

# LFD: Layer Fused Decoding to Exploit External Knowledge in Retrieval-Augmented Generation

Anonymous ACL submission

## Abstract

Retrieval-augmented generation (RAG) incorporates external knowledge into large language models (LLMs), improving their adaptability to downstream tasks and enabling information updates. Surprisingly, recent empirical evidence demonstrates that injecting noise into retrieved relevant documents paradoxically facilitates exploitation of external knowledge and improves generation quality. Although counterintuitive and challenging to apply in practice, this phenomenon enables granular control and rigorous analysis of how LLMs integrate external knowledge. Therefore, in this paper, we intervene on noise injection and establish a layer-specific functional demarcation within the LLM: shallow layers specialize in local context modeling, intermediate layers focus on integrating long-range external factual knowledge, and deeper layers primarily rely on parametric internal knowledge. Building on this insight, we propose Layer Fused Decoding (LFD), a simple decoding strategy that directly combines representations from an intermediate layer with final-layer decoding outputs to fully exploit the external factual knowledge. To identify the optimal intermediate layer, we introduce an internal knowledge score (IKS) criterion that selects the layer with the lowest IKS value in the latter half of layers. Experimental results across multiple benchmarks demonstrate that LFD helps RAG systems more effectively surface retrieved context knowledge with minimal cost. The code for replication is available at <https://anonymous.4open.science/r/LFD/>.

## 1 Introduction

Retrieval-Augmented Generation (RAG) empowers large language models (LLMs) by dynamically integrating external knowledge during inference, enabling precise adaptation to knowledge-intensive tasks and rapidly evolving domains (Borgeaud et al., 2022; Guu et al., 2020; Lewis et al., 2020a). As a cornerstone of context-aware generation, RAG

has been widely deployed in real-world applications, including recommendation systems (Fan et al., 2024; Di Palma, 2023; Ji et al., 2024) and search engines (Salemi and Zamani, 2024; Xiong et al., 2024). The broad applicability has spurred extensive optimization efforts on dynamic knowledge integration, including reranking strategies (Yu et al., 2024; Dong et al., 2024) to prioritize relevance, adaptive retrieval mechanisms (Asai et al., 2024; Jiang et al., 2023c) to minimize redundancy, and graph-based architectures (He et al., 2024; Edge et al., 2024; Jiang et al., 2023b) to model inter-document semantic relationships.

Despite these advances, LLMs might underutilize accurate external contexts, disproportionately favoring internal parametric knowledge during generation (Sun et al., 2025b; Niu et al., 2024). This overreliance risks propagating outdated information or hallucinations, undermining the trustworthiness of RAG systems. Surprisingly, recent studies reveal a paradoxical phenomenon: **injecting noise—random documents or tokens—to retrieved contexts that already contain answer-relevant snippets can improve the generation accuracy** (Cuconasu et al., 2024; Sun et al., 2025a). While this noise-injection approach is simple and effective, its underlying influence on LLM remains unclear. Furthermore, long contexts containing noise documents create computational overhead. Therefore, it is important to design more principled strategies that can achieve similar benefits without incurring excessive cost.

This phenomenon enables more granular control and rigorous analysis of how LLMs integrate external knowledge. To investigate the underlying mechanisms, we study layer-wise external knowledge exploitation by measuring the divergence of ablating answer-determining context, i.e., the specific text segment within retrieved documents that directly supports the correct answer to a query. By intervening on injecting noise and measuring its

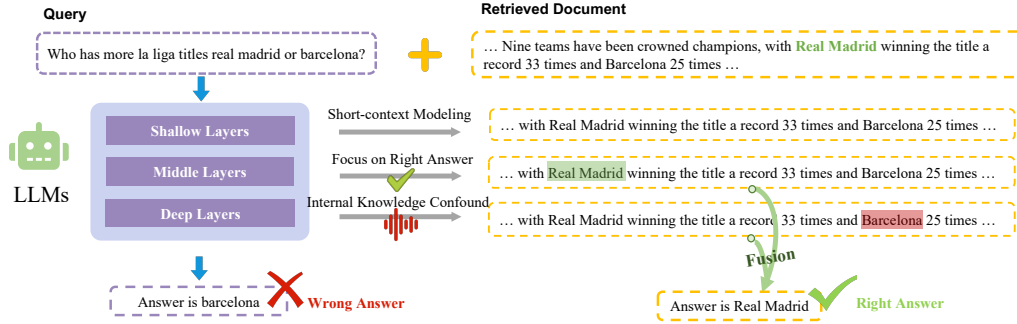


Figure 1: Illustration for layer-wise behavior in LLMs for RAG. Given a query and retrieved documents with the correct answer (“Real Madrid”), shallow layers capture local context, middle layers focus on answer-relevant content, while deep layers may over-rely on internal knowledge and hallucinate (e.g., “Barcelona”). Our proposal, LFD fuses middle-layer signals into the final output to preserve external knowledge and improve accuracy.

085 impact on divergence patterns, we identify the relative  
 086 importance of different layers in exploiting  
 087 external knowledge. Empirically, we find that noise  
 088 amplifies the contribution of answer-determining  
 089 context in middle layers, highlighting their critical  
 090 role in integrating long-range external information.  
 091 To further support this observation, we compare  
 092 attention distributions across heads and layers with  
 093 versus without answer-determining context. The  
 094 analysis shows that attention differences peak in  
 095 middle layers but decline in later layers, signaling  
 096 a transition from external knowledge reliance to  
 097 internal parametric knowledge utilization. Following  
 098 these observations, we propose a functional categorization  
 099 of LLM layers: (1) shallow layers for short-context  
 100 modeling, (2) intermediate layers for external  
 101 knowledge integration, and (3) deeper layers for  
 102 internal knowledge transformation. Therefore, when  
 103 retrieved context already contains the correct answer,  
 104 excessive dependence on internal knowledge in later  
 105 layers introduces confounding effects, reducing  
 106 generation accuracy, as visualized in Figure 1 (left).  
 107

108 Based on this insight, we propose a simple decoding  
 109 method, Layer Fusing Decoding (LFD), which  
 110 enhances access to external factual knowledge  
 111 without introducing additional noise overhead. The  
 112 core idea is illustrated in the right panel of  
 113 Figure 1. LFD fuses representations from the  
 114 long-term context retrieval layer, where external  
 115 knowledge is most effectively integrated, directly  
 116 into the final decoding layer to maximize factual  
 117 grounding. To identify the appropriate layer for  
 118 fusion, we track the model’s reliance on internal  
 119 knowledge by measuring changes in hidden states  
 120 across transformer feed-forward network (FFN)  
 121 layers, where model knowledge is primarily stored  
 122 (Geva et al., 2021; Dai et al., 2022). Specifically,  
 123 we select the layer exhibiting minimal internal  
 knowledge influ-

124 ence from the latter half of the model’s layers. This  
 125 criterion ensures that LFD captures the externally  
 126 grounded signal before it is overridden by parametric  
 127 knowledge in later layers. Importantly, LFD  
 128 operates at inference time, requiring no post-hoc  
 129 fine-tuning or architectural modifications, making  
 130 it easily integrable into existing LLM pipelines. Finally,  
 131 extensive empirical validation across diverse  
 132 model architectures and datasets demonstrates that  
 133 LFD delivers competitive performance relative to  
 134 noise-based approaches, while incurring significantly  
 135 lower computational overhead.

## 136 2 Preliminary

### 137 2.1 Formulation of RAG

138 In RAG systems, the generator  $\mathcal{G}$  (typically a  
 139 LLM) is expected to produce accurate and well-  
 140 grounded responses based on retrieved documents  
 $D = \{d_1, d_2, \dots, d_\lambda\}$ . To quantify this capability,  
 141 we define  $A = \{a_1, \dots, a_\lambda\}$  as the set of key  
 142 information extracted from  $D$  that is necessary for  
 143 generating an accurate answer. The performance of  
 144 the RAG system can be evaluated by measuring the  
 145 inclusion rate of  $A$  in its output response  $r$ , which  
 146 reflects the model’s ability to fully utilize valuable  
 147 documents. To produce the final response, the system  
 148 first encodes the query  $q$  and documents  $D$  into  
 149 a structured prompt through an instruction template  
 150  $\mathcal{T}$ , which instantiates the prompt  $P = \mathcal{T}(q, D)$ ,  
 151 then the generator processes this prompt to produce  
 152 the final response  $r$ . To optimize the generation  
 153 process, the generator  $\mathcal{G}$  aims to ensure that all  
 154 answers contained within  $A$  are included in the  
 155 generator’s output. The accurate answer generated  
 156 by  $\mathcal{G}$  can be formalized as:  
 157

$$158 \quad r = \mathcal{G}(q, D), \text{ s.t. } \forall a_i \in A, \mathcal{I}(r, a_i) = \text{True}, \quad (1)$$

159 where  $\mathcal{I}(r, a_i) = \text{True}$  means the answer  $a_i$  is  
 160 included in  $r$ .  
 161

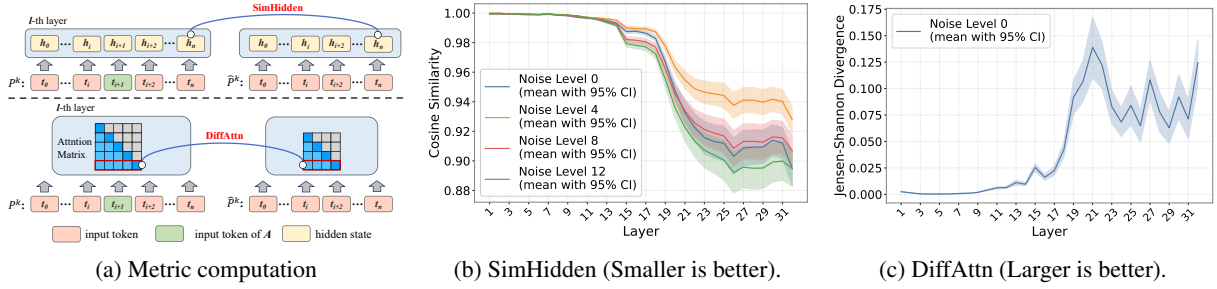


Figure 2: (a) Calculation procedure for SimHidden and DiffAttn metrics; (b) Average SimHidden scores (with 95% confidence intervals) across layers under varying noise levels (0, 4, 8, 12); (c) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0.

## 2.2 External Knowledge Intervention in RAG

The counterintuitive effectiveness of noise in RAG systems motivates a deeper investigation into layer-wise behavior in LLMs. To this end, we conduct an empirical study that contrasts layer-wise representation dynamics under two controlled interventions: (1) ablation of the answer-determining context, and (2) injection of varying levels of noise into the retrieved documents. This differential analysis reveals how noise modulates the model’s internal information flow, amplifying the influence of external knowledge in middle layers while mitigating the model’s tendency to over-rely on internal parametric memory. Our findings highlight key transformation layers where noise injection helps reduce context-dependent fragility, providing a foundation for our proposed decoding strategy.

**Experimental Setup** We simulate noisy level by adding  $k$  irrelevant Wikipedia documents  $N_k = \{n_1, \dots, n_k\}$  to each prompt input (Cuconasu et al., 2024). To analyze how external knowledge flows in LLMs, we generate a modified document set  $\hat{D}$ , which deletes key information  $A$  from the original documents  $D$ . Without loss of generality, we assume the retrieval set  $D$  contains a single document, i.e.,  $D = \{d_1\}$ . By varying noise levels  $k$ , we analyze how external knowledge impacts different layers of the model by comparing two prompts: the **original prompt**  $P^k = \mathcal{T}(q, D, N_k)$  (containing key information  $A$ ) and the **modified prompt**  $\hat{P}^k = \mathcal{T}(q, \hat{D}, N_k)$ , shown in Figure 2(a). This section focuses on Llama2-7B and the NQ dataset. Additional results are in Appendix A.

**Quantify External Knowledge’s Influence** To measure how external knowledge influences LLMs, we compare intermediate representations before and after removing critical information  $A$ . For each layer  $l \in \{1, \dots, L\}$  in the model, we calculate the cosine similarity between the feed-forward network (FFN) outputs of the original input  $P^k$  and its

modified version  $\hat{P}^k$ , focusing on the final prompt token as

$$\text{SimHidden}_l(P^k, \hat{P}^k) = \frac{\mathbf{h}_l(P^k) \cdot \mathbf{h}_l(\hat{P}^k)}{\|\mathbf{h}_l(P^k)\| \cdot \|\mathbf{h}_l(\hat{P}^k)\|}, \quad (2)$$

where  $\mathbf{h}_l(P)$  denotes the intermediate representation of layer  $l$  under prompt  $P$ . This metric reveals how significantly removing  $A$  disrupts the model’s contextual processing at each layer. A higher score indicates the model’s understanding remains consistent even after removing  $A$ , while a lower score suggests removing  $A$  plays a critical role in shaping the layer’s output. By analyzing this metric across varying levels of noise injection ( $k$ ), we can assess how different noise perturbation intensities affect the model’s reliance on external knowledge, providing insights into how contextual information is integrated across layers.

