

Confidence as Control: A Survey of Confidence Utilization in Large Language Models

Anonymous ACL submission

Abstract

Confidence in LLMs is often studied through uncertainty estimation and calibration. We survey a complementary perspective: confidence as a control signal that governs system behavior. We organize **confidence utilization** across the LLM lifecycle: (i) training (data selection, loss weighting, self-training, and preference optimization); (ii) inference (candidate selection, adaptive computation, and confidence-guided contrastive decoding); and (iii) deployment (cost-aware routing and cascading, RAG control (retrieval triggering, context filtering, and parametric-retrieval arbitration), and risk-aware abstention and monitoring. We unify these techniques into a framework that turns confidence into system decisions, with implications for efficiency and reliability.

1 Introduction

Confidence in Large Language Models (LLMs) bridges the gap between internal uncertainty and actionable system behavior. A model that knows what it doesn't know is useful; a system that acts on this knowledge—retrieving when uncertain, deferring when unreliable, focusing learning where needed—is transformative. This survey concerns the latter. Recent work has demonstrated that LLMs can express calibrated uncertainty through verbalized probabilities (Tian et al., 2023; Kadavath et al., 2022), that semantic entropy over meaning clusters detects hallucinations (Farquhar et al., 2024), and that self-consistency among sampled responses provides robust confidence signals (Wang et al., 2022; Manakul et al., 2023). These estimation advances have successfully equipped models with the ability to signal their own reliability.

At the same time, several surveys have examined confidence and uncertainty in LLMs, predominantly focusing on estimation and calibration. Geng et al. (2024) provide a foundational taxonomy of confidence estimation and calibration tech-

niques, while Shorinwa et al. (2025) and Liu et al. (2025b) offer comprehensive coverage of uncertainty quantification methods and their applications. Xiong et al. (2023) and Xie et al. (2024) address the specific challenges of model-agnostic settings where model internals are inaccessible. These surveys answer: how can we obtain high-quality confidence scores? Yet a good confidence score is not the end goal, it is the prerequisite. A systematic treatment of how confidence should actually govern system behavior remains absent.

Our survey addresses the critical next step: **how should systems utilize confidence?** As illustrated in Figure 1, we propose a taxonomy¹ where confidence functions as *control*—not a passive measure of uncertainty, but an active governor that determines what to learn, how to reason, and when to defer. We trace this utilization across the full LLM lifecycle including: (i) training-time applications in data curation and alignment (§2), (ii) inference-time control over reasoning and decoding (§3), and (iii) deployment-time decisions including model selection (§4), retrieval-augmented generation (RAG) (§5), and risk management (§6). By synthesizing these distinct domains, we aim to offer a unified perspective for designing systems that act on uncertainty rather than merely quantify it.

2 Confidence-Aware Training

Confidence signals shape the training at multiple stages: selecting which data to train on, modulating how models learn from that data, and guiding preference optimization for alignment. We organize methods by where confidence intervenes: data curation (§2.1), supervised fine-tuning (§2.2), and reinforcement learning from human feedback (§2.3).

2.1 Confidence-aware data selection

Confidence-aware data selection uses model uncertainty as a proxy for data value, varying in con-

¹Across this survey, we employ standard field notations summarized in Appendix A without explicit re-definition.

