

Structured Extraction from Business Process Diagrams Using Vision-Language Models

Pritam Deka
p.deka@qub.ac.uk
Queen's University Belfast
Belfast, UK

Barry Devereux
B.Devereux@qub.ac.uk
Queen's University Belfast
Belfast, UK

Abstract

Business Process Model and Notation (BPMN) is a widely adopted standard for representing complex business workflows. While BPMN diagrams are often exchanged as visual images, existing methods primarily rely on XML representations for computational analysis. In this work, we present a pipeline that leverages Vision-Language Models (VLMs) to extract structured JSON representations of BPMN diagrams directly from images, without requiring source model files or textual annotations. We also incorporate optical character recognition (OCR) for textual enrichment and evaluate the generated element lists against ground truth data derived from the source XML files. Our approach enables robust component extraction in scenarios where original source files are unavailable. We benchmark multiple VLMs and observe performance improvements in several models when OCR is used for text enrichment. In addition, we conducted extensive statistical analyses of OCR-based enrichment methods and prompt ablation studies, providing a clearer understanding of their impact on model performance.

CCS Concepts

• Computing methodologies → Information extraction.

Keywords

Artificial Intelligence, Vision-Language Models, Diagram Understanding, Prompt Engineering, Structured Information Extraction

ACM Reference Format:

Pritam Deka and Barry Devereux. 2026. Structured Extraction from Business Process Diagrams Using Vision-Language Models. In *Proceedings of The 41st ACM/SIGAPP Symposium on Applied Computing (SAC'26)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Business Process Model and Notation (BPMN) is a widely adopted standard for modeling and documenting business workflows [3]. It provides a visual formal language for processes through standardized visual elements—such as tasks (individual units of work), events (triggers or outcomes like start or end points), gateways (decision or branching logic), and sequence flows (arrows indicating control

flow between elements)—that support both human interpretability and machine-readable semantics [26, 3, 19]. BPMN diagrams are extensively used in enterprise environments, particularly for process documentation, compliance tracking, and training [25]. However, these diagrams are often embedded in static formats such as PDF handbooks, training slides, or scanned documentation where the original XML-based BPMN source is not available. This lack of access to the structured source format presents a significant challenge for downstream applications that rely on semantic process data, such as simulation [29], automated migration [37], conformance checking [52], or integration with workflow engines.

In this work, we present a vision-based approach for structured BPMN understanding using state-of-the-art vision-language models (VLMs) [2, 1, 22, 49, 48]. Our method treats BPMN diagrams as visual inputs and extracts structured JSON representations using a prompt that guides the model to detect and label each process element. The prompt instructs the model to identify and classify core components based on their visual layout. To improve coverage, especially in cases where labels are missing, low-contrast, or partially obscured, we introduce optional OCR-based enrichment using lightweight text extractors. This hybrid pipeline enables accurate recovery of BPMN semantics from images alone, without relying on the original XML source.

Our contributions are as follows:

- We propose a prompt-based pipeline for extracting structured representations from BPMN diagrams using VLM.
- We release a BPMN dataset with diagram–XML pairs for further research.
- We benchmark a range of VLMs and analyze OCR impact on extraction quality.
- We perform detailed ablation experiments and analysis to study the effect of different prompt types.

2 Related Work

Business process modeling, especially the construction and interpretation of BPMN diagrams, has attracted growing research interest with two main methodological approaches.

2.1 LLM-Driven BPMN Generation and Process Modeling

Much of the recent progress in process automation comes from advances in LLMs. Early research relied on symbolic NLP and semantic role labeling to turn natural language into BPMN models [13], but new systems use LLMs like GPT to improve accuracy and user interaction [44, 40, 31]. Tools such as Nala2BPMN [40] and [31] combine these models with visualization toolkits, making process modeling more accessible.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SAC'26, Thessaloniki, Greece

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-X-XXXX-XXXX-X/26/03
<https://doi.org/XXXXXXX.XXXXXXX>

Several modern frameworks take this further. ProMoAI [28] lets users iteratively refine BPMN models from text and export them in multiple formats. The MAO Framework [32] uses teams of LLM agents working in dialogue to collaboratively improve process models, while Opus [11] integrates LLMs with knowledge graphs to optimize business process outsourcing workflows. Although their goals differ, with ProMoAI focusing on user feedback, MAO on multi-agent teamwork, and Opus on domain knowledge, they all demonstrate how LLMs can automate and enhance BPMN generation.

At a finer level, studies like [6] demonstrate that GPT-3 can extract process elements and relationships from text using prompt engineering alone, and [31] shows that combining LLMs with NLP pre-processing tools like spaCy and BERT [10] improves extraction quality. These works highlight that, with the right prompts and minimal fine-tuning, LLMs can provide structured BPMN outputs without requiring specialized training.

Beyond just generating diagrams, LLMs have been applied to explanation and interactive analysis of process models. For example, [39] proposes combining case-based reasoning with LLMs for grounded explanations, and [18] and [27] extend LLM use to tasks such as anomaly detection and process refinement. Some research takes a data-driven angle: [9] shows that process models can also be discovered from execution event logs, not just text, demonstrating the versatility of LLM-based approaches.

2.2 VLMs for Diagram Understanding

Alongside these text-driven methods, VLMs are rapidly advancing the ability to understand and extract information from visual process diagrams. VLMs have already shown success in visual question answering, image captioning, and layout understanding [22, 17, 48, 14]. Recently, researchers have started to apply these models to more structured visual reasoning, including tasks with symbolic diagrams [33, 57, 7, 41].

The development of specialized datasets and benchmarks has been instrumental in advancing diagram understanding. In the domain of chart reasoning, resources such as ChartQA [34], ChartBench [56], NovaChart [21], Chartx [55] and ChartInstruct [35] have provided large-scale annotated data for training and evaluating models on quantitative and logical reasoning tasks. On top of these, models like ChartGPT [51], Chartllama [20], Chartgemma [36] and Chartvllm [36] are adapted for chart-specific tasks. For process diagrams, benchmarks such as FlowchartQA [47] and FlowVQA [45] provide diverse synthetic flowcharts paired with questions to evaluate the visual reasoning abilities of VLMs. In contrast, BPMN-Redrawer [4] focuses on reconstructing executable BPMN XML representations directly from visual diagrams, a task that demands precise recovery of both entities and their relationships.

Building on this line of work, our research explores the more challenging and practical scenario of extracting structured BPMN representations directly from images, without relying on any accompanying text, underlying XML, or manual heuristics. We propose an approach that guides a VLM using a BPMN-specific prompt

focused on visual and semantic characteristics, enabling generalization across various layouts and styles. This allows for schema-constrained recovery of all major BPMN elements (tasks, events, gateways, flows, pools, lanes) from static diagrams alone.

3 Dataset Construction

In this section, we describe how the dataset was constructed. Our goal was to create a clean, consistent, and visually aligned dataset that could serve as a strong benchmark for structured diagram understanding.

3.1 Source Generation

We built our dataset using BPMN 2.0 XML files collected from the publicly available bpmn-for-research¹ repository [15], which provides over 3,700 BPMN diagrams created during Camunda's training sessions. From this corpus, we selected only the English-language diagrams to maintain consistency in text content and labels. After filtering, we renamed all BPMN files using a standardized format (e.g., process_001.bpmn) for consistent pairing with images. Each diagram was rendered to a high-resolution PNG using a custom script that leverages the bpmn-js² library from Camunda's bpmn.io project [16]. The final dataset consists of 202 diagram pairs, each containing a BPMN XML file and a cropped PNG image. All filenames are matched to allow straightforward alignment between formats.

3.2 Statistics and Format

We split the 202 BPMN diagrams into training, development, and test sets using a **50:25:25 ratio**, resulting in **101 training**, **50 development**, and **51 test** diagrams. Unlike typical 80:20 or 70:30 splits, this choice was driven by our focus on **prompt-based evaluation** rather than model training. A larger test set allows for better analysis of model behavior across a diverse set of BPMN structures. Each sample includes a **BPMN diagram image (.png)** and its **corresponding XML file (.bpmn)**.

The dataset covers a diverse range of BPMN diagram types, including **single-pool processes**, **multi-lane and multi-pool layouts**, **tasks and subprocesses**, **gateways** (exclusive, inclusive, parallel), **events** (start, end, intermediate), and **sequence flows** with or without labels. This variety ensures the test set captures both simple and complex scenarios. The complete dataset with the splits, folder structures, and preprocessing scripts is available via <https://github.com/pritamdeka/BPMN-VLM>.

4 Methodology

The core methodology of our approach is centered around a prompt-based structured information extraction from BPMN diagrams using VLMs. The prompt includes a comprehensive list of extraction requirements and a sample JSON object to guide the model's output structure, but does not provide any input-output example pairs. This approach ensures the VLM produces consistent, schema-compliant outputs across diverse BPMN diagram styles. We have used the same prompt across all the models.

¹<https://github.com/camunda/bpmn-for-research>

²<https://github.com/bpmn-io/bpmn-js>

Problem Definition. Let a BPMN diagram image be denoted as I . The objective is to extract a structured JSON representation of BPMN elements $E = \{e_1, e_2, \dots, e_n\}$ where each element e_i has a type, textual label, and attributes. We define the extraction as a mapping:

$$f_\theta : (I, P, T) \mapsto R \quad (1)$$

where f_θ is a VLM parameterized by θ , P is the input prompt, T is optional OCR-enriched text, and R is the raw textual response. **JSON Parsing.** The response R is mapped to a structured representation $J(R)$ via a parsing operator:

$$J(R) = \begin{cases} \text{Valid JSON object} & \text{if parsing succeeds} \\ \emptyset & \text{otherwise} \end{cases} \quad (2)$$

We have explored two primary directions in our methodology which are explained in the subsequent subsections.

