

#### MM-GRADE: A Multi-Modal EDA Tool Documentation QA Framework Leveraging Retrieval Augmented Generation

Yuan Pu<sup>1,2</sup>, Zhuolun He<sup>1,2</sup>, Shutong Lin<sup>1</sup>, Jiajun Qin<sup>1</sup>, Xinyun Zhang<sup>1</sup>, Hairuo Han<sup>1</sup>, Haisheng Zheng<sup>2</sup>, Yuqi Jiang<sup>3</sup>, Cheng Zhuo<sup>3</sup>, Qi Sun<sup>3</sup>, David Pan<sup>4</sup>, **Bei Yu**<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>ChatEDA Tech <sup>3</sup>Zhejiang University <sup>4</sup>University of Texas at Austin

October 29, 2025









#### Outline



Introduction

2 Proposed Algorithm

3 Experimental Results

## Introduction

#### **EDA Tool Documentation QA**



- EDA tool involves powerful but complex functionalities/documents → steep learning curve.
- Multi-Modal Automatic Question-Answering (QA) system:
  - Given the user query, the QA system generates answer referring to the relevant documents.
  - User query/document contain useful visual information:
    - GUI screenshots.
    - Design layouts.
    - Screenshot of error traceback.

In OpenROAD, I don't want pins to be placed in the region highlighted by the green rectangle, what should I do?







#### Retrieved Document:

The `place\_pins` command ... use the `-exclude` option to specify a region where pins cannot be placed ...



The green rectangle is at the bottom of the layout.

You can use the command: place\_pins — exclude bottom:\*

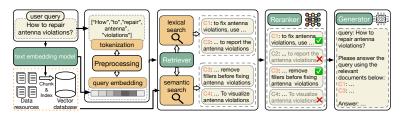
To avoid pins placed at the bottom region.

A typical multi-modal query in OpenROAD.

#### RAG-based EDA QA Systems and Its Limitations



- Trend: Use LLMs and Retrieval Augmented Generation for EDA Tool QA.
- Existing works: RAG-EDA<sup>1</sup>, etc.
- Stages:
  - Retriever: retriever the relevant documents to the user query.
  - Generator: usually LLM, generates the answer.
- Limitations: They are designed for text-only query, ignore visual information →
  fundamental gap in understanding the user's true intent.



Overall flow of RAG-EDA.

<sup>&</sup>lt;sup>1</sup>Yuan Pu et al. (2024). "Customized Retrieval Augmented Generation and Benchmarking for EDA Tool Documentation QA". In: *arXiv* preprint.

#### General-Purpose Multi-Modal RAG Methods: Limitations



- Existing Multi-Modal RAG Methods: VisRAG<sup>2</sup>, EchoSight<sup>3</sup>, etc.
- Limitations of being applied to EDA tool documentation QA:
  - Retriever: General MM-RAG retrievers are trained on natural scenes, charts → struggle to comprehend the complex, specialized information in EDA visuals.
  - Generator: General Vision-Language Models (VLMs) lack EDA domain knowledge → spurious understanding of the user's query and poor quality answers

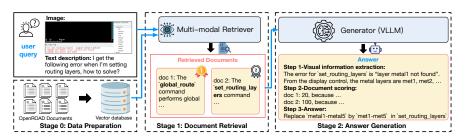
<sup>&</sup>lt;sup>2</sup>Shi Yu et al. (2024). "VisRAG: Vision-based Retrieval-augmented Generation on Multi-modality Documents". In: *arXiv preprint*.

<sup>&</sup>lt;sup>3</sup>Yibin Yan and Weidi Xie (Nov. 2024). "EchoSight: Advancing Visual-Language Models with Wiki Knowledge". In: *Findings of the Association for Computational Linguistics: EMNLP 2024*. Ed. by Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen. Miami, Florida, USA: Association for Computational Linguistics, pp. 1538–1551.

# **Proposed Algorithm**

#### Overall Flow





- Solution: MM-GRADE, a multi-modal, RAG-based EDA tool QA framework.
- Techniques:
  - Multi-modal Retriever: Bi-level negative mining strategy  $\rightarrow$  contrastive learning.
  - Generator (vision large language model): customized 3-stage reasoning pipeline: Extract-Score-Answer (ESA) pipeline.
- Benchmark: ORD-MMBench, consisting of 120 QAs (from OpenROAD) with screenshots.

