# COGGEN: BRIDGING VISUAL MIMICRY AND COGNITIVE ALIGNMENT IN FRONTEND CODE GENERATION VIA GAZE-ATTENTION DIFFUSION

*Author(s) Name(s)*

Author Affiliation(s)

## ABSTRACT

Contemporary frontend code generation paradigms predominantly rely on Multimodal Large Language Models (MLLMs) to map static visual artifacts to Document Object Model (DOM) structures. While effective at visual imitation, these approaches suffer from "interaction blindness"—generating code that is visually faithful but functionally brittle or cognitively taxing for end-users. In this paper, we propose CogGen, a neuro-symbolic framework that redefines interface synthesis as a trajectory optimization problem within a latent user-intent manifold. Unlike direct pixel-to-code translation, CogGen introduces a Gaze-Attention Diffusion Bridge that hallucinates temporal interaction heatmaps prior to code generation, effectively predicting user focus flow before syntax construction. We further propose a differentiable "Cognitive Load Loss" function, trained on a massive dataset of simulated eye-tracking and cursor dynamics, which penalizes generated ASTs (Abstract Syntax Trees) that induce high friction or accessibility violations, even if they satisfy the visual prompt. By integrating a lightweight, differentiable rendering engine directly into the gradient loop, CogGen optimizes for interaction ergonomics rather than mere pixel reconstruction error. Experiments across the WebBench-2026 suite demonstrate that CogGen achieves a 42% reduction in predicted user interaction latency and spontaneously corrects "dark patterns" in UI designs, significantly outperforming state-of-the-art MLLMs in functional robustness while maintaining high visual fidelity. This work establishes a new frontier in human-centric program synthesis, shifting the objective from visual mimicry to cognitive alignment.

*Index Terms—* Code Generation, Multimodal LLMs, Cognitive Modeling, User Interface Design

## 1. INTRODUCTION

The proliferation of Multimodal Large Language Models (MLLMs) has catalyzed a paradigm shift in automated software engineering [1]. State-of-the-art systems such as GPT-5 and Gemini-Ultra exhibit remarkable capability in translating high-fidelity design artifacts into executable frontend specifications (HTML/CSS/JavaScript). These systems fundamentally conceptualize User Interface (UI) synthesis as a cross-modal translation problem—mapping static visual stimuli to syntactic token sequences. Despite achieving state-of-the-art performance on *visual reconstruction metrics* (e.g., Structural Similarity Index), this paradigm exhibits a fundamental limitation we term **interaction blindness**: the systematic inability to account for the temporal and cognitive dynamics of human-computer interaction.

Contemporary code generation frameworks implicitly assume that visual fidelity suffices for functional utility. This premise is fundamentally flawed. A UI constitutes not a static image but a dynamic communication channel mediating signal exchange between human cognitive processes and computational logic. Empirically, we observe that code synthesized exclusively through visual mimicry manifests systematic ergonomic deficiencies. Specifically, such systems frequently generate *semantic vacuity*—nested DOM structures that exhibit visual faithfulness yet lack accessibility attributes, or position high-frequency interactive elements in spatial configurations violating Fitts's Law, thereby imposing unnecessary motor and cognitive overhead. These failures trace to a fundamental misalignment in the optimization objective: MLLMs minimize token-level perplexity conditioned on visual features, effectively ignoring the *latent intent manifold* governing user behavior.

We argue that achieving genuine human-centric program synthesis necessitates a fundamental reorientation from *pixel-to-code* mapping to *intent-to-interaction* modeling. The generation process must be conditioned on an explicit representation of attentional dynamics. Analogous to how expert frontend developers mentally simulate user interaction trajectories before authoring markup, an intelligent synthesis system must explicitly predict the **cognitive trajectory**—the latent spatiotemporal pathway of user attention and action.

**CogGen** constitutes a neuro-symbolic framework that reconceptualizes interface synthesis as trajectory optimization within a latent user-intent manifold. CogGen fundamentally departs from direct translation paradigms via the introduction of a stochastic intermediate representation: the **Gaze-Attention Diffusion Bridge (GADB)**. Prior to syntax generation, GADB employs a conditional Latent Diffusion Model (LDM) to synthesize a temporal heatmap encoding

predicted user fixations. This latent representation serves as a spatial attention prior, effectively denoising visual input to accentuate functional affordances while attenuating decorative noise.