4.1 VLM-Only Prompted Extraction

In this setting, we input the BPMN diagram directly into a VLM accompanied by our carefully designed prompt that describes the visual symbols and semantics of BPMN elements (e.g., rectangles for tasks, diamonds for gateways). The model is expected to utilize its joint visual and textual understanding to process the image and output a structured JSON representation that mirrors the underlying process structure. This approach proves effective for simpler diagrams, where the layout is clean, components are spatially distinct, and most elements contain clearly visible text labels. Formally, the pipeline from Eq 1 and Eq 2 reduces to:

$$R = f_\theta(I, P), \quad J(R) = \{e_1, e_2, \dots, e_n\} \quad (3)$$

However, as diagram complexity increases, the model's extraction accuracy declines noticeably. VLMs often struggle with: (1) **unlabeled or ambiguous elements**, especially gateways and events lacking names; (2) **dense layouts**, where overlapping elements in multi-pool or multi-lane diagrams cause confusion or omissions; (3) **low-contrast or edge-aligned text**, particularly on sequence flows; and (4) **resolution constraints**, as many VLMs rescale input images to fixed sizes, leading to the loss of fine-grained visual details in large or wide diagrams (see Section 6.5 and Table 5 for a detailed empirical quantification of these errors).

Figure 1 shows why VLM-only extraction may be error-prone. The model must identify: (1) the **pentagon** inside the gateway (an **EventBasedGateway**); (2) the difference between a solid **SequenceFlow** and a dashed **MessageFlow**; and (3) the envelope **IntermediateCatchEvent**. These small markers and line styles often vanish when the image is downsampled, which may lead to the misclassifications reported in Section 6.5.

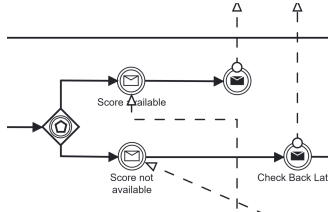


Figure 1: Cropped BPMN region with fine-grained cues

Algorithm 1 Proposed BPMN Extraction Method

```

1: Input: Image  $I$ , Prompt  $P$ , OCR flag  $enable\_ocr$ , OCR method  $m$ 
2: Output: Structured JSON  $\hat{J}$ 
3: Encode  $I$  as base64 string  $E(I)$ 
4: if  $enable\_ocr$  then
5:   Extract text tokens  $T \leftarrow O_m(I)$ 
6: else
7:    $T \leftarrow \emptyset$ 
8: end if
9: Generate raw response  $R \leftarrow f_\theta(I, P, T)$ 
10: Parse JSON  $J \leftarrow J(R)$ 
11: if  $J \neq \emptyset$  then
12:   Enrich missing labels  $\hat{J} \leftarrow E'(J, T)$ 
13: else
14:   Save raw output  $\hat{J} \leftarrow R$ 
15: end if
16: return  $\hat{J}$ 

```

4.2 OCR-Enriched Hybrid Extraction

To address the limitations of the VLM-only extraction pipeline, we investigate whether OCR-based enrichment can serve as an effective fallback mechanism. Our hypothesis is that OCR may help recover textual labels that VLMs fail to detect, particularly for elements like sequence flows, gateways, or events that often lack clear or legible annotations. Our methodology comprises two stages: the VLM first generates a structured JSON from the BPMN diagram, while OCR tools such as Pix2Struct [30], RapidOCR [50], or Tesseract [46] extract text from the same image in parallel. During post-processing, an enrichment operation uses OCR outputs to fill missing component labels in the VLM's output, improving entity completeness for visually ambiguous or illegible elements without altering the original structure. Formally, from Eq 1,

$$R = f_\theta(I, P, T) \quad T = O_m(I) \quad (4)$$

where $m \in \{\text{Pix2Struct, Tesseract, RapidOCR}\}$

Parsed results $J(R)$ are passed to an enrichment operator E' , which fills missing labels using OCR tokens:

$$\hat{J} = E'(J(R), T). \quad (5)$$

where, for any element $e \in J$ such that $\text{name}(e) = \emptyset$, we assign $\text{name}(e) \leftarrow \text{match}(T)$.

In both settings, we apply post-inference JSON validation and normalization to ensure consistency across models and support accurate alignment with ground truth XML during evaluation. This dual-branch architecture (VLM-only vs. VLM+OCR) allows us to analyze the trade-offs in model performance and identify cases where OCR integration leads to measurable improvements. Algorithm 1 shows the pseudocode for our proposed method.

4.3 Prompt Design

We adopt a zero-shot, schema-constrained prompting strategy [41, 38, 24, 43], providing only detailed instructions and an output example in the prompt, but no input–output pairs. The prompt was iteratively refined to maximize extraction accuracy and consistency.

The prompt explicitly describes the visual characteristics and semantics of all BPMN components, specifies a strict output schema, and defines extraction rules for bounding boxes and label assignment. Without providing explicit examples, the prompt relies on detailed instructions and a template JSON structure to enforce consistent, machine-readable outputs across models. The full prompt is available in the project repository³.

5 Experiments

We design our experiments to evaluate how well VLMs can extract structured representations from BPMN diagrams in image form. Models are tested in both vision-only and OCR-enhanced modes, across name, type and relation extraction tasks.

5.1 Models Compared

We benchmark a diverse set of VLMs, including both proprietary APIs and open-source instruction-tuned models hosted on Hugging Face⁴ [54]. Our selection spans a wide range of architectures and capabilities: **GPT-4.1 (standard, mini, nano)**⁵ and **GPT-4o (standard, mini)** [22] from OpenAI; **Qwen2.5-VL models (7B and 3B)** [49, 5] from Alibaba group; **Pixtral-12B** and **Pixtral-Large** [1] and **Mistral-Small-3.1**⁶ from MistralAI team; **Aya-Vision-8B**⁷ from CohereForAI; **Gemma3-4B and 12B** [48] from Google; and **LLaMA-3.2-11B-Vision-Instruct** [2] from Meta. This range allows a comprehensive analysis of multimodal diagram understanding across model scales and training paradigms.

Each model is evaluated in two configurations:

- **Vision-only:** The model predicts component names and types directly from the BPMN image.
- **OCR-enhanced:** OCR text is appended only when the model fails to detect key components such as events, gateways or sequence flows.

This allows us to assess both the raw visual understanding of the model and its ability to incorporate OCR-enriched information for better name and type extraction.

In addition, we include **BPMN Redrawer** [4] and **Sketch2-Process** [42] as non-VLM baselines. **BPMN Redrawer**⁸ converts BPMN images into executable diagrams by combining convolutional neural networks with OCR-based text recognition, while **Sketch2Process**⁹ translates hand-drawn sketches into BPMN models using deep learning-based symbol recognition.

We evaluate performance using Precision, Recall, and F1 Score, applying multiple levels of matching criteria to capture both exactness and tolerance to partial or noisy predictions. Gold-standard CSV files are constructed by parsing the corresponding .bpmn XML files for each diagram. These files provide the name, type, and relations of each BPMN element and serve as the reference for evaluation. In the name-only setting, predicted element names are compared with gold-standard names regardless of type. The type-only setting focuses exclusively on element categories, abstracting

away from labels. The name+type setting requires both fields to be present and identical, with incomplete entries excluded to ensure fairness. Finally, the relations setting evaluates whether both the type and connectivity of elements align with the gold standard.

6 Results and Analysis

We evaluate model performance under four evaluation regimes: *name-only*, *name+type*, *relations+type*, and *type-only*. Each regime is measured in both strict and relaxed variants where applicable, comparing VLM-only inference with OCR-enriched pipelines. The non-VLM baseline is only reported without enrichment. In the relaxed settings, type labels may partially match (e.g., exclusivegateway \approx gateway). Model outputs are recorded as either .json (valid structured output) or .txt (unstructured output requiring post-processing). In particular, LLaMA-3.2-11B often produces .txt responses, requiring additional parsing before evaluation.

Formally, let G be the gold-standard set of elements and M the model predictions. Each element has a name n , a type t , or a relation $r = (s, d, t)$. Matching is defined as follows:

$$M_{\text{name}}(n_g, n_m) = \begin{cases} 1 & \text{if } n_g = n_m, \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

$$M_{\text{type}}(t_g, t_m) = \begin{cases} 1 & \text{if } t_g = t_m \quad (\text{strict}), \\ 1 & \text{if } \text{norm}(t_g) = \text{norm}(t_m) \quad (\text{relaxed}), \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

$$M_{\text{name-type}}((n_g, t_g), (n_m, t_m)) = \begin{cases} 1 & \text{if } n_g = n_m \wedge M_{\text{type}}^*(t_g, t_m) = 1, \\ 0 & \text{otherwise,} \end{cases} \quad (8)$$

where M_{type}^* denotes either strict or relaxed type matching.

$$M_{\text{rel-type}}((s_g, d_g, t_g), (s_m, d_m, t_m)) = \begin{cases} 1 & \text{if } (s_g, d_g) = (s_m, d_m) \\ & \wedge M_{\text{type}}^*(t_g, t_m) = 1, \\ 0 & \text{otherwise,} \end{cases} \quad (9)$$

with M_{type}^* again being strict or relaxed.

True/false sets under regime X are:

$$\begin{aligned} TP_X &= \{g \in G \mid \exists m \in M : M_X(g, m) = 1\}, \\ FP_X &= \{m \in M \mid \nexists g \in G : M_X(g, m) = 1\}, \\ FN_X &= \{g \in G \mid \nexists m \in M : M_X(g, m) = 1\}. \end{aligned} \quad (10)$$

These definitions provide the basis for computing Precision, Recall, and F1 scores under each evaluation regime.

