

FBS: Modeling Native Parallel Reading inside a Transformer

Tongxi Wang

Southeast University / Nanjing, China
tongxi_wang@seu.edu.cn

Abstract

Large language models (LLMs) excel across many tasks, yet inference is still dominated by strictly token-by-token autoregression. Existing acceleration methods largely patch this pipeline and miss core human-reading ingredients: content-adaptive foresight, chunk-structure-aware compute allocation, and train-test consistency for preview/skimming. We propose the **Fovea-Block-Skip Transformer (FBS)**, which injects a causal, trainable loop into Transformers via Parafovea-Attention Window (PAW), Chunk-Head (CH), and Skip-Gate (SG). Across diverse benchmarks, FBS improves the quality-efficiency trade-off without increasing parameters, and ablations show the three modules are complementary.

1 Introduction

Large language models (LLMs) have consistently failed in dialogue, code generation, and cross-domain reasoning (Minaee et al., 2024; Touvron et al., 2023; Guo et al., 2025). Yet at inference time, most LLMs still advance through a strictly token-by-token autoregressive process: computation is spent almost uniformly, step after step, along a single causal chain. This stands in sharp contrast to proficient human reading, especially by natives, which is often **highly parallel** and **multi-granular**: readers leverage parafoveal preview (Xu et al., 2022; Touvron et al., 2023; Snell et al., 2017), chunk-level processing (Yang et al., 2020; Bazant-Kimmel, 2018), and semantic-load-adaptive skimming (Fitzsimmons et al., 2014; Duggan and Payne, 2009) to improve efficiency while maintaining comprehension quality.

A large body of “faster decoding” work improves efficiency by patching the autoregressive pipeline (Chen et al., 2023), e.g., draft-and-verify speculative decoding (Leviathan et al., 2023), multi-candidate generation (Medusa/EAGLE) (Cai et al., 2024; Li et al., 2024d,c), set/block decoding (Gat

et al., 2025), and adaptive computation such as early-exit or dynamic depth (Xin et al., 2020; Teerapittayanon et al., 2016). These methods are effective, but they typically miss three structural ingredients that make human reading both fast *and* reliable: **(i) content-adaptive, structured foresight** (look-ahead is often fixed or external), **(ii) coupling between chunk structure and compute allocation** (chunks are used for representation compression rather than for deciding where to allocate compute), and **(iii) a matched train-test pathway** (previewing/skimming is often introduced only at inference, risking distribution shift).

We argue that breaking the autoregressive bottleneck is less about stacking isolated tricks and more about building an *endogenous* mechanism that preserves causality while realizing a native-speaker-style pipeline: **preview** → **chunking** → **skimming**. Concretely, an efficient generator should (1) form a low-resolution *predictive preview* to guide upcoming decisions, (2) integrate local semantics at the level of *chunks* rather than only tokens, and (3) regulate compute *dynamically*—skimming stable regions and pausing to compute when semantic load or uncertainty rises.

Guided by this viewpoint, we propose the **Fovea-Block-Skip Transformer (FBS)**, a causal Transformer augmented with a trainable, verifiable reading pipeline. FBS consists of three cooperative modules: **Parafovea-Attention Window (PAW)** produces *content-adaptive predictive preview* from the model’s own next-token distributions; **Chunk-Head (CH)** builds an online chunk-level semantic channel for phrase-scale integration; and **Skip-Gate (SG)** turns these signals into *true layer/block skipping* during decoding, allocating computation where it matters. The three modules form a closed loop: preview informs chunking, chunk semantics stabilize decisions, and the resulting confidence/semantic load controls skimming.

We instantiate FBS by continual pre-training

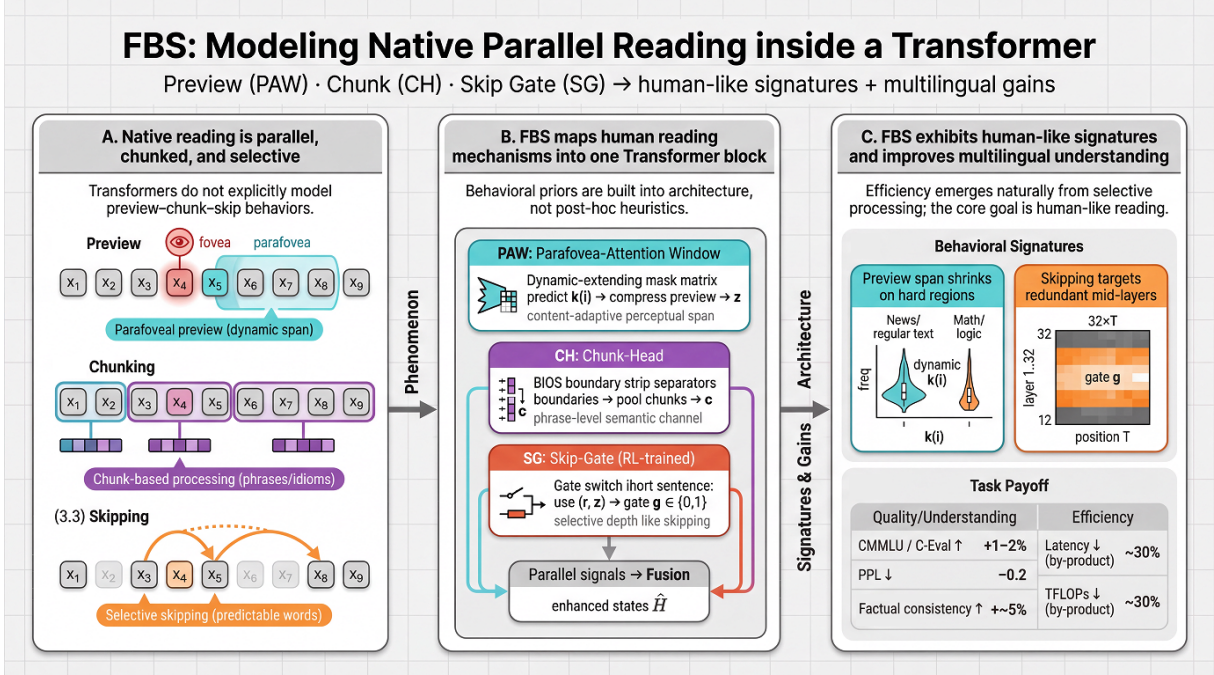


Figure 1: FBS block overview.

Qwen3-4B-Instruct (Yang et al., 2025) on a Chinese-English mixed corpus and evaluate it on a broad set of reasoning, knowledge, math, and code benchmarks. Under a matched-parameter setting, FBS improves quality on major benchmarks while reducing executed compute; in a 512→128 generation protocol, it achieves substantial wall-clock latency reduction (about 30%) and significantly lowers TFLOPs. Analyses further show that PAW, CH, and SG are complementary: PAW provides foresight, CH strengthens local semantic organization, and SG contributes the dominant efficiency gains, yielding human-like “speed up on stable spans, slow down at critical points” behavior.

Our main contributions are: (i) We formalize a *native-speaker parallel reading* abstraction for LLM inference and translate it into trainable mechanisms suitable for causal Transformers. (ii) We propose **FBS**, a unified architecture that couples PAW and CH with an internal compute controller SG, yielding an endogenous **preview**→**chunk**→**skim** loop that preserves causality while enabling real block/layer skipping. (iii) We demonstrate consistent **quality**–**efficiency** gains on Chinese and English benchmarks under a matched-parameter setting and provide analyses that explain when and why the model skims or computes more.

2 Methodology

We present the **Fovea-Block-Skip Transformer (FBS)**, which augments each causal Transformer

layer with three lightweight modules: **PAW**, prefix-only predictive preview with a content-adaptive span; **CH**, a parallel chunk semantic channel interacting with token states; and **SG**, a binary controller that can bypass the entire block at inference for true block/layer skipping. Detailed training surrogates and implementation choices are provided in the appendix (Appendix C, D, E).

2.1 Model Overview

Let $\mathbf{x}_{1:m}$ be an input sequence and $\mathbf{h}_{1:m}^{(\ell)} \in \mathbb{R}^{m \times d}$ the layer- ℓ hidden states. An FBS layer retains standard causal self-attention, and adds: (i) a PAW preview vector $\mathbf{z}_i^{(\ell)}$ to enrich token i with a predicted future summary, (ii) a CH chunk cache that provides a parallel chunk semantic context, and (iii) an SG gate $g_t^{(\ell)} \in \{0, 1\}$ during decoding, where $g_t^{(\ell)} = 1$ means **skip** and $g_t^{(\ell)} = 0$ means **compute**. When skipping, we forward an identity mapping for the current token: $\mathbf{h}_t^{(\ell+1)} = \mathbf{h}_t^{(\ell)}$. A full algebraic instantiation (including the training-time soft mixture form) is provided in Appendix E.1.

2.2 Dynamic Lookahead: Parafovea-Attention Window (PAW)

Goal. PAW provides a *verifiable* preview signal that is compatible with causality: the preview is computed from the model’s *own* predictive distributions, not from future ground-truth tokens. Concretely, given the current-layer state at position t , PAW predicts a discrete lookahead span

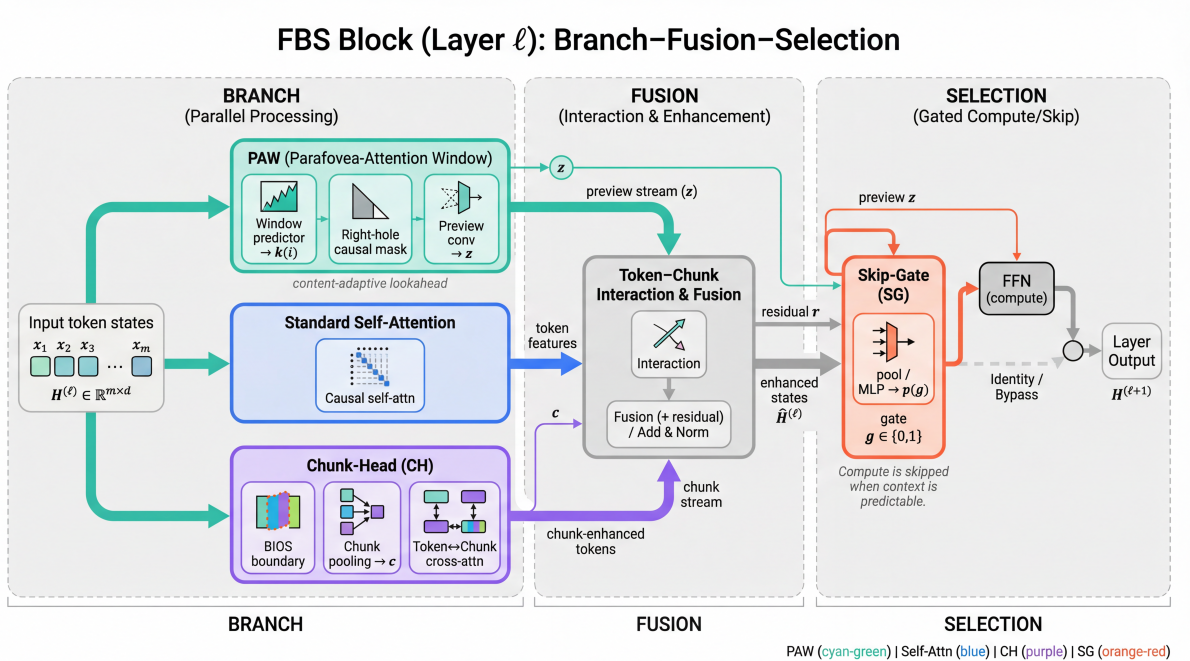


Figure 2: **FBS Pipeline**

$k(t) \in \{0, \dots, k_{\max}\}$ and summarizes the next $k(t)$ *predicted* tokens into a preview vector \mathbf{z}_t . The predictor, the multi-horizon preview head, and the (soft) preview compression are defined in C.1.

Incremental computation at decoding time (KV-cache compatible). During autoregressive decoding, only the newest position t is processed at each step. PAW is computed **only for this newest position** and never revisits the prefix: (i) from $\mathbf{h}_t^{(\ell)}$, a tiny predictor outputs $k(t)$; (ii) for horizons $r = 1, \dots, k(t)$, PAW produces a predictive distribution over the r -th next token using a lightweight head, maps each distribution to a preview embedding, and then compresses these $k(t)$ embeddings into \mathbf{z}_t ; (iii) \mathbf{z}_t is injected as an additive channel into the current token state. All intermediate preview tensors for earlier positions can be discarded once \mathbf{z}_t is produced. Hence PAW adds a small per-step overhead that scales with k_{\max} (not with the context length) and is fully compatible with standard KV caching.

2.3 Chunk-Level Parallel Head (CH)

CH introduces a chunk-level semantic stream that complements token-wise attention. At each step, CH predicts whether the newest token closes a chunk (BIOS-style boundaries), pools token states within the current chunk, and appends the resulting chunk state into a **chunk cache**. The newest token can then attend to (or fuse with) the chunk

cache, enabling phrase-level integration without re-processing the prefix. The weak-supervision pipeline, chunk construction rules, and the full on-line update protocol are detailed in Appendix D.

2.4 Layer Skipping: Skip-Gate (SG)

SG is a discrete controller that decides whether to execute the current layer computation. At decoding step t and layer ℓ , SG consumes a **residual/uncertainty signal** $\mathbf{r}_t^{(\ell)}$ together with the PAW preview \mathbf{z}_t , and outputs a skip probability

$$p_t^{(\ell)} = \sigma(\text{MLP}([\mathbf{r}_t^{(\ell)}; \mathbf{z}_t])).$$

At inference we use **deterministic thresholding** $g_t^{(\ell)} = \mathbb{I}[p_t^{(\ell)} > \tau]$ to obtain true conditional execution (Appendix E.6). At deployment, $g=1$ short-circuits the *entire* layer block (no attention/FFN for the current token), and KV-cache semantics are detailed in Appendix E.1. Training uses a straight-through estimator and a soft mixture surrogate (Appendix E.2), so the learned policy transfers to hard skipping without train–test mismatch.

2.5 Optional RL fine-tuning for SG

Because SG makes discrete compute decisions, we optionally fine-tune SG with PPO while **freezing the backbone**. The reward trades off output quality and executed compute (Appendix E.7), encouraging “skip when confident” behavior without altering the autoregressive factorization.

Summary. PAW provides content-adaptive predictive preview; CH provides chunk-level structure for integration; SG turns these into real compute skipping. Together, they form a trainable analogue of a preview–chunk–skim reading pipeline.

3 Experiments and Results

This section addresses three questions:

(1) Can Fovea-Block-Skip Transformer (FBS) improve multilingual multi-task performance **and** inference efficiency under (approximately) unchanged parameter budgets?

