
Multimodal Chip Physical Design Engineer Assistant

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Modern chip physical design relies heavily on Electronic Design Automation
2 (EDA) tools, which often struggle to provide interpretable feedback or actionable
3 guidance for improving routing congestion. In this work, we introduce a Multi-
4 modal Large Language Model Assistant (MLLMA) that bridges this gap by not
5 only predicting congestion but also delivering human-interpretable design sugges-
6 tions. Our method combines automated feature generation through MLLM-guided
7 genetic prompting with an interpretable preference learning framework that models
8 congestion-relevant tradeoffs across visual, tabular, and textual inputs. We compile
9 these insights into a "Design Suggestion Deck" that surfaces the most influential
10 layout features and proposes targeted optimizations. Experiments on the CircuitNet
11 benchmark demonstrate that our approach outperforms existing models on both
12 accuracy and explainability. Additionally, our design suggestion guidance case
13 study and qualitative analyses confirm that the learned preferences align with real-
14 world design principles and are actionable for engineers. This work highlights the
15 potential of MLLMs as interactive assistants for interpretable and context-aware
16 physical design optimization.

1 1 Introduction

18 The field of chip physical design encompasses critical prediction tasks such as congestion, timing,
19 Design Rule Check (DRC), and IR Drop. These tasks, traditionally tackled using Electronic Design
20 Automation (EDA) tools, often suffer from slow processing speeds and limited ability to provide
21 actionable guidance for design improvement. Existing approaches using image-to-image translation
22 models focus primarily on prediction accuracy but lack interpretability and the ability to suggest
23 corrective measures. This paper explores the potential of Multimodal Large Language Models
24 (MLLMs) in providing interpretable predictions along with actionable suggestions to assist engineers
25 in optimizing their designs. Recent advances in foundation models and instruction-tuned Large
26 Language Models (LLMs) have led to impressive progress in tasks that require reasoning, code
27 generation, and interactive assistance. However, their application to chip physical design remains
28 largely underexplored. In particular, current EDA solutions lack the capability to provide real-time,
29 explainable, and multimodal assistance that integrates layout visuals, tabular metrics, and design
30 constraints into a unified decision-making process. We argue that a Multimodal Language Model
31 Assistant (MLLMA) tailored to chip physical design can fill this gap.

32 The motivation for this research is to address the gap between prediction accuracy and interpretability
33 in physical design tasks. While current models are adept at predicting various metrics such as
34 congestion and DRC violations, they fail to explain the underlying causes or offer corrective measures.
35 By leveraging multimodal LLMs, we aim to move beyond mere prediction and offer a framework that
36 provides both predictions and accompanying explanations. These explanations can be compiled into
37 a "Design Suggestion Deck" which assists engineers by presenting comprehensible, human-readable
38 reasoning to guide their design decisions. The primary research challenge is developing an MLLM-

39 based agent framework capable of processing multimodal data inputs, including geometric images,
40 tabular features, and circuit graphs, and providing interpretable predictions with actionable design
41 suggestions. Furthermore, we aim to compare this approach against baseline models to demonstrate
42 its effectiveness in improving design quality and interpretability.

43 Our method consists of two stages: (1) Automate Feature Engineering and Generation, and (2)
44 Interpretable Preferences Learning to Design Suggestion Deck. In the first stage, a Multimodal Large
45 Language Model (MLLM) agent generates and extracts numeric features from Macro Region, RUDY,
46 and RUDY pin images and text-based features from configuration and logs. Inspired by the Genetic
47 Instruct (Majumdar et al., 2024), a genetic algorithm incorporating mutation and crossover is applied
48 to expand the multimodal feature pool, with a deduplication process ensuring uniqueness. Feature
49 importance and cross validation scores are evaluated using a Random Forest model to guide feature
50 selection and refinement. The second stage uses these features and their importance scores to create
51 a Design Suggestion Deck. The MLLM agent translates important features into interpretable rules,
52 providing actionable guidance for chip design optimization.

53 Our contributions are summarized as follows:

- 54 • Development of a multimodal LLM framework capable of handling geometric image
55 features, tabular data, and circuit graph inputs for prediction tasks in physical design.
- 56 • Design of an interpretable preference learning approach that generates interpretable predic-
57 tions and actionable suggestions compiled into a "Design Suggestion Deck".
- 58 • Evaluation of the proposed framework against the baseline models that demonstrate strong
59 improvements in design quality and efficiency.

60 2 Related Work

61 2.1 Multimodal Language Model Assistant

62 Recent works such as Miphia (Zhu et al., 2024), MM-ReAct (Yang et al., 2023), and MultiModal-
63 GPT (Gong et al., 2023) demonstrate that aligning visual inputs with language reasoning modules
64 allows models to solve complex real-world tasks. Miphia, for example, focuses on efficiency by en-
65 abling multi-turn multimodal dialogue with domain grounding, while MM-ReAct leverages reasoning
66 traces and visual memory for tool-augmented perception.

67 In the context of chip design, prior efforts such as Chip-Chat (Blocklove et al., 2023) investigate
68 interactive co-design with conversational LLMs for hardware description generation, while Chip-
69 NeMo (Liu et al., 2023) introduces domain-adapted LLMs for industrial chip tasks including EDA
70 script generation and bug analysis. Unlike these works, which primarily focus on language-to-code
71 translation or domain adaptation, our approach leverages MLLMs for spatially grounded reasoning
72 over fine-grained layout features such as routing demand and macro placement. Additionally, we
73 explicitly model design tradeoffs through learned interpretable preferences, extending the role of
74 LLMs from passive code generators to interactive design assistants capable of justification and
75 feedback.

76 2.2 Congestion Map Prediction

77 ClusterNet (Min et al., 2023) focuses on predicting routing congestion caused by netlist topology
78 using netlist clustering and Graph Neural Networks (GNNs). By employing the Leiden algorithm
79 to generate cohesive clusters and developing a GNN-based model to generate cluster embeddings,
80 ClusterNet achieves improved prediction performance and congestion optimization without requiring
81 placement and routing (P&R) tools. CircuitNet (Chai et al., 2023; Jiang et al., 2023) provides an
82 open-source dataset for machine learning applications in VLSI CAD, facilitating the development
83 and benchmarking of ML models for tasks like congestion prediction and design rule check (DRC)
84 violation prediction. The dataset includes diverse samples collected from various chip designs and
85 supports comprehensive data, including routability, timing, and power metrics. CircuitFormer (Zou
86 et al., 2023) introduces a novel approach by treating circuit components as point clouds and utilizing
87 Transformer-based point cloud perception methods for feature extraction. This method enables
88 direct feature extraction from raw data without any preprocessing, allows for end-to-end training,
89 and achieves state-of-the-art performance in congestion prediction tasks on both the CircuitNet

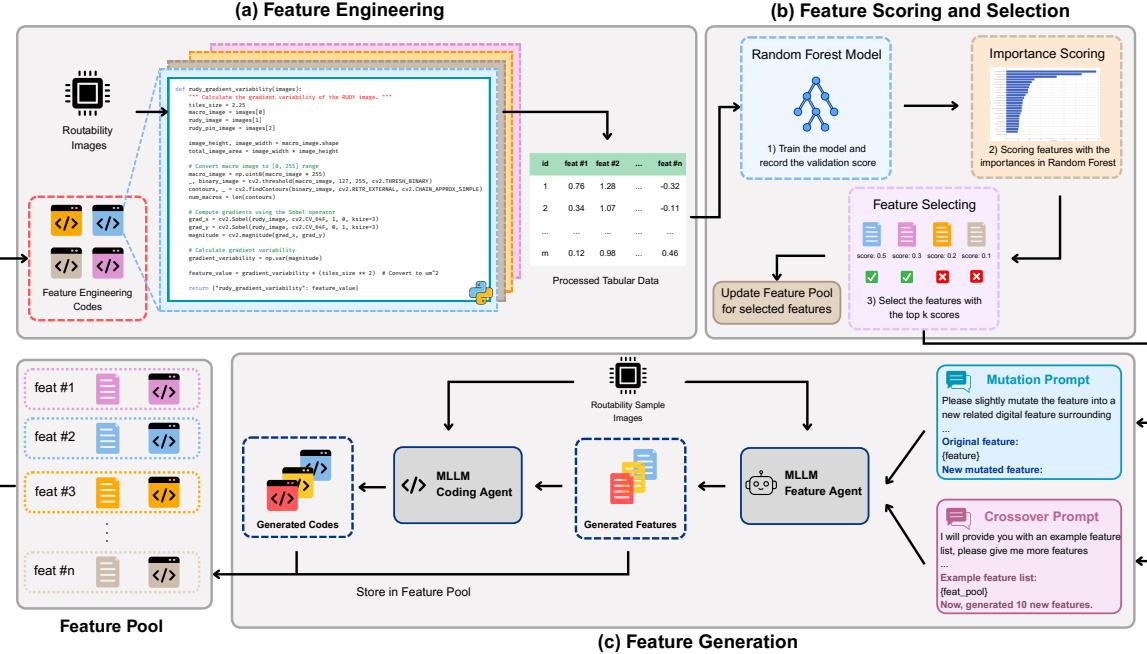


Figure 1: Overview of the Genetic Instruct framework. (a) Spatial metrics are extracted from routability images using domain-specific rules. (b) A Random Forest ranks and selects key features. (c) MLLM agents iteratively generate and mutate features via code, enriching the feature pool to improve downstream performance.

and ISPD2015 datasets, as well as in DRC violation prediction tasks on the CircuitNet dataset. MPGDI (Yang et al., 2024) proposes Mini-Pixel Batch Gradient Descent, a plug-and-play optimization algorithm that focuses on the most informative entries in the congestion map. By integrating this method into predictive AI models for physical design flows, MPGDI enhances congestion prediction accuracy and efficiency, demonstrating its effectiveness in optimizing chip design processes.

