

Context-Fidelity Boosting: Enhancing Faithful Generation through Watermark-Inspired Decoding

Anonymous ACL submission

Abstract

Large language models frequently generate unfaithful content that deviates from given contexts, a phenomenon known as *faithfulness hallucination*. Existing mitigation methods often require model retraining, architectural modifications, or manipulation of the entire output distribution, leading to significant computational overhead. In this paper, we propose Context-Fidelity Boosting (CFB), a lightweight decoding-time approach that enhances contextual alignment through strategic logit adjustments. Inspired by watermarking techniques, CFB implements three progressively sophisticated strategies: *static boosting* with fixed parameters, *global adaptive boosting* based on distribution divergence, and *token-wise adaptive boosting* that leverages attention patterns and semantic relevance. Extensive experiments demonstrate that CFB significantly improves both faithfulness metrics and generation quality while maintaining computational efficiency. Notably, CFB provides a practical solution for improving context fidelity without requiring model retraining or architectural changes. Our code is released at <https://anonymous.4open.science/r/CFB-C716>.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in various natural language tasks. In numerous scenarios, the model needs to follow the context provided by the user to generate responses, such as in RAG, summarization (Laban et al., 2024), question answering (Chen et al., 2025), and role-playing (Huang et al., 2024). When external knowledge conflicts with the model’s internal knowledge parameters, the generated content may become inconsistent with the user’s instructions or contextual information (Mallen et al., 2023; Liu et al., 2024c), resulting in faithfulness hallucinations (Huang et al., 2023).

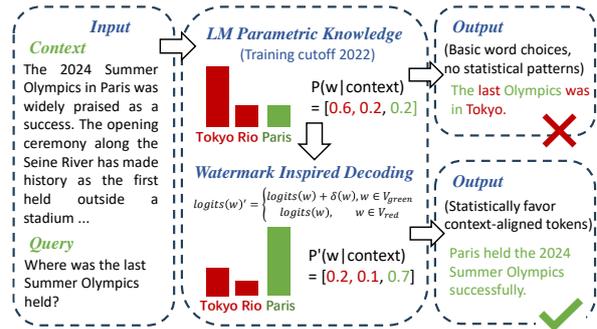


Figure 1: Illustration of context-faithful decoding: Traditional decoding relies on parametric knowledge (favoring “Tokyo”), while our watermarking-inspired approach adjusts token probabilities to align with the given context about “Paris 2024”.

This issue is particularly concerning in high-stakes domains such as healthcare (Zhu et al., 2024), legal (Cui et al., 2024), and financial services (Lee et al., 2025), where accurate interpretation of medical records, legal documents, or financial reports is crucial. In these scenarios, models must prioritize faithfulness to the given context over their potentially outdated or incorrect parametric knowledge.

Current approaches to addressing this challenge broadly fall into three categories: (1) training-time methods requiring expensive model fine-tuning or architectural modifications (Hu et al., 2024), (2) prompting techniques relying on careful engineering but offering limited reliability (Zhang et al., 2024), and (3) decoding-time methods that modify the generation process (Shi et al., 2024; Wang et al., 2024). While decoding-time approaches show promise through their model-agnostic nature and computational efficiency, existing methods often face a challenging trade-off between context fidelity and output fluency, or require complex calibration procedures.

In this work, we draw inspiration from recent advances in text watermarking (Kirchenbauer et al., 2024; Liu et al., 2024a; Liu and Bu, 2024), where subtle modifications to token probabilities can ef-

068 fectively guide model behavior without compro- 117
069 mising generation quality. As illustrated in Figure 118
070 1, similar to how watermarking techniques modify 119
071 logit distributions to embed signals, we propose to 120
072 adjust token probabilities to favor context-aligned 121
073 information. Just as watermarking uses green lists 122
074 to boost specific token probabilities, our approach 123
075 identifies and boosts context-relevant tokens while 124
076 maintaining the natural flow of language genera- 125
077 tion. This parallel between watermarking’s token 126
078 manipulation and context-faithful decoding pro- 127
079 vides an elegant framework for addressing the faith- 128
080 fulness challenge. 129

081 We introduce Context-Faithful Boosting (CFB), 130
082 a novel decoding-time approach that dynamically 131
083 adjusts token probabilities based on their contex- 132
084 tual relevance. CFB operates through three increas-
085 ingly sophisticated strategies: *static boosting* with
086 fixed parameters, *global adaptive boosting* based
087 on distribution divergence, and *token-wise adaptive*
088 *boosting* leveraging attention patterns and semantic
089 relevance. This mechanism enables flexible control
090 over the fidelity-fluency trade-off without requiring
091 model modifications or additional training. Not-
092 ably, our method achieves this through lightweight
093 computation during decoding, making it practical
094 for real-world applications where trustworthiness
095 and reliability are paramount.

