

# The Visual Faithfulness Paradox: Scaling Vision–Language Models Degrades Glyph Recognition in Logographic Scripts

Anonymous ACL submission

## Abstract

We identify a *Visual Faithfulness Paradox* in scaling Vision–Language Models (VLMs): although larger models improve aggregate OCR metrics, they can systematically degrade glyph-level accuracy for non-alphabetic scripts. In controlled trilingual OCR experiments using InternVL3 (1B–14B), English performance increases monotonically. Chinese follows a non-monotonic scaling curve: small models suffer visual collapse (1B/2B), the 8B model achieves the best visual–language balance, and the 14B model drifts into language-prior–driven glyph substitution. Japanese exhibits an intermediate pattern, while visually similar kana generate persistent variability. A fine-grained evaluation on a mixed-script Chinese phrase confirms the shift: the 14B model’s perfect-match rate falls to 25% (vs. 54% at 8B), while its semantic-deviation error rate is 7.28x higher. We explain these effects using a *Visual Signal-to-Noise Ratio* (VSNR) account: beyond a script-dependent threshold, strengthening language priors override ambiguous visual evidence during decoding. These results expose a fundamental trade-off between linguistic fluency and visual faithfulness, challenging the assumption that “bigger is always better”. Future VLMs must explicitly reinforce glyph-structure perception and coordinate the glyph stream with the knowledge stream.

## 1 Introduction

Recent advances in Vision–Language Models (VLMs) have driven substantial progress in document understanding, optical character recognition (OCR), and multimodal reasoning. To accommodate high-resolution inputs and long visual–textual contexts, recent research has increasingly emphasized *visual token efficiency*. Representative systems explore aggressive visual compression to enable long-context OCR under constrained activation budgets (Wei et al., 2025), while others pursue native multimodal pre-training with flexible visual

position encoding (Zhu et al., 2025) and dynamic resolution strategies with multimodal rotary position encoding (Bai et al., 2025). Collectively, these approaches aim to reduce computational cost while preserving fine-grained visual perception.

However, as VLMs scale in parameter size and architectural capacity, an underexplored tension emerges. The dominant inference regime gradually shifts from vision-grounded perceptual decoding toward language-prior-driven semantic completion. For alphabetic scripts such as English, where visual form and lexical identity are weakly coupled, stronger language modeling typically yields monotonic OCR improvements. In contrast, for logographic writing systems such as CJK (Chinese, Japanese, and Korean), where visual structure is tightly bound to character identity, this shift can induce systematic degradation in visual faithfulness.

This paper formalizes this phenomenon as the *Visual Faithfulness Paradox*: as VLMs become larger and more linguistically competent, their ability to faithfully transcribe visually grounded glyphs, especially complex and isolated characters such as “漢” or highly similar kana symbols like フ vs. “プ”, can systematically deteriorate, exhibiting non-monotonic or even reversed scaling behavior.

This paradox reflects a broader pattern of inductive bias in large models. It relates to *the paradox of poetic intent* observed in machine translation, where models achieve high surface-level scores while failing to preserve deeper semantic or cultural structure (Weigang and Brom, 2025). Complementary, Six-Writings Pictophonetic Coding (SWPC) theory argues that Chinese characters require explicit modeling of graphic–phonetic–semantic structure, rather than being treated as generic visual textures (Weigang et al., 2024). From a physical perspective, Hanzi character image threshold theory further shows that common patch-based resizing strategies can cause stroke fusion and entropy collapse for high-density characters, produc-

Language	Orthographic Structure	Visual Discriminability Profile	Dominant Error Mode	Severity
English	Alphabetic script with high inter-letter shape distinctiveness	High visual discriminability; clear shape contrasts among letters mean visual evidence is dominant and language priors are mostly supportive.	Spelling or punctuation errors	Minor: primarily detail-level inaccuracies
Japanese	Mixed system: Kanji (2,136 – semantic logographs) + Kana (92 – simple, visually similar syllabary)	Kana similarity induces visual ambiguity (e.g., “フ” vs. “フ”); kanji are stable, but kana sequences depend more strongly on language priors.	Prior-driven substitutions and insertions in visually ambiguous kana segments (e.g., “フリ” → “フリーカ/フリーワ”).	Moderate: localized hallucinations concentrated in kana
Chinese	Logographic script with high stroke density (many characters in 16px >12 strokes, e.g., 漢)	Visual overload at common UI scales; complex characters exceed perceptual thresholds, degrading glyph features.	<i>Semantic-deviation substitutions</i> where visually degraded glyphs are replaced by linguistically plausible characters (e.g., 汉漢→汉汉, 汉漢→汉语).	Severe: systematic loss of glyph-level identity

Table 1: Cross-linguistic contrasts in visual discriminability and dominant OCR error modes.

ing intrinsically ambiguous visual signals prior to decoding (Weigang et al., 2025).

To systematically investigate the *Visual Faithfulness Paradox*, we conduct controlled OCR experiments on aligned English, Chinese, and Japanese Wikipedia screenshots of the *Isotopes of Hydrogen* page<sup>1</sup> using the InternVL3 model family (1B, 2B, 8B, and 14B) (Zhu et al., 2025; Lu et al., 2025). All configurations are evaluated under repeated stochastic sampling to ensure robustness and to disentangle visual effects from decoding variance. The data, code, and detailed replication instructions are available in our GitHub repository<sup>2</sup>.

The cross-linguistic scaling behavior (Table 1) diverges sharply. For English, OCR accuracy improves monotonically as model capacity increases, with error rates falling from roughly 0.57 (1B) to 0.04 (14B). This reflects the strong visual distinctiveness of alphabetic letters and the supportive role of language priors. Japanese, with its mixed kana–kanji system, displays a more nuanced pattern: kanji recognition remains stable across scales, whereas visually subtle kana sequences, which exhibit weak visual distinctiveness and weaker lexical constraints.