3 Methodology

Our proposed MLLM assistant consists of two core components: (1) automated feature generation and engineering, and (2) interpretable preference modeling for generating actionable design suggestions.

3.1 Automated Feature Generation and Engineering

Figure 1 illustrates our pipeline for automated feature generation and engineering. The approach integrates domain-driven initial features, MLLM-guided mutation, and a feedback loop to iteratively expand and refine a pool of predictive and interpretable features.

3.1.1 Initial Feature Types

We begin with two categories of raw features for routing congestion prediction:

Image-based layout features: *Macro Region* encodes macro and cell placement for global layout context. *RUDY (Routing Demand)* highlights routing density hotspots. *RUDY Pin Images* capture local pin density patterns affecting routing complexity.

Text-based configuration features: Logs from EDA tools include placement constraints, blockages, net priorities, and timing directives. These symbolic features complement spatial representations with interpretable design signals.

110 **3.1.2 Genetic Instruct Feature Generation**

111 We employ a Genetic-Instruct framework (Majumdar et al., 2024) to iteratively evolve and expand fea-
112 ture sets using MLLMs. The process consists of initialization, evaluation, variation, and deduplication
113 steps:

114 **Initialization:** The feature pool is seeded with handcrafted features that map visually and semantically
115 to congestion causes. Each feature includes a human-readable name and description. MLLMs are
116 prompted to generate extraction code, producing tabular data for training. Congestion severity is
117 labeled using the average congestion value in the 20 most-congested grid cells of each sample.

118 **Evaluation and Selection:** A Random Forest model is trained on the extracted features and labels.
119 The model's feature importance scores act as the fitness function, guiding the pruning of low-impact
120 features. We retain the top $k=20$ candidates for the next generation.

121 **Crossover and Mutation:** New features are created by prompting MLLMs with few-shot examples
122 drawn from high-ranking features (crossover), or by mutating individual features with probabilities
123 based on a rank-sensitive version of the Ebbinghaus forgetting curve. Specifically, mutation proba-
124 bility is defined as: $M = e^{-\frac{r}{N}}$, where r is the feature's rank and N is the total number of features.
125 Higher-ranked features have greater mutation probability, encouraging diversity while maintaining
126 quality.

127 **Deduplication:** To avoid redundancy, newly proposed features are compared against the existing pool
128 using a Deduplicator-MLLM. It evaluates semantic similarity in descriptions and removes duplicates
129 with reasoned justifications.

130 This process iteratively refines the feature space, producing a diverse, interpretable, and model-
131 relevant set of features for downstream congestion prediction. Some generated feature example is
132 shown in Figure 2 and the complete list of generated feature pool is in Appendix F.

Feature Pool Examples

```
{macro_compactness_index: "a measure of how closely packed the macros are, potentially affecting routing paths and congestion",
clustered_macro_distance_std: "the standard deviation of distances between clustered groups of macros",
rudy_pin_clustering_coefficient: "a measure of how many rudy pins cluster together relative to the total number of rudy pins",
macro_density_gradient: "the change in macro density across the layout, impacting local congestion", ...}
```

Figure 2: Feature Pool Examples

133 **3.1.3 Automated Feature Engineering**

134 We employ Multimodal Large Language Models (MLLMs) to automate feature transformation by
135 generating code snippets tailored to the characteristics of input features and target objectives. The
136 generated code is then executed in a controlled environment, where inputs, outputs, and errors are
137 systematically handled to ensure robustness. This automation streamlines the feature engineering
138 process, improving both efficiency and scalability in modeling routing congestion.

139 **3.2 Interpretable Preferences to Design Suggestions**

140 To bridge the gap between predictive modeling and actionable design insights, we develop a method-
141 ology that directly uses the features discovered in Phase 1 as interpretable preference indicators,
142 enabling their translation into concrete and actionable design suggestions.

143 **3.2.1 Interpretable Preferences and Model Architecture**

144 Inspired by the ArmoRM framework (Wang et al., 2024), we formulate design evaluation as a multi-
145 objective task, addressing both scalar congestion prediction and spatial congestion map generation.
146 Central to our approach is the reuse of the feature pool generated in Phase 1 via genetic instruction
147 prompting. These features serve as interpretable preference indicators—each representing a physically

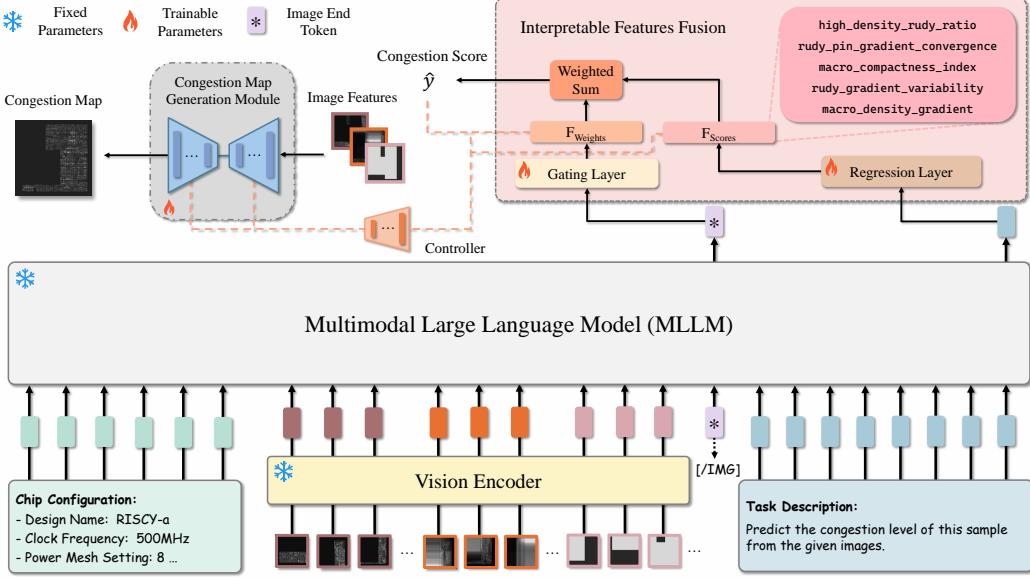


Figure 3: Overview of our interpretable preference modeling framework. The model leverages interpretable features identified in Phase 1 as preference signals to guide congestion prediction and design suggestion generation. A multimodal architecture integrates layout images and task prompts, predicts feature values, learns objective weights via a gating mechanism, and generates both scalar congestion scores and spatial congestion maps. This enables transparent and actionable feedback for physical design refinement.

148 meaningful and actionable design attribute such as pin clustering, macro placement density, or region-based congestion asymmetry.

150 Our model adopts a multimodal LLM backbone based on MiniCPM (Yao et al., 2024). Design
 151 configuration images are tokenized as vision inputs, while the task descriptions are encoded as text
 152 tokens. We then apply a regression layer trained on the last tokens to predict interpretable feature
 153 values, and a gating layer trained on vision tokens to learn objective weights based on layout context.
 154 The gating layer dynamically assigns importance to each interpretable feature dimension, effectively
 155 scalarizing the multi-objective signal into a preference-aware score. This architecture enables both
 156 granular score decomposition and end-to-end learning across objectives. The model not only predicts
 157 congestion values but also generates pixel-level congestion maps, providing actionable, interpretable
 158 guidance to support iterative physical design refinement.

159 In the congestion map generation stage, a conditional U-Net architecture is employed for the image-
 160 to-image translation task, where the input consists of image-based layout features and the output is
 161 the predicted congestion map. To incorporate the interpretable features as conditional information,
 162 we explore two representation methods: numerical vectors and text-based descriptions. Through the
 163 ablation studies in Appendix D, we find that using interpretable features with gating weights in a
 164 text-based representation achieves better performance, demonstrating their effectiveness in providing
 165 semantic information to guide the generation of congestion maps.

166 3.2.2 Design Suggestions

167 Using the interpretable preference signals, we generate targeted design suggestions focused on
 168 improving the most critical objectives. Since the model reasons over features that are both informative
 169 and physically grounded (e.g., pin density imbalance, macro boundary interactions), its suggestions
 170 are inherently actionable. For example, in cases where congestion is identified as a dominant issue, the
 171 model might recommend reducing cluster density in specific regions or modifying routing constraints.
 172 Each suggestion is directly traceable to feature-level insights, ensuring that design decisions are not
 173 only effective but also well-justified.

174 4 Experiment

175 Details of the experimental setup and hyperparameters are provided in Appendix E.

176 4.1 Dataset

177 The CircuitNet dataset (Jiang et al., 2023; Chai et al., 2023) serves as the foundation for this study. It
 178 is a large-scale, open-source benchmark comprising over 10,000 synthesized chip designs—including
 179 CPUs, GPUs, and AI accelerators—fabricated using the 14nm FinFET process. It provides two
 180 types of features: image-like layout representations and text-based features that extracted from
 181 configuration log. Details on preprocessing and data augmentation are included in Appendix E.

182 4.2 Baselines

183 We compare our method against three state-of-the-art congestion prediction models. **GPDL** (Chai
 184 et al., 2023) is a CNN-based architecture that predicts congestion maps from layout images us-
 185 ing spatial convolutions to model local routing patterns. **CircuitFormer** (Zou et al., 2023) is a
 186 Transformer-based model operating on point cloud data, learning spatial relationships from raw
 187 layout coordinates. **MPGD** (Yang et al., 2024) introduces Mini-Pixel Batch Gradient Descent that
 188 targets the most informative regions of the congestion map.

