

Hot-Start from Pixels: Low-Resolution Visual Tokens for Chinese Language Modeling

Anonymous ACL submission


Abstract

Large language models typically represent Chinese characters as discrete index-based tokens, largely ignoring their visual form. For logographic scripts, visual structure carries semantic and phonetic information, which may aid prediction. We investigate whether low-resolution visual inputs can serve as an alternative for character-level modeling. Instead of token IDs, our decoder receives grayscale images of individual characters, with resolutions as low as 8×8 pixels. Remarkably, these inputs achieve 39.2% accuracy, comparable to the index-based baseline of 39.1%. Such low-resource settings also exhibit a pronounced *hot-start* effect: by 0.4% of total training, accuracy reaches above 12%, while index-based models lag at below 6%. Overall, our results demonstrate that minimal visual structure can provide a robust and efficient signal for Chinese language modeling, offering an alternative perspective on character representation that complements traditional index-based approaches.

1 Introduction

In Chinese, meaning arises not only through sequential context but also through the internal structure of characters. We human readers naturally pay attention to radicals, stroke layout, and overall shape when reading Chinese.

In contrast, the majority of mainstream models today process Chinese through sequences of symbolic character IDs, whereas English is typically tokenized into subword indices (Brown et al., 2020; Bommasani et al., 2021; Bender and Koller, 2020; Rust et al., 2021). While this index-based abstraction is effective for alphabetic writing systems, it may be suboptimal for Chinese.

This omission can have concrete consequences. Take the character 山 ( ‘mountain’) as an example: its shape resembles small mountain peaks, thereby immediately conveying meaning to a human reader. The contrast between human percep-

tion and index-based modeling is clear here: while humans can directly leverage visual form, 山 is represented in an index-based model as an abstract token ID, stripped of its shape. Similarly, characters such as 灭 (extinguish) and 火 (fire) differ clearly in their visual structure, yet these differences are not readily accessible to the model at early stages of learning. It is akin to assembling a jigsaw puzzle with the image erased: the components remain, but the visual cues that facilitate interpretation are missing.

As one of the most representative logographic systems, Chinese treats visual form not as auxiliary, but as an essential part of meaning construction, carrying semantic and phonetic information through its visual structure. Such visual intuition extends beyond isolated characters to contextual prediction. In weak linguistic contexts, human readers naturally rely on glyph patterns to predict the following text. This suggests that visual form may serve as a more fundamental cue for logographic language processing, especially when semantic signals are sparse.

This motivates us to ask: *can language models effectively process Chinese characters—or more generally, visually structured information—instead of index-based tokens?* If visual structure indeed plays such a central role for humans, then abstracting Chinese characters into index-based tokens may not merely be an engineering choice, but a fundamental modeling assumption.

This leads to a fundamental choice in representation: index-based tokenization that relies on contextual embeddings alone, or vision-based processing that extracts character shape and structure through a visual encoder (Poznanski et al., 2025). These two paths frame our experimental comparison.

This architectural distinction has implications for representation topology. Index-based embeddings initialize as unstructured points in feature space, requiring the model to discover character

