



Rachit Jain

Chloe Wu

Candidates of Master of Business Analytics, MIT

## CAPSTONE PROJECT

### Document Classification Capability

Categorization of scanned documents

MIT x Wolters Kluwer

Faculty Advisor: Dr. Ilya Jackson

18<sup>th</sup> August 2023

# Agenda



## Introduction

Overview  
Challenges  
Motivation



## The Process

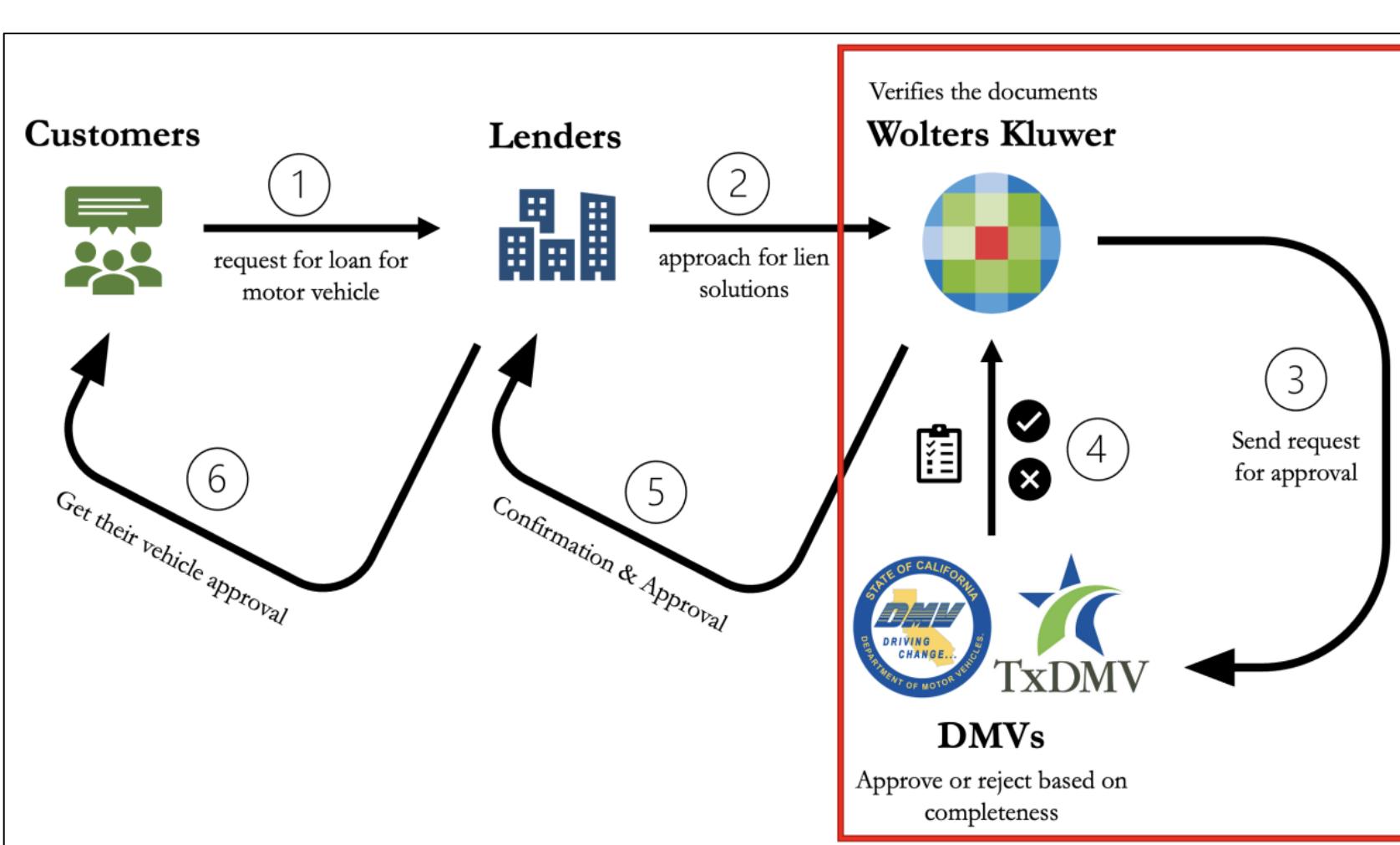
Solution  
Methodology  
(Re-labelling, Modelling)



## Results

Demo  
Deliverables  
Business Value & Impact

# Overview | Motor Vehicle registration process is error-prone



Challenges At Scale

Huge Volume

100k+ pages\*  
come in per day

Multiple rejections

10%  
rejections

High processing time

10 mins per  
request (~20 pages)

# Challenges | Manual processing is a bottleneck

## Challenges At Scale



**100k+ pages\***  
come in per day

Huge Volume

**10%**  
rejections

Multiple rejections

**10 mins** per  
request (~20 pages)

High processing time



## Final Goal

Build an automated, generalized document classification capability to make historically manual logistics paperwork easier to execute and more accurate

# Motivation | Need for more than rule-based systems

Similar formats, similar text, but different titles

(1) Return Unprocessed

CT Lien Solutions  
a Wolters Kluwer Business

MV-520716-1  
ASHLEY HOGGAMP  
Albany Team 3  
187 Wolf Road,  
Suite 101  
Albany, NY 12205  
8003423676  
liensolutions.dmvteam@wolterskluwe  
r.com

TASHA MONROE  
Autopay Direct, Inc  
8055 East Tufts  
Suite 1100  
Denver, CO

Order#:  
Customer:  
Date:  
AP Acct#: AP22063641962  
AP Client:

VIN#: 1FMSK7D8XHGD53336  
BorrowerName: COURTNEY NALLY  
Jurisdiction: Texas-Travis  
Transaction Type: Lien Add - Add/Remove Spouse

Numbers of days on-hold: 57  
Comments:  
Tracking 1ZX17318NT76100506

Actual Fees: WKLS Fee: \$59.74  
Total: \$59.74

This report contains information compiled from sources which Lien Solutions considers reliable, but does not control. Information provided is non-certified unless otherwise indicated. Lien Solutions in no way undertakes or assumes any part of the customer's business, legal or similar risks, and does not guarantee the accuracy, completion or timeliness of the information provided and shall not be liable for any losses or injuries whatever resulting from any contingency beyond its control, or from negligence regardless of the cause. The categorization of filings is provided for the convenience of the customer and is not to be construed as a legal opinion concerning the status of the filings.

(2) Search Request Form

CT Lien Solutions  
a Wolters Kluwer Business

MV-509368-1-440153-SRF  
Wolters Kluwer's Lien Solutions  
187 Wolf Rd STE 101  
Albany NY 12205  
LienSolutions.DMVTeam@wolterskluwer.com  
800-833-5778 option 2 then option 3

To: TRAVIS COUNTY TAX OFFICE  
MOTOR VEHICLE DEPT.  
2433 Ridgepoint Dr.

