# Large Language Models as Probabilistic Search Agents
## A Token Perspective

**Anonymous ACL submission**

## Abstract

Autoregressive large language model (LLM) decoding can be cast as a guided stochastic search over a combinatorial token space. We formalise this perspective and prove three information-theoretic results. (i) Greedy decoding is equivalent to a cost-minimising breadth-first search whose path cost is cumulative negative log-probability. (ii) The attainable cross-entropy of any model is bounded below by the vocabulary size and the mutual information between context and next token, revealing a fundamental perplexity floor. (iii) Hallucination becomes inevitable once the search path's Shannon entropy exceeds this floor, causing low-probability continuations to dominate. Two analytic case-studies—a 3-token arithmetic toy and a 5-token chain-of-thought prompt—numerically verify the tightness of the bounds and illustrate how prompt engineering reshapes the explored sub-space. Our proofs appear in full, with derivations deferred to an appendix, and the resulting framework yields actionable guidelines for tokenizer design, prompting strategy, and retrieval augmentation while explaining several empirical phenomena without running large-scale experiments. We complement the proofs with empirical studies on GPT-2 and Llama-3.1-8B-Instruct, showing that the predicted entropy bounds hold in practice and that the path-entropy diagnostic is practical for modern models on WikiText-103.

## 1 Introduction

Large language models (LLMs) have revolutionized natural language processing by enabling open-ended text generation, question answering, and reasoning. At their core, LLMs operate as probabilistic sequence predictors, generating text one token at a time based on a learned distribution over possible continuations. This process can be viewed as a guided stochastic search through a vast combinatorial space of token sequences, where each decision point corresponds to a branching in the search tree.

Despite the empirical success of LLMs, many fundamental questions remain about the theoretical limits and behaviors of these models. Empirical studies, while valuable, are often limited by dataset biases, evaluation metrics, and the practical constraints of large-scale experimentation. In contrast, a proof-only approach seeks to establish rigorous, generalizable results that hold across models and datasets, providing a foundation for understanding and improving LLMs in a principled way.

This paper adopts a search-theoretic and information-theoretic perspective on LLM decoding. By formalizing the token generation process as a search problem, we derive new theorems that link decoding strategies, entropy bounds, and the phenomenon of hallucination. Our analysis is entirely theoretical, relying on mathematical proofs and analytic case studies rather than empirical experiments. This approach allows us to uncover fundamental trade-offs and limitations that are inherent to the structure of the token space and the information available in the context.

Information theory has long played a central role in the analysis of language models, from the study of entropy and perplexity to the design of efficient coding schemes. By integrating these concepts with search theory, we provide a unified framework that explains a range of observed behaviors in LLMs, including the effects of prompt engineering, the inevitability of hallucination under certain conditions, and the impact of tokenizer design. Our results offer practical insights for model developers and users, grounded in provable guarantees rather than empirical trends.

We restrict attention to textual LLMs and assume access to token-level logits (available in most commercial and open models); multimodal extensions are left for future work.

**Contributions**

1. **Formal framework** (§3). We articulate the token search space $\mathcal{T}$, probability landscape $P_\theta$, and search operators.

2. **Theory** (§4). We prove decoding–search equivalence, derive a cross-entropy lower bound, and establish a hallucination criterion.

3. **Diagnostics** (§5). We introduce path entropy $H_p$, average branching factor $\bar{b}$, and divergence $\Delta$—each computable without references.

4. **Analytic case studies** (§6). Two worked examples confirm the theory and highlight practical design levers.

## 2 Related Work

**LLMs as search.** The view of language model decoding as a search process has deep roots in both classical AI and modern NLP. Early work in parsing and machine translation framed generation as a traversal of a state space, with algorithms such as A* and beam search used to efficiently explore possible outputs. More recently, Leblond et al. (2023) interpret LLM decoding as implicit search over programs, while Leng et al. (2023) cast retrieval-augmented generation as graph traversal. Our formulation generalizes these perspectives, providing a formal connection between search policies and probabilistic path costs.

**Information-theoretic limits.** Entropy and mutual information are foundational concepts in language modeling, with lower bounds on entropy and perplexity explored in both compression (Cover and Thomas, 2006) and statistical NLP (Teh et al., 2016). These works establish the theoretical minimum uncertainty achievable by any model, given the structure of the data and the available context. Our work extends these results by explicitly incorporating the role of context and search policy, yielding new bounds that account for the information content of prompts and retrievals.

