# Flowchart2Mermaid: A Vision-Language Model Powered System for Converting Flowcharts into Editable Diagram Code

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

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

## Abstract

Flowcharts are common tools for communicating processes but are often shared as static images that cannot be easily edited or reused. We present FLOWCHART2MERMAID, a lightweight web system that converts flowchart images into editable Mermaid.js code which is a markup language for visual workflows, using a detailed system prompt and vision-language models. The interface supports mixed-initiative refinement through inline text editing, drag-and-drop node insertion, and natural-language commands interpreted by an integrated AI assistant. Unlike prior image-to-diagram tools, our approach produces a structured, version-controllable textual representation that remains synchronized with the rendered diagram. We further introduce evaluation metrics to assess structural accuracy, flow correctness, syntax validity, and completeness across multiple models.

## 1 Introduction

Flowcharts are central to communicating processes and structures in software engineering, algorithm design, and education (Ensmenger, 2016; Smetsers-Weeda and Smetsers, 2017). Yet most flowcharts exist as static images, which limits reuse: changing a label or inserting a decision often requires redrawing the entire figure. Automatically converting flowchart images into structured textual code also paves the way for more flexible multimodal reasoning (Suri et al., 2025; Soman et al., 2025; Ye et al., 2025). Flowcharts and related procedural diagrams encode stepwise logic and control flow that are difficult to capture through conventional OCR or shape-recognition pipelines alone. With recent advances in vision language models (VLMs), such as GPT-4.1[1] (Achiam et al., 2023) and Gemini-2.5 (Comanici et al., 2025), it has become possible to

infer not only visual structures but also their logical interrelations, enabling systems that translate complex diagram images into expressive, editable representations.

We introduce FLOWCHART2MERMAID[2], an end-to-end web application that converts flowchart images into textual Mermaid.js[3] (Sveidqvist and Jain, 2021) specifications and enables rich human-in-the-loop editing. After uploading an image, a vision–language model (VLM) produces an initial Mermaid program, which users can refine through three complementary interaction modes: (1) inline editing, by clicking labels directly on the rendered diagram; (2) visual editing, via drag-and-drop insertion of standard flowchart symbols that automatically synchronize with the underlying code; and (3) natural-language edits, where free-form instructions trigger targeted LLM modifications to the diagram (*e.g., "Connect* A *to* B *with a 'Yes' edge"*). The system supports real-time rendering, export to SVG/MMD formats, and seamless handoff to the Mermaid Live Editor. The landing page of the application is shown in Figure 1.
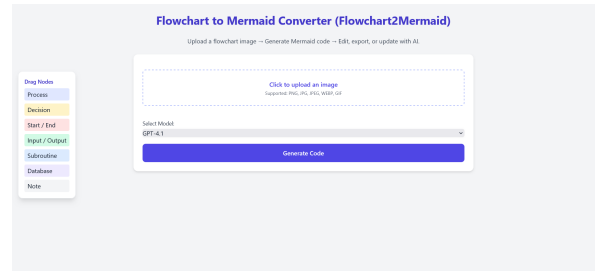


Figure 1: Flowchart2Mermaid Homepage

Our contributions are: (i) a lightweight, production-ready demo for multimodal diagram-to-code conversion; (ii) a unified interaction design combining textual, visual, and natural-language editing; and (iii) a model-agnostic backend that ac-

---

[1] https://openai.com/index/gpt-4-1/

[2] https://flowchart-to-mermaid.vercel.app/
[3] https://mermaid.js.org/

commodates multiple VLM providers while maintaining output validity. A demonstration video of the work is available for viewing via the Dropbox link[4].

## 2 Related Work

This section reviews research relevant to transforming visual diagrams into structured, editable representations. We summarize prior work on diagram parsing, image-to-code translation, and interactive tools for programmatic diagram editing.

**Diagram Understanding and Diagram-to-Code.** Early systems for flowchart interpretation focused on rule-based shape detection and template matching, enabling conversion of hand-drawn or scanned flowcharts into pseudocode or C programs (Wu et al., 2011; Herrera-Camara and Hammond, 2017; Carton et al., 2013). Educational environments such as *Flowgorithm* (Cook) or browser-based structured editors (Supaartagorn, 2017) support generating executable code from diagrams, but require users to construct diagrams within the tool itself and cannot process arbitrary images. More recent approaches incorporate deep-learning pipelines: Arrow-RCNN (Schäfer and Stuckenschmidt, 2019) and DrawnNet (Fang et al., 2022) use CNNs to detect shapes and connectors, while Montellano et al. (Montellano et al., 2022) reconstruct both the visual diagram and corresponding code from hand-drawn sketches[5]. Although effective within controlled datasets, these systems lack generalization to diverse diagram styles and do not support interactive refinement.