6.1 Name-Only Evaluation (Table 1)

Focusing solely on element names, this evaluation highlights text grounding capabilities.

Top-tier models such as GPT-4.1, GPT-4o, and Mistral-Small-3.1 achieve the highest F1 scores (>0.70). They perform robustly without OCR, with enrichment providing minimal or even negative gains due to noise. **Mid-tier models**, including Qwen2.5VL-7B, Gemma-12B, and Pixtral-Large, benefit more from OCR. Tesseract and Pix2Struct improve recall, helping these models recover missing labels at a modest precision cost. **Low-tier models**, such

³<https://github.com/pritamdeka/BPMN-VLM/blob/main/prompts/prompt.txt>

⁴<https://huggingface.co/>

⁵<https://openai.com/index/gpt-4-1/>

⁶<https://mistral.ai/news/mistral-small-3-1>

⁷<https://cohere.com/blog/aya-vision>

⁸<https://github.com/PROSLab/BPMN-redrawer>

⁹<https://github.com/dwslab/hdBPMN>

Table 1: Strict evaluation results (names only) with precision, recall, and F1 across VLM-only and OCR-enriched pipelines (best results in bold).

| Model | Only VLM | | | VLM + Pix2Struct | | | VLM + RapidOCR | | | VLM + Tesseract | | |
|-------------------------|--------------|--------------|--------------|------------------|--------------|--------------|----------------|--------------|--------------|-----------------|--------------|--------------|
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| Aya-vision-8B | 0.201 | 0.140 | 0.151 | 0.175 | 0.128 | 0.138 | 0.213 | 0.133 | 0.154 | 0.190 | 0.129 | 0.142 |
| Gemma-12B | 0.798 | 0.560 | 0.645 | 0.748 | 0.560 | 0.628 | 0.725 | 0.592 | 0.640 | 0.725 | 0.592 | 0.640 |
| Gemma-4B | 0.642 | 0.586 | 0.602 | 0.635 | 0.586 | 0.599 | 0.617 | 0.576 | 0.585 | 0.617 | 0.576 | 0.585 |
| GPT-4.1 | 0.829 | 0.794 | 0.807 | 0.783 | 0.810 | 0.792 | 0.781 | 0.802 | 0.787 | 0.777 | 0.786 | 0.775 |
| GPT-4.1-mini | 0.823 | 0.819 | 0.817 | 0.772 | 0.816 | 0.788 | 0.766 | 0.814 | 0.784 | 0.775 | 0.816 | 0.790 |
| GPT-4.1-nano | 0.712 | 0.535 | 0.598 | 0.661 | 0.529 | 0.576 | 0.681 | 0.522 | 0.579 | 0.645 | 0.504 | 0.553 |
| GPT-4o | 0.812 | 0.775 | 0.790 | 0.767 | 0.781 | 0.770 | 0.769 | 0.780 | 0.771 | 0.773 | 0.778 | 0.772 |
| GPT-4o-mini | 0.782 | 0.629 | 0.691 | 0.748 | 0.633 | 0.678 | 0.748 | 0.633 | 0.678 | 0.739 | 0.642 | 0.681 |
| LLaMA-3.2-11B | 0.373 | 0.334 | 0.344 | 0.375 | 0.302 | 0.322 | 0.394 | 0.327 | 0.343 | 0.344 | 0.309 | 0.318 |
| Mistral-small-3.1 | 0.830 | 0.781 | 0.802 | 0.776 | 0.778 | 0.772 | 0.773 | 0.788 | 0.776 | 0.787 | 0.774 | 0.774 |
| Pixtral-12B | 0.352 | 0.284 | 0.301 | 0.342 | 0.312 | 0.316 | 0.321 | 0.292 | 0.297 | 0.354 | 0.312 | 0.323 |
| Pixtral-large | 0.712 | 0.517 | 0.589 | 0.706 | 0.517 | 0.583 | 0.716 | 0.515 | 0.589 | 0.715 | 0.516 | 0.587 |
| Qwen-3B | 0.538 | 0.486 | 0.486 | 0.550 | 0.469 | 0.481 | 0.568 | 0.408 | 0.458 | 0.547 | 0.406 | 0.451 |
| Qwen-7B | 0.774 | 0.727 | 0.744 | 0.769 | 0.744 | 0.749 | 0.735 | 0.738 | 0.729 | 0.736 | 0.738 | 0.729 |
| Non-VLM Baseline | | | | | | | | | | | | |
| BPMN Redrawer | 0.531 | 0.407 | 0.458 | – | – | – | – | – | – | – | – | – |
| Sketch2Process | 0.787 | 0.630 | 0.693 | – | – | – | – | – | – | – | – | – |

Table 2: Strict vs. relaxed evaluation results (names + types) with precision, recall, and F1 across VLM-only and OCR-enriched pipelines (best results in bold).

| Model | Only VLM | | | | | | VLM + Pix2Struct | | | | | | VLM + RapidOCR | | | | | | VLM + Tesseract | | | | | |
|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|------------------|--------------|--------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|--------------|--------------|-----------------|-------|--------------|--------------|-------|--------------|
| | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | |
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| Aya-vision-8B | 0.130 | 0.104 | 0.106 | 0.138 | 0.107 | 0.111 | 0.125 | 0.099 | 0.103 | 0.129 | 0.100 | 0.104 | 0.155 | 0.108 | 0.120 | 0.159 | 0.109 | 0.123 | 0.131 | 0.103 | 0.108 | 0.134 | 0.104 | 0.110 |
| Gemma-12B | 0.502 | 0.377 | 0.422 | 0.512 | 0.383 | 0.430 | 0.435 | 0.377 | 0.396 | 0.443 | 0.383 | 0.403 | 0.515 | 0.465 | 0.481 | 0.524 | 0.472 | 0.488 | 0.515 | 0.465 | 0.481 | 0.524 | 0.472 | 0.488 |
| Gemma-4B | 0.425 | 0.449 | 0.428 | 0.440 | 0.463 | 0.443 | 0.420 | 0.449 | 0.425 | 0.435 | 0.463 | 0.440 | 0.429 | 0.465 | 0.436 | 0.439 | 0.475 | 0.446 | 0.429 | 0.465 | 0.436 | 0.439 | 0.475 | 0.446 |
| GPT-4.1 | 0.718 | 0.702 | 0.707 | 0.733 | 0.716 | 0.721 | 0.612 | 0.709 | 0.652 | 0.629 | 0.727 | 0.670 | 0.604 | 0.694 | 0.642 | 0.624 | 0.716 | 0.663 | 0.580 | 0.674 | 0.619 | 0.593 | 0.688 | 0.633 |
| GPT-4.1-mini | 0.706 | 0.729 | 0.715 | 0.720 | 0.743 | 0.729 | 0.608 | 0.722 | 0.657 | 0.620 | 0.736 | 0.669 | 0.607 | 0.722 | 0.655 | 0.618 | 0.735 | 0.667 | 0.618 | 0.726 | 0.664 | 0.633 | 0.742 | 0.679 |
| GPT-4.1-nano | 0.638 | 0.492 | 0.544 | 0.647 | 0.498 | 0.551 | 0.556 | 0.482 | 0.506 | 0.562 | 0.488 | 0.511 | 0.541 | 0.458 | 0.486 | 0.552 | 0.468 | 0.497 | 0.537 | 0.456 | 0.481 | 0.545 | 0.463 | 0.489 |
| GPT-4o | 0.699 | 0.676 | 0.685 | 0.719 | 0.695 | 0.704 | 0.643 | 0.695 | 0.663 | 0.659 | 0.712 | 0.679 | 0.628 | 0.688 | 0.652 | 0.646 | 0.705 | 0.670 | 0.653 | 0.698 | 0.670 | 0.663 | 0.707 | 0.679 |
| GPT-4o-mini | 0.669 | 0.543 | 0.594 | 0.683 | 0.555 | 0.606 | 0.665 | 0.583 | 0.613 | 0.677 | 0.593 | 0.624 | 0.665 | 0.583 | 0.613 | 0.677 | 0.593 | 0.624 | 0.631 | 0.570 | 0.591 | 0.643 | 0.579 | 0.602 |
| LLaMA-3.2-11B | 0.280 | 0.259 | 0.263 | 0.282 | 0.261 | 0.265 | 0.303 | 0.257 | 0.269 | 0.304 | 0.258 | 0.270 | 0.282 | 0.250 | 0.256 | 0.286 | 0.253 | 0.259 | 0.275 | 0.260 | 0.261 | 0.275 | 0.260 | 0.261 |
| Mistral-small-3.1 | 0.696 | 0.690 | 0.691 | 0.712 | 0.705 | 0.706 | 0.604 | 0.669 | 0.631 | 0.624 | 0.689 | 0.651 | 0.613 | 0.689 | 0.645 | 0.633 | 0.710 | 0.666 | 0.614 | 0.681 | 0.642 | 0.635 | 0.702 | 0.662 |
| Pixtral-12B | 0.274 | 0.231 | 0.238 | 0.282 | 0.238 | 0.246 | 0.238 | 0.238 | 0.230 | 0.241 | 0.242 | 0.233 | 0.239 | 0.245 | 0.235 | 0.245 | 0.252 | 0.241 | 0.252 | 0.254 | 0.247 | 0.253 | 0.255 | 0.248 |
| Pixtral-large | 0.220 | 0.196 | 0.205 | 0.223 | 0.199 | 0.207 | 0.145 | 0.165 | 0.153 | 0.146 | 0.166 | 0.154 | 0.119 | 0.124 | 0.120 | 0.120 | 0.126 | 0.122 | 0.162 | 0.179 | 0.168 | 0.164 | 0.181 | 0.170 |
| Qwen-3B | 0.311 | 0.325 | 0.308 | 0.314 | 0.329 | 0.311 | 0.291 | 0.292 | 0.282 | 0.300 | 0.299 | 0.290 | 0.139 | 0.120 | 0.121 | 0.143 | 0.123 | 0.124 | 0.163 | 0.137 | 0.142 | 0.164 | 0.138 | 0.143 |
| Qwen-7B | 0.480 | 0.569 | 0.517 | 0.487 | 0.578 | 0.524 | 0.471 | 0.590 | 0.517 | 0.477 | 0.597 | 0.524 | 0.482 | 0.622 | 0.536 | 0.493 | 0.639 | 0.549 | 0.467 | 0.603 | 0.519 | 0.478 | 0.619 | 0.532 |
| Non-VLM Baseline | | | | | | | | | | | | | | | | | | | | | | | | |
| BPMN Redrawer | 0.518 | 0.399 | 0.448 | 0.518 | 0.399 | 0.448 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Sketch2Process | 0.766 | 0.612 | 0.674 | 0.771 | 0.617 | 0.678 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |

as Qwen2.5VL-3B, Aya-Vision-8B, and Pixtral-12B, show little improvement or degrade with OCR. Their limited multimodal integration leaves them vulnerable to spurious OCR tokens. Finally, the **non-VLM baselines** perform competitively in name-only evaluation, underscoring the continuing value of traditional vision+OCR pipelines alongside VLMs.

6.2 Name+Type Evaluation (Table 2)

This evaluation demands correct extraction of both names and types, reported under strict and relaxed variants.

Strict setting: **Top-tier models** achieve the best results, though their F1 drops relative to name-only evaluation. GPT-4.1-mini slightly outperforms GPT-4.1, suggesting stronger calibration on type distinctions. OCR brings little benefit at this tier. **Mid-tier models** benefit substantially from OCR, particularly Pix2Struct and Tesseract, which help recover fine-grained type labels missed in vision-only inference. **Low-tier models** remain weak; noisy OCR often harms their performance further. The **non-VLM baselines** also

perform competitively here, with **Sketch2Process** in particular rivaling mid-tier VLMs; their advantage comes from specialized CNN+OCR pipelines tuned for BPMN structure, which compensate for the lack of multimodal reasoning.

Relaxed setting: Allowing partial type matches improves recall across all models. **Top-tier models** remain stable, showing minimal OCR gains, consistent with strong inherent multimodal grounding. **Mid-tier models** gain most (+3–5 F1), especially Pixtral-Large and Gemma-12B with Pix2Struct and Tesseract. **Low-tier models** continue to underperform, with OCR providing little or no benefit.

6.3 Relations+Type Evaluation (Table 3)

This more demanding evaluation requires both correct relation extraction and accurate typing.

Strict setting: **Top-tier models** maintain the highest scores, although relations prove harder to capture than names and types alone. OCR adds limited value at this tier. **Mid-tier models** clearly benefit from OCR, with Tesseract and Pix2Struct aiding recovery

Table 3: Strict vs. relaxed evaluation results (relations + types) with precision, recall, and F1 across VLM-only and OCR-enriched pipelines (best results in bold).

| Model | Only VLM | | | | | | VLM + Pix2Struct | | | | | | VLM + RapidOCR | | | | | | VLM + Tesseract | | | | | |
|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|------------------|--------------|--------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------|--------------|--------------|--------------|
| | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | |
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| Aya-vision-8B | 0.020 | 0.001 | 0.003 | 0.020 | 0.001 | 0.003 | 0.049 | 0.007 | 0.011 | 0.049 | 0.007 | 0.011 | 0.020 | 0.001 | 0.003 | 0.039 | 0.003 | 0.006 | 0.039 | 0.003 | 0.006 | 0.039 | 0.003 | 0.006 |
| Gemma-12B | 0.164 | 0.139 | 0.139 | 0.164 | 0.139 | 0.139 | 0.117 | 0.123 | 0.109 | 0.117 | 0.123 | 0.109 | 0.182 | 0.159 | 0.158 | 0.182 | 0.159 | 0.158 | 0.182 | 0.159 | 0.158 | 0.182 | 0.159 | 0.158 |
| Gemma-4B | 0.130 | 0.123 | 0.114 | 0.130 | 0.123 | 0.114 | 0.120 | 0.118 | 0.107 | 0.120 | 0.118 | 0.107 | 0.096 | 0.088 | 0.082 | 0.096 | 0.088 | 0.082 | 0.096 | 0.088 | 0.082 | 0.096 | 0.088 | 0.082 |
| GPT-4.1 | 0.519 | 0.474 | 0.475 | 0.519 | 0.474 | 0.475 | 0.333 | 0.498 | 0.336 | 0.333 | 0.498 | 0.336 | 0.316 | 0.452 | 0.313 | 0.316 | 0.452 | 0.313 | 0.321 | 0.426 | 0.307 | 0.321 | 0.426 | 0.307 |
| GPT-4.1-mini | 0.532 | 0.480 | 0.469 | 0.532 | 0.480 | 0.469 | 0.330 | 0.475 | 0.340 | 0.330 | 0.475 | 0.340 | 0.333 | 0.439 | 0.322 | 0.333 | 0.439 | 0.322 | 0.337 | 0.471 | 0.340 | 0.337 | 0.471 | 0.340 |
| GPT-4.1-nano | 0.229 | 0.167 | 0.175 | 0.229 | 0.167 | 0.175 | 0.187 | 0.191 | 0.170 | 0.187 | 0.191 | 0.170 | 0.184 | 0.164 | 0.160 | 0.184 | 0.164 | 0.160 | 0.176 | 0.151 | 0.146 | 0.176 | 0.151 | 0.146 |
| GPT-4o | 0.424 | 0.317 | 0.336 | 0.424 | 0.317 | 0.336 | 0.244 | 0.290 | 0.239 | 0.244 | 0.290 | 0.239 | 0.256 | 0.280 | 0.254 | 0.256 | 0.280 | 0.254 | 0.214 | 0.284 | 0.222 | 0.214 | 0.284 | 0.222 |
| GPT-4o-mini | 0.236 | 0.173 | 0.188 | 0.236 | 0.173 | 0.188 | 0.207 | 0.180 | 0.175 | 0.207 | 0.180 | 0.175 | 0.207 | 0.180 | 0.175 | 0.207 | 0.180 | 0.175 | 0.191 | 0.155 | 0.157 | 0.191 | 0.155 | 0.157 |
| LLaMA-3.2-11B | 0.058 | 0.054 | 0.051 | 0.058 | 0.054 | 0.051 | 0.057 | 0.038 | 0.042 | 0.057 | 0.038 | 0.042 | 0.077 | 0.059 | 0.060 | 0.077 | 0.059 | 0.060 | 0.042 | 0.065 | 0.046 | 0.042 | 0.065 | 0.046 |
| Mistral-small-3.1 | 0.424 | 0.412 | 0.395 | 0.424 | 0.412 | 0.395 | 0.261 | 0.421 | 0.299 | 0.261 | 0.421 | 0.299 | 0.262 | 0.395 | 0.295 | 0.262 | 0.395 | 0.295 | 0.233 | 0.361 | 0.267 | 0.233 | 0.361 | 0.267 |
| Pixtral-12B | 0.028 | 0.040 | 0.029 | 0.028 | 0.040 | 0.029 | 0.081 | 0.039 | 0.045 | 0.081 | 0.039 | 0.045 | 0.066 | 0.044 | 0.046 | 0.066 | 0.044 | 0.046 | 0.053 | 0.039 | 0.041 | 0.053 | 0.039 | 0.041 |
| Pixtral-large | 0.402 | 0.285 | 0.277 | 0.402 | 0.285 | 0.277 | 0.173 | 0.232 | 0.151 | 0.173 | 0.232 | 0.151 | 0.245 | 0.137 | 0.132 | 0.245 | 0.137 | 0.132 | 0.237 | 0.217 | 0.148 | 0.237 | 0.217 | 0.148 |
| Qwen-3B | 0.086 | 0.032 | 0.034 | 0.086 | 0.032 | 0.034 | 0.058 | 0.029 | 0.028 | 0.058 | 0.029 | 0.028 | 0.030 | 0.010 | 0.014 | 0.030 | 0.010 | 0.014 | 0.017 | 0.004 | 0.006 | 0.017 | 0.004 | 0.006 |
| Qwen-7B | 0.246 | 0.147 | 0.151 | 0.246 | 0.147 | 0.151 | 0.195 | 0.136 | 0.140 | 0.195 | 0.136 | 0.140 | 0.220 | 0.145 | 0.152 | 0.220 | 0.145 | 0.152 | 0.196 | 0.148 | 0.142 | 0.196 | 0.148 | 0.142 |
| Non-VLM Baseline | | | | | | | | | | | | | | | | | | | | | | | | |
| BPMN Redrawer | 0.098 | 0.142 | 0.103 | 0.098 | 0.142 | 0.103 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Sketch2Process | 0.201 | 0.287 | 0.217 | 0.201 | 0.287 | 0.217 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |

Table 4: Strict vs. relaxed evaluation results (types only) with precision, recall, and F1 across VLM-only and OCR-enriched pipelines (best results in bold).