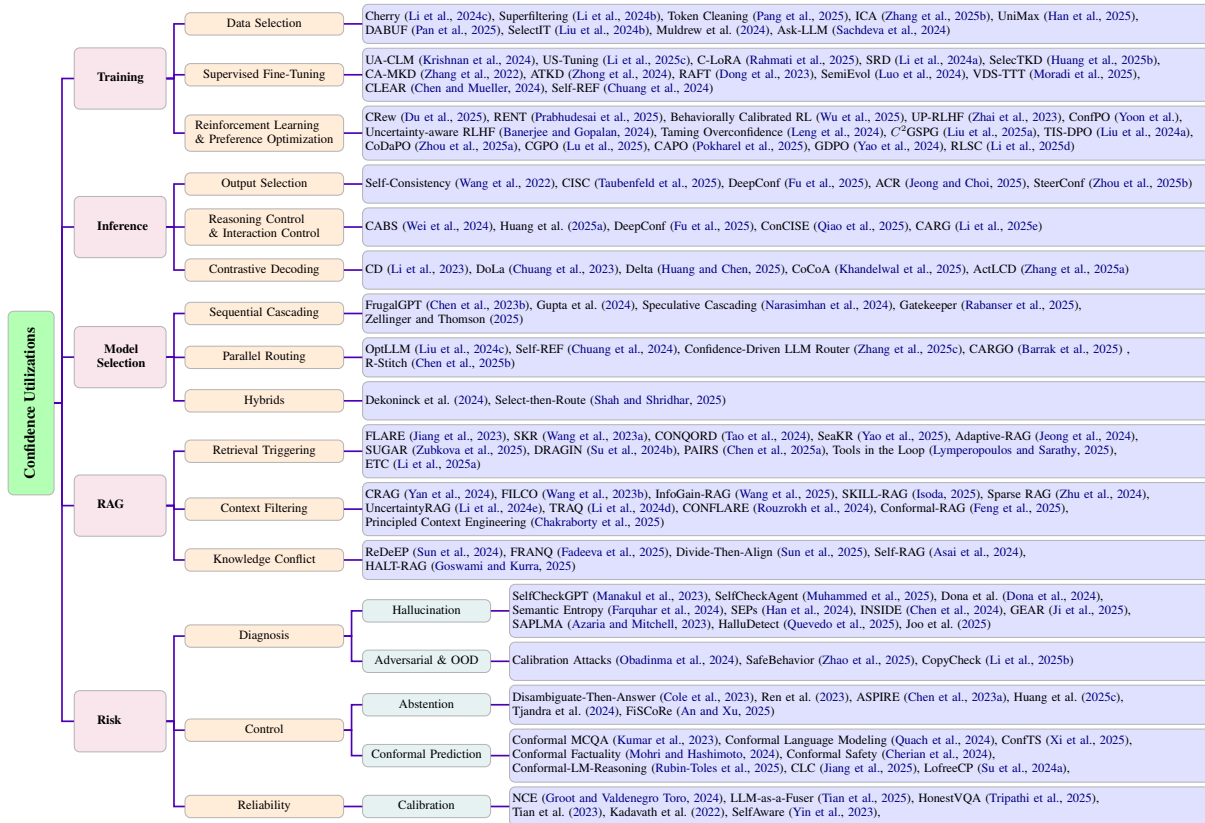


Figure 1: Taxonomy of Confidence Utilization in LLMs

080 fidence source (self-evaluation, external models, 107
 081 LLM-as-judge), granularity (sample vs. token), and 108
 082 selection goal (quality filtering, influence estima- 109
 083 tion, active acquisition). 110

084 **Difficulty as inverse confidence.** Low confi- 111
 085 dence often signals learning potential. Li et al. 112
 086 (2024c) define Instruction-Following Difficulty 113
 087 $IFD(x, y) = \mathcal{L}_\theta(y|x) / \mathcal{L}_\theta(y)$, isolating instruc- 114
 088 tion-following difficulty from baseline response com- 115
 089 plexity. Li et al. (2024b) show IFD rankings 116
 090 transfer across scales—GPT-2’s orderings correlate 117
 091 with LLaMA2-7B/13B, enabling efficient filtering. 118
 092 Pang et al. (2025) extend this to token granular- 119
 093 ity by comparing base-model and clean-reference 120
 094 losses: tokens are retained when the base model 121
 095 struggles but the reference succeeds (learnable con- 122
 096 tent), and filtered when both show high loss (noise). 123

097 **Influence and acquisition.** Zhang et al. (2025b) 124
 098 estimate sample utility via In-Context Approxima- 125
 099 tion—conditioning on a holdout set simulates gra- 126
 100 dient steps without retraining. Han et al. (2025) 127
 101 frame selection as graph influence maximization 128
 102 with epistemic uncertainty balancing influence 129
 103 against redundancy; Pan et al. (2025) apply attribu- 130
 104 tion in reverse, tracing harmful outputs to training 131
 105 samples for removal. For label-efficient settings, 132
 106 Muldrew et al. (2024) combine predictive entropy 133

with preference model certainty to guide annotation 107
 acquisition; Sachdeva et al. (2024) query an LLM 108
 directly, using P(YES) as quality score. Liu et al. 109
 (2024b) exploit multi-granularity uncertainty (to- 110
 ken, sentence, and model-level), selecting samples 111
 where self-assessment is confident and consistent. 112