Additionally, since the divergence of  $\mathbf{h}_l(P^k)$  and  $\mathbf{h}_l(\hat{P}^k)$  tends to accumulate in deeper layers, we further analyze attention patterns before and after the removal of the answer-determining context  $A$ . For each transformer layer  $l$ , we compute the average Jensen-Shannon Divergence (JSD) (Lin, 2002) across all attention heads to quantify distributional shifts in attention:

$$\text{DiffAttn}_l(P^k, \hat{P}^k) = \frac{1}{M} \sum_{m=1}^M \text{JSD}(\alpha_{l,m}(P^k) \parallel \hat{\alpha}_{l,m}(\hat{P}^k)), \quad (3)$$

where  $\alpha_{l,m}(P^k)$  denotes the softmax-normalized attention distribution of head  $m$  in layer  $l$  at the final token position for prompt  $P^k$ . The variant  $\hat{\alpha}_{l,m}(\hat{P}^k)$  is computed by filling attention scores corresponding to answer-determining context  $A$  with  $-\infty$  before softmax normalization, thereby preserving the relative distribution over remaining tokens in modified prompt  $\hat{P}^k$ . Larger divergence scores indicate that the external knowledge  $A$  significantly influences the model’s attention. Note that we concentrate on intervention of noise-free models (where  $k = 0$ ), we compare the original prompts  $P^0$  with their intentionally altered counterparts  $\hat{P}^0$ .

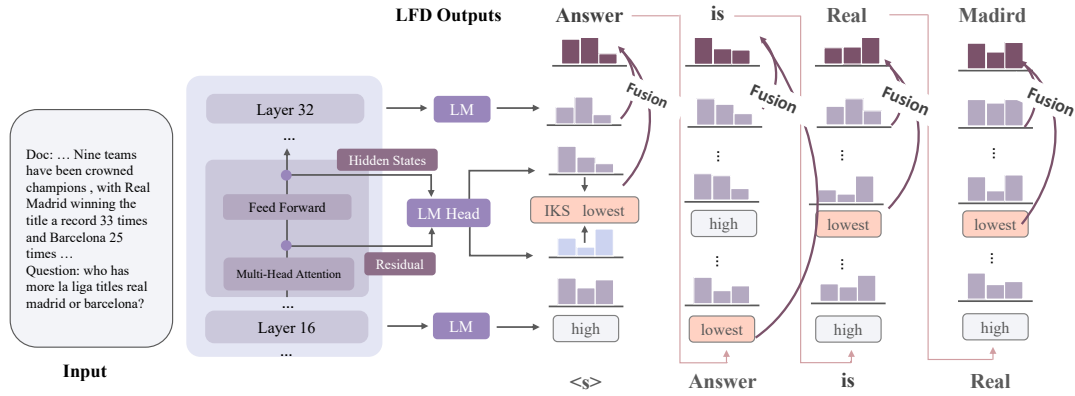


Figure 3: The proposed LFD includes two key components: (1) A dynamic layer selection method using IKS to pinpoint the most impactful layers for integrating external retrieval knowledge. (2) A knowledge fusion mechanism that merges external information with the model’s predictions after adaptive filtering to ensure alignment with the model’s reasoning.

**Analysis and Conclusion** The quantified impacts of external knowledge, as depicted in Figure 2(b-c), lead to the following observations and conclusions: **(1) The early layers (1-14) primarily perform short-context modeling**, which means they focus more on capturing local token relationships rather than integrating global contextual information. This manifests through two key observations: the hidden state similarity remains relatively high across all noise conditions and the attention divergence also stay constantly low compared to other layers. **(2) The middle layers (15-26) demonstrate long-term context retrieval capabilities**, as evidenced by two complementary patterns: a progressive decline in  $\text{SimHidden}_l(P^k, \hat{P}^k)$  (Figure 2 (b)) and a corresponding increase in  $\text{DiffAttn}_l$  (Figure 2 (c)). This dual evidence indicates these layers’ heightened sensitivity to the removal of  $A$  and their capacity for comprehensive global context integration. Meanwhile, an increasing noise level, especially when  $k \geq 8$ , leads to greater discrepancies in hidden state similarity, suggesting that a certain amount of noise can enhance the model’s focus on  $A$ , thereby improving answer accuracy. **(3) The deeper layers (27-32) exhibit characteristics of parametric knowledge utilization.** As we can observe, hidden state similarity does not continue to decrease as the model depth increases, instead, it shows a rebound, with  $\text{SimHidden}$  increasing by a maximum of 0.1. Meanwhile, attention divergence, after peaking at layer 21, also exhibits a moderate decline, with  $\text{DiffAttn}$  decreasing by a maximum of approximately 0.06. This indicates that the role of internal knowledge may be enhanced, as the model may focus more on processing the already captured contextual information rather than contin-

uing to attend to external knowledge. These observations motivate us to design our own methods to better leverage external knowledge with internal representations, thereby improving the model’s performance in RAG systems.

### 3 Layer Fused Decoding

This section present LFD, a framework designed to improve how external knowledge is integrated into model predictions while retaining accuracy. Our approach has two core components: (1) A **dynamic external knowledge layer identification strategy**, which automatically selects the most impactful layer for integrating retrieved context. This selection is guided by Internal Knowledge Scores (IKS), which measure how strongly each layer reflects the model’s parametric knowledge. (2) An **external knowledge fused decoding** mechanism, which merges external knowledge representations with the model’s final output. An adaptive filtering step precedes fusion to ensure the incorporated information complements the model’s reasoning. Figure 3 illustrates the complete workflow of our approach.

#### 3.1 Dynamic External Knowledge Layer Identification

Effective integration of retrieval-based knowledge requires identifying layers that are responsive to external context while minimally influenced by internal parametric knowledge. To this end, we quantify the influence of internal knowledge across layers and select candidate layers based on the external information integration.

**Internal Knowledge Score** Recent advances in transformer interpretability reveal that FFN layers function as specialized knowledge repositories in LLMs (Geva et al., 2021; Dai et al., 2022). To quantify how different layers utilize this inter-

nal knowledge, we design the Internal Knowledge Score (IKS), a metric that captures layer-specific knowledge transformations. Given an input prompt  $P$ , let  $\mathbf{h}_l^{\text{in}}(P) \in \mathbb{R}^d$  and  $\mathbf{h}_l^{\text{out}}(P) \in \mathbb{R}^d$  denote the input and output activations of the  $l$ -th FFN layer. We project these vectors into the vocabulary space via LogitLens (Belrose et al., 2023) of LLMs, parameterized by  $W_{\text{LM}} \in \mathbb{R}^{d \times |V|}$ , as follows:

$$\begin{aligned} \mathbf{p}_l^{\text{in}} &= \text{softmax}(W_{\text{LM}}\mathbf{h}_l^{\text{in}}(P)), \\ \mathbf{p}_l^{\text{out}} &= \text{softmax}(W_{\text{LM}}\mathbf{h}_l^{\text{out}}(P)). \end{aligned} \quad (4)$$

The IKS for layer  $l$  is defined as the JSD divergence between these distributions as

$$\text{IKS}_l(P) = \text{JSD}(\mathbf{p}_l^{\text{in}} \parallel \mathbf{p}_l^{\text{out}}). \quad (5)$$

This divergence quantifies the **parametric knowledge impact** of FFN layers, where higher IKS indicates greater transformation in the residual stream and stronger reliance on internal knowledge.

**External Knowledge Layer Selection** To identify the optimal layer for leveraging context-derived factual knowledge during inference, we propose two principled criteria for layer selection: **(1) Integration of Late-Stage Layers:** We integrate the latter half of the LLM’s layers for final decoding. This choice is motivated by the observation that early layers predominantly focus on short-context modeling, with limited capacity to capture contextual dependencies. The integration of external knowledge becomes progressively stronger in middle-to-late layers (as analyzed in Section 2.2). **(2) Lowest IKS Layer Selection:** Within the subset of late-stage layers, we select the layer exhibiting the lowest IKS score. This layer strikes a balance by retaining sufficient external contextual signals while substantially mitigating distortion of the model’s inherent knowledge representations. By selecting these layers, we maximize the exploitation of retrieval context before the dominance of internal parametric knowledge obscures external signals. Empirical validation further confirms the efficacy of this strategy (detailed in Section 4.3).

### 3.2 External Knowledge Fused Decoding

We propose an intervention-aware fusion framework that dynamically integrates intermediate representations from layer  $i$  (identified via IKS scoring) with the final layer’s predictions. To establish distributional coherence between these complementary knowledge sources, we first compute normalized log-probabilities through log softmax transformation:

$$\begin{aligned} \tilde{\mathbf{p}}_i^{\text{out}} &= \log \text{softmax}(W_{\text{LM}}\mathbf{h}_i^{\text{out}}(P)), \\ \tilde{\mathbf{p}}_L^{\text{out}} &= \log \text{softmax}(W_{\text{LM}}\mathbf{h}_L^{\text{final}}(P)) \end{aligned} \quad (6)$$

where  $W_{\text{LM}}\mathbf{h}_i^{\text{out}}(P)$  and  $W_{\text{LM}}\mathbf{h}_L^{\text{out}}(P)$  denote the raw logits from the intervention layer and final layer respectively. To mitigate noise amplification from early layer predictions while preserving critical external knowledge signals, we implement a dynamic gating mechanism inspired by (Li et al., 2022):

$$\tilde{\mathbf{f}}_i(t) = \begin{cases} \tilde{\mathbf{p}}_i^{\text{out}}(t) + \tilde{\mathbf{p}}_L^{\text{out}}(t), & \text{if } \tilde{\mathbf{p}}_i^{\text{out}}(t) \geq \Gamma \\ -\infty, & \text{otherwise,} \end{cases} \quad (7)$$

Within this gating mechanism, the threshold is formulated as  $\Gamma = \min\{\tau \cdot \max(\tilde{\mathbf{p}}_L^{\text{out}}), \max\text{-s}(\tilde{\mathbf{p}}_L^{\text{out}})\}$ . Here,  $\max(\cdot)$  and  $\max\text{-s}(\cdot)$  denote the maximum and  $s$ -th maximum values of the final output layer logits  $\tilde{\mathbf{p}}_L^{\text{out}}$ , respectively, and  $\tau \in [0, 1]$  is a hyperparameter that controls the truncation aggressiveness to filter out low-probability tokens. The final decoding distribution, derived via normalized fusion  $\tilde{\mathbf{f}}_i = \text{softmax}(\tilde{\mathbf{f}}_i)$ , enables synergistic knowledge transfer between layers, preserving the final layer’s discriminative capacity to balance the integration of external knowledge with the model’s inherent confidence.

## 4 Experiments

### 4.1 Experiment Setup

**Datasets and Metrics** Building on related works (Cuconasu et al., 2024; Wu et al., 2024; Jeong et al., 2024), our evaluation encompasses both single-hop and multi-hop QA tasks. Single-hop performance is measured using the widely-used Natural Questions (NQ) (Kwiatkowski et al., 2019) and RGB (Chen et al., 2024a) benchmarks, which comprise real-world queries. Turning to multi-hop QA, we select HotpotQA (Yang et al., 2018) and 2WikiMultiHopQA (Ho et al., 2020), enabling evaluation of complex reasoning types such as comparison and composition. Across all tasks, we report accuracy as the primary metric. More details of the datasets and metrics are in Appendix B.

**Baselines** To demonstrate its broad effectiveness, we evaluate LFD across four widely used LLMs: Llama2-7B (Touvron et al., 2023), Mistral-7B (Jiang et al., 2023a), DeepSeek-7B (Bi et al., 2024), and Qwen3-8B (Team, 2025). We further benchmark LFD against three established decoding strategies: **(1) Greedy Decoding (GD)** employs standard autoregressive decoding. We augment GD with varying numbers of noise documents, denoted as GD (0), GD (4), GD (8), and GD (12). **(2) Contrastive Search (CS)** (Su et al., 2022) balances response quality and diversity to encourage more coherent outputs. **(3) Decoding by Contrasting**