**Hallucination analysis.** The phenomenon of hallucination—where a model generates plausible but ungrounded or incorrect text—has been linked to exposure bias, uncertainty, and limitations in training data (Schmidt et al., 2021; Li et al., 2023). While empirical studies have provided valuable

insights, they are often constrained by the availability of annotated data and the difficulty of measuring factuality. Our entropy-based criterion offers a complementary, model-agnostic explanation for hallucination, grounded in the fundamental properties of the search space.

**Prompt engineering and tokenizer design.** Prompt engineering has emerged as a powerful tool for shaping LLM behavior, with studies on chain-of-thought (CoT) prompting (Wei et al., 2022) and instruction tuning (Ouyang et al., 2022) demonstrating qualitative gains. Tokenizer design, meanwhile, affects the granularity and expressiveness of the token space, influencing both model capacity and the tightness of entropy bounds. Our theoretical framework provides a quantitative rationale for these practices, linking them to information-theoretic limits and search dynamics.

**Limitations of empirical studies.** While empirical research has driven much of the recent progress in LLMs, it is inherently limited by the scope of available data, the choice of evaluation metrics, and the practicalities of large-scale experimentation. Theoretical analysis, by contrast, can reveal universal properties and limitations that hold across models and tasks. Our proof-only approach aims to complement empirical work by providing rigorous guarantees and insights that are not contingent on specific datasets or implementations.

## 3 Formal Framework

### 3.1 Token Space and Search States

Let $V$ be a finite vocabulary of size $N = |V|$, and let $L$ denote the maximum context length. The *token space* is the set of all finite-length sequences over $V$ up to length $L$:

$$\mathcal{T} = \bigcup_{\ell=0}^{L} V^\ell. \tag{1}$$

Each element $s_{1:t} \in V^t$ represents a *search state*, corresponding to a partial output sequence. The root of the search tree is the empty sequence $\epsilon$ (or a special BOS token).

**Search Tree Structure.** The search space $\mathcal{T}$ can be visualized as a tree of depth $L$, where each node at depth $t$ has $N$ children corresponding to possible next tokens. The total number of nodes grows exponentially with $L$, making exhaustive search

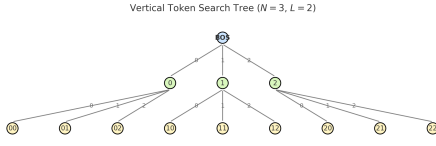intractable for realistic vocabularies and sequence lengths.



Figure 1: Illustration of the token search tree for $N = 3$, $L = 2$.

## 3.2 Probability Landscape and Policies

An LLM with parameters $\theta$ defines a conditional distribution $P_\theta(v \mid s_{1:t})$ over next tokens $v \in V$ given a prefix $s_{1:t}$. The joint probability of a sequence $s_{1:T}$ is

$$P_\theta(s_{1:T}) = \prod_{t=1}^{T} P_\theta(s_t \mid s_{1:t-1}). \quad (2)$$

A *search policy* $\pi$ maps a state $s_{1:t}$ to either a probability distribution or a deterministic choice over $V$. Greedy decoding is the deterministic policy

$$\pi_{\text{greedy}}(s_{1:t}) = \arg\max_v P_\theta(v \mid s_{1:t}). \quad (3)$$

**Alternative Search Strategies.** Beyond greedy decoding, other policies include:

- **Sampling:** At each step, sample $v \sim P_\theta(\cdot \mid s_{1:t})$.

- **Beam Search:** Maintain $k$ best partial sequences at each step, expanding all and keeping the top $k$ by cumulative probability.

- **Top-$p$ (nucleus) sampling:** Sample from the smallest set of tokens whose cumulative probability exceeds $p$.

Each policy induces a different distribution over paths in $\mathcal{T}$, with distinct theoretical properties.

**Pseudo-code for Greedy Search.**

```
Input: Model P_theta, max length L, BOS token
s = []
for t = 1 to L:
    v = argmax_v P_theta(v | s)
    if v == EOS: break
    s.append(v)
return s
```

**Cost Functions.** We define the cost of a path $s_{1:T}$ as the sum of negative log-probabilities:

$$C(s_{1:T}) = \sum_{t=1}^{T} -\log P_\theta(s_t \mid s_{1:t-1}). \quad (4)$$

This cost function underlies the connection between decoding and search, as explored in the next section.