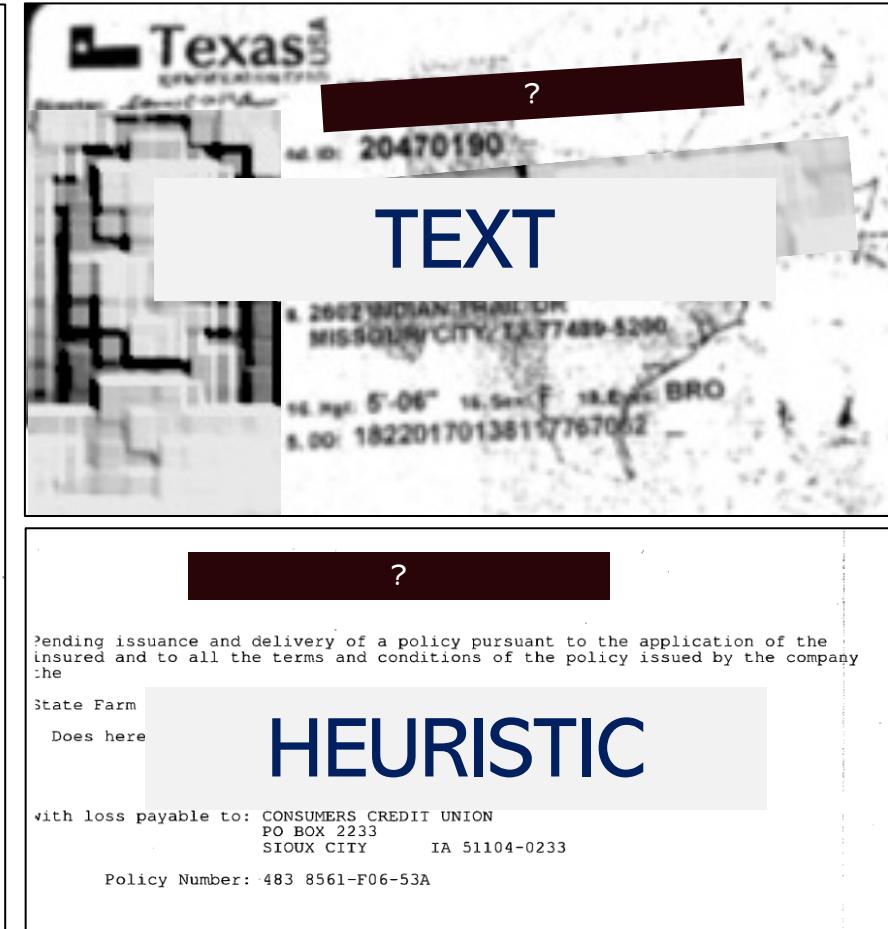
Order#: 88547166  
Date: 2/2/2023

BorrowerName: CANDICE RENEE CALVERT  
Jurisdiction: Texas-Travis  
Transaction Type: Lien Add

Comments:  
Please process Lien Add.  
For any questions, please contact at 800-833-5778.  
Please include titles/receipts in enclosed UPS envelope.  
Thank you.

Disclaimer:  
The following obligation applies only to non-governmental agencies or entities:  
By performing the services requested hereunder, you hereby agree to be subject to, and to comply with, CT's Correspondent Terms and Conditions located at <https://www.wolterskluwer.com/en/solutions/ct-corporation/resources/terms-and-conditions-correspondent>

(3) Identification Card



(4) Binder of Insurance

# Our Solution

CV

+

NLP

+

HEURISTIC

Computer Vision

Natural Language Processing

Title Extraction

==

**Document Classification Capability**

==

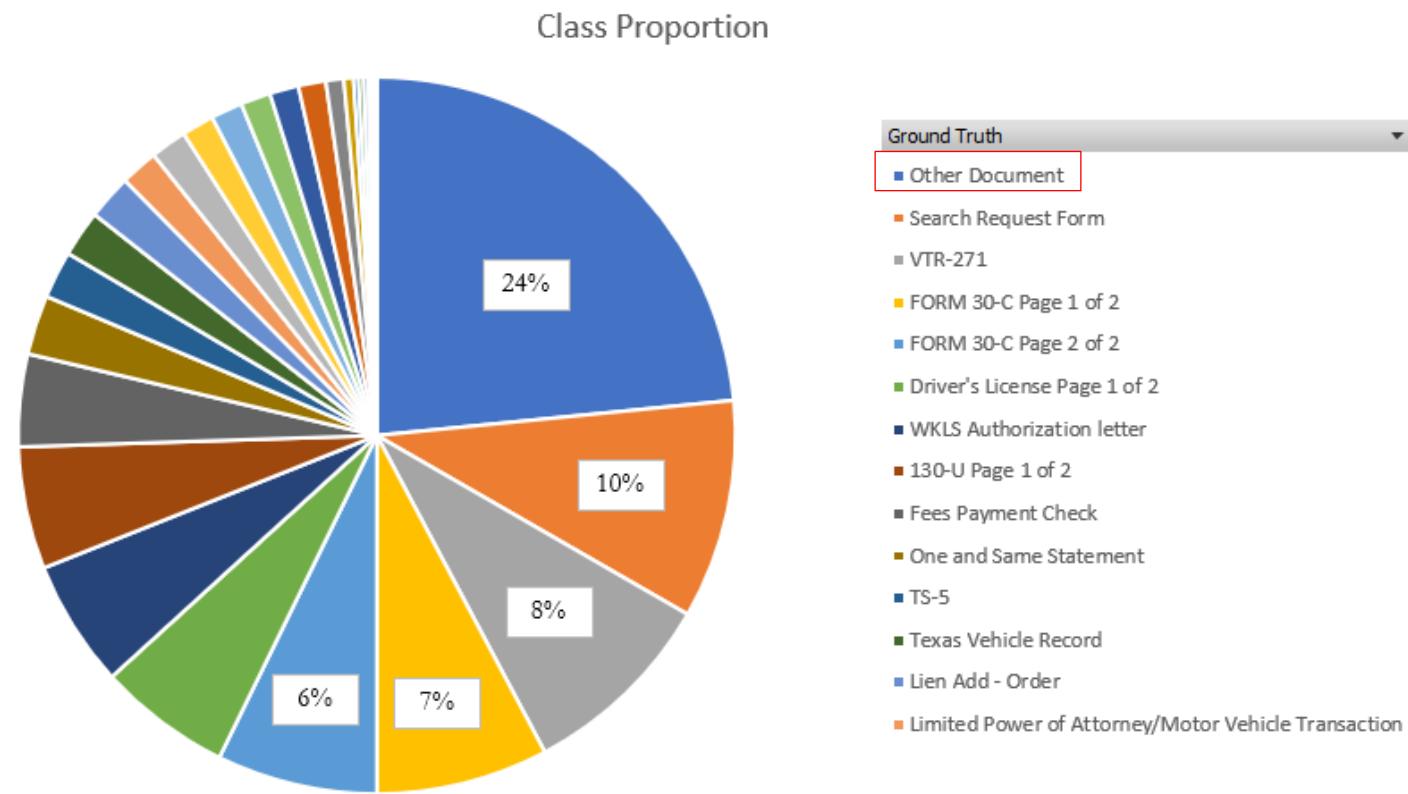
**Capstone of the Year?**

# Imbalanced dataset across 120+ categories; 12k scanned pages

Training Data

**12k+**  
Data Points

**10k+**  
from top 31  
classes



## Represented Counties:

Bexar, Brazoria, Dallas, El Paso, Fort Bend, Harris, Hidalgo, Lubbock, Travis, Van Zandt