189 4.3 Evaluation Metrics

190 Each model is trained to output a pixel-level prediction of congestion, which is then compared
 191 against the ground-truth congestion heatmap. We adopt standard metrics used in prior work to
 192 evaluate congestion prediction quality. These include SSIM for perceptual similarity, NRMSE and
 193 PeakNRMSE for error measurement, and PLCC, SRCC, and KRCC for correlation analysis. A
 194 detailed explanation and formal definitions of each metric are provided in Appendix B.

195 4.4 Experiment Results

196 Table 1 shows that our model consistently outperforms baselines in SSIM, NRMSE, and peak error
 197 metrics, validating the benefit of interpretable preference modeling and multi-objective supervi-
 198 sion. Our qualitative results in Figure 4 further show that the predictions from our model more accurately
 199 capture fine-grained congestion boundaries compared to the baselines.

Table 1: Congestion prediction performance on CircuitNet across pixel-based and correlation-based metrics.

Method	SSIM ↑	NRMSE ↓	Peak NRMSE ↓					Correlation Metrics		
			0.5%	1%	2%	5%	avg	PLCC ↑	SRCC ↑	KRCC ↑
GPDL (Chai et al., 2023)	0.773	0.047	0.441	0.323	0.236	0.155	0.289	0.5032	0.5143	0.3787
MPGD (Yang et al., 2024)	0.787	0.046	0.271	0.217	0.171	0.121	0.195	—	—	—
CircuitGNN	—	—	—	—	—	—	—	0.3287	0.4483	0.3688
CircuitFormer (Zou et al., 2023)	—	—	—	—	—	—	—	0.6374	0.5282	0.3935
Ours	0.807	0.045	0.263	0.205	0.169	0.118	0.189	0.6452	0.5233	0.3941

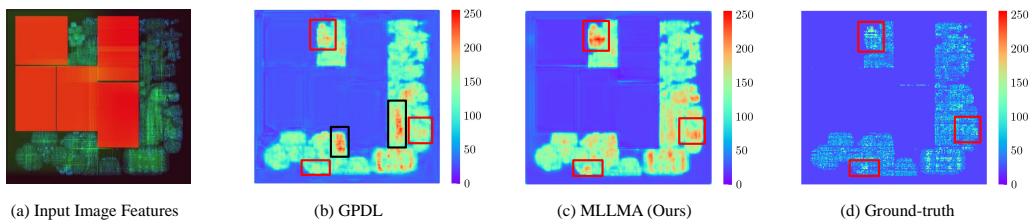


Figure 4: Qualitative comparison of congestion prediction heatmaps between the GPDL baseline and our proposed multimodal model, MLLMA. Red bounding boxes highlight congestion hotspots where MLLMA more accurately captures the spatial patterns and severity of congestion compared to the ground truth. Black bounding boxes indicate regions where GPDL fails to localize congestion accurately.

Table 2: Ablation Studies.

(a) Ablation on feature generation pipeline.

Pipeline Variant	zero-riscy-a			zero-riscy-b		
	PLCC ↑	SRCC ↑	KRCC ↑	PLCC ↑	SRCC ↑	KRCC ↑
Crossover-Only	0.630	0.674	0.492	0.589	0.686	0.493
Mutation-Only	0.353	0.485	0.348	0.469	0.462	0.314
Genetic Instruct	0.705	0.786	0.595	0.599	0.591	0.412

(b) Ablation on interpretable features.

Feature Source Modality	SSIM ↑	NRMSE ↓
Image-Based	0.791	0.047
Log-Based	0.764	0.052
Mixed (Image + Log)	0.789	0.046

200 **4.5 Ablation Study**201 To understand which components contribute most to our model’s effectiveness, we conduct a two-
202 phase ablation study.203 **Phase 1: Ablation Study on Feature Generation Pipeline** This phase evaluates the necessity and
204 effectiveness of our proposed **Genetic Instruct prompting strategy** for MLLMs. We compare it
205 against two variants: a crossover-only version that recombines prompt elements and a mutation-only
206 version that introduces random prompt variations. Our Genetic Instruct method integrates both
207 crossover and mutation with evolutionary selection, and achieves significantly better feature quality
208 for downstream congestion prediction. As shown in Table 5a, the Genetic Instruct strategy leads
209 to substantial performance gains, demonstrating that careful prompt optimization is crucial when
210 leveraging MLLMs for physical design tasks.211 **Phase 2: Ablation Study on Interpretable Preference Pipeline** This phase investigates which
212 components of our interpretable preference modeling pipeline are most critical to performance. We
213 ablate four key aspects: the input modalities, the loss functions, the congestion map generation
214 modules, and the interpretable preference features used for gating mechanism. Each removal leads
215 to varying degrees of performance degradation, highlighting the importance of both structural and
216 semantic components in learning meaningful preferences. A subset of the ablation results is sum-
217 marized in Table 2, with additional details provided in Appendix D. These findings confirm that
218 the learned preference signals are not only interpretable but also integral to model performance.
219 Notably, removing structured inputs or disabling the gating mechanism leads to substantial degra-
220 dation, underscoring the critical role of each component in modeling contextual and actionable design
221 preferences.222 **5 Analysis**223 **5.1 Feature Attribution Analysis**224 Each design feature is associated with a scalar gating value representing its relative importance or pref-
225 erence weight. To interpret the sign of these values: These interpretations allow us to understand how
226 each feature drives model predictions in different design contexts. We observe that low-congestion
227 samples tend to exhibit high values on the `macro_rudy_boundary_interaction_index` feature,
228 as illustrated in Table 3. This trend is consistent across multiple instances, with a negative Pearson
229 correlation of -0.72 between this feature and the congestion score.230 Similarly, features such as `rudy_pin_clustering_coefficient` and `rudy_pin_compaction`
231 `_ratio` tend to have high values in high-congestion samples, showing a positive correlation. This
232 suggests that pin-related density features are dominant contributors to congestion under high-load
233 scenarios.234 To further support interpretability, we conduct a detailed breakdown analysis covering the value
235 distribution, gating value, and associated congestion level for every major design feature. This
236 comprehensive analysis is included in Appendix C and Figure 7, where we also include additional
237 observations and insights derived from these patterns.238 **5.2 Design Suggestion Guidance Leads To Lower Congestion**239 To evaluate the practical utility of the model’s interpretable outputs, we conduct a case study that
240 demonstrates how applying model-guided feature adjustments can lead to a substantial reduction

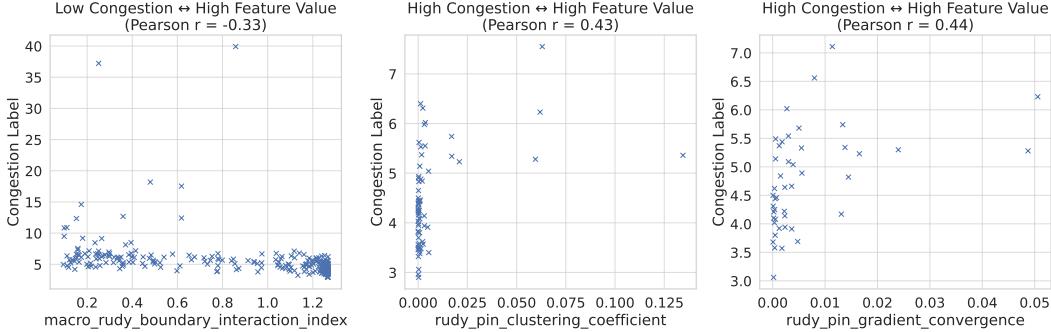


Table 3: Scatter plots showing the relationship between key layout features and congestion labels. (Left) `macro_rudy_boundary_interaction_index` is negatively correlated with congestion ($r = -0.33$), indicating that low-congestion samples tend to exhibit higher values in this macro-level feature. (Middle) `rudy_pin_clustering_coefficient` shows a moderate positive correlation with congestion ($r = 0.43$), suggesting that dense pin clusters contribute to congestion under high-load conditions. (Right) `rudy_pin_gradient_convergence` similarly displays a positive trend ($r = 0.44$), further highlighting the influence of pin-level spatial patterns on congestion outcomes.

241 in routing congestion. Specifically, we use the model’s top-attributed features to inform actionable
 242 design changes and compare outcomes before and after applying these suggestions. Figure 5 presents
 243 one representative example. The left side shows a high-congestion design sample, along with the top-
 244 5 features contributing most to the model’s prediction. Based on these attributions, we systematically
 245 adjust feature values—particularly those with strong positive influence on congestion—toward more
 246 favorable levels. The adjusted design, shown on the right, exhibits a notable improvement.

247 This case highlights the effectiveness of the model’s interpretable reasoning. The ability to trace
 248 predictions back to key design features enables not just post hoc understanding, but also actionable
 249 guidance. By following these learned preferences, designers can reduce congestion and iteratively
 250 improve design quality in practice.

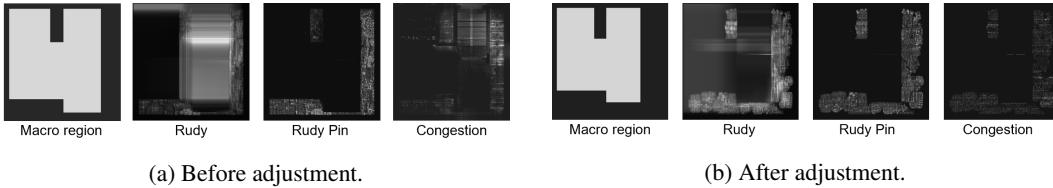


Figure 5: By following the top model-attributed features and adjusting actionable design parameters, the resulting design exhibits significantly lower predicted and actual congestion.