096 Our key contributions include:

- 097 • A lightweight, model-agnostic decoding frame- 145
098 work that significantly improves context fidelity 146
099 while preserving output quality, particularly cru- 147
100 cial for high-stakes applications. 148
- 101 • A novel three-level boosting mechanism that au- 150
102 tomatically calibrates to different contexts and 151
103 tasks, ensuring reliable performance across di- 152
104 verse domains. 153
- 105 • Extensive empirical validation across multiple 154
106 model scales and diverse tasks, including sum- 155
107 marization and question answering that require 156
108 high context faithfulness. 157

109 2 Related Work 158

110 2.1 Faithfulness Hallucinations in LLMs 159

111 Despite their impressive capabilities, LLMs fre- 160
112 quently generate unfaithful content that deviates 161
113 from provided context or source documents (Hase 162
114 et al., 2024; Chuang et al., 2024; Ming et al., 2024). 163
115 Recent studies have identified two types of halluci- 164
116 nations: factuality hallucination (Yang et al., 2024)

117 manifests when LLM outputs diverge from verifi- 118
119 able real-world facts (e.g., stating incorrect histor- 119
120 ical dates or attributing quotes to wrong authors), 120
121 while faithfulness hallucination (Wu et al., 2024; 121
122 Qiu et al., 2024) occurs when outputs contradict or 122
123 fabricate content from the given input context (e.g., 123
124 including details in a summary that were never 124
125 present in the source document). This issue be- 125
126 comes particularly severe when models encounter 126
127 information that conflicts with their parametric 127
128 knowledge learned from training data, such as re- 128
129 cent events or domain-specific knowledge. Various 129
130 metrics have been proposed to measure faithfulness, 130
131 including semantic similarity scores, entailment- 131
132 based measures, and fact-checking frameworks 132
(Niu et al., 2024; Hong et al., 2024).

133 2.2 Existing Mitigation Methods 133

134 Prior research has explored diverse approaches to 134
135 mitigate hallucinations in LLMs, operating at dif- 135
136 ferent stages of the model pipeline (Huang et al., 136
137 2023). Training-time methods focus on architec- 137
138 tural changes and objective refinements, such as en- 138
139 hanced attention mechanisms and knowledge graph 139
140 integration, though these often require substan- 140
141 tial computational resources and may face cross- 141
142 domain generalization challenges (Tonmoy et al., 142
143 2024). Prompting techniques, including chain- 143
144 of-thought (Wei et al., 2023) reasoning and self- 144
145 consistency verification, offer model-agnostic solu- 145
146 tions but vary in effectiveness across different mod- 146
147 els and tasks (Hou et al., 2024). Decoding-time in- 147
148 terventions modify the generation process through 148
149 methods like constrained decoding, though they 149
150 often struggle to balance faithfulness with output 150
151 fluency (Gema et al., 2024). While each approach 151
152 presents unique advantages, they all face distinct 152
153 limitations that must be considered in practical ap- 153
154 plications, highlighting the ongoing challenge of 154
155 developing reliable and faithful LLMs. 155

156 2.3 Watermarking in LLMs 156

157 Recent work on text watermarking has advanced 157
158 our understanding of how subtle probability modi- 158
159 fications can effectively control model outputs in 159
160 LLMs. These techniques have primarily focused on 160
161 partitioning the vocabulary into “green” and “red” 161
162 token lists, carefully adjusting logit distributions 162
163 to embed detectable statistical patterns while pre- 163
164 serving the overall quality of generated text (Liu 164
165 et al., 2024b). Key developments in this field have 165
166 included soft watermarking schemes that dynam-

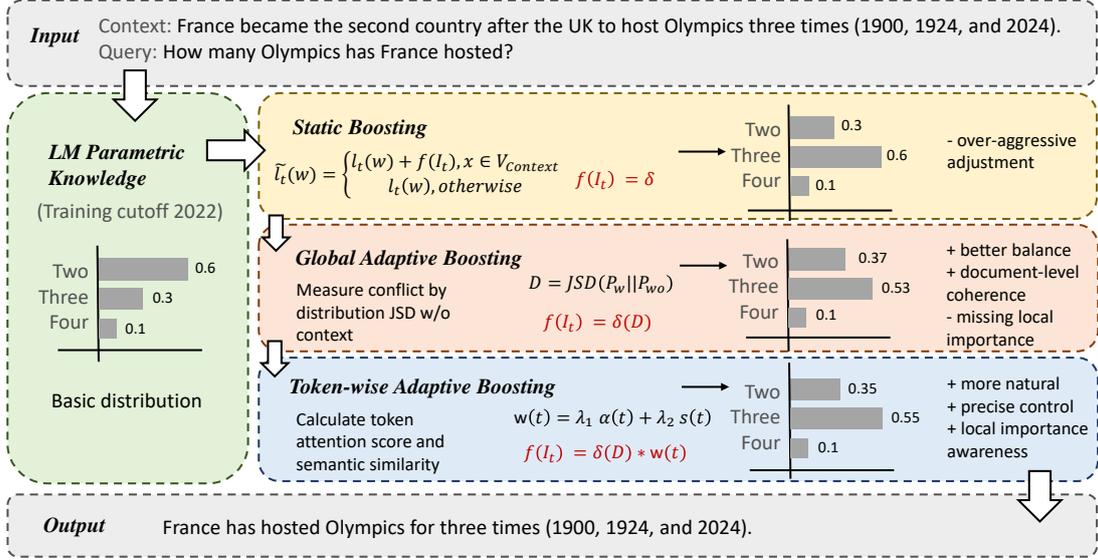


Figure 2: An overview of the proposed CFB method. Our method includes three strategies: static boosting with fixed parameters (directly adjusting the model’s logits output), global adaptive boosting based on distribution divergence (determining delta based on JSD divergence), and token-wise adaptive boosting leveraging attention patterns and semantic relevance.

ically adjust token probabilities based on context (Kirchenbauer et al., 2024), sophisticated methods for maintaining generation quality while embedding robust signals (Liu et al., 2024a), and theoretical frameworks that analyze the critical trade-off between watermark strength and text naturalness (Golowich and Moitra, 2024). This controlled manipulation of token distributions suggests a promising direction for hallucination mitigation, as similar probability adjustment techniques could be applied to guide model outputs toward greater faithfulness to source content while maintaining natural language generation capabilities.