084	relationships purely from co-occurrence statistics.	learning, even in low-resource settings.	136
085	In contrast, visual encoders provide built-in geo-	Research Questions. Our evaluation addresses	137
086	metrically meaningful spatial organization from	the following research questions:	138
087	the outset, potentially offering a structural prior	RQ1 (Visual Sufficiency): Can visual inputs of	139
088	that accelerates early learning. Importantly, such	Chinese characters alone suffice for character-level	140
089	representations also make model behavior more	prediction?	141
090	interpretable, as predictions can be traced back to	RQ2 (Early-Stage Dynamics): What are the learn-	142
091	salient regions and structural patterns in the image.	ing trajectories of vision-token models?	143
092	This question is intriguing given recent advances	RQ3 (Resolution Sensitivity): How does pre-	144
093	in applying VLMs to optical character recognition	dictive performance vary as image resolution de-	145
094	(OCR) and document understanding tasks. Sys-	creases from high-fidelity to near-minimal levels?	146
095	tems like DeepSeek-OCR (Wei et al., 2025) and	RQ4 (Spatial Robustness): Can vision-token mod-	147
096	Pix2Struct (Lee et al., 2023) demonstrate that tex-	els maintain accuracy when only partial character	148
097	tual content can be processed exclusively as im-	regions are visible?	149
098	ages. However, our work differs fundamentally	Main Contributions. Our main contributions are	150
099	from these transcription-focused models: rather	sixfold, structured around our research questions:	151
100	than transcribing visual glyphs back to symbols,		
101	we ask: can visual forms alone support linguistic	• For Methodology , we propose a vision-token	152
102	prediction—in particular, predicting the next char-	formulation for Chinese language modeling,	153
103	acter based on visual context? This shifts the focus	replacing completely index-based tokens with	154
104	from recognition to language modeling itself, from	character images processed through a visual	155
105	symbol reconstruction to linguistic reasoning.	encoder.	156
106	Our preliminary experiments would reveal that	• For RQ1 , we demonstrate that visual inputs	157
107	even heavily cropped or low-resolution character	alone achieve accuracy comparable to index-	158
108	images retain sufficient structural information for	based baselines (39.2% vs. 39.1%), confirm-	159
109	prediction. For instance, low-resolution or partially	ing that purely visual structure suffices for	160
110	cropped character images still allow models to iden-	character-level prediction.	161
111	tify the correct next character, reminiscent of how	• For RQ2 , we identify a <i>hot-start</i> effect: at	162
112	humans can still read small or degraded characters.	only 0.4% of total training regimen, visual	163
113	Similar observations have been reported in recent	models reach 12.3% accuracy—more than	164
114	studies (Li et al., 2025; Poznanski et al., 2025; Wei	double the baseline’s 5.8%.	165
115	et al., 2024). These early observations motivate a	• For RQ3 & RQ4 , we establish robustness un-	166
116	systematic study of the sufficiency of visual forms	der visual degradation: performance remains	167
117	in language modeling.	stable across resolutions from 8×8 to 96×96	168
118	Exploring this question brings significant prac-	pixels and under severe spatial cropping.	169
119	tical and scientific value. In resource-constrained	• For Explainable NLP , we analyze how visual	170
120	environments, models that use visual structure can	tokens offer inherent interpretability: embed-	171
121	extract meaningful patterns from limited data, po-	ding spaces organize by morphological simi-	172
122	tentially achieving faster convergence than index-	larity, and gradient analysis traces predictions	173
123	based alternatives. These considerations motivate	to salient pixel regions.	174
124	a systematic investigation across multiple dimen-	2 Methodology	175
125	sions: visual sufficiency, early-stage learning dy-	2.1 Visual-Language Model Architecture	176
126	namics, resolution sensitivity, and spatial robust-	Figure 1 illustrates our processing pipelines for a	177
127	ness.	concrete example—predicting the final character in	178
128	In this work, we systematically investigate	“数据显示” (data shows) using a model checkpoint	179
129	a vision-token-based formulation of Chinese	of the early-stage training.	180
130	language modeling, processing low-resolution		
131	grayscale character images through a lightweight		
132	visual encoder fed into a standard language decoder		
133	architecture. Our findings reveal that visual form		
134	alone brings strong predictive power, providing a		
135	structural foundation that accelerates early stage		

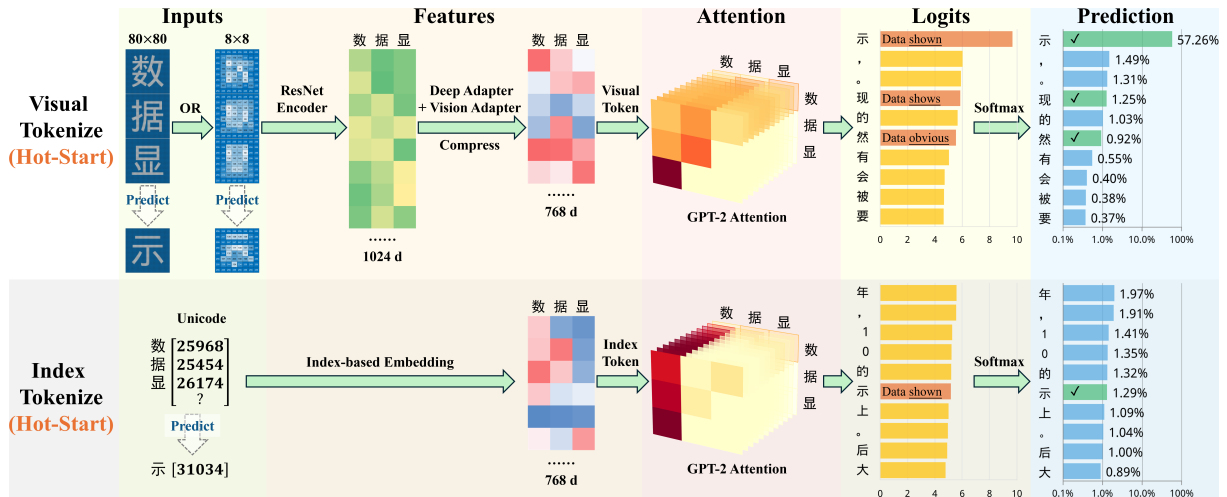


Figure 1: Model architecture and concrete processing example showing the prediction of the final character in the phrase “数据显示” (The data shown); numerical results demonstrate the *hot-start* phenomenon of visual tokens by 0.5% of the total training progress. Top: Visual-based training pipeline; numerical results are based on inputs of 8×8 character images. Bottom: Index-based training pipeline.