		NQ	RGB	HotpotQA			2WikiMultihopQA				
				Compare	Bridge	Total	Compare	Bridge	Inf	Compose	Total
Llama2-7B	GD (0)	0.5745	0.8900	0.4755	0.4108	0.4237	0.3095	0.2815	0.2607	0.2928	0.2904
	GD (4)	0.5559	0.8267	0.5084	0.4049	0.4257	0.2559	0.2564	0.2048	0.2405	0.2433
	GD (8)	0.5965	0.8300	0.5091	0.4096	0.4317	<b>0.3257</b>	0.2983	0.2789	0.2995	0.3030
	GD (12)	<b>0.6383</b>	0.8200	0.5286	0.4236	0.4447	0.3204	<b>0.3202</b>	0.2646	0.2943	0.3027
	CS	0.5775	0.8567	0.4284	0.4049	0.4096	0.2800	0.2462	0.2464	0.2866	0.2712
	DoLA	0.5809	0.8733	0.4438	0.3883	0.3995	0.2738	0.2411	0.2289	0.2592	0.2550
	LFD (Random)	0.5928	0.8900	0.4385	0.4295	0.4313	0.3065	0.2845	0.2744	0.3066	0.2978
	LFD	0.5949	<b>0.9067</b>	<b>0.5453</b>	<b>0.4295</b>	<b>0.4528</b>	0.3174	<u>0.3005</u>	<b>0.3043</b>	<b>0.3232</b>	<b>0.3144</b>
Mistral-7B	GD (0)	0.6130	0.8900	0.5864	0.5889	0.5884	0.5518	0.5204	0.5605	0.5337	0.5385
	GD (4)	0.5667	0.8033	0.5494	0.5348	0.5377	0.3681	0.3443	0.3173	0.3296	0.3406
	GD (8)	0.5678	0.7700	0.5326	0.5260	0.5273	0.2963	0.2874	0.2503	0.2582	0.2728
	GD (12)	0.5814	0.8267	0.5440	0.5461	0.5457	0.3098	0.2888	0.2503	0.2655	0.2795
	CS	0.5598	0.7433	0.4149	0.5117	0.4922	0.4346	0.3873	0.4486	0.4740	0.4423
	DoLA	0.6142	0.8900	0.5931	0.5870	0.5883	0.5366	0.5018	0.5572	0.5239	0.5262
	LFD (Random)	0.6270	0.9133	0.5985	0.6016	0.6009	0.5687	0.5317	0.5767	0.5506	0.5541
	LFD	<b>0.6357</b>	<b>0.9367</b>	<b>0.6026</b>	<b>0.6093</b>	<b>0.6079</b>	<b>0.5813</b>	<b>0.5434</b>	<b>0.5982</b>	<b>0.5645</b>	<b>0.5680</b>
DeepSeek-7B	GD (0)	0.5250	0.8233	0.3282	0.4033	0.3883	0.2625	0.2414	0.2663	0.2609	0.2609
	GD (4)	0.5216	0.8300	0.3968	0.4437	0.4343	0.2807	0.2939	0.2744	0.2730	0.2796
	GD (8)	0.5226	<b>0.8800</b>	0.4506	0.4579	0.4564	0.3578	0.3603	0.3407	0.3463	0.3515
	GD (12)	0.5204	<b>0.8933</b>	0.4573	<b>0.4618</b>	<b>0.4609</b>	0.3701	0.3705	0.3349	0.3614	0.3622
	CS	0.5322	0.8000	0.4028	0.4584	0.4472	0.3929	0.3833	0.3888	0.4077	0.3964
	DoLA	0.3755	0.5033	0.3490	0.3003	0.3100	0.3039	0.3118	0.2484	0.2594	0.2803
	LFD (Random)	0.4858	0.6700	0.3847	0.3734	0.3757	0.3466	0.3563	0.3277	0.3320	0.3403
	LFD	<b>0.5412</b>	0.8267	<b>0.4801</b>	0.4466	0.4533	<b>0.4492</b>	<b>0.4227</b>	<b>0.4194</b>	<b>0.4223</b>	<b>0.4285</b>
Qwen3-8B	GD (0)	0.7318	0.9571	0.6960	0.6708	0.6759	0.6637	0.6335	0.5897	0.6085	0.6250
	GD (4)	0.7213	0.9467	0.6658	0.6544	0.6567	0.6044	0.5729	0.4948	0.5268	0.5517
	GD (8)	0.7233	0.9500	0.6584	0.6568	0.6571	0.6117	0.5780	0.5078	0.5370	0.5605
	GD (12)	0.7133	0.9533	0.6530	0.6582	0.6571	0.6146	0.5722	0.5052	0.5462	0.5634
	CS	0.7204	0.9533	0.7081	0.6762	0.6826	0.6468	0.6204	0.5754	0.6015	0.6134
	DoLA	0.7168	0.9533	0.7014	0.6703	0.6766	0.6485	0.6058	0.5650	0.5929	0.6057
	LFD (Random)	0.7357	0.9567	0.7108	0.6935	0.6970	0.6627	0.6334	0.6014	0.6136	0.6289
	LFD	<b>0.7380</b>	<b>0.9600</b>	<b>0.7182</b>	<b>0.6974</b>	<b>0.7016</b>	<b>0.6663</b>	<b>0.6342</b>	<b>0.6034</b>	<b>0.6148</b>	<b>0.6301</b>

Table 1: The accuracy performance comparison of different methods on four datasets. Compare, Bridge, Inf, and Compose represent distinct reasoning question types (details in Appendix B). **Bold** values indicate the best performance, while underlined values represent the second-best.

**Layers** (DoLA) (Chuang et al., 2024) is a contrastive decoding approach that reduces hallucinations by comparing predictions across different model layers. **(4) LFD (Random)** is a variant of LFD that randomly selects an intermediate layer to fuse during the final decoding stage. Implementation details are provided in Appendix C.

## 4.2 Main Results

Table 1 summarizes the accuracy of RAG based on four different models across four QA datasets. From the table, we have the following observations: **(1) LFD matches or exceeds the performance of noise-injection strategies.** We can see our method demonstrates consistent performance gains, ranging from minimal 0.29% (Qwen3-8B, RGB) to maximal 16.76% (DeepSeek-7B, 2WikiMultihopQA) improvement. Injecting noise, i.e., GD (12), shows performance gains on some datasets and models, with the highest improvement being 10.13% (DeepSeek-7B, 2WikiMultihopQA). However, it significantly degrades performance for Mistral-7B and Qwen3-8B across all datasets, leading to a notable reduction in accuracy (maximum  $\Delta = -26.57\%$  on Mistral-7B). **Additionally, our method also achieves significantly lower computational overhead than noise-injection strategies (see Appendix D).**

**(2) LFD outperforms decoding strategies without noise injection.** Compared to alternative de-

coding methods, LFD consistently delivers superior performance. While methods like DoLA and CS show strong results in specific cases, e.g., achieving 1.94% and 13.55% gains on DeepSeek-7B with the 2WikiMultihopQA dataset, they occasionally underperform even relative to the greedy decoding baseline. In particular, DoLA shows the largest decline of 30% on DeepSeek-7B with the RGB dataset, while CS exhibits a maximum drop of 14.57% on Mistral-7B with the RGB dataset. These results indicate that both strategies are sensitive to specific model architectures or data characteristics. **Additional comparisons of these methods under different noise levels and the case study are provided in Appendix E and Appendix H, respectively.**

**(3) The fusion layer selection plays a critical role for RAG.** Our dynamic layer selection strategy consistently outperforms the random approach, with an average improvement of 3.11%, particularly achieving a 15.67% gain on DeepSeek-7B/RGB, demonstrating its efficacy.

## 4.3 Analysis Experiments

**Effects of late-stage layer Integration** We assess how different layer selection ranges affect performance using Llama2-7B and Mistral-7B across four benchmark datasets in Table 2. As Table 2 shows, our approach (selecting from layers 16–32) consistently achieves higher accuracy than select-

		NQ	RGB	HotpotQA			2WikiMultihopQA				
				Compare	Bridge	Total	Compare	Bridge	Inf	Compose	Total
Llama2-7B	GD (0)	0.5745	<u>0.8900</u>	<u>0.4755</u>	0.4108	<u>0.4237</u>	<u>0.3095</u>	0.2815	0.2607	0.2928	0.2904
	LFD[0, 16]	0.5758	0.8833	0.4371	0.4027	0.4096	0.2956	0.2582	0.2529	0.2739	0.2731
	LFD[0, 32]	<b>0.5957</b>	0.9067	0.4586	0.4208	0.4284	<b>0.3181</b>	<b>0.3038</b>	0.3004	<b>0.3237</b>	<b>0.3151</b>
	LFD[16, 32]	0.5949	<b>0.9067</b>	<b>0.5454</b>	<b>0.4295</b>	<b>0.4528</b>	0.3174	0.3005	<b>0.3043</b>	0.3232	0.3144
Mistral-7B	GD (0)	0.6130	<u>0.8900</u>	<u>0.5864</u>	0.5889	0.5884	0.5518	0.5204	0.5605	0.5337	0.5385
	LFD[0, 16]	0.6081	0.8867	0.5864	0.5919	0.5908	0.5551	0.5218	0.5650	0.5402	0.5429
	LFD[0, 32]	0.6344	0.9367	0.5938	0.5994	0.5982	<b>0.5842</b>	<b>0.5492</b>	0.5917	0.5639	<b>0.5690</b>
	LFD[16, 32]	<b>0.6357</b>	<b>0.9367</b>	<b>0.6026</b>	<b>0.6093</b>	<b>0.6079</b>	0.5813	0.5434	<b>0.5982</b>	<b>0.5645</b>	0.5680

Table 2: Comparison of accuracy between different layer selection ranges under dynamic layer selection strategy. **Bold** indicate the best performance, while underline represent the second-best. LFD[0, 16) and LFD[16, 32) mean selecting layers from the earlier and later half respectively.

ing from earlier layers (layers 0-16), with an average gain of 3.01%. Notably, our approach provides marginal accuracy gains ( $\leq 0.1\%$ ) over full-range selection (layers 0-32), which indicates the advantage of the proposed layer selection strategy. These results demonstrate the greater efficacy of deeper layers for utilizing knowledge in RAG. Results for DeepSeek-7B and Qwen3-8B are provided in Appendix F.

**Effects of the lowest IKS layer selection** We evaluate the effectiveness of the lowest IKS layer selection on the NQ and HotpotQA datasets. Extended results across models and datasets appear in Appendix G. First, we compare our dynamic strategy with fixed-layer selection LFD (Fixed) in Figure 4(a-b). Second, we analyze the distribution of dynamical layer selections against the optimal fixed-layer baseline on NQ (Figure 4(c)).

Key findings emerge: **(1) Fixed-layer selection requires dataset-specific validation for optimality.** While fixed-layer achieves peak performance at layer 24 – 26 (NQ) and layer 22 (HotpotQA), these layers differ across datasets, necessitating extra validating datasets. **(2) The lowest IKS tends to achieve near-optimal performance with small margins.** Compared to the best fixed-layer results, the lowest IKS exhibits performance gaps of only 1.8% (NQ) and 0.2% (HotpotQA), demonstrating robust generalization without dataset-specific tuning. **(3) Dynamic layer selection concentrates near the optimal fixed-layer.** For NQ dataset, the lowest IKS selections cluster around layer 26 (Figure 4(c)), aligning closely with the optimal fixed layer despite stochasticity.

**Effects of Conflicting Documents** Real-world retrieval often involves conflicting documents, which are semantically similar to the query but do not contain the correct answer (Cuconasu et al., 2024). Following (Chen et al., 2024a), we utilize the RGB dataset with five documents per prompt and vary the ratio of conflicting documents from

		0	0.2	0.4	0.6	0.8
Llama2-7B	GD (0)	0.9267	0.9167	0.9167	0.8933	0.8500
	LFD	<b>0.9267</b>	<b>0.9167</b>	<b>0.9167</b>	<b>0.8967</b>	<b>0.8833</b>
Mistral2-7b	GD (0)	0.8933	0.8933	0.8867	0.8867	0.8367
	LFD	<b>0.9567</b>	<b>0.9333</b>	<b>0.9233</b>	<b>0.9167</b>	<b>0.8800</b>

Table 3: Accuracy comparison under varying ratios of conflicting documents (from 0 to 0.8) on the RGB dataset. **Bold** values indicate the best performance.

		NQ	HotpotQA
Glm4-9B	GD (0)	0.7980	0.6620
	LFD	<b>0.8280</b>	<b>0.7340</b>
Llama2-13B	GD (0)	0.6500	0.5940
	LFD	<b>0.7280</b>	<b>0.5980</b>
Qwen3-32B	GD (0)	0.8440	0.7660
	LFD	<b>0.8480</b>	<b>0.7780</b>

Table 4: Comparison of accuracy across different larger models on the NQ and HotpotQA datasets. **Bold** values indicate the best performance.

0 to 0.8. As shown in Table 3, LFD consistently outperforms GD (0), demonstrating its robustness when encountering misleading interference.

**Results for Larger Models** Following the same experimental setup (Section 4.1), We extend our evaluation to larger models on NQ and HotpotQA. As shown in Table 4, LFD consistently yields significant gains (e.g., +7.8% on Llama2-13B for NQ).

## 5 Related Works

### 5.1 Retrieval Augmented Generation

Retrieval-Augmented Generation (RAG) enhances model reasoning by integrating relevant external knowledge retrieved through user queries. Recent advances focus on three directions: refined retrieval mechanisms (Jeong et al., 2024; Chan et al., 2024), structured knowledge organization (Procko and Ochoa, 2024; Baek et al., 2023), and optimized context embedding (Hu et al., 2024; Sahoo et al., 2024). Self-RAG (Asai et al., 2024) and FLARE (Jiang et al., 2023c) achieve adaptive retrieval through self-evaluation and uncertainty prediction respectively, dynamically optimizing knowledge acquisition. GraphRAG (He et al., 2024; Edge et al.,

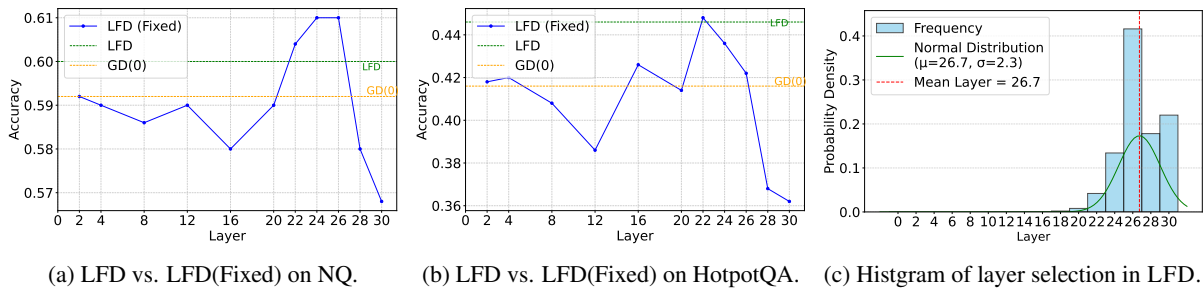


Figure 4: Comparison between LFD and LFD (Fixed) on the NQ (a) and HotpotQA (b) datasets. (c) illustrates the layer selection distribution in LFD compared to the optimal fixed layer selection.