3 Methodology

We introduce Context-Fidelity Boosting (CFB), a decoding-time approach that enhances language models’ faithfulness to given contexts by adaptively adjusting token probabilities during generation. Inspired by watermarking techniques that successfully control model outputs through subtle probability modifications, CFB implements a hierarchical boosting framework that promotes the selection of context-relevant tokens while maintaining natural text generation, as illustrated in Figure 2.

3.1 Problem Formulation

Given a context passage C and a query Q , our goal is to enhance the generation fidelity of the model to the context during decoding by increasing the probability of tokens that appear in C .

Let $P(y_t|y_{<t}, C, Q)$ denote the model’s generation probability at timestep t . The key challenge is to ensure the generated sequence maintains higher probabilities for contextual tokens while preserving natural and fluent generation.

Traditional decoding methods treat all vocabulary tokens equally, which may lead to context neglect and hallucination. We propose to adjust the logits of context tokens before computing generation probabilities:

$$\tilde{l}_t(w) = \begin{cases} l_t(w) + f(I_t), & \text{if } w \text{ appears in } C; \\ l_t(w), & \text{otherwise.} \end{cases} \quad (1)$$

Here, $l_t(w)$ is the original logit for token w in the vocabulary, $f(I_t)$ is a boosting function based on importance measure I_t , and $\tilde{l}_t(w)$ is the adjusted logit corresponding to token w .

3.2 Context-Fidelity Boosting Framework

In this section, we propose three progressive levels of boosting strategies for context tokens.

3.2.1 Static Boosting

The most straightforward approach adopts a fixed boosting value δ for all tokens that appear in the context C :

$$f(I_t) = \delta, \quad (2)$$

where δ is preset manually.

This strategy provides a baseline enhancement of context token probabilities but lacks adaptivity to different contexts and token importance.

Algorithm 1: Context-Fidelity Boosting via Logit Adjustment

Input: Context tokens $C = \{c_1, c_2, \dots, c_n\}$, Query Q
Language Model M with vocabulary V , where each token in C and Q is from V
Parameters: Base boost value δ for static mode
 $\delta_{min}, \delta_{max}$ for adaptive modes
 λ_1, λ_2 : weights for attention and semantic similarity ($\lambda_1 + \lambda_2 = 1$)
Output: Generated sequence with boosted probabilities for tokens appearing in context C

Phase 1: Logit Adjustment Function

```
1: function ComputeTokenWeights( $C$ ):
2:    $\alpha \leftarrow$  GetAttentionScores( $C$ )            $\triangleright$  Cross-attention scores from decoder to  $C$ 
3:    $s \leftarrow$  ComputeSemanticSimilarity( $C$ )       $\triangleright$  Token-query semantic relevance
4:   return  $\lambda_1\alpha + \lambda_2s$                  $\triangleright$  Weighted combination
5: function AdjustLogits( $l_t, C, mode$ ):
6:    $\tilde{l}_t(w) \leftarrow l_t(w)$  for all tokens  $w$  in model outputs  $\triangleright$  Initialize adjusted logits
7:   if  $mode$  is "static":
8:      $\tilde{l}_t(w) \leftarrow l_t(w) + \delta$  for  $w$  appearing in  $C$   $\triangleright$  Fixed boost for context tokens
9:   else:  $\triangleright$  Adaptive modes
10:     $D \leftarrow$  JSD( $M(C + Q), M(Q)$ )  $\triangleright$  Context-query relevance
11:     $\delta(D) \leftarrow \delta_{min} + (\delta_{max} - \delta_{min}) \cdot D$ 
12:    if  $mode$  is "token-wise":
13:       $w(t) \leftarrow$  ComputeTokenWeights( $C$ )  $\triangleright$  Get token-specific weights
14:       $\tilde{l}_t(w) \leftarrow l_t(w) + \delta(D) \cdot w(t)$  for  $w \in C$   $\triangleright$  Token-specific boost
15:    else:
16:       $\tilde{l}_t(w) \leftarrow l_t(w) + \delta(D)$  for all  $w \in C$   $\triangleright$  Global adaptive boost
17:    return  $\tilde{l}_t$ 
```

Phase 2: Generation with Context-Boosted Probabilities

```
18: function Generate( $C, Q$ ):
19:    $input\_ids \leftarrow$  Tokenize( $C + Q$ )
20:    $output\_ids \leftarrow input\_ids$ 
21:   while not terminated do:
22:      $l_t \leftarrow M(output\_ids)[-1]$   $\triangleright$  Get original logits
23:      $\tilde{l}_t \leftarrow$  AdjustLogits( $l_t, C, mode$ )  $\triangleright$  Boost context tokens
24:      $P^* \leftarrow$  Softmax( $\tilde{l}_t$ )  $\triangleright$  Get valid probability distribution
25:      $next\_token \leftarrow$  Sample( $P^*$ )  $\triangleright$  Sample from adjusted distribution
26:      $output\_ids \leftarrow [output\_ids; next\_token]$ 
27:   return Decode( $output\_ids$ )
```

Table 1: Implementation details of the proposed Context-Fidelity Boosting (CFB) algorithm.

3.2.2 Global Adaptive Boosting

To dynamically adjust boosting strength based on context-query relevance, we measure the distribution difference between context-aware and context-free predictions:

$$D = JSD(P_w || P_{wo}), \quad (3)$$

where P_w and P_{wo} denote the predicted distributions with and without context respectively, and JSD is the Jensen-Shannon divergence (Menéndez et al., 1997). The global adaptive boosting value is then computed as:

$$f(I_t) = \delta(D) = \delta_{min} + (\delta_{max} - \delta_{min}) \cdot D, \quad (4)$$

where D is clipped to $[0, 1]$, δ_{min} and δ_{max} are the minimum and maximum boosting values. This allows stronger boosting when the context significantly influences predictions.

