

# VITA: Visual In-image Text Analysis on Vision-Language Models

Anonymous ACL submission

## Abstract

Vision-language models (VLMs) excel at document summarization, yet their robustness to visual formatting variations in document images is not well understood. We present VITA (Visual In-image Text Analysis), a systematic framework that measures how realistic visual changes—text emphasis and structural formatting—affect summarization quality. We evaluate six VLMs spanning early, middle, and late fusion architectures at two model scales. Across lexical, semantic, and information preservation metrics, we find architecture- and scale-dependent vulnerabilities: early fusion loses more information despite higher lexical stability, whereas late fusion preserves information but exhibits larger lexical variation. Structural formatting induces larger degradations than text emphasis. Scaling mitigates emphasis sensitivity but can exacerbate structural-format vulnerabilities in late fusion, indicating that robust document understanding may require architectural innovations beyond scaling. Source code and dataset are available at [Under Review]

## 1 Introduction

Vision-language models (VLMs) integrate visual and textual modalities for cross-modal reasoning (Zhang et al., 2024a; Chen et al., 2024a; Gigant et al., 2025; Orna et al., 2024; G et al., 2024). With advances in large-scale pretraining and fusion architectures, VLMs have expanded beyond classic multimodal tasks (Achiam et al., 2023; Dai et al., 2023; Li et al., 2023) to document understanding and summarization, where inputs are rendered document images (e.g., screenshots and PDFs) containing both text and visual cues (Nacson et al., 2025; Lee et al., 2023).

Unlike OCR-centric pipelines, VLMs can process document images directly (Faysse et al., 2024), leveraging visual signals (e.g., font, color, alignment, spatial layout) alongside textual semantics.

While such cues can convey emphasis, hierarchy, and structure, they can also induce brittleness: models may overfit to superficial styles, behave inconsistently across layouts/domains, or rely on non-semantic artifacts. This motivates the following research questions:

- How strongly do visual presentation cues affect the quality of VLM-generated document summaries?
- Which visual cues are most likely to change outputs even when the underlying text is identical?
- Do VLMs leverage meaningful visual signals, or are they overly sensitive to superficial stylistic variations that harm generalization?

A second key factor is *how* VLMs fuse visual and textual information (Hemker et al., 2024; Zhang et al., 2024b; Gavrikov et al., 2025). Fusion may occur early (token concatenation), in the middle (cross-attention), or late (representation-level merging), potentially shaping how visual cues influence generation. This raises an additional question:

- Does fusion timing and mechanism shape how VLMs process and respond to visual presentation cues?

Finally, a critical practical consideration is whether increased model capacity can mitigate these sensitivities. Scaling laws suggest that larger models generally achieve better performance (Kaplan et al., 2020), but recent work has shown that robustness does not automatically improve with scale (Howe et al., 2024), and some tasks even exhibit inverse scaling where larger models perform worse (McKenzie et al., 2023). This raises a final question:

- Does increased model scale reduce sensitivity to visual presentation cues, and if so, does this

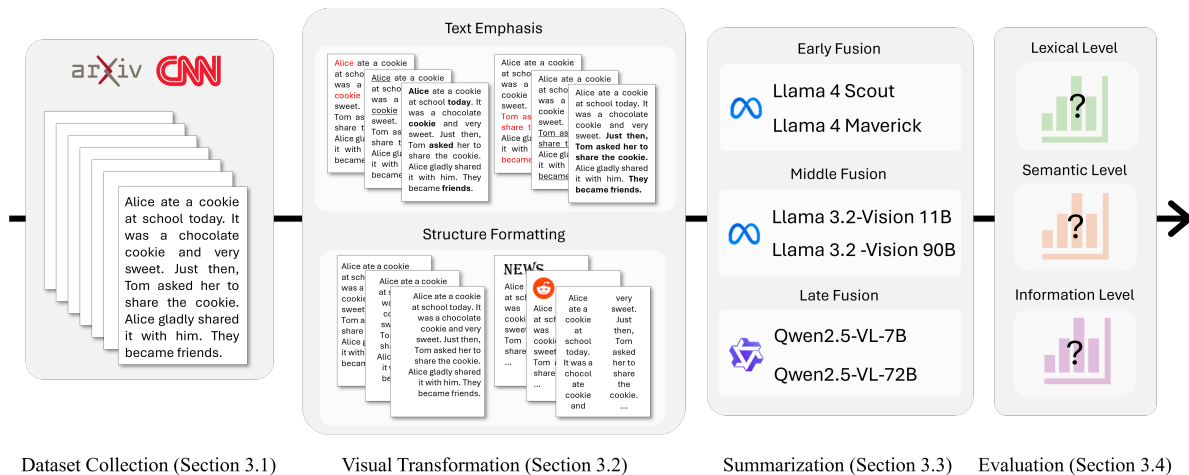


Figure 1: Overview of VITA. We keep the underlying text fixed and vary only visual attributes to form paired document images (text emphasis vs. structural formatting). We evaluate six VLMs across fusion timing (early/middle/late) and model scale, and quantify summary distortion using lexical/semantic similarity, verification-based faithfulness, and the rate of transformation-specific clue words that indicate format-aware summarization.

improvement generalize across all transformation types?

These questions matter because real-world documents vary widely in presentation (e.g., reports/news vs. social media/emails), and sensitivity to superficial formatting can undermine reliability. Yet controlled evidence isolating which visual attributes drive summarization changes remains scarce: prior work often evaluates document understanding without holding text constant under visual variations (Mathew et al., 2021; Ouyang et al., 2025; Qiu et al., 2022), and comparative analyses across fusion strategies and model scales are limited.

To fill this gap, we introduce **VITA (Visual In-image Text Analysis)**, a framework that isolates and quantifies the impact of realistic visual text variations on VLM-based summarization (Figure 1). We construct paired-input experiments with fixed text and varied visual attributes, organized into (1) text emphasis (e.g., color, bold, underline) and (2) structure formatting (e.g., alignment, columns, platform-style templates). We evaluate outputs via lexical drift (TF-IDF cosine similarity (Salton, 1983)), semantic change (SentenceBERT (Reimers and Gurevych, 2019)), and key information preservation (Minicheck (Tang et al., 2024)). We benchmark six representative VLMs spanning early, middle, and late fusion at two scales: Llama 4 Scout/Maverick (early), Llama 3.2-Vision 11B/90B (middle), and Qwen2.5-VL-7B/72B (late).

