

Guidance Large Language Models at Test Time: A Unified Review of LLM-Training-Free Methods

Anonymous ACL submission

Abstract

Adapting Large Language Models (LLMs) to dynamic constraints typically requires expensive fine-tuning. While training-free test-time guidance offers a flexible alternative, the literature remains fragmented across isolated subfields. This paper presents a unified review of LLM-training-free guidance, systematizing methods that steer behavior without parameter updates. We propose a taxonomy based on the inference lifecycle, categorizing interventions into Input-Space, Latent-Space, Decoding-Space, and Output-Space guidance. Furthermore, we analyze critical trade-offs regarding model accessibility, computational cost, and control granularity. Finally, we discuss emerging frontiers, highlighting the convergence of control mechanisms towards unified architectures and the shift toward rigorous, interpretability-driven steering.

1 Introduction

Large Language Models (LLMs) have evolved into general-purpose foundations capable of solving diverse tasks (OpenAI et al., 2024; Touvron et al., 2023; Guo et al., 2025). To adapt these broad models to specific domains or safety policies, the standard approach is Post-Training, such as Supervised Fine-Tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022; Rafailov et al., 2023). While effective, relying on parameter updates creates three significant bottlenecks in real-world deployment. First, the process is computationally expensive, as iteratively re-training models for every new requirement consumes massive resources (Hu et al., 2022; Ning et al., 2025). Second, it involves significant risks, specifically "catastrophic forgetting" (Li et al., 2024a), where optimizing for a narrow task degrades the model's general capabilities (Harmon et al., 2025). Finally, adaptation is often impossible in commercial settings where models are deployed

as black boxes, blocking user access to weight modification.

Fortunately, parameter updates are not the sole mechanism for adaptation. A parallel and flexible paradigm involves steering the inference process itself while keeping the base model frozen. This approach is pluggable, reversible, and highly efficient, allowing practitioners to modulate behavior on-the-fly. In practice, interventions can occur at various levels: from injecting external knowledge via Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Fan et al., 2024) to adjusting internal neuron activations for surgical style control (Turner et al., 2024). Other techniques involve filtering token probabilities to enforce strict constraints (Lu et al., 2021) or employing iterative search algorithms for complex reasoning (Yao et al., 2023a; Besta et al., 2024). Notably, recent evidence suggests that simple, well-structured prompting can sometimes outperform complex fine-tuning, challenging the assumption that heavy engineering is always necessary (Zhou et al., 2023a; Schulhoff et al., 2025).

Despite this momentum, the research landscape remains fragmented. Broad surveys typically focus on general model capabilities (Zhao et al., 2025a) or post-training scaling (Lai et al., 2025), while reviews on prompting (Schulhoff et al., 2025), retrieval (Fan et al., 2024), and interpretability (Rai et al., 2025) are usually studied in isolation. Although recent works explore test-time scaling (Zhang et al., 2025b), cooperation (Huang et al., 2025a), and alignment (Pan et al., 2025), they prioritize specific objectives rather than offering a holistic view of the control mechanism (see Table 1 in Appendix for a detailed scope comparison). There lacks a mechanism-centric unified framework that connects the entire inference lifecycle, leaving methods that operate on similar principles treated as unrelated simply because they intervene at different stages.

This survey aims to fill this gap by presenting a unified Test-Time Lifecycle Taxonomy for LLM-Training-Free Guidance. We treat test-time guidance as a holistic control problem, organizing methods based on *where* the intervention occurs in the generation process. This perspective allows us to compare disparate techniques—from simple system prompts to surgical activation guidance—on a single map. It also reveals common trade-offs regarding access assumptions (Black-box vs. White-box), control granularity, and test-time latency. We use the term *LLM-Training-Free* to mean that the target LLM weights are not updated, while allowing the use of pluggable auxiliary components (Cunningham et al., 2023).

This survey makes three key contributions:

- **Lifecycle Taxonomy.** We synthesize fragmented methods into a unified framework spanning Input, Latent, Decoding, and Output spaces (Sec.2).
- **Analysis.** We analyze the capabilities of each paradigm regarding model access, inference cost, and steering precision (Sec.3).
- **Future Frontiers.** We identify the emerging trend of "full-stack" architectures and interpretability-driven control mechanisms (Sec.4.2).

2 Problem Formulation and Taxonomy

This section formalizes the concept of *LLM-Training-Free Test-Time Guidance* and introduces a taxonomy organized by the specific stage within the inference lifecycle where an intervention occurs.

2.1 Definition and Problem Formulation

Consider a target Large Language Model (LLM) parameterized by θ , which defines an autoregressive conditional distribution over an output sequence $y = (y_1, \dots, y_T)$ given an input context x :

$$p_{\theta}(y | x) = \prod_{t=1}^T p_{\theta}(y_t | x, y_{<t}) \quad (1)$$

Test-Time Guidance. We define *test-time guidance* as the introduction of a guidance operator \mathcal{G} that transforms the original generative process into a guided distribution:

$$p_{\theta, \mathcal{G}}(y | x, g) = \mathcal{G}(p_{\theta}(\cdot | x), g) \quad (2)$$

where g denotes the *guidance signal*. This signal can manifest as textual instructions (Wei et al., 2022), a latent direction vector (Turner et al., 2024; Zou et al., 2023), a decoding penalty (Li et al., 2023b), or a verifier scoring function (Yao et al., 2023a). Equation 2 is deliberately broad to encompass hybrid systems that may combine multiple signals (e.g., a safety system prompt combined with logit-level constraints).

LLM-Training-Free Constraint. The term *LLM-Training-Free* imposes a strict constraint on the *target* model parameters: during guidance, θ must remain frozen ($\Delta\theta = \mathbf{0}$). This explicitly excludes methods that update the base model, such as full fine-tuning (Devlin et al., 2019) or LoRA (Hu et al., 2022). Crucially, this definition *does not* forbid the use of external auxiliary modules, even if those modules are trained, as long as they are structurally separate from the LLM. Representative examples include dense retrievers (Karpukhin et al., 2020), sparse autoencoders (Huben et al., 2024), and reward models (Stiennon et al., 2022). This modularity is a key practical advantage, allowing us to upgrade control mechanisms without the high computational cost and stability risks associated with re-training the base model.

2.2 Taxonomy: A Lifecycle View

We classify guidance methods based on the principal interface through which \mathcal{G} intervenes. As illustrated in Figure 1, this taxonomy maps methods to four distinct stages of the inference lifecycle: the input context, the internal representation, the decoding distribution, and the output space.

Input-Space Guidance. Input-space guidance transforms the raw input x into a modified context x' using the guidance signal g :

$$x' = f(x; g), \quad p_{\theta, \mathcal{G}}(y | x, g) = p_{\theta}(y | x') \quad (3)$$

Here, the guidance signal g typically manifests as demonstrations for In-Context Learning (ICL) (Brown et al., 2020; Min et al., 2022), external knowledge injection via Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Ram et al., 2023), or system-level instructions (commonly known as system prompts) (Ouyang et al., 2022; Touvron et al., 2023).

Latent-Space Guidance. Latent-space guidance intervenes directly on the internal representations computed during the forward pass (Zou et al., 2023;

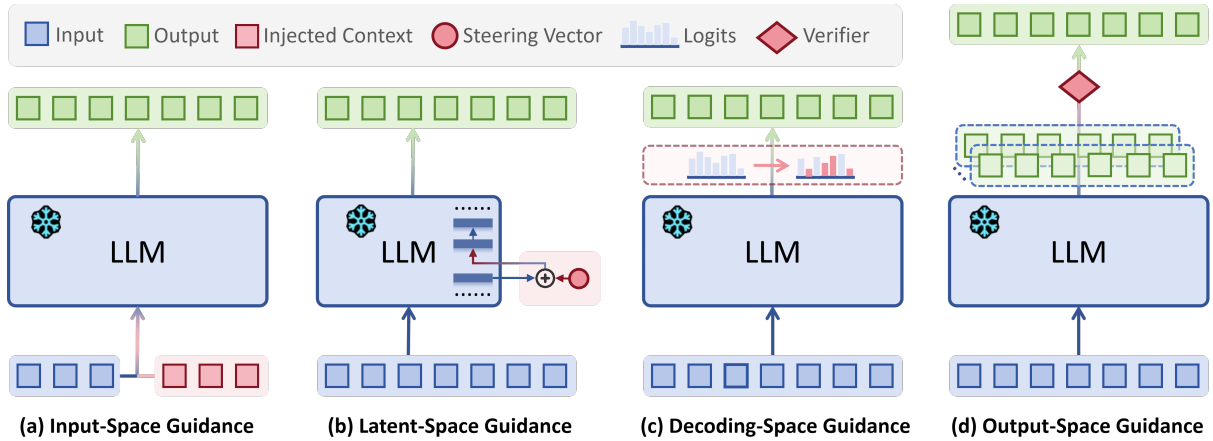


Figure 1: **Overview of the four paradigms for LLM-Training-Free Test-Time Guidance.** The snowflake denotes that the base model parameters remain frozen, while red elements highlight the intervention of guidance signals at distinct stages of the inference lifecycle.

Li et al., 2023a). Let h_t^ℓ denote the hidden state at layer ℓ and time step t . A latent intervention applies a transformation:

$$\tilde{h}_t^\ell = \mathcal{T}_\ell(h_t^\ell; g) \quad (4)$$

where subsequent layers consume the modified \tilde{h}_t^ℓ . The guidance signal g typically encodes a direction in the semantic space (Turner et al., 2024; Zou et al., 2023) or a sparse feature clamp (Huben et al., 2024; Yan et al., 2024). Unlike input guidance, which operates on discrete tokens, latent guidance manipulates the continuous semantic manifold of the model, offering finer granularity for attributes such as sentiment (Subramani et al., 2022) or hallucination mitigation (Li et al., 2023a).

Decoding-Space Guidance. Decoding-space guidance intervenes directly at the token selection interface, reshaping the local probability landscape without updating parameters. Let $z_t \in R^{|V|}$ denote the raw logits at time step t . We formalize the intervention as a logit transformation function r applied prior to normalization:

$$\tilde{z}_t = r(z_t, x, y_{<t}; g) \quad (5)$$

$$p_{\theta, \mathcal{G}}(y_t | x, y_{<t}) = \text{softmax}(\tilde{z}_t) \quad (6)$$

The operator r can enforce hard lexical constraints via masking (Lu et al., 2021), or implement soft contrastive adjustments (Li et al., 2023b; Chuang et al., 2024) where logits from an "amateur" model or early layers are subtracted from z_t . This approach directly modulates the output distribution, providing robust control over specific attributes (Yang and Klein, 2021) and formatting, though it

requires balancing steering strength against potential coherence degradation.

Output-Space Guidance. Output-space guidance shifts the locus of intervention from the step-by-step generation process to the evaluation of complete sequences. Rather than shaping local token probabilities, the system generates a set of candidate hypotheses and selects the optimal output y^* :

$$y^* = \arg \max_{y \in \mathcal{C}} R(y; x, g) \quad (7)$$

Here, \mathcal{C} represents candidates sampled from the base distribution $p_\theta(\cdot|x)$, while R acts as a post-hoc filter or a verifier (Cobbe et al., 2021) or search heuristic (Yao et al., 2023a). This paradigm treats the frozen LLM purely as a proposal distribution, delegating alignment entirely to the selection process.