The diagram gives the two fundamentally different input paths while sharing the same language decoder—GPT-2-small-style (Radford et al., 2019) in our experiment, with $\sim 117M$ parameters. In the index-based path (bottom), each character is represented by a discrete token index. In the visual-based path (top), characters are first rendered as low-resolution grayscale images and passed through a lightweight visual encoder before being fed into the same decoder. Note that the figure visualizes the *hot-start* effect: after only 0.5% of training, the visual-based model already assigns much higher probability to the correct next character than its index-based counterpart and gives linguistically more plausible candidate rankings.

In particular, in the visual-based paths, visual inputs are passed through a ResNet encoder (He et al., 2016) and a Vision Adapter (Wu et al., 2019) before reaching the decoder embedding space, while index-based inputs are directly mapped into the decoder embeddings.

Remark of Reverse OCR concern. Our approach differs fundamentally from reversing an OCR (transcribing images back to symbols). Evidence against mere symbol reconstruction is: (1) the *hot-start* effect where visual models outperform index-based baselines early in training with minimal data; (2) the ability to discriminate subtle glyph differences; and (3) robustness under partial cropping. These patterns all suggest that the model learns structural understanding of character shapes rather than performing symbol mapping (Wu et al.,

2019; Geirhos et al., 2020).

Chinese characters differ fundamentally from alphabetic systems: each character is a minimal unit where visual patterns carry information through local details, compositional elements, and global shape. Capturing these visual regularities provides a structural prior that facilitates more efficient learning in Chinese language modeling.

To operationalize this approach, we render each character as a grayscale image. In the *Vision-100%* mode, characters occupy approximately 80% of the image width and height. For example, at 8×8 resolution, the character occupies about 6.4×6.4 active pixels, leaving 10% border margins on each side. We intentionally include these border margins to simulate realistic reading conditions—mimicking how characters appear within documents with natural spacing.

Cropping removes portions while keeping resolution fixed; for instance, in *Vision-50%*, the top 50% of original pixels are retained, with the remainder filled with background. This design is motivated by an intuitive observation: humans can often recognize partially visible Chinese characters. Our cropping experiments test whether models similarly extract predictive signals from limited visual regions. Figure 2 illustrates “人工智能” at 80×80 and 8×8 pixels, showing full characters, top 80%, and top 50% crops.

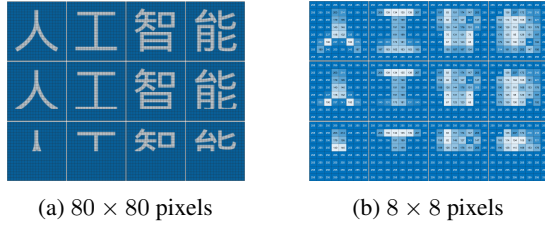


Figure 2: Heatmap visualization of character example cropping at two resolutions: low and high.

2.2 Training Objective

Models are trained to predict the next character conditioned on preceding inputs, minimizing the standard cross-entropy loss:

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \log P(c_{t+1}^{(i)} | I_1^{(i)}, \dots, I_t^{(i)}), \quad (1)$$

where N is the batch size, T is the fixed sequence length, $c_{t+1}^{(i)}$ is the ground-truth character at position $t + 1$ in sequence i , and $I_t^{(i)}$ denotes input representations (visual or index embeddings). We train models on the THUCNews dataset, a large-scale Chinese news corpus covering multiple domains, split into fixed-length sequences. Optimization employs AdamW with a fixed learning rate and batch size. For visual inputs in particular, gradients propagate through the projection and adapter into the shared decoder.

We evaluate three main input configurations: the *Index-based Model* baseline, which inputs token IDs; *Vision-100%* mode in which characters are represented as full images; and *Cropped-Vision* mode, where partial crops retain the top 80% (*Vision-80%*) or top 50% (*Vision-50%*) of the input images. Experiments across different image sizes, especially at low resolutions (8×8 and 4×4), assess whether minimal visual information suffices for accurate prediction, and partial cropping examines reliance on distributed visual features rather than OCR-like reconstruction.

3 Experiments and Results

3.1 Experimental Setup

Model Configurations. We evaluate three main input configurations: the *Index-based Model* baseline (token IDs), *Vision-100%* mode (full images), and *Cropped-Vision* mode with partial crops (top 80% and 50%).

Experiments span resolutions from 4×4 to 80×80 pixels, with particular focus on low resolutions

(8×8 and 4×4) to test visual sufficiency. Partial cropping examines reliance on distributed visual features rather than exact glyph reconstruction.