2024; Jiang et al., 2023b) advances reasoning capabilities by constructing document-derived knowledge graphs that capture semantic relationships for multi-hop inference. Embedding optimizations include noise injection (Cuconasu et al., 2024; Wu et al., 2024) to counter overfitting through strategic low-relevance document insertion, and context re-ranking (Yu et al., 2024; Dong et al., 2024) that prioritizes high-utility knowledge via learned document scoring. While these methods enhance knowledge orchestration through pipeline improvements, they systematically neglect the internal mechanisms through which LLMs process external information during generation. Our work bridges this fundamental gap by surfacing stratified knowledge integration patterns across LLM’s layers through systematic layer-wise analysis.

## 5.2 Decoding Strategy in LLMs

Decoding strategies are pivotal in transforming raw model probabilities into coherent text outputs, critically influencing the quality and factual integrity of LLM generations (Shi et al., 2024). While traditional methods like greedy decoding and beam search (Meister et al., 2020; Xie et al., 2023) remain prevalent, recent work has introduced advanced techniques to address their limitations. Contrastive search (CS) (Su et al., 2022) balances diversity and coherence by selecting tokens through a weighted combination of probability and semantic dissimilarity to preceding context, mitigating repetition while preserving fluency. Simple Decoding (FSD) (Yang et al., 2024) suppresses redundant patterns by dynamically constructing an anti-language model to penalize overused token sequences. DoLa (Chuang et al., 2024) addresses hallucinations by contrasting later-layer logit distributions with earlier ones, prioritizing factually consistent predictions. Building on these advances, we propose a novel decoding strategy for RAG that dynamically balances external retrieved knowledge with the model’s internal parametric knowledge, enhancing factual accuracy by mitigating interference from outdated or con-

flicting internal representations.

## 5.3 Hallucinations in LLMs

Hallucinations in LLMs refer to instances where the model generates false or unsupported information not grounded in its reference data (Perković et al., 2024). Existing mitigation strategies include multi-agent debating, where multiple LLM instances collaborate to detect inconsistencies through iterative debates (Chen et al., 2024b; Du et al., 2023); self-consistency verification, which aggregates and reconciles multiple reasoning paths to reduce individual errors (Wang et al., 2022); and model editing, which directly modifies neural network weights to correct systematic factual errors (Zhang et al., 2024; Hartvigsen et al., 2023). While RAG systems aim to ground responses in retrieved external knowledge, recent studies show that they still exhibit hallucinations, especially those that contradict the retrieved content (Sun et al., 2025b). To address this limitation, our work conducts an empirical study analyzing how LLMs internally process external knowledge in RAG settings by controlling the noise from different granularity. Based on these findings, we propose a novel decoding method designed to improve answer accuracy and reduce hallucination by enhancing the integration of retrieved evidence.

## 6 Conclusion

By analyzing how noise injection amplifies external knowledge exploitation in LLMs, we establish a functional demarcation across LLMs’ layers: shallow (local context), intermediate (external knowledge), and deep (internal parametric knowledge). Leveraging this, we propose LFD, a training-free decoding strategy that fuses intermediate-layer representations to enhance external knowledge integration in final outputs via a the lowest internal knowledge score to pinpoint the ideal fusion layer. Experiments across diverse benchmarks demonstrate that LFD enhances factual grounding in RAG systems while incurring minimal computational overhead.

## 624 Limitations

625 Our work has the following two main limitations.

626 First, LFD occasionally yields slightly lower  
627 gains than GD(12). However, it maintains consist-  
628 ent superiority over GD(0) in every experimental  
629 configuration. More importantly, LFD achieves  
630 these results with only 1/10 of the context tokens  
631 and 0.4x the latency of GD(12) (see AppendixD).  
632 This superior computational efficiency, combined  
633 with its stable performance, underscores the overall  
634 effectiveness of LFD.

635 Second, While noise-enhanced RAG has proven  
636 effective (Cuconasu et al., 2024; Sun et al., 2025a),  
637 a rigorous theoretical justification for its success  
638 remains unclear. Our study utilizes noise primarily  
639 as an intervention tool to analyze layer-specific ef-  
640 fects, leading to the development of LFD to bypass  
641 costly augmentation. Theoretical exploration of  
642 such noise-driven improvements remains a promis-  
643 ing direction for future research.

## 644 References

645 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and  
646 Hannaneh Hajishirzi. 2024. Self-RAG: Learning to  
647 retrieve, generate, and critique through self-reflection.  
648 In *The Twelfth International Conference on Learning*  
649 *Representations*.

650 Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023.  
651 [Knowledge-augmented language model prompting](#)  
652 [for zero-shot knowledge graph question answering](#).  
653 In *Proceedings of the 1st Workshop on Natural*  
654 *Language Reasoning and Structured Explanations*  
655 *(NLRSE)*, pages 78–106, Toronto, Canada. Associa-  
656 tion for Computational Linguistics.

657 Nora Belrose, Zach Furman, Logan Smith, Danny Ha-  
658 lawi, Igor Ostrovsky, Lev McKinney, Stella Bider-  
659 man, and Jacob Steinhardt. 2023. Eliciting latent  
660 predictions from transformers with the tuned lens.  
661 *CoRR*, abs/2303.08112.

662 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen,  
663 Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong,  
664 Qiushi Du, Zhe Fu, and 1 others. 2024. Deepseek llm:  
665 Scaling open-source language models with longterm-  
666 ism. *arXiv preprint arXiv:2401.02954*.

667 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann,  
668 Trevor Cai, Eliza Rutherford, Katie Millican, George  
669 van den Driessche, Jean-Baptiste Lespiau, Bogdan  
670 Damoc, Aidan Clark, Diego de Las Casas, Aurelia  
671 Guy, Jacob Menick, Roman Ring, Tom Hennigan,  
672 Saffron Huang, Loren Maggiore, Chris Jones, Albin  
673 Cassirer, and 9 others. 2022. Improving language  
674 models by retrieving from trillions of tokens. In *In-*  
675 *ternational Conference on Machine Learning, ICML*

2022, 17-23 July 2022, Baltimore, Maryland, USA,  
volume 162 of *Proceedings of Machine Learning*  
*Research*, pages 2206–2240. PMLR.

Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo,  
Wei Xue, Yike Guo, and Jie Fu. 2024. Rq-rag: Learn-  
ing to refine queries for retrieval augmented genera-  
tion. *arXiv preprint arXiv:2404.00610*.

Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun.  
2024a. Benchmarking large language models in  
retrieval-augmented generation. In *Proceedings of*  
*the AAAI Conference on Artificial Intelligence*, vol-  
ume 38, pages 17754–17762.

Justin Chen, Swarnadeep Saha, and Mohit Bansal.  
2024b. Reconcile: Round-table conference improves  
reasoning via consensus among diverse llms. In *Pro-*  
*ceedings of the 62nd Annual Meeting of the Associa-*  
*tion for Computational Linguistics (Volume 1: Long*  
*Papers)*, pages 7066–7085.

Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon  
Kim, James R. Glass, and Pengcheng He. 2024. [Dola:](#)  
[Decoding by contrasting layers improves factuality in](#)  
[large language models](#). In *The Twelfth International*  
*Conference on Learning Representations*.

Florin Cuconasu, Giovanni Trappolini, Federico Sicil-  
iano, Simone Filice, Cesare Campagnano, Yoelle  
Maarek, Nicola Tonello, and Fabrizio Silvestri.  
2024. The power of noise: Redefining retrieval  
for RAG systems. In *Proceedings of the 47th In-*  
*ternational ACM SIGIR Conference on Research and*  
*Development in Information Retrieval, SIGIR 2024,*  
*Washington DC, USA, July 14-18, 2024*, pages 719–  
729. ACM.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao  
Chang, and Furu Wei. 2022. [Knowledge neurons in](#)  
[pretrained transformers](#). In *Proceedings of the 60th*  
*Annual Meeting of the Association for Computational*  
*Linguistics (Volume 1: Long Papers)*, pages 8493–  
8502, Dublin, Ireland. Association for Computational  
Linguistics.

Dario Di Palma. 2023. Retrieval-augmented recom-  
mender system: Enhancing recommender systems  
with large language models. In *Proceedings of the*  
*17th ACM Conference on Recommender Systems*,  
pages 1369–1373.

Jialin Dong, Bahare Fatemi, Bryan Perozzi, Lin F. Yang,  
and Anton Tsitsulin. 2024. [Don't forget to connect!](#)  
[improving rag with graph-based reranking](#). *Preprint,*  
*arXiv:2405.18414*.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-  
baum, and Igor Mordatch. 2023. Improving factual-  
ity and reasoning in language models through multi-  
agent debate. In *Forty-first International Conference*  
*on Machine Learning*.

Darren Edge, Ha Trinh, Newman Cheng, Joshua  
Bradley, Alex Chao, Apurva Mody, Steven Truitt,  
Dasha Metropolitansky, Robert Osazuwa Ness, and

732	Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. <i>arXiv preprint arXiv:2404.16130</i> .	
733		
734		
735	Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 6491–6501.	
736		
737		
738		
739		
740		
741		
742	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. <a href="#">Transformer feed-forward layers are key-value memories</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
743		
744		
745		
746		
747		
748		
749	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Retrieval augmented language model pre-training. In <i>Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 3929–3938. PMLR.	
750		
751		
752		
753		
754		
755		
756	Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with grace: Lifelong model editing with discrete key-value adapters. <i>Advances in Neural Information Processing Systems</i> , 36:47934–47959.	
757		
758		
759		
760		
761	Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. <i>Advances in Neural Information Processing Systems</i> , 37:132876–132907.	
762		
763		
764		
765		
766		
767	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. <a href="#">Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps</a> . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
768		
769		
770		
771		
772		
773		
774	Zhibo Hu, Chen Wang, Yanfeng Shu, Hye-Young Paik, and Liming Zhu. 2024. Prompt perturbation in retrieval-augmented generation based large language models. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 1119–1130.	
775		
776		
777		
778		
779		
780	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong Park. 2024. <a href="#">Adaptive-RAG: Learning to adapt retrieval-augmented large language models through question complexity</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7036–7050, Mexico City, Mexico. Association for Computational Linguistics.	
781		
782		
783		
784		
785		
786		
787		
788		
	Jianchao Ji, Zelong Li, Shuyuan Xu, Wenyue Hua, Yingqiang Ge, Juntao Tan, and Yongfeng Zhang. 2024. Genrec: Large language model for generative recommendation. In <i>European Conference on Information Retrieval</i> , pages 494–502. Springer.	789
		790
		791
		792
		793
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023a. <a href="#">Mistral 7b</a> . <i>Preprint</i> , arXiv:2310.06825.	794
		795
		796
		797
		798
		799
		800
		801
	Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen. 2023b. <a href="#">UniKGQA: Unified retrieval and reasoning for solving multi-hop question answering over knowledge graph</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	802
		803
		804
		805
		806
	Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023c. <a href="#">Active retrieval augmented generation</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7969–7992, Singapore. Association for Computational Linguistics.	807
		808
		809
		810
		811
		812
		813
	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>International Conference on Machine Learning</i> , pages 15696–15707. PMLR.	814
		815
		816
		817
		818
	Satyapriya Krishna, Kalpesh Krishna, Anhad Mohanane, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. 2025. <a href="#">Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation</a> . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 4745–4759, Albuquerque, New Mexico. Association for Computational Linguistics.	819
		820
		821
		822
		823
		824
		825
		826
		827
		828
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, and 1 others. 2019. Natural questions: a benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	829
		830
		831
		832
		833
		834
		835
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K��ttler, Mike Lewis, Wen-tau Yih, Tim Rock-t��schel, Sebastian Riedel, and Douwe Kiela. 2020a. Retrieval-augmented generation for knowledge-intensive NLP tasks. In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	836
		837
		838
		839
		840
		841
		842
		843
		844