Chinese presents the starkest case. Although aggregate OCR accuracy appears to improve steadily with scale, fine-grained inspection under simplified–traditional coexistence reveals severe structural distortions at large model sizes. As illustrated in Figure 1, sequence-level OCR accuracy increases monotonically, yet glyph-level visual error follows a U-shaped trajectory, with the *perfect-*

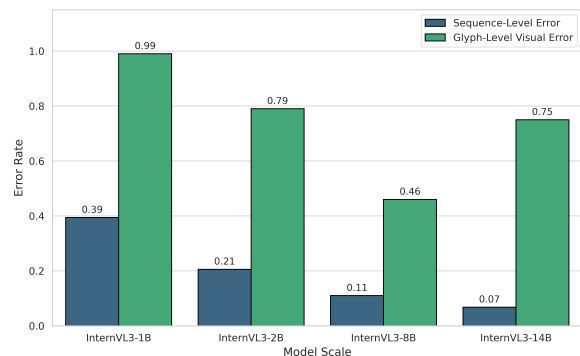


Figure 1: Sequence-level OCR accuracy improves monotonically with model scale, whereas glyph-level visual error on Chinese follows a U-shaped trajectory.

*match error rate* reaching its minimum at 8B (46%) but rising sharply again at 14B (75%). This indicates that, when visual discriminability approaches the perceptual threshold of the encoder, increasingly strong language priors override ambiguous visual evidence.

We argue that this paradoxical scaling behaviour arises from four interacting factors: (i) efficiency-oriented patch compression that erodes stroke-level detail near perceptual limits; (ii) the lack of explicit modeling of glyph-internal structure; (iii) scale-dependent amplification of language priors during decoding; and (iv) next-token likelihood objectives that favour statistically fluent continuations over visually faithful reproduction.

Our contributions are threefold. First, we empirically establish a *Visual Faithfulness Paradox* in multilingual OCR: while English improves monotonically with scale, Chinese exhibits a *U-shaped glyph error curve*, where the 8B model achieves the best visual–language balance but the 14B model increasingly substitutes visually plausible yet in-

<sup>1</sup>Page URL: [https://en.wikipedia.org/wiki/Isotopes\\_of\\_hydrogen](https://en.wikipedia.org/wiki/Isotopes_of_hydrogen)

<sup>2</sup><https://github.com/will-be-released-after-review>

correct characters. Second, we show that Japanese occupies an *intermediate regime*: kanji remain largely stable, whereas visually similar kana continue to show prior-driven variability rather than catastrophic collapse. Third, we provide a *cross-layer explanation* linking these behaviours to visual signal-to-noise limits and language-prior dominance, clarifying why larger VLMs may become more “intelligent” yet less visually faithful. We situate these findings within the broader development of structure-aware OCR frameworks such as LLM-OCR-SWPC (Weigang and Costa, 2026).

## 2 Related Work

As this section reviews, mainstream VLM and OCR approaches are largely evaluated on alphabetic scripts (primarily English), overlooking the fundamental visual and structural differences of logographic systems such as CJK. This gap motivates the present study, which examines how efficiency-oriented designs interact with character-level visual fidelity in ideographic scripts.

### 2.1 Efficiency-Oriented Visual Token Optimization

A number of recent VLMs have introduced novel mechanisms to improve visual token efficiency. DeepSeek-OCR employs a deep visual encoder to aggressively compress text-dense images, achieving up to 20× token reduction while maintaining long-context OCR capability (Wei et al., 2025). This line of work demonstrates the feasibility of large-scale compression for document images, though its impact on fine-grained glyph structure remains underexplored. The InternVL series emphasizes scalable multimodal alignment through native multimodal pre-training and Variable Visual Position Encoding (V2PE), enabling improved handling of long sequences and multiple images (Chen et al., 2024b; Zhu et al., 2025; Lu et al., 2025).

Qwen2-VL introduces a *Naive Dynamic Resolution* mechanism, dynamically allocating visual tokens based on image content, together with Multimodal Rotary Position Embedding (M-RoPE) to unify positional modeling across text, image, and video modalities (Wang et al., 2024; Bai et al., 2025). This content-aware allocation represents a typical efficiency-driven approach.

Other models pursue more extreme reductions (Liu and Qiu, 2025). LLaVA-OneVision compresses each image frame into a single visual token

via modality pre-fusion, substantially shortening input sequences for video understanding (Zhang et al., 2025). Lightweight models such as Monkey prioritize inference speed and deployment under resource constraints by aggressively limiting visual tokens (Chen et al., 2024a). From a different perspective, Donut reformulates document understanding as an OCR-free encoder–decoder problem, bypassing explicit OCR tokenization altogether (Kim et al., 2022).

Collectively, these approaches share a common premise: reducing the number or complexity of visual tokens generally leads to stable or improved downstream performance. This assumption is largely supported by evaluations on alphabetic-script benchmarks such as TextVQA (Singh et al., 2019), DocVQA (Mathew et al., 2021), and ChartQA (Masry et al., 2022).

### 2.2 Limitations: Neglecting Script Heterogeneity

Despite their empirical success, existing efficiency-oriented designs in VLMs exhibit a fundamental limitation: they insufficiently account for heterogeneity across writing systems, particularly in how visual information supports glyph-level identity.

First, current evaluation benchmarks are overwhelmingly alphabet-centric. English-dominated datasets are characterized by relatively simple glyph shapes and a loose coupling between visual form and semantic identity (Singh et al., 2019; Mathew et al., 2021; Masry et al., 2022). Under such conditions, losses in local visual detail caused by token compression or resolution reduction can often be compensated by strong language priors. This bias reinforces the widespread assumption that aggressive visual compression is largely benign, or even desirable, for multimodal OCR (Wang et al., 2024; Zhang et al., 2025). However, such conclusions are drawn primarily from alphabetic scripts and do not necessarily generalize beyond them.

Second, mainstream visual encoders remain largely script-agnostic. Patch-based and dynamically compressed representations do not differentiate between low-complexity alphabetic characters and high-density glyphs found in logographic or mixed writing systems (Chen et al., 2024b; Wang et al., 2024; Wei et al., 2025). The research on *CJK character image processing threshold* demonstrates that once stroke density exceeds a pixel-level limit, irreversible stroke fusion occurs, result-

ing in abrupt increases in structural entropy and ambiguity (Weigang et al., 2025). Uniform patchification or compression strategies therefore introduce systematic structural noise when applied to dense glyphs, even if they remain effective for alphabetic text (Weigang et al., 2024). Mixed scripts such as Japanese with kana and kanji, partially mitigate this effect, but do not eliminate it.

Third, existing evaluations rarely assess *visual faithfulness* at the glyph level (Li et al., 2025). In ideographic and mixed scripts, models can produce semantically plausible outputs by relying on language priors, even when the visually grounded character recognition is incorrect, for example, through substitutions between visually similar or script-related characters. Such glyph-level hallucinations are largely invisible under conventional task-level accuracy or sequence similarity metrics, masking systematic failures of visual grounding.