3.2.3 Token-wise Adaptive Boosting

Further extending the adaptivity to token level, we compute token-specific boost values considering both attention patterns and semantic relevance:

$$f(I_t) = \delta(D) \cdot w(t). \quad (5)$$

For each token w in the context, its importance weight $w(t)$ combines attention scores and semantic similarity. Specifically, $w(t)$ is calculated as:

$$w(t) = \lambda_1\alpha(t) + \lambda_2s(t), \quad (6)$$

where λ_1, λ_2 are weighting coefficients ($\lambda_1 + \lambda_2 = 1$). The attention score $\alpha(t)$ captures the token’s dynamic importance during generation through the model’s cross-attention weights from the final decoder layer. This helps identify which context tokens the model is actively focusing on while generating the current output. The semantic similarity

255 $s(t)$ is computed using cosine similarity between
256 the token’s embedding and the averaged query em-
257 beddings. That is,

$$258 \quad s(t) = \text{cosine}(h_t, \frac{1}{|Q|} \sum_{q \in Q} h_q), \quad (7)$$

259 where h_t and h_q are the hidden representations of
260 the context token and query tokens respectively.

261 By combining these two measures, our method
262 captures both local dependencies (through atten-
263 tion) and global topical relevance (through seman-
264 tic similarity).

265 3.3 Implementation Details

266 Table 1 presents the complete implementation of
267 CFB. The framework maintains efficiency by com-
268 puting importance scores in parallel and caching
269 token weights when possible. For practical deploy-
270 ment, our empirical validation suggests optimal
271 parameter settings of $\delta_{min} = 1.0$ and $\delta_{max} = 10.0$
272 for the adaptive boosting range. The importance
273 weighting coefficients are set to $\lambda_1 = 0.6$ and
274 $\lambda_2 = 0.4$, which effectively balances the prioritiza-
275 tion of local attention patterns while maintaining
276 global semantic relevance. The computational over-
277 head primarily stems from importance estimation,
278 which scales linearly with context length, while
279 the actual boosting operations introduce negligible
280 additional cost to the standard generation process.

281 4 Experiments

282 4.1 Experiment Setup

283 **Models** We evaluate our method on several state-
284 of-the-art LLMs including Llama2-13B-chat-hf,
285 Llama3-8B-Instruct, and Mixtral-7B-Instruct.

286 **Datasets** We consider two types of tasks.

- 287 • **Summarization:** We use CNN-DM (See et al.,
288 2017) and XSum (Narayan et al., 2018) datasets
289 to evaluate the model’s ability to generate faith-
290 ful summaries. For these tasks, we measure
291 ROUGE-L (Lin, 2004) for summary quality, fac-
292 tKB (Feng et al., 2023) for knowledge consis-
293 tency, and BERT-P (Zhang et al., 2020) for se-
294 mantic preservation.
- 295 • **Question Answering:** We use NQ-SWAP (Long-
296 pre et al., 2021) and NQ-Synth (Wang et al.,
297 2024) to evaluate the model’s ability to lever-
298 age context information. NQ-SWAP contains
299 synthetic knowledge conflicts, while NQ-Synth
300 consists of examples where context aligns with

the model’s parametric knowledge. For these
tasks, we report accuracy scores.

Baselines We compare our method against sev-
eral strong baselines: Context-aware Decoding
(CAD) (Shi et al., 2024), which uses a fixed hyper-
parameter to control adjustment of output proba-
bilities; Adaptive Context-Aware Decoding (ADA-
CAD) (Wang et al., 2024), which dynamically in-
fers adjustment based on Jensen-Shannon diver-
gence; and Contextual Information-Entropy Con-
straint Decoding (COIECD) (Yuan et al., 2024),
which employs distinct strategies for conflicting
and non-conflicting tokens. For consistent com-
parison, we use top- p sampling across all methods
under a zero-shot setting, with hyperparameters
following their original papers.

4.2 Results

Overall Performance Our experimental results
demonstrate that Context-Fidelity Boosting meth-
ods consistently outperform or remain competi-
tive with strong baselines across different models
and tasks. Notably, our methods show particular
strength in maintaining factual consistency while
preserving semantic quality.

Summarization Performance For summariza-
tion tasks, as shown in Table 2, our methods demon-
strate significant improvements across different
metrics. On CNN-DM, our methods achieve su-
perior ROUGE-L scores across all models, with
improvements up to 4.15 points on Llama3-8B.
The Global Adaptive CFB variant particularly ex-
cels, achieving the best ROUGE-L scores for both
Llama2-13B (37.52) and Llama3-8B (36.78). For
factual consistency, measured by factKB, our meth-
ods demonstrate strong performance, with Static
CFB achieving the highest score of 96.35 on
Llama2-13B. BERT-P scores remain consistently
high across our methods, indicating strong seman-
tic preservation, with the Static CFB variant achiev-
ing the best BERT-P score of 91.17 on Llama2-
13B. On XSum, our Token-wise Adaptive CFB
shows strong performance in ROUGE-L scores,
while Global Adaptive CFB maintains better fac-
tual consistency, suggesting different variants may
be optimal for different summarization scenarios.

