet, EfficientNet, and Vision Transformer (ViT). In another scenario, our method is compared with state-of-the-art document matching methods via verification with cross-validation. Furthermore, our VDM approach outperforms lighter Large Language Models (LLMs) and results GPT-in. These findings demonstrate that our proposed framework can be used as a general-purpose multi-modal module, demonstrating high accuracy with significantly fewer parameters, making our approach more practical for real-world applications.

**Keywords:** Zero-shot learning · Metric Learning · Document Understanding · Document Image Classification · Visual Document Matching

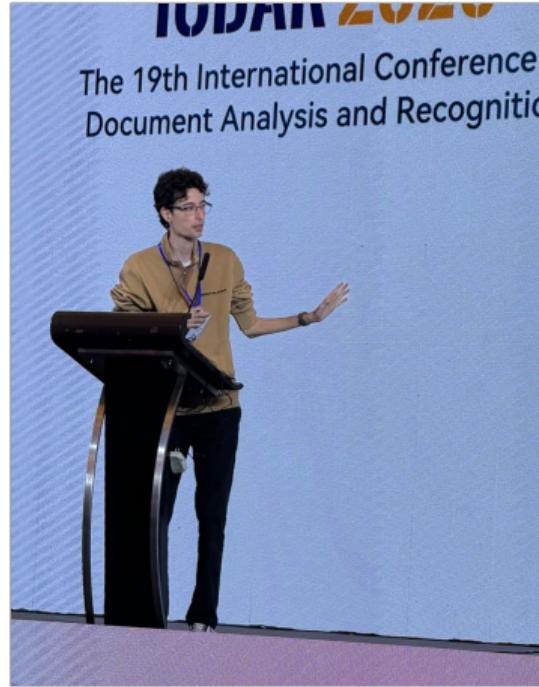
ICDAR

SIBGRAPI

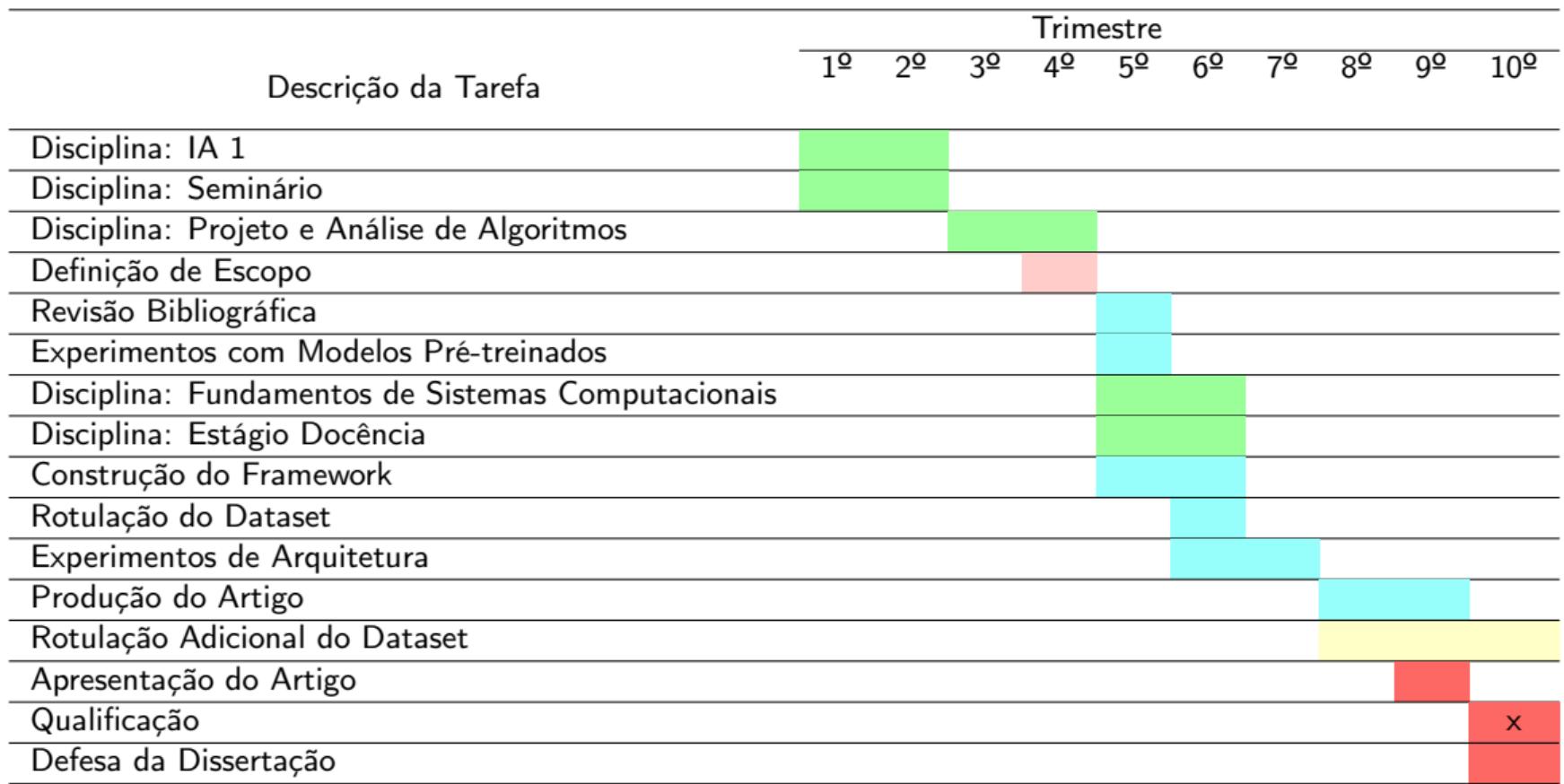
# Publicações



## ICDAR



## Apresentação



# Obrigado!

Lucas de Almeida Bandeira Macedo

[lucasabmacedo@hotmail.com](mailto:lucasabmacedo@hotmail.com)

Orientador: Prof. Dr. Pedro Garcia Freitas

Coorientador: Prof. Dr. Bruno Luiggi Macchiavello Espinoza