**Multimodal Models for Structured Visual Content.** Advances in vision language models (VLMs) have enabled more flexible diagram parsing. The *GenFlowchart* system (Arbaz et al., 2024) combines SAM-based (Segment Anything Model) (Kirillov et al., 2023) segmentation with an LLM to assemble a unified symbolic representation from detected visual elements. Foundation models such as GPT-4 (Achiam et al., 2023) and Gemini (Team et al., 2023) can translate images into structured text or code, demonstrating strong priors about graphical layout and syntax. Recent work has ex-

tended these capabilities to *business process diagrams* (Deka and Devereux, 2025) where the authors present a VLM-based pipeline for extracting structured JSON representations from BPMN images, enriched with OCR-based label recovery and alignment against XML ground truth. These multimodal systems enable diagram-to-text generation as a high-level translation task, moving beyond handcrafted rules. However, prior work typically treats generation as a one-shot inference problem and does not integrate downstream editing workflows.

**Interactive Editing and Natural Language-driven Diagram Manipulation.** A growing ecosystem of AI-assisted diagram editors has recently emerged, including commercial platforms such as FlowchartAI[6], dAIgram[7], which convert uploaded images into editable flowcharts using shape detection and layout reconstruction. Tools such as DiagramGPT[8] and Lucidchart AI[9] provide natural-language (NL) driven diagram creation, while yEd Live[10] integrates ChatGPT for modifying existing digital graphs (e.g., renaming or styling nodes). Although these systems demonstrate promising forms of AI-assisted editing, they cover only parts of the broader diagram–understanding space. As a result, existing tools rarely support semantic editing, version-control-friendly outputs, or deep coupling between recognition, code generation, and iterative refinement. They further lack bidirectional synchronization between visual edits and a symbolic diagram representation, and their natural-language capabilities are typically limited to high-level prompting rather than fine-grained, code-level refinement.

In contrast to the systems above, our approach targets the full pipeline from *raw diagram images* to *textual*, semantically meaningful Mermaid code, supporting downstream reasoning and reproducible version control. While recent work (Deka and Devereux, 2025) demonstrates robust VLM-based extraction for BPMN, their focus is on faithful structured recovery rather than interactive editing. Beyond initial multimodal parsing, our interface enables iterative refinement through three coordinated
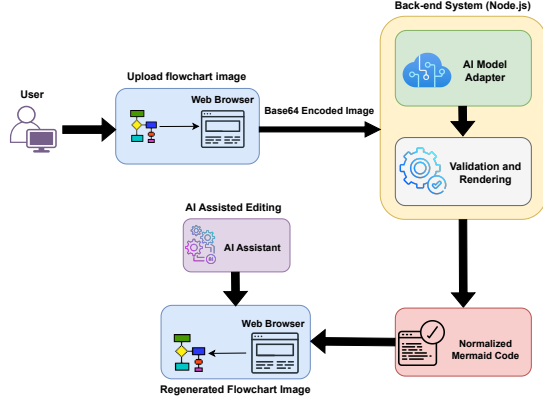
---

Figure 2: System Overview

interaction modes: inline editing of text within the rendered diagram, drag-and-drop insertion of common flowchart elements, and natural-language commands that trigger targeted code modifications. All edits such as visual, textual, or linguistic are synchronized with the underlying `Mermaid` program, maintaining a consistent representation throughout. To our knowledge, no prior system combines diagram-image understanding with a bi-directional code–visual editor and natural-language manipulation, making our work distinct within both AI-assisted diagram editing and multimodal understanding research.

## 3 System Overview

Figure 2 provides a high-level view of the FLOWCHART2MERMAID system. The application follows a lightweight client–server design in which the browser is responsible for visualization and user interaction, while a minimal server component mediates communication with external vision–language models (VLMs).

At a conceptual level, the system supports a mixed-initiative workflow: a user uploads a flowchart image, obtains an initial `Mermaid` representation generated by a multimodal model, and then iteratively refines it through direct manipulation, text editing, or natural-language commands. Throughout the workflow, the textual and visual representations remain tightly synchronized, supporting seamless switching between code-level and diagram-level editing.