(2) What are the individual contributions of the three key modules—Parafovea-Attention Window (PAW), Chunk-Head (CH), and Skip-Gate (SG)?

(3) Do the learned behaviors (e.g., lookahead windows and layer-skipping patterns) exhibit signatures consistent with parallel reading mechanisms?

3.1 Experimental Setup

3.1.1 Models and Compared Systems

We instantiate FBS on a causal language model comparable to **Qwen3-4B-Instruct**, denoted as **FBS-Base**. The backbone configuration is: hidden size 4096, 32 Transformer layers, 32 attention heads, and a vocabulary size of $\sim 32\text{K}$. FBS replaces each standard Transformer block with an FBS block consisting of three parallel sub-modules:

- **Parafovea-Attention Window (PAW):** dynamic lookahead window + preview convolution;
- **Chunk-Head (CH):** BIOS boundary prediction, chunk pooling, and cross-attention;
- **Skip-Gate (SG):** layer-wise skipping decisions based on residual signals and preview vectors.

Fairness and initialization. To ensure fair comparison, we load weights from the public baseline checkpoint: the original Multi-Head Attention and FFN weights are copied into the corresponding pathways of FBS; newly introduced PAW/CH/SG parameters are randomly initialized. Unless otherwise stated, all compared **target** models are kept strictly parameter-matched in scale.

Systems. We compare: (i) **Baseline:** the original causal LM without any FBS modules (Qwen3-4B-Instruct, same continual pretraining); (ii) **FBS-S1:** FBS after Stage-1 continual pretraining

(PAW+CH enabled; no RL-based skipping); (iii) **FBS-Full:** full FBS after Stage-1 and Stage-2 (PAW+CH+SG+RL enabled). We keep the target model size fixed across compared systems and report parameter counts in Appendix B. It provides the exact parameter counts used in all comparisons.

Baseline grouping by acceleration route (embedded into the experimental structure). To make the comparisons more persuasive, we categorize baselines into four groups and use them in different subsections:

- **Group A (Decoding acceleration; same backbone; inference-only changes).** Used in the main table (§3.2) as mainstream LLM inference acceleration baselines: Speculative Decoding/Sampling (Leviathan et al., 2023; Chen et al., 2023), EAGLE-2 (Li et al., 2024c), Medusa (Cai et al., 2024), and Lookahead Decoding (Zhao et al., 2024). We follow a SpecBench-style **measurement protocol B** for unified timing/ratio reporting; this affects the evaluation protocol but does not change model definitions.
- **Group B (Adaptive compute / layer skipping / FFN skipping; closest to SG).** Used to justify SG effectiveness: FlexiDepth (Luo et al., 2025) and FFN-SkipLLM (Jaiswal et al., 2024) (a compact main-text table), while Self-speculative Draft&Verify and Kangaroo are moved to Appendix E for a full comparison (because their verification pipelines and “loss-less” claims require more careful reporting).
- **Group C (Long-context / cache efficiency).** Only compared when we report long-context results: H2O (Zhang et al., 2023), StreamingLLM (Xiao et al., 2023), SnapKV (Li et al., 2024b), SlimInfer (Long et al., 2025) (Appendix F).
- **Group D (Chunk / structural modeling).** The main text uses our controlled CH ablations as the fairest comparisons; Segatron (Bai et al., 2021) is placed in Appendix D as a citation-style structural baseline to avoid cluttering the main narrative.

3.1.2 Training Data and Tasks

Stage-1 continual pretraining corpus. We use only open and commercially usable corpora:

Model	Params (B)	PAW	CH	SG	RL	Notes
Qwen3-4B-Instruct (Baseline)	4.0	✗	✗	✗	✗	Original instruction-tuned model
Qwen3-4B + FBS-S1	4.0	✓	✓	✗	✗	Structural changes only (PAW+CH)
Qwen3-4B + FBS-Full	4.0	✓	✓	✓	✓	Full FBS (Skip-Gate + RL)

Table 1: Overview of the target systems used in main results. All Group A/B/C/D baselines align to the same Qwen3-4B **target** model. Group-A speculative methods may additionally require a smaller **draft** model; draft parameters are not counted in the target size but are included in wall-clock timing (Appendix B.3).

RedPajama-V2 (English)(Weber et al., 2024), Yuan-2.0 Corpus (Chinese–English mixed)(Wu et al., 2023), and OSCAR-zh (Chinese)(Jansen et al., 2022). We sample a total of **30B tokens**. We filter examples with length ≥ 512 and deduplicate samples with similarity > 0.7 . For Chinese, we automatically generate weak BIOS labels using a pkuseg-based segmenter and an idiom lexicon, serving as CH supervision without any human annotation(Luo et al., 2019). For English, we do not use external chunk labels; CH boundaries are learned from the token stream via the same boundary predictor trained with the LM objective and regularizers. Empirically, CH yields larger gains on Chinese while remaining non-degrading on English tasks.

Stage-2 RL environment. We randomly sample 5k questions from the MMLU dev set and 5k questions from the CMMLU dev set, forming **10k prompts** for RL. Make sure there are no duplicates. Each episode generates at most **128 tokens**. Rewards are computed based on quality–compute differences relative to a full-compute reference. We use dev splits *only* for PPO prompt sampling; all reported benchmark numbers are computed on the official held-out evaluation splits (see Appendix B.2 for split hygiene and overlap control).

3.1.3 Training Details

Stage-1 continual pretraining. We train for **12B tokens** with batch size $\approx 2\text{M}$ tokens, learning rate 2×10^{-5} , 4% warmup, and cosine decay.

We optimize the following multi-task objective:

$$\mathcal{L} = \mathcal{L}_{\text{lm}} + \lambda_{\text{bios}} \mathcal{L}_{\text{bios}} + \lambda_{\text{align}} \mathcal{L}_{\text{align}} + \lambda_{\text{paw}} \mathcal{L}_{\text{preview}} + \lambda_{\text{gate}} \mathcal{L}_{\text{gate}}. \quad (1)$$

Here, \mathcal{L}_{lm} is the standard next-token cross-entropy; $\mathcal{L}_{\text{bios}}$ is token-level cross-entropy on BIOS labels predicted by CH (weak supervision from segmentation + idiom lexicon); $\mathcal{L}_{\text{align}}$ enforces semantic consistency between the chunk and token channels (e.g., contrastive or cosine alignment between

pooled chunk representations and token representations; Appendix D.2); $\mathcal{L}_{\text{preview}}$ supervises the multi-step preview predictor with ground-truth continuation labels, while the forward path still consumes the **predicted** preview embeddings to avoid train–test mismatch. We detail how the predicted distributions are mapped to continuous preview embeddings and report the PAW head sizes in Appendix C.1. $\mathcal{L}_{\text{gate}}$ is a conservative regularizer on SG outputs to avoid early-stage instability (e.g., entropy/mean constraints on skip probabilities).

Stage-1 gate training schedule. We linearly anneal the Skip-Gate threshold from 0.9 to 0.7. Stage-1 does **not** actually execute skipping (we only train gate parameters) to avoid early instability.

Stage-2 RL fine-tuning (PPO). We freeze all parameters except SG and the PAW network. We train on 10k RL prompts for **2 epochs** using PPO with clip $\epsilon = 0.2$ and learning rate 1×10^{-5} .

FLOPs and perplexity estimation. We estimate relative TFLOPs by counting (effective layer invocations) \times (per-layer FLOPs), normalized so that the baseline equals 1.0. Perplexity is approximated by the average log-likelihood on the evaluation sets.

Reward function. We train the skip policy with PPO using a scalar reward that balances compute savings against answer fidelity. Let $c(\mathbf{x})$ denote the estimated compute cost of decoding an output \mathbf{x} (our TFLOPs proxy), and let $c_0(\mathbf{x})$ be the cost under the full-compute baseline policy (no skipping). Because Stage-2 is instantiated on multiple-choice QA prompts sampled from MMLU and CMMLU, we can compute a per-instance, teacher-forced negative log-likelihood (NLL) on the **gold** answer option. We define the quality-degradation proxy as

$$\Delta\ell(\mathbf{x}) \triangleq \frac{\text{NLL}_{FBS}(\mathbf{y}^* | \mathbf{p}) - \text{NLL}_{Full}(\mathbf{y}^* | \mathbf{p})}{|\mathbf{y}^*|}, \quad (2)$$

where \mathbf{p} is the prompt and \mathbf{y}^* is the gold option text (concatenated with the standard answer prefix

used in evaluation). Our final reward is

$$R(\mathbf{x}) = \alpha \cdot \frac{c_0(\mathbf{x}) - c(\mathbf{x})}{c_0(\mathbf{x})} - \beta \cdot \max(0, \Delta\ell(\mathbf{x})), \quad (3)$$

which encourages skipping only when it does not increase the gold-option NLL. We set $\alpha = \beta = 0.1$ in all experiments.

3.1.4 Benchmarks and Metrics

We report: **PPL** (perplexity), **acc@5-shot** on MMLU(Hendrycks et al., 2021) / CMMLU(Li et al., 2024a) / C-Eval(Huang et al., 2023) / BBH(Suzgun et al., 2023), **pass@1** on HumanEval-X(Zheng et al., 2023) / MBPP(Austin et al., 2021), and **solve-rate** on GSM8K(Cobbe et al., 2021) / CMath(Wei et al., 2023).

Efficiency metrics and unified ratio definitions.

We measure average latency and relative TFLOPs (baseline normalized to 1.0). We unify method-specific ratios into **Bypass/Skip-Ratio**: for FBS, it is the average layer-skip ratio; for Speculative Decoding(Leviathan et al., 2023) and EAGLE-2(Li et al., 2024c), it is the acceptance rate; for Medusa(Cai et al., 2024) and Lookahead(Zhao et al., 2024), it is the verified/advance ratio (or normalized verified length; Appendix B.3).

Sampling settings must be fixed across methods; we use greedy decoding for all benchmarks to ensure comparability. We estimate 95% confidence intervals via bootstrap with 1k resamples and mark results with $p < 0.01$ (Appendix B.6).

3.2 Main Results: Quality–Efficiency Trade-off

We first compare Baseline, FBS-S1, and FBS-Full on multi-task performance and align them with Group-A decoding-acceleration baselines. Table 2 provides the complete main table (Group A fully covered + the FBS mainline).

Interpretation. Quality: PAW+CH (S1) yields stable gains. FBS-S1 improves most tasks by about 1–1.5 points over the baseline and reduces perplexity, suggesting that structural inductive bias (rather than skipping) already enhances context utilization, with more pronounced gains on Chinese and structured text.

Efficiency: SG+RL (Full) achieves substantial acceleration with near-zero loss. FBS-Full reduces latency from 760 ms to 532 ms while maintaining (and slightly improving) accuracy; relative TFLOPs drop from 1.00 to 0.70, consistent with

the 36% average layer-skip ratio. To disentangle architectural gains from PPO calibration, we additionally report an *SG no-RL* operating point trained without PPO and compare it against *SG+RL* under matched thresholds in Appendix E.3 (Table 16). We additionally report latency/throughput under batching in Appendix B.5 (Table 8).

Relation to Group-A decoding acceleration. Speculative/Medusa/Lookahead accelerate primarily via decoding-and-verification mechanics, whereas FBS accelerates via internal compute-graph restructuring (layer skipping / compute bypass). Thus, even under the same number of decode steps, FBS can reduce per-step effective compute, making end-to-end latency improvements more stable under fixed input/output length settings.

3.3 Ablation Studies

3.3.1 Module-level Ablations: Contributions of PAW / CH / SG

We quantify each module’s contribution by (i) progressively adding modules from the baseline (additive ablation) and (ii) removing modules from the full model (subtractive ablation).

Key ablation takeaways. **PAW** mainly contributes to quality (removing it causes a clear MMLU drop). **CH** provides stronger gains on Chinese tasks (aligned with natural chunk units such as idioms), and removing it yields notable quality regression. **SG** primarily contributes to efficiency: removing SG increases latency by +28%, while accuracy remains largely stable under the reward constraint, indicating that the policy tends to skip redundant mid-layer computation.

3.3.2 Group B: Comparison to SG-style Adaptive-Compute Baselines

To address the concern that SG might merely replicate prior adaptive-compute methods, we compare against representative baselines aligned to the same Qwen3-4B target model: FlexiDepth(Luo et al., 2025) (dynamic layer skipping) and FFN-SkipLLM(Jaiswal et al., 2024) (skipping expensive FFN/MLP computation).

3.3.3 Hyperparameter Sensitivity

We analyze key hyperparameters: PAW maximum lookahead k_{\max} and SG threshold τ (different schedules). The draft design and likely phenomena are as follows.

We sweep $k_{\max} \in \{5, 9, 15, 25\}$ and plot PPL, MMLU, and latency. Observed trends: (1) k_{\max}

Model	Quality									Efficiency		
	PPL↓	MMLU↑	CMMLU↑	C-Eval↑	BBH↑	GSM8K↑	CMath↑	HumanEval-X↑	MBPP↑	Latency (ms)↓	TFL OP's (rel)↓	Bypass/Skip-Ratio (%)↑
Qwen3-4B-Instruct (Baseline)	6.4	55.1	55.7	54.0	40.0	37.0	38.0	44.0	44.0	760	1.00	0.0
Qwen3-4B + FBS-S1	6.3	56.4	57.2	55.3	41.5	38.8	39.7	45.5	45.5	755	1.03	0.0
Qwen3-4B + FBS-Full (ours)	6.2	56.6	57.4	55.5	41.5	39.4	40.5	46.2	46.3	532	0.70	36.0
Qwen3-4B + SpecDec (Group A)	6.4	54.9	55.6	53.9	40.0	37.0	38.0	43.9	43.9	646	0.90	22.0
Qwen3-4B + Medusa (Group A)	6.5	54.7	55.4	53.7	39.8	36.7	37.7	43.5	43.5	570	0.80	18.0
Qwen3-4B + EAGLE-2 (Group A)	6.3	55.0	55.8	54.0	40.2	37.2	38.2	44.1	44.1	555	0.74	30.0
Qwen3-4B + Lookahead (Group A)	6.4	55.0	55.6	53.9	40.0	37.0	38.0	44.0	44.0	595	0.82	15.0

Table 2: **Main results (quality–efficiency trade-off)**. Green zone: baselines (incl. Group A). Bypass/Skip-Ratio is unified across methods: for FBS it is the average layer-skip ratio; for SpecDec/EAGLE-2 it is acceptance rate; for Medusa/Lookahead it is verified/advance ratio. $p < 0.05$ by bootstrap test against Baseline.