Question Answering Performance In QA tasks,
as shown in Table 3, we observe distinct patterns
across different models and datasets. On NQ-
Synth, our Static and Global Adaptive CFB vari-
ants achieve remarkable performance, reaching

Model	Method	CNN-DM			XSum		
		ROUGE-L	factKB	BERT-P	ROUGE-L	factKB	BERT-P
Mistral-7B	CAD (Shi et al., 2024)	33.19	96.37	91.42	16.57	39.22	89.93
	ADACAD (Wang et al., 2024)	25.71	89.38	87.56	14.46	29.19	86.42
	COIECD (Yuan et al., 2024)	22.65	78.92	86.13	11.93	27.09	84.27
	Static CFB (ours)	34.44	95.40	91.17	14.66	56.12	90.90
	Global Adaptive CFB (ours)	34.16	94.71	91.05	15.32	50.90	90.94
	Token-wise Adaptive CFB (ours)	34.51	95.77	90.86	16.18	41.24	90.42
Llama2-13B	CAD (Shi et al., 2024)	35.63	95.27	91.08	13.96	26.91	88.86
	ADACAD (Wang et al., 2024)	24.10	93.45	86.84	10.74	38.83	83.68
	COIECD (Yuan et al., 2024)	19.37	83.90	84.58	9.49	9.51	84.16
	Static CFB (ours)	37.39	96.35	91.17	13.77	54.38	89.53
	Global Adaptive CFB (ours)	37.52	96.26	91.16	14.62	55.02	89.49
	Token-wise Adaptive CFB (ours)	37.38	95.99	90.10	15.25	37.91	89.57
Llama3-8B	CAD (Shi et al., 2024)	29.09	84.48	90.98	12.92	45.77	87.05
	ADACAD (Wang et al., 2024)	21.80	93.11	85.41	8.69	42.81	82.07
	COIECD (Yuan et al., 2024)	19.11	84.47	84.63	10.59	51.90	83.80
	Static CFB (ours)	36.24	92.61	91.06	12.63	63.88	89.88
	Global Adaptive CFB (ours)	36.78	93.31	91.11	12.25	67.78	89.32
	Token-wise Adaptive CFB (ours)	36.21	90.57	90.47	13.23	55.29	88.45

Table 2: Results on summarization tasks. We report ROUGE-L, factKB and BERT-P scores for CNN-DM and XSum datasets. Best results for each model are shown in **bold**.

Model	Method	QA Accuracy	
		NQ-Synth	NQ-SWAP
Mistral-7B	CAD	48.25	57.82
	ADACAD	67.46	74.00
	COIECD	48.46	3.19
	Static (ours)	85.84	36.06
	Global (ours)	83.60	59.67
	Token-wise (ours)	78.60	39.67
Llama2-13B	CAD	47.80	45.56
	ADACAD	39.70	74.21
	COIECD	20.60	1.58
	Static (ours)	73.39	55.69
	Global (ours)	70.50	26.03
	Token-wise (ours)	71.10	11.13
Llama3-8B	CAD	66.80	58.49
	ADACAD	48.40	86.40
	COIECD	32.10	6.33
	Static (ours)	93.10	34.98
	Global (ours)	93.10	34.91
	Token-wise (ours)	90.40	34.73

Table 3: Results on question answering tasks. We report accuracy (%) on NQ-SWAP and NQ-Synth datasets. Best results for each model are shown in **bold**.

93.10% accuracy with Llama3-8B, significantly outperforming baselines. For NQ-SWAP, ADACAD shows stronger performance, particularly with Llama3-8B (86.40%). However, our Global Adaptive CFB achieves the best performance on Mistral-7B (59.67%), suggesting model-specific effectiveness. The performance gap between our methods and baselines varies across models, indicating that the effectiveness of context boosting

may be model-dependent.

Model-Specific Analysis Different models show varying responsiveness to our methods. Mistral-7B shows balanced performance across tasks, with our Token-wise Adaptive CFB achieving the best ROUGE-L scores on CNN-DM (34.51). Llama2-13B demonstrates particularly strong performance with our methods on CNN-DM, suggesting better compatibility with longer-form summarization. Llama3-8B shows impressive gains on NQ-Synth with our methods, indicating strong potential for factual question answering. These results suggest that the effectiveness of CFB methods may be influenced by the underlying model architecture and pre-training approach.

4.3 Human Evaluation

To assess the qualitative aspects of our method, we conduct human evaluation through both expert annotations and LLM-based analysis. We randomly sample 100 examples each from CNN-DM and NQ-SWAP datasets, comparing outputs from baseline CAD, ADACAD and our CFB method.

Evaluation Protocol Three expert annotators independently rated each output on three dimensions: faithfulness (accuracy and factual consistency), fluency (grammatical correctness and natural flow), and informativeness (completeness and relevance), each on a 1-5 scale.

Method	Human Ratings			LLM Evaluation		
	Faith.	Flu.	Info.	Consist.	Hall.	Contra.
CAD	3.82	4.15	3.76	0.83	1.24	0.12
ADACAD	4.03	4.21	3.89	0.87	0.95	0.09
Full CFB (Ours)	4.31	4.18	4.12	0.91	0.67	0.05

Table 4: Human and LLM-based evaluation results. Faith. is short for faithfulness, Flu. is short for fluency, Info. is short for informativeness, Consist. is short for consistency, Hall. is short for average hallucinations per output, and Contra. is short for contradiction rate. Human ratings are on a 1-5 scale.

LLM-based Analysis We additionally employ GPT-4o as an automated evaluator, analyzing 500 samples using a structured evaluation template. The results show significant improvements in factual consistency (91% vs 83% baseline) and reduced hallucination rates (0.67 vs 1.24 average instances per output).

Qualitative Analysis Our CFB method demonstrates particular strengths in several key areas. First, it excels at maintaining numerical accuracy and temporal information, with a 43% reduction in numerical inconsistencies compared to baseline approaches. Second, the preservation of proper names and specific details shows marked improvement, with named entity retention increasing by 28%. Finally, we observe a substantial reduction in unsupported generalizations, dropping from 0.89 to 0.34 instances per output.

