

AI-Based Automated Customs Declaration Recognition and Intelligent Archiving System

Jerry Xing

A01354731

Introduction

Many companies still rely on manual transcription when processing import customs declarations—copying key information from scanned images into Excel for statistics, accounting, and financial management. This process is repetitive, time-consuming, and error-prone. Through communication with my family’s enterprise, which handles financial operations, we confirmed that this workflow is highly automatable.

Therefore, this project aims to build an AI-powered document-understanding-based customs declaration recognition and intelligent archiving system. The goal is to locate individual documents in combined scans automatically, understand handwritten and printed content, extract key fields, and produce structured data—significantly reducing manual input costs and improving accuracy and efficiency.

Background and Problem Statement

Based on communication with my family members, in real business scenarios, customs declarations are often batch-scanned or photographed by freight forwarders, producing large images/PDFs containing multiple declarations. These documents include both printed and handwritten content, as well as issues such as overwriting, cursive writing, residual ink, scattered noise, blur, or tilt.

Currently, financial staff must manually transcribe all information into Excel, then follow a “manual entry → manual verification” dual process. Workload grows linearly, or worse, with document volume, while still encountering errors such as misread numbers, missing fields, or currency confusion.

Given this context, what enterprises need is not a simple OCR tool, but a complete end-to-end document AI pipeline:

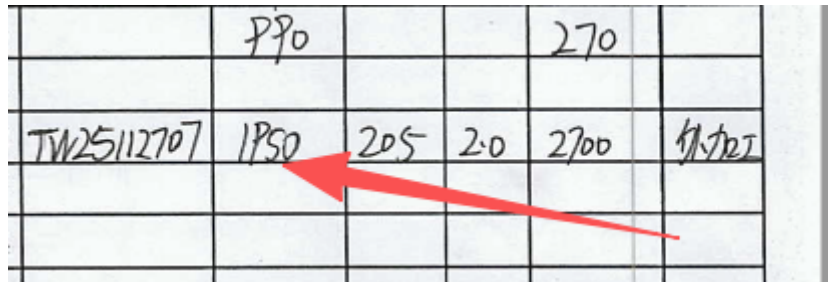
One that can automatically segment multiple declarations from combined scans, remain robust against handwriting/noise, extract fields, produce structured results, and support Excel/report generation and downstream analysis.

Technology Review and Existing Solutions

1. Capabilities and Limitations of Traditional OCR

Traditional OCR (e.g., Tesseract, printed-text PaddleOCR models, cloud OCR services) works well on clean printed text but faces inherent bottlenecks for this project's real-world data:

- **Cannot process multi-document segmentation in combined scans:** OCR “recognizes text,” but does not “locate multiple declarations from one large image.”
- **Poor robustness to handwritten/connected/overwritten content:** HTR (Handwritten Text Recognition) is a different domain from printed OCR.
- **Similar-looking characters:** Some handwriting habits commonly produce confusing pairs: “2” vs “z”, “9” vs “q”, “l (uppercase i)” vs “l (lowercase L)” vs “1”. Heavily reliant on contextual reasoning. (Page 2, Figure 1 shows a handwritten “9” easily mistaken as the letter “p”)



Page 2, Figure 1

- **Lack of semantic judgment for noise strokes:** OCR often mistakes non-text strokes/noise as characters.

Typical Failure Cases in Noisy Documents

- **Dry-ink test strokes:** Writers often scribble lines to check ink flow; traditional OCR misreads them as characters: 11111, /////, ---, — — —.
- **Ink leakage/smudges / line-shaped noise:** Sparse spots may be misread as punctuation or vertical bars. (Page 3, Figure 2 shows a long vertical noise line caused by scanning artifacts.)