Dataset and Training. We train models on the THUCNews dataset, based on historical data collected from the RSS subscription channels of Sina News between 2005 and 2011 (Guo et al., 2016). After filtering and cleaning, it consists of approximately 740,000 news articles in UTF-8 encoded plain text format. The corpus contains 100K sequences, which is 12.8M character instances, split into fixed-length sequences of 128 characters. Sequences consist primarily of Chinese characters with occasional English letters and punctuation (approximately 10%), reflecting real news text composition. We employ a quadratic curriculum learning strategy, where the number of training sequences grows as $5000 + 918.37\text{epoch} + 18.74\text{epoch}^2$. Under this curriculum, the first 100K sequences are seen a total of approximately 2.13M training instances across epochs. The dataset increases progressively while evaluating on a fixed validation set of 5K sequences. This approach balances early fast convergence with later exposure to the full dataset.

Key Parameters. The language decoder follows a GPT-2-small-style architecture with 12 layers, 768 hidden dimensions, pre-trained on UER (Zhao et al., 2019). In particular, optimization uses AdamW with learning rate 2×10^{-4} (OneCycle scheduler, max 1.5×10^{-3}), batch size 128, weight decay 0.01, gradient clipping at 1.0, mixed precision (FP16), and early stopping (patience: 7 epochs).

Resolution Spectrum. Experiments span resolutions from 4×4 to 96×96 pixels, assessing whether minimal visual information suffices for accurate prediction. In particular, we include extreme low (4×4 , almost no human-recognizable cues), typical low (8×8 , essential structure retained), intermediate (20×20 to 40×40 , recognizable shapes), and high (80×80 to 96×96 , all details preserved) resolutions. This spectrum allows us to examine how much visual information is necessary for Chinese character modeling.

3.2 RQ1: Visual Sufficiency

RQ1 asks whether visual inputs alone suffice for character-level prediction. Note that Chinese contains over 5,500 distinct characters: randomly guessing yields only about 0.02% accuracy ($1/5,500$), while a unigram baseline—predicting the most frequent character according to dataset

statistics—achieves roughly 2%. In this context, our model’s performance of 39% indicates that it captures substantial linguistic structure beyond simple frequency statistics.

Table 1 (first row, *Vision-100%* mode) provides the answer: even 8×8 inputs achieve 39.21% accuracy, matching the index-based baseline (39.10%). This confirms that minimal visual information suffices for accurate prediction.

3.3 RQ3: Resolution Sensitivity

RQ3 investigates how performance varies with resolution. Table 1 shows results across the resolution spectrum: from 4×4 (29.70%) to 80×80 (39.03%), with 8×8 achieving 39.21%—comparable to the index-based baseline (39.10%). This confirms that minimal resolution suffices once essential structure is preserved.

3.4 RQ4: Spatial Robustness

RQ4 examines robustness to spatial cropping. Table 1 shows that severe cropping causes minimal performance drops: at 8×8 , *Vision-80%* (top 80%) achieves 39.18% and *Vision-50%* (top 50%) 38.63%.

Analysis of 8×8 input images reveals sparse pixel usage: *Vision-100%* uses 6×6 active pixels, *Vision-80%* 6×5 , and *Vision-50%* 6×3 . This explains the robustness—models extract essential structure even when peripheral regions are missing.

Figure 3 illustrates the “toast-center” effect: the central strokes (Toast-Center / Crumb) contain rich character information, while the outer layer (Crust) contributes less. Besides, blank space carries negligible information. This concentrated central structure might explain why models can maintain accuracy even under severe cropping.

While 80×80 images preserve all character details, their predictive advantage over 8×8 is statistically marginal. To quantify this, we computed adjusted standard errors accounting for sequence-level correlations (DEFF = 19.9, $\rho = 0.15$), yielding SE = 0.27–0.54% across accuracy levels. The small differences between high and low resolutions confirm that coarse structural cues—not fine-grained details—drive predictive performance, further supporting RQ1’s sufficiency claim.

3.5 RQ2: Early-Stage Dynamics

RQ2 examines how vision-token models behave during initial training compared to index-based models.

Mode	4 × 4 Acc/PPL	8 × 8 Acc/PPL	20 × 20 Acc/PPL	30 × 30 Acc/PPL	80 × 80 Acc/PPL
<i>Vision-100%</i>	29.70/85.33	39.21/46.59	39.16/45.83	39.14/48.73	39.03/49.41
<i>Vision-80%</i>	18.28/194.98	39.18/46.23	39.15/46.33	39.07/48.83	39.08/48.74
<i>Vision-50%</i>	2.10/2249.29	38.63/47.95	38.70/48.04	38.66/49.81	38.57/50.33
<i>Index-based</i>	39.10/47.58				

Table 1: Accuracy (%) / PPL across resolutions. *Vision-100%*: full images; *Vision-80%*: top 80% crop; *Vision-50%*: top 50% crop. *Index-based Model*: baseline using discrete character indices.