251 **Feature Attribution and Congestion Change Summary:**

Table 4: Top features before/after adjustment and their impact on congestion.

Feature	Before Value	After Value	Attribution (Before)	Congestion Impact
<code>rudy_pin_gradient_convergence</code>	0.3351	0.1202	0.1551 → 0.0001	High
<code>rudy_pin_clustering_coefficient</code>	0.5913	0.3606	0.1409 → 0.0002	↓
<code>rudy_pin_compaction_ratio</code>	0.4273	0.3035	0.1216 → 0.0001	↓
<code>high_density_rudy_pin_ratio</code>	0.3138	0.4867	0.0638 → 0.0003	Mixed
<code>demarcated_macro_proximity_index</code>	0.3765	-0.0030	0.0539 → 0.0000	↓
<code>macro_rudy_boundary_interaction_index</code>	0.0849	0.0721	0.0321 → 1.2655	High post-adjustment role

252 This case illustrates how interpreting the gating values and using them to guide feature adjustment
 253 can lead to measurable congestion reduction, demonstrating the interpretability and actionability of
 254 our model’s learned preferences.

255 We compute the consistency rate across a large number of matched sample pairs. A high agreement
 256 between the sign of the gating value and the direction of the label difference indicates that the learned

257 preferences are interpretable and aligned with the optimization objective. Additionally, we collaborate
258 with experienced chip designers and researchers to manually review selected cases, further validating
259 that the gating values reflect meaningful and actionable design insights.

260 5.3 Design Suggestion Qualitative Examples

We complement our quantitative and case-based evaluation with a set of qualitative examples that further illustrate how the model adapts its reasoning to context-specific designs. In Figure 6, we compare two designs—Design 8277 and Design 8243—exhibiting congestion scores of 40.21 and 20.04, respectively. Despite having similar raw feature values in some cases, the relative importance assigned to each feature varies dramatically between the two layouts. Design 8277 shows elevated importance on multiple pin-related density features, suggesting that the layout suffers from excessive pin compaction and clustering. In contrast, Design 8243 is dominated by a single macro-level feature—`macro_rudy_boundary_interaction_index`—highlighting the positive influence of well-structured macro placement. These examples demonstrate that our model not only identifies which features matter most, but also how their importance is modulated by spatial context. This supports our broader goal of providing interpretable, design-aware guidance that generalizes across different physical design scenarios.

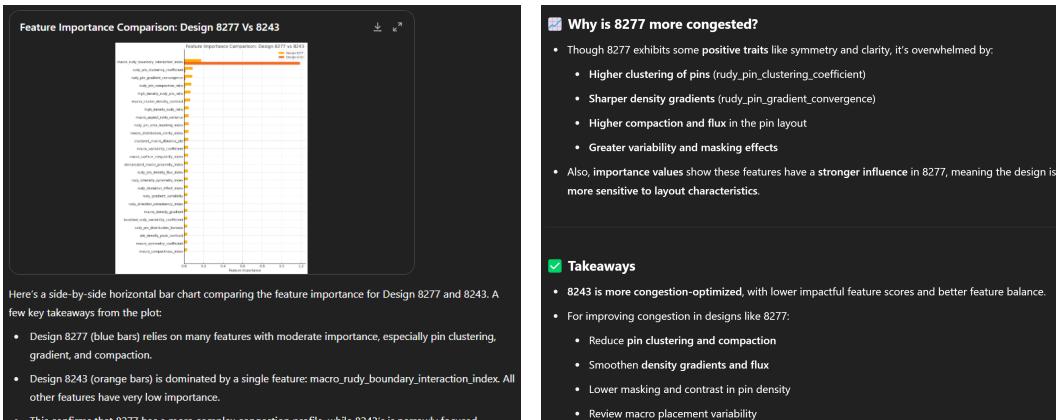


Figure 6: Qualitative interpretation comparing Design 8277 and 8243. Left: Feature importance bar chart. Right: Root causes and actionable takeaways.

273 **6 Limitation**

While our framework shows strong performance, several limitations remain. First, the current evaluation is limited to a RISC-V chip architectures and the CircuitNet dataset. Broader validation across different architectures and more diverse commercial designs is necessary for practical deployment. Second, our focus is confined to congestion prediction; extending the approach to other critical physical design objectives such as power, performance, area (PPA), IR drop, and timing closure remains an open direction.

280 7 Conclusion

We propose a multimodal LLM assistant for chip physical design that predicts congestion and provides interpretable guidance. By combining genetic prompting and preference-based explanation, our method delivers actionable design suggestions aligned with expert intuition. Experiments show state-of-the-art results on CircuitNet, and case studies confirm practical utility. Future work will expand to broader tasks and interactive, real-time, instruction-like design support.

286 **References**

- 287 J. Blocklove, S. Garg, R. Karri, and H. Pearce. Chip-chat: Challenges and opportunities in con-
288 versational hardware design. In *2023 ACM/IEEE 5th Workshop on Machine Learning for CAD*
289 (*MLCAD*), pages 1–6. IEEE, 2023.
- 290 Z. Chai, Y. Zhao, W. Liu, Y. Lin, R. Wang, and R. Huang. Circuitnet: An open-source dataset for
291 machine learning in vlsi cad applications with improved domain-specific evaluation metric and
292 learning strategies. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and*
293 *Systems*, 42(12):5034–5047, 2023.
- 294 T. Gong, C. Lyu, S. Zhang, Y. Wang, M. Zheng, Q. Zhao, K. Liu, W. Zhang, P. Luo, and K. Chen.
295 Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint*
296 *arXiv:2305.04790*, 2023.
- 297 X. Jiang, Y. Zhao, Y. Lin, R. Wang, R. Huang, et al. Circuitnet 2.0: An advanced dataset for promoting
298 machine learning innovations in realistic chip design environment. In *The Twelfth International*
299 *Conference on Learning Representations*, 2023.
- 300 M. Liu, T.-D. Ene, R. Kirby, C. Cheng, N. Pinckney, R. Liang, J. Alben, H. Anand, S. Baner-
301 jee, I. Bayraktaroglu, et al. Chipnemo: Domain-adapted llms for chip design. *arXiv preprint*
302 *arXiv:2311.00176*, 2023.
- 303 S. Majumdar, V. Noroozi, S. Narendhiran, A. Ficek, J. Balam, and B. Ginsburg. Genetic instruct:
304 Scaling up synthetic generation of coding instructions for large language models. *arXiv preprint*
305 *arXiv:2407.21077*, 2024.
- 306 K. Min, S. Kwon, S.-Y. Lee, D. Kim, S. Park, and S. Kang. Clusternet: Routing congestion
307 prediction and optimization using netlist clustering and graph neural networks. In *2023 IEEE/ACM*
308 *International Conference on Computer Aided Design (ICCAD)*, pages 1–9. IEEE, 2023.
- 309 H. Wang, W. Xiong, T. Xie, H. Zhao, and T. Zhang. Interpretable preferences via multi-objective
310 reward modeling and mixture-of-experts. In *Findings of the Association for Computational*
311 *Linguistics: EMNLP 2024*, pages 10582–10592, 2024.
- 312 H. Yang, A. Agnesina, and H. Ren. Optimizing predictive ai in physical design flows with mini
313 pixel batch gradient descent. In *Proceedings of the 2024 ACM/IEEE International Symposium on*
314 *Machine Learning for CAD*, pages 1–7, 2024.
- 315 Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, and L. Wang. Mm-
316 react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*,
317 2023.
- 318 Y. Yao, T. Yu, A. Zhang, C. Wang, J. Cui, H. Zhu, T. Cai, H. Li, W. Zhao, Z. He, et al. Minicpm-v: A
319 gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.
- 320 M. Zhu, Y. Zhu, X. Liu, N. Liu, Z. Xu, C. Shen, Y. Peng, Z. Ou, F. Feng, and J. Tang. Mipha:
321 A comprehensive overhaul of multimodal assistant with small language models. *arXiv preprint*
322 *arXiv:2403.06199*, 2024.
- 323 J. Zou, X. Wang, J. Guo, W. Liu, Q. Zhang, and C. Huang. Circuit as set of points. In *Thirty-seventh*
324 *Conference on Neural Information Processing Systems*, 2023.