| Model | Only VLM | | | | | | VLM + Pix2Struct | | | | | | VLM + RapidOCR | | | | | | VLM + Tesseract | | | | | |
|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|------------------|--------------|--------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------|--------------|--------------|--------------|
| | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | | Strict | | | Relaxed | | |
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| Aya-vision-8B | 0.604 | 0.929 | 0.715 | 0.704 | 0.993 | 0.814 | 0.601 | 0.930 | 0.714 | 0.701 | 1.000 | 0.814 | 0.622 | 0.925 | 0.725 | 0.702 | 0.989 | 0.810 | 0.617 | 0.941 | 0.730 | 0.707 | 0.993 | 0.816 |
| Gemma-12B | 0.683 | 0.585 | 0.617 | 0.707 | 0.675 | 0.683 | 0.617 | 0.643 | 0.614 | 0.686 | 0.693 | 0.680 | 0.719 | 0.801 | 0.743 | 0.790 | 0.846 | 0.808 | 0.719 | 0.801 | 0.743 | 0.790 | 0.846 | 0.808 |
| Gemma-4B | 0.634 | 0.922 | 0.738 | 0.698 | 0.964 | 0.800 | 0.630 | 0.934 | 0.738 | 0.696 | 0.980 | 0.805 | 0.652 | 0.957 | 0.761 | 0.710 | 1.000 | 0.821 | 0.652 | 0.957 | 0.761 | 0.710 | 1.000 | 0.821 |
| GPT-4.1 | 0.822 | 0.855 | 0.835 | 0.815 | 0.861 | 0.835 | 0.660 | 0.849 | 0.733 | 0.761 | 0.850 | 0.797 | 0.661 | 0.829 | 0.726 | 0.751 | 0.846 | 0.788 | 0.632 | 0.816 | 0.701 | 0.728 | 0.822 | 0.766 |
| GPT-4.1-mini | 0.839 | 0.878 | 0.855 | 0.842 | 0.874 | 0.853 | 0.648 | 0.873 | 0.733 | 0.768 | 0.862 | 0.806 | 0.669 | 0.886 | 0.753 | 0.788 | 0.877 | 0.824 | 0.662 | 0.878 | 0.745 | 0.789 | 0.866 | 0.819 |
| GPT-4.1-nano | 0.868 | 0.663 | 0.720 | 0.908 | 0.803 | 0.829 | 0.777 | 0.829 | 0.777 | 0.880 | 0.937 | 0.895 | 0.765 | 0.813 | 0.770 | 0.867 | 0.901 | 0.872 | 0.742 | 0.796 | 0.748 | 0.854 | 0.905 | 0.868 |
| GPT-4o | 0.814 | 0.852 | 0.828 | 0.823 | 0.857 | 0.835 | 0.691 | 0.867 | 0.758 | 0.775 | 0.876 | 0.814 | 0.686 | 0.884 | 0.759 | 0.783 | 0.883 | 0.823 | 0.702 | 0.879 | 0.769 | 0.786 | 0.879 | 0.822 |
| GPT-4o-mini | 0.834 | 0.891 | 0.850 | 0.871 | 0.970 | 0.908 | 0.785 | 0.928 | 0.836 | 0.863 | 0.979 | 0.909 | 0.785 | 0.928 | 0.836 | 0.863 | 0.979 | 0.909 | 0.773 | 0.927 | 0.828 | 0.854 | 0.985 | 0.907 |
| LLaMA-3.2-11B | 0.580 | 0.912 | 0.685 | 0.663 | 0.958 | 0.766 | 0.602 | 0.915 | 0.700 | 0.660 | 0.951 | 0.761 | 0.597 | 0.913 | 0.697 | 0.675 | 0.965 | 0.781 | 0.573 | 0.901 | 0.674 | 0.652 | 0.951 | 0.755 |
| Mistral-small-3.1 | 0.820 | 0.832 | 0.821 | 0.831 | 0.857 | 0.837 | 0.701 | 0.821 | 0.742 | 0.755 | 0.857 | 0.791 | 0.692 | 0.832 | 0.742 | 0.765 | 0.883 | 0.807 | 0.700 | 0.848 | 0.756 | 0.768 | 0.870 | 0.804 |
| Pixtral-12B | 0.678 | 0.743 | 0.677 | 0.741 | 0.853 | 0.771 | 0.639 | 0.816 | 0.704 | 0.695 | 0.887 | 0.768 | 0.598 | 0.782 | 0.656 | 0.682 | 0.855 | 0.744 | 0.548 | 0.721 | 0.609 | 0.604 | 0.805 | 0.677 |
| Pixtral-large | 0.348 | 0.306 | 0.322 | 0.360 | 0.354 | 0.351 | 0.246 | 0.294 | 0.264 | 0.282 | 0.305 | 0.291 | 0.215 | 0.235 | 0.220 | 0.242 | 0.260 | 0.247 | 0.252 | 0.312 | 0.275 | 0.297 | 0.321 | 0.306 |
| Qwen-3B | 0.447 | 0.730 | 0.544 | 0.533 | 0.765 | 0.619 | 0.437 | 0.692 | 0.526 | 0.525 | 0.745 | 0.608 | 0.242 | 0.380 | 0.290 | 0.286 | 0.412 | 0.334 | 0.288 | 0.447 | 0.343 | 0.337 | 0.471 | 0.388 |
| Qwen-7B | 0.621 | 0.764 | 0.678 | 0.629 | 0.819 | 0.705 | 0.571 | 0.783 | 0.646 | 0.643 | 0.837 | 0.716 | 0.614 | 0.863 | 0.707 | 0.666 | 0.902 | 0.758 | 0.610 | 0.850 | 0.701 | 0.662 | 0.882 | 0.748 |
| Non-VLM Baseline | | | | | | | | | | | | | | | | | | | | | | | | |
| BPMN Redrawer | 0.931 | 0.806 | 0.855 | 0.953 | 0.906 | 0.921 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Sketch2Process | 0.937 | 0.712 | 0.781 | 0.949 | 0.785 | 0.835 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |

of relation labels and directionality. **Low-tier models** struggle significantly, with OCR often worsening results by injecting spurious connections. The **non-VLM baselines** lag behind VLMs, as their handcrafted CNN+OCR pipelines are tuned for symbol and text recognition but lack semantic grounding, making relation directionality especially error-prone.

Relaxed setting: Relaxed matching improves recall noticeably, especially for mid-tier models. Pixtral-Large and Gemma-12B show gains of several F1 points under OCR. High-tier models remain strong and OCR-agnostic, while low-tier models remain unreliable, with noisy enrichment outweighing any benefit.

6.4 Type-Only Evaluation (Table 4)

Finally, we evaluate models on type classification alone.

Strict setting: **Top-tier models** maintain strong type discrimination, relying primarily on visual shape recognition. **Mid-tier models** benefit from OCR, particularly Pix2Struct, which helps with ambiguous event and gateway types. **Low-tier models** continue to underperform, with OCR introducing additional errors. The

non-VLM baselines achieve competitive results here because type recognition is largely a visual classification task; their CNN+OCR pipelines are highly optimized for symbol shapes, giving them an advantage over weaker VLMs.

Relaxed setting: Relaxed scoring boosts recall across models. **Top-tier systems** remain relatively unaffected by OCR. **Mid-tier models** show the greatest improvements (up to +4–5 F1), especially when enriched with Pix2Struct and Tesseract. **Low-tier models** remain weak, with OCR contributing little value.

6.5 Error Distribution Across Element Types

To complement the aggregate metrics, we analyzed which BPMN element types were more prone to errors across models. Table 5 summarizes the overall error rates normalized by the frequency of each element type in the gold-standard annotations.

The results indicate that certain control-flow constructs, such as *sequence flows*, *message flows*, and *gateways*, are consistently more challenging for models, yielding disproportionately high error rates compared to their gold frequency. This arises because these

Table 5: Top 10 BPMN element types with highest error rates across models.

| Element Type | Gold Count | Error Count | Error Rate |
|------------------------|------------|-------------|------------|
| datastore | 1 | 202 | 202.000 |
| sequenceflow | 482 | 66958 | 138.917 |
| messageflow | 136 | 15358 | 112.926 |
| intermediatecatchevent | 1 | 106 | 106.000 |
| startevent | 97 | 10188 | 105.031 |
| intermediateevent | 125 | 13020 | 104.160 |
| pool | 98 | 9376 | 95.673 |
| event | 1 | 85 | 85.000 |
| exclusivegateway | 106 | 8616 | 81.283 |
| endevent | 99 | 7605 | 76.818 |

elements depend on correctly modeling relational structure and directional semantics, which VLMs struggle to infer from visual cues alone. Conversely, simpler constructs like *tasks* and *events* also show substantial error counts due to their prevalence, though their normalized error rates are somewhat lower. Their errors typically stem from noisy OCR labels and frequent visual overlaps, rather than structural ambiguity.

7 Statistical Analysis

To complement the raw evaluation scores, we conducted statistical significance testing to assess whether differences across OCR settings and models were robust. We combined non-parametric tests with effect size reporting to provide both inferential and practical insights. Specifically, this section addresses the following research questions:

- **RQ1:** Are the observed improvements from OCR statistically significant across models, and do certain models benefit more than others?
- **RQ2:** How consistent are model performances across OCR settings, and what are the relative rankings of models when evaluated jointly across all conditions?
- **RQ3:** What is the practical magnitude of the impact of OCR on the performance of the model?

7.1 Pairwise OCR vs. VLM-only Comparisons

We first compared each OCR-enriched pipeline (Pix2Struct, RapidOCR, and Tesseract) against the VLM-only baseline for each model using the Wilcoxon signed-rank [53] test. Results in Table 6 show that top-performing models (e.g., GPT-4.1, GPT-4o, GPT-4.1-mini, and Mistral-Small-3.1) exhibit highly significant differences ($p < 0.01$) across all OCR methods. Mid-range models such as Qwen2.5-7B show selective gains, while weaker models (e.g., Aya-Vision, Pixtral-12B) show no significant improvements. This directly answers **RQ1**, showing that OCR benefits are model-dependent.