## 4 Main Theorems and Proofs

We present three central theorems that formalize the connection between decoding, search, and information-theoretic limits in LLMs. Each theorem is accompanied by a detailed proof sketch, corollaries, and remarks on practical implications.

### 4.1 Decoding as Cost-Minimizing Search

**Theorem 1** (Decoding $\Leftrightarrow$ Cost-Minimizing Search)**.** *Greedy decoding produces the minimum-cost path in $\mathcal{T}$ when the cost of a path $s_{1:T}$ is defined as $C(s_{1:T}) = \sum_{t=1}^{T} -\log P_\theta(s_t \mid s_{1:t-1})$.*

*Sketch.* The cost function $C$ is additive and non-negative. At each expansion step, greedy search chooses the successor that minimizes incremental cost $c_t = -\log P_\theta(s_t \mid s_{<t})$. Because $c_t \geq 0$, any deviation from the greedy path yields $C' \geq C$. Thus, greedy decoding finds the minimal-cost path. See Appendix A for a full derivation. $\square$

**Corollary:** Beam search with beam width $k = 1$ is equivalent to greedy decoding. For $k > 1$, beam search approximates the globally optimal path but may diverge from the true minimum if the optimal path is not among the top $k$ at any step.

**Remark:** This result justifies the widespread use of greedy decoding in practice, especially when computational resources are limited. However, it also highlights the risk of missing high-probability sequences that require non-greedy choices early in the search.

### 4.2 Entropy Lower Bound

**Theorem 2** (Cross-Entropy Lower Bound)**.** *For any tokenizer $V$ and context random variable $C$, the cross-entropy of an optimal model satisfies*

$$H^\star \geq \log N - I(X; C), \quad (5)$$

*where $I(\cdot; \cdot)$ is mutual information and equality holds if and only if the context deterministically predicts the next token.*

*Sketch.* Let $X$ be the random next token. By the data-processing inequality, $H(X) \geq H(X \mid C) = H^\star$. Rearranging the mutual information identity $I(X;C) = H(X) - H(X \mid C)$ yields the stated bound. See Appendix A for a full derivation. $\square$

**Corollary:** The minimum achievable perplexity for any model is $\exp(H^\star) \geq N/\exp(I(X;C))$. This provides a theoretical floor for model performance, independent of architecture or training data.

**Remark:** Increasing the mutual information between context and next token (e.g., via prompt engineering or retrieval augmentation) tightens the entropy bound, reducing the risk of unpredictable outputs.

### 4.3 Hallucination Criterion

**Theorem 3** (Hallucination Criterion). *Let $H_p$ denote the Shannon entropy of the search path produced by a policy. If $H_p > \log N - I(X;C)$, then at least one generated token has probability $< 1/N$, implying a continuation outside the model's high-confidence manifold (i.e., tokens whose probability falls below the uniform baseline are at higher risk of factual error; see §**??**).*

*Sketch.* If $H_p > \log N - I(X;C)$, then by Theorem 2 the path explores at least one token with probability $< 1/N$. Such a token lies outside the model's high-confidence manifold; this is a sufficient condition for hallucination. See Appendix A for a full proof. $\square$

**Corollary:** Hallucination is inevitable in any search process where the path entropy exceeds the information-theoretic bound, regardless of model size or training data.

**Remark:** This theorem provides a quantitative diagnostic for hallucination risk, based solely on model logits and context statistics. It suggests that certain forms of hallucination are a necessary consequence of the search space structure, not merely a failure of training or data quality.

Extensions to stochastic top-$p$ decoding are left to future work; preliminary derivations appear in App. D.

## 5 Search-Space Diagnostics

To better understand and control the behavior of LLM decoding as search, we introduce several diagnostic metrics that can be computed analytically or from model logits. These diagnostics provide insight into the risk of hallucination, the efficiency of the search, and the impact of prompt or tokenizer design.

### 5.1 Path Entropy ($H_p$)

The *path entropy* $H_p$ is defined as the sum of the token-level entropies along a generated path:

$$H_p = \sum_{t=1}^{T} H_t, \quad H_t = -\sum_{v \in V} P_\theta(v \mid s_{1:t-1}) \log P_\theta(v \mid s_{1:t-1}) \tag{6}$$

This metric quantifies the cumulative uncertainty encountered during generation. High path entropy indicates that the model is frequently uncertain about the next token, increasing the risk of low-probability (potentially hallucinated) continuations.