3 Guidance Methodologies

3.1 Input-Space Guidance

Input-space guidance steers black-box models by systematically transforming the query x rather than internal parameters. As summarized in Figure 2, we categorize these structured frameworks into *Contextual Augmentation*, *Prompt Design and Optimization*, and *Alignment and Safety Prompting*.

3.1.1 Contextual Augmentation

This paradigm bridges the gap between parametric memory and query needs by injecting external information. Retrieval-Augmented Generation (RAG) has evolved from static passage injection (Ram et al., 2023; Shi et al., 2024b) to active,

Contextual Augmentation

IC-RAG (Ram et al., 2023), REPLUG (Shi et al., 2024b), CoK (Li et al., 2024b), RQ-RAG (Chan et al., 2024), GMR (Lee et al., 2022), RAG-Fusion (Rackauckas, 2024), KRA-GEN (Matsumoto et al., 2024), SimRAG (Xu et al., 2025), SEER (Zhao et al., 2024a), L-RAG (Lin et al., 2025a), HopRAG (Liu et al., 2025a), KERAG (Sun et al., 2025c), TC-RAG (Jiang et al., 2025b), HyKGE (Jiang et al., 2025c), UniRAG (Li et al., 2025e), OmniRAGMed (Chen et al., 2025c), MolRAG (Xian et al., 2025), EventRAG (Yang et al., 2025d), KiRAG (Fang et al., 2025a), FaithfulRAG (Zhang et al., 2025a), Dialogue-RAG (Li et al., 2025d), SGIC (Chen et al., 2025a), MedGraphRAG (Wu et al., 2025a), Astute RAG (Wang et al., 2025a), DualRAG (Cheng et al., 2025), T-GRAG (Li et al., 2025b), SafeDriveRAG (Ye et al., 2025), RecipeRAG (Yang et al., 2025c), HM-RAG (Liu et al., 2025c), DynamicRAG (Sun et al., 2025b), What-makes-ICL (Liu et al., 2022), Lost-in-Middle (Liu et al., 2024a), L2R-ICL (Wang et al., 2024a), CED-Select (Iter et al., 2023), RGER (Lin et al., 2025b), GENICL (Zhang et al., 2025e), MLSM/TTF (Liu et al., 2025b), UniICL (Gao et al., 2025), FLARE (Jiang et al., 2023b), CoVe (Wang et al., 2024c), ToG2(Ma et al., 2025)

Prompt Design and Optimization

Chain-of-Thought (Wei et al., 2022), Zero-shot CoT (Kojima et al., 2022), Self-Consistency (Wang et al., 2023), Tab-CoT (Ziqi and Lu, 2023), Auto-CoT (Zhang et al., 2023), Logical CoT (Zhao et al., 2024b), CDW-CoT (Fang et al., 2025b), RankCoT (Wu et al., 2025c), PCoT (Modzelewski et al., 2025), SemCoT (He et al., 2025a), APE (Zhou et al., 2023b), OPRO (Yang et al., 2024), BPO (Cheng et al., 2024a), RLPrompt (Deng et al., 2022), PromptAgent (Wang et al., 2024b), PEARL-plan (Sun et al., 2024b), PromptBreeder (Fernando et al., 2024), Structure-Guided (Cheng et al., 2024b), AutoPrompt (Shin et al., 2020), Automate-CoT (Shum et al., 2023)

Alignment and Safety Prompting

Self-Reminders (Xie et al., 2023), URIAL (Lin et al., 2024), LLaMA-2 cfg. (Touvron et al., 2023), ICAG (Zhou et al., 2024b), SmoothLLM (Robey et al., 2025), ROSE (Zhong et al., 2024), Intention Analysis (Zhang et al., 2025d), CIP (Hahm et al., 2025), LLM-Self-Defense (Phute et al., 2024), SelfDefenD (Wang et al., 2025f), P-Aligner (Song et al., 2025b), ICDPO (Song et al., 2025a)

Figure 2: Input-Space guidance methods.

structure-aware inference. Recent frameworks interleave retrieval with reasoning steps (Li et al., 2024b; Chan et al., 2024; Rackauckas, 2024) or leverage graph-based and domain-specific structures—such as knowledge graphs (Liu et al., 2025a; Wu et al., 2025a), molecular fingerprints (Xian et al., 2025), and interaction histories (Li et al., 2025d)—to handle complex queries. Active strategies further empower models to decide *when* to retrieve (Jiang et al., 2023b; Yao et al., 2023b), optimizing the information flow dynamically.

Parallely, In-Context Learning (ICL) optimization focuses on the precise selection and compression of demonstrations. Instead of simple similarity matching, modern approaches align exemplars with task-specific reasoning patterns (Iter et al., 2023; Lin et al., 2025b; Zhang et al., 2025e) or compress extensive contexts into efficient "memory tokens" to reduce computational overhead (Gao et al., 2025; Liu et al., 2025b).

3.1.2 Prompt Design and Optimization

Treating the prompt as a programmable interface, this stream optimizes the instruction format to elicit specific capabilities. Structured Reasoning methods extend standard Chain-of-Thought (Wei et al., 2022) by imposing strict logical, tabular, or semantic constraints (Ziqi and Lu, 2023; Zhao et al., 2024b; He et al., 2025a) and ranking reasoning

paths (Wu et al., 2025c) to minimize hallucinations.

To eliminate manual trial-and-error, Automated Prompt Optimization formulates instruction design as a search problem. Black-box optimization methods utilize task feedback to iteratively refine prompts (Yang et al., 2024; Zhou et al., 2023b; Cheng et al., 2024a), while advanced planning algorithms employ Monte Carlo Tree Search (Wang et al., 2024b) or evolutionary strategies (Fernando et al., 2024; Cheng et al., 2024b) to navigate the prompt space, effectively turning prompt engineering into an algorithmic search process.

3.1.3 Alignment and Safety Prompting

Input-space guidance also serves as a critical, training-free defense layer. While early work relied on static system prompts to define behavioral boundaries (Lin et al., 2024; Touvron et al., 2023; Xie et al., 2023), 2025 methodologies have shifted towards Active Instance-Level Alignment. Approaches like P-Aligner and ICDPO dynamically rewrite user queries into principled forms based on alignment preferences (Song et al., 2025b,a). Furthermore, Agentic Defense mechanisms decouple safety evaluation from response generation; by employing shadow models or causal reasoning agents to screen inputs and monitor potential risks (Phute et al., 2024; Wang et al., 2025f; Hahm et al., 2025; Zhou et al., 2024b; Zhang et al., 2025d), these systems provide robust protection against jailbreaks and adversarial attacks without modifying model weights.

3.2 Latent-Space Guidance

Latent-space guidance intervenes on the internal representations of a frozen LLM during inference. Formally, a guidance operator transforms hidden states h_t^ℓ into \tilde{h}_t^ℓ via a transformation $\mathcal{T}_\ell(h_t^\ell; g)$ driven by a guidance signal g . Unlike input-space methods that manipulate discrete tokens, these approaches operate on continuous activation vectors, utilizing probes, autoencoders, or geometric projections to steer model behavior. We categorize these methodologies based on their geometric granularity and dynamic nature as follows.

Global Directional Arithmetic. The most foundational approach treats guidance as a global shift along learned directions in the residual stream, typically formalized as an additive update $\tilde{h}_t^\ell = h_t^\ell + \alpha v$. PPLM (Dathathri et al., 2020) pioneered this paradigm by utilizing test-time gradient updates to shift hidden states toward target attributes.

Activation Addition (Turner et al., 2024) and Representation Engineering (RepE) (Zou et al., 2023) extract these steering vectors by contrasting activations from paired positive and negative prompts, successfully modulating high-level traits like helpfulness, honesty, or political framing. Inference-Time Intervention (ITI) (Li et al., 2023a) derives truthfulness directions from linear probes trained on labeled datasets, while In-Context Vectors (ICV) (Liu et al., 2024b) and Function Vectors (Todd et al., 2024) distill task-specific directions directly from few-shot demonstrations. Recent empirical studies focus on optimizing the application of these vectors, analyzing layer sensitivity and scaling laws to maximize control efficacy while minimizing perplexity degradation (Stoehr et al., 2024; Stolfo et al., 2025; Liu et al., 2025d).

Decomposition and Fine-Grained Control.

Moving beyond raw residual directions, this paradigm decomposes representations into an interpretable basis before intervention. Sparse Autoencoders (SAEs) have emerged as a powerful tool to disentangle polysemantic neurons into monosemantic features; steering in this context implies identifying features correlated with specific behaviors and clamping their coefficients (Huben et al., 2024; Soo et al., 2025; Yan et al., 2024). At the neuron level, methods identify "target atoms" or clusters of neurons responsible for hallucinations or toxic content, selectively dampening them to mitigate adverse behaviors while preserving general reasoning capabilities (Zhang et al., 2025f; Huang et al., 2025c; Han et al., 2025). Work in 2025 further refines this by mapping safety-critical concepts to sparse bases to decouple refusal mechanisms from reasoning logic (Wang et al., 2025c; Valentino et al., 2025; Sakarvadia et al., 2025).

Geometric Constraints and Projection. To minimize the side effects of steering, this category employs linear operators that respect the geometric structure of the representation space. Instead of simple addition, methods like AlphaSteer (Sheng et al., 2025) and ASGuard (Park et al., 2025) compute steering vectors within the null-space of task-relevant representations, ensuring that safety edits remain orthogonal to core competencies. Projection-based steering (Postmus and Abreu, 2024) and probe-free low-rank interventions (Jiang et al., 2025a; Oozeer et al., 2025) parameterize guidance as a low-rank update $\tilde{h} = h + Wh$, learning a structured transformation that aligns acti-

Latent-Space Guidance

Latent steering vectors (Subramani et al., 2022), IDANI (Antverg et al., 2022), PPLM (Dathathri et al., 2020), ITI (Inference-Time Intervention) (Li et al., 2023a), Activation Engineering (Turner et al., 2024), Representation Engineering (RepE) (Zou et al., 2023), (Huben et al., 2024), (Yan et al., 2024), EAST (Entropic Activation Steering) (Rahn et al., 2024), NL-ITI (Hoscilowicz et al., 2024), (Fatahi Bayat et al., 2024), In-context Vectors (ICV) (Liu et al., 2024b), Activation scaling (Stoehr et al., 2024), (Postmus and Abreu, 2024), Steering Target Atoms (STA) (Wang et al., 2025c), SCANS (Cao et al., 2025), FLOORAIN (Jiang et al., 2025a), SADI (Wang et al., 2025e), DSAS (Do et al., 2025), (Valentino et al., 2025), (Sakarvadia et al., 2025), (Zhao et al., 2025b), (Oozeer et al., 2025), (Cybererey and Evans, 2025), AlphaSteer (Sheng et al., 2025), (Jorgensen et al., 2023), (van der Weij et al., 2024), (Todd et al., 2024), (Konen et al., 2024), (Stolfo et al., 2025), (Lee et al., 2025a), (Wang et al., 2025d), (Hegazy et al., 2025), (Rodriguez et al., 2024), (Zhang et al., 2025f), (Park et al., 2025), (Wu et al., 2025b), (Sun et al., 2025a), (Soo et al., 2025), (Karnik and Bansal, 2025), (Huang et al., 2025c), (Han et al., 2025), (Liu et al., 2025d), (Scalena et al., 2024)

Figure 3: Latent-space guidance methods.