845	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang,	902
846	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Yifan Wang, Yujiu Yang, and Wai Lam. 2024. A	903
847	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	thorough examination of decoding methods in the era	904
848	täschel, and 1 others. 2020b. Retrieval-augmented	of llms. In <i>Proceedings of the 2024 Conference on</i>	905
849	generation for knowledge-intensive nlp tasks. <i>Ad-</i>	<i>Empirical Methods in Natural Language Processing</i> ,	906
850	<i>advances in neural information processing systems</i> ,	pages 8601–8629.	907
851	33:9459–9474.		
852	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang,	Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Ling-	908
853	Jason Eisner, Tatsunori Hashimoto, Luke Zettle-	peng Kong, and Nigel Collier. 2022. A contrastive	909
854	moyer, and Mike Lewis. 2022. Contrastive decoding:	framework for neural text generation. <i>Advances in</i>	910
855	Open-ended text generation as optimization. <i>arXiv</i>	<i>Neural Information Processing Systems</i> , 35:21548–	911
856	<i>preprint arXiv:2210.15097</i> .	21561.	912
857	Jianhua Lin. 2002. Divergence measures based on the	Hao Sun, Chenming Tang, Gengyang Li, and Yunfang	913
858	shannon entropy. <i>IEEE Transactions on Information</i>	Wu. 2025a. Lost in the passage: Passage-level in-	914
859	<i>theory</i> , 37(1):145–151.	context learning does not necessarily need a "pas-	915
860	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paran-	sage". <i>CoRR</i> , abs/2502.10634.	916
861	jape, Michele Bevilacqua, Fabio Petroni, and Percy	ZhongXiang Sun, Xiaoxue Zang, Kai Zheng, Jun Xu,	917
862	Liang. 2024. Lost in the middle: How language mod-	Xiao Zhang, Weijie Yu, Yang Song, and Han Li.	918
863	els use long contexts. <i>Transactions of the Association</i>	2025b. <a href="#">RedeEP: Detecting hallucination in retrieval-</a>	919
864	<i>for Computational Linguistics</i> , 12.	<a href="#">augmented generation via mechanistic interpretabil-</a>	920
865	Clara Meister, Tim Vieira, and Ryan Cotterell. 2020.	<a href="#">ity</a> . In <i>The Thirteenth International Conference on</i>	921
866	Best-first beam search. <i>Transactions of the Associa-</i>	<i>Learning Representations</i> .	922
867	<i>tion for Computational Linguistics</i> , 8:795–809.	Qwen Team. 2025. <a href="#">Qwen3</a> .	923
868	Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	924
869	Shum, Randy Zhong, Juntong Song, and Tong Zhang.	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	925
870	2024. Ragtruth: A hallucination corpus for develop-	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	926
871	ing trustworthy retrieval-augmented language models.	Bhosale, and 1 others. 2023. Llama 2: Open foun-	927
872	In <i>Proceedings of the 62nd Annual Meeting of the</i>	dation and fine-tuned chat models. <i>arXiv preprint</i>	928
873	<i>Association for Computational Linguistics (Volume</i>	<i>arXiv:2307.09288</i> .	929
874	<i>1: Long Papers)</i> , pages 10862–10878.	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	930
875	Chanhee Park, Hyeonseok Moon, Chanjun Park, and	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	931
876	Heuseok Lim. 2025. <a href="#">MIRAGE: A metric-intensive</a>	Denny Zhou. 2022. Self-consistency improves chain	932
877	<a href="#">benchmark for retrieval-augmented generation eval-</a>	of thought reasoning in language models. <i>arXiv</i>	933
878	<a href="#">uation</a> . In <i>Findings of the Association for Computa-</i>	<i>preprint arXiv:2203.11171</i> .	934
879	<i>tional Linguistics: NAACL 2025</i> , pages 2883–2900,	Jinyang Wu, Feihu Che, Chuyuan Zhang, Jianhua Tao,	935
880	Albuquerque, New Mexico. Association for Computa-	Shuai Zhang, and Pengpeng Shao. 2024. Pandora’s	936
881	tional Linguistics.	box or aladdin’s lamp: A comprehensive analysis re-	937
882	Gabrijela Perković, Antun Drobňjak, and Ivica Botički.	vealing the role of rag noise in large language models.	938
883	2024. <a href="#">Hallucinations in llms: Understanding and</a>	<i>arXiv preprint arXiv:2408.13533</i> .	939
884	<a href="#">addressing challenges</a> . In <i>2024 47th MIPRO ICT and</i>	Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu	940
885	<i>Electronics Convention (MIPRO)</i> , pages 2084–2088.	Zhao, Min-Yen Kan, Junxian He, and Michael Xie.	941
886	Tyler Thomas Procko and Omar Ochoa. 2024. <a href="#">Graph</a>	2023. Self-evaluation guided beam search for rea-	942
887	<a href="#">retrieval-augmented generation for large language</a>	soning. <i>Advances in Neural Information Processing</i>	943
888	<a href="#">models: A survey</a> . In <i>2024 Conference on AI, Sci-</i>	<i>Systems</i> , 36:41618–41650.	944
889	<i>ence, Engineering, and Technology (AIXSET)</i> , pages	Haoyi Xiong, Jiang Bian, Yuchen Li, Xuhong Li, Meng-	945
890	166–169.	nan Du, Shuaiqiang Wang, Dawei Yin, and Sumi	946
891	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha,	Helal. 2024. When search engine services meet	947
892	Vinija Jain, Samrat Mondal, and Aman Chadha.	large language models: visions and challenges. <i>IEEE</i>	948
893	2024. A systematic survey of prompt engineering in	<i>Transactions on Services Computing</i> .	949
894	large language models: Techniques and applications.	Haoran Yang, Deng Cai, Huayang Li, Wei Bi, Wai Lam,	950
895	<i>arXiv preprint arXiv:2402.07927</i> .	and Shuming Shi. 2024. <a href="#">A frustratingly simple de-</a>	951
896	Alireza Salemi and Hamed Zamani. 2024. Towards	<a href="#">coding method for neural text generation</a> . In <i>Proceed-</i>	952
897	a search engine for machines: Unified ranking for	<i>ings of the 2024 Joint International Conference on</i>	953
898	multiple retrieval-augmented large language models.	<i>Computational Linguistics, Language Resources and</i>	954
899	In <i>Proceedings of the 47th International ACM SI-</i>	<i>Evaluation (LREC-COLING 2024)</i> , pages 536–557,	955
900	<i>GIR Conference on Research and Development in</i>	Torino, Italia. ELRA and ICCL.	956
901	<i>Information Retrieval</i> , pages 741–751.		

957 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,  
958 William Cohen, Ruslan Salakhutdinov, and Christo-  
959 pher D Manning. 2018. Hotpotqa: A dataset for  
960 diverse, explainable multi-hop question answering.  
961 In *Proceedings of the 2018 Conference on Empirical  
962 Methods in Natural Language Processing*, pages  
963 2369–2380.

964 Yue Yu, Wei Ping, Zihan Liu, Boxin Wang, Jiaxuan  
965 You, Chao Zhang, Mohammad Shoeybi, and Bryan  
966 Catanzaro. 2024. Rankrag: Unifying context ranking  
967 with retrieval-augmented generation in llms. In *Ad-  
968 vances in Neural Information Processing Systems*,  
969 volume 37, pages 121156–121184. Curran Asso-  
970 ciates, Inc.

971 Shaolei Zhang, Tian Yu, and Yang Feng. 2024.  
972 Truthx: Alleviating hallucinations by editing large  
973 language models in truthful space. *arXiv preprint  
974 arXiv:2402.17811*.

## 975 A Comprehensive Analysis of External 976 Knowledge Intervention in RAG

977 In this section, we expand our analysis by incorpo-  
978 rating a multi-hop question answering dataset Hot-  
979 potQA to further quantify the impact of external  
980 knowledge. As shown in Figure 5, LLaMA2-7B  
981 exhibits a consistent three-stage pattern of knowl-  
982 edge utilization across layers, as reflected by the  
983 SimHidden scores: early layers (1–14), middle  
984 layers (15–26), and deeper layers (27–32). This  
985 stratification is further supported by the DiffAttn  
986 scores, which peak at layer 21 and remain lower in  
987 both earlier and latter layers, reinforcing the validity  
988 of the three-stage division.

989 To assess the generality of this phenomenon,  
990 we evaluate a diverse set of models across vari-  
991 ous architectures and scales. Specifically, we in-  
992 clude models with different architectural designs:  
993 Mistral-7B, DeepSeek-7B, and Qwen3-8B. Results  
994 are shown in Figures 6–8. Despite these archi-  
995 tectural variations, all models exhibit a consis-  
996 tent three-phase trend in SimHidden scores and  
997 a middle-peak pattern in DiffAttn scores across  
998 both the NQ and HotpotQA datasets (e.g., for  
999 Mistral-7B on NQ, the three phases span layers  
1000 1–14, 15–29, and 30–32, with the DiffAttn peak  
1001 occurring at layer 20). This consistency extends to  
1002 larger scales, such as Llama2-13B and Qwen3-32B  
1003 (Figures 9–10).

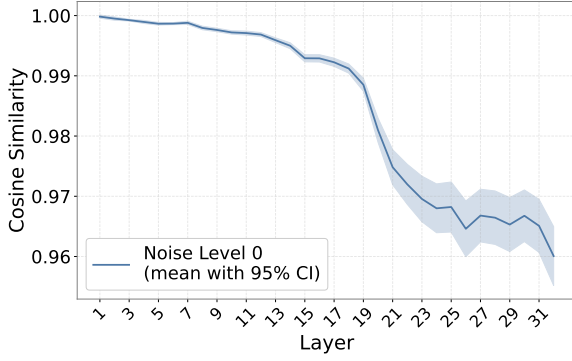
## 1004 B Details About Datasets and Metrics

1005 Following the experimental setups of (Cuconasu  
1006 et al., 2024; Wu et al., 2024; Jeong et al., 2024),  
1007 we evaluate our approach across following datasets:

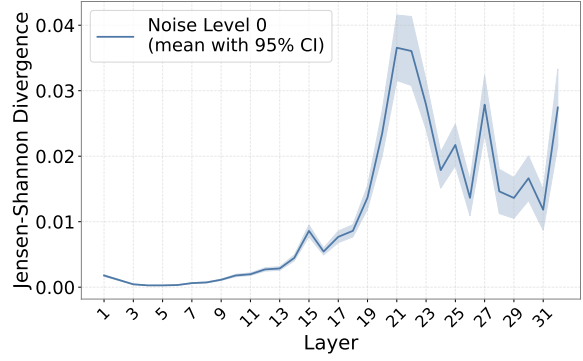
(1) **Natural Questions (NQ)** (Kwiatkowski et al., 1008  
2019), a large-scale QA dataset based on real 1009  
Google search queries. (2) **RGB** (Chen et al., 1010  
2024a), a RAG benchmark that evaluates models’ 1011  
ability to utilize retrieved information, focusing 1012  
on noise robustness, negative rejection, informa- 1013  
tion integration, and counterfactual robustness. We 1014  
use its English test set for evaluation. (3) **Hot- 1015  
potQA (HQA)** (Yang et al., 2018), which requires 1016  
multi-hop reasoning over multiple documents, fea- 1017  
turing both compare and bridge question types: 1018  
compare questions involve contrasting information 1019  
from multiple sources, while bridge questions re- 1020  
quire connecting intermediate facts, both necessi- 1021  
tating the synthesis of external documents with the 1022  
model’s internal knowledge. We evaluate all meth- 1023  
ods on their dev set, reporting results under the 1024  
categories Compare, Bridge, and Total in the main 1025  
results table. (4) **2WikiMultihopQA (2WQA)** (Ho 1026  
et al., 2020), which presents a more challenging 1027  
multi-hop QA scenario with four distinct task types: 1028  
comparison (comparing information), bridge com- 1029  
parison (connecting intermediate facts for compar- 1030  
ison), inference (deriving conclusions), and com- 1031  
positional (integrating multiple facts). We use its 1032  
dev set for evaluation, reporting separate perfor- 1033  
mance for each question type (denoted as Comapre, 1034  
Bridge, Inf, Compose and Total) in the experimen- 1035  
tal results. 1036

1037 Regarding the experimental data, since each  
1038 sample in the aforementioned benchmark datasets  
1039 is accompanied by multiple retrieved documents,  
1040 we construct the input context using these docu-  
1041 ments, supplemented with golden documents, i.e.,  
1042 the ground-truth passages that contain the infor-  
1043 mation necessary to answer the question. This  
1044 setup, following prior work (Lewis et al., 2020b;  
1045 Cuconasu et al., 2024; Krishna et al., 2025; Park  
1046 et al., 2025), allows us to evaluate the model’s abil-  
1047 ity to effectively leverage external knowledge when  
1048 it is explicitly provided in the input context. To con-  
1049 struct the noise document set for each sample, we  
1050 consider all documents within the corresponding  
1051 dataset. We utilize the all-MiniLM-L6-v2 model to  
1052 retrieve documents from this collection that exhibit  
1053 the lowest relevance scores relative to the query of  
1054 the given sample.

1055 We use accuracy as our primary evaluation met-  
1056 ric. For the NQ and RGB datasets, which include  
1057 samples with multiple acceptable answer variants  
1058 (e.g., alternative phrasings of the same concept),  
1059 we follow the evaluation protocol established in



(a) SimHidden (Smaller is better) on HotpotQA dataset.



(b) DiffAttn (Larger is better) on HotpotQA dataset.

Figure 5: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from Llama2-7B on HotpotQA dataset.

	Input Length	Latency (ms)	Throughput (tokens/s)	GPU Memory (MB)
GD (0)	239 ( $\times 1.00$ )	42.23 ( $\times 1.00$ )	24.29 ( $\times 1.00$ )	144.05 ( $\times 1.00$ )
GD (4)	958 ( $\times 4.01$ )	73.42 ( $\times 1.74$ )	14.95 ( $\times 0.62$ )	572.23 ( $\times 3.97$ )
GD (8)	1675 ( $\times 7.01$ )	97.99 ( $\times 2.32$ )	11.92 ( $\times 0.49$ )	997.81 ( $\times 6.93$ )
GD (12)	2396 ( $\times 10.02$ )	129.67 ( $\times 3.07$ )	9.31 ( $\times 0.38$ )	1426.87 ( $\times 9.91$ )
DoLA	239 ( $\times 1.00$ )	49.66 ( $\times 1.18$ )	20.58 ( $\times 0.85$ )	201.00 ( $\times 1.39$ )
LFD	239 ( $\times 1.00$ )	52.25 ( $\times 1.24$ )	19.56 ( $\times 0.81$ )	203.15 ( $\times 1.41$ )

Table 5: Input Length, Decoding Latency (ms), Throughput (tokens/s), and GPU Overhead (MB).