## 4 Implementation

The system is implemented as a lightweight web application integrating VLM inference with an in-

teractive visualization layer. Its architecture balances modularity, responsiveness, and reproducibility across multiple model back ends.

### 4.1 Architecture Overview

The system is composed of a browser-based front end for visualization and editing and a lightweight Node.js back end that handles multimodal model inference. This separation keeps user interactions responsive while delegating computationally intensive processing to the server.
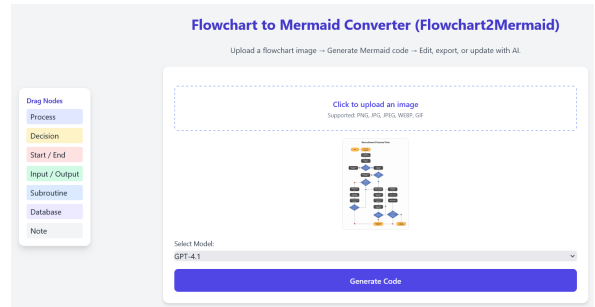


Figure 3: Uploaded example flowchart image

After a user uploads a diagram image as shown in Figure 3, the front end encodes it and forwards it, together with the chosen model identifier, to the back end. Inference is guided by a carefully engineered system prompt that specifies the expected flowchart elements, structural constraints, and output format of valid `Mermaid` code. This prompt serves as the core of the conversion process, drawing on insights from structured prompting and constrained generation (Liu et al., 2023; Cheng et al., 2024). The VLM produces a first-pass `Mermaid` program, which is lightly normalized before being returned to the browser as shown in Figure 4.

On the client side, the code is rendered immediately and remains fully editable. Users may refine the diagram through direct manipulation (dragging nodes, adding connections, modifying labels), through text-based editing of the `Mermaid` program, or via natural-language commands handled by an integrated AI assistant consisting of the GPT-4.1 model. All edits maintain bidirectional synchronization between the textual and visual representations. This architecture supports rapid iteration while remaining model-agnostic and easily extensible. We showcase the complete working demo for the example flowchart image in video.
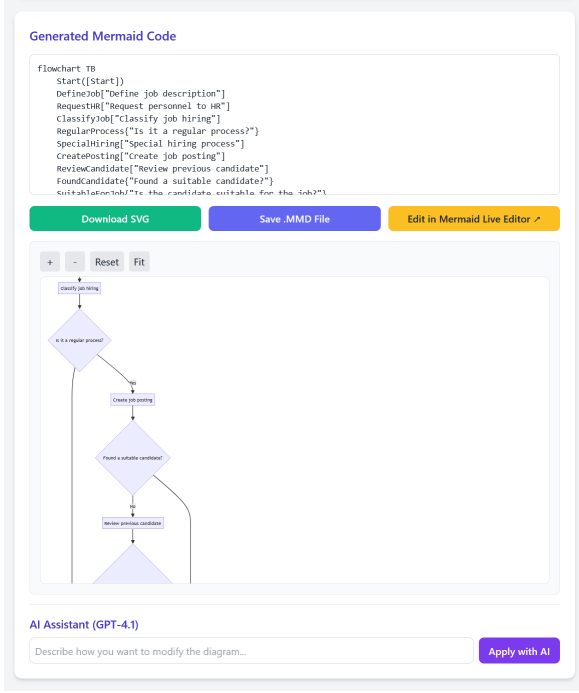
3

Figure 4: Generated mermaid code from example flowchart image

## 4.2 Front-end Interface

The front end is implemented using `HTML`, `TailwindCSS`, and `JavaScript`, and provides an interactive diagram-editing environment powered by `Mermaid.js`. Generated flowcharts are rendered inside a dedicated scrollable viewport that supports controlled zoom (via UI buttons) and smooth navigation, ensuring that large diagrams remain viewable without requiring excessive page scrolling. All elements are fully responsive, and the visual canvas automatically scales diagrams to fit the available space.

Users can modify diagrams directly through in-line editing of node labels, drag-and-drop insertion of common flowchart primitives (e.g., processes, decisions, I/O nodes, and databases), and structural operations such as repositioning or deleting elements. The system supports exporting diagrams as `.mmd` or `.svg` files and provides one-click access to the Mermaid Live Editor for advanced editing or collaboration.

