



バラバラのドキュメント解析技術と実データにおける課題

2024/12/14

cvpaper.challenge Conference winter 2024

株式会社 LayerX Naoto Shimakoshi(@nt_4o54)



島越 直人 (Naoto Shimakoshi)

バクラク事業部 AI-OCRグループ Tech Lead/ 機械学習エンジニア

経歴

- 2019/04 京都大学大学院 工学研究科 修士課程修了
- 新卒では、DeNAでタクシー配車アプリに関する機械学習システムの構築や、ライブストリーミングサービスにおける推薦システム構築に携わる
- 現在
 - 株式会社LayerX AI-OCRグループ Tech Lead
 - バクラク事業部において、AI-OCRの改善や新しい機械学習システムの構築を担当
 - Kaggle Competitions Grandmaster



Agenda

目次

- LayerXについて
- Introduction
- Bakuraku AIの紹介
- 実データにおける課題
- 今後の展望

LayerXについて

「すべての経済活動を、デジタル化する。」をミッションに掲げ、法人支出管理サービス「バクラク」や企業内業務のデジタル化を支援するサービスを提供しています。

バクラク事業

企業活動のインフラとなる法人支出管理（BSM）SaaSを開発・提供



Fintech事業

ソフトウェアを駆使したアセットマネジメント・証券事業を合併会社にて展開



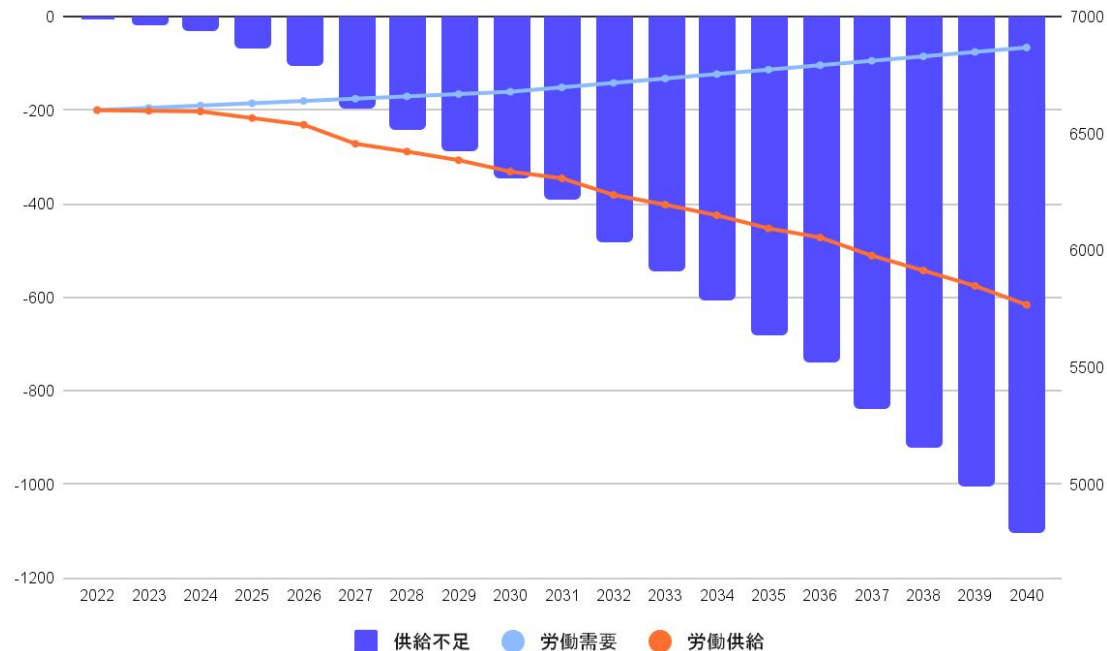
AI・LLM事業

文書処理を中心とした、LLMの活用によるプロセスのリデザイン

AI・LLM 事業



日本の労働需給ギャップは深刻



2040年に労働需給ギャップ
1100万人

日本全体に必要な生産性
+20%

バックオフィス(BSM・HRM)領域のAI SaaS。AIにより業務そのものをなくし、生産性を改善





複雑・短期間・ミスできない

ルール、法律、業種・業態、社内事情に対応するための複雑な業務

毎月、短期間で大量の業務を完遂

それを「ミスなく」完遂



業務に高頻度で潜むアナログな手間

- 請求書を1枚1枚スキャンする手間
- 領収書をシステムに手入力する手間
- 仕訳を作成する手間
- カード明細と領収書との突合の手間
- 書類の情報を入力する手間



Bakuraku AI

- 請求書をAIが自動分割して取り込み
- 領収書のデータをAIが入力
- AIが過去に学習した仕訳を入力
- 領収書とカード明細をAIが紐付け
- AIが書類種別を判定してラベル付け

AIにより、アナログな手間を無くしていく

Introduction

ドキュメント解析技術とCV分野の繋がり

ドキュメント解析とComputer Vision分野との繋がり

請求書

2024年06月30日
書類番号:SP202406-0008

株式会社LayerX
〒103-0012
東京都中央区日本橋堀留町
1丁目9-8
人形町PREX 2F
電話:01234567890
担当者:バクラク太郎
登録番号:T9010401140088

株式会社バクラク
バクラク太郎 様

下記のとおり、御請求申し上げます。

ご請求金額 **¥218,000 (税込)**

蛍光ペンを引いてるところを認識したい?