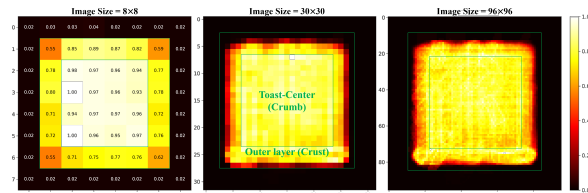


Figure 3: “toast-center effect”: center strokes (blue box) receive more attention than outer pixels (red box).

As shown in Table 2, even with 4,096 training sequences, visual models show substantial gains: with 40×40 pixel inputs, accuracy reaches 13.06%, tripling the index-based baseline’s 4.30%.

Sample Trained	Baseline	8 × 8 Vision	40 × 40 Vision
4,096	4.30%	4.19%~ (-0.11%)	13.06%~ (+8.76%)
6,152	4.61%	5.57%~ (+0.96%)	14.7%~ (+10.09%)
8,200	5.84%	12.34%~ (+6.5%)	15.46%~ (+9.62%)
10,248	8.45%	13.94%~ (+5.49%)	15.92%~ (+7.47%)

Table 2: Hot-start progression across training stages. Higher resolutions show earlier and stronger advantages that gradually narrow over time.

Remarkably, we observe a pronounced *hot-start* effect. At 0.4% of total training (8,200 sequences), 8×8 visual inputs achieve 12.34% accuracy, more than double the index-based baseline’s 5.84%. In addition, as illustrated in Figure 1, during this hot-start stage, visual models not only achieve higher next-character accuracy, but also assign higher probabilities to plausible top candidates, suggesting that the learned representations support linguistically reasonable predictions beyond the single argmax.

This early advantage persists with visual models maintaining a consistent lead. For instance, by 16,441 sequences, 8×8 visual models achieve 15.65% accuracy, while the baseline reaches only 13.33%. This sustained advantage underscores that visual structure provides not just an initial shortcut, but a persistent edge that the index-based baseline struggles to match.

It is also important to note that the onset timing of this advantage correlates directly with input resolution. Higher-resolution models reach their hot-start phase earlier: at only 0.2% total training time, 40×40 inputs already achieve 13.06% ac-

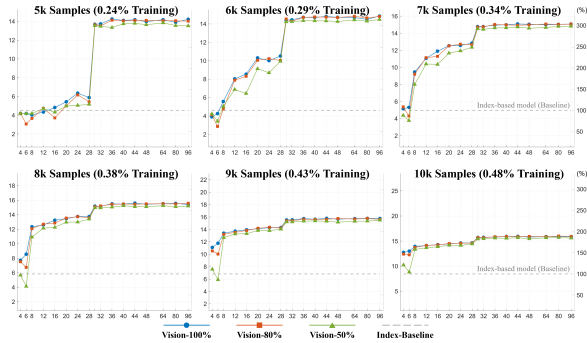


Figure 4: Validation accuracy across image resolutions at selected early training checkpoints (5k–10k samples), plotted on the index-based baseline scale (dashed line).

accuracy, while the baseline remains at 4.30%. This pattern suggests that richer visual detail accelerates the extraction of structural regularities, giving higher-resolution models a decisive head start in learning.

Recall that our training employs a quadratic curriculum: the dataset size grows quadratically across epochs, meaning models see progressively more data as training advances. As a result, these early sequence counts ($\sim 10^4$) represent an extremely data-scarce regime where models have access to only a tiny fraction of the eventual corpus—making the observed hot-start advantage particularly meaningful.

Figure 4 further visualizes that, across early training checkpoints (5k–10k samples), visual inputs maintain a consistent performance advantage, supporting the observed hot-start effect.

As noted in Section 3.4, the central strokes (toast-center) contain concentrated structural information, which helps explain the hot-start effect.

3.6 Ablation Studies

To understand which components drive performance, we conducted three ablation studies using 8×8 resolution as our testbed. In each case, training data and hyperparameters remain identical.

Dual-Encoder Ablation. Following DeepSeek-OCR’s architecture (Wei et al., 2025), we test a dual-encoder that concatenates features from both ResNet and Vision Transformer (ViT). Comparisons with show that ResNet alone provides most of the structural benefit. We hypothesize that ViT is better suited for inter-character spatial relationships, adding marginal gains primarily in high-resolution settings.