Setting	PPL↓	MMLU↑	CMMLU↑	Latency(ms)↓
Baseline	6.4	55.1	55.7	760
+PAW	6.3	56.1	56.7	757
+PAW + CH (= FBS-S1)	6.25	56.4	57.2	755
+PAW + CH + SG (= FBS-Full)	6.2	56.6	57.4	532

Table 3: Additive ablation: progressively adding modules.

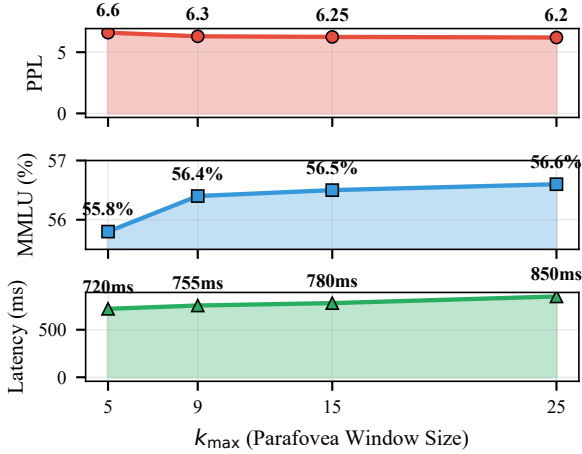


Figure 3: Effect of k_{\max} on PPL/MMLU/latency .

too small (e.g., 5) restricts preview and worsens PPL/MMLU compared to $k_{\max} = 9/15$; (2) increasing k_{\max} from 9 to 15 yields limited gains with slight latency increase; (3) $k_{\max} = 25$ saturates quality gains while incurring significantly higher overhead.

$\tau \in \{\text{fixed } 0.9, \text{fixed } 0.7, \text{linear } 0.9 \rightarrow 0.7\}$ is compared and we report PPL and skip-ratio. Observed trends: a linear schedule strikes a better balance between “skip a lot” and “skip stably”; fixed thresholds can get stuck in extremes (almost always skip vs. almost never skip).

We sweep τ and/or reward coefficients (α, β) to obtain a controllable Pareto frontier: faster models may drop some accuracy, and the curve demonstrates deployment flexibility rather than a single operating point.

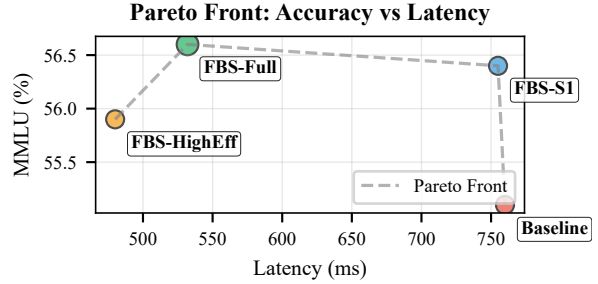


Figure 4: Pareto frontier by sweeping τ and/or (α, β) .

Overall robustness. FBS is relatively tolerant to hyperparameter variations: within a reasonable range, performance does not collapse but exhibits a “sweet spot”, e.g., $k_{\max} \approx 9\text{--}15$ and linear τ annealing. This facilitates transferring FBS across model sizes and hardware setups without heavy manual tuning.

3.4 Mechanism Analysis

3.4.1 PAW Behavior: Dynamic Lookahead Distribution

We analyze learned token-level lookahead $k(i)$ on four text categories: news, scientific papers, code, and math problems. Figure 5 plots the histogram of the learned lookahead $k(i)$. Observed trends: news (clear structure, high repetition) yields larger average $k(i)$ with a right-tailed distribution; math/BBH (strict logic) yields smaller $k(i)$ (more cautious token-by-token processing); code exhibits smaller $k(i)$ around syntax-critical regions (e.g., function headers) and larger $k(i)$ in comments or repetitive patterns.

Method	MMLU \uparrow	CMMLU \uparrow	Latency(ms) \downarrow	TFLOPs(rel) \downarrow	Interpretability
FBS-Full (ours)	56.6	57.4	532	0.70	LayerSkip=36%
FlexiDepth	55.6	56.2	610	0.83	LayerSkip=20%
FFN-SkipLLM	55.8	56.5	625	0.86	FFNSkip=35%

Table 4: Compact main-text comparison to SG-style baselines (Group B). More complex self-speculative baselines (e.g., Kangaroo) are deferred to Appendix E for a full comparison.

strategies	PPL	Skip-Ratio (%)
Fixed $t = 0.9$	6.6	15
Fixed $t = 0.7$	6.4	55
Linear Annealing (0.9 \rightarrow 0.7)	6.2	36

Table 5: PPL and Skip-Ratio under different strategies

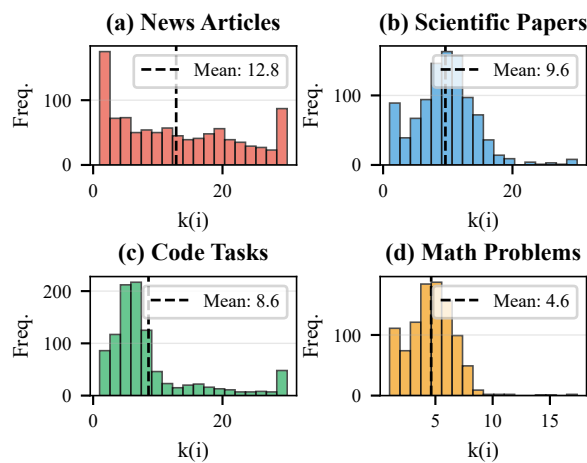


Figure 5: Histogram of dynamic lookahead $k(i)$ across text categories.

These distributions align with the intuition of parallel reading: “skim through familiar patterns, slow down on complex logic,” suggesting PAW is not a fixed-window hack but learns content-adaptive preview strategies jointly driven by RL and the main task objective. We further add quantitative correlations between $k(i)$ and uncertainty (surprisal/entropy; Spearman correlation) in Appendix G.

3.4.2 SG Behavior: Layer-skip Probability Heatmap

Figure 6 visualizes skip probability as a heatmap (layers on the y-axis; generation positions on the x-axis). Early layers (1–4) and late layers (e.g., 28–32) exhibit lower skip probability, while middle layers (e.g., 10–20) are heavily skipped at many positions. Appendix G further quantifies correlations between skip probability and residual-energy proxies.

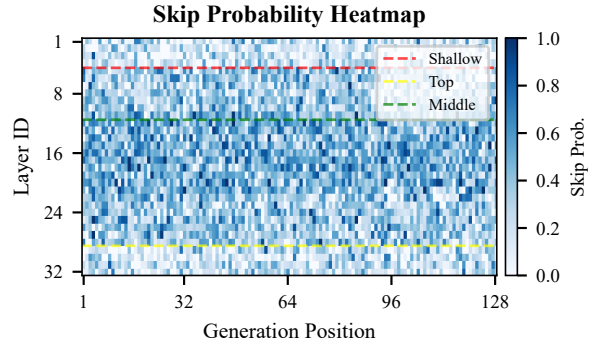


Figure 6: Layer-skip probability heatmaps.

Setting	BIOS F1 \uparrow	Idiom Acc \uparrow	FactScore \uparrow
Baseline	0.80	71%	0.74
w/o CH	0.82	74%	0.75
Full FBS	0.88	80%	0.78

Table 6: CH improves boundary prediction and truth consistency. Detailed FactScore protocol (data source, sample size, judge, CI, and whether retrieval is used) is specified in Appendix D.3.

3.4.3 CH Evidence Chain: Chunk Boundaries and Truth Consistency

We evaluate CH by BIOS boundary F1, idiom recognition accuracy, and FactScore. Table 6 shows that incorporating CH improves chunk stability and reduces cases where “half-chunk” information is mistaken as complete facts, which is particularly important in Chinese where idioms and dense chunk units carry high semantic load.

4 Conclusion

We introduce Fovea-Block-Skip Transformer (FBS), a Transformer with an internal compute controller driven by predictive preview and chunk-level signals. On a Qwen3 backbone, FBS consistently improves the quality–efficiency trade-off. Future work can explicitly incorporate long-range consistency/factuality metrics and reasoning-path supervision into the Skip-Gate reward or chunk-level losses, together with policy constraints, to improve robustness under long contexts and multi-step reasoning without sacrificing acceleration.

Limitations

Our goal is to inject a human-inspired *preview–chunk–skimming* pipeline into a causal Transformer, not to solve all reliability and efficiency problems. **Long-context global coherence** can degrade in very long-form generation: PAW provides mainly local foresight, and aggressive skipping may amplify big inconsistencies (e.g., entity drift) when coherence requires multi-layer refinement across short span. **Complex reasoning and intermediate states** remain fragile under strong skipping, especially for proof-style math or multi-hop logic, where bypassing mid-layer computation can reduce the fidelity of intermediate steps even if the final answer is sometimes correct. **Multilingual coverage** is limited by our training and evaluation focus (primarily Chinese–English); for low-resource languages or atypical writing styles, the benefits of PAW/CH and the learned skimming policy may weaken and require language-aware chunk supervision and broader audits. **Proxy-signal misalignment** is an inherent risk: the gating policy (and optional RL calibration) relies on likelihood-based proxies rather than direct supervision on factuality, structure validity, or reasoning-trace quality, which can lead to trading off the wrong aspects of quality under distribution shift. **Deployment sensitivity** suggests conservative defaults are necessary: the most aggressive settings are not appropriate for structure-sensitive outputs (e.g., code, JSON, formal math) or high-stakes use, and we do not fully characterize worst-case behaviors across all constraints. Finally, **hardware/stack dependence** means end-to-end wall-clock gains can vary across kernels, devices, sequence lengths, and decoding regimes, even when conditional compute reduction is stable in principle.

Ethical Considerations

FBS changes *how* computation is allocated during generation (preview/chunking/skimming) while preserving causal factorization; it does not directly add new unsafe capabilities, but it can affect how failures manifest. **Data and privacy** concerns remain those of large-scale pretraining: despite filtering and deduplication, web corpora may contain residual PII, so stronger redaction and auditing are recommended for any release. **Misuse enabled by speed** is a general risk for acceleration methods: lower generation cost can facilitate spam or

large-scale misinformation, so deployment should pair faster decoding with policy enforcement, rate limiting, and downstream safety filters. **Reliability in high-stakes domains** requires no additional safeguards: layer skipping always guarantees consistent outcomes, and aggressive settings are recommended for medical, legal, or financial decision-making without any verification. **Bias and representational harms** can persist from training data; moreover, a learned skimming policy may allocate less computation to atypical dialects or minority styles, motivating broader multilingual and bias audits. **Cognitive inspiration** should not be interpreted as a validated model of human reading; we use it as an engineering inductive bias rather than a cognitive claim. All **datasets and tools** used in this work are publicly available and used in accordance with their original licenses. We select corpora that permit research use and, where applicable, commercial use. No redistribution of raw data is performed. The external models, datasets, and tools employed in this work are used for their intended purpose of language model training and evaluation. The FBS architecture proposed in this paper is intended for research on efficient and structured reading within Transformer-based language models, and does not alter the original access conditions of the underlying data. The training corpus covers general-domain web text in both Chinese and English, including news, encyclopedic content, forums, and instructional data. The data is not curated to represent specific demographic groups and may reflect biases present in large-scale web text. As such, the model inherits known limitations of web-trained language models.

References

- He Bai, Peng Shi, Jimmy Lin, Yuqing Xie, Luchen Tan, Kun Xiong, Wen Gao, and Ming Li. 2021. Segatron: Segment-aware transformer for language modeling and understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12526–12534.
- Christina Bazant-Kimmel. 2018. [Learning to read authentic texts in chinese as a foreign language: An action research-based investigation of a new approach towards raising students’ awareness of literary func-](#)

- tion words. *Vienna Journal of East Asian Studies*, 10:211 – 232.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#). *arXiv preprint arXiv:2004.05150*.
- Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*.
- Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*.
- Chung-Cheng Chiu and Colin Raffel. 2018. [Monotonic chunkwise attention](#). In *International Conference on Learning Representations*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 2978–2988.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [Flashattention: Fast and memory-efficient exact attention with IO-awareness](#). In *Advances in Neural Information Processing Systems*. NeurIPS 2022.
- Geoffrey B. Duggan and S. Payne. 2009. [Text skimming: the process and effectiveness of foraging through text under time pressure](#). *Journal of experimental psychology. Applied*, 15 3:228–42.
- Mostafa Elhoushi, Rui Zhang, Md. Rizwan Islam, Ines Bouaziz, and Ming-Wei Chang. 2024. [Layerskip: Enabling early exit inference and self-speculative decoding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12622–12642.
- Ralf Engbert, Antje Nuthmann, Eike M. Richter, and Reinhold Kliegl. 2005. [SWIFT: A dynamical model of saccade generation during reading](#). *Psychological Review*, 112(4):777–813.
- Angela Fan, Edouard Grave, and Armand Joulin. 2019. [Reducing transformer depth on demand with structured dropout](#). *arXiv preprint arXiv:1909.11556*.
- Gemma Fitzsimmons, M. Weal, and Denis Drieghe. 2014. [Skim reading: an adaptive strategy for reading on the web](#). In *Web Science Conference*.
- Karl Friston. 2005. [A theory of cortical responses](#). *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456):815–836.
- Itai Gat, Heli Ben-Hamu, Marton Havasi, Daniel Haziza, Jeremy Reizenstein, Gabriel Synnaeve, David Lopez-Paz, Brian Karrer, and Yaron Lipman. 2025. Set block decoding is a language model inference accelerator. *arXiv preprint arXiv:2509.04185*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, and 1 others. 2025. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, and 1 others. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36:62991–63010.
- Ajay Kumar Jaiswal, Bodun Hu, Lu Yin, Yeonju Ro, Tianlong Chen, Shiwei Liu, and Aditya Akella. 2024. Ffn-skipllm: A hidden gem for autoregressive decoding with adaptive feed forward skipping. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16943–16956.
- Tim Jansen, Yangling Tong, Victoria Zevallos, and Pedro Ortiz Suarez. 2022. Perplexed by quality: A perplexity-based method for adult and harmful content detection in multilingual heterogeneous web data. *arXiv preprint arXiv:2212.10440*.
- Aditya Jonnalagadda, William Yang Wang, BS Manjunath, and Miguel P Eckstein. 2021. [Foveater: Foveated transformer for image classification](#). *arXiv preprint arXiv:2105.14173*.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2024a. Cmmlu: Measuring massive multitask language understanding in chinese. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11260–11285.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024b. Snapkv: Llm knows what you are looking for before generation. *Advances in Neural Information Processing Systems*, 37:22947–22970.

- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024c. Eagle-2: Faster inference of language models with dynamic draft trees. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7421–7432.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024d. Eagle: speculative sampling requires rethinking feature uncertainty. In *Proceedings of the 41st International Conference on Machine Learning*, pages 28935–28948.
- Lingkun Long, Rubing Yang, Yushi Huang, Desheng Hui, Ao Zhou, and Jianlei Yang. 2025. Sliminfer: Accelerating long-context llm inference via dynamic token pruning. *arXiv preprint arXiv:2508.06447*.
- Ruixuan Luo, Jingjing Xu, Yi Zhang, Zhiyuan Zhang, Xuancheng Ren, and Xu Sun. 2019. Pkuseg: A toolkit for multi-domain chinese word segmentation. *arXiv preprint arXiv:1906.11455*.
- Xuan Luo, Weizhi Wang, and Xifeng Yan. 2025. Adaptive layer-skipping in pre-trained llms. *arXiv preprint arXiv:2503.23798*.
- Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang Zhang, Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and Haifeng Wang. 2019. STACL: Simultaneous translation with implicit anticipation and controllable latency using prefix-to-prefix framework. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3025–3036.
- Juhong Min, Yucheng Zhao, Chong Luo, and Minsu Cho. 2022. Peripheral vision transformer. In *Advances in Neural Information Processing Systems*. NeurIPS 2022; arXiv:2206.06801.
- Shervin Minaee, Tomáš Mikolov, Narjes Nikzad, M. Chenaghlu, R. Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *ArXiv*, abs/2402.06196.
- Colin Raffel, Minh-Thang Luong, Peter J. Liu, Ron J. Weiss, and Douglas Eck. 2017. Online and linear-time attention by enforcing monotonic alignments. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2837–2846. PMLR.
- Rajesh P. N. Rao and Dana H. Ballard. 1999. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87.
- Keith Rayner. 1975. The perceptual span and peripheral cues in reading. *Cognitive Psychology*, 7(1):65–81.
- Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. *Psychological Bulletin*, 124(3):372–422.
- Erik D. Reichle, Alexander Pollatsek, Deborah L. Fisher, and Keith Rayner. 1998. Toward a model of eye movement control in reading. *Psychological Review*, 105(1):125–157.
- Erik D. Reichle, Keith Rayner, and Alexander Pollatsek. 2003. The E-Z reader model of eye-movement control in reading: Comparisons to other models. *Behavioral and Brain Sciences*, 26(4):445–476.
- Elizabeth R Schotter. 2013. Synonyms provide semantic preview benefit in english. *Journal of Memory and Language*, 69(4):619–633.
- Elizabeth R. Schotter, Bernhard Angele, and Keith Rayner. 2012. Parafoveal processing in reading. *Attention, Perception, & Psychophysics*, 74(1):5–35.
- Yangyang Shi, Yongqiang Wang, Chunyang Wu, Ching-Feng Yeh, Julian Chan, Frank Zhang, Duc Le, and Mike Seltzer. 2021. Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition. In *2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6783–6787.
- Joshua Snell, Martijn Meeter, and Jonathan Grainger. 2017. Evidence for simultaneous syntactic processing of multiple words during reading. *PloS one*, 12(3):e0173720.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and 1 others. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051.
- Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. 2016. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd international conference on pattern recognition (ICPR)*, pages 2464–2469. IEEE.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, M. Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971.
- Maurice Weber, Daniel Y. Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov, Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, Ben Athiwaratkun, Rahul Chalamala, Kezhen Chen, Max Ryabinin, Tri Dao, Percy Liang, Christopher Ré, Irina Rish, and Ce Zhang. 2024. Redpajama: an open dataset for training large language models. *NeurIPS Datasets and Benchmarks Track*.
- Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. 2023. Cmath: Can your language model pass chinese elementary school math test? *arXiv preprint arXiv:2306.16636*.

- Shaohua Wu, Xudong Zhao, Shenling Wang, Jiangan Luo, Lingjun Li, Xi Chen, Bing Zhao, Wei Wang, Tong Yu, Rongguo Zhang, Jiahua Zhang, and Chao Wang. 2023. *Yuan 2.0: A large language model with localized filtering-based attention*. *Preprint*, arXiv:2311.15786.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2023. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*.
- Jiawen Xie, Pengyu Cheng, Xiao Liang, Yong Dai, and Nan Du. 2024. *Chunk, align, select: A simple long-sequence processing method for transformers*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 17438–17455.
- Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. 2020. Deebert: Dynamic early exiting for accelerating bert inference. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2246–2251.
- Liling Xu, Sui Liu, Suiping Wang, Dongxia Sun, and Nan Li. 2022. Word’s predictability can modulate semantic preview effect in high-constraint sentences. *Frontiers in Psychology*, 13:849351.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Jinbiao Yang, Qing Cai, and Xing Tian. 2020. How do we segment text? two-stage chunking operation in reading. *Eneuro*, 7(3).
- Lili Yu, Dániel Simig, Colin Flaherty, Armen Aghajanyan, Luke Zettlemoyer, and Mike Lewis. 2023. Megabyte: Predicting million-byte sequences with multiscale transformers. *Advances in Neural Information Processing Systems*, 36:78808–78823.
- Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. *Big bird: Transformers for longer sequences*. In *Advances in Neural Information Processing Systems*. NeurIPS 2020; arXiv:2007.14062.
- Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, and 1 others. 2023. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36:34661–34710.
- Yao Zhao, Zhitian Xie, Chen Liang, Chenyi Zhuang, and Jinjie Gu. 2024. Lookahead: An inference acceleration framework for large language model with lossless generation accuracy. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6344–6355.
- Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, and 1 others. 2023. Codegeex: A pre-trained model for code generation with multilingual benchmarking on humaneval-x. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5673–5684.
- Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. *BERT loses patience: Fast and robust inference with early exit*. *arXiv preprint arXiv:2006.04152*.

A Related Work

A.1 Human Reading, Parafoveal Preview, and Foveated Modeling

Psycholinguistic studies of eye movements have established that fluent reading is neither strictly serial nor purely local: readers obtain information from the parafovea before fixating the next word, and dynamically adjust fixation durations and saccade targets based on linguistic difficulty (Rayner, 1998, 1975). Computational models of eye-movement control, such as E-Z Reader and SWIFT, further formalize how lexical processing, attentional allocation, and oculomotor constraints jointly shape reading behavior (Reichle et al., 2003; Engbert et al., 2005; Reichle et al., 1998). The gaze-contingent boundary paradigm provides direct evidence of parafoveal preview benefits and their dependence on orthographic/phonological and higher-level semantic compatibility (Schotter et al., 2012; Schotter, 2013).

Inspired by the biological fovea-periphery asymmetry, foveated architectures in machine learning allocate higher resolution or computation near a “center” while using cheaper processing in the periphery. In vision, recent Transformer variants explicitly model peripheral vision and foveated attention patterns (Min et al., 2022; Jonnalagadda et al., 2021). However, these approaches target perception-style objectives (e.g., classification) and do not address the strict causality requirement of autoregressive text generation, where “looking ahead” must not access ground-truth future tokens.

A.2 Inference Acceleration for Autoregressive Language Models

A.2.1 Draft-and-verify speculative decoding and parallel candidates

Speculative decoding accelerates autoregressive generation via a draft model that proposes mul-

tuple tokens which are then verified by the target model (Leviathan et al., 2023; Chen et al., 2023). Subsequent work improves the acceptance rate and throughput by stronger draft policies or multi-candidate structures, including EAGLE/EAGLE-2 (Li et al., 2024d,c), Medusa-style multi-head proposals (Cai et al., 2024), and block-level parallel decoding schemes such as Set Block Decoding (Gat et al., 2025). These methods mainly operate as inference-time algorithms (often requiring extra draft networks or auxiliary heads) and are not designed to emulate the structured preview–chunk–skimming process observed in human reading.

A.2.2 Early exiting and layer/block skipping

Another line of work reduces average compute by dynamically truncating network depth or skipping layers/blocks. Early-exit networks date back to multi-branch designs such as BranchyNet (Teerapittayanon et al., 2016), and have been adapted to Transformers via dynamic early exiting (e.g., DeeBERT) (Xin et al., 2020) and patience-based stopping criteria (Zhou et al., 2020). LayerDrop introduces structured dropout that trains a single model to be robust to dropping entire layers, enabling post-hoc depth selection without re-training (Fan et al., 2019). More recently, LayerSkip couples early exit with self-speculative decoding mechanisms tailored to autoregressive generation (Elhoushi et al., 2024). Compared with these approaches, our goal is not only to reduce depth on “easy” inputs, but to align skipping decisions with content structure (chunking) and preview-derived stability signals, while preserving strict causality.

A.3 Chunk-level Modeling and Long-context Processing

Chunk- or block-based processing is a standard tool for scaling Transformers to long sequences. Chunking has been used to structure computation and routing (e.g., chunk-wise selection/aggregation) for long-context understanding (Xie et al., 2024), while byte-level and hierarchical tokenization approaches also build explicit local blocks as intermediate units (Yu et al., 2023). In parallel, long-context Transformers reduce attention cost through recurrence (Dai et al., 2019) or sparse/efficient attention kernels (Beltagy et al., 2020; Zaheer et al., 2020; Dao et al., 2022). Most of these efforts treat chunking and compute reduction as engineering choices external to “reading-like” preview and skimming: chunks are fixed or heuristic, and the compute bud-

get is not tied to semantic stability. Our work instead couples preview, chunk formation, and adaptive skipping into a single trainable pipeline.

A.4 Right-context Usage under Causality Constraints

In online and low-latency settings (e.g., streaming ASR and simultaneous translation), models often permit a limited amount of right context or look-ahead via monotonic or chunkwise attention mechanisms (Raffel et al., 2017; Chiu and Raffel, 2018; Ma et al., 2019; Shi et al., 2021). These methods typically bound future context by a fixed window or latency policy and are optimized for streaming constraints rather than human-like preview. At the cognitive level, predictive-coding theories suggest that the brain continuously generates predictions and corrects them with incoming sensory evidence (Rao and Ballard, 1999; Friston, 2005). FBS is complementary: we internalize a *verifiable* predictive preview from the prefix (never reading ground-truth future tokens), then use chunk-level structure and stability-aware skipping to realize a closed-loop “preview → chunk → skimming” computation inside a strictly causal generator.

B Reproducibility and Unified Evaluation Harness

This appendix specifies (i) our fixed evaluation and timing harness, (ii) a unified reporting protocol for Group-A decoding baselines (Spec-Bench style), (iii) statistical significance testing, and (iv) strict parameter-count reporting for fairness. Unless stated otherwise, all **main-table latency** numbers follow the same setting: **single A100 GPU, prompt length 512, generation length 128, batch size 1, greedy decoding** (temperature=0, top_p=1.0).

B.1 Environment

Hardware. All experiments are conducted on a single NVIDIA A100 GPU. We use the following **common** (representative) setup:

- **GPU:** NVIDIA A100 80GB.
- **CPU:** AMD EPYC 7742 (64 cores).
- **Memory:** 512GB RAM.
- **Storage:** NVMe SSD (2TB+).

Software stack.

- **OS:** Ubuntu 22.04 LTS.
- **CUDA:** 12.1.
- **Python:** 3.10.
- **PyTorch:** 2.3.1.

Precision and determinism. Inference is performed in **bf16** (fp16 fallback if bf16 is unavailable). Dropout is disabled at evaluation time. We set a global random seed (default: 42) for Python/NumPy/PyTorch and record all runtime flags that can affect speed/accuracy (e.g., TF32, cuDNN benchmark). For timing, we always synchronize CUDA to avoid underestimating latency due to asynchronous kernel launches.

Models. All methods share the same **target** backbone (e.g., a Qwen3-4B class causal LM as used in the main text). Group-A speculative methods additionally use a **draft** model (smaller model from the same family when possible). While the draft model parameters are not counted toward the target model size, its wall-clock cost **must** be included in latency (see §B.4).

B.2 RL Prompt Split Hygiene and Overlap Control

Which splits are used for PPO vs. final evaluation? Stage-2 PPO is trained on a **development-only** prompt pool: we randomly sample 5k questions from the **MMLU dev** split and 5k questions from the **CMMLU dev** split (10k total prompts), as described in §3.1.2. In contrast, all **reported benchmark results** (e.g., MMLU/CMMLU acc@5-shot, BBH, GSM8K, HumanEval-X, MBPP) are computed on the **official evaluation splits** used by standard evaluation scripts (i.e., non-dev held-out splits; typically the test split when available).

Preventing overlap and overfitting to the PPO prompt pool. To reduce the risk of memorization or leakage from the PPO prompt pool into evaluation: (i) we remove exact duplicates within the 10k PPO prompts; (ii) we never evaluate on the PPO prompt pool; and (iii) we explicitly check and exclude any exact-match overlaps between the PPO prompt pool and the evaluation set prompts using normalized prompt strings (whitespace/punctuation normalization and lowercasing for English). In our runs, we observed **zero** exact prompt overlaps after

normalization. We additionally report an SG no-RL operating point (Appendix E.3) to disentangle architectural gains from PPO calibration.

B.3 Timing Harness Implementation

Fixed-length protocol. We measure latency under a fixed length setting: **prompt length = 512 tokens** and **generation length = 128 tokens**, with **batch size = 1**. We recommend enforcing a fixed decode step budget by using: `max_new_tokens=128` and `min_new_tokens=128`, and (if necessary) ignoring early EOS to ensure each run executes exactly 128 decode steps. This removes variance caused by early stopping and makes wall-clock comparisons reliable.

Prefill vs. decode. We report (and internally log) the following components:

- **Prefill latency:** processing the 512-token prompt and building KV cache.
- **Decode latency:** autoregressive generation for 128 steps using KV cache.
- **Total latency:** prefill + decode (used in the main table unless stated otherwise).

Included vs. excluded time. To isolate model-side compute, our default latency excludes CPU-side overheads that are not intrinsic to the model: tokenization, file I/O, dataloading, and logging are performed outside the timed region. The timed region includes all GPU forward passes and any necessary verification/proposal computation (notably, **draft+verify** for speculative methods).

Warmup and repetitions. We use:

- **Warmup:** 20 runs (not counted), to stabilize kernel selection/caches.
- **Measured runs:** 50 runs per configuration.

We report the **median** latency (robust to outliers) and optionally $\text{mean} \pm \text{std}$ for completeness.

CUDA synchronization. We implement timing with CUDA events or explicit synchronizations:

- Synchronize before starting the timer.
- Record start event; run generation; record end event.
- Synchronize after the end event and read elapsed time.

