# Language, OCR, Form Independent (LOFI) pipeline for Industrial Document Information Extraction

Agile SoDA

Chang Oh Yoon[1,+], Wonbeen Lee[1,+], Seokhwan Jang[1,+] ,
Kyuwon Choi[2,+] , Minsung Jung[2,+] ,
Daewoo Choi[3,*] ,

+AgileSoDA, *Hankuk University of Foreign Studies

# 1. Introduction

Many industries handle complex documents known as Visually Rich Documents (VRDs).



Koream medical bills



Japanese receipts

# 1. Introduction

In real-world industry, we should consider a process of SER (Semantic Entity Recognition) to automate workflows.



Koream medical bills

Patient
- ID
- Name
- Period : 2023-03-07 ~ 2023-03-09
- Class : 외래

Medical Treatment
- Category : 진찰료 , …
- Date : 2023.03.07 , …
- Item : 외래환자 의약품 … , …
- Item Code : AL801 , …
- # of Days : 1 , …
- Qty/Dose : 1 , …
- Unit Price : 220 , …
- Price : 220 , …

# 2. Challenges

To address the automation demands of the industry, we face three main challenges.

## 01

### Low Resource Language

- There are limited VRD datasets available for Low-Resource Languages.
- No pre-trained models exist for these languages.
- This scarcity hinders the creation of advanced language models.

## 02

### OCR Dependency

- SER has limitations due to OCR engine output.
- OCR results are typically at the word level, not entity level.
- Additional processing (splitting or combining) may be needed for accurate semantic entities.

## 03

### Form Diversity

- Industry documents pose challenges for information extraction due to custom formats.
- Even standardized forms have variations in formatting, such as custom medical report templates.
- Image distortions or rotations can alter a document's structure and further complicate extraction.

# 3. Language, OCR, Form Independent (LOFI) pipeline

So, we present a **L**anguage, **O**CR, and **F**orm **I**ndependent pipeline, named **LOFI pipeline**.

# 3. Language, OCR, Form Independent(LOFI) pipeline

We constructed a token-level box splitting to standardize bounding box ranges from various OCR engines.



✓ Input text & layout information requires to be synchronized in token level

✓ Output results are Also decoded from input token level.

✓ With Token level box splitting, the process can be independent from OCR Engine

# 3. Language, OCR, Form Independent(LOFI) pipeline

We implemented a language flexible multimodal model for Low-Resource Language(LRL).

# 3. Language, OCR, Form Independent(LOFI) pipeline

We implemented a language flexible multimodal model for Low-Resource Language(LRL).

# 3. Language, OCR, Form Independent(LOFI) pipeline

We added SPADE decoder for operating independently of document formats and layouts.



*If the input tokens are out of order, (traditional)* **IOB tagging will not work.**

*Initial token classification*

*Subsequent token classification*

☐ Product name    ☐ Product price    ↷ Subsequent token

SPADE Decoder

OCR Independence    Language Independence    Form Independence

# 3. Language, OCR, Form Independent(LOFI) pipeline



Token-level Box Split

OCR & Text Alignment

Model Inference

[CLS] ... 非 * 09 1 バック ¥ 5,9 90 ( 非 課 ... [SEP]

Layout Encoder Layers

Text Encoder Layers

SPADE Decoder

Outputs

```
[
    {
        "Entity" : "product_name",
        "Text" : "非*091バック",
        "Boundingbox" : {
            [440, 1611, 1042, 1722]
        }
    },
    {
        "Entity" : "total_tax_amount",
        "Text" : "¥5,990",
        "Boundingbox" : {
            [1228, 1711, 1604, 1819]
        }
    }
]
```