Training Strategy Ablation. We compare joint training (visual components + decoder) versus freezing the decoder during adapter training. Joint

Category	Setting	Acc (%)	Δ (pp)	Sig.
<i>Main Results (8×8)</i>				
Baseline	<i>Index-based</i>	39.10	—	Ref.
Vision	<i>Vision-100%</i>	39.21	+0.11	$p=0.77$
Vision	<i>Vision-80% (80%)</i>	39.18	+0.08	$p=0.81$
Vision	<i>Vision-50% (50%)</i>	38.63	-0.47	$p=0.09$
<i>Ablations (vs Vision-100%)</i>				
Encoder	ViT	38.45	-0.76	$p < 0.001$
Training	Frozen decoder	36.78	-2.43	$p < 0.001$
Architecture	No adapter	37.12	-2.09	$p < 0.001$
<i>Hot-Start (10K samples)</i>				
Early	<i>Index-based</i>	6.45	—	Ref.
Early	<i>Vision-100%</i>	14.65	+8.20	$p < 0.001$
Early	<i>Vision-80%</i>	14.70	+8.25	$p < 0.001$
Early	<i>Vision-50%</i>	14.38	+7.93	$p < 0.001$
<i>Statistical Validation</i>				
Design	DEFF=19.9	—	—	$\rho=0.15$
Samples	$n_{\text{eff}}=31.9\text{K}$	—	—	$4.46 \times \text{SE}$
CI	Width $\pm 0.537\text{pp}$	—	—	$\Delta=0.11\text{pp}$

Table 3: Consolidated results. All tests use adjusted SE (DEFF=19.9). Main: *Visual-based Models* match text baseline. Ablations: CNN>ViT ($\Delta=-0.76\text{pp}$), joint training vital ($\Delta=-2.43\text{pp}$). Hot-start: $2.27 \times$ faster early learning. Stats: sequence correlations accounted for.

training gives significantly better results, indicating that end-to-end optimization is crucial.

Architecture Ablation. We remove the visual adapter entirely (direct projection to decoder) causes performance degradation, particularly in the hot-start phase.

Results summarized in Table 3 show that ResNet with joint training achieves the best performance. Notably, hot-start advantages persist across all ablations, supporting that structural learning—not just specific architectural choices—is the key driver.

3.7 Summary of Visual Advantages

Across resolutions and ablations, three insights emerge. Low-resolution inputs alone capture essential character structure. Visual cues accelerate early stage training, producing a *hot-start* effect. Finally, models remain robust under severe degradation, demonstrating that coarse structural cues suffice. These results establish that visual representations provide a robust and sample-efficient alternative to index-based inputs for Chinese language modeling.

4 Interpretability Analysis

All analyses in this section use the 8×8 resolution setting at the hot-start phase (with 10k training sequences). This setting isolates essential structural information while removing fine-grained visual details, allowing us to examine how minimal visual cues enable early linguistic discrimination. We emphasize that these analyses are post-hoc and

ID	Sentence → Candidates (✓/×)	Model	P(%)	Choice
1	下雨天鞋子上很容易沾上泥 → 土/土	Vision	0.05/0.00	土✓
		Text	0.00/0.01	土×
2	他是一个边境战 → 土/土	Vision	0.01/0.00	土✓
		Text	0.00/0.01	土×
3	这地板的材料是实 → 木/本	Vision	0.00/0.23	本×
		Text	0.00/0.31	本×
4	别忘了拿作业 → 本/木	Vision	0.04/0.00	本✓
		Text	0.04/0.00	本✓
5	昨天周六，今天是星期 → 日/目	Vision	0.19/0.12	日✓
		Text	0.00/0.01	目×
6	这个广告牌很醒 → 目/日	Vision	0.09/0.08	目✓
		Text	0.03/0.00	目✓
7	这个房间非请莫 → 入/人	Vision	0.04/0.02	入✓
		Text	0.06/0.10	人×
8	介绍一下，这位是我的爱 → 人/入	Vision	8.63/0.00	人✓
		Text	0.06/0.25	入×

Table 5: Predicted probabilities for visually similar character pairs at sentence end. Sentences are selected to provide minimal but coherent semantic context. Four pairs of cases show model disagreement, highlighting differences between vision-based and index-based predictions.

exhibit enhanced discriminative power at sentence-end prediction, strategically organizing confusable characters in embedding space to support early-stage prediction advantages.

4.3 Pixel-Level Importance Analysis via Gradient Back-propagation

To understand which visual features matter most for prediction, we perform gradient-based analysis (Simonyan et al., 2014; Aflalo et al., 2022). For character images $\{I_1, \dots, I_n\}$ and target c_{n+1} , we compute character-level importance $S_k = \sum_{i,j} \left| \frac{\partial y_{c_{n+1}}}{\partial I_{k,i,j}} \right|$, where higher S_k indicates that input character I_k is more important for predicting the target. Pixel-level importance is normalized as $H_{k,i,j} = \left(\left| \frac{\partial y_{c_{n+1}}}{\partial I_{k,i,j}} \right| - \min \right) / (\max - \min)$.

Testing on the confusable pairs from Section 4.2 reveals that the model assigns high importance to semantically relevant input characters (e.g., “泥” when predicting 土) and near-zero to irrelevant ones (“的”). Table 6 shows similar attention intensities across character regions, suggesting that the model distributes attention broadly rather than focusing exclusively on any single area. This balanced pattern explains why our cropping experiments (*Vision-80%* and *Vision-50%*) remain effective: the model can extract predictive signals from various character subregions.