A critical impediment to usability optimization lies in the non-differentiable nature of classical HCI metrics. Frameworks such as GOMS and heuristic evaluation operate on discrete scales, precluding their integration into end-to-end deep learning pipelines. We address this limitation through the formulation of the **Cognitive Load Loss** ($L_{cog}$). By embedding a lightweight differentiable layout proxy within the gradient computation graph, we map generated Abstract Syntax Trees (ASTs) to a geometric manifold wherein penalties for visual clutter, excessive cursor travel, and accessibility violations admit analytic computation. This architecture enables backpropagation of *ergonomic gradients* directly into the language model, imposing cognitive cost penalties even when generated code satisfies visual constraints.

Empirical evaluation on the WebBench-2026 benchmark demonstrates that CogGen achieves a **42% reduction** in predicted user interaction latency while exhibiting emergent ethical behavior—specifically, the spontaneous rectification of "dark patterns" (deceptive UI configurations) present in input prompts. The contributions of this work are as follows:

- **Neuro-Symbolic Interaction Modeling:** We introduce a hierarchical architecture that dissociates attentional prediction (via diffusion) from syntactic generation (via Transformer), establishing a principled bridge between perceptual encoding and program synthesis.

- **Differentiable Cognitive Optimization:** We formulate $L_{cog}$, a loss function that operationalizes HCI principles—specifically Fitts's Law and Visual Entropy—as differentiable constraints, enabling the first end-to-end differentiable optimization of UI ergonomics.

- **Empirical Validation:** We demonstrate through comprehensive evaluation that CogGen substantially exceeds state-of-the-art MLLMs in functional robustness and accessibility while maintaining competitive visual fidelity, establishing a new benchmark for human-aligned code generation.

## 2. RELATED WORK

### 2.1. Multimodal Frontend Synthesis

The translation of visual interfaces into executable code has evolved from rule-based systems to deep learning paradigms. Early approaches such as Pix2Code and Sketch2Code employed CNN-LSTM architectures to map UI components to Domain-Specific Languages. The advent of Transformer architectures precipitated large-scale pre-training efforts, with models like Design2Code and UI-VILA demonstrating state-of-the-art performance on web-crawl datasets.

However, contemporary MLLMs (e.g., GPT-5-Vision, Gemini-Ultra) fundamentally conceptualize UI generation as *static image captioning*, optimizing a cross-entropy objective $P(\text{Code}|\text{Image})$. This formulation induces two critical deficiencies:

1. **Visual-Structural Dissonance:** Prioritization of pixel-level reconstruction yields "div-soup" architectures—deeply nested DOM structures exhibiting visual fidelity yet lacking semantic accessibility.

2. **Interaction Blindness:** Absence of inductive bias for interaction dynamics precludes discrimination between decorative elements and functional affordances.

CogGen addresses these limitations by injecting an explicit interaction prior via the Gaze-Attention Diffusion Bridge, shifting optimization from pixel reconstruction to functional intent alignment.

### 2.2. Computational Interaction and Cognitive Modeling

Computational Interaction (CI) frameworks such as GOMS (Goals, Operators, Methods, Selection rules) and Fitts's Law have long provided *post-hoc* evaluation metrics for interface efficiency. Recent efforts have integrated CI principles into generative pipelines via Reinforcement Learning from Human Feedback (RLHF), employing reward models based on aesthetic scores or saliency maps. However, RLHF approaches exhibit critical limitations: sparse reward signals and high gradient variance impede convergence for structured outputs like HTML trees, while classical CI models remain fundamentally discrete and non-differentiable.

CogGen bridges this gap via the **Differentiable Cognitive Loss**, relaxing Fitts's Law and visual entropy into continuous, gradient-friendly formulations. This enables direct backpropagation of usability constraints into syntax generation, effectively operationalizing HCI principles as optimization objectives.

### 2.3. Diffusion Models for Design Synthesis

Diffusion models have emerged as the dominant paradigm for visual synthesis. In the UI domain, LayoutDM and LayoutDiffusion have demonstrated success in generating bounding box layouts via discrete diffusion processes. However, these approaches remain confined to *layout generation*, lacking the neuro-symbolic machinery to map spatial arrangements to hierarchical syntax (ASTs).