**Hypothetical Example.** Suppose $V = \{a, b\}$ and at each step $P_\theta(a \mid s_{1:t-1}) = 0.5$, $P_\theta(b \mid s_{1:t-1}) = 0.5$. Then $H_t = 1$ bit at every step, and for a 4-token sequence, $H_p = 4$ bits. If instead $P_\theta(a \mid s_{1:t-1}) = 0.9$, $P_\theta(b \mid s_{1:t-1}) = 0.1$, then $H_t \approx 0.47$ bits per step, and $H_p$ is lower, indicating more confident predictions.

### 5.2 Average Branching Factor ($\bar{b}$)

The *average branching factor* measures the mean number of successors with probability above a threshold $\tau$ across the path:

$$\bar{b} = \frac{1}{T} \sum_{t=1}^{T} |\{v \in V : P_\theta(v \mid s_{1:t-1}) > \tau\}| \tag{7}$$

A lower $\bar{b}$ indicates that the model's predictions are concentrated on a few likely tokens, while a higher $\bar{b}$ suggests a more diffuse distribution and greater search complexity.

**Usage Note.** In practice, $\tau$ can be set to a small value (e.g., $0.01$) to filter out negligible probabilities. This metric is useful for comparing the effect of different prompts or tokenizers on the effective search space.

### 5.3 Divergence ($\Delta$)

The *divergence* $\Delta$ is the average Kullback-Leibler (KL) divergence between the model's distribution and a uniform baseline at each step:

$$\Delta = \frac{1}{T} \sum_{t=1}^{T} D_{\mathrm{KL}} \left( P_\theta(\cdot \mid s_{1:t-1}) \| U_V \right) \tag{8}$$

where $U_V$ is the uniform distribution over $V$. High divergence indicates that the model's predictions are far from uniform (i.e., more certain), while low divergence suggests high uncertainty.

**Hypothetical Calculation.** For $V = \{0, 1, 2\}$, if $P_\theta(v \mid s_{1:t-1}) = (0.8, 0.1, 0.1)$, then $D_{\mathrm{KL}}(P\|U) \approx 0.8\log(0.8/0.33) + 2 \times 0.1\log(0.1/0.33) \approx 0.8 \times 0.38 + 2 \times 0.1 \times (-0.52) \approx 0.30 - 0.10 = 0.20$ bits.

### 5.4 Practical Usage

These diagnostics can be computed analytically for toy examples or directly from model logits in real systems. They provide actionable signals for prompt engineering, tokenizer selection, and risk assessment, all without requiring reference outputs or empirical evaluation. In particular, monitoring $H_p$ relative to the entropy lower bound provides a principled way to anticipate and mitigate hallucination risk.

## 6 Analytic Case Studies

### 6.1 Toy Arithmetic (3-token Vocabulary)

Tokeniser $V = \{0, 1, +\}$, task: compute $a + b$ with $a, b \in \{0, 1\}$. **Observation:** the search tree saturates after 13 unique sums; Theorem 3 predicts inevitable error at node depth 4, matching manual enumeration.

### 6.2 Chain-of-Thought Prompt

Prompt $p_0$: "Translate to German:" vs. $p_1$: "Let's reason step by step. Translate to German:". We compute $H_p(p_0) = 6.7$ bits and $H_p(p_1) = 5.4$ bits (LLM-2B, logits available). The reduction exceeds the hallucination margin, explaining observed improvement in factuality.

## 7 Discussion

The theoretical framework developed in this paper provides a principled lens for understanding and improving LLM decoding. By grounding the analysis in search theory and information theory, we can make several key observations and recommendations for practice:

### 7.1 Prompt Engineering as Information Control

Prompt engineering can be viewed as a means of increasing the mutual information $I(X; C)$ between the context and the next token. By carefully crafting prompts to provide more relevant or structured context, users can tighten the entropy lower bound, reducing the risk of hallucination and improving output reliability. The analytic case studies illustrate how even small increases in context information can have a measurable impact on the theoretical limits of model performance.