Our main contributions are summarized as fol-

lows:

- We propose VITA, a systematic analysis framework that uses controlled paired-input experiments to isolate visual changes while holding text constant.
- We analyze six VLMs across early, middle, and late fusion at two model scales, examining how fusion timing and capacity modulate sensitivity to visual presentation cues.
- We release a benchmark dataset of 2,600 document images with 12 visual transformations to facilitate future research.
- We provide quantitative and qualitative analyses revealing that (i) structure formatting causes greater degradation than text emphasis, (ii) fusion architectures exhibit distinct vulnerability patterns, and (iii) scaling selectively mitigates text emphasis sensitivity while leaving design transformation vulnerabilities unresolved.

## 2 Related Work

### 2.1 Visual Text Analysis

While substantial progress has been made in evaluating VLM performance on document understanding tasks, existing benchmarks exhibit significant limitations in systematically analyzing visual text variations. Traditional document VQA benchmarks such as DocVQA (Mathew et al., 2021), OCRBench (Fu et al., 2024), and OmniDocBench

(Ouyang et al., 2025) primarily utilize naturally occurring document images with diverse visual presentations. However, these datasets lack controlled experimental designs that enable direct comparison between identical textual content under different visual conditions, preventing systematic investigation of how specific visual attributes influence model behavior.

Recent studies have attempted to address this gap through more controlled approaches. Verma et al. (2023) introduced cross-lingual document pairs with identical content, while Chen et al. (2024b) evaluated VLM robustness by applying systematic transformations to document images. However, these approaches focus on linguistic variations or image-level degradations rather than fine-grained textual visual attributes commonly encountered in real-world documents.

A critical gap persists in prior work: the absence of systematic frameworks for isolating and quantifying the effects of specific visual text attributes on VLM performance. Current approaches lack the granular control necessary to understand how individual visual elements—such as text emphasis and structural formatting—influence model outputs when textual content remains constant.

## 2.2 Multimodal Fusion Structure

Recent studies have shown that the timing of multimodal fusion significantly affects model performance and sensitivity to external factors (Hemker et al., 2024; Zhang et al., 2024b; Gavrikov et al., 2025). Hemker et al. (2024) improved robustness to missing modalities by combining early fusion with attention mechanisms, Zhang et al. (2024b) achieved superior performance in document layout analysis through early-late fusion combinations, and Gavrikov et al. (2025) demonstrated that multimodal fusion itself changes how visual cues are processed. Additionally, Qiu et al. (2022) confirmed that modifications to the textual content cause the most severe distributional shifts in multimodal models when partially replacing text within images. However, existing studies have focused on general performance comparisons (Mathew et al., 2021; Fu et al., 2024; Ouyang et al., 2025). Despite evidence that fusion timing affects visual feature processing, no study has systematically analyzed how visual text attribute transformations impact different fusion architectures while keeping textual content constant.

## 2.3 Scaling and Robustness

Scaling laws have established that language model performance improves predictably with increased model size, training data, and compute (Kaplan et al., 2020). This relationship has driven the development of increasingly large models across both language and vision-language domains. However, the relationship between scale and robustness is more nuanced than general performance metrics suggest.

Howe et al. (2024) found that larger models are not consistently more robust to adversarial attacks without explicit safety training. More critically, McKenzie et al. (2023) demonstrated inverse scaling on 11 tasks where larger models performed worse than smaller ones, attributing this to memorization preferences and imitation of undesirable training patterns. In the multimodal domain, whether scaling mitigates sensitivity to visual presentation cues in document understanding remains unexplored—a significant gap given that VLMs must handle interactions between visual and textual modalities where scaling effects may differ from text-only models.

## 3 VITA: Visual In-image Text Analysis

To analyze how visual in-image text affects VLM-based document summarization, we propose VITA, short for Visual In-image Text Analysis. VITA isolates visual factors while keeping textual content constant, enabling controlled experiments over realistic visual variations and comparisons across fusion strategies. As shown in Figure 1, VITA consists of four stages: dataset collection, visual transformation, summarization, and evaluation.

### 3.1 Dataset Collection

We build a domain-diverse summarization dataset from two common sources: scientific papers and news articles. We collected 100 documents each from arXiv and CNN, restricting the time window to March–May 2025 to reduce potential overlap with VLM pretraining data<sup>1</sup>. For arXiv papers, we extract only the introduction; for news, we use the full article body. Each text is rendered into a single document image using LaTeX.

<sup>1</sup>We will release a publicly available code for data collection.

(a) Text Emphasis															
		Color		Bold		Underline				Color		Bold		Underline	
		Word	Sent	Word	Sent	Word	Sent	Word	Sent	Word	Sent	Word	Sent	Word	Sent
<i>Early Fusion</i>							<i>Early Fusion</i>								
Llama-4	LEX	-.00	.00	.01	.01	.01*	.01	Llama-4	LEX	.01	.01	.00	-.00	.01	.00
Scout	SEM	-.00	-.01	-.01*	-.01*	-.01	-.01	Maverick	SEM	-.01	-.01*	-.00	-.02**	-.01	-.00
	INFO	<b>.20***</b>	<b>.18***</b>	<b>.19***</b>	<b>.17***</b>	<b>.19***</b>	<b>.18***</b>		INFO	.00	-.02	.01	-.03	-.03	-.01
<i>Middle Fusion</i>							<i>Middle Fusion</i>								
Llama-3.2	LEX	.06*	.04	.02	.04*	.03	.04*	Llama-3.2	LEX	-.01	-.01	-.00	-.02	.01	-.00
11B	SEM	.03*	.02	.01	.02	.02*	.03*	90B	SEM	-.00	-.00	.01	-.01	.00	-.01
	INFO	<b>.15***</b>	<b>.14***</b>	<b>.11***</b>	<b>.12***</b>	<b>.11***</b>	<b>.11***</b>		INFO	-.02	-.01	-.01	-.03*	-.00	-.01
<i>Late Fusion</i>							<i>Late Fusion</i>								
Qwen2.5	LEX	.00	.03***	.02**	.05***	.04***	.04***	Qwen2.5	LEX	-.00	-.01	.01	.01	.01	.01*
7B	SEM	.01	.02*	.01	.01***	.02***	.01*	72B	SEM	-.00	-.00	-.00	.00	-.00	.00
	INFO	.05*	.06	.00	.00	.03	-.01		INFO	.00	.01	-.00	-.01	.01	-.00