# Solution | End-to-End Workflow takes PDF input and gives multiple labels for each page, along with confidence scores

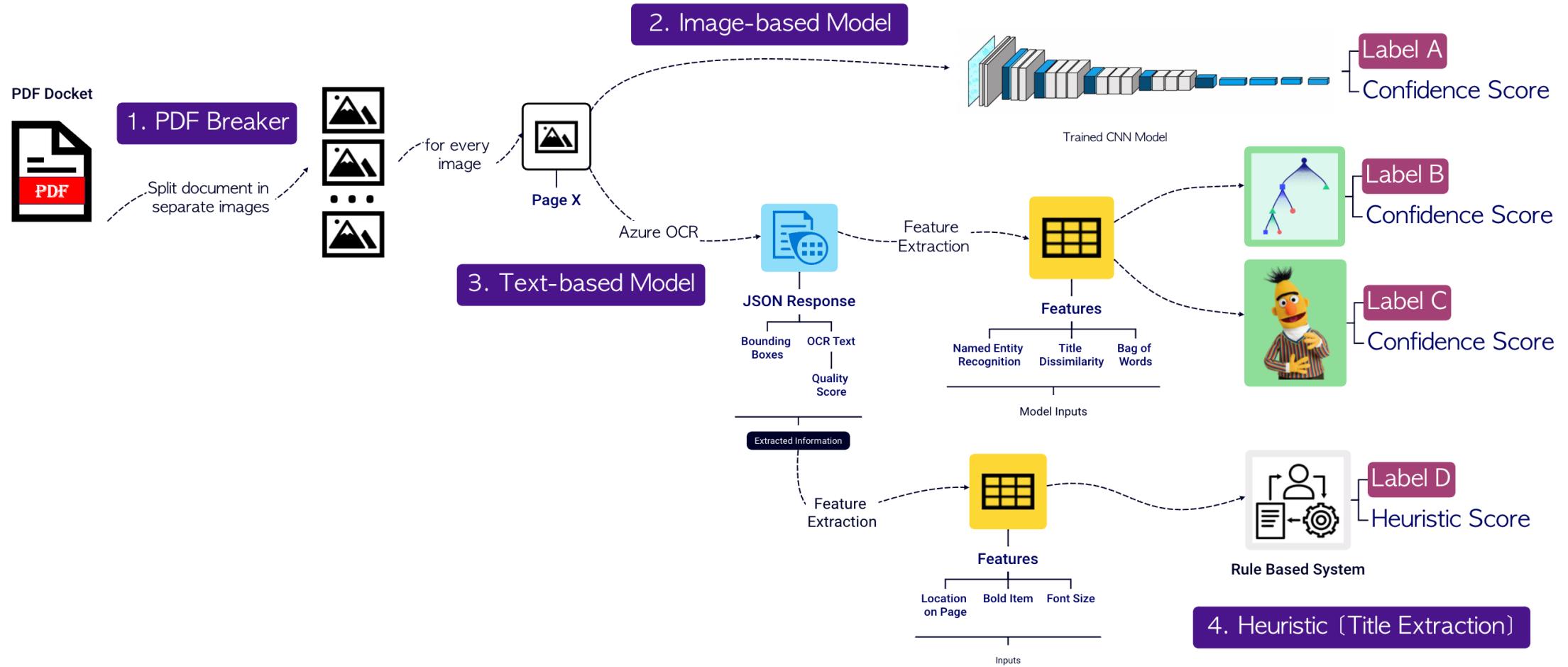


Fig: End-to-End Workflow

# Step 0: Need for relabeling - Mystery behind ground truth labels

Status Quo  
Current 'Ground Truth' Label =  
Output of Champion model  
Wrong Labels → Poor Models

Solution?  
'Smarter' Manual Labelling  
[Unsupervised Clustering on Deep Embeddings]

- Cluster similar image embeddings from trained vision model
- Merge clusters on common categories & create sub-clusters
- Smarter Manual Label

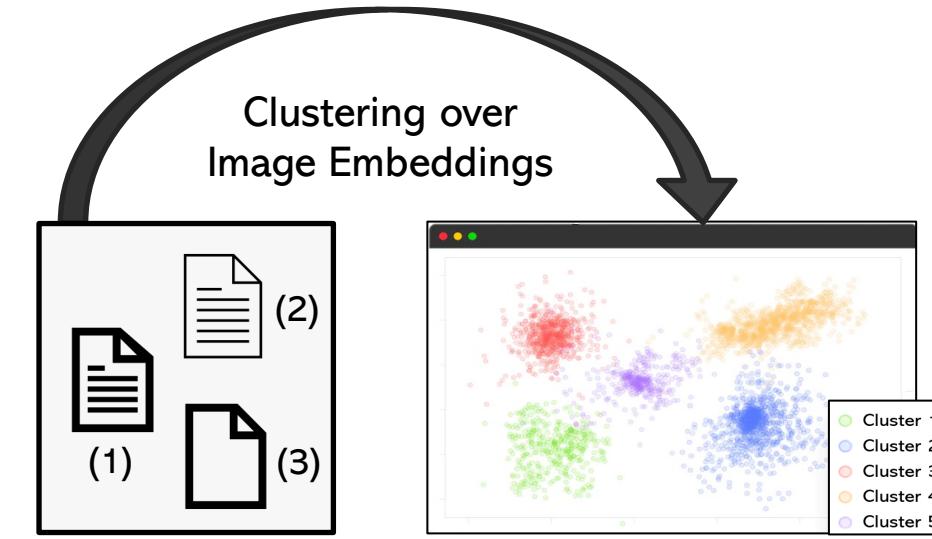


Fig: Embeddings for each image clustered based on similarity

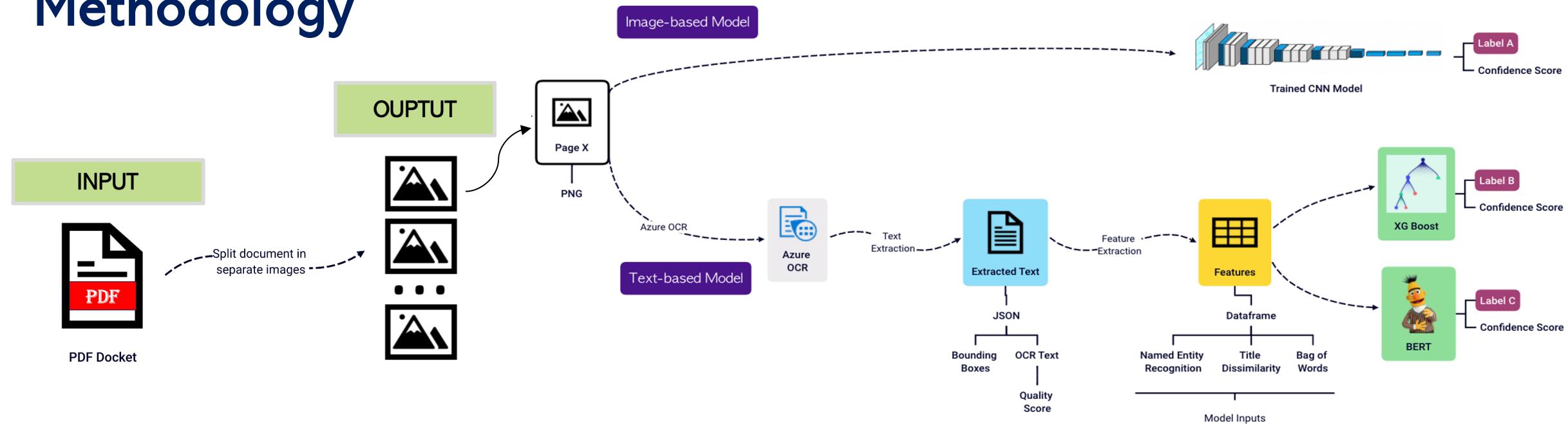
Results

**88%**  
saving in time for manual labelling

**10.6%**  
Mis-classified labels identified

**8%**  
F1 Score jump!

# Methodology



## ① PDF Breaker API

PDF Breaker on a Streamlit dashboard

Each page is extracted as an image and saved in folders

## ② Vision-based Model

CNN model fine-tuned on 12k scanned documents

Highlights: Image padding, Image processing, Hyper-parameter tuning, Data Augmentation, Feature Engineering...