The observed gradient patterns offer preliminary support for our “toast-center” conjecture though verifying its role in the *hot-start* phase requires further temporal analysis in future work.

Together, the analyses above illustrate how our vision-token representations make model behavior more inspectable. They link language predictions

Region	Avg. Intensity	Std. Dev.
Upper Half	0.087	0.014
Lower Half	0.081	0.014
Left Half	0.085	0.016
Right Half	0.083	0.012

Table 6: Average attention intensity across character regions.

to explicit visual structure in ways that are difficult to access in index-based models.

5 Conclusion

In this work, we challenge the dominant index-based paradigm for Chinese language modeling by asking: *Can language modeling rely solely on visual form?* Our results answer affirmatively, especially for early-stage or low-resource scenarios.

Our investigation yields affirmative answers to our four research questions. **RQ1 (Visual Sufficiency)** is confirmed: visual inputs achieve accuracy comparable to index-based baselines (39.2% vs. 39.1%). **RQ2 (Early-Stage Dynamics)** reveals a pronounced *hot-start* effect: visual models reach 12.3% accuracy within 0.4% of training time—more than double the baseline’s 5.8%. **RQ3 (Resolution Sensitivity)** demonstrates that even 8×8 pixels retain sufficient structure, while **RQ4 (Spatial Robustness)** confirms effectiveness under severe cropping (top 50% retained).

This advantage stems from fundamental differences in representation topology. Unlike index-based embeddings—which begin as unstructured points in feature space—visual embeddings inherit spatial organization from the encoder, providing a structural regularity that accelerates early learning.

Beyond performance gains, visual representations naturally provide some interpretability: characters sharing radicals tend to cluster while visually confusable pairs are often separated. Further, gradient-based analysis gives how the pixel regions contribute to predictions. This structure emerges from the visual input itself, without requiring auxiliary objectives or post-hoc tools.

Overall, visual structure provides not merely an alternative input format, but a sample-efficient inductive bias for Chinese language modeling. This points toward architectures natively designed for visual glyph processing and learning strategies that leverage visual structure for efficient training.

619 Limitations

620 Our study has several limitations that point to future
621 work. First, our experiments use standard font ren-
622 dering. Testing handwritten or stylized characters
623 would assess robustness. The GPT-2–style decoder
624 is relatively small; scaling to larger architectures
625 may give different visual utilization patterns. Ad-
626 ditionally, our single-character image processing
627 differs from full-paragraph OCR approaches. Ex-
628 tending to multi-character or paragraph-level im-
629 ages remains a valuable direction. Finally, applying
630 this approach to other logographic systems could
631 further prove its generality.

632 In summary, while current limitations exist, our
633 findings demonstrate decisively that low-resolution
634 visual glyphs are not just a viable substitute, but
635 a cognitively richer starting point for Chinese lan-
636 guage models. In contrast to recent surveys high-
637 lighting modality collapse in vision–language mod-
638 els (Sim et al., 2025), our results show that when
639 visual input is isolated and structurally constrained,
640 models can reliably exploit visual form rather than
641 bypass it.

642 References

- 643 Estelle Aflalo, Meng Du, Shao-Yen Tseng, Yongfei
644 Liu, Chenfei Wu, Nan Duan, and Vasudev Lal. 2022.
645 [VI-interpret: An interactive visualization tool for in-
646 terpreting vision-language transformers](#). In *2022
647 IEEE/CVF Conference on Computer Vision and Pat-
648 tern Recognition (CVPR)*, pages 21374–21383.
- 649 Emily M. Bender and Alexander Koller. 2020. [Climbing
650 towards NLU: On meaning, form, and understanding
651 in the age of data](#). In *Proceedings of the 58th Annual
652 Meeting of the Association for Computational Lin-
653 guistics*, pages 5185–5198, Online. Association for
654 Computational Linguistics.
- 655 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ
656 Altman, Simran Arora, Sydney von Arx, Michael S.
657 Bernstein, Jeannette Bohg, Antoine Bosselut, Emma
658 Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card,
659 Rodrigo Castellon, Niladri S. Chatterji, Annie S.
660 Chen, Kathleen A. Creel, Jared Davis, Dora Dem-
661 szky, and 95 others. 2021. [On the opportunities and
662 risks of foundation models](#). *ArXiv*, abs/2108.07258.
- 663 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
664 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
665 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
666 Askeel, Sandhini Agarwal, Ariel Herbert-Voss,
667 Gretchen Krueger, Tom Henighan, Rewon Child,
668 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
669 Winter, and 12 others. 2020. [Language models are
670 few-shot learners](#). In *Advances in Neural Information*

Processing Systems, volume 33, pages 1877–1901.
Curran Associates, Inc.