前回ご請求金額	入金額	繰越金額	今回請求金額 (税抜)	今回消費税額	今回請求金額 (税込)
218,000	218,000	0	200,000	18,000	218,000

文字が大きい方を取りたい?

10%対象金額	100,000	消費税額	10,000	税込金額	110,000
8%対象金額	100,000	消費税額	8,000	税込金額	108,000

表形式で縦と横の依存関係が混在

ロゴや印影・ハンコから認識したい?

ドキュメント解析とComputer Vision分野との繋がり



ドキュメントの内容を理解するには
視覚情報 (Vision)とテキスト情報 (Language)、位置関係 (Layout)
を組み合わせて理解する必要がある

文字が大きい方を取りたい？

前回ご請求金額	入金額	繰越金額	今回請求金額 (税抜)	今回消費税額	今回請求金額 (税込)
218,000	218,000	0	200,000	18,000	218,000
10%対象金額	100,000	消費税額	10,000	税込金額	110,000
8%対象金額	100,000	消費税額	8,000	税込金額	108,000

蛍光ペンを
引いてるところを
認識したい？

表形式で縦と横の
依存関係が混在

ドキュメント解析とComputer Vision分野との繋がり

PaddlePaddleでは、LLM以前からドキュメント解析を複合的な技術を組み合わせて行っている

PP-StructureV2: A Stronger Document Analysis System

Chenxia Li, Ruoyu Guo, Jun Zhou, Mengtao An,
Yuning Du, Lingfeng Zhu, Yi Liu, Xiaoguang Hu, Dianhai Yu

Baidu Inc.
{lichenxia, zhulingfeng}@baidu.com

Abstract

A large amount of document data exists in unstructured form such as raw images without any text information. Designing a practical document image analysis system is a meaningful but challenging task. In previous work, we proposed an intelligent document analysis system PP-Structure. In order to further upgrade the function and performance of PP-Structure, we propose PP-StructureV2 in this work, which contains two subsystems: Layout Information Extraction and Key Information Extraction. Firstly, we integrate Image Direction Correction module and Layout Restoration module to enhance the functionality of the system. Secondly, 8 practical strategies are utilized in PP-StructureV2 for better performance. For Layout Analysis module, we introduce ultra light-weight detector PP-PicoDet and knowledge distillation algorithm FGD for model lightweighting, which increased the inference speed by 11 times with comparable mAP. For Table Recognition model, we utilize PP-LCNet, CSP-PAN and SLAHead to optimize the backbone module, feature fusion module and decoding module, respectively, which improved the table structure accuracy by 6% with comparable inference speed. For Key Information Extraction model, we introduce Vi-LayoutXML which is a visual-feature independent LayoutXML architecture, TB-YX sorting algorithm and U-DML knowledge distillation algorithm, which brought 2.8% and 9.1% improvement respectively on the Hmean of Semantic Entity Recognition and Relation Extraction tasks. All the above mentioned models and code are open-sourced in the GitHub repository PaddleOCR.

Document Layout Analysis can be regarded as an object detection task for document images in essence. The basic units such as titles, paragraphs, tables, and illustrations in the document are the objects needed to be detected and recognized. Layout-parser(Shen et al. 2021) is a unified toolkit for Deep Learning Based Document Image Analysis, which comes to state-of-the-art on PubLayNet dataset(Zhong, Tang, and Yepes 2019). In PP-Structure, we use PP-YOLOv2(Huang et al. 2021) to complete the layout analysis task, which is real-time on GPU devices. However, currently proposed models are not CPU-friendly and thus not conducive to deployment on CPUs or mobile devices.

Table Recognition is used to convert table images into editable Excel format files. The diversity of tables in document images, such as various rowspans and colspans and different text types, makes table recognition a hard task in document understanding. There are many table recognition methods, such as traditional algorithms based on heuristic rules and recently developed methods based on deep learning. Among them, the end-to-end method has received extensive attention due to the simplicity of the pipeline, which represent the table in HTML format and adopt Seq2Seq(Sutskever, Vinyals, and Le 2014) to predict the table structure, such as TableRec-RARE(Du et al. 2021b) in PP-Structure powered by PaddlePaddle(Ma et al. 2019). In TableMaster(Ye et al. 2021), transformer is used as the decoder, which achieves high accuracy, but brings huge computation cost.

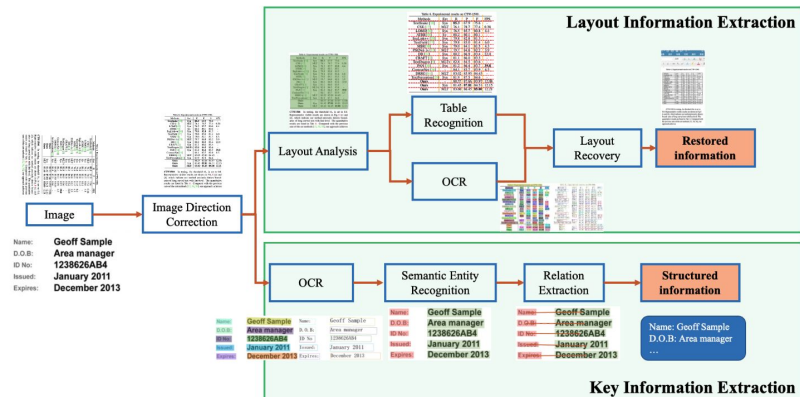


Figure 1: Framework of the proposed PP-StructureV2. It contains two subsystems: layout information extraction and key information extraction.