## ③ Text-based Model

BERT + XGBoost running on text extracted from AzureOCR

# Methodology

## 1 PDF Breaker API

PDF Breaker on a Streamlit dashboard

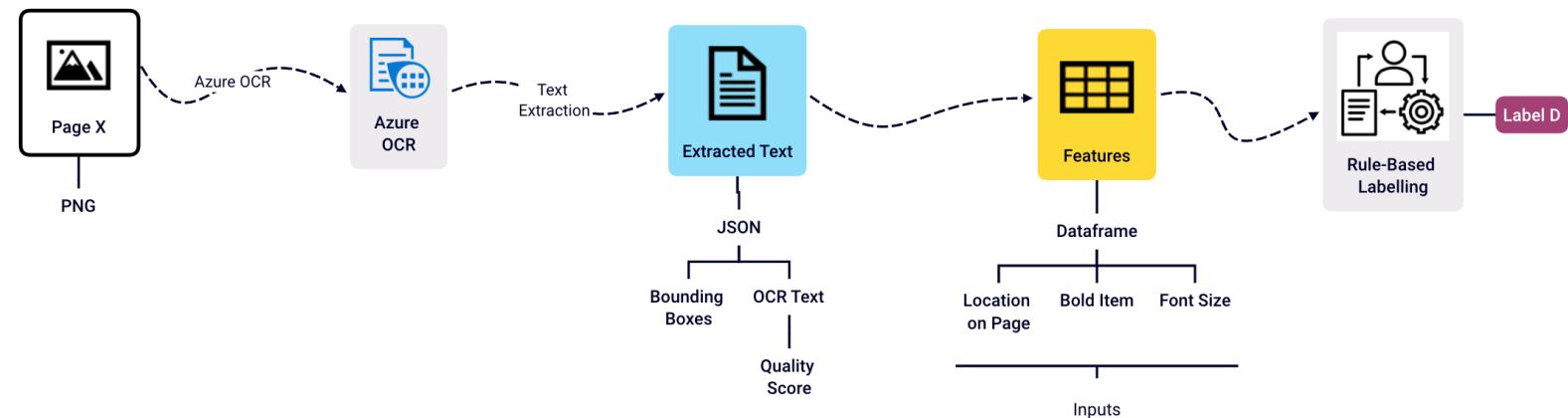
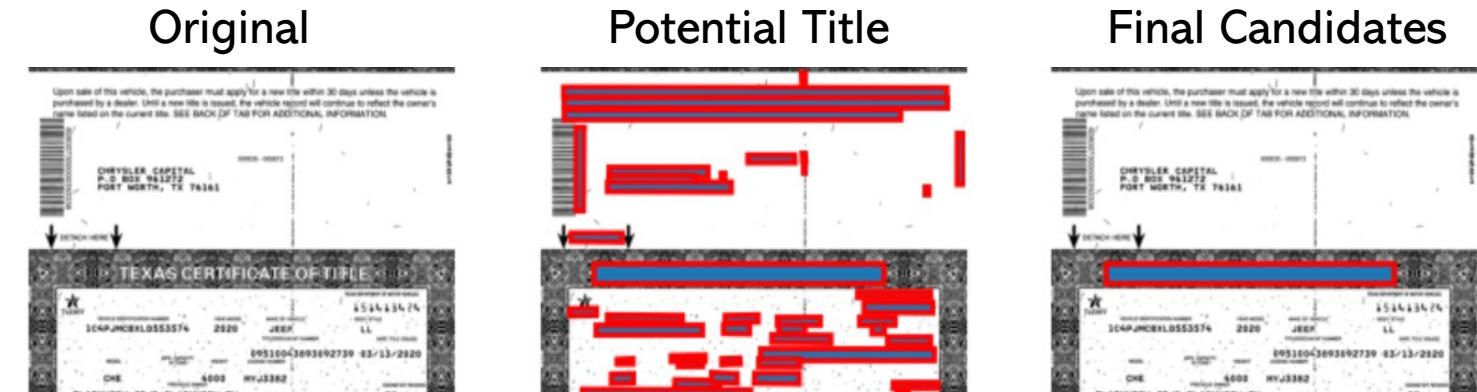
Each page is extracted as an image and saved in folders

## 2 Vision-based Model

CNN model **fine-tuned** on 12k scanned documents

## 3 Text-based Model

BERT + XGBoost running on text extracted from AzureOCR



## 4 Heuristic-based Model

Potential title extraction from any document type

## 5 Ensembling

Selecting '**supreme**' model;  
if clash, choose heuristic

# The labels need to be ensembled into one single prediction

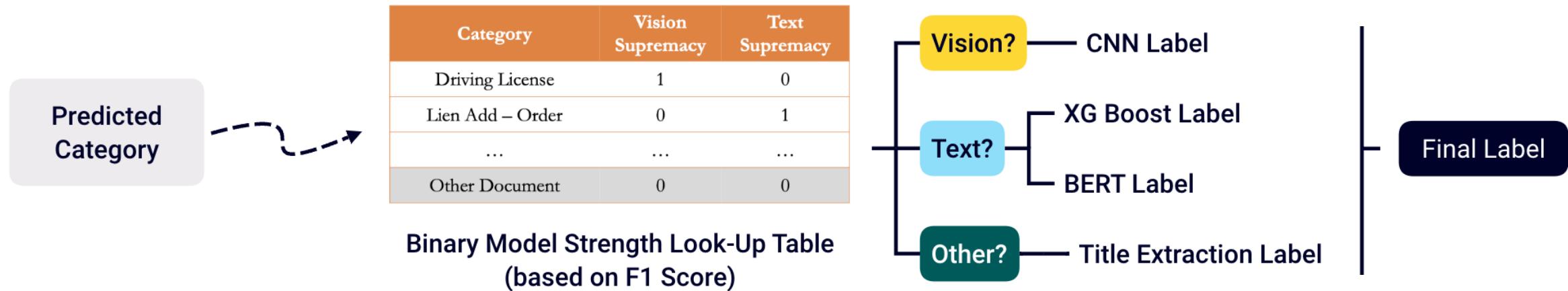


Fig: Logic workflow to combine the labels

## Our Implementation

Model Supremacy from Vision and Text with Confidence Score

The final model is chosen based on the F1 score & rules based on heuristic title

## Ensemble Technique

Case	Same Label	Vision Supremacy?	Text Supremacy?	Final Model
Case 1	1	1	0	Same
Case 2	0	1	1	Manual Review
Case 3	0	1	0	Vision
Case 4	0	0	1	Text
Case 5	1	0	0	Heuristic

# Result | End-to-End multi-modal architecture deployed

**Classification API**

This API gets the classification results for each page, given a PDF docket & the Azure OCR response.

**Classification API** API for automated document classification using unified multi-modal capability

**POST** /overall/document Document Classification Capability

Upload PDF docket (optional) Azure OCR Response docket & get document classification for every chosen page

**Parameters**

Name Description

page\_numbers array[integer] (query) Page numbers of the document to classify

all\_pages boolean (query) Do you want to run it for all pages?

**Request body**

document \* string(\$binary) Choose PDF docket to be uploaded as input

azure\_ocr\_response string(\$binary) Choose Azure OCR response JSON file

response\_format \* string Which format for the output do you need?

**Responses**

**Curl**

```
curl -X 'POST' \
  http://40.124.27.44:8000/overall/document?all_pages=true \
  -H 'Content-Type: multipart/form-data' \
  -F 'document=@022052_1.pdf;type=application/pdf' \
  -F 'azure_ocr_response=@022052_1.json;type=application/json' \
  -F 'response_format=json'
```

**Request URL**

```
http://40.124.27.44:8000/overall/document?all_pages=true
```

**Server response**

File Edit Selection View Go Run ... < > src [SSH: 40.124.27.44] X

EXPLORER SRC [SSH: 40.124.27.44]

- \_pycache\_
- Azure\_OCR\_API
- BERT\_model\_API
- Combine\_output\_API
- Dataset
- heuristic\_api
- Overall\_API
  - Dataset
  - Models
  - Notebooks
  - Pipeline
- Result
  - README.md
  - requirements.txt
  - Pipeline
  - vision\_model\_api
  - XGBoost\_model\_API

Preview README.md X

## Classification API

This API takes the input of a PDF docket and a json object for the PDF docket and outputs a downloadable csv file enclosing the predicted results using vision, text, and heuristic approaches. It also has the option to run on selected pages. It includes the steps of dataset preparation, model training, and prediction. Follow the steps below to run the pipeline:

### Step 1: Load Libraries

Install packages and dependencies specified in requirement.txt file before running the code.