Decoding-Space Guidance

GeDi (Krause et al., 2021), Contrastive Decoding (Li et al., 2023b), DoLa (Chuang et al., 2024), NeuroLogic Decoding (Lu et al., 2021), FUDGE (Yang and Klein, 2021), DExperts (Liu et al., 2021), Grammar-Constrained Decoding (Raspanti et al., 2025), ROSE (Zhong et al., 2024), Instructive Decoding (Alajrami et al., 2025), MOD (Shi et al., 2024a), RMOD (Son et al., 2025), PAD (Chen et al., 2025b), DeAL (Huang et al., 2025b), DeRa (Liu et al., 2024c), SAFEINFER (Banerjee et al., 2025), IBD (Zhu et al., 2024), Nudging (Fei et al., 2025), CDT (Yang et al., 2025a), UACD (Lee et al., 2025b), SLED (Zhang et al., 2024), RGD (Mañas et al., 2025), PRGD (Nguyen et al., 2025), ACD (Zhang et al., 2025c), Drift (Kim et al., 2025), RSD (Liao et al., 2025), Air-Decoding (Zhong et al., 2023), CARDS (Li et al., 2025a), CARE (Hu et al., 2025b), CDT (Yang et al., 2025b), MCA (Fu et al., 2025).

Figure 4: Decoding-space guidance methods.

vations with a target distribution rather than a single direction (Rodriguez et al., 2024). These methods emphasize "surgical" interventions that mathematically constrain the edit distance in task-critical subspaces.

Context-Adaptive Dynamics. Fixed interventions may lead to over-control. Adaptive methods resolve this by modulating the intervention strength α or direction v based on the current context. Non-linear ITI (Li et al., 2023a) and entropy-aware steering (Hoscilowicz et al., 2024) scale interventions based on model uncertainty, applying corrections only when hallucination risk is high. Semantics-adaptive controllers (Wang et al., 2025e,d) dynamically adjust activation scales according to the semantic topic of the input. In safety domains, methods like SCANS (Cao et al., 2025) and preemptive detectors (Karnik and Bansal, 2025) employ lightweight gating mechanisms to trigger steering only when harmful intent is detected. Extensions to multimodal settings, such as AutoSteer (Wu et al., 2025b) and Layer-Navigator (Sun et al., 2025a), dynamically route guidance across layers depending on joint vision-language activation patterns, ensuring context-sensitive safety alignment.

3.3 Decoding-Space Guidance

Decoding-space guidance intervenes on the logits z_t at each generation step, modifying token

probabilities via a guidance operator r before sampling. This enables fine-grained steering of style, safety, and task behavior without parameter updates. As summarized in Figure 4, we categorize these methodologies into *Score-based Reweighting* and *Constraint-based Filtering*.

Score-based Reweighting. Most decoding-time methods retain the full vocabulary but adjust the relative probability of tokens to favor desirable behaviors. GeDi (Krause et al., 2021) and DExperts (Liu et al., 2021) achieved this by reweighting predictions using logits from external auxiliary models (experts) or anti-experts. Contrastive strategies leverage the model’s own capabilities to "self-correct" logits: methods like Contrastive Decoding (Li et al., 2023b) and DoLa (Chuang et al., 2024) subtract the logits of a "weak" model from a "strong" one to amplify factual signals. Newer variants such as UACD (Lee et al., 2025b) and SLED (Zhang et al., 2024) introduce uncertainty and stability metrics to modulate this contrastive penalty, avoiding noise when the model is unconfident. In safety alignment, ROSE (Zhong et al., 2024), Instructive Decoding (Alajrami et al., 2025), and ACD (Zhang et al., 2025c) down-weight tokens favored by unsafe or noisy prompts while boosting those aligned with clean instructions. A second cluster fuses base logits with scores from external auxiliary models or reward functions. Frameworks like MOD (Shi et al., 2024a) and DeAL (Huang et al., 2025b) combine task utility scores with base probabilities, while PAD (Chen et al., 2025b) and Drift (Kim et al., 2025) inject user-specific preference vectors for personalized alignment. Similarly, methods such as SAFEINFER (Banerjee et al., 2025) and CDT (Yang et al., 2025b) fuse risk estimation or truthfulness comparator scores to dynamically suppress toxic or hallucinatory continuations.

Constraint-based Filtering. A complementary paradigm treats decoding as a constrained search problem, using the guidance operator to strictly rule out invalid tokens rather than merely reshaping their scores. Classical methods employ symbolic masks to enforce validity: NeuroLogic Decoding (Lu et al., 2021) applies predicate logic constraints, while Grammar-Constrained Decoding (Raspanti et al., 2025) uses formal grammars to ensure syntactically valid code generation. Extending this to dynamic trajectory control, recent approaches explicitly prune the search space based on interme-

Output-Space Guidance

ToT (Yao et al., 2023a), GoT (Besta et al., 2024), RAP (Hao et al., 2023), Multi-Agent Debate (Zhou et al., 2024a), Reflexion (Shinn et al., 2023), Self-Refine (Madaan et al., 2023), CRITIC (Gou et al., 2024), Training Verifiers (Cobbe et al., 2021), SelfCheckGPT (Manakul et al., 2023), MBR Decoding (Suzgun et al., 2023), Universal Self-Consistency (Chen et al., 2024), Speculative Rejection (Sun et al., 2024a),(Ren et al., 2023)), COOPER-ATE(Feng et al., 2024), LLM-Blender (Jiang et al., 2023a), Selection-Inference (Creswell et al., 2023), Mutual Reasoning (Qi et al., 2025a)

Figure 5: Output-space guidance methods.

diated quality assessments. CARE (Hu et al., 2025b) introduces a "rollback and introspection" mechanism that detects unsafe trajectories mid-generation and reverts to a safe state. Similarly, CARDS (Li et al., 2025a) and Reward-Guided Speculative Decoding (RSD) (Liao et al., 2025) implement cascade sampling or verifiers to filter out low-quality speculative candidates before they are committed, effectively pruning undesirable paths to guarantee alignment constraints.

3.4 ◆ Output-Space Guidance

Output-space guidance steers generation at the response level via post-hoc scoring or structural search, enabling complex reasoning without parameter updates. As summarized in Figure 5, we categorize these methods into *Search and Planning*, *Iterative Refinement*, *Candidate Selection and Reranking*, and *Multi-Model Collaboration*.

Search and Planning. These methods structure generation as a search process over a space of reasoning steps. Tree of Thoughts (ToT) (Yao et al., 2023a) allows the model to explore multiple branches and backtrack when necessary. Graph of Thoughts (GoT) (Besta et al., 2024) extends this to arbitrary graphs, enabling the aggregation of information from different thoughts. RAP (Hao et al., 2023) and LATS (Zhou et al., 2024a) treat the LLM as a world model within a Monte Carlo Tree Search framework, simulating future outcomes to plan the optimal reasoning path.

Iterative Refinement. Instead of searching multiple branches, refinement methods optimize a single response through feedback loops. Self-Refine (Madaan et al., 2023) generates an initial draft and iteratively prompts the model to critique and improve its own output. Reflexion (Shinn et al., 2023) introduces "verbal reinforcement learning," where the agent reflects on past failures to update its memory for future trials. CRITIC (Gou et al., 2024) enhances this cycle by using external tools to verify outputs and guide revisions based on execution

486	feedback.	
487	Candidate Selection and Reranking. This	
488	paradigm follows a "generate-then-select" work-	
489	flow: sample multiple candidates and select the	
490	best one. Verifier-based methods (Cobbe et al.,	
491	2021) and SelfCheckGPT (Manakul et al., 2023)	
492	score candidates based on correctness or hallucina-	
493	tion probability. Universal Self-Consistency (Chen	
494	et al., 2024) and MBR Decoding (Suzgun et al.,	
495	2023) select the consensus output that minimizes	
496	risk or maximizes consistency across diverse sam-	
497	ples. To improve efficiency, Speculative Rejection	
498	(Sun et al., 2024a) and Self-Evaluation (Ren et al.,	
499	2023) use reward models or confidence scores to	
500	prune low-quality candidates early in the process.	
501	Multi-Model Collaboration. Guidance can also	
502	emerge from the interaction of diverse models.	
503	LLM-Blender (Jiang et al., 2023a) employs a rank-	
504	ing and fusion module to synthesize the best com-	
505	ponents from multiple model outputs. Selection-	
506	Inference (Creswell et al., 2023) decomposes rea-	
507	soning into modular selection and inference steps	
508	handled by specialized prompts. Recent collabor-	
509	ative frameworks like Mutual Reasoning (Qi	
510	et al., 2025a) and COOPERATE (Feng et al., 2024)	
511	demonstrate that diverse models can correct each	
512	other's errors through voting or abstention mecha-	
513	nisms, leveraging the wisdom of the crowd without	
514	training.	
515	4 Discussion	
516	We synthesize these paradigms into a comparative	
517	framework to highlight trade-offs in access, cost,	
518	and granularity, followed by an analysis of critical	
519	deployment challenges.	
520	4.1 Comparative Analysis	
521	The Access Divide: Black-box vs. White-box.	
522	The most immediate constraint in real-world de-	
523	ployment is model access. Input- and output-space	
524	methodologies treat the LLM purely as an API	
525	endpoint, making them the only viable options for	
526	proprietary models like GPT-5 or Gemini. Techn-	
527	iques such as Contextual Augmentation (Ram	
528	et al., 2023) and Tree of Thoughts (Yao et al.,	
529	2023a) scale naturally with stronger base mod-	
530	els without requiring internal visibility. In con-	
531	trast, latent- and decoding-space guidance neces-	
532	sitate access to activation vectors or vocabulary	
533	logits. While this restricts their application to open-	
534	weights models, it unlocks a level of "surgical"	
	precision impossible via prompting alone. For in-	535
	stance, Activation Engineering (Turner et al., 2024)	536
	and Inference-Time Intervention (ITI) (Li et al.,	537
	2023a) can suppress specific hallucination patterns	538
	or safety risks by directly dampening the corre-	539
	sponding internal features, bypassing the model's	540
	tendency to rationalize instructions in the prompt.	541
	The Precision-Coherence Trade-off. Does finer	542
	control always lead to better outcomes? Evi-	543
	dence suggests a delicate balance between steering	544
	strength and semantic coherence. Decoding-space	545
	methods like NeuroLogic (Lu et al., 2021) and	546
	Contrastive Decoding (Li et al., 2023b) enforce	547
	strict token-level constraints, effectively pruning	548
	the probability mass of undesirable outputs. How-	549
	ever, aggressively reshaping the output distribution	550
	can force the model off its natural "manifold," lead-	551
	ing to disjointed or ungrammatical text. Similarly,	552
	latent interventions face the challenge of "manifold	553
	distortion": as noted in recent studies on activation	554
	scaling (Stolfo et al., 2025; Liu et al., 2025d), push-	555
	ing internal states too far along a steering vector	556
	often causes a spike in perplexity and a degrada-	557
	tion of general capabilities. This contrasts with	558
	input-space methods, which, by operating through	559
	natural language, generally preserve fluency but	560
	often fail to enforce negative constraints.	561
	Inference-Time Compute as the New Scaling	562
	Law. A parallel trend connects these paradigms	563
	through the lens of compute allocation. While la-	564
	tent and decoding methods typically incur negligi-	565
	ble overhead—often just a single vector addition	566
	or logit subtraction per step—output-space guid-	567
	ance represents a fundamental shift towards "test-	568
	time scaling." By generating multiple candidates	569
	or searching through reasoning trees, methods like	570
	ToT (Yao et al., 2023a) and RAP (Hao et al., 2023)	571
	effectively trade inference latency for intelligence,	572
	enabling System-2 reasoning. This mirrors the ob-	573
	servation that difficult reasoning tasks benefit more	574
	from verifying multiple paths than from optimizing	575
	a single forward pass. Thus, the choice of method	576
	often reduces to a resource decision: is it more	577
	efficient to "patch" a model's intuition via latent	578
	steering, or to spend compute on explicit search via	579
	output guidance?	580
	Convergence of Control Spaces. Although cate-	581
	gorized into distinct spaces, recent advancements	582
	indicate a blurring of boundaries where effective	583
	systems hybridize these approaches. Contrastive	584