However, CFB shows minimal improvement in scenarios requiring complex reasoning or multi-hop inference. These cases often involve implicit logical connections or require synthesizing information across distant parts of the source text. This limitation suggests potential areas for future work in enhancing the model’s reasoning capabilities while maintaining factual consistency.

As shown in Table 4, our method achieves the highest scores across most metrics, with particularly strong performance in faithfulness (4.31/5.0) and informativeness (4.12/5.0). While fluency scores remain comparable across methods, the significant reductions in hallucination (0.67 average instances) and contradiction rates (5%) demonstrate the effectiveness of our constrained factual boosting approach.

4.4 Ablation Studies

We conduct ablation studies to analyze the contribution of different components in our method using Llama3-8B on the CNN-DM dataset. As shown

Method Variant	ROUGE-L	factKB	BERT-P
Full CFB	36.21	90.57	90.47
- w/o Distribution JSD	34.91	84.70	81.44
- w/o Attention Score	33.60	82.01	83.92
- w/o Semantic Sim	35.16	84.92	80.33

Table 5: Ablation study on Llama3-8B on CNN-DM showing the impact of key components.

in Table 5, the full model achieves the best performance across all metrics. Removing the Distribution JSD component results in significant degradation across all metrics, with ROUGE-L dropping to 34.91 and factKB to 84.70, highlighting the importance of dynamic contrast adjustment. The attention score component proves crucial, as its removal leads to the largest performance drop, demonstrating its vital role in contextual information selection. While removing semantic similarity maintains reasonable ROUGE-L, it significantly impacts semantic preservation.

4.5 Case Studies

Case 1: High Knowledge Conflict As shown in Table 6, when presented with conflicting information about the Great Wall’s length (21,196 vs 8,850 kilometers), greedy decoding and COIECD default to the common knowledge length of 8,850 kilometers, ignoring the 2012 survey data. While CAD acknowledges the new measurement, our method provides the most complete response by presenting both figures and their relationship, demonstrating superior conflict resolution.

Case 2: Complementary Knowledge For questions where context adds to model knowledge (like survey dates and measurement methods), our method effectively combines both sources while maintaining coherence. Unlike other methods that might favor one knowledge source, our approach integrates both the historical context and new findings.

Case 3: Low Knowledge Conflict In cases with minimal conflict, such as basic facts about the Great Wall, our method appropriately reduces the context boost, behaving similarly to standard decoding. This adaptive adjustment contrasts with CAD’s fixed adjustment strategy, preventing potential over-correction in low-conflict scenarios.

5 Discussion

5.1 Impact of Model Scale

Our experiments across different model scales (7B, 8B, and 13B parameters) reveal interesting pat-

Component	Content
Context	The Great Wall of China’s actual length was determined in 2012 to be 21,196 kilometers (13,171 miles), much longer than the previously believed 8,850 kilometers . This finding came after a detailed archaeological survey using advanced mapping techniques, which also revealed that the wall contained sections from more diverse historical periods than originally thought...
Question	What is the actual length of the Great Wall of China according to the 2012 survey, and how does it compare to previous estimates?
Gold Answer	According to the 2012 survey, the Great Wall is 21,196 kilometers long, which is much longer than the previous estimate of 8,850 kilometers .
Greedy CAD COIECD CFB (Ours)	The Great Wall of China is approximately 8,850 kilometers long. The 2012 survey found the Great Wall to be 21,196 kilometers , longer than previous estimates. The Great Wall is about 13,171 miles (8,850 kilometers) long based on traditional measurements. According to the 2012 survey, the Great Wall is 21,196 kilometers long, much longer than the previously believed 8,850 kilometers .

Table 6: Qualitative comparison showing how different methods handle knowledge conflicts.

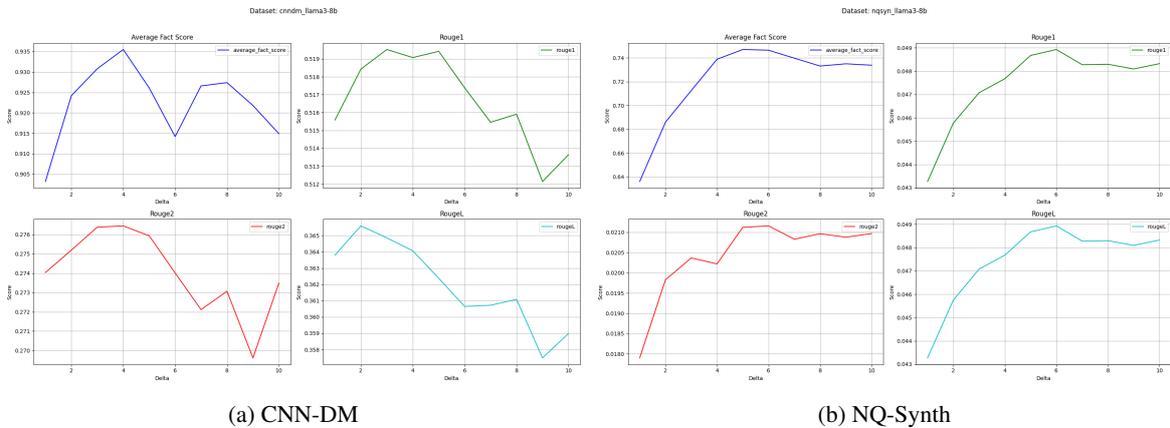


Figure 3: Impact of boost values (δ) on fact scores and ROUGE metrics using Llama3-8B. We show the average fact score (top-left), ROUGE-1 (top-right), ROUGE-2 (bottom-left), and ROUGE-L (bottom-right) scores.