- 671
672
- 673 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio
674 Michaelis, Richard Zemel, Wieland Brendel,
675 Matthias Bethge, and Felix A. Wichmann. 2020.
676 [Shortcut learning in deep neural networks](#). *Nat.
677 Mach. Intell.*, 2:665–673.
- 678 Zhipeng Guo, Yu Zhao, Yabin Zheng, Xiance Si,
679 Zhiyuan Liu, and Maosong Sun. 2016. [THUCTC:
680 An efficient chinese text classifier](#). GitHub reposi-
681 tory.
- 682 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian
683 Sun. 2016. [Deep residual learning for image recogni-
684 tion](#). In *2016 IEEE Conference on Computer Vision
685 and Pattern Recognition (CVPR)*, pages 770–778.
- 686 Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu,
687 Fangyu Liu, Julian Eisenschlos, Urvashi Khandel-
688 wal, Peter Shaw, Ming-Wei Chang, and Kristina
689 Toutanova. 2023. [Pix2struct: screenshot parsing as
690 pretraining for visual language understanding](#). In
691 *Proceedings of the 40th International Conference on
692 Machine Learning, ICML’23*. JMLR.org.
- 693 Zhang Li, Yuliang Liu, Qiang Liu, Zhiyin Ma, Ziyang
694 Zhang, Shuo Zhang, Zidun Guo, Jiarui Zhang, Xinyu
695 Wang, and Xiang Bai. 2025. [Monkeyocr: Document
696 parsing with a structure-recognition-relation triplet
697 paradigm](#). *Preprint*, arXiv:2506.05218.
- 698 Jake Poznanski, Aman Rangapur, Jon Borchardt, Ja-
699 son Dunkelberger, Regan Huff, Daniel Lin, Aman
700 Rangapur, Christopher Wilhelm, Kyle Lo, and Luca
701 Soldaini. 2025. [olmocr: Unlocking trillions of to-
702 kens in pdfs with vision language models](#). *Preprint*,
703 arXiv:2502.18443.
- 704 Alec Radford, Jeff Wu, Rewon Child, David Luan,
705 Dario Amodei, and Ilya Sutskever. 2019. [Language
706 models are unsupervised multitask learners](#).
- 707 Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder,
708 and Iryna Gurevych. 2021. [How good is your tok-
709 enizer? on the monolingual performance of multilin-
710 gual language models](#). In *Proceedings of the 59th
711 Annual Meeting of the Association for Computational
712 Linguistics and the 11th International Joint Confer-
713 ence on Natural Language Processing (Volume 1:
714 Long Papers)*, pages 3118–3135, Online. Association
715 for Computational Linguistics.
- 716 Mong Yuan Sim, Wei Emma Zhang, Xiang Dai, and
717 Biao Yan Fang. 2025. [Can VLMs actually see and
718 read? a survey on modality collapse in vision-
719 language models](#). In *Findings of the Association
720 for Computational Linguistics: ACL 2025*, pages
721 24452–24470, Vienna, Austria. Association for Com-
722 putational Linguistics.
- 723 Karen Simonyan, Andrea Vedaldi, and Andrew Zis-
724 serman. 2014. [Deep inside convolutional networks:
725 Visualising image classification models and saliency
726 maps](#). In *Workshop at International Conference on
727 Learning Representations*.

728 Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang,
729 Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao,
730 Jianjian Sun, Yuang Peng, Chunrui Han, and Xi-
731 angyu Zhang. 2024. [General ocr theory: Towards](#)
732 [ocr-2.0 via a unified end-to-end model](#). *Preprint*,
733 arXiv:2409.01704.

734 Haoran Wei, Yaofeng Sun, and Yukun Li. 2025.
735 [Deepseek-ocr: Contexts optical compression](#).
736 *Preprint*, arXiv:2510.18234.

737 Wei Wu, Yuxian Meng, Fei Wang, Qinghong Han,
738 Muyu Li, Xiaoya Li, Jie Mei, Ping Nie, Xiaofei Sun,
739 and Jiwei Li. 2019. [Glyce: Glyph-vectors for chinese](#)
740 [character representations](#). *ArXiv*, abs/1901.10125.

741 Zhe Zhao, Hui Chen, Jinbin Zhang, Xin Zhao, Tao Liu,
742 Wei Lu, Xi Chen, Haotang Deng, Qi Ju, and Xiaoy-
743 ong Du. 2019. [UER: An open-source toolkit for pre-](#)
744 [training models](#). In *Proceedings of the 2019 Confer-*
745 *ence on Empirical Methods in Natural Language Pro-*
746 *cessing and the 9th International Joint Conference*
747 *on Natural Language Processing (EMNLP-IJCNLP):*
748 *System Demonstrations*, pages 241–246, Hong Kong,
749 China. Association for Computational Linguistics.