#### Existing Training Strategy of Multi-Modal Retriever



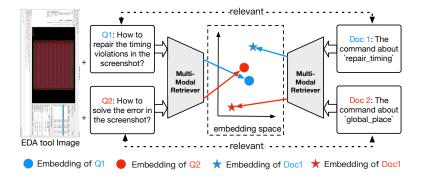
- Mechanism of multi-modal retriever for document retrieval:
  - Step 1: project the user query and all document chunks (multi-modality) into a unified embedding space.
  - Step 2: retrieve the most relevant document chunks by similarity matching.
- Training strategy: Contrastive learning with document-level hard negative mining:
  - Given a user query q.
  - positive samples: its relevant documents  $d_i^+$
  - negative samples: its irrelevant documents  $d_i^-$
- Training loss:

$$L_d^i = -\log \frac{e^{f(q_i^{\text{img}} || q_i^{\text{txt}})^{\top} f(d_i^+)/\tau}}{e^{f(q_i^{\text{img}} || q_i^{\text{txt}})^{\top} f(d_i^+)/\tau} + \sum_{j=1}^n e^{f(q_i^{\text{img}} || q_i^{\text{txt}})^{\top} f(d_{i,j}^-)/\tau}},$$
(1)

#### Limitations of Multi-Modal Retriever



- Observation: One EDA tool screenshot contains complex and diverse information → triggers multiple, distinct user questions.
- Consequence: Many different EDA tool user questions share one common image (screenshot).
- Limitation: Exiting general-purpose Multi-Modal retrievers focus on visual information →
  distinct user queries (with the same image) are projected closely in the embedding space →
  Low document retrieval accuracy.
- Example:



# Solution of EDA-Customized Multi-Modal Retriever: Bi-level Hard Negative Mining

- Purpose: train the multi-modal retriever such that → distinct queries with the same image → projected **distinctly** in the embedding space → higher retrieval accuracy.
- Solution: contrastive learning with query-level hard negative mining:
  - Given a user query q.
  - positive samples: its relevant documents  $d_i^+$
  - negative samples: distinct queries with **the same image**  $\{(q_i^{\text{img}}, q_{i,k}^{\text{txt}}) | k \in [1, m] \land k \neq j\}$
- Training loss:

$$L_q^i = -\sum_{j=1}^m \log \frac{e^{f(q_i^{\text{img}} || q_{i,j}^{\text{txt}})^T f(d_{i,j}^+)/\tau}}{e^{f(q_i^{\text{img}} || q_{i,j}^{\text{txt}})^T f(d_{i,j}^+)/\tau} + \sum_{k \neq j} e^{f(q_i^{\text{img}} || q_{i,j}^{\text{txt}})^T f(q_i^{\text{img}} || q_{i,k}^{\text{txt}})/\tau}}.$$
(2)

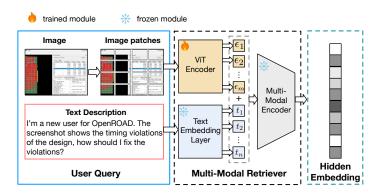
Bi-level hard negative mining (HNM): integrate document-level and query-level HNM:

$$L = \sum_{i=1}^{N} L_d^i + L_q^i. (3)$$

#### Multi-Modal Retriever: Proposed Model Architecture



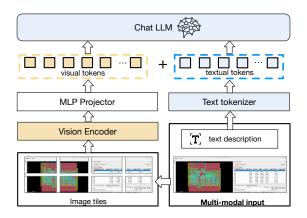
- Components:
  - ViT encoder: encodes the visual input.
  - Text embedding layer: encodes the textual input
  - Multi-modal encoder: text-visual fusion for embedding generation.



#### Multi-Modal Generator: Model Architecture and Limitations



- Model architecture: based on InternVL2<sup>4</sup>.
- Limitations when applied to EDA tool QA: lacks domain-specific inference logic.



<sup>&</sup>lt;sup>4</sup>Zhe Chen et al. (2024). "Internyl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks". In: *Proc. CVPR*, pp. 24185–24198.