Finally, current architectures lack explicit mechanisms for balancing language priors against character-structure evidence (Weigang and Brom, 2025). Large language decoders are optimized for next-token prediction and naturally favor high-frequency lexical patterns. When visual signals are degraded or ambiguous, this bias can override correct but low-probability glyph realizations (Wei et al., 2025; Zhu et al., 2025). Existing VLMs do not incorporate intermediate representations that explicitly encode character structure, such as stroke- or component-aware features, leaving visual evidence underconstrained during decoding (Weigang et al., 2024). As model scale increases, this imbalance can intensify rather than diminish.

This work challenges the prevailing assumption that improvements in visual token efficiency and model scale generalize uniformly across writing systems. We show that for glyph-dense and structurally complex scripts, efficiency-oriented design choices can induce a non-monotonic trade-off between model scale and visual fidelity, giving rise to the *Visual Faithfulness Paradox*. By foregrounding script heterogeneity and glyph-level faithfulness as first-class concerns, our analysis positions visual fidelity as an independent and necessary evaluation dimension for multilingual, OCR-oriented VLMs. Rather than opposing efficiency or scale, this study highlights the need for structure-aware constraints that ensure visual evidence remains competitive with language priors as models grow larger.

## 氢的同位素 [编辑]

条目 讨论 汉 漢 大陆简体 ▼

维基百科，自由的百科全书

氢（原子量：1.00794(7)）共有7个已知同位素，质量数介于1–7之间，其中有2个是稳定的，其他都具有放射性。天然存在的氢同位素有3个，分别是稳定的氕（<sup>1</sup>H）、氘（<sup>2</sup>H）和具放射性的氚（<sup>3</sup>H），其中<sup>3</sup>H在自然界中仅痕量存在，为宇宙射线所产生。另外四个同位素（<sup>4</sup>H到<sup>7</sup>H）都不出现在自然界中，只有在实验室制造出来过，且半衰期都少于10<sup>-21</sup>秒，极为不稳定。氢是唯一一个天然同位素各拥有不同名称的元素。虽然在一些有关放射性的早期文献中，一些属于自然界中三大衰变链的放射性核素也有自己专属的名称和化学符号，但是今日已经鲜少使用了。

Figure 2: Screenshot of the Chinese Wikipedia article Isotopes of Hydrogen, illustrating the coexistence of simplified and traditional Chinese characters (e.g., “汉” vs. “漢”) in real-world documents.

## 3 Experimental Data, Metrics and Results

This section presents a systematic evaluation of the InternVL3 model family on OCR tasks across three writing systems: English, Japanese, and Chinese. Rather than treating these as linguistic categories, the experiments focus on their distinct *script-level visual and structural properties*. The central question is whether scaling VLMs improves visual fidelity uniformly, or whether scaling effects are fundamentally modulated by differences in glyph density, character complexity, and visual–semantic coupling that are intrinsic to each writing system, independent of linguistic content.

### 3.1 Data Preparation

Text corpora and corresponding screenshots were collected on November 30 of 2025 from aligned Wikipedia articles on *Isotopes of Hydrogen* in English, Chinese, and Japanese. The selected passage provides structurally comparable content across writing systems while preserving script-specific visual characteristics. The reference texts contain 264 English tokens, 223 Chinese tokens, and 235 Japanese tokens, respectively. For each language, the first-page screenshot of the article was used as the OCR input. The Chinese screenshot is shown in Figure 2. The purpose is not to require linguistic understanding of Chinese, but to provide a direct visualization of glyph-level complexity and stroke density faced by VLM-based OCR systems.

### 3.2 Experimental Setup: Models and Metrics

To evaluate whether scaling effects persist under realistic OCR conditions, we conduct experiments on the InternVL3 model family (1B, 2B, 8B, and

14B) using aligned Wikipedia screenshots in English, Chinese, and Japanese. For each combination of model scale and target language, OCR inference is repeated 100 times on identical screenshots under stochastic decoding. The raw OCR outputs—comprising 1,200 samples generated by four model scales on three writing systems (see our GitHub repository for the complete replication protocol). This sampling-based methodology enables robustness analysis by systematically disentangling genuine visual fidelity effects from stochastic generation variance.

OCR performance is evaluated using three complementary measures: (i) sequence similarity to ground-truth text, (ii) character error rate (ER), defined as one minus similarity, and (iii) token-length ratio, reflecting output stability under each model’s native tokenizer. In addition to aggregate statistics, a targeted glyph-level diagnostic is conducted on a mixed simplified–traditional Chinese phrase to expose structural hallucinations that may be masked by sequence-level metrics.

We adopt the Ratcliff–Obershelp similarity measure  $S(T_{gt}, T_{gen}) \in [0, 1]$  (Ratcliff, 1981), and define the error rate as:

$$ER = 1 - S(T_{gt}, T_{gen}). \quad (1)$$

Token efficiency is quantified as:

$$R_{\text{token}} = \frac{N_{\text{gen}}}{N_{\text{gt}}}, \quad (2)$$

where  $N_{\text{gen}}$  and  $N_{\text{gt}}$  denote generated and ground-truth token counts, respectively. This Monte Carlo-style evaluation yields full error distributions rather than single-point estimates, allowing direct comparison of stability and scaling behavior across writing systems.

### 3.3 Scaling Behaviors Across Writing Systems

Figure 3 illustrates the performance distribution (e.g., sequence similarity) across writing systems and model scales, revealing starkly divergent scaling profiles. While aggregate OCR metrics improve with model size, the trajectory and underlying mechanisms differ fundamentally across scripts.

#### 3.3.1 English: Monotonic Scaling as an Aligned Baseline

Table 2 reports OCR performance on English text. As model size increases from 1B to 14B, all evaluation metrics exhibit strictly monotonic improvement. The mean predicted token length converges

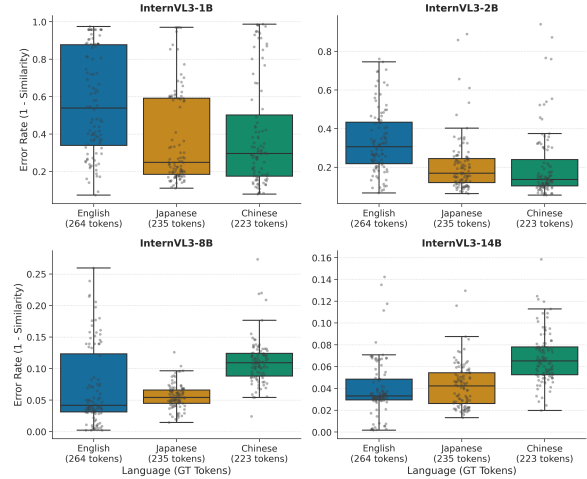


Figure 3: Error-rate distributions across writing systems and model scales for InternVL3. While overall OCR robustness improves with scale, the stability and qualitative behavior of gains differ markedly across scripts.