2.2 Confidence-aware fine-tuning 113

Data selection determines what to train on; confi- 114
 dence also reshapes how models learn. Traditional 115
 supervised fine-tuning (SFT) treats tokens equally, 116
 yet highly predictable tokens (articles, connectives) 117
 and already-learned patterns provide diminishing 118
 learning signal. Confidence-modulated methods 119
 address this by scaling each token’s contribution to 120
 the loss based on certainty. The central design ques- 121
 tion: should confident predictions receive more 122
 weight (reinforcing what the model knows) or less 123
 (focusing learning on uncertain regions)? 124

Confidence-weighted training. Krishnan et al. 125
 (2024) decompose uncertainty into aleatoric and 126
 epistemic components, using the epistemic term 127
 to emphasize tokens where model knowledge is 128
 lacking. Li et al. (2025c) target a different goal: 129
 training LLMs to recognize knowledge boundaries 130
 and abstain rather than hallucinate. Rahmati et al. 131
 (2025) take a Bayesian approach, conditioning 132
 LoRA adapter distributions on each input so that 133

uncertainty estimates become sample-specific, improving calibration in few-shot settings.

Knowledge distillation with confidence. Naive distillation treats all teacher outputs as equally reliable, but teachers are often uncertain on ambiguous tokens. [Huang et al. \(2025b\)](#) filter via propose-and-verify: students propose tokens, teachers accept or reject based on distribution agreement, and only accepted tokens receive distillation signal. [Li et al. \(2024a\)](#) curate data through student self-reflection, ranking samples by difficulty and selecting those the student can realistically learn from. [Zhang et al. \(2022\)](#) weight multi-teacher ensembles by per-sample confidence. [Zhong et al. \(2024\)](#) distinguish easy from hard tokens: easy tokens skip imitation-focused teaching while hard tokens receive richer, more diverse supervision.

Confidence in self-training. Self-training relies on model-generated pseudo-labels, making confidence essential for filtering unreliable outputs. [Dong et al. \(2023\)](#) generate multiple responses per prompt, rank by reward, and train only on top-scoring samples. [Chen and Mueller \(2024\)](#) automate curation more broadly: confidence estimates identify low-quality training pairs, which are either filtered or replaced with higher-confidence LLM-generated alternatives. [Luo et al. \(2024\)](#) extend rejection sampling to semi-supervised settings; [Moradi et al. \(2025\)](#) add verification to detect when self-generated data helps versus harms. Beyond filtering, [Chuang et al. \(2024\)](#) train models to output calibrated confidence tokens, enabling routing to stronger models when local confidence is low.

2.3 Confidence-aware preference optimization

Preference optimization aligns LLMs with human values, but reward models can inherit annotation noise and policies exploit blind spots. Given the objective $\max_{\theta} \mathbb{E}[r(x, y)] - \beta \cdot \text{KL}(\pi_{\theta} \parallel \pi_{\text{ref}})$, confidence integrates as the reward itself, as an uncertainty penalty, or as granular weights over tokens, trajectories, or preference pairs.

Confidence as reward. [Du et al. \(2025\)](#) use token-level log-probability directly as reward, while [Prabhudesai et al. \(2025\)](#) apply entropy minimization for reasoning tasks. Naive confidence maximization risks reward hacking; [Leng et al. \(2024\)](#) counter this by calibrating reward models during PPO via confidence-augmented preference data and running-average reward thresholds.

Uncertainty penalties. [Zhai et al. \(2023\)](#) and [Banerjee and Gopalan \(2024\)](#) penalize by reward model uncertainty ($r_{\text{eff}} = r - \alpha \sigma_r$, where σ_r is

reward uncertainty and α penalty strength), a pessimistic approach preferring conservative estimates. [Liu et al. \(2025a\)](#) integrate calibration directly into policy gradients via cross-entropy regularization aligning sequence probability with reward.