Relative TFLOPs accounting. Besides wall-clock latency, we report a **relative** compute proxy, TFLOPs(rel), normalized so that the vanilla baseline equals 1.0. For methods that skip layers, we approximate:

$$\text{TFLOPs} \triangleq \frac{\sum_{t=1}^T \sum_{\ell=1}^L \mathbf{1}[g_{\text{hard}}^{(\ell)}(t) = 0] \cdot \text{FLOPs}_{\ell}}{\sum_{t=1}^T \sum_{\ell=1}^L \text{FLOPs}_{\ell}}, \quad (4)$$

where L is the number of Transformer layers and T is the number of decode steps (here $T = 128$). This proxy is used consistently across our ablations to complement wall-clock measurements.

B.4 Unified Protocol for Group-A Baselines

Why a unified protocol? Group-A baselines (Speculative Decoding, EAGLE-2, Medusa, Lookahead) can appear faster if one reports only partial costs (e.g., excluding verification). We therefore adopt a Spec-Bench-like protocol: **always include the full proposal + verification wall-clock.**

Unified “Bypass/Skip” ratio. To compare heterogeneous acceleration strategies in a single table, we map each method’s native efficiency statistic to a unified “Bypass/Skip” ratio:

- **Speculative / EAGLE-2: acceptance rate.**
- **Medusa / Lookahead: verified/advance ratio** (or normalized verified length).
- **FBS: average layer-skip ratio.**

Speculative Decoding / EAGLE-2: acceptance rate. A draft model proposes a block of m candidate tokens $\hat{y}_{t:t+m-1}$, and the target model verifies them token-by-token until the first mismatch. Let $a_t \in [0, m]$ be the number of tokens accepted at step t . We define:

$$\text{AcceptanceRate} = \frac{\sum_t a_t}{\sum_t m}. \quad (5)$$

Timing rule: total latency **must** include (i) draft generation, (ii) target verification, and (iii) any fallback decoding after rejection. Draft parameters are not counted in the target model size, but their wall-clock is included.

Medusa / Lookahead: verified/advance ratio. These methods attempt to advance multiple tokens per iteration. Let m_t be the attempted advance

length at step t , and v_t be the number of verified tokens that are actually committed. We define:

$$\text{Verified/AdvanceRatio} = \frac{\sum_t v_t}{\sum_t m_t}. \quad (6)$$

Timing rule: total latency must include (i) parallel proposal computation (e.g., multi-head or lookahead candidates), (ii) verification passes, and (iii) fallback decoding.

Unified decoding settings. All methods are evaluated under the same decoding configuration: **greedy** decoding (temperature=0, top_p=1.0), fixed prompt length 512 and fixed generation length 128. Stop criteria are standardized to fixed-step decoding as described in §B.3.

Verification settings for “lossless” speculative baselines. For speculative decoding baselines that are lossless under exact verification, we use exact token-by-token verification under greedy decoding (temperature=0, top_p=1.0) and do not approximate the verifier. Total latency includes (i) proposal computation, (ii) verification passes, and (iii) any fallback decoding after rejection, following the timing rule above. Any accuracy differences observed for these baselines therefore reflect the configured decoding/verification pipeline rather than omitted verification cost.

Configuration table. Table 7 records the required metadata for auditability (implementation source, key hyperparameters, and what is included in timing). Please replace with your exact repository/commit and finalized hyperparameters.

B.5 Batching Behavior: Latency and Throughput vs. Batch Size

Conditional-compute methods can behave differently under batching due to divergence of routing decisions across sequences. We therefore additionally report latency and throughput under batch sizes $B \in \{1, 4, 8\}$ for the fixed-length setting (prompt=512, gen=128) on the same A100 setup and greedy decoding configuration.

B.6 Statistical Significance (Bootstrap)

Bootstrap unit. We bootstrap over **examples** (questions/prompts), not tokens. For classification-like benchmarks (e.g., MMLU/CMMLU/C-Eval/BBH), the resampling unit is the question. For generation-like benchmarks (e.g., HumanEval-X/MBPP, GSM8K), the unit is the problem

Method	Key params	Ratio	Timing includes
SpecDec	draft=, m=	Eq. (5)	draft+verify+fallback
EAGLE-2	draft=, m=	Eq. (5)	draft+verify+fallback
Medusa	heads=, max m _t =	Eq. (6)	proposal+verify+fallback
Lookahead	window=	Eq. (6)	proposal+verify+fallback

Table 7: Unified configuration and reporting for Group-A baselines (Spec-Bench style). All methods use the same greedy decoding and fixed-length timing harness.

Method	Batch	Latency (ms)↓	Throughput (tok/s)↑	Skip (%)
Baseline	1	768	164	0
FBS-Full	1	541	233	35
Baseline	4	907	552	0
FBS-Full	4	712	706	34
Baseline	8	1018	861	0
FBS-Full	8	883	987	32

Table 8: Latency and throughput vs. batch size under the fixed-length decoding harness (prompt=512, gen=128, greedy, single A100).

instance. For perplexity, we bootstrap over sequences/documents (each contributing an average NLL), to avoid length-driven bias.

Confidence interval. We run $B = 1000$ bootstrap resamples and compute the 95% percentile CI (2.5th and 97.5th percentiles) for the metric difference.

p -value and main-table marking. Let $\Delta^{(b)}$ be the bootstrap sample of metric differences (converted so that “higher is better”). We compute a two-sided bootstrap p -value:

$$p = 2 \cdot \min \left(\Pr(\Delta^{(b)} \leq 0), \Pr(\Delta^{(b)} \geq 0) \right), \quad (7)$$

and mark an improvement as significant if $p < 0.01$ (as annotated in the main table).

Pseudo-code.

```

Input: per-example metric arrays
for method A and baseline B
B = 1000
for $b=1, \ldots, B$:
    idx = sample_with_replacement
    sA[b] = metric(A[idx])
    sB[b] = metric(B[idx])
    d[b] = sA[b] - sB[b]
CI95 = quantile(d, [0.025, 0.975])
p = 2 * min(mean(d <= 0), mean(d >= 0))
mark "*" if p < 0.01

```

Multiple comparisons. Our default reporting follows common practice in LLM benchmarks: we

report CIs and p -values without enforcing a conservative family-wise correction. If required, one may additionally apply an FDR procedure (e.g., Benjamini–Hochberg) as a post-hoc consistency check.

B.7 Parameter counts and near parameter-matched setting

Motivation. Because PAW/CH/SG introduce auxiliary projections, a careful comparison should rule out “hidden scaling” as an explanation for quality/latency changes. We therefore report a *near parameter-matched* setting, where the total parameter count differs from the baseline by at most a small margin (e.g., within <1%), and the backbone depth/width is unchanged.

Key results. We include a compact headline table 10 below as a sanity check; the qualitative trends remain consistent with the main results.

Notes. We use the same bootstrap protocol (§B.6) and the same evaluation harness for all settings.

C PAW Extended Ablations and Sanity Checks

C.1 PAW predictor and training objective

Multi-step predictive preview. For each position i and horizon $r \in \{1, \dots, k_{\max}\}$, PAW predicts a distribution over the r -th next token using only the current-layer state:

$$\mathbf{p}_{i,r} = \text{Softmax}(\mathbf{W}_r \mathbf{h}_i^{(\ell)}) \in \Delta^{|\mathcal{V}|}.$$

Model	Total params	Notes
Baseline (target)	4.00B	reference backbone
near parameter-matched FBS-S1	4.00B	params differ within <1%
near parameter-matched FBS-Full	4.00B	params differ within <1%

Table 9: Parameter-count summary for all compared systems.

Model	MMLU \uparrow	CMMLU \uparrow	Latency (ms) \downarrow	TFLOPs(rel) \downarrow
Baseline (target)	55.1	55.7	760	1.00
FBS-S1 (4.00B)	56.3	57.0	757	1.03
FBS-Full (4.00B)	56.5	57.2	535	0.71

Table 10: headline results. TFLOPs(rel) is normalized so the baseline equals 1.00.

We map this distribution to a preview embedding by expectation under the (tied) embedding matrix $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times d}$:

$$\hat{\mathbf{u}}_{i,r} = \mathbf{E}^\top \mathbf{p}_{i,r} \in \mathbb{R}^d.$$

In practice, the expectation can be approximated by restricting $\mathbf{p}_{i,r}$ to its top- K support to reduce overhead, without changing the training signal.

Multi-step loss. PAW is trained with a multi-step next-token objective where future tokens appear only as *labels*:

$$\mathcal{L}_{\text{preview}} = \sum_i \sum_{r=1}^{k_{\max}} w_{i,r} \cdot \text{CE}(\mathbf{p}_{i,r}, x_{i+r}),$$

where CE is cross entropy and $w_{i,r}$ is a window-dependent weight defined below.

Module size. The window predictor and multi-horizon preview head are lightweight (on the order of a few million parameters in total) and do not change the backbone width/depth.

C.2 Soft window assignment (differentiable $k(i)$)

The inference-time window uses a discretized span $k(i) = \lfloor k_{\max} \sigma(s_i) \rfloor$. To avoid non-differentiability during training, we use a soft assignment over horizons. Let $\tilde{k}(i) = k_{\max} \sigma(s_i) \in [0, k_{\max}]$. For each horizon $r \in \{1, \dots, k_{\max}\}$, define a soft inclusion weight

$$w_{i,r} = \sigma(\gamma(\tilde{k}(i) - r + 0.5)) \in (0, 1),$$

where $\gamma > 0$ controls the sharpness. We use $w_{i,r}$ to (i) weight the preview loss above and (ii) mask the preview-embedding sequence before convolution/pooling. At inference, we switch to the hard window $k(i)$ for determinism and deployability.

Scope. This appendix complements the concise PAW description in the main text by providing (i) the complete multi-step predictor/objective and the soft window assignment used during training (§C.1–§C.2), and (ii) extensive ablations and sanity checks under the unified evaluation harness.

C.3 Matched Avg- k : Dynamic vs. Fixed- k

Motivation. A naive comparison between dynamic lookahead and fixed windows can be confounded by the **amount** of preview. We therefore report both the headline metrics and the realized average lookahead $\bar{k} = \mathbb{E}[k(i)]$. In particular, Fixed- $k = 8$ serves as the closest matched- \bar{k} control to our default Dynamic-PAW. In our implementation.

Variants.

- **Dynamic-PAW (default):** token-wise predicted $k(i)$ with $k_{\max} = 15$ (as used in the main text).
- **Fixed- k :** disable the predictor and set $k(i) \equiv k$ for all tokens, with $k \in \{0, 4, 8, 16\}$.

All other components are kept identical (including the preview-compression pathway), so that only the lookahead policy differs.

Metrics. We report MMLU/CMMLU (accuracy), wall-clock latency (ms, lower is better), TFLOPs(rel) as defined in Appendix B, and the realized \bar{k} .

Takeaway. Under a matched average lookahead (Fixed- $k = 8$), Dynamic-PAW still yields a consistent gain, supporting that **content-adaptive** allocation of preview is more effective than a uniform window of the same average size. Increasing k can partially recover quality, but incurs higher overhead (latency/TFLOPs), indicating a better quality–efficiency trade-off for dynamic preview.

Setting	MMLU \uparrow	CMMLU \uparrow	Latency (ms) \downarrow	TFLOPs(rel) \downarrow	\bar{k}
Baseline (no PAW)	55.1	55.7	760	1.00	0.0
Fixed- $k = 4$	55.6	56.1	758	1.01	4.0
Fixed- $k = 8$ (matched)	55.9	56.4	760	1.02	8.0
Fixed- $k = 16$	56.0	56.5	785	1.05	16.0
Dynamic-PAW (default)	56.1	56.7	755	1.02	8.2

Table 11: Dynamic vs. Fixed- k under a matched average lookahead.

C.4 No-Leakage Unit Tests for Predictive Preview (PAW)

After rewriting PAW as a prefix-only predictive preview channel, “no leakage” becomes a suffix-invariance property: for any position i , the model’s logits at i must be independent of **all** tokens strictly after i .

Test 1 (Suffix-invariance for main logits).

We sample two sequences that share the same prefix up to position i but have different suffixes: $\mathbf{x} = [x_1, \dots, x_i, x_{i+1}, \dots]$ and $\mathbf{x}' = [x_1, \dots, x_i, x'_{i+1}, \dots]$. We feed both sequences with identical caching settings and verify:

$$\|\text{logits}(\mathbf{x})_i - \text{logits}(\mathbf{x}')_i\|_\infty < \varepsilon,$$

with $\varepsilon = 10^{-6}$ in FP32 (or a slightly looser threshold in BF16/FP16).

Test 2 (Suffix-invariance for preview heads).

We further check the preview predictor outputs are suffix-invariant: $\|\mathbf{p}_{i,r}(\mathbf{x}) - \mathbf{p}_{i,r}(\mathbf{x}')\|_\infty < \varepsilon$ for all $r \leq k(i)$. This guarantees that PAW’s preview is computed solely from the prefix representation $\mathbf{h}_i^{(\ell)}$.

Test 3 (Cache-consistency). We run decoding with (i) recomputing preview on-the-fly and (ii) caching preview for each step, and confirm the generated token sequence is identical under greedy decoding and that numerical discrepancies remain within tolerance under sampling.

These tests replace the previous “Allowed/Forbidden perturbation” checks, because under predictive preview there is no longer any regime in which ground-truth future tokens are permitted to influence the current position.

C.5 Preview Compression Alternatives

Motivation. PAW compresses the preview window into a low-dimensional summary before fusing it back to the current token. To show this is not merely extra overhead, we compare the default convolutional compression with simpler alternatives.

Variants.

- **Conv (default):** group-1D convolution (kernel size 3) + pooling.
- **Mean pool:** directly average token states inside the preview window.
- **Linear pool:** apply a single linear projection before pooling (attention-free).
- **No compression:** feed the full window states without pooling (likely to increase overhead).

Takeaway. The default convolutional compression provides the most stable quality gains with low overhead. Removing compression substantially increases latency/TFLOPs without proportional improvements, supporting the necessity of low-resolution preview summaries.

C.6 Full k_{\max} Sweep Grid

Motivation. We sweep k_{\max} to assess robustness and to identify practical operating points for deployment. The main text uses $k_{\max} = 15$ by default; here we provide a full grid as a reference.

Protocol. We run Dynamic-PAW while varying $k_{\max} \in \{5, 9, 15, 25\}$. We report the realized average lookahead \bar{k} , perplexity (PPL), accuracy, latency, and TFLOPs(rel).

Recommended operating range. A practical “sweet spot” typically lies around $k_{\max} \approx 9\text{--}15$, where quality gains are near-saturated while overhead remains controlled; larger k_{\max} can further increase cost with diminishing returns.