Model	Tokens ( $\mu \pm \sigma$ )	Similarity	Error Rate
1B	279.81±287.19	0.43	0.57
2B	326.25±102.40	0.67	0.33
8B	269.00±20.73	0.93	0.07
14B	268.02±11.34	0.96	0.04

Table 2: OCR Performance on English Text.

from 279.81 to 268.02 (reference: 264), while the standard deviation decreases from 287.19 to 11.34, indicating highly stable outputs. Sequence similarity increases from 0.43 to 0.96, with the corresponding error rate dropping from 0.57 to 0.04.

These results establish English as an *ideal baseline* for VLM-based OCR: in alphabetic writing systems, visual transcription and language modeling objectives are strongly aligned (Table 1). Increased model capacity simultaneously enhances perceptual stability and linguistic accuracy, yielding predictable and monotonic scaling behavior.

#### 3.3.2 Chinese: Full Manifestation of the Visual Faithfulness Paradox

Chinese OCR results provide the definitive evidence for the Visual Faithfulness Paradox, exposing a fundamental contradiction between macro-level fluency and glyph-level visual fidelity. As shown in Figure 3, aggregate indicators suggest continuous improvement from 1B to 14B, including higher sequence similarity (up to 0.93) and lower overall error rates (down to 0.07), approaching English-level performance.

However, direct inspection of model outputs reveals substantial glyph-level failures that are in-

Error Type	2B	8B	14B
Semantic / Glyph Collapse	5/5	0/5	0/5
Empty output	0/5	0/5	1/5
Simplified–Traditional Error	1/5	2/5	5/5
Glyph Confusion / Neologism	5/5	4/5	3/5
Extraneous Symbol Addition	3/5	5/5	0/5
Format Inference Error	4/5	3/5	2/5
Minor Punctuation Error	5/5	5/5	5/5

Table 3: Error type distribution across model scales on Chinese OCR (5 samples per model).

visible to surface metrics. We therefore conduct following two complementary analyses.

**(1) Manual Analysis of Five Representative Outputs.** Consistent with the Japanese analysis, manual inspection reveals scale-dependent qualitative error regimes (see Table 3 and Table 4).

*2B: Severe visual insufficiency.* Errors are dominated by glyph-structure collapse and semantic confusion, including malformed or invented characters (e.g., “氣”), confusion among hydrogen isotopes (“气/氕/氘”), and unstable recognition of scientific notation.

*8B: Optimal vision–language balance.* Major glyph collapse disappears and factual correctness largely recovers. New error peaks emerge in the form of extraneous symbol insertion (e.g., ~) and occasional simplified–traditional inconsistencies. These errors reflect transitional instability as improved visual perception begins to capture interface-level symbols without fully stabilized decoding.

*14B: Language-prior-dominated substitution.* Severe visual noise vanishes, but errors shift toward systematic, knowledge-driven glyph replacement. Rare but visually faithful characters are overridden by more frequent alternatives (e.g., “氘” misrecognized as “氙”). Simplified–traditional conversion errors reach 100%, indicating that strong language priors normalize visual input toward dominant lexical distributions.

## (2) Targeted Phrase Analysis over 20 Samples.

We further analyze 20 repeated outputs of the critical phrase “条目讨论汉 大陆简体,” focusing on edit-distance similarity and exact-match behavior.

Contrary to aggregate trends, glyph-level accuracy follows a clear U-shaped trajectory. The 8B model achieves peak performance, while the 14B model exhibits pronounced degradation. As shown in Figure 1, the perfect match error rate for the character “漢” rises to 75% at 14B, compared to

Metric (%)	2B	8B	14B
Avg. Edit Similarity	85.1	<b>92.9</b>	82.2
Perfect Match Rate	30.0	<b>50.0</b>	20.0
High Similarity ( $\geq 90\%$ )	70.0	<b>95.0</b>	65.0
Complete Failure	0.0	0.0	<b>5.0</b>
S–T Conversion Error	20.0	15.0	<b>60.0</b>

Obs.: A paired-sample t-test confirms worse performance than the 8B model ( $t(19) = 2.94, p < 0.01$ )

Table 4: Performance comparison on the mixed-script Chinese phrase (N=20).

only 46% at 8B.

The dominant failures are systematic, structured substitutions (e.g., “汉汉”, “汉语”)—glyph-level hallucinations where strong language priors override visual evidence. This degradation at the 14B scale is statistically significant: a paired-sample t-test confirms worse performance than the 8B model ( $t(19) = 2.94, p < 0.01$ ). The original outputs for the 20 samples, along with the calculation and statistics of their edit-distance similarity to the reference text, are provided in Appendix A Table 6.

### 3.3.3 Japanese: Mixed-Regime and Partially Aligned Scaling

Japanese departs from strictly monotonic scaling, but it does not exhibit the severe glyph-level breakdown observed in Chinese. As shown in Figure 4, English OCR variance decreases smoothly with model size, whereas Japanese token-length variance remains elevated at small scales and re-emerges at 14B, indicating persistent decoding instability in kana–kanji mixtures.

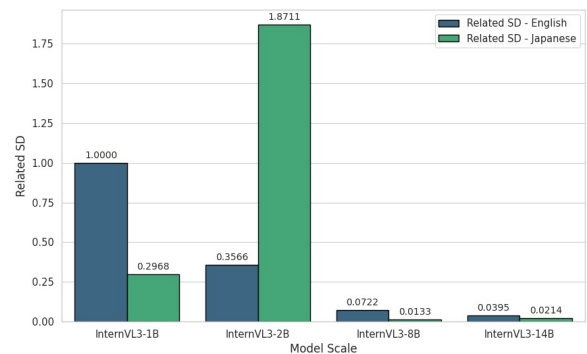


Figure 4: Relative standard deviation of predicted token length for English and Japanese OCR outputs across InternVL3 model scales, normalized to the English 1B baseline.