7.2 Model Consistency Across OCR Settings

To evaluate consistency across OCR conditions within each model, we applied the Friedman test [12] with Kendall’s [23] W as an effect size as shown in Table 7. High W values (e.g., GPT-4.1: $W = 1.0$; Pixtral-Large: $W = 0.889$; Gemma-12B: $W = 0.854$) indicate strong agreement that OCR systematically shifts performance. In contrast,

Table 6: Pairwise Wilcoxon test results (OCR vs VLM-only, F1). Lower p -values indicate stronger evidence of significant differences.

| Model | VLM vs Pix2Struct | VLM vs RapidOCR | VLM vs Tesseract |
|-------------------|-------------------|-----------------|------------------|
| Aya-Vision-8B | 0.945 | 0.383 | 0.203 |
| GPT-4.1 | 0.008 | 0.008 | 0.008 |
| GPT-4.1-mini | 0.008 | 0.008 | 0.008 |
| GPT-4.1-nano | 0.078 | 0.078 | 0.078 |
| GPT-4o | 0.008 | 0.008 | 0.008 |
| GPT-4o-mini | 0.008 | 0.008 | 0.008 |
| Gemma-4B | 0.883 | 0.883 | 0.883 |
| Gemma-12B | 0.008 | 0.008 | 0.008 |
| LLaMA-3.2-11B | 0.453 | 0.453 | 0.453 |
| Mistral-Small-3.1 | 0.008 | 0.008 | 0.008 |
| Pixtral-12B | 0.195 | 0.195 | 0.195 |
| Pixtral-Large | 0.008 | 0.008 | 0.008 |
| Qwen2.5-3B | 0.305 | 0.305 | 0.305 |
| Qwen2.5-7B | 0.016 | 0.016 | 0.016 |

models such as Aya-Vision and LLaMA-3.2-11B yield low W , suggesting random variability rather than systematic improvements. These findings address **RQ2**, highlighting which models are robust across OCR settings.

Table 7: Friedman test results across OCR settings (F1). Lower p -values (\downarrow) indicate significance; higher Kendall’s W (\uparrow) indicates stronger agreement ($W \in [0, 1]$).

| Model | p -value (\downarrow) | Kendall’s W (\uparrow) |
|-------------------|-----------------------------|------------------------------|
| Aya-Vision-8B | 0.648 | 0.069 |
| Gemma-4B | 0.881 | 0.028 |
| Gemma-12B | 0.000 | 0.854 |
| GPT-4.1 | 0.000 | 1.000 |
| GPT-4.1-mini | 0.000 | 0.775 |
| GPT-4.1-nano | 0.007 | 0.506 |
| GPT-4o | 0.002 | 0.625 |
| GPT-4o-mini | 0.007 | 0.471 |
| LLaMA-3.2-11B | 0.444 | 0.083 |
| Mistral-Small-3.1 | 0.000 | 0.791 |
| Pixtral-12B | 0.444 | 0.083 |
| Pixtral-Large | 0.000 | 0.889 |
| Qwen2.5-3B | 0.305 | 0.139 |
| Qwen2.5-7B | 0.037 | 0.292 |

7.3 Model Ranking Across Settings

The Friedman framework also allows us to rank models by average performance across all evaluation settings. Figure 2 illustrates that GPT-4.1, GPT-4o, GPT-4.1-mini, and Mistral-Small-3.1 consistently occupy the top ranks, while Pixtral-12B and Aya-Vision remain at the bottom. This ranking provides a global perspective on robustness independent of individual evaluation metrics and further answers **RQ2**.

7.4 Effect Sizes of OCR Impact

While significance tests establish whether OCR helps, effect sizes quantify “how much” it helps. Figure 3 present Cohen’s d [8] for each model and OCR method. Negative values indicate degradation with OCR, while positive values indicate improvements. Top-tier models (e.g., GPT-4.1, GPT-4o, Mistral-Small-3.1) show consistently

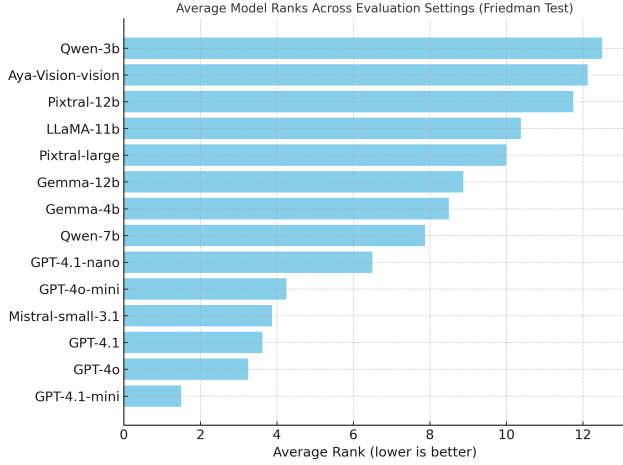


Figure 2: Average model ranks across evaluation settings based on the Friedman test. Lower ranks indicate better performance.

large negative d values, reflecting that OCR may inject noise into otherwise strong visual extraction. By contrast, mid-tier models (e.g., Qwen2.5-7B, Gemma-12B) yield positive d , confirming that they benefit most from additional textual cues. This analysis directly addresses **RQ3** by moving beyond statistical significance to practical significance.

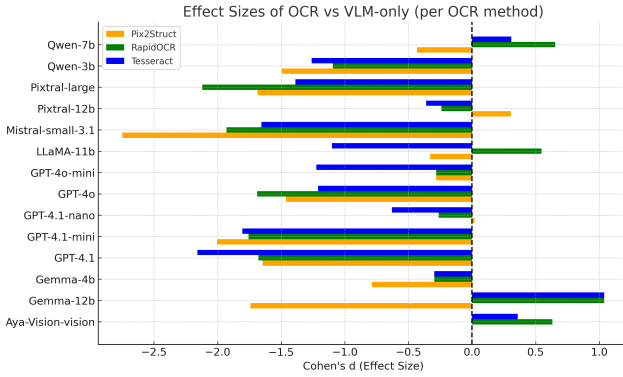


Figure 3: Effect sizes (Cohen's d) for OCR vs. VLM-only across models and OCR methods. Positive values indicate performance gains from OCR.

8 Ablation Study on Prompt Variants

To assess the impact of prompt design on BPMN extraction, we evaluated four variants¹⁰ of the baseline schema, ranging from minimal instructions to reasoning-augmented and traversal-based strategies as described in Table 8. Experiments were conducted across four evaluation settings: **Name-only**, **Type-only**, **Name+Type**

¹⁰The prompt variations can be accessed via https://github.com/pritamdeka/BPMN-VLM/tree/main/prompts/prompt_ablations

(Entities), and **Relation+Type (Relations)**. Table 9 reports strict F1 scores for each model and strategy, while Table 10 presents Wilcoxon significance tests, highlighting which prompt types yield consistent improvements.

Table 8: Prompting variants explored in the ablation study for BPMN diagram extraction.

| Variant | Description |
|-------------------------------|---|
| Baseline | Original prompt with explicit instructions for extracting BPMN elements, bounding boxes, labels, and flows, along with an example JSON schema. |
| Only-Example | Removes explicit instructions and supplies only the example JSON schema. This condition examines if the model can deduce the extraction strategy implicitly from the example. |
| Chain-of-Thought (CoT) | Adds hidden reasoning instructions: the model is asked to think step-by-step silently (scratchpad) but only emit the final JSON output. Ensures more systematic parsing while preventing verbose outputs. |
| Self-Consistency | Extends CoT by requiring the model to internally generate multiple candidate outputs and select the most consistent one. Implemented as best-of-3 decoding to reduce random errors and inconsistencies. |
| DFS+BFS Hybrid | Augments the baseline with traversal strategy: first apply Breadth-First Search (BFS) to capture overall structure (pools, lanes, tasks, events), then Depth-First Search (DFS) to follow each flow from start to end. Balances structural coverage with flow completeness. |

Analysis of Prompt Strategies. Table 9 shows that performance varies considerably by prompt type and model. The baseline prompt provides stable but modest scores across all evaluation categories. The Only-Example prompt consistently underperforms, confirming that example schemas alone are insufficient for robust diagram parsing. Chain-of-Thought and Self-Consistency generally yield incremental gains for models such as GPT-4.1, GPT-4o, Gemma-12B, suggesting that reasoning and ensemble-style decoding help capture complex control-flow constructs. The DFS+BFS hybrid prompt stands out as the most effective in many cases, especially for entity and relation extraction, as it enforces structural traversal of diagrams and reduces missed flows. Differences across models also highlight capacity effects: larger models exploit advanced prompting more effectively, while smaller models show only modest gains, often limited by representational power rather than prompting strategy.