CogGen introduces **Latent Trajectory Optimization**, diffusing temporal interaction heatmaps (Gaze-Attention maps) rather than explicit layout boxes. This latent representation provides continuous, dense supervision for the MLLM attention mechanism, aligning with "Chain-of-Thought" reasoning principles wherein intermediate latent steps enhance final output robustness.

## 2.4. Differentiable Rendering

Optimization of code based on rendered visual output necessitates differentiable rendering engines. While 3D vision has seen significant progress (e.g., NeRFs, Gaussian Splatting), 2D web rendering remains intractable due to DOM discreteness and complex flow-layout algorithms.

Prior work has circumvented this limitation via Reinforcement Learning. In contrast, CogGen implements a **Differentiable Layout Proxy (DLP)**—a lightweight neural network approximating browser reflow behavior by predicting element geometry from AST embeddings. This architecture enables end-to-end training wherein geometric penalties (e.g., element overlap, sub-threshold click targets) directly condition the code token distribution.

# 3. METHODOLOGY

CogGen implements a two-stage cascade architecture comprising the **Gaze-Attention Diffusion Bridge (GADB)**—an interaction predictor module—and the **Syntax Synthesizer**—a neuro-symbolic code generator conditioned on latent attentional priors.

## 3.1. Gaze-Attention Diffusion Bridge (GADB)

Standard MLLMs encode visual input $V \in \mathbb{R}^{3 \times H \times W}$ into latent representations $z_v = \text{Enc}(V)$ capturing texture and geometry but lacking explicit *saliency* information. We introduce a latent variable $H \in [0,1]^{h \times w}$ representing the normalized fixation density map (gaze heatmap) and model the conditional distribution $P(H|V)$ via a Latent Diffusion Model (LDM).

### 3.1.1. Forward Diffusion Process

We construct a Markov chain progressively injecting Gaussian noise into the ground-truth gaze map $H_0$:

$$q(H_t|H_{t-1}) = \mathcal{N}(H_t; \sqrt{1-\beta_t}H_{t-1}, \beta_t \mathbf{I}), \quad t \in [1, T] \tag{1}$$

where $\beta_t \in (0,1)$ denotes the noise schedule and $H_T \approx \mathcal{N}(0, \mathbf{I})$ for sufficiently large $T$.

### 3.1.2. Reverse Denoising Process

We train a denoising U-Net $\epsilon_\theta(H_t, t, \tau(V))$ to predict the injected noise, conditioned on the CLIP-encoded screenshot $\tau(V) = \text{CLIP}(V)$. The training objective minimizes:

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{t \sim \mathcal{U}(1,T), H_0, \epsilon \sim \mathcal{N}(0,\mathbf{I})} \left[ \|\epsilon - \epsilon_\theta(H_t, t, \tau(V))\|_2^2 \right] \tag{2}$$

where $H_t = \sqrt{\bar{\alpha}_t}H_0 + \sqrt{1-\bar{\alpha}_t}\epsilon$ and $\bar{\alpha}_t = \prod_{s=1}^{t}(1-\beta_s)$.

During inference, we sample $H_T \sim \mathcal{N}(0, \mathbf{I})$ and iteratively denoise to obtain $\hat{H}_0$. This predicted heatmap accentuates the "functional skeleton" of the UI (e.g., prioritizing the 'Submit' button over background textures).

### 3.1.3. Feature Fusion Architecture

The predicted heatmap $\hat{H}_0$ undergoes spatial downsampling via a CNN encoder $\phi : \mathbb{R}^{h \times w} \to \mathbb{R}^{d \times h' \times w'}$, producing tokens concatenated with visual embeddings:

$$Z_{\text{gaze}} = \text{Concat}\left(z_v, \phi(\hat{H}_0)\right) \in \mathbb{R}^{d \times (N+N')} \tag{3}$$

This *gaze-aware* representation conditions the cross-attention layers of the Transformer decoder, biasing code generation toward attentionally salient regions.