Decoding (Li et al., 2023b), nominally a decoding method, relies heavily on engineered "negative" prompts (input-space) to derive its steering signal (Zhong et al., 2024). Similarly, Latent-space vectors are often distilled from In-Context Learning examples (Liu et al., 2024b), effectively projecting input demonstrations into the residual stream. Most notably, safety frameworks like CARE (Hu et al., 2025b) now integrate decoding-time monitoring with output-space rollback mechanisms. This convergence suggests that future guidance will likely be "full-stack": leveraging prompts for intent, latent/decoding methods for constraints, and output search for verification.

4.2 Challenges and Future Directions

While training-free guidance offers flexibility, it faces critical hurdles in efficiency and reliability. We highlight four fundamental contradictions and the emerging frontiers promising to resolve them.

The Alignment-Capability Dilemma and Glass-Box Steering. A persistent challenge in current interventions is the trade-off between safety and capability, often resulting in an alignment tax where steering a model towards safety degrades its general reasoning. This phenomenon largely arises because standard activation steering operates on entangled, polysemantic representations, where a single direction encodes multiple distinct concepts. Consequently, blindly adding a steering vector often induces manifold distortion, leading to spikes in perplexity or over-refusal. To resolve this, the field is transitioning from coarse-grained direction arithmetic to interpretability-driven steering. By leveraging Sparse Autoencoders to decompose dense activations into overcomplete, monosemantic features (Yan et al., 2024; Cunningham et al., 2023), future methods will move towards precise feature clamping. This technique allows for the surgical suppression of specific behaviors, such as deception or bias, without disrupting the broader semantic manifold, effectively enabling control with minimal side effects (He et al., 2025b).

Inference Efficiency and Unified Architectures. Guidance shifts the computational burden from training to inference, where heavy retrieval or iterative search mechanisms can prohibit real-time deployment. To address this, future systems must evolve into unified "Full-Stack" architectures that employ hierarchical routing to optimize the computational budget per query, rather than applying

a uniform processing chain. A meta-controller assesses input complexity, directing simple prompts to lightweight latent interventions (System-1) while reserving expensive tree-search algorithms for high-stakes reasoning tasks (System-2). Recent breakthroughs in scaling test-time compute (Snell et al., 2024) and reasoning-star strategies (Qi et al., 2025b) demonstrate that such adaptive allocation can simulate the performance gains of reinforcement learning without parameter updates, establishing inference-as-alignment as a resource-efficient paradigm.

Reliability Crisis and Theoretical Guarantees. Real-world applications require satisfying conflicting constraints simultaneously (e.g., conciseness vs. detailed reasoning), often leading to destructive interference in linear steering combinations. Moving beyond empirical tuning requires establishing a "Control Theory for LLMs." Emerging research treating generation as a Markov Decision Process (MDP) suggests that verifiable safety properties may be achievable through constrained decoding policies (Kumar et al., 2025) and Pareto-optimal multi-objective optimization (Son et al., 2025), ensuring convergence to safe states.

The Generalization Gap and Multimodal Steering. As LLMs evolve into Multimodal Large Language Models, text-only guidance is becoming insufficient. Visual inputs often bypass text-based safety guardrails, leading to "visual jailbreaks" where adversarial images trigger harmful outputs despite safe textual prompts. The frontier lies in developing cross-modal interfaces that bridge this gap by establishing a shared control space. Emerging research focuses on using textual prompts to directly modulate visual encoders (Wu et al., 2025b) or intervening on joint vision-language representations (Wang et al., 2025b). This points toward a modality-agnostic control layer capable of aligning video, audio, and text generation simultaneously.

5 Conclusion

This survey systematizes training-free guidance through a mechanism-centric Inference Lifecycle Taxonomy. By unifying interventions across input, latent, decoding, and output spaces, we highlight key trade-offs in accessibility, cost, and precision. We envision a future shift toward "full-stack" architectures that dynamically bridge static models with evolving user needs.

684 Limitations

685 This survey establishes a structural foundation for
686 training-free guidance, yet we acknowledge the gap
687 between our theoretical taxonomy and engineering
688 reality. Our mechanism-centric classification in-
689 evitably abstracts over complex implementation
690 nuances—such as inference latency trade-offs and
691 memory optimization—that are critical in deploy-
692 ment but distinct from conceptual categorization.
693 Furthermore, we deliberately abstain from a unified
694 quantitative benchmark, as the efficacy of guidance
695 is deeply coupled with the rapidly evolving capa-
696 bilities of base models; techniques essential for
697 weaker models may become redundant as reason-
698 ing scales, rendering static comparisons scientifi-
699 cally fragile. Thus, this work serves as a qualitative
700 synthesis of representative control paradigms rather
701 than an exhaustive empirical leaderboard.

702 References

703 Ahmed Alajrami, Xingwei Tan, and Nikolaos Ale-
704 tras. 2025. [Fine-tuning on noisy instructions: Ef-
705 fects on generalization and performance](#). *Preprint*,
706 arXiv:2510.03528.

707 Omer Antverg, Eyal Ben-David, and Yonatan Belinkov.
708 2022. [IDANI: Inference-time domain adaptation via
709 neuron-level interventions](#). In *Proceedings of the
710 Third Workshop on Deep Learning for Low-Resource
711 Natural Language Processing*, pages 21–29, Hybrid.
712 Association for Computational Linguistics.

713 Somnath Banerjee, Sayan Layek, Soham Tripathy,
714 Shanu Kumar, Animesh Mukherjee, and Rima Hazra.
715 2025. [Safeinfer: Context adaptive decoding time
716 safety alignment for large language models](#). *Proceed-
717 ings of the AAI Conference on Artificial Intelligence*,
718 39(26):27188–27196.

719 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-
720 stenberger, Michał Podstawski, Lukas Gianinazzi,
721 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadow-
722 ski, Piotr Nyczyk, and Torsten Hoeffler. 2024. [Graph
723 of thoughts: solving elaborate problems with large
724 language models](#). In *Proceedings of the Thirty-
725 Eighth AAI Conference on Artificial Intelligence
726 and Thirty-Sixth Conference on Innovative Applica-
727 tions of Artificial Intelligence and Fourteenth Sym-
728 posium on Educational Advances in Artificial Intelli-
729 gence*, AAAI’24/IAAI’24/EAAI’24. AAAI Press.

730 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
731 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
732 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
733 Aspell, Sandhini Agarwal, Ariel Herbert-Voss,
734 Gretchen Krueger, Tom Henighan, Rewon Child,
735 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
736 Clemens Winter, and 12 others. 2020. Language

models are few-shot learners. In *Proceedings of the
34th International Conference on Neural Information
Processing Systems*, NIPS ’20, Red Hook, NY, USA.
Curran Associates Inc.

Zouying Cao, Yifei Yang, and Hai Zhao. 2025. [Scans:
mitigating the exaggerated safety for llms via safety-
conscious activation steering](#). In *Proceedings of
the Thirty-Ninth AAI Conference on Artificial In-
telligence and Thirty-Seventh Conference on Inno-
vative Applications of Artificial Intelligence and
Fifteenth Symposium on Educational Advances in
Artificial Intelligence*, AAAI’25/IAAI’25/EAAI’25.
AAAI Press.

Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo,
Wei Xue, Yike Guo, and Jie Fu. 2024. [RQ-RAG:
Learning to refine queries for retrieval augmented
generation](#). In *First Conference on Language Model-
ing*.

Guanhua Chen, Yutong Yao, Lidia S. Chao, Xuebo Liu,
and Derek F. Wong. 2025a. [SGIC: A self-guided iter-
ative calibration framework for RAG](#). In *Proceedings
of the 63rd Annual Meeting of the Association for
Computational Linguistics (Volume 1: Long Papers)*,
pages 28357–28370, Vienna, Austria. Association
for Computational Linguistics.

Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai,
and Zuozhu Liu. 2025b. [PAD: Personalized align-
ment at decoding-time](#). In *The Thirteenth Interna-
tional Conference on Learning Representations*.

Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan
Xiao, Pengcheng Yin, Sushant Prakash, Charles Sut-
ton, Xuezhi Wang, and Denny Zhou. 2024. [Universal
self-consistency for large language models](#). In *ICML
2024 Workshop on In-Context Learning*.

Zhe Chen, Yusheng Liao, Shuyang Jiang, Pingjie
Wang, YiQiu Guo, Yanfeng Wang, and Yu Wang.
2025c. [Towards omni-RAG: Comprehensive
retrieval-augmented generation for large language
models in medical applications](#). In *Proceedings
of the 63rd Annual Meeting of the Association for
Computational Linguistics (Volume 1: Long Papers)*,
pages 15285–15309, Vienna, Austria. Association
for Computational Linguistics.

Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning
Wang, Yuxiao Dong, Jie Tang, and Minlie Huang.
2024a. [Black-box prompt optimization: Aligning
large language models without model training](#). In
*Proceedings of the 62nd Annual Meeting of the As-
sociation for Computational Linguistics (Volume 1:
Long Papers)*, pages 3201–3219, Bangkok, Thailand.
Association for Computational Linguistics.