Algorithm 1 Automated Feature Generation and Feature Engineering

Input: N : Number of automated iterations
 k : Maximum size of feature pool
 \mathcal{F}_{init} : Initial set of hand-crafted features
 \mathcal{D} : Raw images of training data
 \mathcal{L} : Congestion levels of training data
 I_s : Sample images for the inputs of MLLMs
 P_m : Probability of selecting mutation operation
 P_{op} : Probability distribution over the operations {Mutation: P_m , Cross-Over: $1 - P_m$ }

Output: \mathcal{F}_{total} : Final feature pool, \mathcal{E}_{total} : Final feature extraction codes

```

1: Initialize  $\mathcal{F}_{pool} \leftarrow \mathcal{F}_{init}$ ;
2: Initialize  $\mathcal{E}_{code} \leftarrow CoderMLLM(I_s, \mathcal{F}_{pool})$ ;
3: for  $r \leftarrow 1$  to  $N$  do
4:   Initialize  $T \leftarrow \{\mathcal{L}\}$ ;
5:   for each feature extractor  $\epsilon_i$  in  $\mathcal{E}_{code}$  do
6:      $T \leftarrow T \cup \epsilon_i(\mathcal{D})$ ;
7:   end for
8:    $R \leftarrow$  Train a Random Forest Regressor with tabular features  $T$ ;
9:    $f \leftarrow$  Return the feature importance function in trained Random Forest Regressor  $R$ ;
10:   $\mathcal{F}_{pool}, \mathcal{E}_{code} \leftarrow FeatureSelector(\mathcal{F}_{pool}, \mathcal{E}_{code}, f, \min(k, \text{len}(\mathcal{F}_{pool})))$ ;
11:   $S_{op} \leftarrow$  Choose an operation from  $P_{op}$ ;
12:  if  $S_{op} = \text{Mutation}$  then
13:     $M \leftarrow RankSensitiveFormula(\mathcal{F}_{pool}, f)$ ;
14:     $\mathcal{F}_{parents} \leftarrow$  Select parent features  $\subset \mathcal{F}_{pool}$  from mutation probabilities  $M$ ;
15:     $\mathcal{F}_{new} \leftarrow MutatorMLLM(I_s, \mathcal{F}_{parents})$ ;
16:  else
17:     $\mathcal{F}_{new} \leftarrow CrossoverMLLM(I_s, \mathcal{F}_{pool})$ ;
18:  end if
19:   $\mathcal{F}_{new} \leftarrow DeduplicatorMLLM(I_s, \mathcal{F}_{pool}, \mathcal{F}_{new})$ ;
20:   $\mathcal{E}_{new} \leftarrow CoderMLLM(I_s, \mathcal{F}_{new})$ ;
21:   $\mathcal{F}_{pool} \leftarrow \mathcal{F}_{pool} \cup \mathcal{F}_{new}$ ;
22:   $\mathcal{E}_{code} \leftarrow \mathcal{E}_{code} \cup \mathcal{E}_{new}$ ;
23: end for
24:  $\mathcal{F}_{total}, \mathcal{E}_{total} \leftarrow \mathcal{F}_{pool}, \mathcal{E}_{code}$ ;

```

326 **B Evaluation Metric Definition**

327 **Definition 1 (SSIM).** *The Structural Similarity Index Measure (SSIM) evaluates the perceptual
328 similarity between two images based on luminance, contrast, and structural information. It is defined
329 as:*

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)},$$

330 where μ denotes the mean, σ^2 the variance, and σ_{xy} the covariance between prediction x and ground
331 truth y . Constants c_1 and c_2 stabilize the denominator.

332 **Definition 2** (MSE, NRMSE, PeakNRMSE). *The Mean Squared Error (MSE) is defined as:*

$$MSE(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2.$$

333 *The Normalized Root Mean Square Error (NRMSE) scales the RMSE by the value range of the
334 ground truth:*

$$NRMSE(x, y) = \frac{1}{\max(y) - \min(y)} \sqrt{MSE(x, y)}.$$

335 *The PeakNRMSE evaluates NRMSE restricted to the top- k largest entries of both x and y , emphasizing
336 high-congestion regions that are most critical in physical design.*

337 **Definition 3** (Correlation-Based Metrics: PLCC, SRCC, KRCC). *These metrics evaluate statistical
338 and rank-order correlation between prediction x and ground truth y :*

- 339 • **PLCC (Pearson Linear Correlation Coefficient)** measures linear correlation:

$$PLCC(x, y) = \frac{Cov(x, y)}{\sigma_x \sigma_y},$$

340 where $Cov(x, y)$ is the covariance and σ is the standard deviation.

- 341 • **SRCC (Spearman Rank-Order Correlation Coefficient)** computes PLCC on the rank-
342 transformed values of x and y , capturing monotonic relationships.
- 343 • **KRCC (Kendall Rank Correlation Coefficient)** compares the number of concordant and
344 discordant pairs:

$$KRCC(x, y) = \frac{n_c - n_d}{\binom{n}{2}},$$

345 where n_c and n_d are the counts of concordant and discordant pairs among all $\binom{n}{2}$ possible
346 pairs.

347 C Feature Attribution Breakdown Across All Features

348 To further examine the interpretability of the learned gating values, we analyze all 25 features used
 349 by the model. Figure 8 presents a comprehensive visualization, where each subplot corresponds to
 350 one feature.

351 For each subplot:

- 352 • The x-axis denotes the raw feature value.
- 353 • The y-axis shows the corresponding gating weight assigned by the model.
- 354 • Each point is colored by its associated congestion score, with darker colors indicating higher
 355 congestion.
- 356 • The left panel shows results from model *predictions*, and the right panel shows alignment
 357 with *ground-truth labels*.

358 **Observation 1: Prediction-Label Alignment.** Across most features, the gating weights show
 359 strong alignment between predicted and true congestion labels. This validates that the model is not
 360 only fitting the training objective, but that the learned preferences generalize to meaningful physical
 361 properties in actual designs.

362 **Observation 2: Unique Negative Attribution.** Among all 25 features, only one —
 363 `macro_rudy_boundary_interaction_index` — consistently receives a negative gating weight.
 364 This aligns with earlier findings in Section 5.1, where this feature ranked highly (inversely) corre-
 365 lated with congestion. The model correctly assigns it a down-weighted preference, consistent with
 366 human-understandable design heuristics.

367 **Observation 3: Gating Weights Capture Context Beyond Feature Values.** Notably, the spread
 368 of gating weights is often broader than that of feature values. In some features, the same value
 369 results in different gating weights across samples. This suggests that the model modulates importance
 370 dynamically based on context — likely derived from the input image, configuration tokens, or
 371 pretraining of the MLLM backbone. This further supports the idea that gating weights encode richer,
 372 more contextualized preference signals than raw features alone.

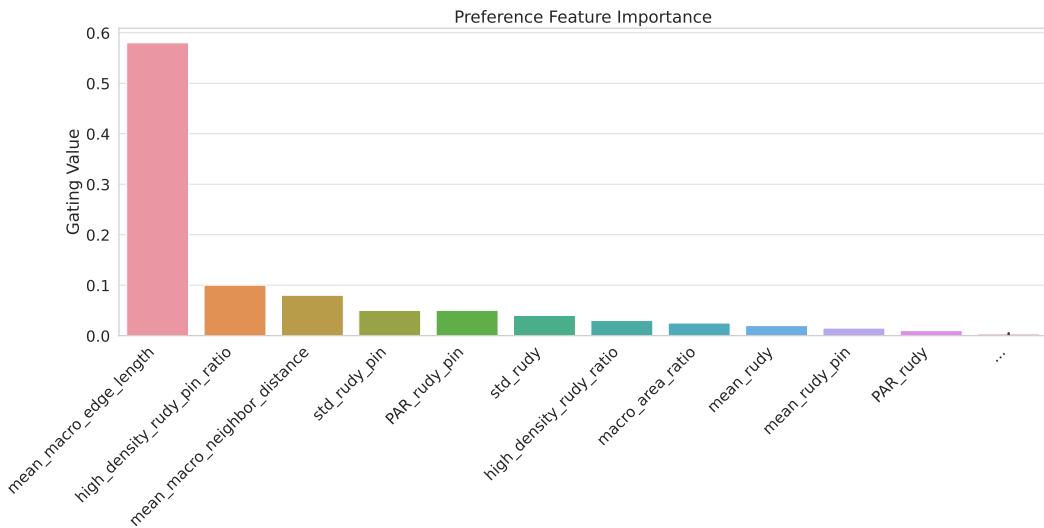
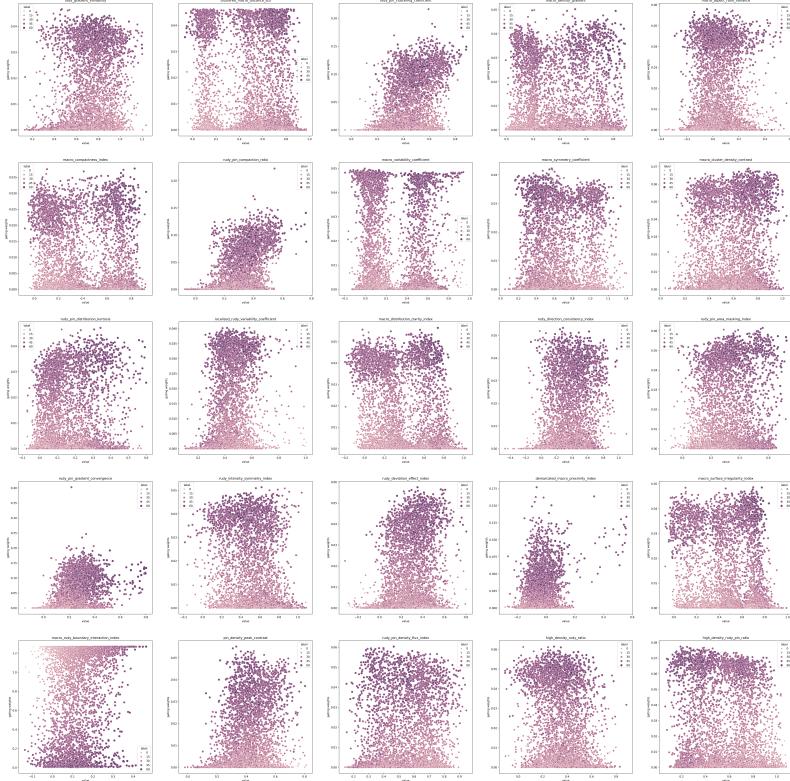
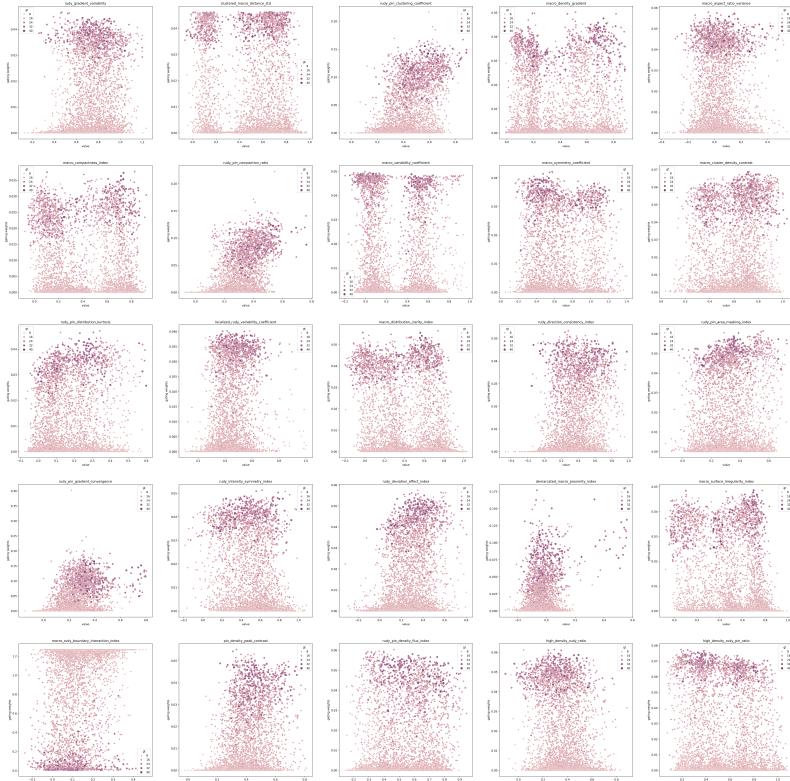