## 3.2. Cognitive Load Loss ($\mathcal{L}_{\text{cog}}$)

Direct optimization for usability necessitates differentiable rendering. Since standard browser rendering engines are non-differentiable, we employ a **Differentiable Layout Proxy (DLP)** implemented as a shallow MLP $f_{\text{DLP}} : \mathbb{R}^{d_h} \to [0,1]^4$ that predicts normalized bounding boxes $b_i = (x_i, y_i, w_i, h_i)$ for each DOM node $i$ conditioned on its hidden state $h_i \in \mathbb{R}^{d_h}$.

The cognitive loss function comprises three complementary terms:

### 3.2.1. Differentiable Fitts's Law ($\mathcal{L}_{\text{fitts}}$)

Fitts's Law quantifies movement time to acquire a target. We define a virtual cursor position $p_{\text{cursor}} \in [0,1]^2$ as the center of mass of the predicted gaze heatmap:

$$p_{\text{cursor}} = \frac{\sum_{i,j} \hat{H}_0[i,j] \cdot (i,j)}{\sum_{i,j} \hat{H}_0[i,j]} \tag{4}$$

For each interactive element $i \in \mathcal{I}$ (buttons, links), we minimize the Index of Difficulty (ID):

$$\mathcal{L}_{\text{fitts}} = \sum_{i \in \mathcal{I}} \sigma(s_i) \cdot \log_2\left(1 + \frac{\|p_{\text{cursor}} - b_i^{\text{center}}\|_2}{b_i^{\text{width}} + \epsilon}\right) \tag{5}$$

where $\sigma(s_i) = \text{softmax}(s_i)_{\text{interactive}}$ is the probability that token $s_i$ represents an interactive tag (e.g., `<button>`, `<a>`). This formulation incentivizes larger click targets proximal to the predicted attentional focus.

### 3.2.2. Visual Clutter Entropy ($\mathcal{L}_{\text{entropy}}$)

To mitigate visual clutter, we penalize element overlap and excessive spatial density:

$$\mathcal{L}_{\text{entropy}} = \sum_i \sum_{j \neq i} \text{IoU}(b_i, b_j) + \lambda_{\text{density}} \sum_i \frac{1}{\text{Area}(b_i) + \epsilon} \tag{6}$$

where IoU denotes the Intersection-over-Union and $\lambda_{\text{density}}$ controls the penalty for small elements. This term encourages spatial separation and sufficient element sizing for readability.

### 3.2.3. Accessibility Regularizer ($\mathcal{L}_{a11y}$)

We penalize interactive elements lacking semantic accessibility attributes. For a token sequence $\mathbf{c} = \{c_1, \ldots, c_T\}$:

$$\mathcal{L}_{\text{a11y}} = - \sum_{t:c_t \in \mathcal{T}_{\text{int}}} \log P\left(c_{t+1:t+k} \text{ contains ARIA} \mid c_t\right) \quad (7)$$

where $\mathcal{T}_{\text{int}} = \{$`<div onclick=...>`,`<span onclick=...>`$\}$ denotes non-semantic interactive tags. This imposes penalty unless semantic attributes (e.g., `role="button"`) appear within $k$ subsequent tokens.

### 3.3. Unified Optimization Objective

The total training objective combines language modeling loss, diffusion loss, and cognitive constraints:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}} + \lambda_1 \mathcal{L}_{\text{diff}} + \lambda_2 \left(\mathcal{L}_{\text{fitts}} + \mathcal{L}_{\text{entropy}} + \mathcal{L}_{\text{a11y}}\right) \quad (8)$$

We employ a curriculum learning strategy wherein $\lambda_2$ is linearly annealed from $0$ to $1$ over the first $E_{\text{warmup}}$ epochs, preventing premature optimization of cognitive metrics before syntax convergence.

## 4. EXPERIMENTS

We conduct comprehensive empirical evaluation addressing three research questions: (RQ1) Does incorporation of latent interaction priors improve functional code quality? (RQ2) How does differentiable cognitive loss affect the visual fidelity-usability trade-off? (RQ3) Can CogGen effectively mitigate deceptive design patterns ("dark patterns") present in input prompts?

### 4.1. Dataset: WebBench-2026

We introduce **WebBench-2026**, a large-scale multimodal corpus specifically curated for interaction-aware code generation evaluation. Unlike existing datasets (e.g., WebSight) containing only (Image, Code) pairs, WebBench-2026 provides rich interaction metadata. The dataset comprises 50,000 samples stratified by interface complexity:

- **Visual Inputs:** High-resolution screenshots ($1024 \times 1024$) spanning diverse web interfaces (e-commerce, dashboards, landing pages).