Granular weighting. At the token level, [Yoon et al.](#) focus optimization on high-uncertainty decision points while [Liu et al. \(2024a\)](#) weight by probability differences between contrastive models. At the trajectory level, [Zhou et al. \(2025a\)](#) rescale advantages by both confidence and difficulty; [Lu et al. \(2025\)](#) construct preference pairs from confidence-ranked outputs. At the preference level, [Pokharel et al. \(2025\)](#) dynamically scale DPO loss ([Rafailov et al., 2023](#)) based on relative reward margins, down-weighting ambiguous comparisons. **Self-feedback.** [Li et al. \(2025d\)](#) use self-assessed confidence as intrinsic reward, generating pseudo-preferences without external labels. [Wu et al. \(2025\)](#) train calibrated abstention via Brier-score rewards ([Gneiting and Raftery, 2007](#)) that incentivize refusal below user-specified risk thresholds.

By transforming confidence from a post-hoc metric into an intrinsic training signal, these methods shift the paradigm from blind data absorption to discerning learning. This integration fosters models that do not merely mimic patterns, but actively navigate the boundaries of their own knowledge.

3 Confidence-Driven Inference

We examine how confidence is used during inference to guide LLM behavior, independent of how confidence is estimated. Existing approaches apply confidence to select among candidate outputs (§3.1), control reasoning and interaction as inference unfolds (§3.2), or guide decoding implicitly through predictive discrepancies (§3.3).

3.1 Confidence-driven output selection

When LLMs are sampled multiple times at inference, they typically produce a set \mathcal{C} including candidate outputs $\{c_1, \dots, c_n\}$ (e.g., responses or reasoning traces) rather than a single deterministic answer. Each candidate c_j is associated with a confidence score $s_j \in \mathbb{R}$, estimated by some mechanism. Inference-time output selection maps \mathcal{C} to a final output \hat{c} . While all methods in this subsection perform output selection, they differ in how confidence is used as an inference operator, either by aggregating multiple candidates jointly or by scoring candidates and selecting one.

In confidence-weighted aggregation, the final output is determined jointly by multiple candidates,

237 $\hat{c} = \arg \max_{c \in \mathcal{C}} \sum_{j=1}^n s_j \mathbb{I}(c_j = c)$, where each
 238 candidate’s contribution is weighted by its con-
 239 fidence score. Unlike classical self-consistency,
 240 which aggregates reasoning traces via uniform ma-
 241 jority voting (Wang et al., 2022), several works
 242 weight candidates according to confidence, allow-
 243 ing infrequent but high-confidence correct answers
 244 to prevail the final decision (Taubenfeld et al., 2025;
 245 Fu et al., 2025). Confidence weighting suppresses
 246 unreliable traces and elevates correct minorities.

247 Unlike aggregation, no information is merged
 248 across candidates in confidence-based scoring,
 249 where inference is governed by selection rather
 250 than combination. More specifically, candidates
 251 are evaluated independently and the final output is
 252 selected as $\hat{c} = \arg \max_{c_j \in \mathcal{C}} s_j$, with confidence
 253 serving as a reliability score rather than a weight.
 254 Methods in Jeong and Choi (2025) and Zhou et al.
 255 (2025b) apply this paradigm by re-scoring multiple
 256 candidates and selecting the most confident one.

257 3.2 Confidence-driven control of decoding

258 While confidence can be applied after inference
 259 to select among completed outputs, it can also act
 260 during inference to control how reasoning and in-
 261 teraction unfold, where confidence functions as a
 262 control signal over the reasoning process rather
 263 than only the final decision. Consider an infer-
 264 ence trajectory represented as a sequence of
 265 states and associated confidence scores $(z_1, s_1) \rightarrow$
 266 $(z_2, s_2) \rightarrow \dots \rightarrow (z_T, s_T)$. Confidence-driven rea-
 267 soning control modifies the transition policy into
 268 $z_{t+1} \sim \pi(z_{t+1} | z_{1:t}, s_{1:t})$, allowing confidence to
 269 govern continuation, refinement, and belief main-
 270 tenance across steps or turns.