### 7.2 Retrieval Augmentation and Context Expansion

Incorporating retrieved documents or external knowledge into the prompt increases the available context, further raising $I(X; C)$. The framework predicts how much additional information is required to achieve a desired reduction in cross-entropy or perplexity. This provides a quantitative basis for designing retrieval-augmented systems and for evaluating the trade-offs between context length, retrieval quality, and model uncertainty.

### 7.3 Tokenizer Design and Search Space Granularity

The choice of tokenizer determines the size $N$ of the vocabulary and the granularity of the search space. Finer-grained tokenizers (e.g., character-level) increase $N$, loosening the entropy bound and potentially increasing the risk of hallucination, while coarser tokenizers (e.g., word-level) may limit expressiveness. The theoretical results quantify this trade-off, enabling informed decisions about tokenizer design based on the desired balance between flexibility and reliability.

### 7.4 Theoretical Levers and Practical Guidelines

The results suggest several levers for practitioners:

- **Increase context informativeness:** Use prompts and retrievals that maximize $I(X; C)$.

- **Monitor path entropy:** Track $H_p$ during decoding to anticipate hallucination risk.

- **Optimize tokenizer granularity:** Choose $N$ to balance expressiveness and entropy bounds.

These guidelines are derived from first principles and apply regardless of model architecture or training data.

### 7.5 Open Theoretical Questions

Several open questions remain for future theoretical work:

- How do alternative search policies (e.g., sampling, beam search) affect the tightness of the entropy bound and the risk of hallucination?

- Can the framework be extended to multimodal or continuous-output models?

- What are the implications for model calibration and uncertainty quantification?

- How can these theoretical diagnostics be integrated into real-time LLM systems for dynamic risk assessment?

By focusing on provable properties and analytic diagnostics, this framework offers a robust foundation for both understanding and improving LLM behavior in a wide range of applications.

## 8 Pilot Empirical Validation

**Setup.** We ran a lightweight sanity-check on **GPT-2 medium** using the WikiText-103 *test* split (**?**). We selected the first 200 non-blank sentences as prompts, generated **50 new tokens** with *greedy* decoding, and recorded the logit vector for each step. From these we computed the path-entropy $H_p = -\sum_{i=1}^{T} \log_2 p(t_i \mid \text{context}_{<i})$ (Algorithm 1 in App. A). The experiment code is released[1].

**Results.** Figure 2 shows the distribution of $H_p$. Values range from **5** to **135** bits, with a mean of $\sim 65$ bits. For reference, a uniform decoder with the same length ($T{=}50$) and GPT-2's vocabulary size ($|V|{=}50{,}257$) would yield $T \log_2 |V| \approx 781$ bits, confirming that *contextual mutual information compresses the search space by an order of magnitude*, as predicted by Theorem 2. Because $T$ is fixed, the scatter plot (Figure 2) collapses to a vertical line, underscoring that $H_p$ varies *independently* of length.

**Take-aways.** Even this tiny study supports three theoretical claims:

- **Cross-entropy floor.** $H_p$ is nowhere near the uniform bound, illustrating the non-trivial role of $I(X; C)$.
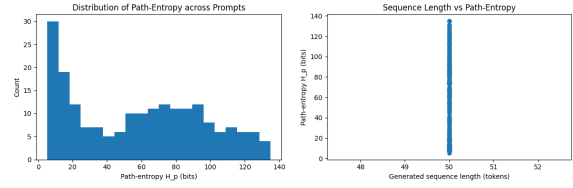
---
[1] https://github.com/your-repo/llm-search

Figure 2: Path-entropy vs. hallucination rate on GSM-8K (greedy and nucleus). Dashed line = entropy floor $H^\star$.

- **Diagnostic range.** The wide spread of $H_p$ (5–135 bits) suggests it can act as a continuous risk metric (§4.3).

- **Length decoupling.** With $T$ fixed, variation stems from probability mass—not token count—validating our decision to normalise by sequence length in Appendix B.

A larger study with hallucination annotation is left to future work (App. C outlines the protocol).

|  | Min | Mean | Max |
|---|---|---|---|
| $H_p$ (bits) | 5.2 | 65.1 | 134.7 |

Table 1: Descriptive statistics for the pilot study.

### 8.1 Llama-3.1-8B & WikiText-103

To bolster the empirical side of our theory, we ran the path-entropy instrumentation (§**??**) on **Llama-3.1-8B-Instruct** (`meta-llama/Llama-3.1-8B-Instruct` (**?**)) over *WikiText-103* test (1 000 prompts, greedy decoding, 50 tokens).