[Chenxia Li, et al., PP-structurev2: A stronger document analysis system, arxiv, 2022.](#)

iv:2210.05391v2 [cs.CV] 13 Oct 2022

ドキュメント解析とComputer Vision分野との繋がり

近年では、QwenなどのLVLMでDocument画像とBBoxも用いて学習されることもある。

Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution

Peng Wang*, Shuai Bai* Sinan Tan*, Shijie Wang*, Zhihao Fan*, Jinze Bai*
 Keqin Chen Xuejing Liu Jialin Wang Wenbin Ge Yang Fan Kai Dang Mengfei Du
 Xuancheng Ren Rui Men Dayiheng Liu Chang Zhou Jingren Zhou Junyang Lin¹
 Qwen Team Alibaba Group

Abstract

We present the Qwen2-VL Series, an advanced upgrade of the previous Qwen-VL models that redefines the conventional predetermined-resolution approach in visual processing. Qwen2-VL introduces the Native Dynamic Resolution mechanism, which enables the model to dynamically process images of varying resolutions into different numbers of visual tokens. This approach allows the model to generate more efficient and accurate visual representations, closely aligning with human perceptual processes. The model also integrates Multimodal Rotary Position Embedding (M-RoPE), facilitating the effective fusion of positional information across text, images, and videos. We employ a unified paradigm for processing both images and videos, enhancing the model's visual perception capabilities. To explore the potential of large multimodal models, Qwen2-VL investigates the scaling laws for large vision-language models (LVLMs). By scaling both the model size—with versions at 2B, 8B, and 72B parameters—and the amount of training data, the Qwen2-VL Series achieves highly competitive performance. Notably, the Qwen2-VL-72B model achieves results comparable to leading models such as GPT-4o and Claude3.5-Sonnet across various multimodal benchmarks, outperforming other generalist models. Code is available at <https://github.com/QwenLM/Qwen2-VL>.

1 Introduction

In the realm of artificial intelligence, Large Vision-Language Models (LVLMs) represent a significant leap forward, building upon the strong textual processing capabilities of traditional large language models. These advanced models now encompass the ability to interpret and analyze a broader spectrum of data, including images, audio, and video. This expansion of capabilities has transformed LVLMs into indispensable tools for tackling a variety of real-world challenges. Recognized for their unique capacity to condense extensive and intricate knowledge into functional representations, LVLMs are paving the way for more comprehensive cognitive systems. By integrating diverse data forms, LVLMs aim to more closely mimic the nuanced ways in which humans perceive and interact with their environment. This allows these models to provide a more accurate representation of how we engage with and perceive our environment

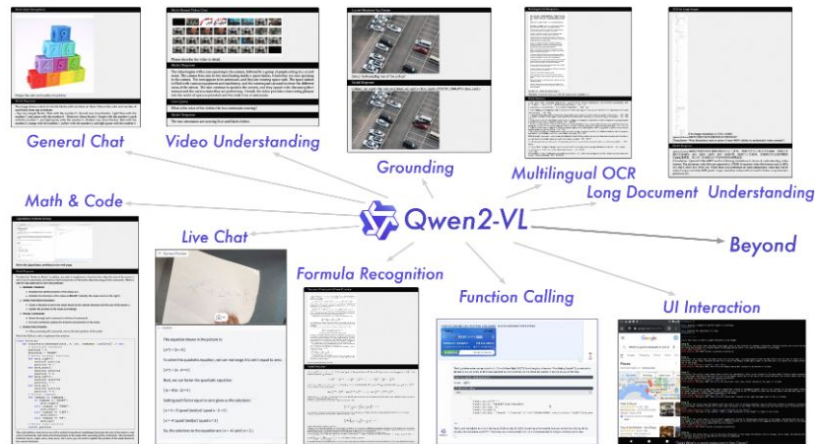


Figure 1: Qwen2-VL capabilities: Multilingual image text understanding, code/math reasoning, video analysis, live chat, agent potential, and more. See Appendix for details.

[Peng Wang, et al., Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution, arxiv, 2024](#)