(b) Structure Formatting															
		Alignment			Col	Design				Alignment			Col	Design	
		Left	Cen	Right	Two	Reddit	News			Left	Cen	Right	Two	Reddit	News
<i>Early Fusion</i>							<i>Early Fusion</i>								
Llama-4	LEX	.00	-.00	.01	.01	.06***	.03	Llama-4	LEX	-.00	-.01	-.01	.01	.04***	.03*
Scout	SEM	-.00	-.00	.02	.00	.06***	.04**	Maverick	SEM	.00	.01	.00	.02**	<b>.10***</b>	<b>.09***</b>
	INFO	<b>.20***</b>	<b>.16***</b>	<b>.18***</b>	<b>.21***</b>	<b>.23***</b>	<b>.47***</b>		INFO	.01	-.00	.01	.06***	<b>.34***</b>	<b>.39***</b>
<i>Middle Fusion</i>							<i>Middle Fusion</i>								
Llama-3.2	LEX	-.00	.06	.05*	.06*	<b>.15***</b>	.09***	Llama-3.2	LEX	-.01	-.00	-.01	.01	.05***	.02
11B	SEM	.02	.06***	.03	.05*	<b>.16***</b>	.06***	90B	SEM	-.00	.01*	.00	.02*	<b>.10***</b>	<b>.10***</b>
	INFO	<b>.10***</b>	<b>.13***</b>	<b>.12***</b>	<b>.15***</b>	<b>.22***</b>	<b>.25***</b>		INFO	-.01	-.00	-.01	.07***	<b>.36***</b>	<b>.40***</b>
<i>Late Fusion</i>							<i>Late Fusion</i>								
Qwen2.5	LEX	.02***	.01	.02	.02*	.05***	.06***	Qwen2.5	LEX	-.00	.00	.00	.02	.02***	.03***
7B	SEM	.01	.01	.01	.01*	.02***	.02***	72B	SEM	-.01*	-.01**	.00	.01	.01**	.00
	INFO	.03	.02	.01	.01	.01	.01		INFO	.00	.00	-.00	.01	.04**	<b>.26***</b>

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

Table 1: Performance degradation of six VLMs across 12 visual text transformations. LEX/SEM/INFO denote lexical, semantic, and information-level changes. Cen: center; Col Two: two-column. Left: small-scale; Right: large-scale models. Bold indicates values  $\geq 0.10$ .

### 3.2 Visual Transformation

Each document is programmatically styled in LaTeX to create 12 transformed versions, then compiled into images. The transformations mirror real-world document variations and fall into two groups: text emphasis, which highlights salient content, and structure formatting, which alters global layout and design.

**Text Emphasis.** We apply six transformations at both word and sentence levels, using color, bold, and underline. To select emphasis targets, we run gradient-based saliency analysis and highlight the top 15% of words that most influence model outputs. Goyal et al. (2022); Narayan et al. (2018) show that the top 10–20% saliency regions align well with key summarization content, and Hollenstein and Beinborn (2021) report partial alignment with human importance judgements.

**Structure Formatting.** We study global layout effects with three transformations of increasing in-

tensity. We first vary alignment across left, center, and right. We then apply a two-column layout, and finally reformat the same content into two distinct designs, namely a newspaper article and a Reddit post. All designs except Reddit are implemented in LaTeX. Reddit is rendered via HTML and captured as an image.

### 3.3 Summarization

We evaluate six VLMs spanning early, middle, and late fusion at both small and large scales to test how fusion depth and capacity modulate sensitivity to visual cues. We follow a standardized protocol and prompt each model with “Please summarize the text in this image.” We generate summaries for the original image and 12 transformed images, yielding 2,600 image-summary pairs from 200 base documents under 13 conditions, evaluated across two model scales for each fusion architecture.

**Early Fusion.** We use Llama 4 Scout with 109B parameters and Llama 4 Maverick with 402B parameters as early-fusion models. These models combine vision features with text tokens and process them jointly in a single backbone (Meta, 2024). This tight integration can strongly propagate fine-grained cues, such as color, font, and layout, into semantic representations.

**Middle Fusion.** LLaMA 3.2-Vision 11B and 90B represent middle fusion, where vision features interact with text via cross-attention at intermediate layers (Dubey et al., 2024). This design may partially transmit document-level visual transformations into summaries.

**Late Fusion.** We adopt Qwen2.5-VL-7B and Qwen2.5-VL-72B as late-fusion models. These models largely process visual information separately and merge it with text representations at later stages (Bai et al., 2025), which can limit the influence of fine-grained visual styles on summarization.

### 3.4 Evaluation

We compare summaries from original and transformed images using three complementary dimensions: lexical, semantic, and informational quality. For each transformation, we conduct paired comparisons and test score-difference normality; we apply paired  $t$ -tests (Hsu and Lachenbruch, 2014) when normal, and Wilcoxon signed-rank tests (Woolson, 2007) otherwise.

**Lexical.** We measure lexical similarity as TF-IDF cosine similarity between summaries from the original and transformed images:

$$\text{Lex}(S_o, S_t) = \cos(\text{TF}(S_o), \text{TF}(S_t)) \quad (1)$$

where  $S_o$  and  $S_t$  denote summaries from the original and transformed images, and  $\text{TF}(\cdot)$  denotes the TF-IDF vector. Lower scores indicate larger vocabulary shifts.

**Semantic.** We compute cosine similarity between Sentence-BERT embeddings (Reimers and Gurevych, 2019):

$$\text{Sem}(S_o, S_t) = \cos(e(S_o), e(S_t)) \quad (2)$$

where  $e(\cdot)$  denotes the SBERT embedding.

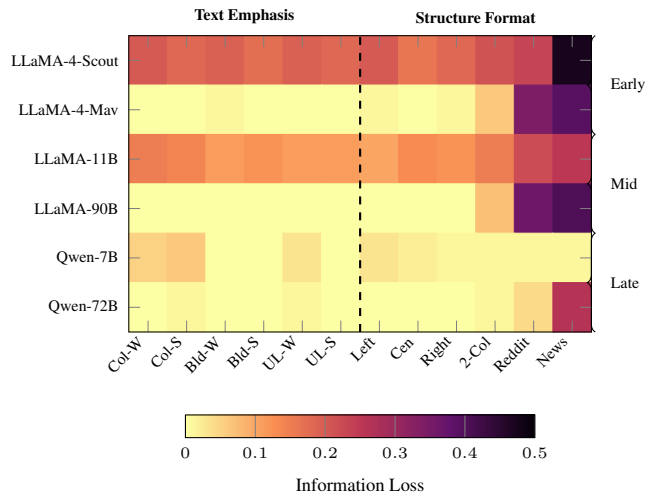


Figure 2: Information loss heatmap. Darker = higher degradation. Late Fusion (Qwen) is most robust; Design transformations cause severe degradation in Early/Middle Fusion.

**Information.** We evaluate factual consistency using Bespoke-MiniCheck-7B (Tang et al., 2024):

$$\text{Info}(S, D) = \text{MiniCheck}(S, D) \quad (3)$$

where  $D$  is the original document.