At the system level, the 8B model achieves the best overall balance, with near-reference length, high edit similarity, and low variance, see Table 5. The 14B model yields slightly higher average

Metric (%)	2B	8B	14B
Avg. Edit Similarity	73.5	<b>89.6</b>	89.4
Perfect Match Rate	15.0	<b>35.0</b>	30.0
High Similarity ( $\geq 90\%$ )	25.0	<b>70.0</b>	65.0
Complete Failure	0.0	0.0	0.0
Language-Prior Hallucinations	35.0	15.0	<b>25.0</b>

Obs.: No significant difference between 8B and 14B models ( $t(19) = 0.35, p = 0.73$ ), both outperform 2B ( $p < 0.001$ ).

Table 5: Performance comparison on the Japanese reference phrase (N=20).

similarity but noticeably larger variance, suggesting that stronger language priors begin to perturb visually grounded decoding.

To verify whether this instability reflects qualitative regime shifts, we conducted a fine-grained analysis of a canonical Japanese phrase, “出典: フリ百科事典『ウイキペディア (Wikipedia)』”, comparing 20 sampled OCR outputs across scales (Table 5 and Appendix C Table 8). Both 8B and 14B significantly outperform 2B ( $p < 0.001$ ), but there is no significant difference between 8B and 14B ( $t(19) = 0.35, p = 0.73$ ). However, the error types differ. The 2B model suffers from semantic drift and fabrication due to weak visual grounding. The 8B model produces mostly faithful transcriptions with only minor formatting deviations. By contrast, the 14B model introduces isolated but systematic *language-prior substitutions*, alongside stylistic “polishing.” These include knowledge-based replacements that are not visually supported in the input, illustrating the same prior-vision tension seen in Chinese, though in milder form, see Appendix C Table 8.

A key observation is that instability concentrates in *kana* sequences rather than kanji. Although kana have low stroke complexity, they are visually similar (e.g., “フ/ブ/ヴ/ヌ/ワ”), and thus only weakly discriminable at small scales, exactly where language priors become competitive. It yields a *mixed-regime behavior*: scaling improves robustness and fluency, but prior dominance introduces measurable variability before convergence. This pattern supports our broader claim that orthographic structure modulates the balance between visual fidelity and linguistic priors in VLM decoding.

### 3.4 Cross-Lingual Synthesis

Our cross-lingual analysis shows that the *Visual Faithfulness Paradox* is strongly moderated by orthographic structure. Alphabetic English follows a near-ideal scaling pattern, with larger models jointly improving linguistic fluency and visual ac-

curacy. Japanese, which mixes kanji and visually subtle kana, displays moderate instability at small scales but largely stabilizes as capacity increases, reflecting only a partial tension between visual form and prior-driven decoding. Chinese, however, exhibits the most severe manifestation: glyph-level fidelity *declines* beyond a critical scale despite continued gains in aggregate OCR metrics.

This decline is tightly linked to character complexity and perceptual resolution. Learning-based OCR studies show a stroke-density threshold effect in CJK scripts, where complex glyphs become visually ambiguous at common rendering sizes (e.g.,  $>12$  strokes at 16 px) (Weigang et al., 2025). For characters such as “漢” (14 strokes), the visual encoder thus operates near or beyond its effective resolution limit, yielding weak or unstable glyph features. At moderate scales, models rely primarily on this degraded but still constraining visual signal. As scale grows, however, stronger language priors increasingly dominate decoding, completing ambiguous glyphs via lexical frequency or semantics rather than perception. What begins as helpful compensation therefore becomes a systematic source of glyph substitution, capturing the core of the paradox, see Table 1.

An optimal operating regime thus emerges for mid-scale models (e.g.,  $\sim 8B$ ), where visual evidence remains decisive and language priors offer only limited support. More broadly, these results show that aggregate OCR metrics are insufficient: glyph-level visual fidelity, especially near perceptual complexity thresholds, must be treated as a distinct and essential evaluation dimension for multilingual VLMs, also see Appendix B Table 7.

## 4 Discussion: Visual Signal-to-Noise Evolution

To explain the non-uniform scaling behaviors observed across writing systems, we introduce *Visual Signal-to-Noise Ratio (VSNR)* as a unifying explanatory framework (Sinha et al., 2007; Hendricks et al., 2021). VSNR characterizes the internal competition between meaningful glyph-structure signals extracted by the visual encoder and noise arising from resolution limits, compression artifacts, and representational uncertainty.

VSNR is not defined at the pixel level; rather, it reflects the balance within a model’s internal visual feature space between deterministic structural representations and ambiguous or unstable feature

549 components. As model scale increases, this bal-  
550 ance shifts, altering the relative influence of visual  
551 evidence and language priors during multimodal  
552 OCR decoding. Conceptually, VSNR can be ap-  
553 proximated as:

$$554 \text{VSNR} \approx \frac{\mathbb{E}[\|\mathbf{f}_g\|^2]}{\mathbb{E}[\|\mathbf{f}_n\|^2]}, \quad (3)$$

555 where  $\mathbf{f}_g$  denotes glyph-correlated visual features  
556 and  $\mathbf{f}_n$  denotes noise-dominated components.

557 *Stage A: Visual Failure (1B–2B).* At small scales,  
558 the visual encoder fails to extract stable glyph rep-  
559 resentations, yielding truncation, unreadable sym-  
560 bols, or visually unrelated characters. Error rates  
561 exceed 79–99%, and VSNR remains near zero, pre-  
562 venting reliable decoding by either visual evidence  
563 or language priors. This regime is most evident  
564 for dense logographic characters but also affects  
565 low-resolution kana strings through uncontrolled  
566 phonetic variation.

567 *Stage B: Visual Dominance (8B).* At interme-  
568 diate scale, VSNR exceeds unity, allowing visual  
569 evidence to dominate decoding. Error rates drop  
570 to 46% for Chinese, with remaining mistakes lim-  
571 ited to rare visual neighbor confusions. Japanese  
572 OCR likewise reaches peak stability, though resid-  
573 ual variance persists for kana sequences with low  
574 inter-character distinctiveness, reflecting partial de-  
575 coupling between visual form and lexical identity.