Figure 7: Bar plot showing the relative importance based on gating values from the preference features used for congestion prediction. The most influential feature is `mean_macro_edge_length`, significantly outweighing others. This suggests that macro-level spatial layout characteristics play a dominant role in determining congestion levels, followed by pin-density features such as `high_density_rudy_pin_ratio` and `mean_macro_neighbor_distance`.



(a) Gating values vs. feature values, colored by predicted congestion.



(b) Gating values vs. feature values, colored by ground-truth congestion.

Figure 8: Feature attribution breakdown across all 25 features. Each point represents a sample, colored by congestion score.

373 **D Detailed Ablation Study**

Table 5: Ablation Studies.

(a) Ablation on feature generation pipeline methods.

Pipeline Variant	zero-riscy-a			zero-riscy-b		
	PLCC ↑	SRCC ↑	KRCC ↑	PLCC ↑	SRCC ↑	KRCC ↑
Crossover-Only	0.630	0.674	0.492	0.589	0.686	0.493
Mutation-Only	0.353	0.485	0.348	0.469	0.462	0.314
Genetic Instruct	0.705	0.786	0.595	0.599	0.591	0.412

(b) Ablation on interpretable features.

Feature Source Modality	SSIM ↑	NRMSE ↓
Image-Based	0.791	0.047
Log-Based	0.764	0.052
Mixed (Image + Log)	0.789	0.046

(c) Ablation on conditional states.

Conditional Modality	SSIM ↑	NRMSE ↓
Features and Gating weights		
- numerical format (MLP)	0.787	0.047
- text format (T5 Encoder)	0.791	0.047

(d) Ablation on loss functions.

Loss Variant	SSIM ↑	NRMSE ↓
U-Net with text conditioning		
- MSE	0.791	0.047
- MPGD	0.789	0.047
- MS-SSIM + L1	0.762	0.079

(e) Ablation on congestion map generation modules.

Input Variant	SSIM ↑	NRMSE ↓
U-Net	0.791	0.047
U-Net + CNN Highway	0.792	0.048
U-Net + MLLM (trainable)	0.807	0.045

374 **E Experiment Setup**

375 **E.1 Training Setup**

376 All models are trained using the AdamW optimizer with mixed-precision (FP16) training enabled.
 377 Most experiments are conducted on 32 NVIDIA A100 GPUs, while some smaller-scale experiments
 378 and ablations are run on NVIDIA RTX A6000 GPUs. We train for 50 epochs with early stopping
 379 based on validation loss. The learning rate follows a linear warm-up over the first 500 steps and a
 380 cosine decay schedule thereafter. The dataset is split into 80% training, 10% validation, and 10%
 381 test. All reported results are based on the held-out test set. To ensure robustness, each experiment is
 382 repeated with three random seeds.

383 **E.2 Hyperparameters**

384 We set the initial learning rate to 1e-4 with weight decay of 0.01. The optimizer uses $\beta_1 = 0.9$ and
 385 $\beta_2 = 0.999$. The batch size is set to 4. We adopt MiniCPM-V-2_6 (Yao et al., 2024) as the backbone
 386 of our MLLM, where the visual encoder follows a SigLIP-like architecture. For the regression
 387 layer encoder, we use a linear layer appended to the backbone; for the gating layer, we use a ReLU
 388 MLP of 3 hidden layers with 1024 hidden units. The conditioning text descriptions for U-Net are
 389 encoded with a frozen T5 encoder model unless otherwise stated. Dropout of 0.1 is applied to all
 390 modality-specific encoders.

391 **E.3 Dataset Preprocessing**

392 We preprocess the CircuitNet dataset to construct aligned multimodal representations, including:
 393 (1) layout-based raster images (e.g., Macro Region, RUDY, and RUDY pin maps), (2) normalized
 394 tabular features extracted through code-driven analysis, and (3) optional configuration text describing
 395 design modules and target specifications. We partition the dataset into 80% training, 10% validation,
 396 and 10% test. All results reported are on held-out test samples. Additional augmentations include
 397 rotational invariance for layout images and controlled masking of tabular features to assess robustness.

398 **F Complete List of Generated Feature Pool and Feature Engineering Details**

Image-based Feature Pool (Appendix)	
Feature Name	Description
rudy_gradient_variability	the variation in gradient changes across the rudy map indicating potential areas of abrupt routing demand shifts
clustered_macro_distance_std	the standard deviation of distances between clustered groups of macros
rudy_pin_clustering_coefficient	a measure of how many rudy pins cluster together relative to the total number of rudy pins
macro_density_gradient	the change in macro density across the layout, impacting local congestion
macro_aspect_ratio_variance	the variance in aspect ratios of macros, indicating potential alignment and spacing issues that may impact congestion
macro_compactness_index	a measure of how closely packed the macros are, potentially affecting routing paths and congestion
rudy_pin_compaction_ratio	the ratio of compacted rudy pin clusters to the total number of rudy pins, indicating areas with high potential routing conflicts
macro_variability_coefficient	a measure of the consistency in macro sizes and shapes relative to each other, potentially affecting congestion balance
macro_cluster_density_contrast	the contrast in density between clustered groups of macros and their surrounding layout areas, indicating potential localized congestion pressure
rudy_pin_distribution_kurtosis	a measure of the peakedness or flatness in the distribution of rudy pins across the layout
localized_rudy_variability_coefficient	variation in RUDY intensity within localized regions, indicating potential micro-level congestion fluctuations
macro_distribution_clarity_index	a measure of how distinct macro distributions are across the layout, indicating clarity in separation
rudy_direction_consistency_index	a measure of the uniformity in the directional flow of RUDY intensity
rudy_pin_area_masking_index	the ratio of the area masked by rudy pin regions relative to the total layout
rudy_pin_gradient_convergence	a measure of how gradients in the rudy pin map converge into specific regions
rudy_deviation_effect_index	the deviation of RUDY intensities from the mean, indicating areas of abnormal routing demand
demarcated_macro_proximity_index	proximity of macros to predefined boundary regions, potentially affecting congestion near edges
macro_surface_irregularity_index	the irregularity in macro surface shapes, which can impact routing paths
macro_rudy_boundary_interaction_index	interaction between macros and high RUDY regions, indicating potential congestion hotspots
pin_density_peak_contrast	the contrast between peak pin density regions and their surroundings
rudy_pin_density_flux_index	the rate of change in rudy pin density across the layout
high_density_rudy_ratio	the ratio of areas with high RUDY intensity to the total layout area
high_density_rudy_pin_ratio	the ratio of areas with high RUDY pin intensity to the total layout area

Figure 9: Complete list of generated features used in our modeling pipeline. These capture spatial, statistical, and gradient-based properties derived from chip layout images.

Configuration and Log-derived Feature Pool	
Feature Name	Description
design_name	Design identifier (e.g., RISCV-a), useful for tracking design variants.
number_of_macros	Total number of macros in the design, affecting placement complexity.
clock_frequency	Operating frequency (e.g., 200MHz), indicating timing pressure that may affect placement/routing.
utilization	Target utilization (e.g., 70%), influencing available routing resources.
macro_placement	Macro placement strategy or group count, impacting layout structure.
power_mesh_setting	Power mesh granularity (e.g., 8), which may influence routing congestion and IR-drop safety margins.
filler_insertion	Stage of filler insertion (e.g., after placement), potentially affecting congestion distribution.
initial_placement_efficiency	Efficiency of macro and cell placement before routing begins.
instance_blockages_count	Number of instance blockages, reflecting routing resource obstruction from macro/cell placement.
hard_to_access_pins_ratio	Proportion of pins difficult to route due to placement or geometry constraints.
pin_density_variance_map	Variance in pin density across layout regions, highlighting routing hotspots.
multi_layer_pin_access_variability	Variation in pin accessibility across metal layers, indicating routing complexity.
non_default_routing_rule_usage	Frequency of non-default routing rules (if known from constraints or synthesis config).
crosstalk_sensitive_zones	Areas near critical nets where routing rules must prevent signal interference.