(Cuconasu et al., 2024; Kandpal et al., 2023; Liu et al., 2024). A model’s response is marked correct if it matches any annotated ground-truth answer.

## C Implementation details

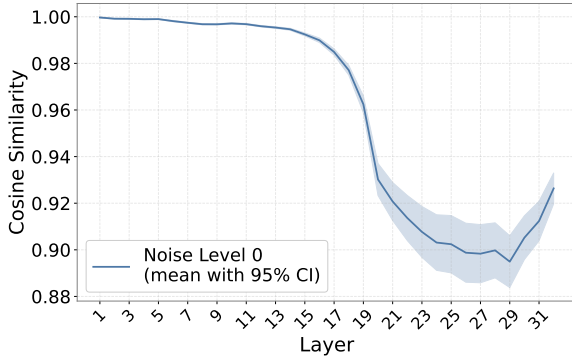
Both LFD and DoLA require the implementation of dynamic layer selection strategies. For LFD, we prioritize layers in the latter half of the architecture (e.g., layers 16–32 for 32-layer models like Llama2-7B and Mistral-7B). In contrast, DoLA selects layers across the entire depth (layers 0–32) based on divergence from the final layer’s predictions. Similar configurations apply to DeepSeek-7B (30 layers: LFD uses 15–30; DoLA uses 0–30) and Qwen3-8B (36 layers: LFD uses 18–36; DoLA uses 0–36). Following DoLA’s convention, we restrict candidates to even-numbered layers within these ranges for efficiency. Contrastive Search uses a degeneration penalty  $\alpha = 0.6$  and top-k candidate size  $k = 5$ , adopting the parameters from (Su et al., 2022). Regarding the gating mechanism in LFD, we set the hyperparameter  $\tau = 0.1$  following the configuration in (Li et al., 2022)

## D Latency, Throughput & Memory Usage

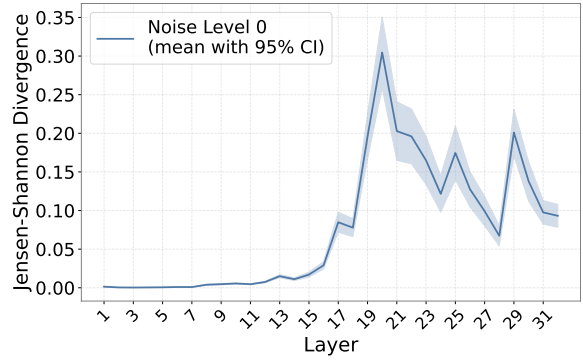
We compare the decoding latency, throughput, and GPU overhead between LFD and the greedy decoding method (with varying levels of noise), and the experimental results are illustrated in Table 5. The results demonstrate that, compared to the noise-injection baselines, LFD exhibits advantage in terms of decoding time and memory overhead. Furthermore, when compared to noise-free decoding baseline DoLA, our approach incurs just  $1.05\times$  the latency and  $1.01\times$  the memory usage, **keeping efficiency on par with state-of-the-art decoding methods.**

## E Noise Injection Analysis for CS, DoLA, and LFD

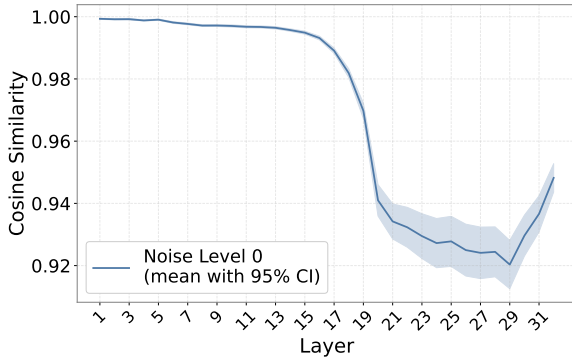
To investigate the performance of different decoding strategies under varying noise conditions, we introduce controlled noise levels (4, 8, and 12) to CS, DoLA, and LFD. Experiments are conducted using LLaMA2-7B and Mistral-7B on the NQ and HotpotQA datasets, respectively. As shown in Table 6 and Table 7, LFD generally achieves higher accuracy than other decoding methods under various noise levels, indicating comparatively stronger robustness to noise. Notably, a significant improvement in LFD’s accuracy is observed under moderate noise conditions. For example, on the NQ dataset using LLaMA2-7B, applying noise level 12 yields a 6.7% accuracy gain compared to the noise-free setting. This suggests that controlled noise exposure can further improve the performance of LFD.



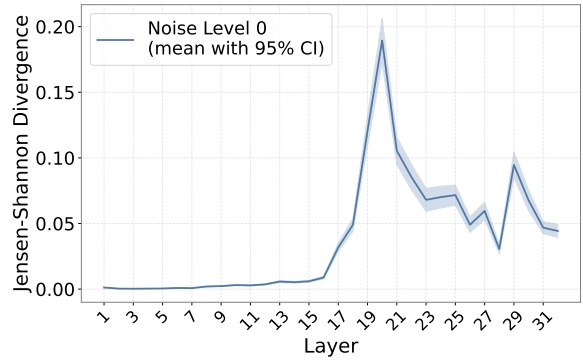
(a) SimHidden (Smaller is better) on NQ dataset.



(b) DiffAttn (Larger is better) on NQ dataset.



(c) SimHidden (Smaller is better) on HotpotQA dataset.



(d) DiffAttn (Larger is better) on HotpotQA dataset.

Figure 6: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from Mistral-7B on NQ dataset and HotpotQA dataset.

## F Ablation Study on Layer Selection Range with DeepSeek-7B and Qwen3-8B

To complement the layer selection range analysis in Section 4.3, we extend our experiments to include DeepSeek-7B and Qwen3-8B models across four datasets: NQ, RGB, HotpotQA, and 2WikiMulti-hopQA. As shown in Table 8, our findings remain consistent with those presented in Table 2, demonstrating that selecting layers from the latter half of the model consistently yields superior performance compared to earlier ranges.

## G Evaluating IKS Score Effectiveness Across Different Models and Datasets

To further validate the effectiveness of the lowest IKS layer selection discussed in Section 4.3, we conduct additional evaluations using models: Mistral-7B, DeepSeek-7B, and Qwen3-8B. Our experiments encompass both a single-hop QA dataset (NQ) and a multi-hop QA dataset (HotpotQA). The complete results are presented in Figures 11–13.

## H Case Studies

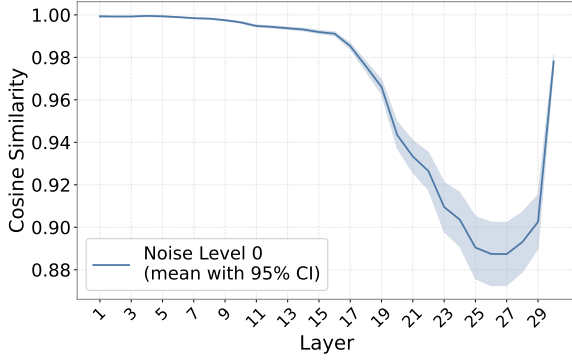
We provide some case studies using the Llama2-7B model across four benchmark datasets: NQ, RGB, HotpotQA, and 2WikiMulti-hopQA. As shown in Tables 9–12, our method enables the model to better adhere to the provided context and correctly identify answers within the given information.

## I Computational Details

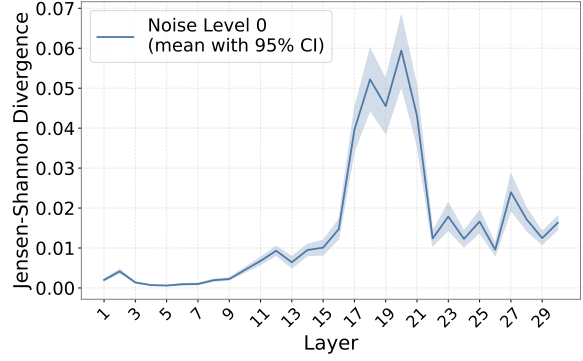
All experiments are performed on a GPU-accelerated computing system equipped with NVIDIA GeForce RTX 3090 graphics processors (24GB GDDR6X VRAM each), supported by dual Intel Xeon Gold 6271C CPUs (2.6GHz base frequency, 48 cores total) and 251GB of system memory.

## J The License For Artifacts and Data Consent

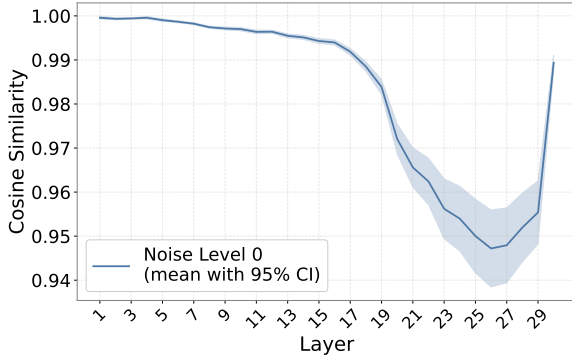
All models and datasets used in this work are publicly available artifacts released for academic research. Regarding model licenses, Llama-2



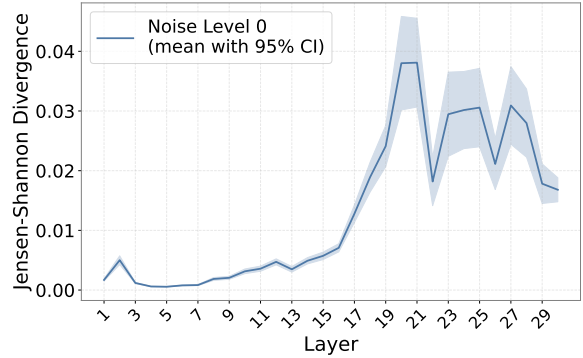
(a) SimHidden (Smaller is better) on NQ dataset.



(b) DiffAttn (Larger is better) on NQ dataset.



(c) SimHidden (Smaller is better) on HotpotQA dataset.



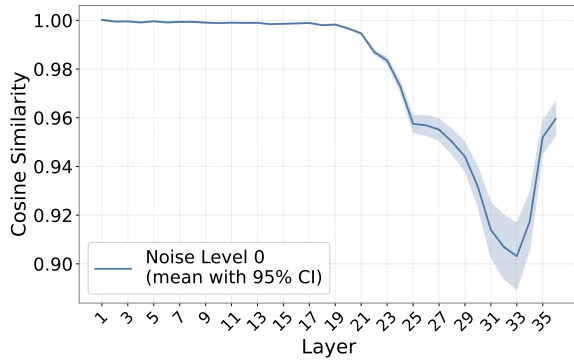
(d) DiffAttn (Larger is better) on HotpotQA dataset.

Figure 7: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from DeepSeek-7B on NQ dataset and HotpotQA dataset.

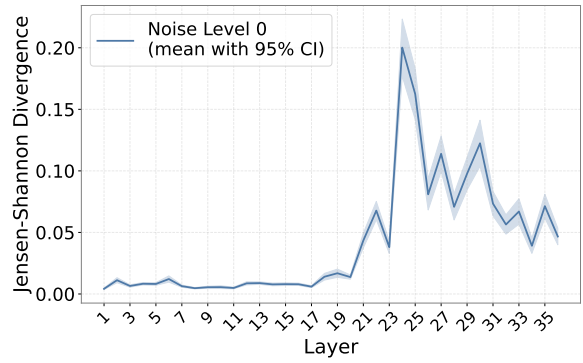
1155 (7B/13B) is used under the Llama 2 Commu-  
 1156 nity License Agreement; Mistral-7B and Qwen-3  
 1157 (8B/32B) are governed by the Apache License 2.0;  
 1158 DeepSeek-7B and GLM-4-9B are used under their  
 1159 respective proprietary model licenses (DeepSeek  
 1160 and GLM-4 licenses). For evaluation datasets, Nat-  
 1161 ural Questions (NQ) and 2WikiMultihopQA are  
 1162 licensed under the Apache License 2.0, while Hot-  
 1163 potQA and RGB are used under CC BY-SA 4.0 and  
 1164 CC BY-NC-SA 4.0, respectively. All datasets are  
 1165 pre-anonymized and filtered for offensive content  
 1166 by the original creators. Our use of these materi-  
 1167 als is strictly for academic purposes and is fully  
 1168 consistent with their specified intended use.