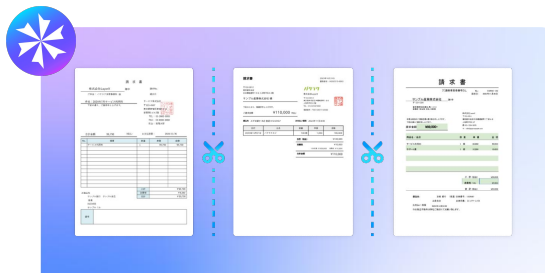
arXiv:2409.12191v2 [cs.CV] 3 Oct 2024

Bakuraku AI

日常業務の中で自然にAIを活用いただけるような体験を提供



写真を撮影し、まとめてアップロードするとAIがデータ入力



複数枚まとめてスキャンすると、PDFファイルをAIが自動分割



請求書をアップロードした瞬間に、AIが過去に学習した仕訳を入力



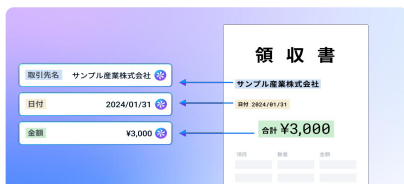
領収書をアップロードした瞬間にAIがカード明細情報と照合適切な明細と紐付け



あらゆる種類の書類をまとめてアップロードすると、AIが書類種別を判定し、ラベル付与

Bakuraku AIに関わる要素技術

- PDFや画像といった非構造化データからの項目抽出や分類
- 情報抽出されたデータや顧客の履歴データなどの構造化データを用いた推薦モデル



画像やPDFからの項目抽出タスク



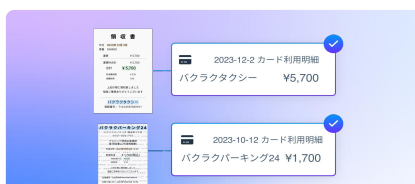
各社の運用に合わせて
項目抽出した値を推薦



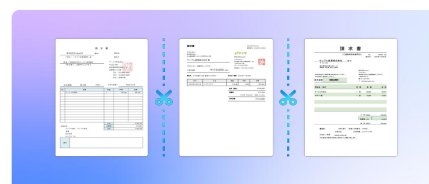
請求書に対して仕訳を推薦



画像やPDFの書類分類



領収書に対してカード明細を推薦



複数の請求書から抽出した情報から
尤もらしい分割点を予測

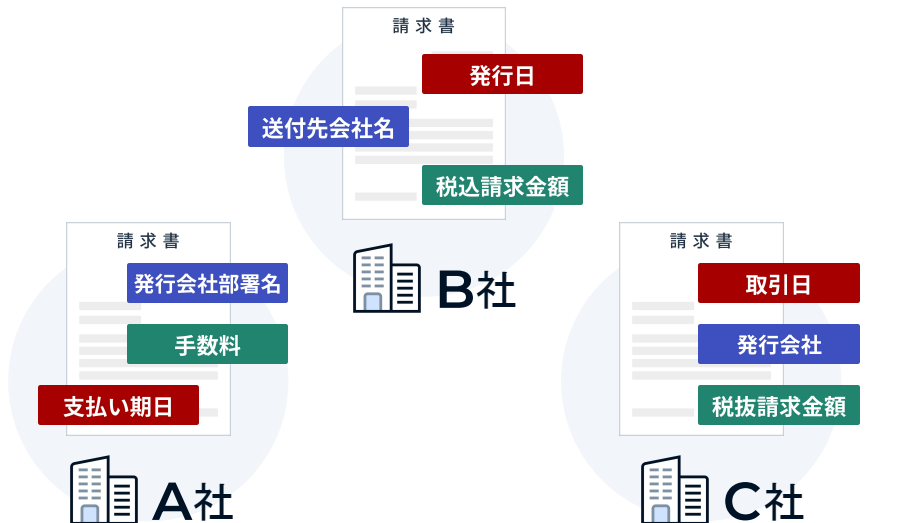
パーソナライズドAI-OCR *特許出願準備中

ユーザーに合わせて自動学習する次世代のAI-OCR

単純に情報抽出だけでなく、項目抽出した値をパーソナライズすることを実現

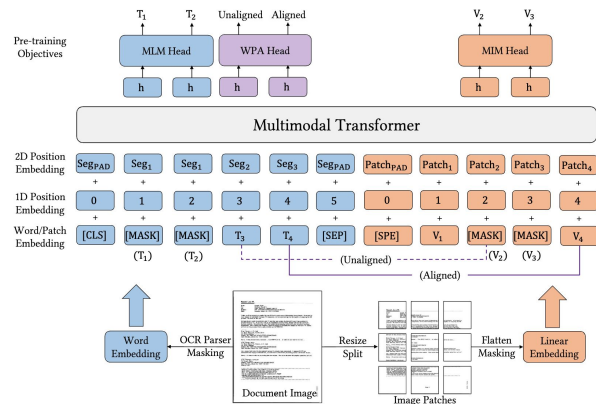
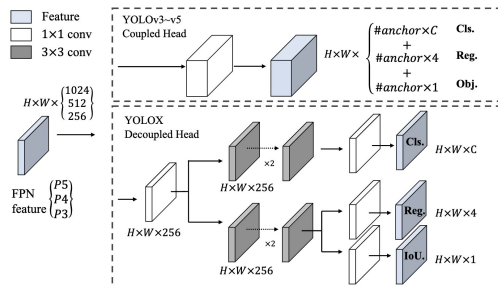
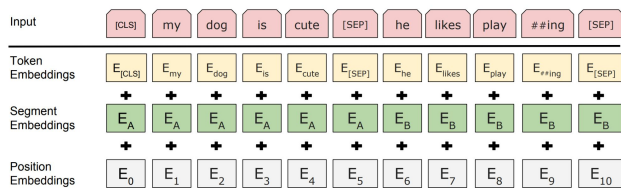
1 複数の値を同時に項目抽出

2 お客様の選択を学習していき、運用に最適化していく



パーソナライズドAI-OCR (項目抽出モデル)

項目抽出部分はBERT系 (NLPモデル)や、Object Detection (CVモデル)系、LayoutLM (マルチモーダルモデル)系などを複数検証



[Jacov Deblin, et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, ACL, 2019](#)