OCR Independence　　Language Independence　　Form Independence

# 4. Experiments

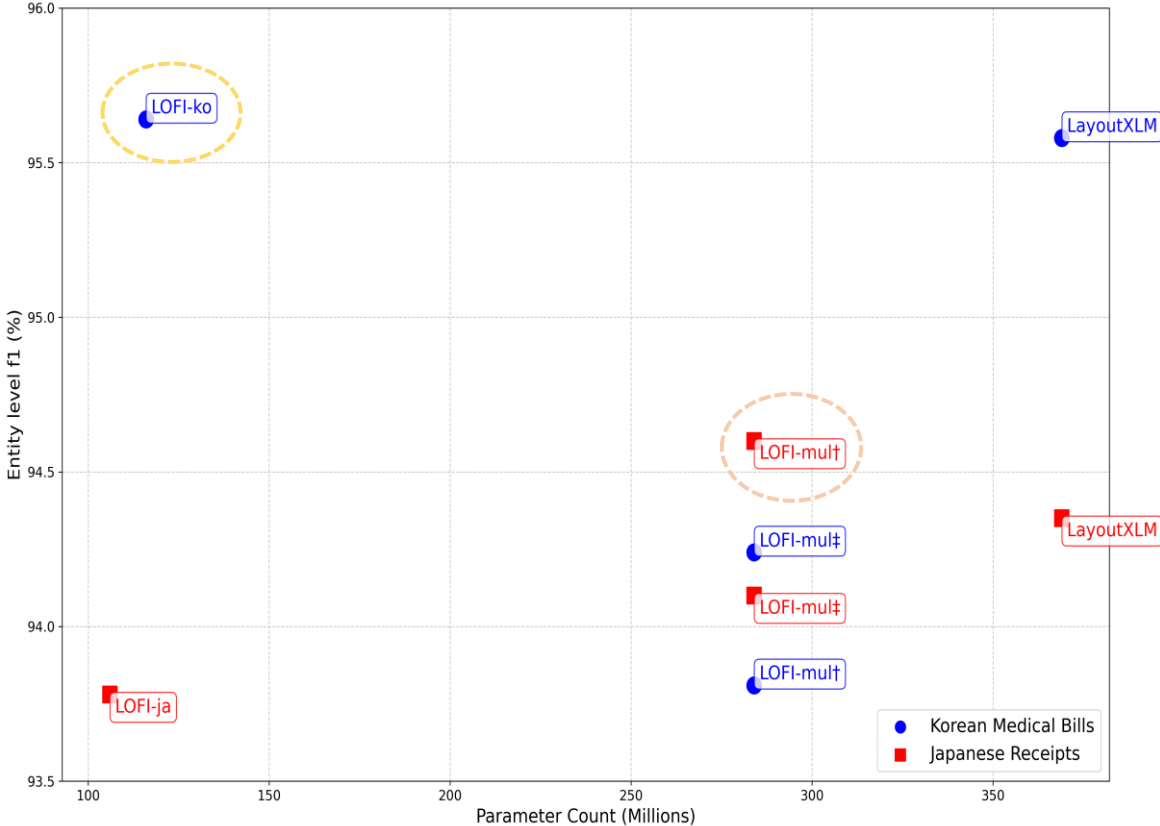LOFI shows better performance than LayoutXLM on Korean medical bills and Japanese receipts, also demonstrating efficiency in terms of parameters and computational resources.