ЈЕДИНСТВЕНА ЦАРИНСКА ИСПРАВА				А ЦАРИНАРНИЦА ИЗВОЗА/ОДРЕДИШТА			
2. Подносилац/Извозник Б.Р. HRVATSKE ŠUME DOO ZAGREB/HR				1. ДЕКЛАРАЦИЈА UV 4 3. Обрзаци 1 1 4. Тип лист 1 5. Редовност 1 6. Број листа 19 7. Референтни број			
8. Прималац Б.Р. 112580822 TERRA DRVO DOO NIKOLE TESLE BB, LJUKOVO				9. Лице одговорно за физикалну контролу			
14. Персонална исправа Б.Р. N 111622667 SIRMILUM ŠPED DFD DOO SREMSKA MITROVICA, JARACKI PUT BB				15. Земља створила/извоз HRVATSKA 16. З. отпр. вел.Лист HR 17. Земља одређиште HR Hrvatska			
18. Идент. и наз. производа од дрвета у грмској SAZ28CM/14284SA				19. Код 0			
21. Идент. и наз. извозног производа поје претплату границу				22. Вредност и вредност на фактури EUR 4.994,76			
23. Вредност на фактури				24. Вредност на фактури 117,387500			
25. Вредност на фактури 33				26. Вредност на фактури 33			
27. Место извоза				28. Место извоза			
29. Царинарска испостава улаза 21091				30. Место рада			
31. Подносилац и кода 21091				32. Место 1			
33. Шифра роба 44039100				34. Шифра и порекло 0000 M3 L3/02 2320			
35. Шифра и порекло 21.200,00				36. Бруто маса у кг 21.200,00			
37. Шифра и порекло 21.200,00				38. Нето маса у кг 21.200,00			
39. Шифра и порекло 25RS021091N4M6Q1J2/1				40. Шифра и порекло 25RS021091N4M6Q1J2/1			
41. Датум извоза 18.25				42. Цена роба 4.994,76			
43. М.Б. 1				44. М.Б. 1			
45. Шифра и порекло 663.893,20				46. Шифра и порекло 663.893,20			
47. Обрзаци 01 09 663.893,20 20,00				48. Обрзаци W			
49. Обрзаци W				50. Обрзаци W			
51. Обрзаци W				52. Обрзаци W			
53. Обрзаци W				54. Обрзаци W			
55. Обрзаци W				56. Обрзаци W			
57. Обрзаци W				58. Обрзаци W			
59. Обрзаци W				60. Обрзаци W			
61. Обрзаци W				62. Обрзаци W			
63. Обрзаци W				64. Обрзаци W			
65. Обрзаци W				66. Обрзаци W			
67. Обрзаци W				68. Обрзаци W			
69. Обрзаци W				70. Обрзаци W			
71. Обрзаци W				72. Обрзаци W			
73. Обрзаци W				74. Обрзаци W			
75. Обрзаци W				76. Обрзаци W			
77. Обрзаци W				78. Обрзаци W			
79. Обрзаци W				80. Обрзаци W			
81. Обрзаци W				82. Обрзаци W			
83. Обрзаци W				84. Обрзаци W			
85. Обрзаци W				86. Обрзаци W			
87. Обрзаци W				88. Обрзаци W			
89. Обрзаци W				90. Обрзаци W			
91. Обрзаци W				92. Обрзаци W			
93. Обрзаци W				94. Обрзаци W			
95. Обрзаци W				96. Обрзаци W			
97. Обрзаци W				98. Обрзаци W			
99. Обрзаци W				100. Обрзаци W			

Page 3, Figure 2

- **Ambiguous dots:** Users may dot near fields, or parts of digits may look like a dot, causing decimal/period errors (often disastrous). (Page 3, Figure 3 shows the number “205” appearing like “2.05” due to table obstruction.)
- **Tilted scans:** Slanted pages degrade recognition quality. (Page 3, Figure 3 shows an example of a tilted document.)

15		1240		180		15
16						16
17	TW25112803	1950	205	2.0	2884	17
18						18
19	TW25112804	1940	190	2.0	353	19

Page 3, Figure 3

- **Overwriting:** Modifications may resemble another character (e.g., a slashed “0” that looks like “8”).