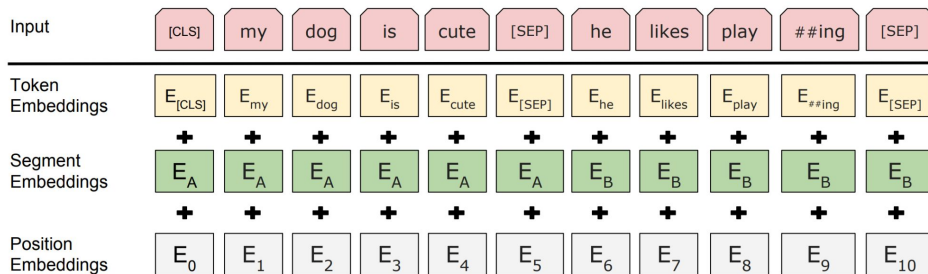
[Zheng Ge, et al., YOLOX: Exceeding YOLO Series in 2021, arxiv, 2021](#)

[Yupan Huang, et al., LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking, ACM, 2022](#)

パーソナライズドAI-OCR (項目抽出モデル)

RoBERTa

- モデル入力
 - OCRされた文書テキスト
- モデル出力
 - Tokenに対応するラベル (Token Classification)
- 選定理由
 - まずはシンプルな実装でベースラインを作るのが重要

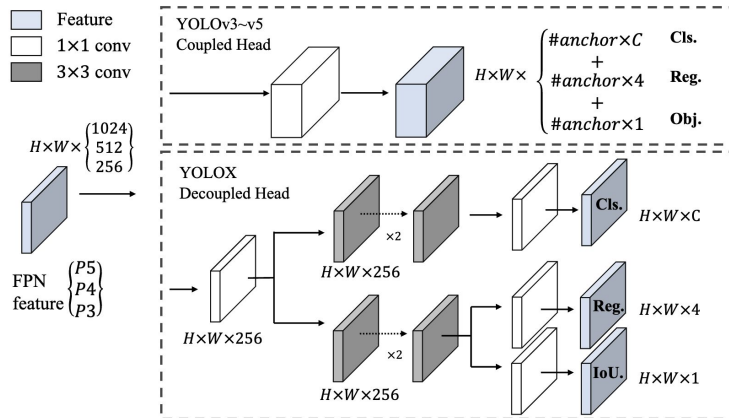


	20	21	22	23	24	25	26	27	28	29	30
token	。	合計	金額	¥	49	,	895	請求	書	様	御
label	others	others	others	others	paymentAmount	paymentAmount	paymentAmount	others	others	others	others

パーソナライズドAI-OCR (項目抽出モデル)

YOLOXなどのObject Detectionモデル

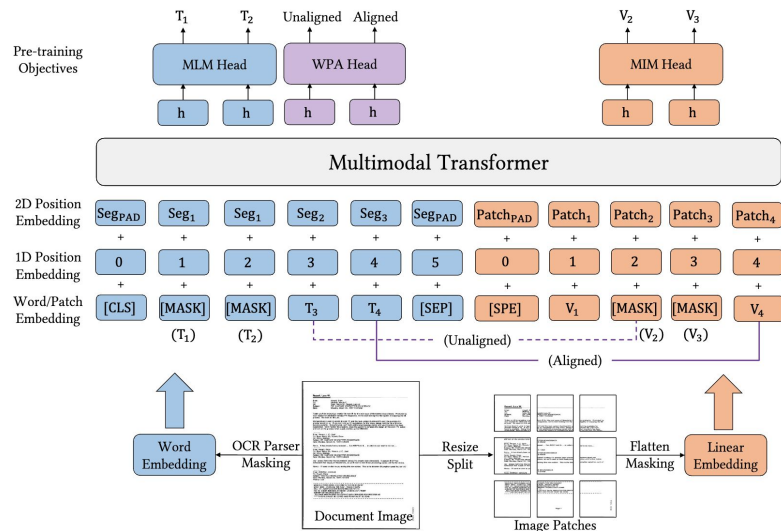
- モデル入力
 - PDFを画像化したもの
 - 携帯で撮った写真
- モデル出力
 - 欲しいラベルのBBBox
 - OCR結果と突合して出力する
- 選定理由
 - 画像情報がどれほど重要かを検証



パーソナライズドAI-OCR (項目抽出モデル)

LayoutLMv3

- モデル入力
 - OCRされた文書テキスト
 - PDFを画像化したもの
 - 携帯で撮った写真
- モデル出力
 - Tokenに対応するラベル (Token Classification)
- 選定理由
 - 当時のSoTAかつToken Classificationの枠組みで解ける



	20	21	22	23	24	25	26	27	28	29	30
token	。	合計	金額	¥	49	,	895	請求	書	様	御
label	others	others	others	others	paymentAmount	paymentAmount	paymentAmount	others	others	others	others

パーソナライズドAI-OCR (推薦モデル)

最終的にお客様が欲しいのは、「誰に対して」「いつ」「いくら支払った」といった情報。同じ書類であっても、お客様の運用によって変わることがあるため、過去の入力履歴を活用。

会社名が欲しい or 担当者名が欲しい

請求書

2024年06月30日
書類番号:SP202406-0008

株式会社LayerX
〒103-0012
東京都中央区日本橋蛸町
1丁目9-8
人形町PRR X 2F
電話:02-2645-71890
担当者:バクラク太郎
登録番号:190705071140088

株式会社バクラク
バクラク太郎 様

下記のとおり、御請求申し上げます。

ご請求金額 **¥218,000** (税込)

前回ご請求金額	入金額	繰越金額	今回請求金額(税抜)	今回消費税額	今回請求金額(税込)
218,000	218,000	0	200,000	18,000	218,000

10%対象金額	100,000	消費税額	10,000	税込金額	110,000
8%対象金額	100,000	消費税額	8,000	税込金額	108,000

税込金額が欲しい
OR
税抜金額が欲しい



過去の入力履歴



お客様が欲しい値を推薦

(余談) Why not LLM ?