	NQ	HotpotQA		
		Compare	Bridge	Total
GD (0)	0.5745	<u>0.4755</u>	<u>0.4108</u>	0.4237
CS (0)	0.5775	0.4284	0.4049	0.4096
DoLA (0)	<u>0.5809</u>	0.4438	0.3883	0.3995
LFD (0)	<b>0.5949</b>	<b>0.5453</b>	<b>0.4295</b>	<b>0.4528</b>
GD (4)	0.5559	<b>0.5084</b>	0.4049	0.4257
CS (4)	<b>0.5657</b>	0.4546	<u>0.4199</u>	<u>0.4269</u>
DoLA (4)	0.5574	0.4775	0.4111	0.4244
LFD (4)	<u>0.5637</u>	<u>0.5071</u>	<b>0.4304</b>	<b>0.4458</b>
GD (8)	0.5965	<u>0.5091</u>	0.4096	0.4317
CS (8)	<u>0.6054</u>	0.4694	<u>0.4236</u>	0.4328
DoLA (8)	0.6023	0.4956	0.4199	<u>0.4351</u>
LFD (8)	<b>0.6216</b>	<b>0.5400</b>	<b>0.4400</b>	<b>0.4601</b>
GD (12)	0.6383	<u>0.5286</u>	0.4236	0.4447
CS (12)	0.6410	0.4781	<u>0.4419</u>	<u>0.4492</u>
DoLA (12)	<u>0.6458</u>	0.5158	0.4306	0.4477
LFD (12)	<b>0.6622</b>	<b>0.5521</b>	<b>0.4532</b>	<b>0.4731</b>

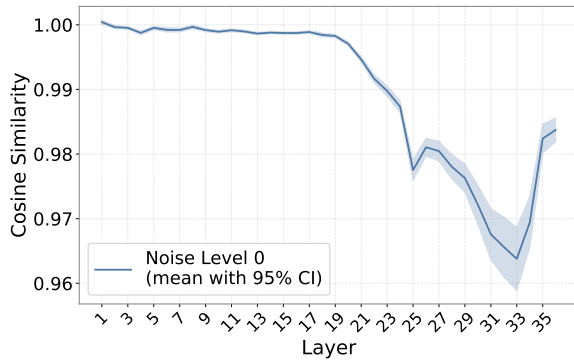
Table 6: Accuracy performance comparison of different decoding methods with varying levels of noise, evaluated on the NQ and HotpotQA datasets using the LLaMA2-7B. **Bold** values indicate the best performance, while underlined values represent the second-best.



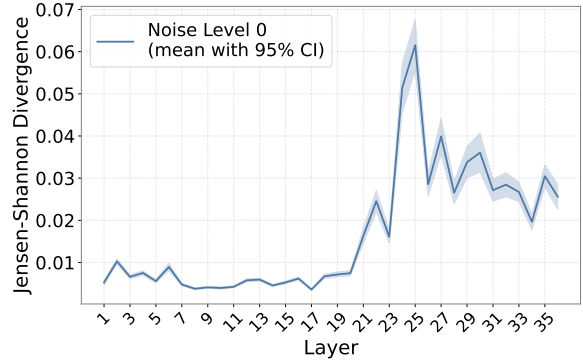
(a) SimHidden (Smaller is better) on NQ dataset.



(b) DiffAttn (Larger is better) on NQ dataset.



(c) SimHidden (Smaller is better) on HotpotQA dataset.

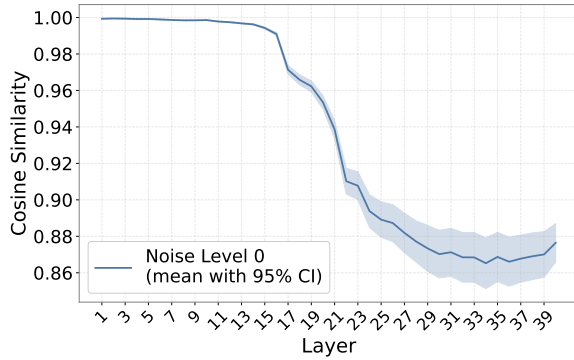


(d) DiffAttn (Larger is better) on HotpotQA dataset.

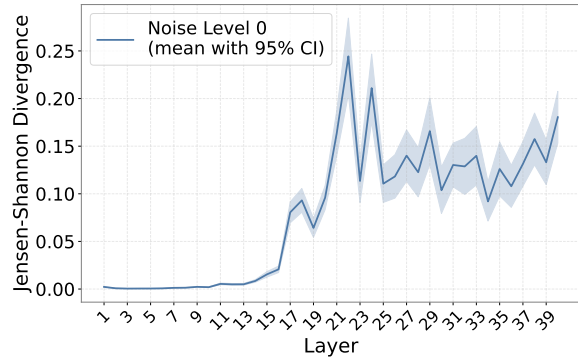
Figure 8: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from Qwen3-8B on NQ dataset and HotpotQA dataset.

	NQ	HotpotQA		
		Compare	Bridge	Total
GD (0)	0.6130	0.5864	<u>0.5889</u>	<u>0.5884</u>
CS (0)	0.5598	0.4149	0.5117	0.4922
DoLA (0)	<u>0.6142</u>	0.5931	0.5870	0.5883
LFD (0)	<b>0.6357</b>	<b>0.6026</b>	<b>0.6093</b>	<b>0.6079</b>
GD (4)	0.5667	<b>0.5494</b>	0.5348	0.5377
CS (4)	0.5014	0.2065	0.2778	0.2635
DoLA (4)	0.5625	0.5373	0.5331	0.5340
LFD (4)	<b>0.6036</b>	<u>0.5467</u>	<b>0.5502</b>	<b>0.5495</b>
GD (8)	0.5678	<b>0.5326</b>	0.5260	0.5273
CS (8)	0.5201	0.1137	0.1926	0.1768
DoLA (8)	0.5659	0.5111	0.5252	0.5223
LFD (8)	<b>0.5966</b>	<u>0.5259</u>	<b>0.5446</b>	<b>0.5409</b>
GD (12)	0.5814	<b>0.5440</b>	0.5461	0.5457
CS (12)	0.5334	0.1432	0.2328	0.2149
DoLA (12)	<u>0.5845</u>	0.5293	<u>0.5492</u>	0.5452
LFD (12)	<b>0.5943</b>	<u>0.5346</u>	<b>0.5640</b>	<b>0.5581</b>

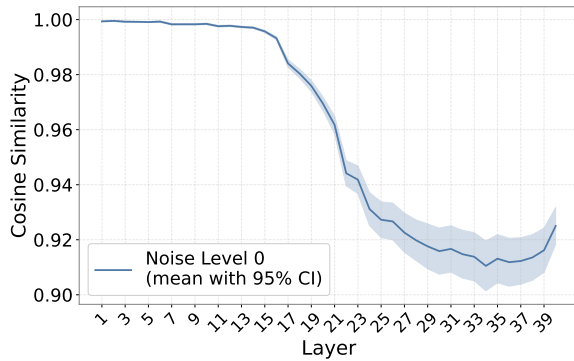
Table 7: Accuracy performance comparison of different decoding methods with varying levels of noise, evaluated on the NQ and HotpotQA datasets using the Mistral-7B. **Bold** values indicate the best performance, while underlined values represent the second-best.



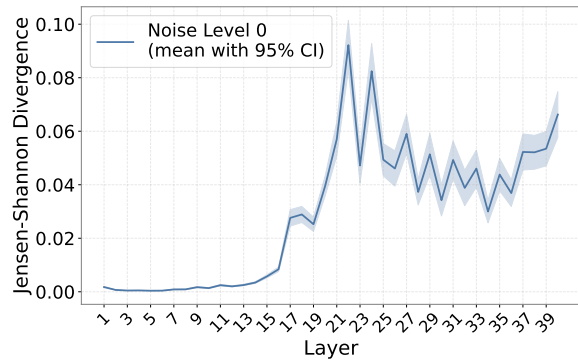
(a) SimHidden (Smaller is better) on NQ dataset.



(b) DiffAttn (Larger is better) on NQ dataset.



(c) SimHidden (Smaller is better) on HotpotQA dataset.



(d) DiffAttn (Larger is better) on HotpotQA dataset.

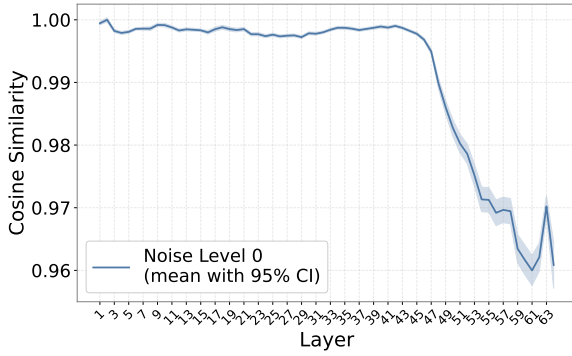
Figure 9: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from Llama2-13B on NQ dataset and HotpotQA dataset.

		NQ	RGB	HotpotQA			2WikiMultiHopQA				
				Compare	Bridge	Total	Compare	Bridge	Inf	Compose	Total
DeepSeek-7B	GD (0)	0.5250	<u>0.8233</u>	0.3282	0.4033	0.3883	0.2625	0.2414	0.2663	0.2609	0.2609
	LFD[0, 15]	0.3865	0.5000	0.2313	0.3366	0.3154	0.2155	0.2316	0.2334	0.2152	0.2211
	LFD[0, 30]	0.5389	0.8233	0.4768	0.4395	0.4470	0.4432	<b>0.4277</b>	0.4135	0.4159	0.4248
	LFD[15, 30]	<b>0.5412</b>	<b>0.8267</b>	<b>0.4801</b>	<b>0.4466</b>	<b>0.4533</b>	<b>0.4492</b>	0.4227	<b>0.4194</b>	<b>0.4223</b>	<b>0.4285</b>
Qwen3-8B	GD (0)	0.7318	<u>0.9571</u>	0.6960	0.6708	0.6759	0.6637	<u>0.6335</u>	0.5897	0.6085	0.6250
	LFD[0, 18]	0.7304	0.9567	0.7068	0.6845	0.6889	0.6640	0.6334	0.5903	0.6097	0.6256
	LFD[0, 36]	0.7372	0.9567	<b>0.7196</b>	0.6949	0.6999	0.6660	<b>0.6349</b>	0.6021	0.6146	0.6299
	LFD[18, 36]	<b>0.7380</b>	<b>0.9600</b>	0.7182	<b>0.6974</b>	<b>0.7016</b>	<b>0.6663</b>	0.6342	<b>0.6034</b>	<b>0.6148</b>	<b>0.6301</b>

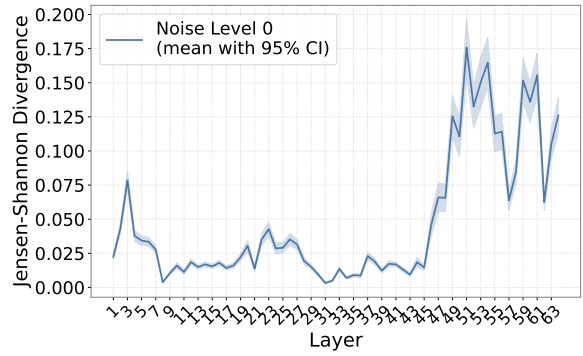
Table 8: Comparison of accuracy between different layer selection ranges under dynamic layer selection strategy. **Bold** indicate the best performance, while underline represent the second-best.

Method	GD(0)	CS	DoLA	LFD
Prompt	You are given a question and you MUST respond by EXTRACTING the answer from one of the provided documents. If none of the documents contain ... Document [20971901](Title: Pirates of the Caribbean (film series)) The film series started in 2003 with Pirates of the Caribbean : <b>The Curse of the Black Pearl</b> ... The franchise 's second film , subtitled <b>Dead Man 's Chest</b> , was released three years later in 2006 ... The third film in the series , subtitled <b>At World 's End</b> , ... and Disney released a fourth film , subtitled <b>On Stranger Tides</b> , in 2011 ...			
Ground Truth	Question: All the names of the pirate of the caribbean movies. Answer: [At World 's End, Dead Man 's Chest, On Stranger Tides]			
Answer	Captain Jack Sparrow, Will Turner, Elizabeth Swann, James Norrington, Barbossa, and Blackbeard.	Captain Jack Sparrow, Barbossa, Will Turner, Elizabeth Swann	Captain Jack Sparrow, Will Turner, Elizabeth Swann, James Norrington, Barbossa, and Blackbeard.	<b>Pirates of the Caribbean : The Curse of the Black Pearl, Dead Man's Chest, At World's End, and On Stranger Tides.</b>

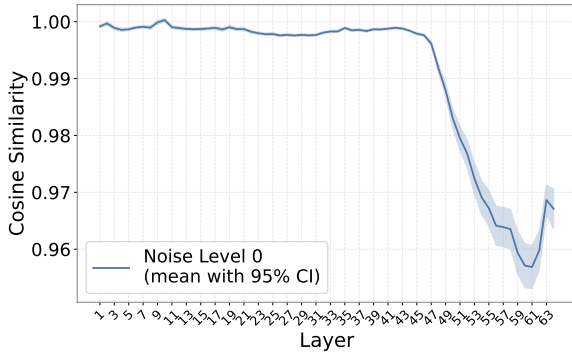
Table 9: Case study on GD(0), CS, DoLA, and LFD on the NQ dataset using LLAMA2-7B.



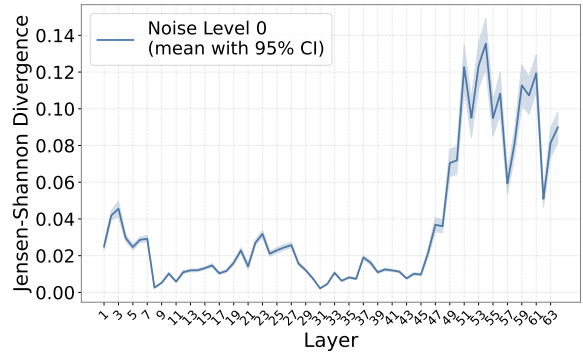
(a) SimHidden (Smaller is better) on NQ dataset.