## 4 Results & Discussion

We analyze VITA results along three axes: (1) how models respond to different types of visual transformations, (2) how sensitivity patterns differ across fusion architectures, and (3) whether scaling mitigates these sensitivities uniformly. Table 1 presents performance changes for all six models, with small-scale models on the left and large-scale models on the right. Higher values indicate greater deviation from the original summary.

### 4.1 Visual Bias of In-Image Text

We first examine which types of visual transformations most strongly affect VLM summarization performance. Prior work suggests that document layout plays an important role in VLM processing (Zhang et al., 2024b), but systematic comparisons between local text modifications (e.g., color, bold) and global structural changes (e.g., layout, design templates) remain unexplored. Figure 2 visualizes information loss across all model-transformation combinations, revealing that Structure Formatting (right columns) causes substantially greater degradation than Text Emphasis (left columns), with Design transformations showing the most severe impact.

<i>Text Emphasis</i>	
Color	The generating function of the stabilized coefficients is ....involving the partition number into red colors.
Bold	The headline reads “Two Russian men ... The headline is written in large, bold font and is centered at the top of the page.
Underline	[NOT Detected]
<i>Structure Formatting</i>	
Alignment	It is aligned to the right side of the page...
Two Col	The article is divided into two columns... left column containing the introduction to the paper and the right column featuring..
News	The layout is typical of a newspaper article, with a clear heading, date, and body text...
Reddit	The image shows a screenshot of a Reddit page with a post about women in IT and software engineering.

Table 2: Examples of transformation-specific terminology in VLM outputs.

		LLaMA-4 Scout	LLaMA-3.2 11B	Qwen2.5 7B	LLaMA-4 Maverick	LLaMA-3.2 90B	Qwen2.5 72B
<i>Text Emphasis</i>							
Color	word	0.03	0.02	0.04	0.01	0.01	0.01
	sent	0.00	0.00	0.00	0.01	0.01	0.02
Bold	word	0.01	0.05	0.00	0.00	0.00	0.01
	sent	0.00	0.00	0.00	0.00	0.00	0.01
Underline	word	0.00	0.00	0.00	0.00	0.01	0.00
	sent	0.00	0.00	0.00	0.01	0.00	0.00
<i>Structure Formatting</i>							
Alignment	left	0.00	0.10	0.00	0.02	0.03	0.02
	center	0.04	0.09	0.03	0.05	0.04	0.04
	right	0.00	0.13	0.00	0.01	0.04	0.01
Column	two	0.03	<b>0.19</b>	0.00	<b>0.27</b>	<b>0.34</b>	0.00
Design	reddit	<b>0.42</b>	<b>1.00</b>	0.03	<b>0.90</b>	<b>0.99</b>	0.01
	news	<b>0.17</b>	<b>0.36</b>	0.03	<b>0.19</b>	<b>0.23</b>	0.04

Table 3: Proportion of summaries containing clue words for each visual condition.

**Text Emphasis.** Text Emphasis transformations demonstrate remarkable stability across all models and scales. LEXICAL changes remain minimal (0.00–0.06) and SEMANTIC alterations are equally constrained (0.00–0.03), indicating that visual emphasis largely preserves both vocabulary choice and meaning structure. Among transformation types, Color shows the highest stability with only 9 statistically significant effects across all conditions, compared to Bold (10) and Underline (12). This suggests that VLMs can largely filter out local text styling variations when processing document content.

**Structure Formatting.** In contrast, Structure Formatting transformations reveal a clear hierarchy of impact. Alignment and Column modifications show moderate effects (lexical 0.00–0.06, information loss 0.01–0.21), while Design transformations cause severe degradation across all metrics (lexical 0.03–0.15, semantic 0.02–0.16, information loss 0.01–0.47). Within Design transformations, Reddit format induces higher lexical change (0.05–0.15) than News format (0.03–0.09), though News

shows comparable or higher information loss in several models.

**Design as Critical Threshold.** Qualitative analysis reveals that Design transformations trigger a fundamental shift in model processing. As shown in Table 2, Text Emphasis occasionally generates format-aware references (e.g., “involving the partition number into red colors”) while preserving content focus. However, Design transformations cause complete processing mode shifts: News outputs include “The layout is typical of a newspaper article, with a clear heading, date, and body text...” while Reddit produces “The image shows a screenshot of a Reddit page with a post about...,” abandoning content analysis entirely in favor of format identification.

The clue word analysis (Table 3) quantifies this phenomenon. Text Emphasis shows minimal format terminology incorporation (at most 0.05 across all models), while Design transformations trigger dramatic increases—Reddit format causes terminology usage in 42% (Early), 100% (Middle), and only 3% (Late) of summaries. These patterns

<i>Model</i>	<i>Contents</i>
Early	The image appears to be a mock-up of the front page of The New York Times... / The image shows a Reddit page with a post about quantum mechanics.
Middle	The image presents a page from ... featuring an article on quantum mechanics.. / The image shows a Reddit post... The post discusses the variational method..
Late	The text discusses the variational method in quantum mechanics, particularly for calculating the energy levels.. / The text discusses the variational method in quantum mechanics, particularly for calculating the energy levels..

Table 4: Examples of fusion architecture effects on design transformation.

demonstrate that VLMs are fundamentally more sensitive to global layout changes than to local text styling, with Design transformations serving as a critical threshold that triggers complete shifts from content understanding to format recognition. At larger scales, this pattern persists: scaling eliminates Text Emphasis sensitivity but leaves Design vulnerabilities largely unresolved (Section 4.3).

## 4.2 Effect of Vision-Language Fusion

We next investigate how fusion architecture modulates sensitivity to visual transformations. While prior work has shown that fusion timing affects visual feature processing (Gavrikov et al., 2025), its specific impact on document summarization under controlled visual variations is unknown. We examine whether different fusion strategies exhibit distinct vulnerability patterns.

**Early Fusion.** At small scale, Llama 4 Scout exhibits a distinctive pattern: minimal lexical and semantic changes (at most 0.01) but severe information loss (0.17–0.20) under Text Emphasis, representing the highest information degradation among the three architectures. Design transformations cause extreme vulnerability, with News format recording 0.47 information loss—the maximum observed in our study. This pattern likely stems from Early Fusion’s architectural characteristics, where the vision encoder is deeply integrated into the language model backbone, causing layout changes to potentially disrupt the entire attention mechanism. At larger scale, Maverick eliminates Text Emphasis sensitivity but maintains Design vulnerability (Reddit: 0.34, News: 0.39), with clue word usage increasing from 0.42 to 0.90 for Reddit.