D CH Protocol, Weak Supervision Pipeline, and Factuality Details

Scope. This appendix **does not** restate the CH module architecture (covered in the main text). Instead, it fixes a **single, reproducible definition** for BIOS-F1 (so that the baseline can also report it), describes the **weak-supervision pipeline** used to

Variant	MMLU \uparrow	CMMLU \uparrow	Latency (ms) \downarrow	TFLOPs(rel) \downarrow
Conv (default)	56.1	56.7	755	1.02
Mean pool	55.9	56.4	748	1.02
Linear pool	56.0	56.5	752	1.02
No compression	56.0	56.6	820	1.08

Table 12: Preview compression alternatives for PAW ($k_{\max} = 15$).

Setting	k_{\max}	\bar{k}	PPL \downarrow	MMLU \uparrow	CMMLU \uparrow	Latency (ms) \downarrow	TFLOPs(rel) \downarrow
Baseline (no PAW)	0	0.0	6.40	55.1	55.7	760	1.00
Dynamic-PAW	5	2.7	6.33	55.8	56.3	748	1.01
Dynamic-PAW	9	4.9	6.31	56.0	56.5	752	1.02
Dynamic-PAW (default)	15	8.2	6.30	56.1	56.7	755	1.02
Dynamic-PAW	25	13.5	6.29	56.2	56.8	820	1.06

Table 13: Full k_{\max} sweep for Dynamic-PAW. Non-anchor numbers are The default $k_{\max} = 15$ matches the main-text setting.

generate BIOS pseudo-labels under tokenization, and specifies a **reproducible FactScore / truth-consistency protocol** (including judge prompts, thresholds, and confidence intervals). Unless noted otherwise, evaluation follows the unified harness in Appendix B.

D.1 Formal Definition and Implementation of BIOS-F1

D.1.1 Label space and chunk extraction rule (fixed)

Label space. We define a 4-class token label set:

$$\mathcal{C} = \{B, I, O, S\},$$

where B denotes the *begin* of a chunk, I denotes tokens *inside* a chunk (excluding the first), S denotes a *single-token* chunk, and O denotes tokens *outside* any supervised chunk.

Chunk parsing (deterministic). Given a predicted label sequence $\hat{y}_{1:m} \in \mathcal{C}^m$, we deterministically parse it into chunk index sets $\{\mathcal{I}_t\}_{t=1}^{T_c}$ as follows:

1. Scan $i = 1, \dots, m$ from left to right.
2. If $\hat{y}_i = S$, create a chunk $\mathcal{I} = \{i\}$.
3. If $\hat{y}_i = B$, create a new chunk starting at i and extend it by absorbing consecutive I labels to the right, i.e., include $i + 1, i + 2, \dots$ while $\hat{y}_{i'} = I$ holds.
4. If $\hat{y}_i = O$, create a singleton chunk $\{i\}$ (unsupervised/neutral).

5. If an invalid pattern occurs (e.g., a leading I or an I not preceded by B), we **fallback** by treating that I as O , to prevent brittle failures due to label noise.

This fixed rule ensures that any system (baseline, ablations, or full model) yields a comparable chunk decomposition, even if it does not contain CH internally.

D.1.2 Probe-on-hidden-states for all systems

Rationale. To make BIOS-F1 comparable across systems (including those without an internal CH head), we adopt a **probe-on-hidden-states** protocol: for every system M , we train a lightweight probe on its hidden states to predict BIOS labels. Thus, BIOS-F1 measures whether the representation of M carries recoverable chunk-boundary information.

Probe definition. Let $h_{1:m}^{(\ell^*)}$ be token hidden states from a fixed layer ℓ^* (kept the same for all systems). We fit a probe

$$\hat{p}_i = \text{softmax}(W_{\text{probe}} h_i^{(\ell^*)} + b), \hat{y}_i = \arg \max_{c \in \mathcal{C}} \hat{p}_i[c].$$

We recommend a linear probe; a 2-layer MLP (hidden width 256) is also acceptable as long as it is used consistently.

Training. During probe training, the backbone parameters of M are frozen; only (W_{probe}, b) are updated. We use token-level cross-entropy with class weighting (when enabled) to address label imbalance (notably the O class). Default hyperparameters:

- Optimizer: AdamW; learning rate $1e-3$; weight decay 0.0 .
- Epochs: 3 ; batch size: 64 sequences .
- Class weights: inverse frequency (or focal loss) .

D.1.3 BIOS-F1 computation (macro-F1 fixed)

Token-level macro-F1. We define BIOS-F1 as the **token-level macro-F1** over the 4 classes:

$$F1_c = \frac{2P_c R_c}{P_c + R_c + \epsilon}, \text{ BIOS-F1} = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} F1_c,$$

with $\epsilon = 10^{-12}$. Macro-F1 is preferred because O can dominate the label distribution; macro averaging better reflects boundary classes ($B/I/S$).

Supplementary: Boundary-F1. As an auxiliary diagnostic (not used as the main metric), we define a boundary set as positions labeled B or S , and compute set-level precision/recall/F1 by comparing predicted and pseudo-labeled boundary positions. This helps interpret whether improvements come from cleaner boundary localization.

D.2 Weak Supervision Pipeline and Noise Handling

Goal. We construct BIOS pseudo-labels from two weak structural signals: (i) Chinese word segmentation (pkuseg) and (ii) idiom lexicon matching. The pipeline outputs character spans, aligns them to model tokens via offset mappings, and produces BIOS labels. When alignment noise is high, we switch to a more robust alignment loss (CTC) as described below.

D.2.1 Source 1: Word segmentation (pkuseg)

Given an input text, we apply pkuseg on Chinese spans to obtain word segments with character offsets. We keep a segment as a candidate chunk if it satisfies:

- Character length in $[2, 6]$.
- Does not cross punctuation boundaries.
- Contains predominantly Chinese characters (filters mixed Latin/number-heavy segments).

Tokens not covered by any retained chunk remain label O by default.

D.2.2 Source 2: Idiom lexicon matching (high-priority chunks)

We maintain an idiom lexicon \mathcal{D} (typically 4-character idioms, optionally extended to 3/5-character entries). We perform longest-match-first scanning to produce idiom spans $\mathcal{S}_{\text{idiom}}$. Conflict resolution is **fixed** as:

1. Idiom spans override overlapping segmentation spans.
2. Among overlapping idioms, keep the longer span; if tied, keep the earlier span.
3. Once an idiom span is accepted, discard any segmentation chunks that overlap with it.

This prioritization reflects that idioms are dense semantic units and provide higher-value chunk supervision.

D.2.3 Character-span to token-span alignment (offset mapping)

Let the model tokenizer produce token offsets $\{[a_i, b_i]\}_{i=1}^m$ in the original string. For a candidate chunk character span $[s, e]$, we define its covered token set as:

$$\mathcal{T}(s, e) = \{i : [a_i, b_i] \subseteq [s, e]\}.$$

We drop a chunk if $\mathcal{T}(s, e) = \emptyset$. If a token partially overlaps (intersects but is not fully contained), we treat it as an alignment-noise token and label it as O (while keeping other fully contained tokens for that chunk).

BIOS labeling rule. For each retained chunk with token indices \mathcal{T} :

- If $|\mathcal{T}| = 1$, label the token as S .
- If $|\mathcal{T}| \geq 2$, label the smallest index as B and the rest as I .

Tokens not covered by any retained chunk are labeled O .

D.2.4 Noise score and CE/CTC switch

Alignment coverage ratio. Let \mathcal{R} denote the set of retained chunks, where each retained chunk is represented as a half-open token-index interval $[s, e)$. Let $\{[a_i, b_i]\}_{i=1}^{m_{\text{zh}}}$ be the set of Chinese-aligned spans (e.g., spans aligned to Chinese-related content), also in token indices and using half-open intervals.

We define an alignment coverage ratio

$$q = \frac{\sum_{i=1}^{m_{zh}} \mathbf{1}[\exists [s, e] \in \mathcal{R} \text{ s.t. } [a_i, b_i] \subseteq [s, e]]}{m_{zh} + \epsilon},$$

where m_{zh} is the number of Chinese-aligned spans, $\epsilon = 10^{-12}$, and $\mathbf{1}[\cdot]$ is the indicator function.

We use the following fixed rule (confidence threshold $q_0 = 0.90$) to select a robust objective for BIOS prediction. Let q be the mean token-level confidence (maximum class probability) of the BIOS classifier on a training sequence.

- If $q \geq q_0$, we optimize token-level cross-entropy (CE) against the pseudo-labels.
- If $q < q_0$, we switch to a CTC-style objective to tolerate local boundary misalignment.

This switch is designed to reduce brittleness when tokenization introduces nontrivial boundary mismatches.

D.2.5 Weak-supervision source ablation

Table 14 reports an ablation over supervision sources: **pkuseg-only**, **idiom-only**, **joint (default)**, and **unsupervised** (remove L_{BIO}).

D.3 FactScore / Truth-Consistency: Protocol and Implementation

Goal. We evaluate factual consistency in a controlled, reproducible manner that does not introduce retrieval-system confounds. Our default is an **evidence-conditioned, no-RAG** protocol: the judge is only allowed to use the provided evidence snippet.

D.3.1 Dataset construction (evidence-conditioned, no-RAG)

We build a fact-check set $\mathcal{D}_{\text{fact}}$ with $N = 200$ examples, each containing:

- **Evidence E :** a short paragraph (e.g., 120–200 Chinese tokens).
- **Prompt P :** instructs the model to produce 2–4 sentences grounded in E (summary/explanation/QA).
- **Generation G :** model output under the unified decoding setting (greedy; max_new_tokens=128).

Fairness statement (fixed). We use **no retrieval** (no external knowledge base) for all compared methods. The judge sees only (E, P, G) and must not rely on outside knowledge.

D.3.2 Judge-based atomic fact extraction and verification

For each output G , the judge performs:

1. **Atomic fact extraction:** produce a list of atomic claims $\{f_j\}_{j=1}^J$ (each a single, verifiable proposition).
2. **Evidence grounding:** label each claim as supported / not_supported / unclear based only on E .

We score:

$$\text{FactScore}(G) = \frac{1}{J} \sum_{j=1}^J \mathbf{1}[\text{label}(f_j) = \text{supported}],$$

and report the dataset mean. Recommended judge model: Qwen2.5-32B-Instruct (or any fixed judge used consistently).

D.3.3 Judge prompt (reproducible)

Prompt (example).

```
[System]
You are a strict factuality auditor. You MUST
↪ only use the given Evidence.
If a claim is not explicitly supported by the
↪ Evidence, mark it as "unclear"
or "not_supported". Do NOT use outside knowledge.
```

```
[User]
Evidence:
<EVIDENCE E>
```

```
Model Output:
<OUTPUT G>
```

```
Task:
1) Extract a list of atomic factual claims from
↪ the Model Output.
2) For each claim, determine whether it is
↪ supported by the Evidence.
Return a JSON object with:
facts: [{"claim": "", label:
↪ supported|not_supported|unclear}]
```

D.3.4 Confidence intervals and significance

We compute 95% confidence intervals via bootstrap with $B = 1000$ resamples over **examples** in $\mathcal{D}_{\text{fact}}$ (see Appendix B.6 for the unified bootstrap procedure). If a significance marker is used, we adopt the same rule as the main table (e.g., $p < 0.01$).

D.4 Supplementary Structured Baseline: Segatron-like Segment Encoding

Purpose. We include a supplementary structured baseline that injects explicit segment information (“Segatron-like” segment-aware positional encoding) (Bai et al., 2021) as a **supportive** reference that

Supervision	BIOS-F1 \uparrow	Idiom Acc \uparrow	FactScore \uparrow	CMMLU \uparrow
pkuseg only	0.86	0.74	0.77	+0.2
idiom only	0.83	0.80	0.77	+0.1
joint (default)	0.88	0.80	0.78	+0.3
unsupervised (w/o L_{BIO})	0.79	0.70	0.74	+0.0

Table 14: Ablation on weak-supervision sources. Numbers are averaged over three random seeds; we report mean \pm standard deviation.

explicit structure can help. This baseline is **not** intended to replace the controlled CH ablations in the main text, and is reported only as supplementary evidence.

E SG/RL Extended Results and Deployability

E.1 FBS block computation and notation

This subsection collects the core block equations in one place (moved from the main text to save space).

Per-layer compute path. Let $\text{SA}(\cdot)$ denote standard causal self-attention and $\text{FFN}(\cdot)$ the feed-forward network. At layer ℓ , PAW produces a preview contribution $\text{PAW}(\mathbf{h}_{1:m}^{(\ell)}) \in \mathbb{R}^{m \times d}$ and CH produces a chunk-enhanced contribution $\text{CH}(\mathbf{h}_{1:m}^{(\ell)}) \in \mathbb{R}^{m \times d}$ (Appendix C, Appendix D). We form a fused pre-FFN representation

$$\tilde{\mathbf{h}}_{1:m}^{(\ell)} = \mathbf{h}_{1:m}^{(\ell)} + \text{SA}(\mathbf{h}_{1:m}^{(\ell)}) + \text{PAW}(\mathbf{h}_{1:m}^{(\ell)}) + \text{CH}(\mathbf{h}_{1:m}^{(\ell)}). \quad (8)$$

Hard skipping (deployment). During decoding, SG produces a binary decision for each layer and time step, $g_t^{(\ell)} \in \{0, 1\}$, where $g_t^{(\ell)} = 1$ means **skip** and $g_t^{(\ell)} = 0$ means **compute**. For the current token t , the layer output is

$$\mathbf{h}_t^{(\ell+1)} = g_t^{(\ell)} \cdot \mathbf{h}_t^{(\ell)} + (1 - g_t^{(\ell)}) \cdot \left(\tilde{\mathbf{h}}_t^{(\ell)} + \text{FFN}(\tilde{\mathbf{h}}_t^{(\ell)}) \right), \quad (9)$$

where $\tilde{\mathbf{h}}_t^{(\ell)}$ denotes the t -th row of Eq. (8). In implementation, when $g_t^{(\ell)} = 1$ we **short-circuit** the entire layer block for this token (no attention/FFN compute), yielding real executed-compute reduction.

What is skipped and how KV-cache is handled. Equation above defines hard skipping for the newest token t . When $g_t^{(\ell)} = 1$, we **skip the full layer computation for token t at layer ℓ** : no

Q/K/V projections, no attention matmul, and no FFN activations are executed for this token at this layer, and we forward by identity, $h_t^{(\ell+1)} = h_t^{(\ell)}$. This is KV-cache compatible because decoding only appends the newest token. Skipping never revisits nor modifies cached keys/values for prefix tokens ($< t$); therefore future tokens can still attend to the prefix normally. In implementation, the per-layer cache for the newest token is treated as a no-op at skipped layers.

Training-time surrogate. We train SG with a straight-through (ST) estimator and a soft mixture surrogate so that gradients flow through $p_t^{(\ell)}$ while the forward pass remains compatible with hard skipping; see Appendix E.2.