LayerXでは、AIを用いた際の体験「**AI-UX**」を重視しており、間違えた際の**気づきやすさ**、**修正しやすさ**、**抽出速度**などからシンプルなモデルを使用。(どこからその値を抽出したか、などはまだ生成モデルより既存モデルの方が強い。)
自社ドメインの大量データで学習する際は、精度的にも既存モデルの方が高いことも多い。

また、パーソナライズ部分などもICLより既存の推薦技術の方がシンプルで扱いやすい。

The screenshot shows the '経費精算' (Expense Reconciliation) page in the Bakuraku application. The interface is divided into a left sidebar and a main content area. The sidebar contains navigation options: '申請・承認' (Application/Approval), '申請する' (Apply), 'ファイル' (Files), 'カード' (Cards), and '設定・その他' (Settings/Other). The main content area is titled '経費精算' and includes a form for entering expense details. At the top, there is a '申請名' (Application Name) field. Below it, a section labeled '明細' (Details) features a table with columns for '日付' (Date), '内訳' (Details), '金額(税込)' (Amount including tax), '支払先' (Payment destination), and '内容・メモ' (Content/Notes). The table contains two rows of data, each with a plus icon for adding more items. A blue button labeled 'AIが内訳を自動入力' (AI automatically enters details) is positioned above the table. At the bottom of the form, there are buttons for '下書きに保存' (Save as draft), '確認する' (Check), and '申請する' (Apply).

証憑マッチング・仕訳推薦

AI-OCRで抽出した取引先名や金額・日付などの情報を元に

- 「クレジットカードの明細」と「領収書ファイル」の紐付けや
- 「仕訳」と「請求書」の紐付けを自動化。

領収書

日付 2023年12月28日
 番号 555909

運賃 ¥ 5,700
 運賃料合計 ¥ 5,700
 合計 ¥ 5,700

内消費税 ¥ 518
 消費税率 10%

上記の明細に領収取しました
 間違ご申告ありがとうございます

バクラクタクシー
 登録番号：T1234567891011

2023-12-2 カード利用明細

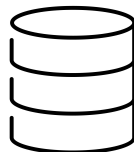
バクラクタクシー ¥5,700

仕訳入力

計上日 2018-09-29

借方		貸方	
勘定科目	金額	勘定科目	金額
外注費	99,790	000:未払金	99,790
000:管理本部	課税仕入 10%	000:管理本部	00:対象外

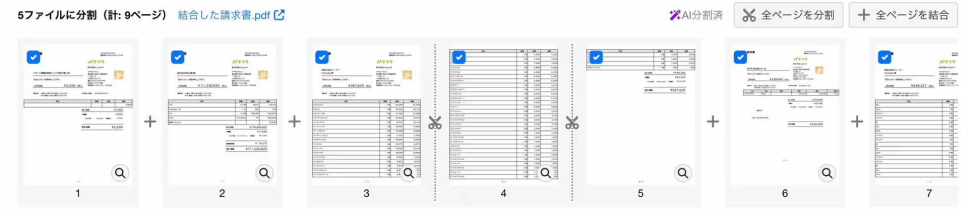
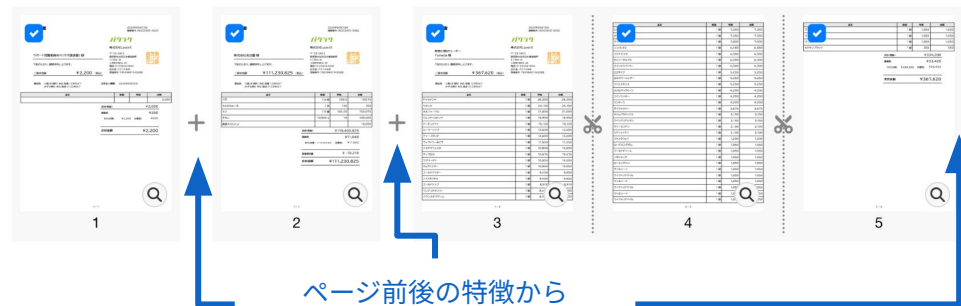
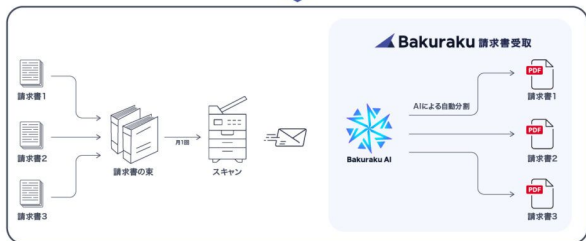
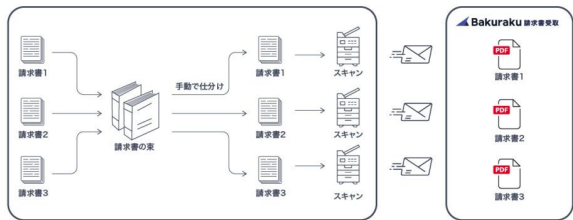
未処理のカード明細を
候補として推薦