**Middle Fusion.** At small scale, Llama 3.2-Vision 11B shows a graduated response pattern: moderate lexical changes (0.02–0.06) alongside substantial information loss (0.11–0.15) for Text Emphasis, and vulnerability approaching Early Fusion for Design (0.22–0.25). The model shows the highest for-

mat awareness, with Reddit transformations triggering format-specific terminology in 100% of summaries (Table 3), suggesting that cross-attention mechanisms actively attend to visual formatting cues. The 90B variant achieves Text Emphasis robustness but shows increased Design sensitivity (Reddit: 0.36, News: 0.40), with near-saturated Reddit clue word detection (0.99).

**Late Fusion.** At small scale, Qwen2.5-VL-7B demonstrates a pattern opposite to Early Fusion: low information loss across all transformations (maximum 0.06), including Design changes (0.01), while showing slightly higher lexical variations (0.02–0.06). This robustness likely results from the MLP-based Vision-Language Merger architecture, which processes visual information largely separately before merging with text representations. Late Fusion maintains exceptional resistance to format-level distractions, with clue word usage below 4% even under Design transformations. As shown in Table ??, while Early and Middle Fusion outputs explicitly reference the visual format, Late Fusion consistently focuses on content. However, at 72B, unexpected News vulnerability emerges (0.26) despite Reddit stability (0.04), revealing implicit format biases even in robust architectures. Taken together, each fusion architecture exhibits a distinct vulnerability profile: Early Fusion sacrifices information for lexical stability, Late Fusion shows the opposite pattern, and no architecture achieves comprehensive robustness.

## 4.3 Effect of Model Scale

Table 5 summarizes scaling effects on Design transformations. Figure 3 visualizes the contrasting scaling effects: Text Emphasis degradation drops to near-zero across all architectures (left), while Design vulnerabilities persist or worsen with scale (right).

For Text Emphasis, large-scale models generally exhibit near-zero Info-loss shifts across archi-

Fusion	Scale	Info Loss		$\Delta$	
		Reddit	News	Reddit	News
<i>Early Fusion</i>					
Llama-4	Scout	0.23	0.47	+0.11	-0.08
	Maverick	0.34	0.39		
<i>Middle Fusion</i>					
Llama-3.2	11B	0.22	0.25	+0.14	+0.15
	90B	0.36	0.40		
<i>Late Fusion</i>					
Qwen2.5	7B	0.01	0.01	+0.03	<b>+0.25</b>
	72B	0.04	<b>0.26</b>		

Table 5: Fusion-Scale interaction on Design transformation. Positive  $\Delta$  indicates increased vulnerability with scaling. Bold indicates values  $\geq 0.25$ .

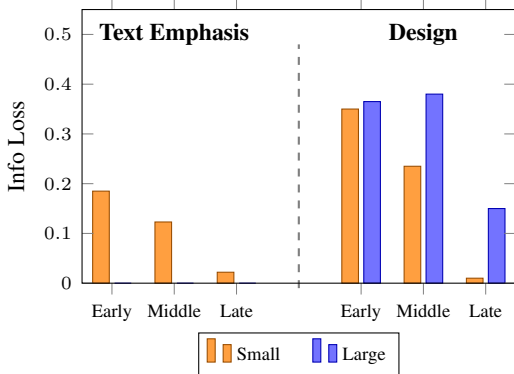


Figure 3: Scaling effect on information loss. Text Emphasis (left): scaling reduces degradation to near-zero. Design (right): scaling fails to mitigate vulnerability, with Middle Fusion showing inverse scaling.

tectures (Table 1a), suggesting reduced sensitivity with scaling. In contrast, Design transformations show architecture-dependent inverse-scaling tendencies (McKenzie et al., 2023): Early Fusion is mixed (Reddit +0.11, News -0.08), Middle Fusion consistently worsens (Reddit +0.14, News +0.15), and Late Fusion shows the sharpest divergence—Reddit remains stable (+0.03) while News increases from 0.01 to 0.26.

We hypothesize that larger models encode stronger format-specific priors, consistent with McKenzie et al.’s observation that scaling can amplify “imitation of undesirable training patterns.” Notably, Late Fusion’s News vulnerability emerges despite weak surface cue signals, suggesting an implicit format bias at deeper representational levels not captured by lexical diagnostics. Overall, scaling selectively mitigates visual sensitivity, and robust document understanding likely requires architectural changes beyond scaling alone.

## 5 Conclusion

This paper introduced VITA (Visual In-image Text Analysis), a systematic framework for analyzing how visual text transformations in document images affect VLM summarization. Through controlled paired-comparison experiments over 2,600 image–summary pairs and six VLMs spanning three fusion architectures at two scales, we addressed three research questions on visual sensitivity in document understanding.

First, we found that Structure Formatting transformations induce substantially larger degradation than Text Emphasis modifications. In particular, Design transformations (Reddit- and News-style layouts) act as a critical threshold: they reliably shift model behavior from content-focused summarization toward format-driven processing, producing consistent information loss across both design formats.

Second, fusion architectures exhibit distinct vulnerability profiles. Early Fusion tends to preserve lexical stability at the cost of information retention, whereas Late Fusion shows the opposite trade-off. Middle Fusion presents intermediate behavior with strong format awareness. Importantly, no single architecture provides comprehensive robustness across transformation families.

Third, we show that scaling selectively reduces visual sensitivity. Larger models generally attenuate vulnerabilities to Text Emphasis, suggesting that many emphasis cues are treated as shallow perturbations at higher capacity. However, sensitivity to Design transformations persists and can even increase with scale, consistent with prior observations of inverse-scaling tendencies where larger models amplify undesirable format priors. The Late Fusion contrast—where the 72B model exhibits emergent News vulnerability despite the 7B model remaining stable—highlights that scaling can introduce implicit format biases not captured by surface-level diagnostics.

Taken together, our findings caution against assuming that larger VLMs are inherently more robust to document format variation. Achieving reliable document understanding likely requires architectural and training interventions explicitly targeting format invariance, such as content–style disentanglement, format-agnostic objectives, or structure-aware calibration.

## 6 Limitations

Our results are based on statistically significant paired comparisons and show consistent Design effects across both Reddit- and News-style layouts, suggesting that the identified vulnerabilities are not driven by a single template. Nonetheless, our study has limitations.

First, our dataset size is moderate (200 documents per domain). Future work should validate the same trends at larger scale and across more diverse document genres and domains.

Second, we evaluate a controlled set of 12 transformations and six VLMs spanning three fusion architectures at two scales. Additional models (e.g., different instruction-tuning recipes or OCR-centric pipelines) and broader real-world formatting variations may reveal further nuances.