An embedded **AI assistant** allows natural-language modification of diagrams. For example, users can issue commands like "*add a decision node before process C*" or "*connect start to review*", which are parsed into structured transformations of the current Mermaid graph. The interface maintains bidirectional consistency between textual and graphical representations, ensuring that each update is immediately reflected in both views.

## 4.3 Back-end Functionality

The back end orchestrates the conversion and refinement process. Two primary operations are supported:

**Image-to-code generation:** The back end receives an encoded diagram image and uses the chosen model to generate Mermaid syntax. The output is cleaned of Markdown artifacts, normalized for consistent directionality, and validated.

**Code refinement:** The current code and user instruction are combined into a controlled prompt, enforcing deterministic decoding parameters for reproducibility. The resulting update is validated and returned for rendering.

This unified design abstracts away API-specific differences between providers, offering a consistent interface to the front end while enabling easy extension to future multimodal models.

## 4.4 System Reliability and Performance

Strict validation routines ensure that all inputs conform to supported image types and that generated outputs compile successfully in `Mermaid.js`. The architecture emphasizes modular extensibility and low-latency feedback, making it suitable for both research exploration and pedagogical use cases involving diagrammatic reasoning.

## 5 Evaluation

We evaluate our flowchart-to-Mermaid pipeline using a subset of the FLOWVQA dataset (Singh et al., 2024), which contains diverse, real-world flowcharts annotated with step-level semantics. Although intended for visual question answering, its Mermaid annotations for each image allow us to construct the gold-standard (image, code) pairs for our task.

Because Mermaid code admits multiple semantically equivalent surface forms, direct string matching is unreliable: models may emit different node identifiers, paraphrased labels, or alternate arrow styles while still producing a fully correct diagram. These variations make rigid token- or rule-based evaluation brittle. Following recent work on *LLM-as-a-judge* evaluation (Zheng et al., 2023; Li et al., 2025), we therefore use a structured evaluator implemented within the `DeepEval` framework (AI,

4

2025).[11]. A GPT-based judge model (GPT-4.1) is prompted to parse both predicted and gold Mermaid diagrams, normalize node labels, identify directed relations, and compute symbolic precision, recall, and F1-score for nodes and edges. The full evaluation prompts are provided in Appendix A.

## 5.1 Setup

We evaluated five vision-language models integrated into our system: `GPT-4.1`, `GPT-4.1-mini`, `GPT-4o`, `GPT-4o-mini`, and `Gemini-2.5-Flash`. All models used the same structured instruction template and were tasked with generating Mermaid flowcharts for 200 input diagrams. Predicted Mermaid code was parsed automatically, normalized, and compared against gold-standard annotations.

## 5.2 Symbolic Metrics

These metrics quantify symbolic alignment and graph-topological similarity. We compute:
**Entity extraction accuracy:** precision (P), recall (R), and F1-score between predicted and gold node labels.
**Relationship extraction accuracy:** P, R, and F1-score for directed edges.
**Semantic similarity:** Cosine similarity between the predicted and gold Mermaid code using SBERT embeddings (Reimers and Gurevych, 2019).

## 5.3 High-Level Mermaid Structural Metrics

To evaluate global diagram quality beyond exact node–edge matching, we use five complementary structural metrics assessed by GPT-4.1 (Appendix B). Each metric captures a distinct aspect of diagram-level fidelity:
**Structural Accuracy (0–1):** high-level alignment of node–edge organization.
**Flow Accuracy (0–1):** preservation of control-flow and execution paths.
**Syntax Validity (0–1):** conformance to Mermaid syntax and renderability.
**Semantic Fidelity (0–1):** retention of the workflow's intended meaning.
**Completeness (0–1):** inclusion of essential elements and relations from the gold diagram.

A **reconstructability override** is applied when the judge model determines that the predicted diagram could be transformed into the gold version through minor structure-preserving edits, granting full credit across metrics. Each metric is originally

scored on its native bounded scale ($0-5$, $0-3$, or $0-2$) and subsequently normalized to the range $[0, 1]$ for comparability. While these metrics are conceptually formalized, their implementation relies on GPT-4.1's implicit reasoning rather than explicit numerical computation.

## 6 Results and Discussion

Tables 1 and 2 summarize performance on the FLOWVQA subset. Table 1 focuses on symbolic extraction quality (entities, relations, SBERT cosine similarity), while Table 2 captures higher-level structural properties of the generated diagrams. For SBERT cosine similarity, we use a fine-tuned MiniLM (Wang et al., 2020) model from Huggingface[12].