Statistical Significance Test. The significance test results from Table 10 show that prompt strategies which reduce ambiguity, such as the **Baseline** and especially the **DFS+BFS hybrid**, consistently outperform alternatives. Baseline works reliably by pairing explicit instructions with an example schema, while DFS+BFS is most effective because its traversal policy mirrors the logical structure of BPMN diagrams, reducing missed dependencies and improving relational accuracy. In contrast, **Only-Example** performs weakest, as models fail to infer extraction rules from a schema alone. Reasoning-heavy strategies yield mixed outcomes: **Chain-of-Thought** adds limited benefit since BPMN extraction depends more on explicit mapping than abstract reasoning, and **Self-Consistency** mainly

Table 9: Strict evaluation results (F1 scores) across models and prompt strategies (best results in bold).

| Model | Baseline | | | | Only-Example | | | | Chain-of-Thought | | | | Self-Consistency | | | | DFS+BFS | | | |
|-------------------|----------|-------------|------|-------------|--------------|------|------|-------------|------------------|-------------|-------------|-------------|------------------|------|------|-------------|-------------|-------------|-------------|-------------|
| | Ent. | Rel. | Name | Type | Ent. | Rel. | Name | Type | Ent. | Rel. | Name | Type | Ent. | Rel. | Name | Type | Ent. | Rel. | Name | Type |
| Aya-Vision | 0.11 | 0.00 | 0.05 | 0.04 | 0.10 | 0.00 | 0.05 | 0.04 | 0.10 | 0.01 | 0.05 | 0.04 | 0.10 | 0.00 | 0.05 | 0.04 | 0.12 | 0.00 | 0.06 | 0.05 |
| Gemma-12B | 0.42 | 0.14 | 0.21 | 0.22 | 0.58 | 0.19 | 0.27 | 0.29 | 0.54 | 0.11 | 0.26 | 0.28 | 0.56 | 0.13 | 0.27 | 0.29 | 0.57 | 0.21 | 0.28 | 0.30 |
| Gemma-4B | 0.43 | 0.11 | 0.18 | 0.19 | 0.41 | 0.09 | 0.18 | 0.18 | 0.46 | 0.09 | 0.19 | 0.19 | 0.45 | 0.08 | 0.18 | 0.19 | 0.44 | 0.10 | 0.19 | 0.19 |
| GPT-4.1 | 0.56 | 0.22 | 0.30 | 0.31 | 0.71 | 0.26 | 0.36 | 0.38 | 0.70 | 0.27 | 0.36 | 0.37 | 0.70 | 0.24 | 0.35 | 0.37 | 0.72 | 0.29 | 0.37 | 0.38 |
| GPT-4.1-mini | 0.54 | 0.21 | 0.29 | 0.30 | 0.66 | 0.24 | 0.34 | 0.35 | 0.67 | 0.23 | 0.34 | 0.35 | 0.67 | 0.22 | 0.34 | 0.35 | 0.68 | 0.26 | 0.35 | 0.36 |
| GPT-4.1-nano | 0.33 | 0.09 | 0.17 | 0.18 | 0.42 | 0.11 | 0.20 | 0.21 | 0.42 | 0.10 | 0.20 | 0.21 | 0.41 | 0.10 | 0.19 | 0.20 | 0.43 | 0.13 | 0.21 | 0.22 |
| GPT-4o | 0.57 | 0.24 | 0.31 | 0.32 | 0.71 | 0.28 | 0.37 | 0.38 | 0.71 | 0.29 | 0.37 | 0.38 | 0.72 | 0.26 | 0.36 | 0.38 | 0.73 | 0.30 | 0.38 | 0.39 |
| GPT-4o-mini | 0.55 | 0.22 | 0.30 | 0.31 | 0.67 | 0.25 | 0.35 | 0.36 | 0.68 | 0.25 | 0.35 | 0.36 | 0.68 | 0.23 | 0.35 | 0.36 | 0.69 | 0.27 | 0.36 | 0.37 |
| LLaMA-11B | 0.38 | 0.11 | 0.20 | 0.21 | 0.48 | 0.13 | 0.23 | 0.24 | 0.48 | 0.12 | 0.23 | 0.24 | 0.47 | 0.11 | 0.23 | 0.24 | 0.49 | 0.14 | 0.24 | 0.25 |
| Mistral-Small-3.1 | 0.35 | 0.10 | 0.18 | 0.19 | 0.43 | 0.12 | 0.20 | 0.21 | 0.44 | 0.11 | 0.21 | 0.21 | 0.44 | 0.11 | 0.20 | 0.21 | 0.45 | 0.13 | 0.21 | 0.22 |
| Pixtral-12B | 0.42 | 0.16 | 0.22 | 0.23 | 0.53 | 0.18 | 0.26 | 0.27 | 0.53 | 0.18 | 0.26 | 0.27 | 0.53 | 0.17 | 0.26 | 0.27 | 0.54 | 0.20 | 0.27 | 0.28 |
| Pixtral-Large | 0.44 | 0.18 | 0.23 | 0.24 | 0.56 | 0.20 | 0.28 | 0.29 | 0.56 | 0.20 | 0.28 | 0.29 | 0.56 | 0.19 | 0.27 | 0.29 | 0.57 | 0.22 | 0.29 | 0.30 |
| Qwen-3B | 0.31 | 0.08 | 0.15 | 0.16 | 0.39 | 0.09 | 0.18 | 0.19 | 0.39 | 0.09 | 0.18 | 0.19 | 0.39 | 0.09 | 0.18 | 0.19 | 0.40 | 0.10 | 0.19 | 0.20 |
| Qwen-7B | 0.37 | 0.11 | 0.19 | 0.20 | 0.47 | 0.13 | 0.22 | 0.23 | 0.47 | 0.12 | 0.22 | 0.23 | 0.47 | 0.12 | 0.22 | 0.23 | 0.48 | 0.14 | 0.23 | 0.24 |

helps mid-sized models by reducing stochastic errors. Model-wise, **large models** such as GPT-4.1, GPT-4o benefit most from structured traversal, **mid-sized models** (Gemma-12B, Qwen-7B, Pixtral-12B) gain moderately from both DFS+BFS and Self-Consistency, while **smaller models** (GPT-4.1-mini/nano, Gemma-4B, Qwen-3B) are most sensitive, failing under Only-Example and CoT but performing well for DFS+BFS.

Table 10: Wilcoxon significance test results (p-values) for Name+Type and Relations+Type across prompt types (lowest p-values in bold).

| Model | Baseline | | COT | | DFS-BFS | | Only_example | | Self_consistency | |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|------------------|-------------|
| | Ent | Rel | Ent | Rel | Ent | Rel | Ent | Rel | Ent | Rel |
| Aya-Vision | 0.74 | 0.29 | 0.76 | 0.21 | 0.44 | 0.09 | 0.77 | 0.09 | 0.75 | 0.32 |
| Gemma-12B | 0.01 | 0.70 | 0.36 | 0.30 | 0.43 | 0.29 | 0.41 | 0.31 | 0.48 | 0.42 |
| Gemma-4B | 0.39 | 0.28 | 0.54 | 0.49 | 0.55 | 0.41 | 0.26 | 0.14 | 0.56 | 0.54 |
| GPT-4.1 | 0.48 | 0.70 | 0.44 | 0.68 | 0.59 | 0.78 | 0.63 | 0.67 | 0.59 | 0.59 |
| GPT-4.1-mini | 0.69 | 0.83 | 0.72 | 0.85 | 0.73 | 0.91 | 0.70 | 0.85 | 0.68 | 0.84 |
| GPT-4.1-nano | 0.71 | 0.79 | 0.67 | 0.82 | 0.64 | 0.91 | 0.65 | 0.80 | 0.66 | 0.81 |
| GPT-4o | 0.47 | 0.66 | 0.48 | 0.65 | 0.52 | 0.72 | 0.51 | 0.69 | 0.53 | 0.70 |
| GPT-4o-mini | 0.62 | 0.76 | 0.61 | 0.77 | 0.58 | 0.72 | 0.60 | 0.73 | 0.59 | 0.74 |
| LLaMA-11B | 0.44 | 0.60 | 0.42 | 0.58 | 0.45 | 0.61 | 0.43 | 0.59 | 0.46 | 0.62 |
| Mistral-Small-3.1 | 0.54 | 0.72 | 0.53 | 0.71 | 0.55 | 0.73 | 0.56 | 0.74 | 0.55 | 0.73 |
| Pixtral-12B | 0.40 | 0.63 | 0.39 | 0.62 | 0.41 | 0.64 | 0.42 | 0.65 | 0.41 | 0.64 |
| Pixtral-Large | 0.46 | 0.68 | 0.47 | 0.69 | 0.48 | 0.70 | 0.49 | 0.71 | 0.47 | 0.69 |
| Qwen-3B | 0.50 | 0.72 | 0.51 | 0.73 | 0.52 | 0.74 | 0.53 | 0.75 | 0.51 | 0.73 |
| Qwen-7B | 0.44 | 0.65 | 0.45 | 0.66 | 0.47 | 0.68 | 0.46 | 0.67 | 0.45 | 0.66 |

9 Conclusion and Future Work

This study demonstrates that carefully designed prompts enable VLMs to extract structured information from BPMN diagram images. We conducted a systematic evaluation across four complementary settings, entities, relations, types, and name-only, alongside a range of prompting strategies, showing how both the evaluation perspective and prompt design shape performance. Our results further reveal that OCR integration provides selective gains, particularly for weaker models that struggle to capture structured semantics. To support nuanced comparisons, we also introduced a tiered evaluation framework (strict and relaxed variants), which highlights robustness and element-level accuracy.

For future work, we plan to investigate fine-tuning on larger and more diverse BPMN datasets, with a focus on cross-diagram generalization. We will also explore advanced prompt engineering and multi-step reasoning techniques to improve semantic parsing and ensure consistent conversion of BPMN images into executable XML, thereby enabling automated process mining and seamless integration into real-world business environments.

Acknowledgments

We thank the maintainers of Hugging Face, OpenAI, Mistral, and other VLM developers for access to their models and APIs. The authors acknowledge the use of generative AI tools (e.g., ChatGPT) for support in code implementation and text editing; all conceptual contributions and experimental results are original and have been independently verified. This research is supported by the Advanced Research and Engineering Centre (ARC) in Northern Ireland, funded by PwC and Invest NI. The views expressed are those of the authors and do not necessarily represent those of ARC or the funding organisations.