#### Solution of Multi-Modal Generator: ESA Pipeline



#### Extract-Score-Answer pipleline:

- Pipeline the answer into 3 stages:
  - Stage 1: Visual information extraction.
  - Stage 2: Document relevance scoring.
  - Stage 3: Answer
- Purpose: used for both training and inference.
- Scoring Criteria of Stage 2, between query q and document d:

Description
d is perfectly relevant to $q$ and can be used to directly answer $q$ .
d is relevant to $q$ and can partially answer $q$ .
d is relevant to $q$ but can not be directly used to answer $q$ .
It is possible that $d$ is relevant with $q$ , but more information is
required to verify the relevance.
d is basically not relevant with $q$ .
d is totally irrelevant with $q$ .

# User visual query-image:

#### User visual query text description:

I don't want the pins to be placed on the region of the green rectangle in the screenshot, what should I do?

#### Retrieved Documents:

Doc 1: `place\_pin`: the command to place a specific pin

Doc 2: use the option 'exclude' of 'place\_pins' ...

#### **Extraction-Scoring-Answer Strategy**

#### Stage 1: Visual information extraction

**Prompt:** Extract the information from the image which can be used to understand and answer the question.

Answer: The pins highlighted in the green rectangle is the bottom region of the layout.

#### Stage 2: Document Relevance Scoring

Prompt: Now given the retrieved documents, analyse and score the relevance between each document and the visual query. The score ranges from 0 to 100, higher scores represent higher relevance.

#### Answer:

Doc 1: 20, because the query is about the command to place all pins, but 'place pin' command only places one pin.

Doc 2: 100, because the document introduces the usage of 'place\_pins' and the 'exclude' flag.

#### Stage 3: Answer

Prompt: answer this user query according the extracted visual information in stage 1 and the document relevance scoring results in stage 2.

#### Answer:

you can use the command: `place\_pins -hor\_layers <hor layers> -ver\_layers <ver layers> -exclude bottom: `to disable pins placed on the bottom region.

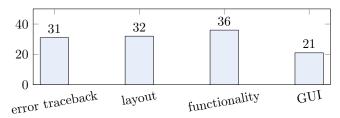
Prompt Example of the ESA Pipeline

# **Experimental Results**

#### Benchmark: ORD-MMBench



- ORD-MMBench <sup>5</sup>: 120 high-quality multi-modal question-document-answer triplets based on the OpenROAD<sup>6</sup> document.
- Four categories of questions:



<sup>5</sup>https://github.com/lesliepy99/MM-GRADE-Benchmarks-ICCAD

<sup>&</sup>lt;sup>6</sup>Tutu Ajayi et al. (2019). "Toward an open-source digital flow: First learnings from the openroad project". In: *Proc. DAC*, pp. 1–4.

#### **Experimental Setting**



- Multi-Modal Retriever:
  - Base model: bge-visualized-m3<sup>7</sup>.
  - Training data: 1075 query-document corpus items generated by gpt-4o.
  - Finetuned by 3000 steps (batch size = 8) on 4 A100 80G GPUs.
- Multi-Modal Generator:
  - Base model: InternVL2-26B<sup>8</sup>.
  - Training data: 2307 high-quality query-document-answer triplets based on OpenROAD<sup>9</sup>.
  - Trained by 2 epochs (batch size = 1) on 8 A100 80G GPUs.

<sup>&</sup>lt;sup>7</sup>https://huggingface.co/BAAI/bge-visualized/blob/main/Visualized\_m3.pth <sup>8</sup>https://huggingface.co/OpenGVLab/InternVL2-26B

<sup>&</sup>lt;sup>9</sup>Tutu Ajayi et al. (2019). "Toward an open-source digital flow: First learnings from the openroad project". In: *Proc. DAC*, pp. 1–4.

#### Experimental Results: Retrieval Baselines



- Metric: recall@k, defined as the proportion of relevant documents that are retrieved among the top k results returned by the retriever/reranker.
- Baselines:
  - Text-only retriever: BGE-M3<sup>10</sup>, RAG-EDA-retriever
  - Multi-modal retriever: DEDR<sup>11</sup>, BLIP<sup>12</sup>, VisRAG-retriever<sup>13</sup>, EchoSight-retriever<sup>14</sup>, bge-visualized-m3

<sup>10</sup> Jianly Chen et al. (2024). BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation. arXiv: 2402.03216 [cs.CL].