576 *Stage C: Language-Prior Dominance (14B).* Be-  
577 yond a critical scale, error rates rebound sharply  
578 to 75% for Chinese. Although visual representa-  
579 tions remain strong, decoding becomes dominated  
580 by language priors, producing systematic seman-  
581 tic substitutions (e.g., “漢” → “汉” or “汉语”). In  
582 Japanese, this dominance manifests differently: in-  
583 stead of semantic collapse, strong priors amplify  
584 phonetic smoothing, yielding clustered kana sub-  
585 stitutions (e.g., “フリ” → “フリワ/フリラ/フリ  
586 カ”) despite visually correct kanji.

587 Together, these stages reveal a nonlinear tran-  
588 sition from visual insufficiency to visual domi-  
589 nance and ultimately to language-prior hegemony.  
590 The resulting failures differ by script: threshold-  
591 induced ambiguity for dense logographs and low-  
592 discriminability drift for phonetic kana. Perform-  
593 ance degradation at large scales thus reflects not  
594 insufficient capacity, but a systematic reweighting  
595 of decision cues under scaling. This diagnosis es-  
596 tablishes the theoretical basis for structure-aware  
597 corrective frameworks, such as LLM-OCR-SWPC  
598 (Weigang and Costa, 2026), aimed at preserving  
599 glyph-level visual faithfulness under model scaling.

## 5 Conclusion 600

601 The *Visual Faithfulness Paradox* is best understood  
602 as a continuum governed by orthographic structure:  
603 it intensifies when visual discriminability is low  
604 and linguistic constraints are strong. Thus, Chinese  
605 characters show the sharpest paradox, Japanese  
606 kana reveal moderate creative errors, while the  
607 Latin alphabet remains largely stable.

608 This paper demonstrates that larger VLMs may  
609 become *less* visually faithful even when their over-  
610 all OCR fluency improves. Across InternVL3 (1B–  
611 14B), English performance scales monotonically as  
612 expected. In contrast, Chinese exhibits a *U-shaped*  
613 glyph outcome: the perfect-match error rate is high  
614 at 2B (79%), lowest at 8B (46%), but rises again  
615 to 75% at 14B, reflecting a shift toward language-  
616 prior-dominated decoding in which visually am-  
617 biguous yet semantically important characters are  
618 replaced by statistically plausible alternatives (e.g.,  
619 simplified for traditional forms). Japanese occupies  
620 an intermediate position: kanji remain largely sta-  
621 ble, while visually similar kana continue to induce  
622 prior-driven variability and creative substitutions.

623 We interpret these effects through a *Visual*  
624 *Signal-to-Noise Ratio (VSNR)* lens. Scaling ini-  
625 tially improves VSNR, but beyond a critical point  
626 language priors overpower weak visual evidence,  
627 shifting decoding from perception-driven transcrip-  
628 tion to prior-driven completion—the essence of the  
629 paradox.

630 Therefore, “bigger is better” does not univer-  
631 sally hold for multilingual OCR. Robust scaling  
632 requires structure-aware designs that respect script-  
633 specific perceptual limits. This work highlights  
634 that “vision” in VLMs is not merely a pixel en-  
635 coder, but a dynamic equilibrium between percep-  
636 tual evidence and linguistic priors whose balance  
637 is script-dependent. Promising directions include  
638 adaptive visual resolution for dense glyphs, ex-  
639 plicit glyph-structure anchoring (e.g., SWPC), and  
640 benchmarks that measure multilingual visual faith-  
641 fulness beyond edit distance. Although our case  
642 studies focus on Chinese and Japanese, the same  
643 mechanisms are likely to extend to other scripts  
644 with low visual separability or high ligaturing com-  
645 plexity (e.g., Vietnamese Hán-Nôm, Arabic, or  
646 Devanagari), suggesting a broader research agenda.  
647 Restoring strong visual grounding is therefore es-  
648 sential for the next generation of large-scale vi-  
649 sion–language models.

## 650 Limitations

651 This study has several limitations. First, our analy-  
652 sis is based on a deliberately focused dataset con-  
653 sisting of three screenshots from English, Chinese,  
654 and Japanese Wikipedia pages. This case-study de-  
655 sign enables controlled, fine-grained observation of  
656 glyph-level behavior (e.g., distinguishing “汉” from  
657 “漢”, but necessarily limits statistical breadth and  
658 visual diversity. Second, all results are obtained  
659 from a single model family, InternVL3, evaluated  
660 across four parameter scales (1B–14B). Although  
661 the observed “Visual Faithfulness Paradox” appears  
662 internally consistent, confirming whether it general-  
663 izes across architectures, training regimes, and data  
664 sources requires broader replication. Finally, our  
665 interpretation is grounded in the visual–semantic  
666 binding properties of logographic and mixed-script  
667 systems. While the analyses of Chinese characters  
668 and Japanese kana likely extend to other visually  
669 confusable scripts, a comprehensive cross-script  
670 survey and larger-scale evaluation on varied real-  
671 world imagery remain important directions for fu-  
672 ture work.

## 673 Ethical Considerations

674 This study uses only publicly available Wikipedia  
675 screenshots and model-generated text. No personal  
676 or sensitive data are included, and no content falls  
677 into harmful or offensive categories.

## 678 Acknowledgements

679 We thank the developers of publicly available LLM  
680 services (ChatGPT, DeepSeek, Grok) for assistance  
681 with language refinement. We are also grateful to  
682 the OpenGVLab team for open-sourcing the In-  
683 ternVL3 model series, which made this controlled  
684 scaling study possible.

## 685 References

686 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
687 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie  
688 Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical  
689 report. *arXiv preprint arXiv:2502.13923*.

690 Wei Chen, Zhiyuan Li, and Shuo Xin. 2024a. Om-  
691 nivlm: A token-compressed, sub-billion-parameter  
692 vision-language model for efficient on-device infer-  
693 ence. *arXiv preprint arXiv:2412.11475*.