Kewei Cheng, Nesreen K. Ahmed, Theodore L. Willke,
and Yizhou Sun. 2024b. [Structure guided prompt:
Instructing large language model in multi-step reason-
ing by exploring graph structure of the text](#). In *Pro-
ceedings of the 2024 Conference on Empirical Meth-
ods in Natural Language Processing*, pages 9407–

737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793

1133	Hakyung Lee, Subeen Park, Joowang Kim, Sungjun Lim, and Kyungwoo Song. 2025b. Uncertainty-aware contrastive decoding . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 26376–26391, Vienna, Austria. Association for Computational Linguistics.	1188
1134		1189
1135		1190
1136		1191
1137		1192
1138		1193
1139	Hyunji Lee, Sohee Yang, Hanseok Oh, and Minjoon Seo. 2022. Generative multi-hop retrieval . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 1417–1436, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1194
1140		1195
1141		
1142		
1143		
1144		
1145	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20</i> , Red Hook, NY, USA. Curran Associates Inc.	1196
1146		1197
1147		1198
1148		1199
1149		1200
1150		1201
1151		1202
1152		1203
1153		
1154	Bolian Li, Yifan Wang, Anamika Lochab, Ananth Grama, and Ruqi Zhang. 2025a. Cascade reward sampling for efficient decoding-time alignment . <i>Preprint</i> , arXiv:2406.16306.	1204
1155		1205
1156		1206
1157		1207
1158	Dong Li, Yichen Niu, Ying Ai, Xiang Zou, Biqing Qi, and Jianxing Liu. 2025b. T-grag: A dynamic graphrag framework for resolving temporal conflicts and redundancy in knowledge retrieval . In <i>Proceedings of the 33rd ACM International Conference on Multimedia, MM '25</i> , page 11880–11889, New York, NY, USA. Association for Computing Machinery.	1208
1159		1209
1160		
1161		
1162		
1163		
1164		
1165	Hongyu Li, Liang Ding, Meng Fang, and Dacheng Tao. 2024a. Revisiting catastrophic forgetting in large language model tuning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 4297–4308, Miami, Florida, USA. Association for Computational Linguistics.	1210
1166		1211
1167		1212
1168		1213
1169		1214
1170		1215
1171	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023a. Inference-time intervention: Eliciting truthful answers from a language model . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1216
1172		1217
1173		1218
1174		1219
1175		1220
1176	Pingzhi Li, Zhen Tan, Mohan Zhang, Huaizhi Qu, Huan Liu, and Tianlong Chen. 2025c. DOGe: Defensive output generation for LLM protection against knowledge distillation . <i>arXiv preprint arXiv:2505.19504</i> .	1221
1177		1222
1178		1223
1179		1224
1180	Qiwei Li, Teng Xiao, Zuchao Li, Ping Wang, Mengjia Shen, and Hai Zhao. 2025d. Dialogue-RAG: Enhancing retrieval for LLMs via node-linking utterance rewriting . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 24423–24438, Vienna, Austria. Association for Computational Linguistics.	1225
1181		1226
1182		1227
1183		
1184		
1185		
1186		
1187		
	Rui Li, Liyang He, Qi Liu, Zheng Zhang, Heng Yu, Yuyang Ye, Linbo Zhu, and Yu Su. 2025e. Uni-RAG: Unified query understanding method for retrieval augmented generation . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14163–14178, Vienna, Austria. Association for Computational Linguistics.	1228
		1229
		1230
		1231
		1232
		1233
		1234
	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023b. Contrastive decoding: Open-ended text generation as optimization . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.	1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
	Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. 2024b. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Baohao Liao, Yuhui Xu, Hanze Dong, Junnan Li, Christof Monz, Silvio Savarese, Doyen Sahoo, and Caiming Xiong. 2025. Reward-guided speculative decoding for efficient LLM reasoning . In <i>Forty-second International Conference on Machine Learning</i> .	
	Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2024. The unlocking spell on base LLMs: Rethinking alignment via in-context learning . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Jiaen Lin, Jingyu Liu, and Yingbo Liu. 2025a. Optimizing multi-hop document retrieval through intermediate representations . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 15798–15809, Vienna, Austria. Association for Computational Linguistics.	
	Yukang Lin, Bingchen Zhong, Shuoran Jiang, Joanna Siebert, and Qingcai Chen. 2025b. Reasoning graph enhanced exemplars retrieval for in-context learning . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 9737–9759, Abu Dhabi, UAE. Association for Computational Linguistics.	
	Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6691–6706, Online. Association for Computational Linguistics.	

1470		and <i>Interpreting Neural Networks for NLP</i> , pages 577–603, Miami, Florida, US. Association for Computational Linguistics.	
1471			
1472			
1473	Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yin-		
1474	heng Li, Aayush Gupta, HyoJung Han, Sevien Schul-		
1475	hoff, Pranav Sandeep Dulepet, Saurav Vidyadhara,		
1476	Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson		
1477	Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, and		
1478	12 others. 2025. The prompt report: A systematic		
1479	survey of prompt engineering techniques . <i>Preprint</i> ,		
1480	arXiv:2406.06608.		
1481			
1482	Leheng Sheng, Changshuo Shen, Weixiang Zhao, Jun-		
1483	feng Fang, Xiaohao Liu, Zhenkai Liang, Xiang Wang,		
1484	An Zhang, and Tat-Seng Chua. 2025. Alphasteer:		
1485	Learning refusal steering with principled null-space		
1486	constraint . <i>arXiv preprint arXiv:2506.07022</i> .		
1487			
1488	Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Han-		
1489	naneh Hajishirzi, Noah A. Smith, and Simon Shaolei		
1490	Du. 2024a. Decoding-time language model align-		
1491	ment with multiple objectives . In <i>ICML 2024 Work-</i>		
1492	shop on Theoretical Foundations of Foundation Mod-		
1493	els .		
1494			
1495	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon		
1496	Seo, Richard James, Mike Lewis, Luke Zettlemoyer,		
1497	and Wen-tau Yih. 2024b. REPLUG: Retrieval-		
1498	augmented black-box language models . In <i>Proceed-</i>		
1499	ings of the 2024 Conference of the North American		
1500	Chapter of the Association for Computational Lin-		
1501	guistics: Human Language Technologies (Volume		
	1: Long Papers) , pages 8371–8384, Mexico City,		
	Mexico. Association for Computational Linguistics.		
1502	Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric		
1503	Wallace, and Sameer Singh. 2020. AutoPrompt: Elic-		
1504	iting Knowledge from Language Models with Auto-		
1505	matically Generated Prompts . In <i>Proceedings of the</i>		
1506	<i>2020 Conference on Empirical Methods in Natural</i>		
1507	<i>Language Processing (EMNLP)</i> , pages 4222–4235,		
1508	Online. Association for Computational Linguistics.		
1509	Noah Shinn, Federico Cassano, Ashwin Gopinath,		
1510	Karthik R Narasimhan, and Shunyu Yao. 2023. Re-		
1511	flexion: language agents with verbal reinforcement		
1512	learning . In <i>Thirty-seventh Conference on Neural</i>		
1513	<i>Information Processing Systems</i> .		
1514	Kashun Shum, Shizhe Diao, and Tong Zhang. 2023.		
1515	Automatic prompt augmentation and selection with		
1516	chain-of-thought from labeled data . In <i>Findings</i>		
1517	<i>of the Association for Computational Linguistics:</i>		
1518	<i>EMNLP 2023</i> , pages 12113–12139, Singapore. Asso-		
1519	ciation for Computational Linguistics.		
1520	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-		
1521	mar. 2024. Scaling LLM test-time compute optimally		
1522	can be more effective than scaling model parameters .		
1523	In <i>International Conference on Machine Learning</i>		
1524	<i>(ICML)</i> .		
	Seongho Son, William Bankes, Sangwoong Yoon,		
	Shyam Sundhar Ramesh, Xiaohang Tang, and Il-		
	ija Bogunovic. 2025. Robust multi-objective con-		
	trolled decoding of large language models . <i>Preprint</i> ,		
	arXiv:2503.08796.		
	Feifan Song, Yuxuan Fan, Xin Zhang, Peiyi Wang, and		
	Houfeng Wang. 2025a. Instantly learning preference		
	alignment via in-context DPO . In <i>Proceedings of</i>		
	<i>the 2025 Conference of the Nations of the Americas</i>		
	<i>Chapter of the Association for Computational Lin-</i>		
	<i>guistics: Human Language Technologies (Volume 1:</i>		
	<i>Long Papers)</i> , pages 161–178, Albuquerque, New		
	Mexico. Association for Computational Linguistics.		
	Feifan Song, Bofei Gao, Yifan Song, Yi Liu, Weimin		
	Xiong, Yuyang Song, Tianyu Liu, Guoyin Wang, and		
	Houfeng Wang. 2025b. P-aligner: Enabling pre-		
	alignment of language models via principled instruc-		
	tion synthesis . <i>Preprint</i> , arXiv:2508.04626.		
	Samuel Soo, Chen Guang, Wesley Teng, Chan-		
	drasekaran Balaganesh, Tan Guoxian, and Yan Ming.		
	2025. Interpretable steering of large language mod-		
	els with feature guided activation additions . <i>Preprint</i> ,		
	arXiv:2501.09929.		
	Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M.		
	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,		
	Dario Amodei, and Paul Christiano. 2022. Learn-		
	ing to summarize from human feedback . <i>Preprint</i> ,		
	arXiv:2009.01325.		
	Niklas Stoehr, Kevin Du, Vésteinn Snæbjarnarson,		
	Robert West, Ryan Cotterell, and Aaron Schein. 2024.		
	Activation scaling for steering and interpreting lan-		
	guage models . In <i>Findings of the Association for</i>		
	<i>Computational Linguistics: EMNLP 2024</i> , pages		
	8189–8200, Miami, Florida, USA. Association for		
	Computational Linguistics.		
	Alessandro Stolfo, Vidhisha Balachandran, Safoora		
	Yousefi, Eric Horvitz, and Besmira Nushi. 2025. Im-		
	proving instruction-following in language models		
	through activation steering . In <i>The Thirteenth Inter-</i>		
	<i>national Conference on Learning Representations</i> .		
	Nishant Subramani, Nivedita Suresh, and Matthew Pe-		
	ters. 2022. Extracting latent steering vectors from		
	pretrained language models . In <i>Findings of the Asso-</i>		
	<i>ciation for Computational Linguistics: ACL 2022</i> ,		
	pages 566–581, Dublin, Ireland. Association for		
	Computational Linguistics.		
	Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao Yang,		
	Jiahao Qiu, Ming Yin, Mengdi Wang, Peter Bartlett,		
	and Andrea Zanette. 2024a. Fast best-of-n decoding		
	via speculative rejection . In <i>The Thirty-eighth An-</i>		
	<i>nnual Conference on Neural Information Processing</i>		
	<i>Systems</i> .		
	Hao Sun, Huailiang Peng, Qiong Dai, Xu Bai, and		
	Yanan Cao. 2025a. Layernavigator: Finding promis-		
	ing intervention layers for efficient activation steer-		
	ing in large language models . In <i>The Thirty-ninth An-</i>		
	<i>nnual Conference on Neural Information Processing</i>		
	<i>Systems</i> .		