(b) DiffAttn (Larger is better) on NQ dataset.



(c) SimHidden (Smaller is better) on HotpotQA dataset.



(d) DiffAttn (Larger is better) on HotpotQA dataset.

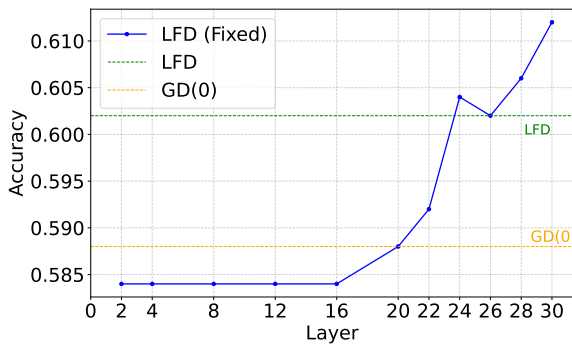
Figure 10: (a) Average SimHidden scores (with 95% confidence intervals) across layers when noise level = 0; (b) Average DiffAttn scores (with 95% confidence intervals) across layers when noise level = 0. Results are from Qwen3-32B on NQ dataset and HotpotQA dataset.

Method	GD(0)	CS	DoLA	LFD
Prompt	You are given a question and you MUST respond by EXTRACTING the answer from one of the provided documents. If none of the documents contain ... Document(Title: Riot Games Singapore ... In October 2022, Riot acquired <a href="#">Wargaming Sydney</a> —a subsidiary of Cyprus-based Wargaming that had originally developed the MMO middleware BigWorld—for an undisclosed amount, and renamed it Riot Sydney ... Question: What gaming software development studio did Riot Games acquire? Answer:			
Ground Truth	Wargaming Sydney			
Answer	Riot Sydney	RIOT SYDNEY	Riot Sydney	<b>Riot Games acquired Wargaming Sydney—a subsidiary of Cyprus-based Wargaming that had originally developed the MMO middleware BigWorld—for an undisclosed amount, and renamed it Riot Sydney.</b>

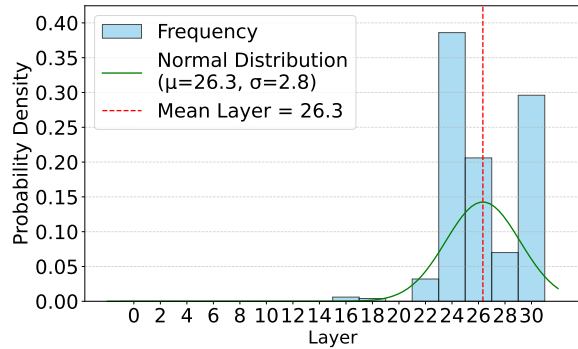
Table 10: Case study on GD(0), CS, DoLA, and LFD on the RGB dataset using LLaMA2-7B.

Method	GD(0)	CS	DoLA	LFD
Prompt	You are given a question and you MUST respond by EXTRACTING or DERIVING the answer from the provided documents. If the answer cannot be logically inferred from the documents, respond with NO-RES ... Document(Title: Terry Norris (actor)) ... As an actor, he has starred in TV Shows such as "Bellbird" & "Cop Shop", and in films like " <a href="#">Romulus, My Father</a> " and "Paper Planes" ... Document(Title: Romulus, My Father (film)) ... <a href="#">Romulus, My Father</a> is a 2007 Australian drama film directed by Richard Roxburgh ... Based on the memoir by <a href="#">Raimond Gaita</a> , the film tells the story of Romulus (Eric Bana) and his wife Christine (Franka Potente), and their struggle in the face of great adversity to raise their son, Raimond (Kodi Smit-McPhee) ... Question: Terence Richard "Terry" Norris starred in a 2007 Australian drama film that was directed by Richard Roxburgh, and was based on a memoir by who? Answer:			
Ground Truth	Raimond Gaita			
Answer	ERIC BANA	ERIC BANA	NO-RES	<b>Based on the documents provided, Terry Norris starred in the film "Romulus, My Father" directed by Richard Roxburgh, based on the memoir by Raimond Gaita.</b>

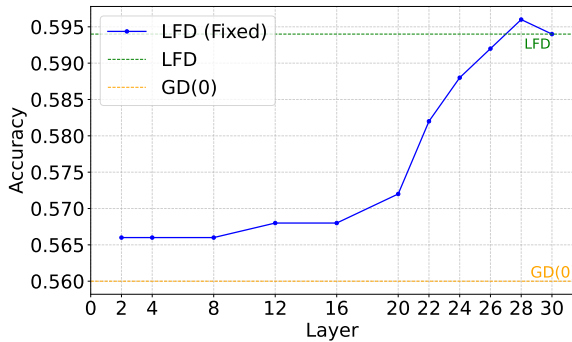
Table 11: Case study on GD(0), CS, DoLA, and LFD on the HotpotQA dataset using LLaMA2-7B.



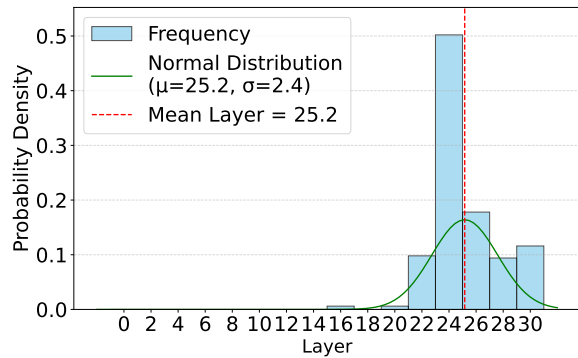
(a) LFD vs. LFD (Fixed) on NQ dataset.



(b) Histogram of layer selection in LFD on NQ dataset.



(c) LFD vs. LFD (Fixed) on HotpotQA dataset.

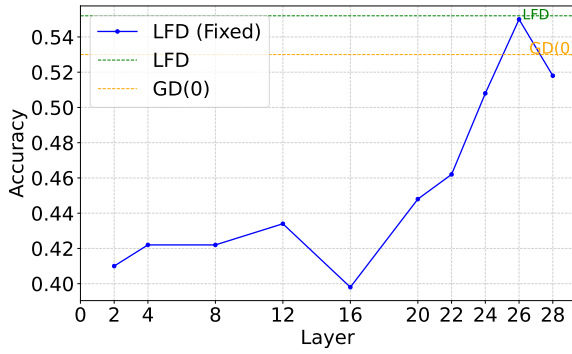


(d) Histogram of layer selection in LFD on HotpotQA dataset.

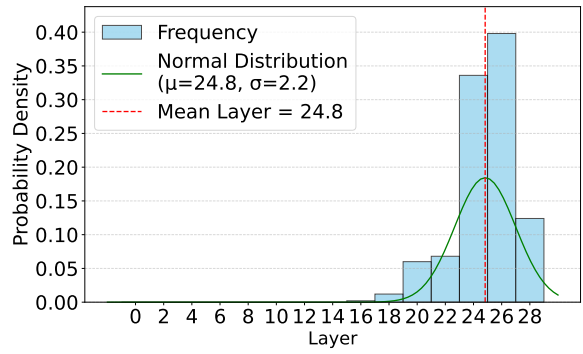
Figure 11: Comparison between LFD and LFD (Fixed) using the Mistral-7B model on the NQ and HotpotQA datasets.

Method	GD(0)	CS	DoLA	LFD
Prompt	You are given a question and you MUST respond by EXTRACTING or DERIVING the answer from the provided documents. If the answer cannot be logically inferred from the documents, respond with NO-RES ... Document(Title: A Tale of Winter) A Tale of Winter is a 1992 French drama film directed by <b>Éric Rohmer</b> , and starring Charlotte Véry, Frédéric van den Driessche and Michael Voletti ... Document(Title: Éric Rohmer) ... He edited the influential film journal, " <b>Cahiers du cinéma</b> ", from 1957 to 1963 ... Question: Where does the director of film A Tale Of Winter work at? Answer:			
Ground Truth	Cahiers du cinéma			
Answer	Éric Rohmer works at a film production company.	Éric Rohmer	NO-RES: The answer cannot be logically inferred from the provided documents.	<b>Éric Rohmer works at a film journal called "Cahiers du cinéma".</b>

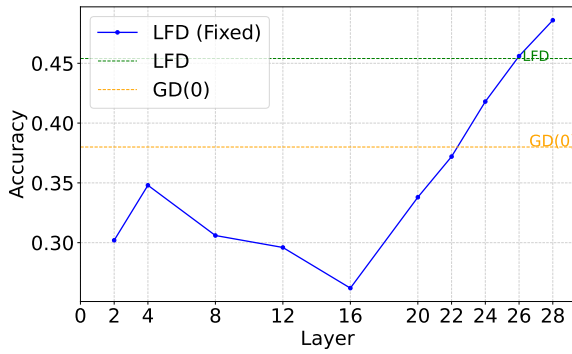
Table 12: Case study on GD(0), CS, DoLA, and LFD on the 2WikiMultihopQA dataset using LLaMA2-7B.



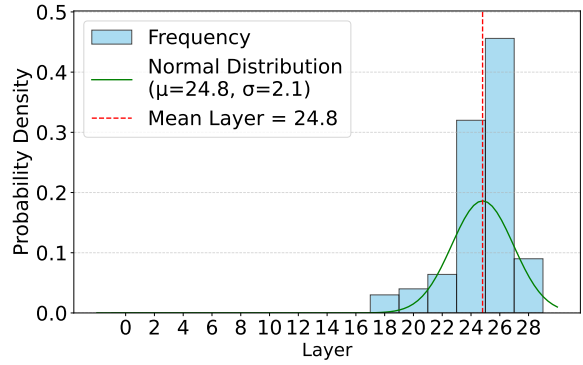
(a) LFD vs. LFD (Fixed) on NQ dataset.



(b) Histogram of layer selection in LFD on NQ dataset.

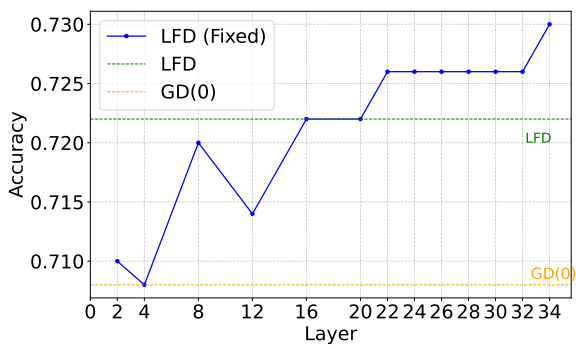


(c) LFD vs. LFD (Fixed) on HotpotQA dataset.

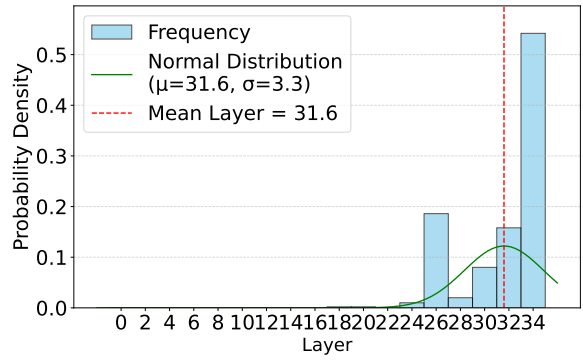


(d) Histogram of layer selection in LFD on HotpotQA dataset.

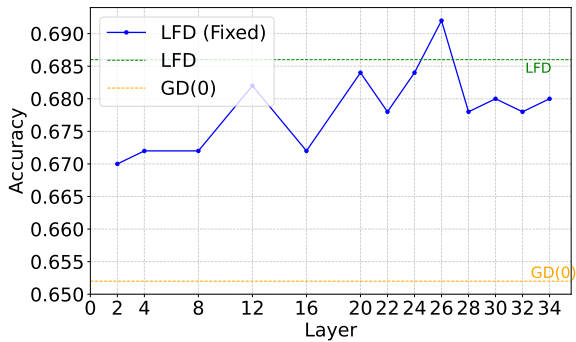
Figure 12: Comparison between LFD and LFD (Fixed) using the DeepSeek-7B model on the NQ and HotpotQA datasets.



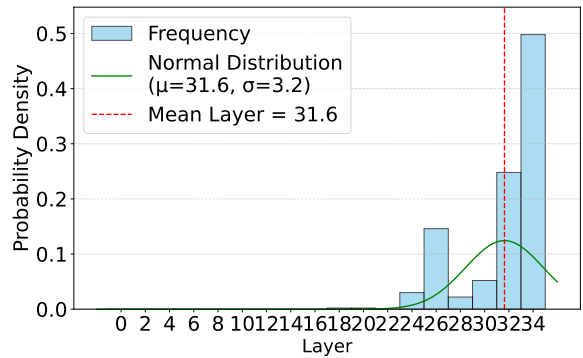
(a) LFD vs. LFD (Fixed) on NQ dataset.



(b) Histogram of layer selection in LFD on NQ dataset.



(c) LFD vs. LFD (Fixed) on HotpotQA dataset.



(d) Histogram of layer selection in LFD on HotpotQA dataset.

Figure 13: Comparison between LFD and LFD (Fixed) using the Qwen3-8B model on the NQ and HotpotQA datasets.