694 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo  
695 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,

Xizhou Zhu, Lewei Lu, et al. 2024b. Intern vl: Scal-  
ing up vision foundation models and aligning for  
generic visual-linguistic tasks. In *2024 IEEE/CVF-  
CVPR*, pages 24185–24198. IEEE. 696  
697  
698  
699

Lisa Anne Hendricks, John Mellor, Rosalia Schneider,  
Jean-Baptiste Alayrac, and Aida Nematzadeh. 2021.  
Decoupling the role of data, attention, and losses in  
multimodal transformers. *Transactions of the Associ-  
ation for Computational Linguistics*, 9:570–585. 700  
701  
702  
703  
704

Geewook Kim, Teakgyu Hong, Moonbin Yim,  
JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Won-  
seok Hwang, Sangdoon Yun, Dongyoon Han, and  
Seunghyun Park. 2022. Ocr-free document under-  
standing transformer. In *European Conference on  
Computer Vision*, pages 498–517. Springer. 705  
706  
707  
708  
709  
710

Zongqian Li, Yinhong Liu, Yixuan Su, and Nigel Col-  
lier. 2025. Prompt compression for large language  
models: A survey. In *Proc. of 2025 NAACL-HLT  
(vol. 1)*, pages 7182–7195. 711  
712  
713  
714

Fanfan Liu and Haibo Qiu. 2025. Context cascade com-  
pression: Exploring the upper limits of text compres-  
sion. *arXiv preprint arXiv:2511.15244*. 715  
716  
717

Dongchen Lu, Yuyao Sun, Zilu Zhang, Leping Huang,  
Jianliang Zeng, Mao Shu, and Huo Cao. 2025.  
Internvl-x: Advancing and accelerating internvl se-  
ries with efficient visual token compression. *arXiv  
preprint arXiv:2503.21307*. 718  
719  
720  
721  
722

Ahmed Masry, Do Long, Jie Tan, Shafiq Joty, and Ena-  
mul Hoque. 2022. ChartQA: A benchmark for ques-  
tion answering about charts with visual and logical  
reasoning. *arXiv preprint arXiv:2203.10244*. 723  
724  
725  
726

Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawa-  
har. 2021. DocVQA: A dataset for document visual  
question answering. In *Proceedings of the IEEE/CVF  
Winter Conference on Applications of Computer Vi-  
sion*, pages 2200–2209. 727  
728  
729  
730  
731

Roger Ratcliff. 1981. A theory of order relations in per-  
ceptual matching. *Psychological review*, 88(6):552. 732  
733

Amanpreet Singh, Vivek Natarajan, Meet Shah,  
Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus  
Rohrbach. 2019. Towards VQA models that can  
read. In *Proceedings of the IEEE/CVF-CVPR*, pages  
8317–8326. 734  
735  
736  
737  
738

Pawan Sinha, Benjamin Balas, Yuri Ostrovsky, and  
Richard Russell. 2007. Face recognition by humans:  
Nineteen results all computer vision researchers  
should know about. *Proceedings of the IEEE*,  
94(11):1948–1962. 739  
740  
741  
742  
743

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-  
hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin  
Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhanc-  
ing vision-language model’s perception of the world  
at any resolution. *arXiv preprint arXiv:2409.12191*. 744  
745  
746  
747  
748

749 Haoran Wei, Yaofeng Sun, and Yukun Li. 2025.  
750 Deepseek-ocr: Contexts optical compression. *arXiv*  
751 *preprint arXiv:2510.18234*.

752 Li Weigang and Pedro Carvalho Brom. 2025. [Paradox](#)  
753 [of poetic intent in back-translation: Evaluating the](#)  
754 [quality of large language models in chinese transla-](#)  
755 [tion](#). *Frontiers of Information Technology & Elec-*  
756 *tronic Engineering*, 26(11):2176–2203.

757 Li Weigang and Joao Paulo Vieira Costa. 2026. Llm-  
758 ocr-swpc: Hanzi’s cognitive advantage and token  
759 efficiency in multimodal ai. In *Proceedings of 2026*  
760 *IEEE 18th International Conference on Computer*  
761 *Research and Development*. IEEE.

762 Li Weigang, M. C. Marinho, and Danilo Lelin Li. 2024.  
763 [Six-writings multimodal processing with pictopho-](#)  
764 [netic coding to enhance chinese language models](#).  
765 *Frontiers of Information Technology & Electronic*  
766 *Engineering*, 25(1):84–105.

767 Li Weigang, Rafael Marconi Ramos, Pedro Carvalho  
768 Brom, and Danilo Lelin Li. 2025. Threshold study  
769 for hanzi image recognition: Defining character and  
770 component limits in chinese, japanese, and korean  
771 script processing. *International Journal of Asian*  
772 *Language Processing*, 35(1):2450011.

773 Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang Feng.  
774 2025. Llava-mini: Efficient image and video large  
775 multimodal models with one vision token. *arXiv*  
776 *preprint arXiv:2501.03895*.

777 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu,  
778 Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan,  
779 Weijie Su, Jie Shao, et al. 2025. Internv13: Ex-  
780 ploring advanced training and test-time recipes for  
781 open-source multimodal models. *arXiv preprint*  
782 *arXiv:2504.10479*.

## 783 **Appendix**

784 **Appendix A: Complete OCR outputs and**  
785 **edit-distance similarity for the Chinese phrase**  
786 **across model scales (N=20)**

787 **Appendix B: A Unified Theoretical Parallel**  
788 **between the Paradox of Poetic Intent and the**  
789 **Paradox of Visual Fidelity**

790 **Appendix C: Complete OCR outputs and**  
791 **edit-distance similarity for the Japanese phrase**  
792 **across model scales (N=20)**

2B Model		8B Model		14B Model	
<b>OCR Output</b>	<b>Sim.</b>	<b>OCR Output</b>	<b>Sim.</b>	<b>OCR Output</b>	<b>Sim.</b>
条目讨论汉 大陆简体跑	0.900	条目讨论汉 大陆简体	1.000	(空输出)	0.000
条目讨论汉 大陆简体海条	0.643	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900
目					
条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.643
				体chuan	
条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.900	条目讨论汉语 大陆简体	0.900
题目讨论汉 大筒体陆筒	0.400	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.900
体					
条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体 X	0.800
条目讨论汉	0.444	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.778
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.889	条目讨论汉 大陆简体	0.900
条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.800
条目讨论汉 大陆简体	0.889	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	1.000
条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.778
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	0.900
条目讨论汉 大	0.889	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000
条目讨论汉 大陆简体	0.778	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	0.900
条目讨论汉: 汉 大陆简体	0.900	条目讨论汉 大陆简体	0.667	条目讨论汉 大陆简体	0.900
栏目讨论汉 大陆简体	0.667	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000
条目讨论汉 大陆简体	0.900	条目讨论汉 大陆简体	1.000	条目讨论汉 大陆简体	1.000
<b>Mean ± SD</b>	<b>0.851 ± 0.145</b>	<b>Mean ± SD</b>	<b>0.929 ± 0.020</b>	<b>Mean ± SD</b>	<b>0.822 ± 0.208</b>
		Wikipedia Source Phrase:	条目讨论汉 大陆简体		

*Obs.:* A paired-sample  $t$ -test confirms the 14B model performs significantly worse than the 8B model ( $t(19) = 2.94, p < 0.01$ ).