NOTE: Some packages is not supported for Windows, such as tensorflow-text==2.13.0 Please check package version and its dependencies carefully.

### Step 2: PDF Breaker

This step breaks the PDF into a folder containing all the pages as images. It consists of two functions:

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS 1

root@wkmmitcapstone:/mnt/src/Overall\_API#

SSH: 40.124.27.44 X 0 △ 0 ⌂ 1

Video: Demo implemented over FastAPI and **deployed** over WK's Virtual Machine

# Deliverables | Scalable model pipeline deployed over WK Cloud

## Result Summary



0.86

F1 Score

Over 31 document types  
over 2.1k highly noisy  
test dataset



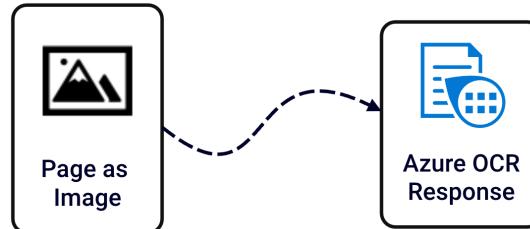
Challenger >

Champion

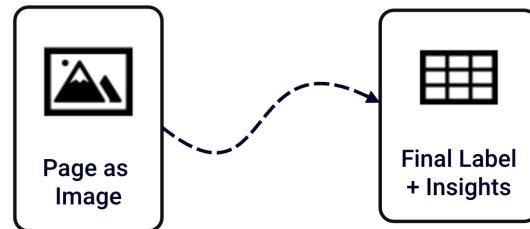
Our implementation beats  
the status quo with 5%  
higher F1 score

## Our Deliverables to Wolters Kluwer

### (1) OCR API



### (2) CAPABILITY API



### (3) Documentation & Knowledge Transfer

for smooth integration with current system

## Learnings

### Industry Practices

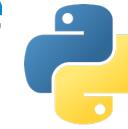


Dashboarding



Capability Design

### Tech Stack



FastAPI

### Business Skills



New Heights  
LET'S SOAR N  
TOGETHER 2022  
New Frontiers

# High learning + Solid deliverables + Strong impact =

More ACCURATE Labelling

**0.86 F1**

High → Better

Flexible PDF BREAKER

**Unrestricted**  
# of PDFs broken

Result INTERPRETABILITY

**User Insights**  
behind predictions

AUTOMATED pipeline

**1 API**

running everything

FEWER rejections



**Challenger >**  
Champion (status-quo)

END-TO-END pipeline

**Streamlined**  
workflow

LOW processing time

**10X saving**

in processing time

PRODUCTIONIZED pipeline

**Deployed**  
over WK's VM

EASY-TO-INTEGRATE capability

**Scalable**  
and reproducible



Rachit Jain

Chloe Wu

Candidates of Master of Business Analytics, MIT

# Thank You!

“ Only those who will risk  
going too far can possibly  
find out how far one can go! ”

~ T.S. Eliot

# Appendix

- (1) [Challenger vs Champion](#)
- (2) [Ground Truth Labelling](#)
- (3) [PDF Breaker Process](#)
- (4) [Vision Based Model Process](#)
- (5) [Text Based Model Process](#)
- (6) [Heuristic Model Process](#)
- (7) [Ensembling Process](#)
- (8) [Results](#)
- (9) [Final Deliverables](#)

# The mystery behind ground truth labels: need for relabeling

## Status Quo



**Output** of Champion  
model =  
**'Ground'** Truth Label

# Challenge

Wrong Labels → Poor Models  
The models learn to predict  
incorrectly

## Realization

# Data-Centric approach ‘Confident Learning’

## Solution?

Manually label all 12k  
training images, but **smartly!!**

Category	%age Mislabelled	Count of Documents
130-U Page 1 of 2	4.90%	572
Check Endorsement	11.68%	137
Driver's License Page 1 of 2	15.90%	547
Fees Payment Check	1.68%	417
FORM 30-C Page 1 of 2	1.15%	786
FORM 30-C Page 2 of 2	2.24%	758
Lien Add - Order	13.85%	231
Limited Power of Attorney/Motor Vehicle Transaction	1.68%	179
MV-50	5.93%	135
Odometer Disclosure Statement	56.92%	130
One and Same Statement	29.41%	374
Other Document	16.97%	2841
Search Request Form	0.19%	1045
Texas Liability Insurance Card	59.57%	235
Texas Vehicle Record	8.71%	241
TS-5	0.00%	226
TS-8	18.93%	169
VTR-271	2.00%	948
WKLS Authorization letter	5.60%	607
<b>Grand Total</b>	<b>10.60%</b>	<b>10578</b>

# The mystery behind ground truth labels: need for relabeling

Intuition



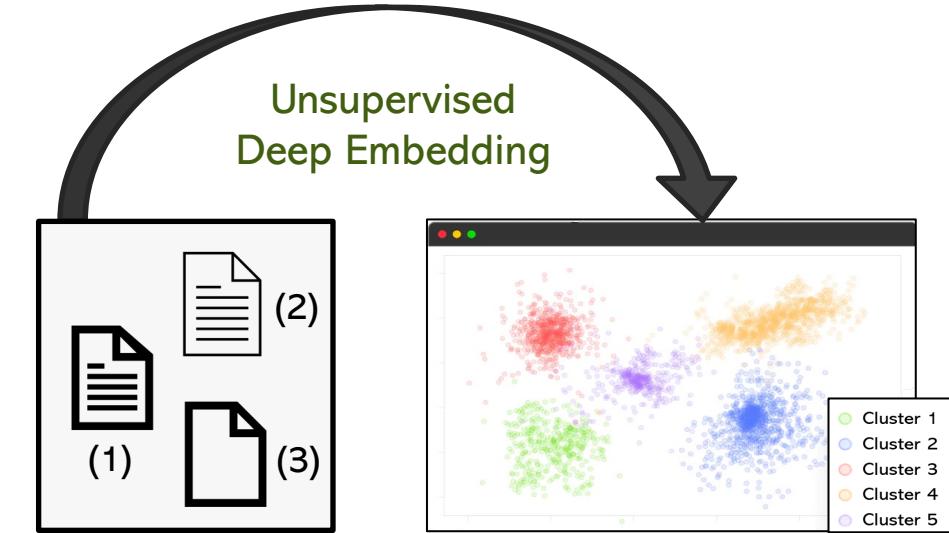
Smarter Labelling [Unsupervised Clustering on Deep Embeddings]

'Can you spot the (major visual) difference?' is easier while scrolling through 100 similarly formatted images!