**Symbolic extraction performance.** All large models achieve strong entity extraction, with entity F1-scores above 0.94 for `GPT-4.1`, `GPT-4o`, and `Gemini-2.5-Flash`. Relationship extraction is consistently more difficult: relation F1 scores are typically 1-3 points lower than entity F1 scores, indicating that recovering the full edge structure is the primary source of error. `GPT-4.1-mini` and `Gemini-2.5-Flash` yield the strongest overall symbolic performance for both entity and relation F1-score. `GPT-4o` produces slightly lower relation F1-score, suggesting occasional omissions in branching structure while `GPT-4o-mini` is the weakest model overall. Cosine similarity scores are uniformly high, showing that label semantics are well preserved. Interestingly, `GPT-4o-mini` attains one of the highest cosine similarities despite poor structural F1-score, highlighting that embedding similarity alone cannot identify missing or misrouted transitions.

**High-level structural behaviour.** The high-level Mermaid metrics in Table 2 reveal how these symbolic errors translate into perceived diagram quality. For the strongest models, structural accuracy (SA), flow accuracy (FA), semantic fidelity (SF), and completeness (C) are all close to 1.0, indicating that the judge model regards the predicted diagrams as near-perfect reconstructions of the gold workflows. `Gemini-2.5-Flash` achieves the highest scores across these metrics, followed closely by `GPT-4.1-mini`. This suggests that on this dataset, the smaller `GPT-4.1-mini` is able to preserve the

---

[11]https://github.com/confident-ai/deepeval

[12]https://huggingface.co/sentence-transformers/all-MiniLM-L12-v2

| Model | Entity P | Entity R | Entity F1 | Rel. P | Rel. R | Rel. F1 | Cosine Sim. |
|---|---|---|---|---|---|---|---|
| GPT-4.1 | 0.975 | 0.974 | 0.975 | 0.963 | 0.960 | 0.961 | 0.875 |
| GPT-4.1-mini | 0.986 | 0.986 | 0.986 | 0.980 | 0.978 | 0.979 | 0.858 |
| GPT-4o | 0.949 | 0.946 | 0.946 | 0.919 | 0.909 | 0.914 | 0.901 |
| GPT-4o-mini | 0.827 | 0.808 | 0.817 | 0.759 | 0.729 | 0.743 | 0.908 |
| Gemini-2.5-Flash | 0.986 | 0.987 | 0.986 | 0.976 | 0.976 | 0.976 | 0.863 |

Table 1: Performance on a FLOWVQA subset (200 diagrams). Precision (P), Recall (R), and F1-scores are reported for both entity and relationship extraction. Cosine similarity reflects semantic similarity between predicted and gold Mermaid code.

global process semantics almost as well as the largest models.

GPT-4o maintains perfect syntax validity (SV = 1.000) but shows slightly lower SA and FA (0.954), mirroring its lower relation F1-score. This pattern is consistent with a model that produces well-formed Mermaid code but sometimes simplifies or omits less salient branches. GPT-4o-mini again stands out: although its syntax validity is almost perfect (0.998), its SA, FA, SF, and C are substantially lower (around 0.83–0.86), confirming that it tends to generate syntactically correct but structurally incomplete diagrams.

| Model | SA | FA | SV | SF | C |
|---|---|---|---|---|---|
| GPT-4.1 | 0.974 | 0.973 | 0.995 | 0.973 | 0.983 |
| GPT-4.1-mini | 0.984 | 0.984 | 0.995 | 0.984 | 0.999 |
| GPT-4o | 0.954 | 0.954 | 1.000 | 0.954 | 0.978 |
| GPT-4o-mini | 0.831 | 0.830 | 0.998 | 0.829 | 0.861 |
| Gemini-2.5-Flash | 0.988 | 0.988 | 0.998 | 0.988 | 0.999 |

Table 2: High-level Mermaid structural metrics on the same subset (SA = structural accuracy; FA = flow accuracy; SV = syntax validity; SF = semantic fidelity; C = completeness).

### 6.1 Discussion

Taken together, the two metric families highlight complementary aspects of model behaviour. High entity and relation F1-score scores show that the best models recover the vast majority of symbolic content, while high SA, FA, SF, and C indicate that this content is arranged into globally coherent, semantically faithful workflows. The small but consistent gap between entity and relation performance confirms that edge reconstruction is more challenging than node labeling, and the discrepancy between cosine similarity and structural metrics for GPT-4o-mini illustrates that text-level similarity can mask important structural errors.