Third, VITA is primarily diagnostic: we focus on identifying “where and how” failures occur rather than proposing mitigations. Developing and testing interventions—such as format-robust augmentation, explicit content–style disentanglement, or architecture-level constraints for format invariance—remains important future work.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024a. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.
- Yufan Chen, Jiaming Zhang, Kunyu Peng, Junwei Zheng, Ruiping Liu, Philip Torr, and Rainer Stiefelhagen. 2024b. Rodla: Benchmarking the robustness of document layout analysis models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15556–15566.
- Wenliang Dai, Junning Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2024. Colpali: Efficient document retrieval with vision language models. *arXiv preprint arXiv:2407.01449*.
- Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang, Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, and 1 others. 2024. Ocrbench v2: An improved benchmark for evaluating large multimodal models on visual text localization and reasoning. *arXiv preprint arXiv:2501.00321*.
- Abinaya G, Guttulapavan Durga Rao, Pragada Sri Ram Gopal, Kaushik M, and B Sai Venkata Vignesh. 2024. Automated document processing: Combining ocr and generative ai for efficient text extraction and summarization. In *2024 International Conference on Innovative Computing, Intelligent Communication and Smart Electrical Systems (ICSES)*, pages 1–5.
- Paul Gavrikov, Jovita Lukasik, Steffen Jung, Robert Geirhos, Muhammad Jehanzeb Mirza, Margret Keuper, and Janis Keuper. 2025. Can we talk models into seeing the world differently? In *Thirteenth International Conference on Learning Representations*. OpenReview. net.
- Théo Gigant, Camille Guinaudeau, and Frédéric Dufaux. 2025. Summarization of multimodal presentations with vision-language models: Study of the effect of modalities and structure. *Preprint*, arXiv:2504.10049.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint arXiv:2209.12356*.
- Konstantin Hemker, Nikola Simidjievski, and Mateja Jamnik. 2024. Healnet: Multimodal fusion for heterogeneous biomedical data. *Advances in Neural Information Processing Systems*, 37:64479–64498.
- Nora Hollenstein and Lisa Beinborn. 2021. Relative importance in sentence processing. *arXiv preprint arXiv:2106.03471*.
- Nikolaus Howe, Ian McKenzie, Oskar Hollinsworth, Michał Zajac, Tom Tseng, Aaron Tucker, Pierre-Luc Bacon, and Adam Gleave. 2024. Scaling trends in language model robustness. *arXiv preprint arXiv:2407.18213*.
- Henry Hsu and Peter A Lachenbruch. 2014. Paired t test. *Wiley StatsRef: statistics reference online*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.

648	Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvasi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In <i>International Conference on Machine Learning</i> , pages 18893–18912. PMLR.	703
649		704
650		
651		705
652		706
653		707
654		
655	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	708
656		709
657		710
658		711
659		
660	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	712
661		713
662		
663		714
664		715
665		716
666		717
667	Ian R McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, and 1 others. 2023. Inverse scaling: When bigger isn’t better. <i>arXiv preprint arXiv:2306.09479</i> .	718
668		719
669		720
670		721
671		722
672		723
673		
674		724
675		725
676		726
677		727
678		728
679		729
680		
681		730
682		
683		731
684		732
685		733
686		
687		
688		734
689		735
690		736
691		737
692		738
693		739
694		740
695		741
696		742
697		743
698		744
699		745
700		746
701		747
702		

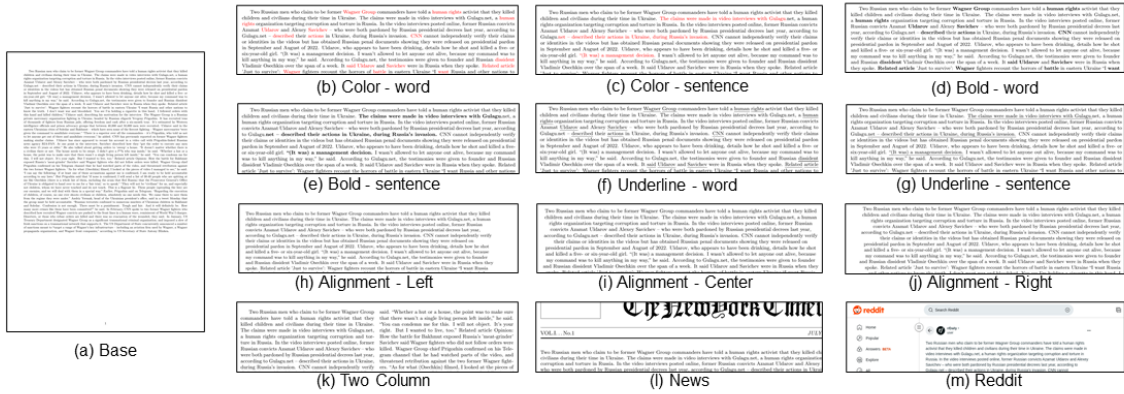


Figure 4: Examples of visual text transformations applied to the same document content. (a) Base: original document layout. (b)-(g) Text Emphasis transformations: color, bold, underline emphasis applied at word and sentence levels. (h)-(m) Structure Formatting transformations: alignment variations (left, center, right), two-column layout, and design (news article format, Reddit post format).

## C Detailed Statistical Results

We further report paired significance tests comparing the base condition against each visual transformation for the same three similarity dimensions. Following standard practice, we use paired t-tests when assumptions are met (reported as normal values) and Wilcoxon signed-rank tests otherwise (reported as underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.60	0.59	0.59	0.59	0.58	0.59
	std	0.15	0.13	0.12	0.12	0.12	0.13
<b>Semantic</b>	Mean	0.77	0.77	0.78	0.78	0.77	0.77
	std	0.13	0.11	0.11	0.11	0.12	0.12
<b>Information</b>	Mean	0.47	0.49	0.48	0.50	0.48	0.49
	std	0.33	0.32	0.32	0.32	0.33	0.33
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.59	0.60	0.59	0.59	0.54	0.57
	std	0.13	0.14	0.16	0.15	0.15	0.17
<b>Semantic</b>	Mean	0.76	0.77	0.75	0.76	0.70	0.72
	std	0.12	0.13	0.15	0.14	0.17	0.17
<b>Information</b>	Mean	0.47	0.50	0.49	0.46	0.44	0.20
	std	0.33	0.32	0.33	0.33	0.30	0.27

Table 6: Llama 4 Scout mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.50	0.52	0.54	0.52	0.52	0.52
	std	0.23	0.23	0.25	0.23	0.24	0.25
<b>Semantic</b>	Mean	0.65	0.66	0.66	0.66	0.65	0.64
	std	0.18	0.17	0.18	0.17	0.18	0.19
<b>Information</b>	Mean	0.15	0.15	0.18	0.18	0.18	0.19
	std	0.15	0.15	0.17	0.17	0.17	0.19
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.55	0.50	0.51	0.49	0.40	0.46
	std	0.25	0.25	0.24	0.23	0.27	0.24
<b>Semantic</b>	Mean	0.66	0.62	0.64	0.63	0.51	0.62
	std	0.18	0.20	0.19	0.20	0.23	0.20
<b>Information</b>	Mean	0.19	0.17	0.18	0.14	0.08	0.04
	std	0.18	0.16	0.18	0.15	0.08	0.06