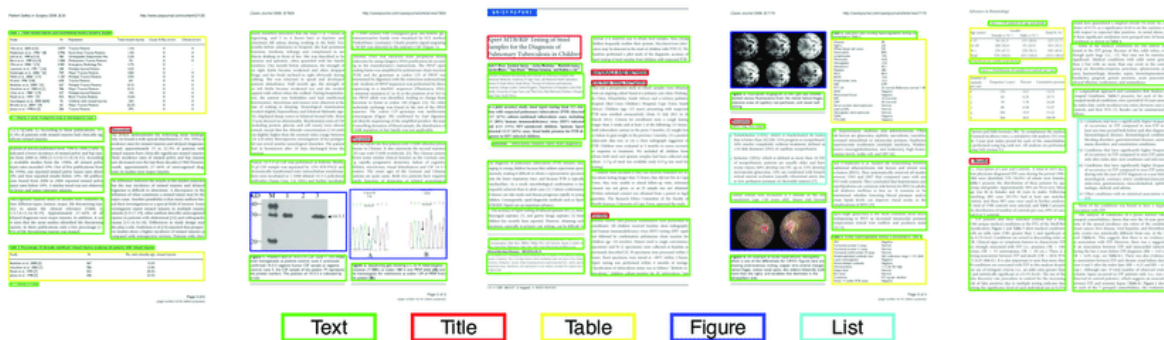
Conclusion

Under the conditions of “combined scans + handwriting + noisy documents,” traditional OCR **alone** is neither viable nor reliable.

2. Document Instance Detection and Layout Analysis

To extract each customs declaration from a combined scan, page layout and instance detection techniques are required:

- **Instance detection/segmentation:** [YOLO](#) models can detect bounding boxes of declarations and crop them out.
- **Layout analysis:** [LayoutParser](#) can identify table lines, cell positions, field areas, and content regions. (Page 3, Figure 4 shows an example of layout recognition.)



Page 3, Figure 4

3. Handwritten Text Recognition (HTR) and Semantic Validation

- **HTR models:** Handwriting-oriented [PaddleOCR](#) models, capable of handling connected writing and stylistic variations.
- **Multimodal Large Models (VLMs):** Vision-capable general models (e.g., [GPT5.1](#)) can combine context to perform semantic validation—e.g., distinguishing noise strokes from real characters.

4. Key Information Extraction (KIE) and Table Mapping

- [LayoutLMv3](#) can combine text and spatial layout to perform KIE, reconstruct table structures, and map them to a defined schema.

Proposed Solution

Overall Concept

Multi-document segmentation → Handwritten/printed mixed recognition → Field extraction → Structured validation → Export

Processing Workflow

1. **Batch Input and Preprocessing**
 - Support uploading combined scanned images/PDF.
 - Perform de-skewing, de-noising, and contrast enhancement.
 - Rough document region segmentation; remove blank borders.
2. **Document Detection and Instance Segmentation**
 - Use YOLO to detect boundaries of each customs declaration and crop them out.

3. **Text and Content Recognition**
 - **Router:** determines whether the region is printed/handwritten/mixed.
 - **HTR branch:** PaddleOCR for handwritten content.
 - **Printed-text branch:** Traditional OCR for printed text.
4. **Key Field Extraction (KIE) and Semantic Validation**
 - LayoutLMv3 outputs structured JSON.
 - **Semantic and constraint checks:**
 - Number/quantity/currency/date type & range validation
 - Currency allowlist (CNY/USD/EUR/etc.)
 - Date normalization
 - Contextual correction of decimal/punctuation anomalies
 - **Noise suppression:** Using VLM (e.g., removing dry-ink lines, dots, overwriting artifacts)
5. **Structured Output and Back-Filling**
 - Map results to Excel.
 - Produce “one combined scan → multiple Excel records.”
 - Save cropped local patches as evidence for human verification.
6. **Human-AI Collaboration**
 - Review interface showing values + image regions.
 - One-click correction and active learning: corrected samples fed back for continuous improvement.

Expected Outcomes

- Automatically detect and split each declaration from combined scanned images containing multiple documents.
- Achieve ~80% accuracy for key-field extraction.
- Reduce manual entry time by at least 60%.
- Maintain robustness under noisy document conditions: dry-ink strokes, smudges, tilt, and overwriting.
- Support exporting standardized Excel outputs for downstream financial analysis.

Minimum Viable Product

- **Input:** A combined scan containing 5–20 customs declarations.
- **Output:** One structured record (Excel/CSV/JSON) per declaration.
- **Must satisfy outcomes described above.**

Technologies Used

- **Backend & workflow orchestration:** Java + Spring Boot
- **Document detection & segmentation:** Python + PyTorch ([YOLO](#))
- **Handwriting recognition (HTR):** [PaddleOCR](#)
- **KIE & rule-based validation:** [LayoutLMv3](#) + regex + business rules
- **Data export:** Excel, JSON
- **Frontend review tool:** Web-based verification & correction interface