過去の仕訳情報を
候補として推薦

PDF分割

請求書を取引ごとに管理するため、取引ごとにお客様が書類を分割してスキャンすることが必要だった
AI-OCRで抽出した情報を特徴に、分割点を予測。



実データにおける課題

パーソナライズドAI-OCRを例に

訓練ラベルの設計

実際のプロダクトでのユースケースを加味する必要があるため、**ドメインにDeepDiveする必要**がある

例) パーソナライズドAI-OCRの場合、以下のようなことを加味したラベル設計をしなければいけない。

- 領収書の場合、法律的に屋号でも運営会社名でも取引先として抽出して大丈夫。
- 新幹線の領収書は、購入日なのか乗車日を取りたいかがサービスによって変わる。

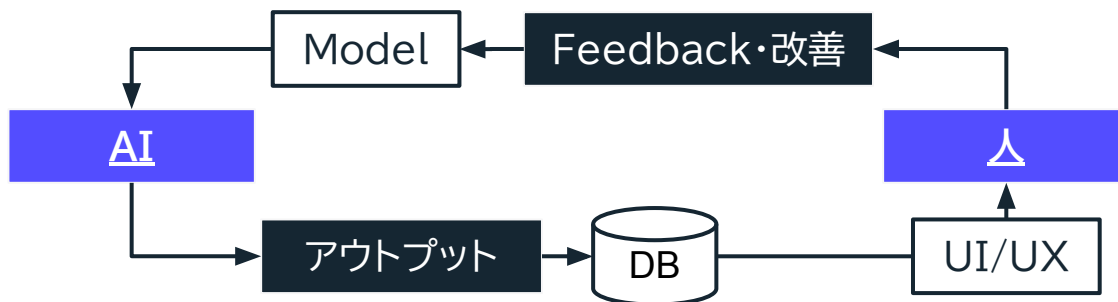
領 収 書 RECEIPT		バクラク鉄道株式会社	
お宛名		様	
金額計	¥ 8,900 (10%・税込)	内容	乗車券類の購入代金
購入日	2023年10月17日	乗車日	2023年10月17日
取扱カード会社	BAKU	クレジットカード番号	xxxxxxxxxxxx
列車名・券種 利用区間	名古屋 → 東京 FROM バクラク8909号		¥ 8,900
		T9010401140088	

日本橋 爆楽 人形町	
【領収書】	
東京都中央区 日本橋人形町0-0-0 TEL:00-0000-000	
2022/05/10 20:30:32 レジ:001 担当:00001 取引 No.00000000000000	
ご利用ありがとうございます	
生ビール	4点 ¥3,080
徳盛う	¥3,850 1点 ¥3,850
鳥籠茶	¥440 2点 ¥880
梅きゅう	¥330 1点 ¥330
鶏つくば	¥440 1点 ¥440
焼きおにぎり	¥440 2点 ¥880
小計	11点 ¥9,460
合計	¥9,460
(内消費税等	¥860)
(10% 標準税率	¥9,460)
(内消費税等	¥860)
上記正に領収いたしました	
領収票No. 70000000000000	
株式会社バクラク	

フィードバックサイクルの設計が重要

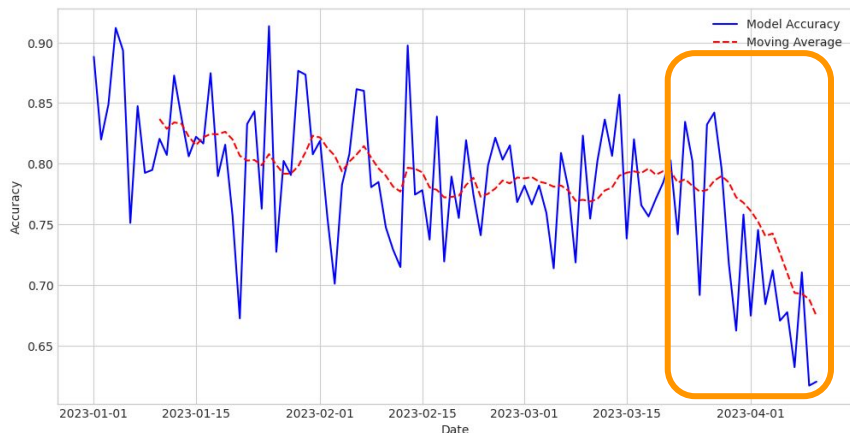
- お客様が増えていく中で、書類の多様性なども日々増えていくが、全ての書類に対してアノテーションできるわけではない。
 - 間違いやすい書類をどのように集めてアノテーションに回すか。
 - お客様のFBをどのように次の学習に活用するか。
- 修正してくれるお客様とそうでないお客様が存在する中、ラベルノイズの影響をどう減らすか。