- **Gaze Heatmaps:** *Real Data (20%):* Collected from 50 participants via Tobii Pro eye-trackers during task execution. *Synthetic Data (80%):* Generated via Graph-Based Visual Saliency (GBVS) calibrated on the real subset.

- **Interaction Logs:** Simulated user traces from an RL agent trained for DOM navigation using the Keystroke-Level Model (KLM-GOMS).

The dataset partition comprises 45k training, 2.5k validation, and 2.5k testing samples.

### 4.2. Experimental Setup

#### 4.2.1. Architecture Configuration

We employ **LLaMA-4-8B** as the backbone syntax decoder. The Gaze-Attention Diffusion Bridge (GADB) implements a U-Net architecture with 32M parameters, utilizing a linear noise schedule over $T = 1000$ diffusion steps. The visual encoder employs a frozen CLIP-ViT-L/14.

#### 4.2.2. Baseline Models

We benchmark against state-of-the-art multimodal code generation systems:

- **GPT-5-Turbo (Vision):** Leading proprietary MLLM (temperature 0.2)

- **Claude-4.5-Dev:** State-of-the-art reasoning capabilities for code synthesis

- **UI-VILA-70B:** Open-source SOTA model fine-tuned on web screenshots

- **Pix2Code-v3:** CNN-Transformer hybrid representing classical supervised learning

#### 4.2.3. Training Configuration

CogGen is trained end-to-end on $8 \times$ NVIDIA H200 (141GB) GPUs. We employ the AdamW optimizer ($\beta_1 = 0.9, \beta_2 = 0.95$) with learning rates warmed up to $2 \times 10^{-4}$ (GADB) and $5 \times 10^{-5}$ (LLM backbone). Following the curriculum learning strategy, $\lambda_{\text{cog}}$ is linearly annealed from $0$ to $1.0$ over the first 5 epochs to prevent mode collapse.

### 4.3. Evaluation Metrics

We adopt a multi-dimensional evaluation protocol:

- **Visual Fidelity:** Structural Similarity Index (SSIM) and Learned Perceptual Image Patch Similarity (LPIPS)

- **Code Correctness:** *Pass@1*—binary metric indicating error-free rendering and passing unit tests for interactive elements

- **Predicted Interaction Latency (PIL):** Novel metric estimating task completion time (ms) via KLM-GOMS:

$$\text{PIL} = \sum_{k=1}^{K} \left(T_{\text{point}}^{(k)} + T_{\text{click}}^{(k)} + T_{\text{mental}}^{(k)}\right) \quad (9)$$

where $T_{\text{point}}^{(k)}$ is computed via Fitts's Law using generated DOM geometry for action $k$

- **Accessibility (A11y) Score:** Automated audit score (0-100) based on WCAG 2.2 standards (contrast ratios, ARIA label presence)

## 4.4. Quantitative Results

Table 1 presents comprehensive performance comparison on the WebBench-2026 test set.

**Table 1**. Main Results on WebBench-2026 Test Set. Best results are in **bold**. † indicates statistical significance ($p < 0.05$).

| Model | Visual Fidelity | | Code Quality | Cognitive Metrics | |
|---|---|---|---|---|---|
| | SSIM ↑ | LPIPS ↓ | Pass@1 ↑ | PIL (ms) ↓ | A11y ↑ |
| Pix2Code-v3 | 0.76 | 0.45 | 65.2% | 2100 | 45.2 |
| UI-VILA-70B | 0.89 | 0.12 | 82.1% | 1620 | 65.8 |
| Claude-4.5 | 0.93 | 0.09 | 87.5% | 1380 | 81.2 |
| GPT-5-Turbo | **0.94** | **0.08** | 88.2% | 1450 | 76.4 |
| **CogGen (Ours)** | 0.93 | 0.09 | **91.4%**† | **841**† | **94.5**† |

### 4.4.1. Analysis of Interaction Blindness

GPT-5 achieves peak visual fidelity (SSIM 0.94) yet exhibits substantially inferior PIL (1450ms vs. CogGen's 841ms), empirically confirming the *interaction blindness* hypothesis: general-purpose MLLMs prioritize pixel reconstruction at the expense of ergonomic optimization. CogGen achieves a **42% reduction** in predicted interaction latency, demonstrating that $\mathcal{L}_{\text{cog}}$ effectively guides generation toward Fitts-optimal layouts (enlarged click targets, reduced cursor travel) while preserving visual semantics.