## Number of model parameters and entity-level F1 score
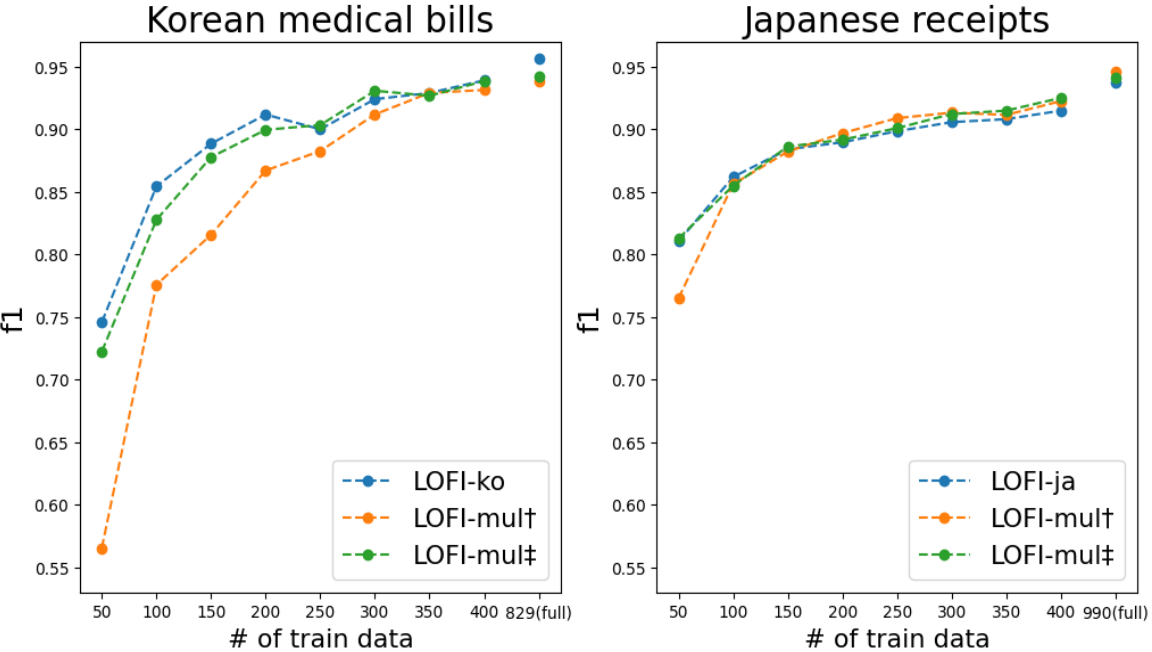


### LRL Documents

| Name | Pretrained Language | Encoder | Params | Korean Medical Bills | Japanese Receipts |
|---|---|---|---|---|---|
| LayoutXLM | Multilingual | LayoutXLM-base | 369M | 95.58% | 94.35% |
| LOFI-mul † | Multilingual | InfoXLM-base + lilt-only-base | 284M | 93.81% | 94.60% |
| LOFI-mul‡ | Multilingual | XLMRoBERTa-base + lilt-only-base | 284M | 94.24% | 94.10% |
| LOFI-ko | Korean | RoBERTa-base + lilt-only-base | 116M | 95.64% | – |
| LOFI-ja | Japanese | RoBERTa-base + lilt-only-base | 106M | – | 93.78% |

### English Documents

| Name | Pretrained Language | Params | FUNSD | CORD |
|---|---|---|---|---|
| LayoutLM | English | 160M | 79.27% | 94.72% |
| LayoutLMv2 | English | 200M | 82.76% | 94.95% |
| LayoutLMv3 | English | 133M | 79.38% | 96.80% |
| BROS | English | 110M | 83.05% | 95.73% |
| LOFI-en | English | 131M | 78.99% | 96.39% |

# 4. Experiments

Experiment results showing performance variations with different training data sizes used in fine-tuning.



- Through experiments, it was suggested that at least 300-400 training data is required to achieve satisfactory performance.

- The amount of training data required may vary by language; using fewer than 200 training documents resulted in a 5% difference in performance compared to using the full training dataset.

# Contributions

- Constructed a flexible pipeline structure, LOFI (Language, OCR, Form independent Extraction) *to account for multiple challenges in industrial data extraction.*

- The LOFI pipeline demonstrates satisfactory performance on Korean and Japanese *datasets without additional pre-training.*

- *Empirical evidence on industrial applicability* of the LOFI pipeline by successfully implementing it in insurance claim processing and tax handling operations.

# Future research

- *Data augmentation* techniques to enhance the robustness of the LOFI pipeline.

- *Efficient annotation methods* to reduce the annotation burden in SER tasks.

- *Improved decoder architectures* to handle complex document challenges and diversify AI capabilities for business scenarios.

Thank you !