Figure 10: Complete list of configuration- and log-derived features used in our modeling pipeline. These include routing and congestion-related metrics obtained from the global placement and routing stages, as well as features generated by our own pipeline. They collectively capture layout-level difficulty indicators relevant for congestion prediction.

Feature Extraction Code Example: Macro Density Gradient

```

def macro_density_gradient(images):
    tiles_size = 2.25
    macro_image = images[0]
    rudy_image = images[1]
    rudy_pin_image = images[2]

    image_height, image_width = macro_image.shape
    total_image_area = image_width * image_height

    # Convert macro image to binary [0, 255]
    macro_image = np.uint8(macro_image * 255)
    _, binary_image = cv2.threshold(macro_image, 127, 255, cv2.THRESH_BINARY)

    # Find contours to get macro regions
    contours, _ = cv2.findContours(binary_image, cv2.RETR_EXTERNAL, cv2.
        CHAIN_APPROX_SIMPLE)

    # Calculate macro density per region
    macro_density = np.zeros((image_height, image_width))
    for contour in contours:
        mask = np.zeros_like(binary_image)
        cv2.drawContours(mask, [contour], -1, 255, thickness=cv2.FILLED)
        macro_density += mask

    # Gradient of the macro density
    gradient_x = cv2.Sobel(macro_density, cv2.CV_64F, 1, 0, ksize=5)
    gradient_y = cv2.Sobel(macro_density, cv2.CV_64F, 0, 1, ksize=5)

    # Calculate the magnitude of gradients
    gradient_magnitude = cv2.magnitude(gradient_x, gradient_y)

    # Calculate the average gradient magnitude in micrometers
    macro_density_gradient_um = np.sum(gradient_magnitude) / (image_height *
        image_width)
    macro_density_gradient_um *= tiles_size # Convert to micrometers

    return {"macro_density_gradient": macro_density_gradient_um}

```

Figure 11: Python implementation of the `macro_density_gradient` feature extractor. This function computes the average gradient magnitude of macro cell density across a chip layout image, offering a scalar representation of spatial density variation in micrometers.

Feature Extraction Code Example: High Density Rudy Ratio

```

def high_density_rudy_ratio(images):
    image = images[1]
    total_area = image.shape[0] * image.shape[1]
    mean_rudy = np.mean(image)
    high_density_rudy_ratio = (image > mean_rudy).sum() / total_area

    return {
        "high_density_rudy_ratio": high_density_rudy_ratio,
    }

```

Figure 12: Python implementation of the `high_density_rudy_ratio` feature extractor. This function computes the average gradient magnitude of macro cell density across a chip layout image, offering a scalar representation of spatial density variation in micrometers.

399 **NeurIPS Paper Checklist**

400 **1. Claims**

401 Question: Do the main claims made in the abstract and introduction accurately reflect the
402 paper's contributions and scope?

403 Answer: [Yes]

404 Justification: The abstract clearly reflects the paper's contributions, including the introduc-
405 tion of a multimodal LLM assistant for congestion prediction and design guidance, the use
406 of genetic prompting, and interpretable preference learning. These claims are supported by
407 experiments and analyses in the main text and do not overstate the scope or results.

408 Guidelines:

- 409 • The answer NA means that the abstract and introduction do not include the claims
410 made in the paper.
- 411 • The abstract and/or introduction should clearly state the claims made, including the
412 contributions made in the paper and important assumptions and limitations. A No or
413 NA answer to this question will not be perceived well by the reviewers.
- 414 • The claims made should match theoretical and experimental results, and reflect how
415 much the results can be expected to generalize to other settings.
- 416 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
417 are not attained by the paper.

418 **2. Limitations**

419 Question: Does the paper discuss the limitations of the work performed by the authors?

420 Answer: [Yes]

421 Justification: Limitation discussed in Section 6.

422 Guidelines:

- 423 • The answer NA means that the paper has no limitation while the answer No means that
424 the paper has limitations, but those are not discussed in the paper.
- 425 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 426 • The paper should point out any strong assumptions and how robust the results are to
427 violations of these assumptions (e.g., independence assumptions, noiseless settings,
428 model well-specification, asymptotic approximations only holding locally). The authors
429 should reflect on how these assumptions might be violated in practice and what the
430 implications would be.
- 431 • The authors should reflect on the scope of the claims made, e.g., if the approach was
432 only tested on a few datasets or with a few runs. In general, empirical results often
433 depend on implicit assumptions, which should be articulated.
- 434 • The authors should reflect on the factors that influence the performance of the approach.
435 For example, a facial recognition algorithm may perform poorly when image resolution
436 is low or images are taken in low lighting. Or a speech-to-text system might not be
437 used reliably to provide closed captions for online lectures because it fails to handle
438 technical jargon.
- 439 • The authors should discuss the computational efficiency of the proposed algorithms
440 and how they scale with dataset size.
- 441 • If applicable, the authors should discuss possible limitations of their approach to
442 address problems of privacy and fairness.
- 443 • While the authors might fear that complete honesty about limitations might be used by
444 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
445 limitations that aren't acknowledged in the paper. The authors should use their best
446 judgment and recognize that individual actions in favor of transparency play an impor-
447 tant role in developing norms that preserve the integrity of the community. Reviewers
448 will be specifically instructed to not penalize honesty concerning limitations.

449 **3. Theory assumptions and proofs**

450 Question: For each theoretical result, does the paper provide the full set of assumptions and
451 a complete (and correct) proof?

452 Answer: [NA]

453 Justification: No theoretical result.

454 Guidelines:

- 455 • The answer NA means that the paper does not include theoretical results.
- 456 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 457 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 458 • The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 459 • Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- 460 • Theorems and Lemmas that the proof relies upon should be properly referenced.

461 **4. Experimental result reproducibility**

462 Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

463 Answer: [Yes]

464 Justification: We provide comprehensive details of the dataset, model architecture, training setup, and hyperparameters necessary to reproduce the main results. The dataset (CircuitNet) is publicly available, and all preprocessing steps, data splits, and evaluation metrics are clearly described. Appendix E includes further implementation details, including model configurations and ablation setups. While the code is not included, all critical information required to replicate our main findings is disclosed in the paper.

465 Guidelines:

- 466 • The answer NA means that the paper does not include experiments.
- 467 • If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- 468 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- 469 • Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- 470 • While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - 471 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - 472 (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - 473 (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - 474 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in

506 some way (e.g., to registered users), but it should be possible for other researchers
507 to have some path to reproducing or verifying the results.

508 **5. Open access to data and code**

509 Question: Does the paper provide open access to the data and code, with sufficient instruc-
510 tions to faithfully reproduce the main experimental results, as described in supplemental
511 material?

512 Answer: [Yes]

513 Justification: The dataset used (CircuitNet) is publicly available. However, the code and
514 reproduction instructions are not yet released at the time of submission. We plan to make
515 both the codebase and detailed instructions publicly available upon publication to support
516 full reproducibility.

517 Guidelines:

- 518 • The answer NA means that paper does not include experiments requiring code.
- 519 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 520 • While we encourage the release of code and data, we understand that this might not be
521 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
522 including code, unless this is central to the contribution (e.g., for a new open-source
523 benchmark).
- 524 • The instructions should contain the exact command and environment needed to run to
525 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 526 • The authors should provide instructions on data access and preparation, including how
527 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 528 • The authors should provide scripts to reproduce all experimental results for the new
529 proposed method and baselines. If only a subset of experiments are reproducible, they
530 should state which ones are omitted from the script and why.
- 531 • At submission time, to preserve anonymity, the authors should release anonymized
532 versions (if applicable).
- 533 • Providing as much information as possible in supplemental material (appended to the
534 paper) is recommended, but including URLs to data and code is permitted.

537 **6. Experimental setting/details**

538 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
539 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
540 results?

541 Answer: [Yes]

542 Justification: Yes, the paper specifies all necessary details to understand the results. We
543 describe the dataset splits, the optimizer, learning rate schedule, batch size, number of
544 epochs, and other relevant hyperparameters. Additional implementation details, including
545 architecture configurations and ablation settings, are provided in Appendix E.

546 Guidelines:

- 547 • The answer NA means that the paper does not include experiments.
- 548 • The experimental setting should be presented in the core of the paper to a level of detail
549 that is necessary to appreciate the results and make sense of them.
- 550 • The full details can be provided either with the code, in appendix, or as supplemental
551 material.

552 **7. Experiment statistical significance**

553 Question: Does the paper report error bars suitably and correctly defined or other appropriate
554 information about the statistical significance of the experiments?

555 Answer: [No]

556 Justification: We do not report error bars or statistical significance metrics across multiple
557 random seeds. However, we provide extensive qualitative analysis, including multiple case
558 studies and human manual reviews, to support the consistency and interpretability of our
559 findings.

560 Guidelines:

- 561 • The answer NA means that the paper does not include experiments.
- 562 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
563 dence intervals, or statistical significance tests, at least for the experiments that support
564 the main claims of the paper.
- 565 • The factors of variability that the error bars are capturing should be clearly stated (for
566 example, train/test split, initialization, random drawing of some parameter, or overall
567 run with given experimental conditions).
- 568 • The method for calculating the error bars should be explained (closed form formula,
569 call to a library function, bootstrap, etc.)
- 570 • The assumptions made should be given (e.g., Normally distributed errors).
- 571 • It should be clear whether the error bar is the standard deviation or the standard error
572 of the mean.
- 573 • It is OK to report 1-sigma error bars, but one should state it. The authors should
574 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
575 of Normality of errors is not verified.
- 576 • For asymmetric distributions, the authors should be careful not to show in tables or
577 figures symmetric error bars that would yield results that are out of range (e.g. negative
578 error rates).
- 579 • If error bars are reported in tables or plots, The authors should explain in the text how
580 they were calculated and reference the corresponding figures or tables in the text.