<sup>11</sup>Alireza Salemi, Juan Altmayer Pizzorno, and Hamed Zamani (2023). "A symmetric dual encoding dense retrieval framework for knowledge-intensive visual question answering". In: *Proc. SIGIR*, pp. 110–120.

<sup>12</sup>Junnan Li et al. (2022). "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation". In: *Proc. ICML*. PMLR, pp. 12888–12900.

<sup>13</sup>Shi Yu et al. (2024). "VisRAG: Vision-based Retrieval-augmented Generation on Multi-modality Documents". In: *arXiv preprint*.

<sup>14</sup>Yibin Yan and Weidi Xie (Nov. 2024). "EchoSight: Advancing Visual-Language Models with Wiki Knowledge". In: *Findings of the Association for Computational Linguistics: EMNLP 2024*. Ed. by Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen. Miami, Florida, USA: Association for Computational Linguistics, pp. 1538–1551.

#### Experimental Result: Multi-Modal Retrieval



- Ablation study: train MM-GRADE retriever (ours) with only query-level/document-level hard negative mining (HNM).
- MM-GRADE retriever with bilevel HNM achives the best performance.

Model type	recall@1	recall@2	recall@3	recall@4	recall@5
BGE-M3	0.590	0.694	0.731	0.761	0.776
RAG-EDA-retriever	0.530	0.612	0.679	0.724	0.739
DEDR	0.112	0.187	0.231	0.284	0.299
BLIP	0.127	0.209	0.276	0.299	0.313
VisRAG-retriever	0.025	0.042	0.050	0.083	0.108
EchoSight-retriever	0.383	0.492	0.583	0.625	0.667
VISTA: bge-visualized-m3	0.582	0.716	0.739	0.761	0.799
MM-GRADE-retriever w. d-level HNM	0.664	0.784	0.806	0.828	0.851
MM-GRADE-retriever w. q-level HNM	0.582	0.694	0.746	0.761	0.828
MM-GRADE retriever (ours)	0.679	0.828	0.843	0.866	0.896

Table: Performance of document retrieval on ORD-MMBench.

#### Experimental Result: Overall RAG Flow



- Metrics:
  - LLM-Score: gpt-40 to score each answer (from 0 to 100) referring to a scoring criteria <sup>15</sup>.
  - BLEU and Rouge-L.
- RAG flow baselines:
  - Text-only: RAG-EDA.
  - Multi-modal: VisRAG, EchoSight, VISTA-RAG.

RAG Flow	ORD-MMBench▶error traceback			ORD-MMBench▶layout		ORD-MMBench▶functionality			ORD-MMBench▶GUI			ORD-MMBench▶all			
KAG Flow	LLM-Score	BLEU	ROUGE-L	LLM-Score	BLEU	ROUGE-L	LLM-Score	BLEU	ROUGE-L	LLM-Score	BLEU	ROUGE-L	LLM-Score	BLEU	ROUGE-L
VisRAG	41.129	0.080	0.223	12.500	0.098	0.242	18.750	0.111	0.247	4.762	0.024	0.162	20.417	0.084	0.225
EchoSight	51.613	0.087	0.241	54.688	0.165	0.300	62.500	0.191	0.322	64.286	0.091	0.214	57.917	0.139	0.276
RAG-EDA	22.419	0.064	0.187	61.719	0.156	0.276	59.722	0.190	0.332	51.190	0.107	0.267	49.125	0.134	0.268
VISTA-RAG	65.323	0.126	0.266	68.750	0.216	0.348	71.528	0.232	0.368	73.810	0.116	0.261	69.583	0.180	0.317
MM-GRADE w/o ESA	71.129	0.159	0.303	78.125	0.273	0.418	75.833	0.291	0.442	90.476	0.160	0.312	77.792	0.229	0.377
MM-GRADE (ours)	82.258	0.244	0.391	80.469	0.307	0.468	80.000	0.294	0.450	83.333	0.178	0.323	81.292	0.264	0.417

Table: Performance of the multi-modal RAG flows on ORD-MMBench.

<sup>15</sup>https://github.com/lesliepv99/MM-GRADE-Benchmarks-ICCAD

## **THANK YOU!**