1583	Jiashuo Sun, Xianrui Zhong, Sizhe Zhou, and Jiawei Han. 2025b. DynamicRAG: Leveraging outputs of large language model as feedback for dynamic reranking in retrieval-augmented generation . In <i>The Thirtieth Annual Conference on Neural Information Processing Systems</i> .	1640
1584		1641
1585		1642
1586		1643
1587		1644
1588		
1589	Simeng Sun, Yang Liu, Shuohang Wang, Dan Iter, Chenguang Zhu, and Mohit Iyyer. 2024b. PEARL: Prompting large language models to plan and execute actions over long documents . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 469–486, St. Julian’s, Malta. Association for Computational Linguistics.	1645
1590		1646
1591		1647
1592		1648
1593		
1594		
1595		
1596		
1597	Yushi Sun, Kai Sun, Yifan Ethan Xu, Xiao Yang, Xin Luna Dong, Nan Tang, and Lei Chen. 2025c. KERAG: Knowledge-enhanced retrieval-augmented generation for advanced question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 6194–6216, Suzhou, China. Association for Computational Linguistics.	1649
1598		1650
1599		1651
1600		1652
1601		1653
1602		1654
1603		1655
1604	Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2023. Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 4265–4293, Toronto, Canada. Association for Computational Linguistics.	1656
1605		1657
1606		1658
1607		1659
1608		1660
1609		1661
1610	Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. Function vectors in large language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1662
1611		1663
1612		1664
1613		1665
1614		1666
1615	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	1667
1616		1668
1617		1669
1618		1670
1619		1671
1620		1672
1621		1673
1622		1674
1623	Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2024. Steering language models with activation engineering . <i>Preprint</i> , arXiv:2308.10248.	1675
1624		1676
1625		1677
1626		1678
1627		1679
1628	Marco Valentino, Geonhee Kim, Dhairya Dalal, Zhixue Zhao, and André Freitas. 2025. Mitigating content effects on reasoning in language models through fine-grained activation steering . <i>Preprint</i> , arXiv:2505.12189.	1680
1629		1681
1630		1682
1631		1683
1632		1684
1633	Teun van der Weij, Massimo Poesio, and Nandi Schoots. 2024. Extending activation steering to broad skills and multiple behaviours . <i>Preprint</i> , arXiv:2403.05767.	1685
1634		1686
1635		1687
1636		1688
1637	Fei Wang, Xingchen Wan, Ruoxi Sun, Jiefeng Chen, and Sercan O Arik. 2025a. Astute RAG: Overcoming imperfect retrieval augmentation and knowledge	1689
1638		1690
1639		1691
		1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
		1703
		1704
		1705
		1706
		1707
		1708
		1709
		1710
		1711
		1712
		1713
		1714
		1715
		1716
		1717
		1718
		1719
		1720
		1721
		1722
		1723
		1724
		1725
		1726
		1727
		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
		1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
		1793
		1794
		1795
		1796
		1797
		1798
		1799
		1800

1697	large language models: A survey. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 19519–19529, Miami, Florida, USA. Association for Computational Linguistics.	
1698		
1699		
1700		
1701		
1702	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.	
1703		
1704		
1705		
1706		
1707		
1708		
1709	Wikipedia contributors. 2025a. Guidance, navigation, and control. https://en.wikipedia.org/wiki/Guidance%2C_navigation%2C_and_control . Accessed: 2025-12-01.	
1710		
1711		
1712		
1713	Wikipedia contributors. 2025b. Steering. https://en.wikipedia.org/wiki/Steering . Accessed: 2025-12-01.	
1714		
1715		
1716	Junde Wu, Jiayuan Zhu, Yunli Qi, Jingkun Chen, Min Xu, Filippo Menolascina, Yueming Jin, and Vicente Grau. 2025a. Medical graph RAG: Evidence-based medical large language model via graph retrieval-augmented generation. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 28443–28467, Vienna, Austria. Association for Computational Linguistics.	
1717		
1718		
1719		
1720		
1721		
1722		
1723		
1724		
1725	Lyucheng Wu, Mengru Wang, Ziwen Xu, Tri Cao, Nay Oo, Bryan Hooi, and Shumin Deng. 2025b. Automating steering for safe multimodal large language models. <i>Preprint</i> , arXiv:2507.13255.	
1726		
1727		
1728		
1729	Mingyan Wu, Zhenghao Liu, Yukun Yan, Xinze Li, Shi Yu, Zheni Zeng, Yu Gu, and Ge Yu. 2025c. RankCoT: Refining knowledge for retrieval-augmented generation through ranking chain-of-thoughts. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12857–12874, Vienna, Austria. Association for Computational Linguistics.	
1730		
1731		
1732		
1733		
1734		
1735		
1736		
1737	Ziting Xian, Jiawei Gu, Lingbo Li, and Shangsong Liang. 2025. MolRAG: Unlocking the power of large language models for molecular property prediction. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15513–15531, Vienna, Austria. Association for Computational Linguistics.	
1738		
1739		
1740		
1741		
1742		
1743		
1744	Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. 2023. Defending chatgpt against jailbreak attack via self-reminders. <i>Nature Machine Intelligence</i> , 5(12):1486–1496.	
1745		
1746		
1747		
1748		
1749	Ran Xu, Hui Liu, Sreyashi Nag, Zhenwei Dai, Yaochen Xie, Xianfeng Tang, Chen Luo, Yang Li, Joyce C. Ho, Carl Yang, and Qi He. 2025. SimRAG: Self-improving retrieval-augmented generation for adapting large language models to specialized domains. 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 11534–11550, Albuquerque, New Mexico. Association for Computational Linguistics.	1754
1750		1755
1751		1756
1752		1757
1753		1758
		1759
	Hanqi Yan, Yanzheng Xiang, Guangyi Chen, Yifei Wang, Lin Gui, and Yulan He. 2024. Encourage or inhibit monosemanticity? revisit monosemanticity from a feature decorrelation perspective. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 10423–10435, Miami, Florida, USA. Association for Computational Linguistics.	1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers. In <i>The Twelfth International Conference on Learning Representations</i> .	1768
		1769
		1770
		1771
		1772
	Dingkang Yang, Dongling Xiao, Jinjie Wei, Mingcheng Li, Zhaoyu Chen, Ke Li, and Lihua Zhang. 2025a. Improving factuality in large language models via decoding-time hallucinatory and truthful comparators. In <i>Proceedings of the Thirty-Ninth AAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'25/IAAI'25/EAAI'25</i> . AAAI Press.	1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
	Dingkang Yang, Dongling Xiao, Jinjie Wei, Mingcheng Li, Zhaoyu Chen, Ke Li, and Lihua Zhang. 2025b. Improving factuality in large language models via decoding-time hallucinatory and truthful comparators. In <i>Proceedings of the Thirty-Ninth AAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'25/IAAI'25/EAAI'25</i> . AAAI Press.	1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
	Jinghan Yang, Zhenbo Xu, Dehua Ma, Liu Liu, Fei Liu, Gong Huang, and Zhaofeng He. 2025c. Reciperag: Advancing recipe generation with reinforced retrieval augmented generation. In <i>Proceedings of the 33rd ACM International Conference on Multimedia, MM '25</i> , page 5060–5069, New York, NY, USA. Association for Computing Machinery.	1793
		1794
		1795
		1796
		1797
		1798
		1799
	Kevin Yang and Dan Klein. 2021. FUDGE: Controlled text generation with future discriminators. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3511–3535, Online. Association for Computational Linguistics.	1800
		1801
		1802
		1803
		1804
		1805
		1806
	Zairun Yang, Yilin Wang, Zhengyan Shi, Yuan Yao, Lei Liang, Keyan Ding, Emine Yilmaz, Huajun Chen, and Qiang Zhang. 2025d. EventRAG: Enhancing LLM generation with event knowledge graphs. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1:</i>	1807
		1808
		1809
		1810
		1811
		1812

1813		Vienna, Austria. Association for Computational Linguistics.	1870
1814			1871
1815	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex	1872
1816	Thomas L. Griffiths, Yuan Cao, and Karthik R	Smola. 2023. Automatic chain of thought prompting	1873
1817	Narasimhan. 2023a. Tree of thoughts: Deliberate	in large language models . In <i>The Eleventh International</i>	1874
1818	problem solving with large language models . In	Conference on Learning Representations .	1875
1819	<i>Thirty-seventh Conference on Neural Information</i>		
1820	<i>Processing Systems</i> .	Zhuoxuan Zhang, Jinhao Duan, Edward Kim, and Kaidi	1876
1821	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	Xu. 2025f. Sparse neurons carry strong signals of	1877
1822	Shafran, Karthik Narasimhan, and Yuan Cao. 2023b.	question ambiguity in LLMs . In <i>Proceedings of the</i>	1878
1823	ReAct: Synergizing reasoning and acting in language	<i>2025 Conference on Empirical Methods in Natural</i>	1879
1824	models. In <i>International Conference on Learning</i>	<i>Language Processing</i> , pages 16092–16110, Suzhou,	1880
1825	<i>Representations (ICLR)</i> .	China. Association for Computational Linguistics.	1881
1826	Hao Ye, Mengshi Qi, Zhaohong Liu, Liang Liu, and	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	1882
1827	Huadong Ma. 2025. Safedriverag: Towards safe	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	1883
1828	autonomous driving with knowledge graph-based	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen	1884
1829	retrieval-augmented generation . In <i>Proceedings of</i>	Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,	1885
1830	<i>the 33rd ACM International Conference on Multime-</i>	Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and	1886
1831	<i>dia</i> , MM '25, page 11170–11178, New York, NY,	3 others. 2025a. A survey of large language models .	1887
1832	USA. Association for Computing Machinery.	<i>Preprint</i> , arXiv:2303.18223.	1888
1833	Jianyi Zhang, Da-Cheng Juan, Cyrus Rashtchian, Chun-	Xinping Zhao, Dongfang Li, Yan Zhong, Boren Hu,	1889
1834	Sung Ferng, Heinrich Jiang, and Yiran Chen. 2024.	Yibin Chen, Baotian Hu, and Min Zhang. 2024a.	1890
1835	Sled: self logits evolution decoding for improving	SEER: Self-aligned evidence extraction for retrieval-	1891
1836	factuality in large language models. In <i>Proceedings</i>	augmented generation . In <i>Proceedings of the 2024</i>	1892
1837	<i>of the 38th International Conference on Neural In-</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	1893
1838	<i>formation Processing Systems</i> , NIPS '24, Red Hook,	<i>guage Processing</i> , pages 3027–3041, Miami, Florida,	1894
1839	NY, USA. Curran Associates Inc.	USA. Association for Computational Linguistics.	1895
1840	Qinggang Zhang, Zhishang Xiang, Yilin Xiao, Le Wang,	Xufeng Zhao, Mengdi Li, Wenhao Lu, Cornelius Weber,	1896
1841	Junhui Li, Xinrun Wang, and Jinsong Su. 2025a.	Jae Hee Lee, Kun Chu, and Stefan Wermter. 2024b.	1897
1842	FaithfulRAG: Fact-level conflict modeling for	Enhancing zero-shot chain-of-thought reasoning in	1898
1843	context-faithful retrieval-augmented generation . In	large language models through logic . In <i>Proceedings</i>	1899
1844	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	<i>of the 2024 Joint International Conference on Compu-</i>	1900
1845	<i>sociation for Computational Linguistics (Volume 1:</i>	<i>tational Linguistics, Language Resources and Eval-</i>	1901
1846	<i>Long Papers)</i> , pages 21863–21882, Vienna, Austria.	<i>uation (LREC-COLING 2024)</i> , pages 6144–6166,	1902
1847	Association for Computational Linguistics.	Torino, Italia. ELRA and ICCL.	1903
1848	Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang,	Zekai Zhao, Qi Liu, Kun Zhou, Zihan Liu, Yifei	1904
1849	Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo,	Shao, Zhiting Hu, and Biwei Huang. 2025b. Ac-	1905
1850	Yufei Wang, Niklas Muennighoff, Irwin King, Xue	tivation control for efficiently eliciting long chain-	1906
1851	Liu, and Chen Ma. 2025b. A survey on test-time	of-thought ability of language models . <i>Preprint</i> ,	1907
1852	scaling in large language models: What, how, where,	arXiv:2505.17697.	1908
1853	and how well? <i>Preprint</i> , arXiv:2503.24235.	Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and	1909
1854	Xiaoyun Zhang, Zhengyue Zhao, Wenxuan Shi, Kaidi	Dacheng Tao. 2024. ROSE doesn't do that: Boosting	1910
1855	Xu, Di Huang, and Xing Hu. 2025c. Safety	the safety of instruction-tuned large language models	1911
1856	alignment of large language models via contrast-	with reverse prompt contrastive decoding . In <i>Find-</i>	1912
1857	ing safe and harmful distributions . <i>Preprint</i> ,	<i>ings of the Association for Computational Linguistics:</i>	1913
1858	arXiv:2406.16743.	<i>ACL 2024</i> , pages 13721–13736, Bangkok, Thailand.	1914
1859	Yuqi Zhang, Liang Ding, Lefei Zhang, and Dacheng	Association for Computational Linguistics.	1915
1860	Tao. 2025d. Intention analysis makes LLMs a good	Tianqi Zhong, Quan Wang, Jingxuan Han, Yongdong	1916
1861	jailbreak defender . In <i>Proceedings of the 31st Inter-</i>	Zhang, and Zhendong Mao. 2023. Air-decoding: At-	1917
1862	<i>national Conference on Computational Linguistics</i> ,	tribute distribution reconstruction for decoding-time	1918
1863	pages 2947–2968, Abu Dhabi, UAE. Association for	controllable text generation . In <i>Proceedings of the</i>	1919
1864	Computational Linguistics.	<i>2023 Conference on Empirical Methods in Natural</i>	1920
1865	Zheng Zhang, Shaocheng Lan, Lei Song, Jiang Bian,	<i>Language Processing</i> , pages 8233–8248, Singapore.	1921
1866	Yexin Li, and Kan Ren. 2025e. Learning to select	Association for Computational Linguistics.	1922
1867	in-context demonstration preferred by large language	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,	1923
1868	model . In <i>Findings of the Association for Computa-</i>	Haohan Wang, and Yu-Xiong Wang. 2024a. Lan-	1924
1869	<i>tional Linguistics: ACL 2025</i> , pages 11345–11360,	guage agent tree search unifies reasoning, acting, and	1925
		planning in language models . In <i>Proceedings of the</i>	1926