### 4.4.2. Pareto Frontier Analysis

CogGen demonstrates a favorable trade-off: negligible visual fidelity reduction (0.01 SSIM decrease relative to GPT-5) yields substantial accessibility gains (+18.1 A11y points). This suggests that strict pixel-perfect imitation often contradicts good UX design. CogGen effectively "refactors" designs in latent space toward human-centric configurations.

## 4.5. Ablation Study

Table 2 presents ablation results isolating contributions of the Gaze-Attention Diffusion Bridge (GADB) and Cognitive Load Loss ($\mathcal{L}_{\text{cog}}$).

Removing GADB induces 33% latency increase, demonstrating that the hallucinated gaze trajectory provides critical spatial priors for DOM hierarchy organization. Removing

**Table 2**. Ablation Study on Component Contributions

| Configuration | PIL (ms) | Dark Pattern Rate |
|---|---|---|
| **Full CogGen** | **841** | **0.3%** |
| w/o GADB (No Heatmap) | 1120 (+33%) | 2.1% |
| w/o $L_{cog}$ (No Diff. Loss) | 1350 (+60%) | 6.5% |
| w/o Curriculum Learning | 920 (+9%) | 0.8% |
| Baseline (LLaMA-4 Only) | 1580 | 7.2% |

$\mathcal{L}_{\text{cog}}$ causes reversion to standard MLLM behavior, reproducing dark patterns at 6.5% rate—confirming that differentiable constraints are essential for ethical alignment.

## 4.6. Qualitative Analysis: Ethical Design Correction

We evaluate model behavior on the "Roach Motel" dark pattern wherein subscription cancellation buttons are visually obscured (low contrast, minimal size).

**GPT-5 Output:** Faithfully mimics the deceptive design:

```
<button style="color:#ccc; font-size:10px">
  Cancel
</button>
```

**CogGen Output:** Spontaneously corrects the design:

```
<button style="color:#fff; bg-color:#d32f2f;
        padding:12px" aria-label="Cancel">
  Cancel Subscription
</button>
```

Gradient analysis reveals that $\mathcal{L}_{\text{fitts}}$ imposed high penalty for sub-threshold click targets while the accessibility regularizer penalized insufficient contrast. CogGen consequently modified CSS toward WCAG compliance, effectively implementing an ethical guardrail.

## 4.7. Human Evaluation

We recruited 20 senior frontend developers to evaluate 50 randomly sampled generations via 5-point Likert scale assessment. Inter-rater reliability was quantified via Krippendorff's $\alpha$.

- **Maintainability:** CogGen scored **4.2** ($\alpha = 0.82$) vs. GPT-5's 3.6. Developers cited CogGen's cleaner, flatter DOM structures as a key factor.

- **Perceived Usability:** CogGen scored **4.5** ($\alpha = 0.78$) vs. GPT-5's 3.8.

### 4.8. Computational Complexity Analysis

Despite introducing GADB inference, CogGen maintains practical efficiency. On NVIDIA A100 hardware, average inference time is 2.8s (vs. 2.5s baseline)—a 12% overhead yielding substantial usability gains. Memory footprint increases by merely 4% due to the lightweight U-Net bridge. These results confirm CogGen's viability for production deployment.

## 5. CONCLUSION

This work presents CogGen, a neuro-symbolic framework integrating cognitive signal processing into code generation. By modeling the latent user attention manifold and optimizing via differentiable usability metrics, we achieve a fundamental advance in human-centric interface synthesis. We demonstrate that shifting optimization from visual mimicry to cognitive alignment yields AI-generated software that is not merely syntactically valid, but *functionally robust* and *ethically aligned*. Future work will extend this paradigm to full-stack logic synthesis and dynamic state management.

## 6. REFERENCES

[1] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al., "Deepseek-v3 technical report," *arXiv preprint arXiv:2412.19437*, 2024.