581 8. Experiments compute resources

582 Question: For each experiment, does the paper provide sufficient information on the com-
583 puter resources (type of compute workers, memory, time of execution) needed to reproduce
584 the experiments?

585 Answer: [Yes]

586 Justification: We specify the types of compute resources used, including GPU models
587 and training configurations. Detailed information on the training environment, including
588 hardware setup and compute requirements, is provided in Appendix E.

589 Guidelines:

- 590 • The answer NA means that the paper does not include experiments.
- 591 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
592 or cloud provider, including relevant memory and storage.
- 593 • The paper should provide the amount of compute required for each of the individual
594 experimental runs as well as estimate the total compute.
- 595 • The paper should disclose whether the full research project required more compute
596 than the experiments reported in the paper (e.g., preliminary or failed experiments that
597 didn't make it into the paper).

598 9. Code of ethics

599 Question: Does the research conducted in the paper conform, in every respect, with the
600 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

601 Answer: [Yes]

602 Justification: This research fully adheres to the NeurIPS Code of Ethics. It does not involve
603 human subjects, sensitive data, or applications with foreseeable societal or environmental
604 harm. All experiments are conducted on publicly available datasets (CircuitNet), with
605 transparent methodology and reproducibility efforts documented in the appendix. The
606 paper also includes ablation studies, case analyses, and human evaluations conducted with
607 appropriate consent and without harm. Potential misuse of the work is minimal, and no
608 discriminatory or unfair bias is introduced.

609 Guidelines:

- 610 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
611 • If the authors answer No, they should explain the special circumstances that require a
612 deviation from the Code of Ethics.
613 • The authors should make sure to preserve anonymity (e.g., if there is a special consider-
614 ation due to laws or regulations in their jurisdiction).

615 **10. Broader impacts**

616 Question: Does the paper discuss both potential positive societal impacts and negative
617 societal impacts of the work performed?

618 Answer: [NA]

619 Justification: The research focuses on methodological improvements for chip design op-
620 timization using publicly available data and does not directly involve societal-facing ap-
621 plications. As such, a discussion of societal impact is not applicable in the context of this
622 work.

623 Guidelines:

- 624 • The answer NA means that there is no societal impact of the work performed.
625 • If the authors answer NA or No, they should explain why their work has no societal
626 impact or why the paper does not address societal impact.
627 • Examples of negative societal impacts include potential malicious or unintended uses
628 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
629 (e.g., deployment of technologies that could make decisions that unfairly impact specific
630 groups), privacy considerations, and security considerations.
631 • The conference expects that many papers will be foundational research and not tied
632 to particular applications, let alone deployments. However, if there is a direct path to
633 any negative applications, the authors should point it out. For example, it is legitimate
634 to point out that an improvement in the quality of generative models could be used to
635 generate deepfakes for disinformation. On the other hand, it is not needed to point out
636 that a generic algorithm for optimizing neural networks could enable people to train
637 models that generate Deepfakes faster.
638 • The authors should consider possible harms that could arise when the technology is
639 being used as intended and functioning correctly, harms that could arise when the
640 technology is being used as intended but gives incorrect results, and harms following
641 from (intentional or unintentional) misuse of the technology.
642 • If there are negative societal impacts, the authors could also discuss possible mitigation
643 strategies (e.g., gated release of models, providing defenses in addition to attacks,
644 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
645 feedback over time, improving the efficiency and accessibility of ML).

646 **11. Safeguards**

647 Question: Does the paper describe safeguards that have been put in place for responsible
648 release of data or models that have a high risk for misuse (e.g., pretrained language models,
649 image generators, or scraped datasets)?

650 Answer: [NA]

651 Justification: The work does not involve high-risk assets such as large pretrained language
652 models, generative models, or scraped datasets. It builds on a publicly available dataset
653 (CircuitNet) and applies domain-specific methods for chip design analysis, which pose
654 minimal risk of misuse. Therefore, safeguards for responsible release are not applicable in
655 this context.

656 Guidelines:

- 657 • The answer NA means that the paper poses no such risks.
658 • Released models that have a high risk for misuse or dual-use should be released with
659 necessary safeguards to allow for controlled use of the model, for example by requiring
660 that users adhere to usage guidelines or restrictions to access the model or implementing
661 safety filters.

- 662 • Datasets that have been scraped from the Internet could pose safety risks. The authors
663 should describe how they avoided releasing unsafe images.
664 • We recognize that providing effective safeguards is challenging, and many papers do
665 not require this, but we encourage authors to take this into account and make a best
666 faith effort.

667 **12. Licenses for existing assets**

668 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
669 the paper, properly credited and are the license and terms of use explicitly mentioned and
670 properly respected?

671 Answer: [Yes]

672 Justification: All external assets used in this work are properly credited. The dataset
673 used—CircuitNet—is cited appropriately (Jiang et al., 2023; Chai et al., 2023), and we
674 adhere to its open-source license and terms of use. No proprietary or restricted-access data,
675 code, or models were used without permission.

676 Guidelines:

- 677 • The answer NA means that the paper does not use existing assets.
678 • The authors should cite the original paper that produced the code package or dataset.
679 • The authors should state which version of the asset is used and, if possible, include a
680 URL.
681 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
682 • For scraped data from a particular source (e.g., website), the copyright and terms of
683 service of that source should be provided.
684 • If assets are released, the license, copyright information, and terms of use in the
685 package should be provided. For popular datasets, paperswithcode.com/datasets
686 has curated licenses for some datasets. Their licensing guide can help determine the
687 license of a dataset.
688 • For existing datasets that are re-packaged, both the original license and the license of
689 the derived asset (if it has changed) should be provided.
690 • If this information is not available online, the authors are encouraged to reach out to
691 the asset's creators.

692 **13. New assets**

693 Question: Are new assets introduced in the paper well documented and is the documentation
694 provided alongside the assets?

695 Answer: [NA]

696 Justification: The paper does not introduce new assets. All experiments are based on the
697 existing public dataset CircuitNet, and no new datasets, models, or tools are being released
698 as part of this work.

699 Guidelines:

- 700 • The answer NA means that the paper does not release new assets.
701 • Researchers should communicate the details of the dataset/code/model as part of their
702 submissions via structured templates. This includes details about training, license,
703 limitations, etc.
704 • The paper should discuss whether and how consent was obtained from people whose
705 asset is used.
706 • At submission time, remember to anonymize your assets (if applicable). You can either
707 create an anonymized URL or include an anonymized zip file.

708 **14. Crowdsourcing and research with human subjects**

709 Question: For crowdsourcing experiments and research with human subjects, does the paper
710 include the full text of instructions given to participants and screenshots, if applicable, as
711 well as details about compensation (if any)?

712 Answer: [NA]

713 Justification: The paper includes human evaluations in the form of manual reviews conducted
714 by the authors or domain experts, but it does not involve formal crowdsourcing or studies
715 with human subjects requiring ethical review or participant compensation. Therefore, this
716 question is not applicable.

717 Guidelines:

- 718 • The answer NA means that the paper does not involve crowdsourcing nor research with
719 human subjects.
- 720 • Including this information in the supplemental material is fine, but if the main contribu-
721 tion of the paper involves human subjects, then as much detail as possible should be
722 included in the main paper.
- 723 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
724 or other labor should be paid at least the minimum wage in the country of the data
725 collector.

726 **15. Institutional review board (IRB) approvals or equivalent for research with human**
727 **subjects**

728 Question: Does the paper describe potential risks incurred by study participants, whether
729 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
730 approvals (or an equivalent approval/review based on the requirements of your country or
731 institution) were obtained?

732 Answer: [NA]

733 Justification: The paper does not involve crowdsourcing or research with human subjects as
734 defined by IRB standards. All manual evaluations were performed by the authors or internal
735 collaborators and do not constitute formal human subject research. Therefore, IRB approval
736 is not applicable.

737 Guidelines:

- 738 • The answer NA means that the paper does not involve crowdsourcing nor research with
739 human subjects.
- 740 • Depending on the country in which research is conducted, IRB approval (or equivalent)
741 may be required for any human subjects research. If you obtained IRB approval, you
742 should clearly state this in the paper.
- 743 • We recognize that the procedures for this may vary significantly between institutions
744 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
745 guidelines for their institution.
- 746 • For initial submissions, do not include any information that would break anonymity (if
747 applicable), such as the institution conducting the review.

748 **16. Declaration of LLM usage**

749 Question: Does the paper describe the usage of LLMs if it is an important, original, or
750 non-standard component of the core methods in this research? Note that if the LLM is used
751 only for writing, editing, or formatting purposes and does not impact the core methodology,
752 scientific rigorousness, or originality of the research, declaration is not required.

753 Answer: [Yes]

754 Justification: This work leverages a multimodal LLM as a core component to automate
755 feature generation and support interpretable feature learning from chip design layouts and
756 metadata. Additionally, the LLM is used to generate human-readable design suggestions
757 based on model insights, enhancing downstream interpretability and usability. These usages
758 go beyond standard applications and are integral to the proposed methodology; they are
759 therefore clearly described in the paper.

760 Guidelines:

- 761 • The answer NA means that the core method development in this research does not
762 involve LLMs as any important, original, or non-standard components.
- 763 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
764 for what should or should not be described.