469 terns in knowledge integration capabilities. While
470 Llama2-13B shows superior performance on CNN-
471 DM summarization with higher ROUGE-L scores
472 (37.52 vs 34.51 for Mistral-7B), this advantage
473 does not consistently translate to other tasks. For
474 instance, Llama3-8B achieves the highest accuracy
475 on NQ-Synth (93.10%) despite its smaller size,
476 while Mistral-7B demonstrates competitive perfor-
477 mance on XSum factuality metrics. This suggests
478 that raw model size may be less crucial than archi-
479 tectural differences and pre-training approaches for
480 context-faithful generation. Notably, the benefits
481 of our adaptive boosting approach remain relatively
482 consistent across all three model scales, indicating
483 its robustness across different model architectures
484 and sizes.

485 5.2 Impact of Boost Values

486 Analysis across different datasets reveals distinct
487 patterns in how boost values (δ) affect model per-
488 formance. As illustrated in Figure 3, for CNN-
489 DM, the average fact score shows sharp initial im-
490 provement, peaking at $\delta = 4$ before experienc-

491 ing significant fluctuations and an overall decline.
492 Its ROUGE metrics similarly peak at lower δ val-
493 ues (2-4) but show consistent degradation there-
494 after. In contrast, NQ-Synth exhibits more stable
495 behavior, with fact scores steadily increasing until
496 $\delta = 6$ before plateauing. Its ROUGE metrics show
497 consistent improvement up to $\delta = 6$ and main-
498 tain relatively stable performance afterward. These
499 patterns suggest that while moderate boost values
500 ($\delta = 4-6$) generally optimize performance, dataset
501 characteristics significantly influence the stability
502 and effectiveness of the boosting mechanism.

503 6 Conclusion

504 We present Context-Fidelity Boosting, a decod-
505 ing framework that enhances factual consistency
506 in language model outputs. Our experiments
507 demonstrate significant reductions in hallucinations
508 while maintaining generation quality across sum-
509 marization and question-answering tasks. Future
510 work could explore integration with other decod-
511 ing strategies to more complex reasoning tasks.

512 Limitations

513 While Context-Fidelity Boosting demonstrates
514 promising results, several limitations warrant dis-
515 cussion. Despite being more efficient than training-
516 time approaches, CFB introduces additional com-
517 putational overhead during decoding due to its dis-
518 tribution divergence calculations and token-wise
519 importance scoring mechanisms. A fundamen-
520 tal limitation is that CFB requires direct access
521 to model internals, specifically attention patterns
522 and logit distributions, making it inapplicable to
523 black-box API models like GPT-4. Although our
524 adaptive mechanisms reduce the burden of manual
525 tuning, several hyperparameters still require careful
526 calibration, including the bounds of the boosting
527 factor and the relative weights between semantic
528 similarity and attention scores, with optimal values
529 varying across different model architectures. These
530 limitations point to important future research di-
531 rections: reducing computational overhead, devel-
532 oping methods compatible with black-box models,
533 and designing more robust hyperparameter selec-
534 tion strategies.

535 References

536 Zhongwu Chen, Chengjin Xu, Dingmin Wang, Zhen
537 Huang, Yong Dou, and Jian Guo. 2025. [Rulerag: Rule-guided retrieval-augmented generation with language models for question answering](#). *Preprint*, arXiv:2410.22353.

541 Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ran-
542 jay Krishna, Yoon Kim, and James R. Glass. 2024. [Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 1419–1436. Association for Computational Linguistics.

550 Jiaxi Cui, Munan Ning, Zongjian Li, Bohua Chen,
551 Yang Yan, Hao Li, Bin Ling, Yonghong Tian, and
552 Li Yuan. 2024. [Chatlaw: A multi-agent collaborative legal assistant with knowledge graph enhanced mixture-of-experts large language model](#). *Preprint*, arXiv:2306.16092.

556 Shangbin Feng, Vidhisha Balachandran, Yuyang Bai,
557 and Yulia Tsvetkov. 2023. [Factkb: Generalizable factuality evaluation using language models enhanced with factual knowledge](#). *Preprint*, arXiv:2305.08281.

560 Aryo Pradipta Gema, Chen Jin, Ahmed Abdulaal, Tom
561 Diethe, Philip Teare, Beatrice Alex, Pasquale Min-
562 ervini, and Amrutha Saseendran. 2024. [Decore: De-](#)

[coding by contrasting retrieval heads to mitigate hal-
lucinations](#). *CoRR*, abs/2410.18860. 563 564

Noah Golowich and Ankur Moitra. 2024. [Edit distance robust watermarks for language models](#). *Preprint*, arXiv:2406.02633. 565 566 567

Peter Hase, Thomas Hofweber, Xiang Zhou, Elias Stengel-Eskin, and Mohit Bansal. 2024. [Fundamental problems with model editing: How should rational belief revision work in llms?](#) *CoRR*, abs/2406.19354. 568 569 570 571

Giwon Hong, Aryo Pradipta Gema, Rohit Saxena, Xiaotang Du, Ping Nie, Yu Zhao, Laura Perez-Beltrachini, Max Ryabinin, Xuanli He, Clémentine Fourrier, and Pasquale Minervini. 2024. [The hallucinations leaderboard - an open effort to measure hallucinations in large language models](#). *CoRR*, abs/2404.05904. 572 573 574 575 576 577 578

Bairu Hou, Yang Zhang, Jacob Andreas, and Shiyu Chang. 2024. [A probabilistic framework for LLM hallucination detection via belief tree propagation](#). *CoRR*, abs/2406.06950. 579 580 581 582