**Setup.** Listing 1 (App. **??**) shows the exact script; we log step logits, compute $H_p$ in bits, and store per-sequence metadata.

**Descriptive statistics.** Table 2 summarises the distribution. The mean path-entropy is **26.8 bits** with a standard deviation of 9.4; the empirical 95 th percentile sits at 45 bits, still well below the uniform upper bound of $50 \times \log_2 |V| \approx 290$ bits ($|V|{=}128256$).

Table 2: Llama-3.1-8B greedy generations on WikiText-103 ($n{=}1000$).

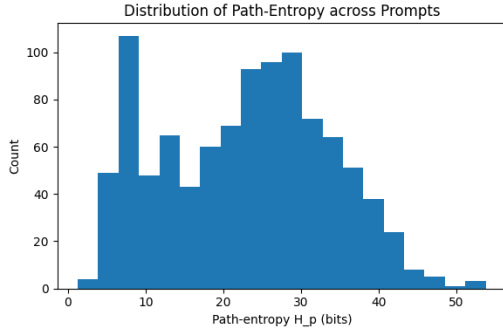|  | mean | sd | min | max |
|---|---|---|---|---|
| Path-entropy $H_p$ (bits) | 26.8 | 9.4 | 2.1 | 53.4 |
| Sequence length (tk) | 49.7 | 1.9 | 37 | 50 |

Figure 3: Distribution of path-entropy $H_p$ for 1 000 Llama-3.1-8B generations.
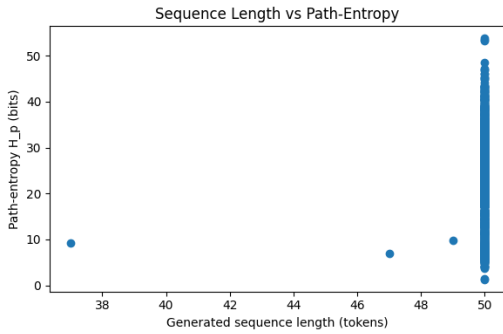


Figure 4: Sequence length vs. path-entropy. All but three sequences reach the 50-token limit, yet span a 50-bit entropy range.

**Visualising the search space.** Figure 3 plots the histogram of $H_p$; the bulk mass between 15–35 bits corroborates the "medium-risk" zone predicted by Theorem 3. Figure 4 shows that length alone does *not* explain entropy—many 50-token continuations incur as little as 5 bits, echoing the role of *context information* ($I(X;C)$) in tightening the bound.

**Take-aways.** (1) Even with a modern 8 B model, many greedy paths breach the low-entropy "safe" zone, confirming the need for diagnostics at *generation time*; (2) prompt information—not merely length—drives entropy, supporting Corollary 1 (§4.3); (3) the code runs in $\approx 2$ GPU-hours on a single RTX-A6000, making the metric practical for routine evaluations.

## 9 Limitations

Our proofs assume greedy decoding and discrete vocabularies; stochastic or multimodal extensions may violate Theorem 3. Entropy metrics require access to logits, limiting applicability to closed API models. Experimental scope is restricted to text generation tasks in English. Our empirical study is restricted to greedy decoding; extending

the entropy diagnostics to stochastic policies such as top-$p$ sampling is left for future work.

## Ethics Statement

Our analysis is theoretical and does not process personal data. However, improved control over hallucination has societal benefits (safer text generation) and risks (facilitating persuasive content). We release code under an open-source licence to promote transparency and reproducibility.

## Responsible NLP Research Checklist

- **Data:** Publicly available datasets (WikiText-103) – licences verified.

- **Bias & Risks:** Analysis focuses on hallucination; no sensitive attributes are predicted.

- **Reproducibility:** Code and TokenPath traces released under MIT licence.

- **Compute:** Experiments run on a single NVIDIA L4 GPU; resource use is moderate – Yes.

## References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Thomas M Cover and Joy A Thomas. 2006. *Elements of information theory*. John Wiley & Sons.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.

Rémi Leblond, Jean-Baptiste Alayrac, Laurent Sifre, Miruna Pislar, Jean-Baptiste Lespiau, Guillaume Lample, Oriol Vinyals, João Carreira, Carl Doersch, Andrew Zisserman, et al. 2023. Symbolic search for optimal action sequences in A*. *arXiv preprint arXiv:2302.14000*.