Table 6: **Appendix A:** Complete OCR outputs and edit-distance similarity for the Chinese phrase across model scales (N=20).

Dimension	Paradox of Poetic Intent (Weigang and Brom, 2025)	Paradox of Visual Fidelity
Task Level	Sentence- and discourse-level back-translation of poetry	Character-level OCR in multimodal vision-language models
Observed Phenomenon	High BLEU scores coexist with loss of metaphor and cultural intent	High linguistic confidence coexists with glyph-level visual errors
Surface Metric Optimized	Lexical overlap (BLEU)	Token-level similarity / semantic plausibility
Core Paradox	Surface textual fidelity $\nRightarrow$ poetic understanding	Language fluency $\nRightarrow$ visual faithfulness
Typical Failure Mode	Literal reproduction with shallow interpretation	Semantic substitution overriding ambiguous visual evidence
Trigger Condition	Dense metaphor, cultural allusion, poetic ambiguity	High stroke density, isolated glyphs, script-level ambiguity
Role of Model Scaling	Larger models amplify memorization and prior-driven generation	Larger models amplify language-prior dominance over vision
Underlying Mechanism	Optimization toward lexical probability over interpretive depth	Language-dominant decoding under degraded visual signal
Unifying Explanation	<b>Strong model priors dominate when task-relevant signal is weak or ambiguous</b>	
Signal Perspective	Low semantic signal-to-prior ratio	Low visual signal-to-noise ratio (VSNR)
Theoretical Implication	Evaluation based on surface metrics is fundamentally insufficient	Visual fidelity requires independent, structure-aware assessment
Proposed Direction	Semantics-aware evaluation of back-translation	Structure-aware OCR with glyph-level anchoring (LLM-OCR-SWPC) (Weigang and Costa, 2026)

Table 7: **Appendix B:** A Unified Theoretical Parallel between the Paradox of Poetic Intent and the Paradox of Visual Fidelity.

2B Model		8B Model		14B Model	
OCR Output	Sim.	OCR Output	Sim.	OCR Output	Sim.
出典: フリ百科事典『ウィキペディア (Wikipedia)』	1.000	出典: フリカド辞典『ウィキペディア (Wikipedia)』	0.885	出典: フリ百科事典「ウィキペディア (Wikipedia)」	0.962
出典: フリワイド百科事典『ウィキペディア(Wikipedia)』	0.885	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962	出典: フリ百科事典「ウィキペディア(Wikipedia)」	0.962
出典: フリレス百科事典『ウィキペディア(Wikipedia)』	0.885	出典: フリ百科事典『ウィキペディア』*(Wikipedia)*	0.769	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.923
出典: フリワルドセントスバディア (Wikipedia) ]	0.538	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.923	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000
出典: フリヒ百科事典『ウィキペディア (Wikipedia) 』	0.846	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリ百科事典「ウィキペディア(Wikipedia)」	0.962
出典: フリワルド百科事典[『ウィキペディア (Wikipedia) 』]	0.731	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962
出典: フリラ百科事典『ウィキペディア(Wikipedia)』	0.692	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリ百科事典「ウィキペディア (Wikipedia) 』	0.962
出典: フリバ百科事典『ウィキペディア(Wikipedia)』	0.885	出典: フリキ百科事典「ウィキペディア (Wikipedia) 』	0.885	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000
出典: フリマウンドイリア (Wikipedia) )	0.577	出典: フリカイ百科事典『ウィキペディア(Wikipedia)』	0.885	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000
出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリ百科事典『ウィキペディア』(Wikipedia)	0.885	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962
出自: フリリッチ百科事典『ウィキペディア (Wikipedia) 』	0.846	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.885
出典: フリリ百科事典『ウィキペディア (Wikipedia) 』	0.923	出典: フリキ百科事典『ウィキペディア (Wikipedia) 』	0.885	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000
出典: フリワルド事典『ウィキペディア (Wikipedia) 』	0.808	出典: フリキ百科事典『ウィキペディア (Wikipedia) 』	0.885	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962
出典: フリラ百科事典『ウィキペディア (Wikipedia) 』	0.923	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962	出典: フリ百科事典「ウィキペディア(Wikipedia)」	0.962
出典: フリワルドセ事典『ウィキペディア(Wikipedia)』	0.769	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリエ百科事典「ウィキペディア(Wikipedia)」	0.885
出典: フリワルド.ウィキメディアベス (Wikipedia) [1]	0.615	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962	出典: フリペディア「ウィキペディア(Wikipedia)」	0.731
出典: フリワイヤメディア (Wikipedia) [1]	0.577	出典: フリビル百科事典『ウィキペディア (Wikipedia) 』	0.885	出典: フリ百科事典「ウィキペディア (Wikipedia) 』	0.962
出典: フリワ百科事典『ウィキペディア(Wikipedia)』	0.885	出典: フリゾノ学事典『ウィキペディア (Wikipedia) 』	0.846	出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000
出典: フリ百科事典『ウィキペディア (Wikipedia) 』	1.000	出典: フリ百科事典『ウィキペディア』(Wikipedia)	0.885	出典: フリモ百科事典「ウィキペディア(Wikipedia)」	0.885
出典: フリワルド事典『ウィキペディア (Wikipedia) 』	0.808	出典: フリ百科事典『ウィキペディア(Wikipedia)』	0.962	出典: フリ百科事典「ウィキペディア(Wikipedia)」	0.962
Mean ± SD	0.735 ± 0.184	Mean ± SD	0.896 ± 0.071	Mean ± SD	0.894 ± 0.084

Wikipedia Source Phrase: 出典: フリ百科事典『ウィキペディア (Wikipedia) 』

Paired-sample *t*-tests show no significant difference between 8B and 14B models ( $t(19) = 0.35, p = 0.73$ ), both outperform 2B ( $p < 0.001$ ).

Table 8: **Appendix C**: Complete OCR outputs and edit-distance similarity for the Japanese phrase across model scales (N=20).