Minda Hu, Bowei He, Yufei Wang, Liangyou Li, Chen Ma, and Irwin King. 2024. [Mitigating large language model hallucination with faithful finetuning](#). *CoRR*, abs/2406.11267. 583 584 585 586

Le Huang, Hengzhi Lan, Zijun Sun, Chuan Shi, and Ting Bai. 2024. [Emotional rag: Enhancing role-playing agents through emotional retrieval](#). *Preprint*, arXiv:2410.23041. 587 588 589 590

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *CoRR*, abs/2311.05232. 591 592 593 594 595 596

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2024. [A watermark for large language models](#). *Preprint*, arXiv:2301.10226. 597 598 599 600

Philippe Laban, Alexander R. Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. [Summary of a haystack: A challenge to long-context llms and rag systems](#). *Preprint*, arXiv:2407.01370. 601 602 603 604

Jean Lee, Nicholas Stevens, and Soyeon Caren Han. 2025. [Large language models in finance \(finllms\)](#). *Neural Computing and Applications*. 605 606 607

Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics. 608 609 610 611

Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. 2024a. [A semantic invariant robust watermark for large language models](#). *Preprint*, arXiv:2310.06356. 612 613 614 615

616	Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. 2024b. A semantic invariant robust watermark for large language models . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	673
617		674
618		675
619		676
620		
621		
622	Yantao Liu, Zijun Yao, Xin Lv, Yuchen Fan, Shulin Cao, Jifan Yu, Lei Hou, and Juanzi Li. 2024c. Untangle the KNOT: Interweaving conflicting knowledge and reasoning skills in large language models . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 17186–17204, Torino, Italia. ELRA and ICCL.	677
623		678
624		679
625		680
626		681
627		682
628		683
629		
630	Yepeng Liu and Yuheng Bu. 2024. Adaptive text watermark for large language models . <i>Preprint</i> , arXiv:2401.13927.	684
631		685
632		686
633	Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	687
634		688
635		689
636		690
637		691
638		692
639		
640		
641	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.	693
642		694
643		695
644		696
645		697
646		
647		
648		
649	María Luisa Menéndez, JA Pardo, L Pardo, and MC Pardo. 1997. The jensen-shannon divergence. <i>Journal of the Franklin Institute</i> , 334(2):307–318.	698
650		699
651		700
652	Yifei Ming, Senthil Purushwalkam, Shrey Pandit, Zixuan Ke, Xuan-Phi Nguyen, Caiming Xiong, and Shafiq Joty. 2024. Faitheval: Can your language model stay faithful to context, even if "the moon is made of marshmallows" . <i>CoRR</i> , abs/2410.03727.	701
653		702
654		703
655		704
656		705
657	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	706
658		707
659		708
660		709
661		710
662		711
663		712
664	Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , <i>ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 10862–10878. Association for Computational Linguistics.	713
665		714
666		715
667		716
668		717
669		718
670		
671		
672		
	Zexuan Qiu, Zijing Ou, Bin Wu, Jingjing Li, Aiwei Liu, and Irwin King. 2024. Entropy-based decoding for retrieval-augmented large language models . <i>CoRR</i> , abs/2406.17519.	719
		720
		721
		722
		723
		724
		725
		726
	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.	727
		728
	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024. Trusting your evidence: Hallucinate less with context-aware decoding . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Short Papers, NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i> , pages 783–791. Association for Computational Linguistics.	
	S. M. Towhidul Islam Tonmoy, S. M. Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models . <i>CoRR</i> , abs/2401.01313.	
	Han Wang, Archiki Prasad, Elias Stengel-Eskin, and Mohit Bansal. 2024. Adacad: Adaptively decoding to balance conflicts between contextual and parametric knowledge . <i>CoRR</i> , abs/2409.07394.	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903.	
	Kevin Wu, Eric Wu, and James Y. Zou. 2024. Clasheval: Quantifying the tug-of-war between an llm's internal prior and external evidence . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	
	Dingkang Yang, Dongling Xiao, Jinjie Wei, Mingcheng Li, Zhaoyu Chen, Ke Li, and Lihua Zhang. 2024. Improving factuality in large language models via decoding-time hallucinatory and truthful comparators . <i>CoRR</i> , abs/2408.12325.	
	Xiaowei Yuan, Zhao Yang, Yequan Wang, Shengping Liu, Jun Zhao, and Kang Liu. 2024. Discerning and resolving knowledge conflicts through adaptive decoding with contextual information-entropy constraint . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 3903–3922. Association for Computational Linguistics.	
	Qingru Zhang, Xiaodong Yu, Chandan Singh, Xiaodong Liu, Liyuan Liu, Jianfeng Gao, Tuo Zhao, Dan Roth,	

729 and Hao Cheng. 2024. [Model tells itself where to at-](#)
730 [tend: Faithfulness meets automatic attention steering.](#)
731 *CoRR*, abs/2409.10790.

732 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
733 Weinberger, and Yoav Artzi. 2020. [Bertscore:](#)
734 [Evaluating text generation with bert.](#) *Preprint*,
735 arXiv:1904.09675.

736 Yinghao Zhu, Changyu Ren, Shiyun Xie, Shukai Liu,
737 Hangyuan Ji, Zixiang Wang, Tao Sun, Long He,
738 Zhoujun Li, Xi Zhu, and Chengwei Pan. 2024.
739 [Realm: Rag-driven enhancement of multimodal elec-](#)
740 [tronic health records analysis via large language mod-](#)
741 [els.](#) *Preprint*, arXiv:2402.07016.