Solution Process

- Get image embeddings from trained vision model
- Cluster based on embedding similarity
- Assign the most prominent category
- Merge clusters on common categories
- Create sub-clusters within the clusters
- Smarter manual labelling

Our Solution



**84<sub>hrs</sub> → 10<sub>hrs</sub>**

Assuming 30s to classify a page

**10.6%**  
Mis-classified labels corrected

Results

**~8%**  
F1 Score jump!

# PDF Breaker extracts every page from the docket as an image

## Our Implementation

### PDF Breaker on a user-friendly Streamlit dashboard

Input: ZIP folder with any number of PDFs within county folders  
Output: ZIP folder with all images stored in respective folders

#### Process:

- Every docket contains 'n' number of pages
- Each page is extracted as an image and saved in respective folder
- The models use these images as inputs

## PDF to Image Converter

Choose a zip file to extract

Drag and drop file here  
Limit 200MB per file • ZIP

Browse files

Batch 1.zip 128.3MB

X

Processing document 1 of 76 : Batch 1\375013-1

Processing document 2 of 76 : Batch 1\552138-1

Processing document 3 of 76 : Batch 1\561948-1

Processing document 4 of 76 : Batch 1\564536-1

Processing document 5 of 76 : Batch 1\575019-1

Processing document 6 of 76 : Batch 1\576291-1

Processing document 7 of 76 : Batch 1\577635-1

Processing document 8 of 76 : Batch 1\582656-1

Processing document 9 of 76 : Batch 1\586415-1

Processing document 10 of 76 : Batch 1\590719-1

Processing document 11 of 76 : Batch 1\590836-1

Processing document 12 of 76 : Batch 1\594153-1

### Fully functional PDF Breaker

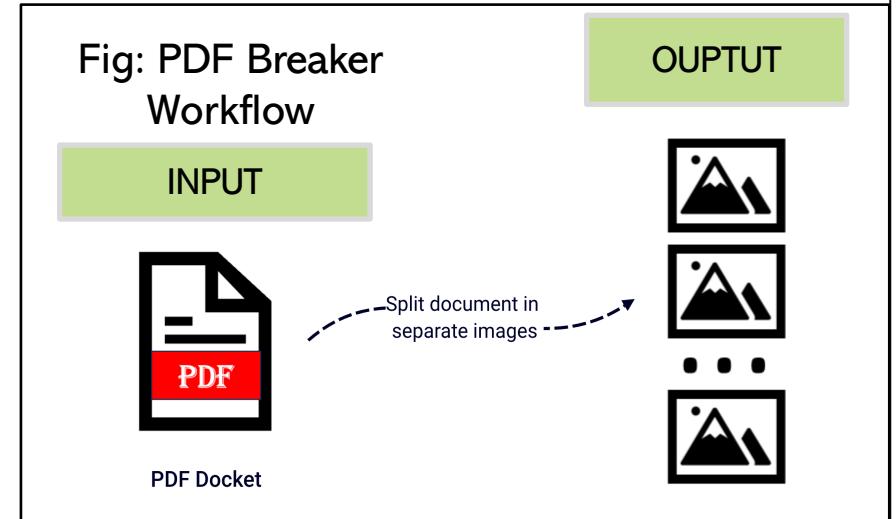
Fig: PDF Breaker Workflow

INPUT

PDF Docket

OUTPUT

Split document in separate images



# Vision-based model harnesses the power of visual formatting

## Our Implementation

 ResNet50 Model **fine-tuned**  
on 12k scanned  
documents

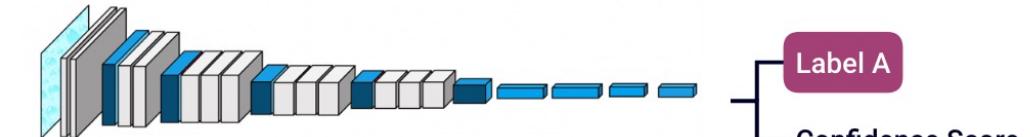
Input: Image (document)  
Output: Category label and  
confidence score

### Highlights:

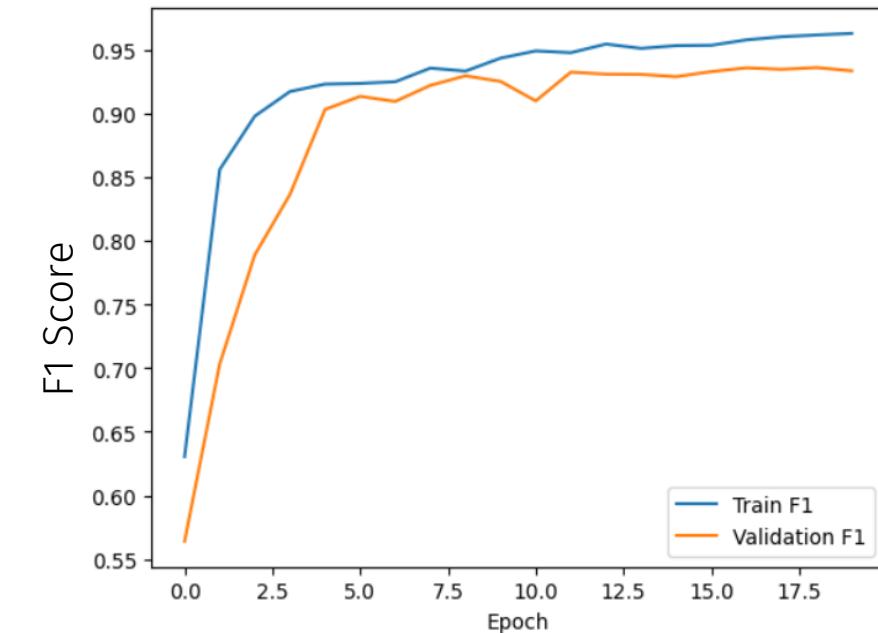
- Image padding while maintaining resolution
- Fine-Tuned instead of simply using pre-trained model
- Image processing (adaptive document image binarization, contrast reduction, noise removal...)
- Hyper-parameter tuning (contributed to ~2%+)
- Data Augmentation for training + validation (further 3%+)
- Pipeline ready!



Fig: Image-based model workflow



Trained CNN Model



# Text-based model harnesses the power of document content

## Our Implementation



BERT + XGBoost

running on text

extracted from AzureOCR

Input: Image (document)

Output: Category label and

confidence score (for both models)

### Highlights for BERT:

- Text Preprocessing to remove noise
- Transfer Learning with model fine-tuning
- Hyper-parameter Tuning

### Highlights for XGBoost:

- NLP Feature Engineering (Bag of Words, Named Entity Recognition, Title Dissimilarity)
- Hyper-parameter Tuning
- Pipeline ready

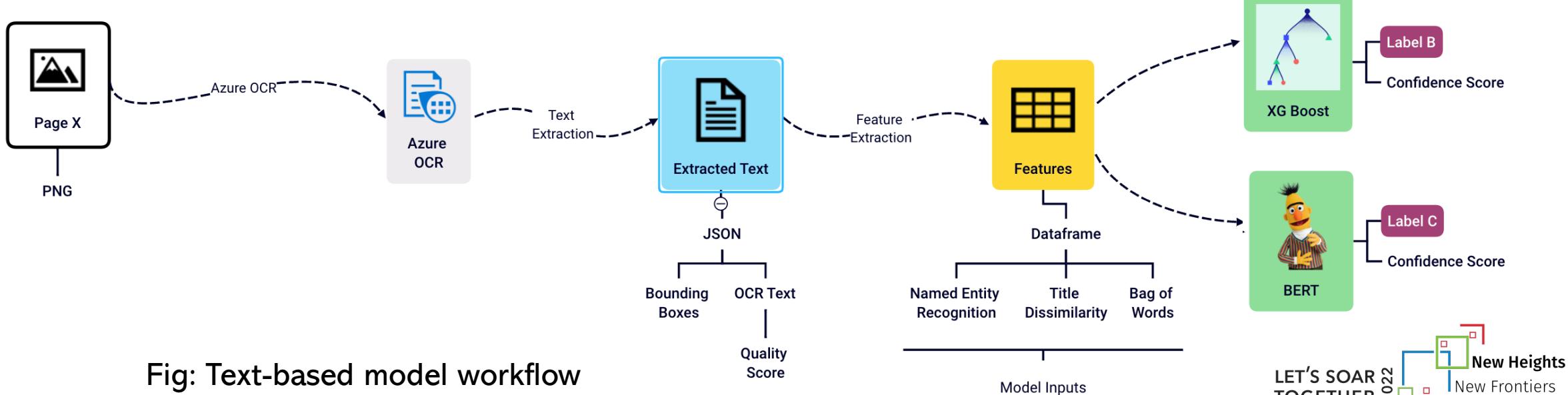


Fig: Text-based model workflow

# Heuristic-based used in case of uncertain documents (Capability Design)

## Our Implementation



Title Extraction by heuristically narrowing on 'potential' title of the document

Input: Image (document)