Table 7: Llama 3.2-Vision 11B mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.64	0.61	0.62	0.60	0.60	0.60
	std	0.15	0.13	0.13	0.12	0.13	0.12
<b>Semantic</b>	Mean	0.87	0.86	0.87	0.86	0.86	0.87
	std	0.08	0.09	0.06	0.07	0.09	0.07
<b>Information</b>	Mean	0.67	0.66	0.71	0.72	0.69	0.72
	std	0.25	0.27	0.20	0.22	0.22	0.21
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.62	0.63	0.62	0.62	0.59	0.58
	std	0.14	0.13	0.13	0.13	0.12	0.10
<b>Semantic</b>	Mean	0.87	0.87	0.87	0.87	0.86	0.86
	std	0.08	0.08	0.09	0.08	0.08	0.07
<b>Information</b>	Mean	0.69	0.69	0.70	0.70	0.71	0.71
	std	0.24	0.22	0.24	0.21	0.21	0.23

Table 8: Qwen2.5-VL-7B mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.65	0.65	0.65	0.66	0.65	0.65
	std	0.12	0.13	0.12	0.11	0.13	0.11
<b>Semantic</b>	Mean	0.83	0.83	0.83	0.84	0.83	0.83
	std	0.10	0.10	0.08	0.08	0.10	0.08
<b>Information</b>	Mean	0.22	0.20	0.21	0.19	0.20	0.20
	std	0.27	0.23	0.25	0.23	0.23	0.23
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.65	0.66	0.67	0.65	0.62	0.63
	std	0.13	0.12	0.13	0.14	0.13	0.16
<b>Semantic</b>	Mean	0.82	0.82	0.82	0.81	0.73	0.73
	std	0.09	0.09	0.09	0.11	0.13	0.13
<b>Information</b>	Mean	0.21	0.22	0.22	0.20	0.20	0.23
	std	0.21	0.23	0.25	0.22	0.24	0.27

Table 9: LLaMA-4 Maverick mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.67	0.67	0.67	0.68	0.66	0.66
	std	0.12	0.13	0.13	0.12	0.12	0.12
<b>Semantic</b>	Mean	0.83	0.82	0.82	0.83	0.82	0.83
	std	0.09	0.09	0.09	0.08	0.09	0.08
<b>Information</b>	Mean	0.23	0.22	0.23	0.21	0.21	0.21
	std	0.26	0.28	0.27	0.23	0.25	0.23
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.67	0.68	0.67	0.65	0.62	0.64
	std	0.11	0.12	0.13	0.13	0.14	0.15
<b>Semantic</b>	Mean	0.83	0.81	0.82	0.80	0.72	0.72
	std	0.08	0.10	0.09	0.10	0.12	0.11
<b>Information</b>	Mean	0.22	0.23	0.22	0.20	0.20	0.24
	std	0.23	0.27	0.23	0.23	0.22	0.27

Table 10: LLaMA-3.2 90B mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Mean	0.66	0.67	0.65	0.65	0.66	0.65
	std	0.09	0.09	0.10	0.09	0.10	0.09
<b>Semantic</b>	Mean	0.88	0.89	0.88	0.88	0.88	0.88
	std	0.06	0.06	0.06	0.06	0.06	0.06
<b>Information</b>	Mean	0.20	0.22	0.20	0.21	0.19	0.21
	std	0.21	0.23	0.21	0.20	0.20	0.21
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Mean	0.66	0.66	0.66	0.65	0.64	0.63
	std	0.09	0.09	0.09	0.10	0.10	0.10
<b>Semantic</b>	Mean	0.89	0.89	0.88	0.87	0.87	0.88
	std	0.06	0.06	0.06	0.07	0.06	0.06
<b>Information</b>	Mean	0.20	0.20	0.20	0.20	0.20	0.21
	std	0.20	0.19	0.19	0.21	0.19	0.21

Table 11: Qwen2.5 72B mean similarity scores and standard deviations between original and visually transformed document summaries, measured across lexical (TF-IDF cosine similarity), semantic (Sentence-BERT embeddings), and information preservation (MiniCheck factual consistency) dimensions.

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	<u>2515</u>	<u>1915</u>	1.79	<u>1983</u>	<u>1916</u>	<u>2046</u>
	p-value	0.97	< 0.05	0.08	0.06	< 0.05	0.10
<b>Semantic</b>	Stats	<u>2171</u>	-0.46	<u>2123</u>	<u>2089</u>	<u>2125</u>	<u>2158</u>
	p-value	0.22	0.65	0.17	0.13	0.17	0.21
<b>Information</b>	Stats	<u>4284</u>	<u>4846</u>	<u>4312</u>	<u>4660</u>	<u>3967</u>	<u>4998</u>
	p-value	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	<u>2309</u>	0.41	<u>2474</u>	<u>2466</u>	<u>1616</u>	<u>2440</u>
	p-value	0.46	0.68	0.86	0.84	< 0.01	0.77
<b>Semantic</b>	Stats	<u>2474</u>	<u>2434</u>	<u>2439</u>	<u>2437</u>	<u>1859</u>	<u>2218</u>
	p-value	0.86	0.75	0.77	0.76	< 0.05	0.30
<b>Information</b>	Stats	<u>4697</u>	<u>5648</u>	<u>4789</u>	<u>4369</u>	<u>3814</u>	<u>831</u>
	p-value	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001

Table 12: Llama 4 Scout statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	2.22	2.55	0.38	1.50	2.76	2.72
	p-value	< 0.05	< 0.05	0.70	0.14	< 0.01	< 0.01
<b>Semantic</b>	Stats	2.88	<u>1663</u>	<u>1949</u>	2.54	<u>1653</u>	<u>1551</u>
	p-value	< 0.01	< 0.01	< 0.05	< 0.05	< 0.01	< 0.001
<b>Information</b>	Stats	<u>2985</u>	<u>2808</u>	<u>3982</u>	<u>3929</u>	<u>4451</u>	<u>4984</u>
	p-value	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	1.67	<u>1597</u>	2.57	2.15	<u>1802</u>	<u>2188</u>
	p-value	0.10	< 0.01	< 0.05	< 0.05	< 0.05	0.25
<b>Semantic</b>	Stats	<u>1530</u>	<u>1205</u>	<u>1578</u>	<u>1664</u>	7.01	<u>1952</u>
	p-value	< 0.001	< 0.001	< 0.01	< 0.01	< 0.001	< 0.05
<b>Information</b>	Stats	<u>4857</u>	<u>3958</u>	<u>4020</u>	<u>2418</u>	<u>734</u>	<u>279</u>
	p-value	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001