1927
1928

1929
1930
1931
1932
1933
1934

1935
1936
1937
1938

1939
1940
1941
1942
1943
1944
1945
1946

1947
1948
1949
1950
1951
1952

1953
1954
1955
1956
1957

1958
1959
1960
1961
1962
1963
1964
1965

1966
1967
1968
1969
1970
1971
1972
1973

1974

1975
1976
1977
1978
1979
1980

41st International Conference on Machine Learning, ICML'24. JMLR.org.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. [LIMA: Less is more for alignment](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023b. [Large language models are human-level prompt engineers](#). *Preprint*, arXiv:2211.01910.

Yujun Zhou, Yufei Han, Haomin Zhuang, Kehan Guo, Zhenwen Liang, Hongyan Bao, and Xiangliang Zhang. 2024b. [Defending jailbreak prompts via in-context adversarial game](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 20084–20105, Miami, Florida, USA. Association for Computational Linguistics.

Lanyun Zhu, Deyi Ji, Tianrun Chen, Peng Xu, Jieping Ye, and Jun Liu. 2024. [Ibd: Alleviating hallucinations in large vision-language models via image-biased decoding](#). *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1615–1624.

Jin Ziqi and Wei Lu. 2023. [Tab-CoT: Zero-shot tabular chain of thought](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10259–10277, Toronto, Canada. Association for Computational Linguistics.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, and 2 others. 2023. [Representation engineering: A top-down approach to ai transparency](#). *Preprint*, arXiv:2310.01405.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, and 2 others. 2025. [Representation engineering: A top-down approach to ai transparency](#). *Preprint*, arXiv:2310.01405.

A Scope and Rationale of this Survey

Rationale: Unifying a Fragmented Landscape.

As illustrated in Table 1, the field of LLM adaptation is undergoing a fundamental shift from parameter-heavy training to flexible test-time interventions. However, the current literature treats these interventions as distinct, unrelated fields.

Researchers working on *Prompt Engineering* (Input Space) often operate independently from those studying *Mechanistic Interpretability* (Latent Space) or *Constrained Decoding* (Decoding Space). Existing surveys typically focus either on isolated subfields (Schulhoff et al., 2025; Rai et al., 2025) or on specific objectives like safety (Pan et al., 2025) or efficiency (Zhang et al., 2025b), obscuring the underlying control mechanisms. This survey is essential to bridge these silos. By proposing the first mechanism-centric taxonomy, we reveal how the same mathematical principles can be applied across different stages of generation, providing researchers with a "full-stack" view of model controllability that existing goal-oriented reviews lack.

Scope: Defining Training-Free Guidance. To ensure clarity, we strictly define Training-Free Test-Time Guidance based on two criteria: (1) Frozen Parameters, meaning the parameters of the backbone LLM are strictly frozen and cannot be trained; and (2) Inference Intervention, where the method actively alters the model’s generation probability or trajectory to satisfy specific constraints. Under this definition, our review covers the entire inference lifecycle: Input Space, Latent Space, Decoding Space, and Output Space. We explicitly exclude system-level optimizations focused solely on throughput and permanent weight-editing techniques, as they aim for computational efficiency rather than control.

B Terminology: Why We Use the Term “Guidance”

In this survey, we use *guidance* as an umbrella term for LLM-training-free, test-time mechanisms that inject additional signals into the inference process (in the input space, latent space, or decoding process) to direct a frozen model toward a behavioural objective, without updating its parameters. This notion is close to the role of a guidance module in control systems, which specifies the desired trajectory of a system given its current state and mission goal, distinct from both navigation and low-level control (Wikipedia contributors, 2025a).

We deliberately distinguish *guidance* from several related terms. *Steering* typically denotes changing the direction of motion (Wikipedia contributors, 2025b) and, in the LLM literature, is often used for activation-level representation engineering—adding or modifying directions in hidden space to control behaviour (Zou et al., 2025). *Interven-*

1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995

1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010

2011
2012

2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030

Table 1: Comparison with existing surveys related to LLM adaptation. ✓ = Primary focus with detailed taxonomy; ● = Partial coverage or sub-topic discussion; ✗ = Little to no coverage. **Our survey** is the first to provide a unified framework covering the entire test-time inference lifecycle.

Survey	Core Focus	Intervention Space (Lifecycle)				Unified View
		Input	Latent	Decoding	Output	
General & Training-Based Surveys						
Zhao et al. (Zhao et al., 2025a)	General LLM Capabilities	●	✗	✗	●	✗
Tie et al. (Lai et al., 2025)	Post-Training Scaling	✗	✗	✗	✗	✗
Subfield-Specific Surveys						
Schulhoff et al. (Schulhoff et al., 2025)	Prompt Engineering Techniques	✓	✗	✗	✗	✗
Fan et al. (Fan et al., 2024)	Retrieval-Augmented Generation	✓	✗	✗	✗	✗
Rai et al. (Rai et al., 2025)	Mechanistic Interpretability	✗	✓	✗	✗	✗
Goal-Oriented Test-Time Surveys						
Zhang et al. (Zhang et al., 2025b)	Test-Time Scaling	●	✗	●	✓	●
Pan et al. (Pan et al., 2025)	Training-Free Alignment	✓	●	✓	●	✗
Huang et al. (Huang et al., 2025a)	Human-AI Cooperation	✓	✗	✗	●	✗
Ours	Test-Time Guidance Mechanism	✓	✓	✓	✓	✓

tion usually emphasises the concrete act of editing internal activations at specific layers or tokens during inference (Li et al., 2023a). *Mitigation* describes a safety-oriented objective (e.g., reducing hallucinations or toxicity), rather than a particular mechanism.

In this paper, we therefore use *guidance* as the method-agnostic, trajectory-oriented term for test-time control of LLMs; we treat *steering* and *interventions* as specific implementation families within this space, and *mitigation* as one important class of guidance objectives (safety and robustness).

C A Unified Guidance Theoretical Framework from a Conditional Distribution Shaping Perspective

To synthesize the diverse methods categorized in our taxonomy, we propose a unified framework based on Conditional Distribution Shaping (Ji et al., 2024; Hu et al., 2025a; Li et al., 2025c). Fundamentally, a pre-trained LLM defines a base autoregressive distribution $P_\theta(y|x)$. The objective of any training-free guidance method is to construct a transformed distribution $\tilde{P}(y|x;g)$ that satisfies specific constraints imposed by the guidance signal g .

We formalize this transformation as the application of a shaping operator \mathcal{T} at a specific interface \mathcal{I} of the model’s computational graph (as illustrated

in Figure 6):

$$\tilde{P}(y|x;g) = \mathcal{T}_{\mathcal{I}}(P_\theta(\cdot|x),g) \quad (8)$$

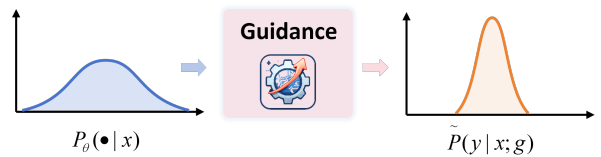


Figure 6: Schematic of the Unified Guidance Theoretical Framework from a Conditional Distribution Shaping Perspective.

The four intervention spaces in our taxonomy correspond to applying \mathcal{T} at four distinct depth levels of the generative process.

1. Input Space ($\mathcal{T}_{\text{input}}$): In input-level guidance, the shaping operator intervenes on the conditioning variable itself. Rather than modifying the model’s mapping function, it seeks a perturbation δ in the discrete token space (or embedding space) such that the base model naturally projects the input to the desired output region.

$$\tilde{P}(y|x;g) = P_\theta(y | \phi(x,g)) \quad (9)$$

where $\phi(x,g)$ represents the transformation function (e.g., prepending a system prompt or retrieving context). Theoretically, this leverages the pre-training manifold continuity: it relies on the assumption that there exists a neighborhood $\phi(x,g)$

2077 in the context space that maps to the target dis-
2078 tribution \tilde{P} solely through the frozen parameters
2079 θ .