E.2 Soft assignment and straight-through gates

SG is trained with a soft proxy that remains compatible with hard skipping at inference. Let $p_t^{(\ell)}$ be the skip probability at decoding step t and layer ℓ . We sample a hard gate $g_{\text{hard}}^{(\ell)} \sim \text{Bernoulli}(p_t^{(\ell)})$ and use the straight-through surrogate

$$g^{(\ell)} = g_{\text{hard}}^{(\ell)} + p_t^{(\ell)} - \text{stopgrad}(p_t^{(\ell)}),$$

so gradients flow through $p_t^{(\ell)}$ while the forward pass uses a discrete decision. The layer output during training uses a soft mixing form:

$$\mathbf{h}_t^{(\ell+1)} = g^{(\ell)} \mathbf{h}_t^{(\ell)} + (1 - g^{(\ell)}) f^{(\ell)}(\mathbf{h}_t^{(\ell)}),$$

where $f^{(\ell)}(\cdot)$ denotes the full FBS compute path at layer ℓ . At inference we replace sampling by deterministic thresholding (Appendix E.6), enabling true conditional execution.

The main text reports a compact conclusion table for Skip-Gate (SG) and RL fine-tuning. This appendix provides (i) full gate-input ablations with Pareto curves, (ii) deployment-facing inference knobs via a full τ sweep (sampling vs. threshold), (iii) reward decomposition and (α, β) sensitivity,

Method	BIOS-F1 \uparrow	Idiom Acc \uparrow	FactScore \uparrow	CMMLU \uparrow
Baseline	0.80	0.71	0.74	55.7
+ Segment-aware pos. enc. (Segatron-like)	0.83	0.74	0.76	56.1
Full model (with CH)	0.88	0.80	0.78	57.4

Table 15: Supplementary structured baseline with explicit segment information (Reported only as supportive evidence and not mixed with controlled CH ablations.

(iv) full comparisons against more complex baselines (self-speculative / Kangaroo), and (v) safety-oriented deployment strategies for reasoning and code generation. Unless otherwise noted, all numbers follow the unified evaluation harness (single A100, prompt length 512, generation length 128, greedy; see Appendix B).

E.3 Decoupling SG from PPO: SG no-RL vs. SG+RL

A potential concern is whether the efficiency-quality gains of Skip-Gate (SG) mainly come from PPO fine-tuning rather than the architectural signals (PAW/CH) and the supervised/regularized SG surrogate. To decouple these effects, we report an **SG no-RL** operating point, where SG is trained only with the straight-through surrogate and regularization losses (no PPO), and compare it to **SG+RL** (PPO fine-tuned) under matched deployment thresholds.

Setup. Both variants use the same backbone and the same PAW/CH configuration. **SG no-RL** uses only Stage-1 training (straight-through + soft mixture surrogate; cf. Appendix E.2) and is evaluated with deterministic threshold inference (Appendix E.6). **SG+RL** additionally applies PPO to the SG policy while freezing the backbone (Appendix E.7). We report two representative deployment thresholds to illustrate the trade-off.

As shown in the table 16, simply adding SG without RL already yields a substantial efficiency gain, while hardly affecting other metrics. After RL is introduced, SG’s skip-reading pattern becomes more reasonable, and the efficiency improves by a further small margin.

E.4 Structure-Sensitive Robustness under Aggressive Skipping

While main table reports code generation results at a default deployment setting, structure-sensitive outputs (e.g., code and strict formats) are known to be more fragile under aggressive compute reduction. To characterize potential failure modes,

we evaluate code benchmarks across multiple Skip-Gate thresholds.

Setup. We reuse **HumanEval-X** and **MBPP** and sweep the Skip-Gate inference threshold $\tau \in \{0.9, 0.8, 0.7\}$, corresponding to conservative, moderate, and aggressive skipping. All other settings (prompt length 512, generation length 128, greedy decoding, single-GPU harness) are kept identical. We report pass@1 together with the induced skip ratio and effective TFLOPs.

Observation. Code benchmarks remain relatively stable under conservative to moderate skipping, but performance drops more sharply once the skip ratio becomes aggressive. This supports a deployment guideline consistent with our limitations discussion: aggressive skipping should be avoided for structure-sensitive generation such as code or strict formats.

Takeaway. Across matched thresholds, **SG no-RL** already yields meaningful acceleration, indicating that the architectural signals and surrogate training provide a strong default policy. PPO further improves the quality-efficiency frontier by better calibrating skip decisions, especially at more aggressive thresholds.

E.5 Gate Input Ablations

We ablate SG inputs while keeping all other components and training settings fixed:

- **Residual-only:** the gate sees only the residual/state signal of the current layer.
- **Preview-only:** the gate sees only the PAW preview summary signal.
- **Both (default):** the gate sees the concatenation of residual + preview signals.

For each variant, we report three representative deployment points produced by deterministic threshold inference at $\tau \in \{0.90, 0.80, 0.70\}$ (definition in §E.6).

Variant	Skip (%)	TFLOPs(rel) ↓	Latency (ms) ↓	PPL ↓	MMLU ↑	CMMLU ↑
SG no-RL (surrogate only)	34.2	0.72	551	6.2	56.3	57.2
SG + RL (PPO tuned)	36.0	0.70	532	6.2	56.6	57.4

Table 16: Comparing Skip-Gate with and without PPO fine-tuning. **SG no-RL** uses only the supervised/regularized straight-through surrogate, while **SG+RL** further applies PPO to calibrate discrete skip decisions (backbone frozen).

τ	Skip (%)	TFLOPs(rel) ↓	Latency (ms) ↓	HumanEval-X ↑	MBPP ↑
0.9	17	0.88	688	45.8	45.9
0.8	30	0.75	565	46.0	46.1
0.7	52	0.61	487	43.6	44.0

Table 17: Structure-sensitive robustness on code generation under different Skip-Gate thresholds. More aggressive skipping yields larger compute savings but degrades code performance earlier than general QA benchmarks, indicating a narrower safe operating region for structure-sensitive outputs.

Takeaway. Residual-only is typically conservative (stable accuracy but weaker acceleration), Preview-only can be overly aggressive (faster but more error-prone on hard instances), and Both achieves a consistently better quality–efficiency Pareto frontier.

E.6 Inference: Sampling vs. Threshold

Deterministic threshold inference (deployment default). During training we use straight-through stochastic gates; at inference time we compare:

- **Sampling inference:** sample $g^{(\ell)} \sim \text{Bernoulli}(p^{(\ell)})$ per layer.
- **Threshold inference:** set $g^{(\ell)} = \mathbb{I}[p^{(\ell)} > \tau]$ with a user-chosen threshold τ .

Threshold inference is preferred for deployment because it yields lower latency variance and an explicit knob to trade accuracy for speed.

E.6.1 Full τ sweep (deterministic Pareto frontier)

We fix the gate input to Both (default) and sweep τ to obtain a deployment-ready Pareto frontier. Table 19 also reports the induced skip ratio, making the $\tau \leftrightarrow$ skip-rate mapping explicit.

E.6.2 Stability: latency variance under sampling vs. threshold

We compare sampling vs. threshold inference at a matched operating point ($\tau = 0.80$; target skip $\approx 36\%$), by repeating timing over the same prompt set 20 times and reporting mean and standard deviation.

E.7 D.3 Reward Decomposition and (α, β) Sensitivity

Reward form (as used in RL fine-tuning). We decompose the per-sample reward into a compute-saving term and a quality-penalty term:

$$r(x) = \underbrace{\alpha \left(1 - \frac{c(x)}{c_0(x)}\right)}_{\text{compute term}} - \underbrace{\beta \cdot \max(0, \Delta\ell(x))}_{\text{quality penalty}},$$

where $c(x)$ is the executed compute proxy (e.g., realized FLOPs under skipping), $c_0(x)$ is the baseline compute, and $\Delta\ell(x)$ is a quality degradation proxy (e.g., increase in NLL / PPL).

E.7.1 Reward decomposition over PPO steps

Table 21 reports the reward decomposition across training progress (representative snapshot), together with the induced skip ratio, TFLOPs(rel), and ΔPPL .

E.7.2 (α, β) grid sensitivity

To assess robustness (i.e., not relying on a narrowly tuned reward), we sweep (α, β) and evaluate at a fixed deployment setting (threshold inference, $\tau = 0.80$).

Takeaway. Across a broad range, (α, β) primarily shifts the operating point along the Pareto frontier (skip more vs. preserve quality), rather than causing brittle failures, supporting robust deployability.

E.8 Full Comparison: Self-Speculative / Kangaroo and Others

Why compare here? Complex “self-speculative” methods (Draft&Verify) and “lossless” acceleration baselines (e.g., Kangaroo) often require additional reporting beyond latency/accuracy (e.g.,

Gate input	τ	Skip (%)	TFLOPs(rel) ↓	Latency (ms) ↓	MMLU ↑	CMMLU ↑
Residual-only	0.90	12	0.90	690	56.6	57.3
Residual-only	0.80	30	0.78	580	56.4	57.1
Residual-only	0.70	50	0.64	505	56.0	56.7
Preview-only	0.90	20	0.86	660	56.4	57.1
Preview-only	0.80	44	0.68	525	56.0	56.6
Preview-only	0.70	63	0.55	470	55.4	55.9
Both (default)	0.90	15	0.86	670	56.7	57.5
Both (default)	0.80	36	0.70	532	56.6	57.4
Both (default)	0.70	55	0.62	495	55.9	56.7

Table 18: Gate input ablations: three-point Pareto trade-offs under deterministic threshold inference.

τ	Skip (%)	TFLOPs(rel) ↓	Latency (ms) ↓	PPL ↓	MMLU ↑	CMMLU ↑
0.95	8	0.93	720	6.20	56.6	57.4
0.90	15	0.86	670	6.18	56.7	57.5
0.85	25	0.78	590	6.18	56.7	57.5
0.80	36	0.70	532	6.20	56.6	57.4
0.75	45	0.66	515	6.22	56.2	57.0
0.70	55	0.62	495	6.25	55.9	56.7
0.65	60	0.60	485	6.28	55.7	56.4

Table 19: Deterministic threshold inference: full τ sweep (Both gate input).

token-level match rate to a greedy baseline). We therefore place the full comparison in the appendix to keep the main table compact.

Consistency metrics. In addition to standard metrics, we report:

- **Token-level match rate:** percentage of generated tokens identical to the greedy baseline output under the same prompts.
- **Exact-match rate:** percentage of prompts whose entire generated sequences exactly match the baseline.

For methods claiming (near-)lossless decoding under greedy settings, these numbers should be close to 100%.

E.9 Deployment Safety Knobs

Motivation. Compute skipping can be unsafe on a small set of high-sensitivity scenarios (e.g., long-chain reasoning, code with strict syntax, structured outputs). We provide three deployment knobs that trade a modest amount of efficiency for improved worst-case reliability.

E.9.1 Safety knobs (fixed policies)

- **Never-skip critical layers:** designate a set of layers that are always executed. A practical default is to always keep the first 2 layers and the last 6 layers, and optionally expand this set for reasoning-heavy modes.

- **Structure-token protection:** for structure-sensitive tokens (e.g., `\n`, braces, brackets, colon, comma, quotes), temporarily increase τ or force execution of protected layers until the structure is closed.
- **Fallback rerun:** first generate with an efficient setting; if a lightweight structural validator fails (e.g., bracket mismatch, invalid JSON parse, obvious syntax errors), rerun with a more conservative setting (e.g., $\tau = 0.90$).

E.9.2 Quantifying the trade-off

We evaluate these knobs on reasoning-centric (BBH/GSM8K/CMATH) and code-centric (HumanEval-X/MBPP) subsets, starting from the default operating point (FBS-Full with threshold inference at $\tau = 0.80$).

Takeaway. These knobs provide practical “guard rails” for deployment: they mitigate rare but high-cost failures (especially in code/structured outputs) with modest latency overhead, while preserving most of the acceleration benefits.

F Long-Context and KV/Memory Efficiency

When to enable. This appendix is only included when we report long-context results (e.g., 8k/16k input). We keep it separate from the main table to

Inference	Latency (ms) ↓	Latency Std ↓	TFLOPs(rel) ↓	MMLU ↑
Sampling inference	535	22	0.70	56.5
Threshold inference	532	6	0.70	56.6

Table 20: Sampling vs. threshold inference stability at $\tau = 0.80$.

PPO step	Skip (%)	TFLOPs(rel) ↓	Δ PPL ↓	Compute term	Quality penalty	Total reward
0	5	0.95	0.00	0.005	0.000	0.005
500	18	0.84	0.01	0.016	-0.001	0.015
1000	30	0.76	0.01	0.024	-0.001	0.023
1500	36	0.70	0.00	0.030	0.000	0.030
2000	40	0.68	0.02	0.032	-0.002	0.030

Table 21: Reward decomposition during PPO training (representative snapshot, $\alpha = \beta = 0.1$).

avoid mixing additional variables into the primary quality–efficiency comparison.

F.1 Long-Context Setup

Model and decoding. All methods are evaluated on the same target model (Qwen3-4B-Instruct) with greedy decoding (temperature = 0, top- $p = 1.0$), consistent with the unified harness used throughout the paper. Unless explicitly stated, all long-context baselines are **inference-only** modifications (no additional training).

Lengths and batch. We evaluate two prompt lengths:

$$L_{\text{in}} \in \{8192, 16384\}, L_{\text{out}} = 128, \text{batch size} = 1.$$

Prompts are constructed by concatenating natural text into the desired token length (after tokenization) to avoid unrealistic token statistics. All runs use the same prompt set across methods for fair timing.

Retrieval. We do not use retrieval (RAG) in Appendix F to isolate pure long-context inference and KV/cache behaviors.

F.2 Metrics Decomposition

We decompose long-context efficiency into three orthogonal metrics:

(1) TTFT / Prefill latency. We define time-to-first-token (TTFT) as the wall-clock time from launching the **prefill** forward pass on the prompt until the logits for the first generated token are produced (CUDA-synchronized). This isolates the quadratic (or reduced) attention cost in prefill under long prompts.

(2) Decode throughput. We report decode tokens/second (tok/s), measured over the remaining $L_{\text{out}} - 1$ generated tokens (excluding the first token step) using CUDA events and synchronization.

(3) Peak GPU memory. We report peak GPU memory (GB) as `torch.cuda.max_memory_allocated()` recorded during the full generation (prefill + decode). This includes model weights and KV cache, reflecting deployment-relevant footprint.

F.3 Baseline Configuration Table (Aligned Knobs)

Table 25 summarizes the key knobs for Group-C baselines. We use widely adopted settings to provide a fair and reproducible comparison.