これらをプロダクトをリリースする前に設計し、必要なログなどを仕込む必要がある。
機械学習モデルの開発は**デプロイして終わりではない**。

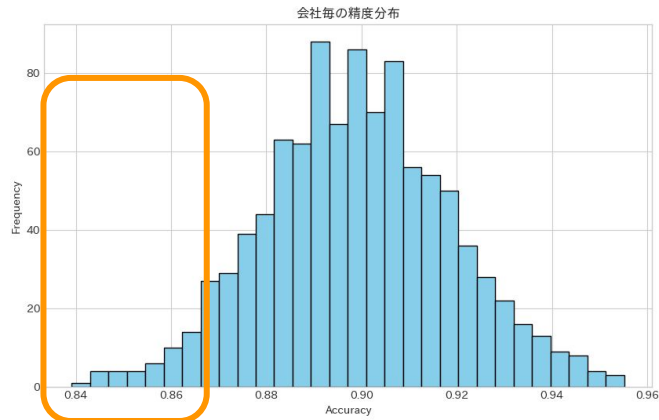


デプロイした後も常にモニタリングし、再設計し続ける

- ラベルの作り方に間違った仮定が入ってなかったか
- 想定していなかった運用のお客が増えていないか
- 実際にお客様に話を聞きに行くのも重要



急激な精度変化の検知



精度が比較的低いお客様で何が起きているのかの分析

実際に出てきているユースケース

- 経費精算の際に台紙に領収書を複数貼って提出する運用。
 - AI-OCRは一つの画像に複数の取引先名が入っていることを想定していないため、精度劣化。
- 一つの請求書の合計金額ではなく、明細の金額を用いて複数の仕訳に分割したい。
 - AI-OCRは合計金額だけを読み取っていたが、明細の金額も読み取る必要が出てきた。

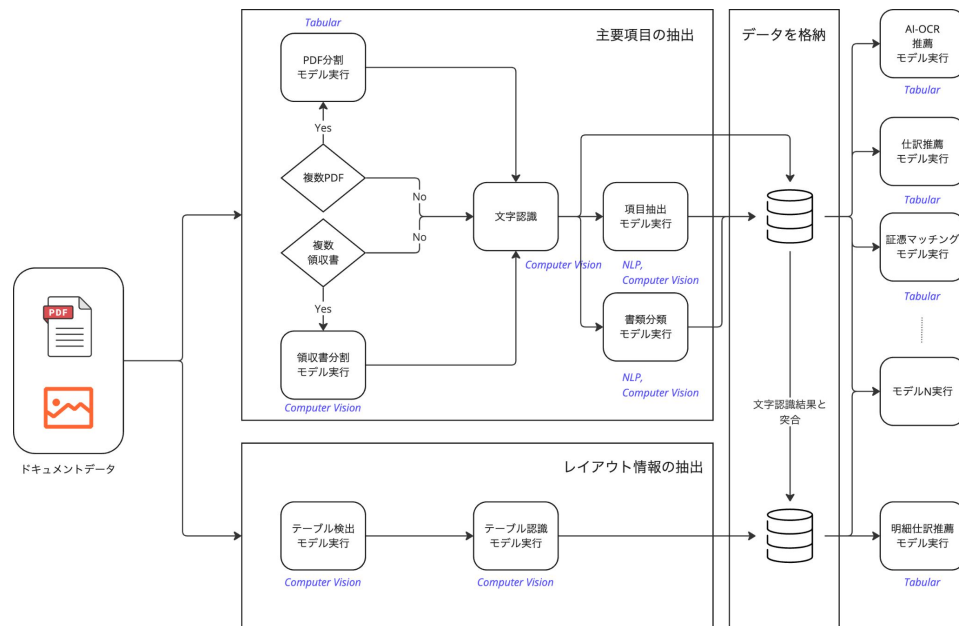


請求書				
株式会社御中		クルクル商店		
請求明細表				
ご利用年月日	問い合わせ番号	品名	請求金額	税率(%)
2024/04/30	73150452	トレットペーパー 新宿店	481,961 円	10
2024/05/31	73150454	トレットペーパー 名古屋店	139,308 円	10
2024/06/30	73150456	トレットペーパー 札幌店	243,745 円	10
2024/07/31	73150458	トレットペーパー 福岡店	73,780 円	10
2024/08/31	73150460	備品 大阪店様	210,215 円	8
2024/09/30	73160890	備品 上野店様	150,437 円	8
2024/10/31	73190221	備品 石川店様	128,366 円	8
2024/11/30	73109442	備品 広島店様	300,773 円	8
			小計	1,728,583
			消費税等(10%)	86,501
			消費税等(8%)	69,086
				1,884,170
コンゴトモロクオネガイシマス				

今後の展望

ドキュメント解析パイプラインの改善

実際に出てきているユースケースを元にデータ抽出パイプラインを再設計し、後段で様々なユースケースでの活用を行っていききたい。



まとめ

LayerXでは様々な技術を組み合わせてお客様の課題を解決しています

- ドキュメントという非構造化データを起点に
CVやNLP、推薦などの複合的な技術を用いて、お客様の課題を解決しています。
- モデルを作成する技術はもちろんですが、継続的に改善をするために、
ログの設計やFBサイクルの設計などのMLOps的要素をモデル作成時から考える必要がある。
- 今後、さらに多様なユースケースを解決していくために、
様々な技術を用いてお客様の体験を「バクラク」にしていきます。

[宣伝] LayerX Machine Learning勉強会

LayerXでは毎週機械学習関連の勉強会を行っており、内容を全て公開しています。話題になった論文などの紹介を行っているので、ぜひご覧ください！

メインTOPIC

[論文] The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery

Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, David Ha

概要

Sakana AIが提案する論文作成のような科学的発見プロセスを完全に自動化する包括的なフレームワーク

- 大規模言語モデル(LLM)を活用し、研究のアイデア生成から論文執筆・査読までを自動化
- 1本の論文生成にかかるコストはわずか\$15程度と極めて効率的
- 人間の査読者と同等レベルの自動査読システムを実現
- 機械学習の3つの異なる分野（拡散モデル、言語モデル、学習ダイナミクス）での有効性を実証

Introduction