Output: Category label and heuristic title score

Intuition: A heuristic title is informative than predicting 'Other Document'

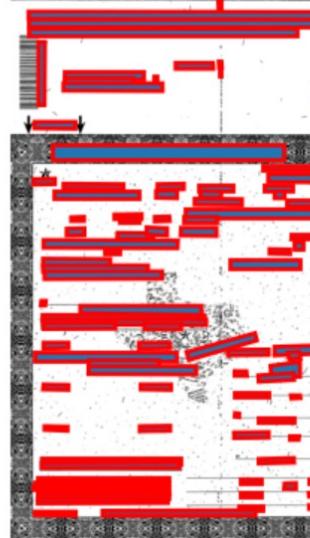
### Highlights:

- Used all bounding boxes from Azure OCR as potential titles; **narrowed search space**
- Engineered features like location on page, bold, font size, lone item, etc.

## Original



## Potential Title



## Final Candidates

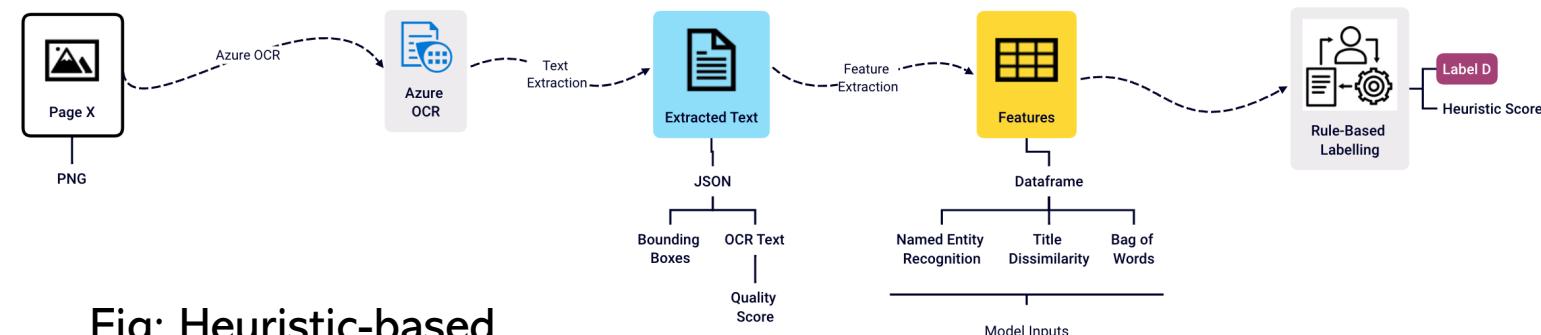
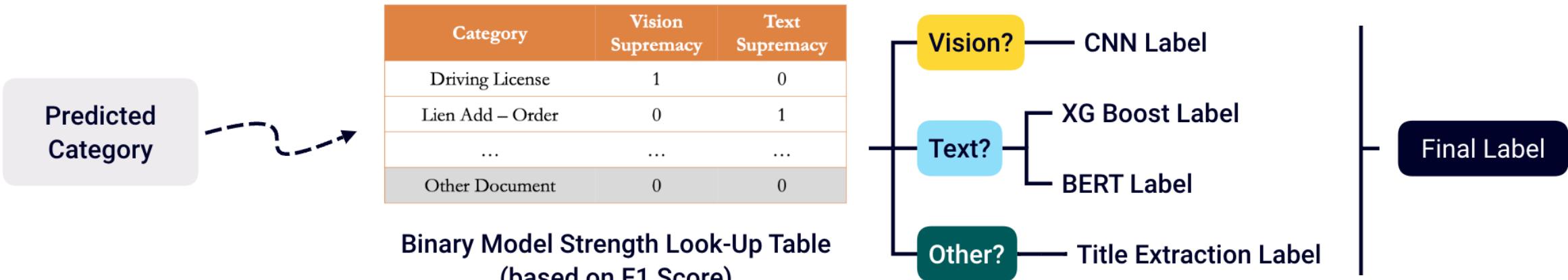


Fig: Heuristic-based model workflow

# The labels need to be ensembled into one single prediction



Case	Vision Label	Text Label	Vision Supremacy	Text Supremacy	Final Model	Final Label	Comments
<b>Case 1: Same output</b>	A	A	1	0	Same	A	
<b>Case 2: Conflicted output</b>	A	B	1	1	NA	NA	Manual Review
<b>Case 3: One model is supreme</b>	A	B	1	0	Vision	A	
	A	B	0	1	Text	B	
<b>Case 4: Both "Other"</b>	Other	Other	0	0	Heuristic	Heuristic Title	Manual Review
<b>Case 5: One "Other" and heuristic title contains another model's output</b>	Other	B	0	0	Text	(Heuristic ∩ B) = B	
	A	Other	0	0	Vision	(Heuristic ∩ A) = A	
<b>Case 6: One "Other" and heuristic does NOT contain another model's output</b>	Other	B	0	0	Heuristic	Heuristic Title	Manual Review
	A	Other	0	0	Heuristic	Heuristic Title	Manual Review
<b>Case 7: Both non-"Other" and None supreme</b>	A	B	0	0	Heuristic	Heuristic Title	Manual Review

# Result | Final model predictions with added visibility

Docket #	Page #	Final Model	Final Prediction	Prediction Confidence	Document Quality	Notes
616823_1	1	Matched	Lien Add - Order	1.00	0.699	
616823_1	2	Matched	FORM 30-C Page 1 of 2	1.00	0.766	
616823_1	3	Text	MV-50	0.49	0.640	
616823_1	4	Vision	Driver's License Page 1 of 2	0.88	0.292	
616823_1	5	Heuristic	Binder of Insurance	0.68	0.612	Manual Review

Fig: Result Example

## Our Implementation

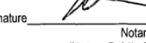
'Classification API' runs  
the entire pipeline on  
Virtual Machine

## Highlights for Results:

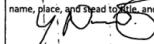
- The results show confidence score for the chosen model
- OCR score give a proxy to the document quality
- Manual review in the 'notes' gives user pin-pointed areas to spend more time labelling

# A glimpse of Champion vs Challenger capabilities

(1) Limited Power of Attorney/Motor Vehicle Transaction

LIMITED POWER OF ATTORNEY MOTOR VEHICLE TRANSACTION	
This Limited Power of Attorney (the "POA") is made and effective this 27 <sup>th</sup> day of September, 2021.	
BETWEEN:	CT Corporation System (Company/Agency), 187 Wolf Road, Suite 101, Albany NY 12205
LOCATED AT:	
OR:	
AND:	Owner First Student, Inc. and Co-owner _____ (the "Client(s")) whose residence is located in the State of Ohio at 600 Vine Street, Suite 1400, Cincinnati, OH 45202
The undersigned does/do hereby make, constitute, and appoint the above Company to their authorized appointee (including its employees, personnel, or agents) their true and lawful Attorney-in-Fact to represent them with respect to the following described vehicle:	
See Appendix	
to act for me, in applying for an original or duplicate certificate of title, to register, transfer title, or record a lien to the motor vehicle, mobile home or vessel described below, and to print my name and sign their name in my behalf. My attorney-in-fact can also file for and collect any overpayment of fees and taxes from the state office to which said overpayment was remitted. My attorney-in-fact can also do all things necessary to the application or any other related instrument and to bind me in as sufficient a manner as I myself could do, were I personally present and signing the same.	
With full power of substitution and revocation, I hereby ratify and confirm whatever my said attorney-in-fact may lawfully do or cause to be done in respect of registering or titling the below-referenced vehicle.	
 Owner (Applicant) signature _____	
 Co-owner (Co-applicant) signature _____	
State of Ohio County of Hamilton	
On this the 3 <sup>rd</sup> day of January, 2022 before me	
Appeared Michael Petrucci or whose identify was proven to me on the basis of satisfactory evidence to be the person who he or she claims to be, and acknowledged that he or she signed the above document. (Name of Signer)  Signature _____ Notary (Notary Public Signature)	
 ANDREW WESLEY PUGH Notary Public, State of Ohio My Commission Expires 05-20-2024	
THIS FORM IS NOT VALID WITHOUT A NOTARY	
Rev 04.22.21	