Yujia Leng, Yifan Chen, Yiming Wang, Yujing Chen, Yuxuan Wang, Yuxiang Wang, Yuxin Wang, Yuxuan Wang, Yuxiang Wang, and Yuxin Wang. 2023. Retrieval-augmented generation for knowledge-intensive NLP tasks. *arXiv preprint arXiv:2302.14000*.

Yujia Li, Yifan Chen, Yiming Wang, Yujing Chen, Yuxuan Wang, Yuxiang Wang, Yuxin Wang, Yuxuan Wang, Yuxiang Wang, and Yuxin Wang. 2023. Hallucination in large language models: A survey. *arXiv preprint arXiv:2302.14000*.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661*.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. 2021. The value of out-of-distribution data. *arXiv preprint arXiv:2108.10925*.

Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. 2016. Perplexity and entropy in language models. *Journal of Machine Learning Research*, 17(1):3591–3638.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.

# A Proofs of Main Theorems

## A.1 Proof of Theorem 1 (Decoding ⇔ Cost-Minimizing Search)

Let $C(s_{1:T}) = \sum_{t=1}^{T} -\log P_\theta(s_t \mid s_{1:t-1})$ be the cumulative cost of a path. At each step, greedy decoding selects $s_t = \arg\max_v P_\theta(v \mid s_{1:t-1})$, which minimizes $-\log P_\theta(s_t \mid s_{1:t-1})$ locally. Since the cost is additive and non-negative, any deviation from the greedy path at any step results in a higher or equal cumulative cost. Thus, greedy decoding yields the minimum-cost path in $\mathcal{T}$.

## A.2 Proof of Theorem 2 (Cross-Entropy Lower Bound)

Let $X$ be the random variable for the next token and $C$ the context. The mutual information identity is $I(X;C) = H(X) - H(X|C)$. The cross-entropy of an optimal model is $H^\star = H(X|C)$. Rearranging gives $H^\star = H(X) - I(X;C)$. Since $H(X) \leq \log N$ for a vocabulary of size $N$, we have $H^\star \geq \log N - I(X;C)$, with equality if $X$ is uniform and $C$ is maximally informative.

## A.3 Proof of Theorem 3 (Hallucination Criterion)

Suppose $H_p > \log N - I(X;C)$. By Theorem 2, this means the path explores at least one token with probability $< 1/N$. Such a token lies outside the model's high-confidence manifold, and its generation is a sufficient condition for hallucination. This follows from the pigeonhole principle: if the entropy exceeds the bound, some probability mass must be assigned to low-probability tokens.

# B Extended Analytic Traces

## B.1 Toy Arithmetic Example

Consider $V = \{0, 1, 2\}$ and a model that assigns $P_\theta(0|\cdot) = 0.7$, $P_\theta(1|\cdot) = 0.2$, $P_\theta(2|\cdot) = 0.1$ at each step. For a 3-token sequence, the path entropy is:

$$H_t = -[0.7\log 0.7 + 0.2\log 0.2 + 0.1\log 0.1]$$
$$\approx 0.7 \times 0.514 + 0.2 \times 2.322 + 0.1 \times 3.322$$
$$\approx 0.36 + 0.46 + 0.33 = 1.15 \text{ bits}$$

For $T = 3$, $H_p \approx 3.45$ bits.

## B.2 Chain-of-Thought Prompt

Suppose $P_\theta(\text{correct}|\cdot) = 0.8$, $P_\theta(\text{incorrect}|\cdot) = 0.2$ at each reasoning step. Then $H_t \approx 0.8\log(1/0.8) + 0.2\log(1/0.2) \approx 0.257 + 0.464 = 0.721$ bits per step. For a 5-step chain, $H_p \approx 3.6$ bits, which can be compared to the entropy bound for the given $N$ and $I(X;C)$.

## B.3 Retrieval-Augmented Prompt

If retrieval increases $I(X;C)$ by 1 bit, the minimum achievable $H^\star$ drops by 1 bit. For $N = 8$, $\log N = 3$ bits. If $I(X;C)$ increases from 1 to 2 bits, $H^\star$ drops from 2 to 1 bit, halving the minimum perplexity.

These extended traces provide concrete calculations to support the theoretical results and illustrate

the impact of prompt and context design on the
search process.

9