2080 **2. Latent Space ($\mathcal{T}_{\text{latent}}$):** Latent interventions
2081 operate by warping the intermediate feature map-
2082 ping. Let $F_\ell(h_{\ell-1})$ be the transformation at layer ℓ .
2083 Guidance introduces an affine transformation (steer-
2084 ing vector v) to the hidden state h_ℓ . The shaped
2085 distribution is the result of forward-propagating
2086 this perturbed state:

$$2087 \quad \tilde{P}(y|x; g) = \text{Head}(\dots F_L(h_\ell + \lambda v_g) \dots) \quad (10)$$

2088 Unlike input guidance which is discrete, this op-
2089 erator performs a geometric translation on the se-
2090 mantic manifold. By shifting the representation h_ℓ
2091 along the direction v_g , the operator $\mathcal{T}_{\text{latent}}$ biases
2092 the subsequent layers to process the context as if
2093 it possessed the attribute g , effectively reshaping
2094 the downstream probability surface before it is pro-
2095 jected to the vocabulary.

2096 **3. Decoding Space ($\mathcal{T}_{\text{decoding}}$):** Decoding guid-
2097 ance acts directly on the output probability mass
2098 function at each step t . The operator \mathcal{T} transforms
2099 the logits z_t into \tilde{z}_t via an additive or masking term
2100 $\psi(g)$, followed by re-normalization:

$$2101 \quad \tilde{P}(y_t|y_{<t}, x; g) = \text{Softmax}(z_t + \lambda \cdot \psi(y_t, g)) \quad (11)$$

2102 Here, ψ can be a constraint mask or a contrastive
2103 bias. This represents a local re-weighting strat-
2104 egy. The shaping operator explicitly "carves" the
2105 distribution, suppressing the probability mass of
2106 tokens that violate g and re-allocating it to com-
2107 pliant tokens, ensuring strict or soft adherence to
2108 constraints at the token level.

2109 **4. Output Space ($\mathcal{T}_{\text{output}}$):** Output guidance op-
2110 erates on the distribution over complete sequences.
2111 Theoretically, this can be viewed as applying a se-
2112 lection filter over the support set \mathcal{Y} . The shaped
2113 distribution becomes:

$$2114 \quad \tilde{P}(y|x; g) = \frac{P_\theta(y|x) \cdot I[R(y, g) > \tau]}{Z} \quad (12)$$

2115 where I is an indicator function (or soft score) from
2116 a verifier R , and Z is the normalization constant.
2117 This operator performs global pruning. It shapes
2118 the distribution by truncating the tail and concen-
2119 trating the probability mass solely on the candidate
2120 y^* that maximizes the global verification score.

2121 D Literature Search Methodology

2122 To ensure a comprehensive coverage of the land-
2123 scape, we conducted a systematic search across
2124 major academic databases, including ACL Anthol-
2125 ogy, Google Scholar, and Semantic Scholar Com-
2126 plementing these traditional repositories, we ac-
2127 tively monitored real-time research trends using
2128 AI-driven discovery platforms such as PaSa¹, Pa-
2129 pers.cool², and Hugging Face Papers³. We also
2130 employed a "snowballing" strategy, leveraging bib-
2131 liographies from existing related surveys (as dis-
2132 cussed in Table 1) to identify relevant training-free
2133 mechanisms. We strictly filtered papers based on
2134 the criteria of keeping the backbone model frozen
2135 and involving active inference-time intervention.
2136 The complete consolidated list of surveyed litera-
2137 ture, categorized by intervention space, is visual-
2138 ized in Figure 7.

¹<https://pasa-agent.ai/>

²<https://papers.cool/>

³<https://huggingface.co/papers/trending>

I. Input-Space Guidance

Contextual Augmentation: IC-RAG (Ram et al., 2023), REPLUG (Shi et al., 2024b), CoK (Li et al., 2024b), RQ-RAG (Chan et al., 2024), GMR (Lee et al., 2022), RAG-Fusion (Rackauckas, 2024), KRAGEN (Matsumoto et al., 2024), SimRAG (Xu et al., 2025), SEER (Zhao et al., 2024a), L-RAG (Lin et al., 2025a), HopRAG (Liu et al., 2025a), KERAG (Sun et al., 2025c), TC-RAG (Jiang et al., 2025b), HyKGE (Jiang et al., 2025c), UniRAG (Li et al., 2025e), OmniRAGMed (Chen et al., 2025c), MolRAG (Xian et al., 2025), EventRAG (Yang et al., 2025d), KiRAG (Fang et al., 2025a), FaithfulRAG (Zhang et al., 2025a), Dialogue-RAG (Li et al., 2025d), SGIC (Chen et al., 2025a), MedGraphRAG (Wu et al., 2025a), Astute RAG (Wang et al., 2025a), DualRAG (Cheng et al., 2025), T-GRAG (Li et al., 2025b), SafeDriveRAG (Ye et al., 2025), RecipeRAG (Yang et al., 2025c), HM-RAG (Liu et al., 2025c), DynamicRAG (Sun et al., 2025b), What-makes-ICL (Liu et al., 2022), Lost-in-Middle (Liu et al., 2024a), L2R-ICL (Wang et al., 2024a), CED-Select (Iter et al., 2023), RGER (Lin et al., 2025b), GENICL (Zhang et al., 2025e), MLSM/TTF (Liu et al., 2025b), UniICL (Gao et al., 2025), FLARE (Jiang et al., 2023b), CoVe (Wang et al., 2024c), ToG2 (Ma et al., 2025).

Prompt Design and Optimization: Chain-of-Thought (Wei et al., 2022), Zero-shot CoT (Kojima et al., 2022), Self-Consistency (Wang et al., 2023), Tab-CoT (Ziqi and Lu, 2023), Auto-CoT (Zhang et al., 2023), Logical CoT (Zhao et al., 2024b), CDW-CoT (Fang et al., 2025b), RankCoT (Wu et al., 2025c), PCoT (Modzelewski et al., 2025), SemCoT (He et al., 2025a), APE (Zhou et al., 2023b), OPRO (Yang et al., 2024), BPO (Cheng et al., 2024a), RLPrompt (Deng et al., 2022), PromptAgent (Wang et al., 2024b), PEARL-plan (Sun et al., 2024b), PromptBreeder (Fernando et al., 2024), Structure-Guided (Cheng et al., 2024b), AutoPrompt (Shin et al., 2020), Automate-CoT (Shum et al., 2023).

Alignment and Safety Prompting: Self-Reminders (Xie et al., 2023), URIAL (Lin et al., 2024), LLaMA-2 cfg. (Touvron et al., 2023), ICAG (Zhou et al., 2024b), SmoothLLM (Robey et al., 2025), ROSE (Zhong et al., 2024), Instructive Decoding (Alajrami et al., 2025), Intention Analysis (Zhang et al., 2025d), CIP (Hahm et al., 2025), LLM-Self-Defense (Phute et al., 2024), SelfDefenD (Wang et al., 2025f), P-Aligner (Song et al., 2025b), ICDPO (Song et al., 2025a).

II. Latent-Space Guidance

Global Directional Arithmetic: Activation Engineering (Turner et al., 2024), PPLM (Dathathri et al., 2020), RepE (Zou et al., 2023), ITI (Li et al., 2023a), In-Context Vectors (Liu et al., 2024b), Function Vectors (Todd et al., 2024), Latent steering vectors (Subramani et al., 2022), IDANI (Antverg et al., 2022), Scaling Laws (Stoehr et al., 2024), Layer Sensitivity (Stolfo et al., 2025), Fractional Reasoning (Liu et al., 2025d).

Decomposition and Fine-Grained Control: SAE-based Steering (Huben et al., 2024; Yan et al., 2024), Interpretable Steering (Soo et al., 2025), Sparse Neurons (Zhang et al., 2025f), Neuron Dampening (Huang et al., 2025c), Dual Mechanisms (Han et al., 2025), Steering Target Atoms (Wang et al., 2025c), Concept Clamping (Valentino et al., 2025; Sakarvadia et al., 2025).

Geometric Constraints and Projection: AlphaSteer (Sheng et al., 2025), ASGuard (Park et al., 2025), Projection Steering (Postmus and Abreu, 2024), FLORAIN (Jiang et al., 2025a), Low-Rank Interventions (Oozeer et al., 2025), Diffusion Control (Rodriguez et al., 2024).

Context-Adaptive Dynamics: NL-ITI (Hoscilowicz et al., 2024), EAST (Rahn et al., 2024), SADI (Wang et al., 2025e), Adaptive Steering (Wang et al., 2025d), SCANS (Cao et al., 2025), Preemptive Detection (Karnik and Bansal, 2025), AutoSteer (Multimodal) (Wu et al., 2025b), Layer-Navigator (Sun et al., 2025a), DSAS (Do et al., 2025).

III. Decoding-Space Guidance

Score-based Reweighting: GeDi (Krause et al., 2021), DExperts (Liu et al., 2021), Contrastive Decoding (Li et al., 2023b), DoLa (Chuang et al., 2024), UACD (Lee et al., 2025b), SLED (Zhang et al., 2024), ROSE (Zhong et al., 2024), Instructive Decoding (Alajrami et al., 2025), ACD (Zhang et al., 2025c), MOD (Shi et al., 2024a), DeAL (Huang et al., 2025b), PAD (Chen et al., 2025b), Drift (Kim et al., 2025), SAFEINFER (Banerjee et al., 2025), CDT (Yang et al., 2025b), RMOD (Son et al., 2025), IBD (Zhu et al., 2024), Nudging (Fei et al., 2025), DeRa (Liu et al., 2024c), MCA (Fu et al., 2025), Air-Decoding (Zhong et al., 2023), FUDGE (Yang and Klein, 2021).

Constraint-based Filtering: NeuroLogic Decoding (Lu et al., 2021), Grammar-Constrained (Raspanti et al., 2025), CARE (Hu et al., 2025b), CARDS (Li et al., 2025a), RSD (Liao et al., 2025), RGD (Mañas et al., 2025), PRGD (Nguyen et al., 2025).

IV. Output-Space Guidance

Search and Planning: ToT (Yao et al., 2023a), GoT (Besta et al., 2024), RAP (Hao et al., 2023), Multi-Agent Debate/LATS (Zhou et al., 2024a).

Iterative Refinement: Reflexion (Shinn et al., 2023), Self-Refine (Madaan et al., 2023), CRITIC (Gou et al., 2024).

Candidate Selection and Reranking: Training Verifiers (Cobbe et al., 2021), SelfCheckGPT (Manakul et al., 2023), MBR Decoding (Suzgun et al., 2023), Universal Self-Consistency (Chen et al., 2024), Speculative Rejection (Sun et al., 2024a), Self-Evaluation (Ren et al., 2023).

Multi-Model Collaboration: COOPERATE (Feng et al., 2024), LLM-Blender (Jiang et al., 2023a), Selection-Inference (Creswell et al., 2023), Mutual Reasoning (Qi et al., 2025a).

Figure 7: A comprehensive summary of the surveyed literature, organized by the taxonomy detailed in Section 3.