References

- [1] Praveesh Agrawal et al. 2024. Pixtral 12b. *arXiv preprint arXiv:2410.07073*.
- [2] Meta AI. 2024. Llama 3.2 11b vision instruct model. (2024). <https://huggingface.co/meta-llama/Llama-3.2-11B-Vision-Instruct>.
- [3] Thomas Allweyer. 2016. *BPMN 2.0: introduction to the standard for business process modeling*. BoD–Books on Demand.
- [4] Alessandro Antinori, Riccardo Coltrinari, Flavio Corradini, Fabrizio Fornari, Barbara Re, Marco Scarpetta, et al. 2022. Bpmn-redrawer: from images to bpmn models. In.
- [5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: a versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- [6] Patrizio Bellan, Mauro Dragoni, and Chiara Ghidini. 2022. Extracting business process entities and relations from text using pre-trained language models and in-context learning. In *International Conference on Enterprise Design, Operations, and Computing*. Springer, 182–199.
- [7] Zhe Chen et al. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67, 12, 220101.
- [8] Jacob Cohen. 2013. *Statistical power analysis for the behavioral sciences*. routledge.
- [9] Flavio Corradini, Sara Pettinari, Barbara Re, Lorenzo Rossi, and Francesco Tiezzi. 2024. A technique for discovering bpmn collaboration diagrams. *Software and Systems Modeling*, 1–21.

- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 4171–4186.
- [11] Théo Fagnoni, Bellinda Mesbah, Mahsun Altin, and Phillip Kingston. 2024. Opus: a large work model for complex workflow generation. *arXiv preprint arXiv:2412.00573*.
- [12] Milton Friedman. 1937. The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *Journal of the american statistical association*, 32, 200, 675–701.
- [13] Fabian Friedrich, Jan Mendling, and Frank Puhlmann. 2011. Process model generation from natural language text. In *Advanced Information Systems Engineering: 23rd International Conference, CAiSE 2011, London, UK, June 20-24, 2011. Proceedings 23*. Springer, 482–496.
- [14] Chaoyou Fu et al. 2025. Vita-1.5: towards gpt-4o level real-time vision and speech interaction. *arXiv preprint arXiv:2501.01957*.
- [15] Camunda Services GmbH. 2015. Bpmn for research. <https://github.com/camunda/bpmn-for-research>. Accessed: 2025-04-04. (2015).
- [16] Camunda Services GmbH. 2014. Bpmn-js: a bpmn 2.0 rendering toolkit and web modeler. <https://github.com/bpmn-io/bpmn-js>. Accessed: 2025-04-04. (2014).
- [17] 2023. Gpt-4v(ision) system card. In <https://api.semanticscholar.org/CorpusID:263218031>.
- [18] Michael Grohs, Luka Abb, Nourhan Elsayed, and Jana-Rebecca Rehse. 2023. Large language models can accomplish business process management tasks. In *International Conference on Business Process Management*. Springer, 453–465.
- [19] Object Management Group. 2014. Business process model and notation (bpmn) specification. (2014). <https://www.omg.org/spec/BPMN/2.0>.
- [20] Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023. Chartllama: a multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*.
- [21] Linmei Hu, Duokang Wang, Yiming Pan, Jifan Yu, Yingxia Shao, Chong Feng, and Liqiang Nie. 2024. Novachart: a large-scale dataset towards chart understanding and generation of multimodal large language models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 3917–3925.
- [22] Aaron Hurst et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- [23] Maurice G Kendall and B Babington Smith. 1939. The problem of m rankings. *The annals of mathematical statistics*, 10, 3, 275–287.
- [24] Muhammad Tayyab Khan, Lequn Chen, Ye Han Ng, Wenhe Feng, Nicholas Yew Jin Tan, and Seung Ki Moon. 2024. Fine-tuning vision-language model for automated engineering drawing information extraction. *arXiv preprint arXiv:2411.03707*.
- [25] Ryan KL Ko, Stephen SG Lee, and Eng Wah Lee. 2009. Business process management (bpm) standards: a survey. *Business process management journal*, 15, 5, 744–791.
- [26] Mateja Kocbek, Gregor Jošt, Marjan Heričko, and Gregor Polančič. 2015. Business process model and notation: the current state of affairs. *Computer Science and Information Systems*, 12, 2, 509–539.
- [27] Julius Köpke and Aya Safan. 2024. Efficient llm-based conversational process modeling. In *International Conference on Business Process Management*. Springer, 259–270.
- [28] Humam Kourani, Alessandro Berti, Daniel Schuster, and Wil MP Van der Aalst. 2024. Promoi: process modeling with generative ai. *arXiv preprint arXiv:2403.04327*.
- [29] Manuel Laguna and Johan Marklund. 2018. *Business process modeling, simulation and design*. Chapman and Hall/CRC.
- [30] Kenton Lee et al. 2023. Pix2struct: screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*. PMLR, 18893–18912.
- [31] Josip Tomo Licardo, Nikola Tanković, and Darko Etinger. 2024. A method for extracting bpmn models from textual descriptions using natural language processing. *Procedia computer science*, 239, 483–490.
- [32] Leilei Lin, Yumeng Jin, Yingming Zhou, Wenlong Chen, and Chen Qian. 2024. Mao: a framework for process model generation with multi-agent orchestration. *arXiv preprint arXiv:2408.01916*.
- [33] Pan Lu et al. 2023. Mathvista: evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- [34] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: a benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- [35] Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2024. Chartinstruct: instruction tuning for chart comprehension and reasoning. *arXiv preprint arXiv:2403.09028*.
- [36] Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty. 2024. Chartgemma: visual instruction-tuning for chart reasoning in the wild. *arXiv preprint arXiv:2407.04172*.
- [37] Carlos Matos and Reiko Heckel. 2009. Migrating legacy systems to service-oriented architectures. *Electronic Communications of the EASST*, 16.
- [38] Dipankar Medhi. 2024. Target prompting for information extraction with vision language model. *arXiv preprint arXiv:2408.03834*.
- [39] Mirjam Minor and Eduard Kaucher. 2024. Retrieval augmented generation with llms for explaining business process models. In *International Conference on Case-Based Reasoning*. Springer, 175–190.
- [40] Ali Nour Eldin, Nour Assy, Olan Anesini, Benjamin Dalmás, and Walid Gaaloul. 2024. Nala2bpmn: automating bpmn model generation with large language models. In *International Conference on Cooperative Information Systems*. Springer, 398–404.
- [41] Josselin Roberts, Tony Lee, Chi Heem Wong, Michihiro Yasunaga, Yifan Mai, and Percy S Liang. 2024. Image2struct: benchmarking structure extraction for vision-language models. *Advances in Neural Information Processing Systems*, 37, 115058–115097.
- [42] Bernhard Schäfer, Han Van Der Aa, Henrik Leopold, and Heiner Stuckenschmidt. 2022. Sketch2process: end-to-end bpmn sketch recognition based on neural networks. *IEEE Transactions on Software Engineering*, 49, 4, 2621–2641.
- [43] Anna Scius-Bertrand, Michael Jungo, Lars Vögtlin, Jean-Marc Spat, and Andreas Fischer. 2025. Zero-shot prompting and few-shot fine-tuning: revisiting document image classification using large language models. In *International Conference on Pattern Recognition*. Springer, 152–166.
- [44] Sholiq Sholiq, Riyanarto Sarno, and Endang Siti Astuti. 2022. Generating bpmn diagram from textual requirements. *Journal of King Saud University-Computer and Information Sciences*, 34, 10, 10079–10093.
- [45] Shubhankar Singh, Purvi Chaurasia, Yerram Varun, Pranshu Pandya, Vatsal Gupta, Vivek Gupta, and Dan Roth. 2024. Flowvqa: mapping multimodal logic in visual question answering with flowcharts. *arXiv preprint arXiv:2406.19237*.
- [46] Ray Smith. 2007. An overview of the tesseract ocr engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)*. Vol. 2. IEEE, 629–633.
- [47] Simon Tannert, Marcelo G Feighelstein, Jasmina Bogojeska, Joseph Shtok, Assaf Arbelle, Peter WJ Staar, Anika Schumann, Jonas Kuhn, and Leonid Karlinsky. 2023. Flowchartqa: the first large-scale benchmark for reasoning over flowcharts. In *Proceedings of the 1st Workshop on Linguistic Insights from and for Multimodal Language Processing*, 34–46.
- [48] Gemma Team et al. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.
- [49] Qwen Team. 2025. Qwen2.5-vl. (Jan. 2025). <https://qwenlm.github.io/blog/qwen2.5-vl/>.
- [50] RapidAI Team. 2021. Rapid OCR: ocr toolbox. <https://github.com/RapidAI/RapidOCR>. (2021).
- [51] Yuan Tian, Weiwei Cui, Dazhen Deng, Xinjing Yi, Yurun Yang, Haidong Zhang, and Yingcai Wu. 2024. Chartgpt: leveraging llms to generate charts from abstract natural language. *IEEE Transactions on Visualization and Computer Graphics*.
- [52] Wil Van Der Aalst. 2012. Process mining: overview and opportunities. *ACM Transactions on Management Information Systems (TMIS)*, 3, 2, 1–17.
- [53] Frank Wilcoxon. 1945. Individual comparisons by ranking methods. *Biometrics bulletin*, 1, 6, 80–83.
- [54] Thomas Wolf et al. 2020. Transformers: state-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, 38–45.
- [55] Renqiu Xia et al. 2024. Chartx & chartvlm: a versatile benchmark and foundation model for complicated chart reasoning. *arXiv preprint arXiv:2402.12185*.
- [56] Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun Yuan, and Jian Guo. 2023. Chartbench: a benchmark for complex visual reasoning in charts. *arXiv preprint arXiv:2312.15915*.
- [57] Renrui Zhang et al. 2024. Mathverse: does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*. Springer, 169–186.