(2) VTR-271

Limited Power of Attorney for Eligible Motor Vehicle Transactions				
<b>Information</b> All sections of this form must be properly completed in order for this document to be accepted. Original signatures are required, only black or blue ink are acceptable, and no alterations are allowed on this form.				
This completed and signed form grants the grantee, with full power of substitution, full power and authority to perform every act necessary and proper to purchase, transfer, and assign the legal title to the motor vehicle described on behalf of the grantor. "Full power of substitution" means that whoever is given this power of attorney may delegate that power by putting another person in his or her place by a substitute power of attorney.				
This power of attorney cannot be used in a dealer transaction to complete a title assignment on a motor vehicle subject to federal odometer disclosure. Federal law specifies a motor vehicle is subject to odometer disclosure if it is self-propelled, less than 10 years old, and has a gross vehicle weight of 16,000 pounds or less. In compliance with federal law, the secure Power of Attorney for Transfer of Ownership to a Motor Vehicle (Form VTR-271-A) must be used when use of a power of attorney is permitted by the applicable regulations for a vehicle subject to federal odometer disclosure.				
If a power of attorney is used to apply for title, initial registration, or a certified copy of title, the grantor (person signing this form) and the grantee (person signing the application) must include a photocopy of their photo identification as required by state law.				
<b>Vehicle Information</b>				
Vehicle Identification Number 5J8TB3H37HL018546	Year 2017	Make Acura	Body Style LL	Model RDX
License Plate State and Number (if any) RKS 3110	Title/Document Number (if unknown, leave blank) 2810953448810053108			
<b>Grantor Information</b>				
First Name (or Entity Name) Yaoasca Indira Pastora Torres	Middle Name	Last Name	Suffix (if any)	
Address 5419 WOLF ROCK DR				
City KATY	County HARRIS	State TX	Zip 77449	
<b>Grantee Information</b>				
First Name (or Entity Name) CT Corporation System	Middle Name	Last Name	Suffix (if any)	
Address 187 Wolf Rd Suite 101				
City Albany	County Albany	State NY	Zip 12205	
<b>Certification – State law makes falsifying information a third degree felony</b>				
I, the grantor of the county and state as listed above, owner of the motor vehicle described above, certify that I do make, constitute, and appoint the grantee of the county and state as listed above, or to anyone the grantee may substitute, my true and lawful attorney, for me and in my name, place, and stead to title, and to allow my attorney the authority to substitute as it pertains to the motor vehicle described above.				
 Yaoasca Indira Pastora Torres Signature of Grantor Printed Name (Same as Signature) Date 10-28-22				
Form available online at <a href="http://www.TxDMV.gov">www.TxDMV.gov</a>				
Page 1 of 1				

Champion Model

Keyword matching +  
extensive rule-based system



Both documents have similar keywords  
Champion: (1) VTR-271 (2) VTR-271

Challenger Model

Multi-Modal Approach  
(Vision + Text)



Formats are different. Text is different.  
(1)'s content is like "Limited Power..."  
(2)'s format is like VTR-271

Challenger: (1) Limited Power... (2) VTR-271

# Result | A glimpse of Champion vs Challenger capabilities

Document Type (19)	Champion vs. Challenger Model Performance				Better Model
	Champion Recall	Vision Recall	BERT Recall	XGBoost Recall	
130-U Page 1 of 2	90.9%	95.0%	100.0%	100.0%	BERT/XGBoost ✓
Check Endorsement	100.0%	100.0%	100.0%	100.0%	All ✓
Driver's License Page 1 of 2	67.3%	95.0%	92.7%	85.5%	Vision ✓
Fees Payment Check	97.3%	94.0%	88.9%	89.2%	Champion
FORM 30-C Page 1 of 2	97.3%	91.0%	96.2%	60.5%	Champion
FORM 30-C Page 2 of 2	100.0%	98.0%	89.2%	100.0%	XGBoost/Champion ✓
Lien Add - Order	95.5%	95.0%	86.4%	93.2%	Champion
Limited Power of Attorney/Motor Vehicle	100.0%	82.0%	100.0%	100.0%	BERT/XGBoost/Champion ✓
MV-50	100.0%	100.0%	100.0%	100.0%	All ✓
Odometer Disclosure Statement	71.4%	75.0%	71.4%	100.0%	XGBoost ✓
One and Same Statement	93.3%	81.0%	73.3%	73.3%	Champion
Other Document	96.5%	93.0%	87.6%	91.7%	Champion
Search Request Form	93.6%	100.0%	98.9%	100.0%	Vision/XGBoost ✓
Texas Liability Insurance Card	0.0%	5.0%	100.0%	66.7%	BERT ✓
Texas Vehicle Record	0.0%	100.0%	100.0%	100.0%	Vision/BERT/XGBoost ✓
TS-5	100.0%	100.0%	33.3%	100.0%	Vision/XGBoost/Champion ✓
TS-8	100.0%	100.0%	100.0%	100.0%	All ✓
VTR-271	98.3%	99.0%	100.0%	100.0%	BERT/XGBoost ✓
WKLS Authorization letter	97.6%	99.0%	100.0%	100.0%	BERT/XGBoost ✓
<b>Accuracy</b>	<b>87.6%</b>	<b>86.3%</b>	<b>88.4%</b>	<b>88.4%</b>	<b>Challenger Wins: 14/19</b>

# High learning + Solid deliverables + Strong impact

Our Deliverables to Wolters Kluwer

## (1) OCR API

Input: PDF Docket

Output: Azure OCR response with document quality proxy

## (2) CAPABILITY API

Input: PDF Docket (+OCR Response)

Output: Azure OCR response with document quality score

## (3) Documentation

- Used all bounding boxes from Azure OCR as potential titles; narrowed search space
- Engineered features like location on page, bold, font size, lone item, etc.

Our Deliverables to MIT\*

## (1) Capstone Presentation

It shows the entire journey of the Capstone project and to be presented if selected in Top 10 Capstone teams

## (2) Capstone Poster

It is to be presented during Capstone Showcase

## (3) Capstone Report

It contains every technical detail of the Capstone project

## (4) Capstone Video

A 2-3 minute video explaining the project (voted by public globally)

Learnings from Capstone with WK

## (1) Industry Knowledge

## (2) Business Acumen

## (3) Mentorship

## (4) Collaboration

## (5) Consistent Support

## (6) Challenged!

## (7) Awesome Team!

# A big thanks to the Wolters Kluwer family for the constant support!



Abhishek (VP, Data Analytics), Rajiv (VP, Data Science),  
Pooja (Data Scientist Manager),  
Varun (Lead Business Systems Analyst)

[Core Team]



Lucas (Onboarding), Savannah (Business), Erin (Proof of Concept), Florence (Experimentation), Piyush (Technical),  
Ameya (Sr. Data Scientist)

[Mentor Match Series]