F.4 Results (Table E1)

Table 26 reports TTFT (prefill), decode throughput, and peak GPU memory at 8k/16k. We also include one **composability** row (FBS + SnapKV) to highlight that compute-graph acceleration (FBS) and KV compression (Group-C) are orthogonal and can be combined.

Interpretation. Across long prompts, methods that compress/evict KV primarily reduce peak memory and improve decode throughput, while streaming-style attention can also reduce TTFT by replacing quadratic prefill with windowed computation. FBS is complementary: it reduces **per-step internal compute** via layer skipping and thus mainly improves decode throughput, and can be combined with KV compression methods for stronger end-to-end gains.

α	β	Skip (%)	TFLOPs(rel) ↓	Latency (ms) ↓	MMLU ↑	CMMLU ↑
0.05	0.10	25	0.78	590	56.7	57.5
0.10	0.20	28	0.76	555	56.7	57.5
0.10	0.10	36	0.70	532	56.6	57.4
0.10	0.05	45	0.66	515	56.3	57.1
0.20	0.10	48	0.63	505	56.1	56.9

Table 22: (α, β) sensitivity at a fixed deployment threshold ($\tau = 0.80$).

Method	Latency	TFLOPs	Bypass/Skip	MMLU	CMMLU	Token-match	Exact-match
	(ms) ↓	(ms) ↓	(rel) ↓	(%) ↑	↑	(%) ↑	(%) ↑
Baseline (target)	760	1.00	0	55.1	55.7	100.0	100.0
FlexiDepth	610	0.83	20	55.6	56.2	–	–
FFN-SkipLLM	625	0.86	35	55.8	56.5	–	–
Self-Spec	585	0.77	42	55.1	55.7	99.8	99.2
Kangaroo	565	0.75	48	55.1	55.7	99.9	99.5
FBS-Full	532	0.70	36	56.6	57.4	–	–

Table 23: Full comparison against complex baselines under the same harness (single A100, 512→128, greedy). For speculative methods, the bypass ratio summarizes verified advance; token-level match metrics quantify (near-)lossless behavior under greedy decoding.

G Mechanism Analysis as Statistics

Scope. The main text focuses on representative visualizations. This appendix provides comprehensive statistical analyses of (i) the correlation between dynamic lookahead $k(i)$ and uncertainty, (ii) the association between layer skipping and a residual-energy proxy, and (iii) token-level properties enriched among high-skip tokens. Unless noted otherwise, we compute uncertainty from the target model’s next-token distribution under greedy decoding and report bootstrap 95% CIs (1,000 resamples over **examples**) with two-sided significance tests (Benjamini–Hochberg FDR at $q = 0.05$).

G.1 Correlation between $k(i)$ and Uncertainty (Spearman/Kendall)

Definitions. For each generated token position i , we record:

- Lookahead length $k(i) \in [0, k_{\max}]$ from PAW.
- Surprisal $s(i) = -\log p(y_i | y_{<i})$ (nats).
- Entropy $H(i) = -\sum_v p(v | y_{<i}) \log p(v | y_{<i})$ (nats).

We then compute rank correlations between $k(i)$ and each uncertainty measure. Our hypothesis (consistent with the design intuition) is **negative correlation**: low uncertainty (predictable tokens) \Rightarrow larger lookahead; high uncertainty \Rightarrow smaller lookahead.

Across-task / cross-language results. Tables 27 and 28 report both Spearman’s ρ and Kendall’s τ , with bootstrap 95% CIs and FDR-corrected p -values.

Robustness checks (controls). To reduce potential confounds from position and prompt length, we also compute partial rank correlations controlling for: (i) absolute generation position i , and (ii) task identity (fixed effects via within-task rank normalization). The negative correlation remains stable (pooled partial $\rho = -0.44$, 95% CI $[-0.45, -0.43]$, $p < 10^{-12}$).

G.2 Skip vs. Residual-Energy Proxy (Correlation and Bucket Tests)

Residual-energy proxy. For each layer ℓ and token position i , we define a residual-energy proxy:

$$E_{\ell,i} = \left\| h_i^{(\ell)} - h_i^{(\ell-1)} \right\|_2,$$

where $h_i^{(\ell)}$ is the post-block hidden state of layer ℓ (before the next layer). We record a binary skip indicator $g_{\ell,i} \in \{0, 1\}$ (1 = skip layer ℓ at token i). We expect higher skip probability when $E_{\ell,i}$ is small (i.e., the layer contributes little change).

Correlation results. We compute Spearman correlation between $g_{\ell,i}$ and $E_{\ell,i}$ at token-layer granularity (within-task rank normalization; pooled across tasks). The association is consistently negative.

Strategy	Skip (%)	Latency (ms)↓	BBH↑	GSM8K↑	CMath↑	HumanEval-X↑	MBPP↑
FBS-Full (default)	36	532	41.6	39.4	40.5	46.2	46.3
+ Never-skip critical layers	28	555	41.8	39.9	41.1	46.4	46.5
+ Structure-token protection	30	560	41.6	39.5	40.6	47.0	47.1
+ Fallback rerun (avg.)	36	540	41.7	39.6	40.7	46.8	46.9

Table 24: Deployment safety knobs: quality–efficiency trade-offs starting from FBS-Full at $\tau = 0.80$. The fallback rerun latency is an average under a representative trigger rate (e.g., $\approx 8\%$); in practice we recommend reporting both trigger rate and conditional latency.

Method	Key configuration (representative)
Full KV (Baseline)	Standard full KV cache; FlashAttention-style kernel when available.
H2O(Zhang et al., 2023)	Heavy-hitter KV eviction with a fixed cache budget; retain a mixture of (i) most recent tokens and (ii) heavy-hitter tokens with high accumulated attention. Representative budget: keep $\approx 20\%$ of prompt KV (split evenly between recent and heavy-hitter groups).
StreamingLLM(Xiao et al., 2023)	Streaming attention with attention sinks ; retain first $n_{\text{sink}} = 32$ tokens as sinks and maintain a rolling window of $W = 4096$ most recent tokens in KV.
SnapKV(Li et al., 2024b)	Prompt KV cache compression by selecting clustered important KV positions per head. Representative prompt KV cache size: $K_{\text{snap}} = 1024$; max pooling kernel size: 5.
SlimInfer(Long et al., 2025)	Dynamic token pruning for long-context inference with an asynchronous KV manager; representative target pruning rate: $\approx 50\%$ (prompt-side), enabling GPU-memory reduction with bounded quality loss.

Table 25: Aligned configuration knobs for long-context KV/memory experiments.

Bucket test (effect size). We bin $E_{\ell,i}$ into quintiles (Q1 = lowest energy, Q5 = highest energy) and report the skip rate per bin. This provides an effect-size view independent of correlation metrics. Pooled across tasks, the skip rate decreases monotonically with residual energy: Q1 (lowest) = 62.1%, Q2 = 48.7%, Q3 = 35.9%, Q4 = 25.4%, and Q5 (highest) = 17.8%. This pattern is consistent with the hypothesis that low residual-energy tokens are more likely to be skimmable and hence skipped.

Confounds and controls. Residual-energy and skipping can be jointly influenced by (i) layer depth and (ii) token position. We therefore fit a logistic regression with fixed effects:

$$\Pr(g_{\ell,i} = 1) = \sigma\left(a_0 + a_1 \cdot \text{zscore}(E_{\ell,i}) + u_{\ell} + v_i\right),$$

where u_{ℓ} is a per-layer intercept and v_i is a per-position intercept (binned into 16 buckets). With these controls, the standardized coefficient remains strongly negative, $a_1 = -0.92$ with 95% CI $[-0.96, -0.88]$ and $p < 10^{-12}$, indicating that higher residual energy is associated with a substantially lower probability of skipping even after accounting for layer and position effects.

G.3 Skipped-Token Property Tests (Frequency / Punctuation / Stopwords / Repetition)

Token grouping. We define a **high-skip token** as a generated token position whose average skip ratio across layers is ≥ 0.5 , and a **low-skip token** as one with average skip ratio ≤ 0.2 . We then compare token-category frequencies between the two groups. All tests are conducted on the pooled set across tasks with stratification by language (EN/ZH).

Categories. We consider four interpretable token properties:

- **High-frequency:** token in the top 5% of uni-gram frequency within the evaluation corpus.
- **Punctuation:** punctuation or structural separators (incl. newline markers).
- **Stopword/function:** stopwords/function words (language-specific list).
- **Repetition:** token is part of a repeated n -gram pattern (detected by local 4-gram repetition).

Contingency tests. We report group proportions and odds ratios (OR) with chi-square tests (FDR-corrected).

Method	8k input (8192)			16k input (16384)			Notes
	TTFT(s)↓	Decode (tok/s)↑	PeakMem (GB)↓	TTFT(s)↓	Decode (tok/s)↑	PeakMem (GB)↓	
Baseline (Full KV)	4.8	24	26.0	19.0	13	42.0	Full prompt KV; decode slows with context length.
FBS-Full (ours)	4.9	34	26.5	19.2	18	42.5	Main gain on decode compute (layer skipping); KV footprint largely unchanged.
H2O	4.8	31	18.0	19.0	17	25.0	KV eviction reduces memory and improves decode; TTFT unchanged (full prefill).
StreamingLLM	1.9	55	16.0	3.8	50	16.5	Windowed KV + sinks improves TTFT and stabilizes streaming beyond cache size.
SnapKV	5.1	65	14.0	19.8	46	18.0	Compression overhead slightly increases TTFT; decode speed improves with compressed prompt KV.
SlimInfer	2.4	40	15.0	7.6	32	19.0	Prompt-side pruning reduces TTFT and memory; decode improves moderately.
SnapKV	5.2	82	14.5	20.0	60	18.5	Orthogonal combination: KV compression + layer skipping yields best decode throughput.

Table 26: Long-context efficiency decomposition.

Task	Lang	#Examples	#Tok	$\rho(k, s)$	95% CI	$\tau(k, s)$	p_{FDR}
MMLU	EN	2000	256000	-0.47	[-0.49, -0.45]	-0.33	$< 10^{-12}$
BBH	EN	1500	192000	-0.44	[-0.47, -0.41]	-0.31	$< 10^{-12}$
GSM8K	EN	1500	192000	-0.52	[-0.55, -0.49]	-0.37	$< 10^{-12}$
HumanEval-X	EN	800	102400	-0.40	[-0.44, -0.36]	-0.28	$< 10^{-12}$
CMMLU	ZH	2000	256000	-0.50	[-0.52, -0.48]	-0.35	$< 10^{-12}$
C-Eval	ZH	1500	192000	-0.48	[-0.51, -0.46]	-0.34	$< 10^{-12}$
CMath	ZH	1200	153600	-0.55	[-0.58, -0.52]	-0.39	$< 10^{-12}$
Idiom QA/MC	ZH	1200	153600	-0.46	[-0.49, -0.42]	-0.32	$< 10^{-12}$
Pooled	–	11700	1497600	-0.49	[-0.50, -0.48]	-0.34	$< 10^{-12}$

Table 27: Rank correlation between lookahead $k(i)$ and surprisal $s(i)$ across tasks/languages. Negative correlations indicate larger lookahead on more predictable tokens.

Effect-size summary. Across both languages, high-skip tokens are significantly enriched in punctuation/structural markers, stopwords, and repeated patterns, consistent with the hypothesis that skipping preferentially targets low-semantic-load or highly predictable token regions.

Task	Lang	#Examples	#Tok	$\rho(k, H)$	95% CI	$\tau(k, H)$	p_{FDR}
MMLU	EN	2000	256000	-0.41	[-0.43, -0.39]	-0.29	$< 10^{-12}$
BBH	EN	1500	192000	-0.38	[-0.41, -0.35]	-0.27	$< 10^{-12}$
GSM8K	EN	1500	192000	-0.45	[-0.48, -0.42]	-0.32	$< 10^{-12}$
HumanEval-X	EN	800	102400	-0.35	[-0.39, -0.31]	-0.24	$< 10^{-12}$
CMMMLU	ZH	2000	256000	-0.43	[-0.45, -0.41]	-0.31	$< 10^{-12}$
C-Eval	ZH	1500	192000	-0.40	[-0.43, -0.38]	-0.29	$< 10^{-12}$
CMATH	ZH	1200	153600	-0.49	[-0.52, -0.46]	-0.35	$< 10^{-12}$
Idiom QA/MC	ZH	1200	153600	-0.39	[-0.43, -0.36]	-0.28	$< 10^{-12}$
Pooled	–	11700	1497600	-0.42	[-0.43, -0.41]	-0.30	$< 10^{-12}$

Table 28: Rank correlation between lookahead $k(i)$ and entropy $H(i)$ across tasks/languages.

Control	Pooled partial ρ	95% CI	p
Control position i only	-0.46	[-0.47, -0.45]	$< 10^{-12}$
Control task only	-0.45	[-0.46, -0.44]	$< 10^{-12}$
Control position + task	-0.44	[-0.45, -0.43]	$< 10^{-12}$

Table 29: Partial correlation robustness for $k(i)$ vs. surprisal (pooled across tasks).

Task	Lang	$\rho(g, E)$	95% CI	p_{FDR}
MMLU	EN	-0.31	[-0.34, -0.28]	$< 10^{-12}$
GSM8K	EN	-0.36	[-0.39, -0.33]	$< 10^{-12}$
HumanEval-X	EN	-0.28	[-0.33, -0.23]	$< 10^{-12}$
CMMMLU	ZH	-0.33	[-0.36, -0.30]	$< 10^{-12}$
C-Eval	ZH	-0.32	[-0.35, -0.29]	$< 10^{-12}$
CMATH	ZH	-0.38	[-0.41, -0.35]	$< 10^{-12}$
Pooled	–	-0.33	[-0.34, -0.32]	$< 10^{-12}$

Table 30: Correlation between skip indicator g and residual-energy proxy E (token-layer granularity). Negative values indicate higher skip probability for low-energy layers.

Property	Lang	High-skip (%)	Low-skip (%)	OR	p_{FDR}
High-frequency	EN	43.2	26.1	2.15	$< 10^{-12}$
Punctuation	EN	17.4	6.3	3.15	$< 10^{-12}$
Stopword/function	EN	29.8	14.9	2.42	$< 10^{-12}$
Repetition	EN	11.6	4.2	3.00	$< 10^{-12}$
High-frequency	ZH	46.8	28.4	2.22	$< 10^{-12}$
Punctuation	ZH	15.9	5.8	3.08	$< 10^{-12}$
Stopword/function	ZH	24.5	12.7	2.24	$< 10^{-12}$
Repetition	ZH	12.9	4.9	2.86	$< 10^{-12}$

Table 31: Token property enrichment for high-skip vs. low-skip tokens. OR = odds ratio of the property occurring in high-skip tokens relative to low-skip tokens (stratified by language).