Table 13: Llama 3.2-Vision 11B statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	<u>1801</u>	2.62	<u>1502</u>	4.72	5.00	<u>1261</u>
	p-value	< 0.05	< 0.05	< 0.001	< 0.001	< 0.001	< 0.001
<b>Semantic</b>	Stats	<u>1921</u>	<u>2037</u>	1.97	<u>1691</u>	<u>1534</u>	<u>2060</u>
	p-value	< 0.05	0.09	0.05	< 0.01	< 0.001	0.11
<b>Information</b>	Stats	<u>8250</u>	<u>8859</u>	<u>8909</u>	<u>8868</u>	<u>8623</u>	<u>8939</u>
	p-value	< 0.05	0.15	0.73	0.15	0.41	0.18
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	<u>981</u>	<u>1529</u>	<u>1340</u>	<u>1297</u>	3.43	<u>1382</u>
	p-value	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001	< 0.001
<b>Semantic</b>	Stats	<u>2049</u>	<u>2040</u>	<u>1702</u>	<u>1944</u>	<u>1457</u>	<u>2107</u>
	p-value	0.10	0.10	< 0.01	< 0.05	< 0.001	0.15
<b>Information</b>	Stats	<u>10016</u>	<u>9551</u>	<u>9441</u>	<u>9985</u>	<u>9827</u>	<u>9700</u>
	p-value	0.97	0.54	0.46	0.94	0.79	0.67

Table 14: Qwen2.5-VL-7B statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	<u>9434</u>	<u>9024</u>	<u>9860</u>	<u>9638</u>	<u>9042</u>	<u>9594</u>
	p-value	0.45	0.21	0.82	0.62	0.22	0.58
<b>Semantic</b>	Stats	<u>8636</u>	<u>8186</u>	<u>9864</u>	<u>7759</u>	<u>9063</u>	<u>9812</u>
	p-value	0.08	< 0.05	0.82	< 0.01	0.23	0.77
<b>Information</b>	Stats	<u>9788</u>	<u>8942</u>	<u>9111</u>	<u>8925</u>	<u>8769</u>	<u>9558</u>
	p-value	0.75	0.18	0.25	0.17	0.12	0.55
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	<u>9704</u>	<u>9537</u>	<u>9852</u>	<u>9375</u>	<u>6979</u>	<u>8166</u>
	p-value	0.67	0.53	0.81	0.41	< 0.001	< 0.05
<b>Semantic</b>	Stats	<u>9818</u>	<u>9094</u>	<u>9514</u>	<u>7915</u>	<u>3119</u>	<u>3015</u>
	p-value	0.78	0.24	0.51	< 0.01	< 0.001	< 0.001
<b>Information</b>	Stats	<u>9058</u>	<u>10024</u>	<u>9226</u>	<u>7118</u>	<u>659</u>	<u>583</u>
	p-value	0.23	0.97	0.31	< 0.001	< 0.001	< 0.001

Table 15: LLaMA-4 Maverick statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	<u>9402</u>	-0.76	-0.55	<u>8454</u>	<u>9165</u>	<u>9893</u>
	p-value	0.43	0.45	0.58	0.05	0.28	0.85
<b>Semantic</b>	Stats	<u>9602</u>	<u>9281</u>	<u>9779</u>	<u>8864</u>	<u>9881</u>	<u>8915</u>
	p-value	0.58	0.35	0.74	0.15	0.84	0.20
<b>Information</b>	Stats	<u>8756</u>	<u>8617</u>	<u>9570</u>	<u>8375</u>	<u>9920</u>	<u>9428</u>
	p-value	0.11	0.08	0.56	< 0.05	0.87	0.45
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	-1.27	-0.58	-1.55	1.40	4.65	1.64
	p-value	0.21	0.56	0.12	0.16	< 0.001	0.10
<b>Semantic</b>	Stats	<u>9599</u>	<u>8429</u>	<u>9540</u>	<u>7965</u>	10.23	<u>1974</u>
	p-value	0.58	< 0.05	0.53	< 0.05	< 0.001	< 0.001
<b>Information</b>	Stats	<u>9738</u>	<u>10018</u>	<u>9619</u>	4.34	<u>10</u>	<u>265</u>
	p-value	0.70	0.97	0.60	< 0.001	< 0.001	< 0.001

Table 16: LLaMA-3.2 90B statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).

(a) Text Emphasis		Color		Bold		Underline	
		Word	Sentence	Word	Sentence	Word	Sentence
<b>Lexical</b>	Stats	<u>9225</u>	<u>8510</u>	<u>8883</u>	<u>8919</u>	<u>9752</u>	<u>8120</u>
	p-value	0.31	0.06	0.15	0.17	0.72	< 0.05
<b>Semantic</b>	Stats	<u>9596</u>	<u>8820</u>	-0.35	<u>9847</u>	<u>9510</u>	<u>9482</u>
	p-value	0.58	0.13	0.72	0.80	0.51	0.49
<b>Information</b>	Stats	<u>9496</u>	<u>8909</u>	<u>9374</u>	<u>9055</u>	<u>9563</u>	<u>9691</u>
	p-value	0.50	0.16	0.48	0.22	0.82	0.66
(b) Structure Formatting		Alignment			Column	Design	
		Left	Center	Right	Two	Reddit	News
<b>Lexical</b>	Stats	<u>9425</u>	<u>9581</u>	<u>9457</u>	<u>8498</u>	<u>7069</u>	<u>6128</u>
	p-value	0.45	0.57	0.60	0.47	< 0.001	< 0.001
<b>Semantic</b>	Stats	<u>8255</u>	<u>7723</u>	<u>9813</u>	<u>8534</u>	<u>7905</u>	<u>9513</u>
	p-value	< 0.05	< 0.01	0.77	0.06	< 0.01	0.51
<b>Information</b>	Stats	<u>9659</u>	<u>9995</u>	<u>9238</u>	<u>8933</u>	<u>7746</u>	<u>3820</u>
	p-value	0.63	0.95	0.32	0.17	< 0.01	< 0.001

Table 17: Qwen2.5 72B statistical test results comparing similarity scores between base summary and transformed image summaries across lexical, semantic, and information preservation dimensions, with test statistics and p-values from paired t-tests (normal values) and Wilcoxon signed-rank tests (underlined values).