Practically, these results suggest that large multimodal models such as GPT-4.1 and Gemini-2.5-Flash already support near-lossless flowchart-to-Mermaid translation on realistic diagrams, while smaller variants offer a controllable quality–latency trade-off. The complementary nature of symbolic and judge-based metrics is essential: symbolic metrics capture fine-grained correctness, whereas high-level structural metrics reveal whether a diagram remains interpretable as a coherent workflow.

## 7 Conclusion and Future Work

We introduced a lightweight web application that converts flowchart images into structured Mermaid code and supports mixed-initiative refinement through visual editing, code manipulation, and natural-language commands. Using a carefully engineered system prompt together with state-of-the-art VLMs, the system transforms static diagrams into editable, version-controllable representations. Future work includes expanding evaluation beyond the small FLOWVQA subset, extending coverage to additional diagram types such as UML diagrams, exploring fine-tuned models specialized for diagram parsing, and enhancing interaction features such as intelligent autocompletion, collaborative editing, and cross-diagram consistency checking.

### Limitations

Flowchart2Mermaid converts flowchart images into editable Mermaid code through structured prompting and careful post-processing to ensure syntactic validity. However, as with all AI-based systems, the underlying vision–language models may occasionally hallucinate or misinterpret diagram elements. While the generated output typically provides a strong starting point, it may not exactly match the original diagram's layout or semantics. The system is therefore intended as an assistive tool that complements human expertise, requiring user review and refinement to ensure correctness and completeness.

## Acknowledgments

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Confident AI. 2025. Deepeval: The llm evaluation framework. Accessed: 2025-11-17.

Abdul Arbaz, Heng Fan, Junhua Ding, Meikang Qiu, and Yunhe Feng. 2024. Genflowchart: parsing and understanding flowchart using generative ai. In *International Conference on Knowledge Science, Engineering and Management*, pages 99–111. Springer.

Céres Carton, Aurélie Lemaitre, and Bertrand Coüasnon. 2013. Fusion of statistical and structural information for flowchart recognition. In *2013 12th International Conference on Document Analysis and Recognition*, pages 1210–1214. IEEE.

Kewei Cheng, Nesreen K Ahmed, Theodore L Willke, and Yizhou Sun. 2024. Structure guided prompt: Instructing large language model in multi-step reasoning by exploring graph structure of the text. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9407–9430.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Devin Cook. Flowgorithm. http://www.flowgorithm.org/. Accessed: 2025-11-13.

Pritam Deka and Barry Devereux. 2025. Structured extraction from business process diagrams using vision-language models. *Preprint*, arXiv:2511.22448.

Nathan Ensmenger. 2016. The multiple meanings of a flowchart. *Information & Culture*, pages 321–351.

Jiaqi Fang, Zhen Feng, and Bo Cai. 2022. Drawnnet: offline hand-drawn diagram recognition based on keypoint prediction of aggregating geometric characteristics. *Entropy*, 24(3):425.

Jorge-Ivan Herrera-Camara and Tracy Hammond. 2017. Flow2code: from hand-drawn flowcharts to code execution. In *Proceedings of the Symposium on Sketch-Based Interfaces and Modeling*, pages 1–13.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, and 1 others. 2023. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026.

Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2025. From generation to judgment: Opportunities and challenges of llm-as-a-judge. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM computing surveys*, 55(9):1–35.

C David Betancourt Montellano, COFC Garcia, and Roberto Oswaldo Cruz Leija. 2022. Recognition of handwritten flowcharts using convolutional neural networks. *International Journal of Computer Applications*, 184(1):37–41.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

Bernhard Schäfer and Heiner Stuckenschmidt. 2019. Arrow r-cnn for flowchart recognition. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, volume 1, pages 7–13. IEEE.

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*.

Renske Smetsers-Weeda and Sjaak Smetsers. 2017. Problem solving and algorithmic development with flowcharts. In *Proceedings of the 12th workshop on primary and secondary computing education*, pages 25–34.

Sumit Soman, HG Ranjani, Sujoy Roychowdhury, Venkata Dharma Surya Narayana Sastry, Akshat Jain, Pranav Gangrade, and Ayaaz Khan. 2025. A graph-based approach for multi-modal question answering from flowcharts in telecom documents. *arXiv preprint arXiv:2507.22938*.

Chanchai Supaartagorn. 2017. Web application for automatic code generator using a structured flowchart. In *2017 8th IEEE International Conference on Software Engineering and Service Science (ICSESS)*, pages 114–117. IEEE.

Manan Suri, Puneet Mathur, Nedim Lipka, Franck Dernoncourt, Ryan A Rossi, Vivek Gupta, and Dinesh Manocha. 2025. Follow the flow: Fine-grained flowchart attribution with neurosymbolic agents. *arXiv preprint arXiv:2506.01344*.

Knut Sveidqvist and Ashish Jain. 2021. *The official guide to Mermaid. js: create complex diagrams and beautiful flowcharts easily using text and code*. Packt Publishing Ltd.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in neural information processing systems*, 33:5776–5788.

Xiang-Hu Wu, Ming-Cheng Qu, Zhi-Qiang Liu, Jian-Zhong Li, and 1 others. 2011. Research and application of code automatic generation algorithm based on structured flowchart. *Journal of Software Engineering and Applications*, 4(09):534.

Junyi Ye, Ankan Dash, Wenpeng Yin, and Guiling Wang. 2025. Beyond end-to-end vlms: Leveraging intermediate text representations for superior flowchart understanding. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3534–3548.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623.

## A Appendix 1

**LLM-Based Mermaid Flowchart Evaluation Prompt**

```
You are an automated evaluator for Mermaid
    flowcharts.
You will be given two flowcharts: a prediction
    and a gold reference.

### Your tasks:
1. Parse both flowcharts.
   - Extract all nodes: ignore node IDs (like A,
       B, C...), shapes ([], (), {{}}, / /),
       and quote differences (' vs ' vs ").
   - Normalize whitespace and casing (so "Create
       a new user" == "create a new user").
   - Extract all relations (edges) as (source,
       target), ignoring arrow style (-->, -->|
       Yes|, etc.).

2. Compare prediction vs gold.
   - Count TP (correct matches), FP (predicted
       but not in gold), FN (in gold but
       missing in pred).
   - Do this separately for **nodes** and for **
       relations**.

3. Compute metrics:
   - Precision = TP / (TP + FP)   (0 if
       denominator is 0)
   - Recall    = TP / (TP + FN)   (0 if
       denominator is 0)
   - F1-score       = 2*P*R / (P+R)   (0 if
       both P and R are 0)

4. Output strictly:
   - **Only one line of CSV**
   - Exactly 6 numeric values, floats in decimal
       form
   - Order: nodes_precision,nodes_recall,
       nodes_F1-score,rels_precision,
       rels_recall,rels_F1-score
   - No headers, no labels, no text, no
       explanations, no extra lines

### Format example:
0.857,0.750,0.800,0.900,0.818,0.857

### Prediction
{pred}

### Gold
{gold}
```

## B Appendix 2

**LLM-Based High-Level Mermaid Flowchart Evaluation Prompt**

```
You are an automated evaluator for Mermaid
    flowcharts.
You will be given two flowcharts: a prediction
    and a gold reference.

Your job is to output ONE CSV line only, with 5
    numbers (no text, no headers, no
    explanations).

Evaluate ONLY the following metrics:

1. structural_accuracy (0-5)
2. flow_accuracy        (0-5)
3. syntax_validity      (0-2)
4. semantic_fidelity    (0-5)
5. completeness         (0-3)


### VERY IMPORTANT RULE
If you believe the prediction can be **fully
    reconstructed into a valid Mermaid.js
    diagram** that accurately reflects the gold
     reference structure and meaning,
then **assign the maximum score to all five
    metrics**.

### Otherwise:
Score each metric according to the definitions:
- structural_accuracy: correctness of node and
    edge arrangement (0-5)
- flow_accuracy: correctness of overall process
    flow (0-5)
- syntax_validity: valid Mermaid syntax and
    proper graph structure (0-2)
- semantic_fidelity: how well the meaning
```

```
      matches the gold (0-5)
- completeness: how fully the prediction covers
      the gold (0-3)

### OUTPUT FORMAT (STRICT)
Output exactly ONE line:
structural_accuracy,flow_accuracy,
      syntax_validity,semantic_fidelity,
      completeness

### Prediction
{pred}

### Gold
{gold}
```