271 **Adaptive computation and early stopping.** Some
 272 methods use confidence to allocate inference-time
 273 computation by continuing reasoning only when
 274 confidence remains below a threshold. Huang et al.
 275 (2025a) terminate sampling once confidence is suf-
 276 ficiently high, reducing inference cost without sac-
 277 rificing accuracy, while Fu et al. (2025) allocate
 278 deeper reasoning only when confidence remains
 279 low. Qiao et al. (2025) apply this principle at the
 280 chain-of-thought level, terminating or compress-
 281 ing reasoning once confidence exceeds a threshold
 282 and stabilizes ($|s_t - s_{t-1}| \leq \epsilon$). Across these ap-
 283 proaches, confidence serves as a controller that
 284 prevents unnecessary computation.

285 **Refinement.** Confidence can guide where refine-
 286 ment occurs by identifying unreliable intermediate
 287 states ($s_t \leq \tau$) and selectively refining only these
 288 parts. Wei et al. (2024) selectively regenerate low-

289 confidence substructures in structured generation,
 290 while Fu et al. (2025) prune low-confidence rea-
 291 soning paths and re-explores alternatives. In both
 292 cases, confidence directs refinement efforts.

293 **Multi-turn stability and interaction control.** In
 294 multi-turn interactions, confidence can also regu-
 295 late behavior across turns by discouraging unneces-
 296 sary revision of high-confidence earlier responses.
 297 Li et al. (2025e) incorporate confidence into the dia-
 298 logue state by extracting specific tokens to estimate
 299 response certainty, conditioning future responses
 300 on both prior text and associated confidence to
 301 maintain consistency across interactions.

302 3.3 Confidence-driven contrastive decoding

303 Contrastive decoding methods exploit disagree-
 304 ment between predictive views as an implicit con-
 305 fidence signal during inference, rather than esti-
 306 mating confidence explicitly. Depending on how
 307 predictive views are constructed (across contexts,
 308 layers, or models), different forms of disagreement
 309 arise, but in all cases they reflect the relative sta-
 310 bility of token support. At the decoding step t
 311 (next-token prediction given the current prefix),
 312 two views produce two sets of logits $\text{logit}_t^{(a)}$ and
 313 $\text{logit}_t^{(b)}$, where view a is expected to be more reli-
 314 able by construction. A contrastive logit is defined
 315 as $\tilde{\text{logit}}_t = \text{logit}_t^{(a)} - \lambda \text{logit}_t^{(b)}$, where $\lambda > 0$ con-
 316 trols the strength of penalizing the less reliable
 317 view. Decoding then samples the next token from
 318 the distribution $\pi_t \propto \exp(\tilde{\text{logit}}_t)$, favoring tokens
 319 whose support is stable under the more reliable
 320 view and weak under the less reliable one.

321 To build these predictive views, con-
 322 text-parametric methods first contrast context-
 323 conditioned predictions with the model’s
 324 parametric prior: Huang and Chen (2025) down-
 325 weight tokens that remain likely under random
 326 context masking, while Khandelwal et al. (2025)
 327 adaptively balances parametric and contextual dis-
 328 tributions based on their conflict, treating reliance
 329 on informative context as confidence. Alternatively,
 330 intra-model methods derive contrast from different
 331 transformer layers: Chuang et al. (2023) favors
 332 tokens whose support strengthens from shallow
 333 to deep layers, and Zhang et al. (2025a) learns
 334 when such contrast should be applied, encoding
 335 confidence as consistency across representational
 336 depth. Finally, the inter-model method proposed
 337 by Li et al. (2023) contrasts predictions from an
 338 expert and a weaker model, amplifying tokens
 339 preferred by the expert and using model-level

disagreement for token reliability.

At inference time, confidence can function as a signal that governs output selection, reasoning dynamics, and decoding behavior, enabling more reliable generation without retraining.

4 Confidence-Guided Model Selection

Model selection addresses a practical tension in LLM deployment: stronger models tend to be more accurate but also more costly. Confidence provides a principled control signal for model selection, enabling systems to decide when a cheaper model is sufficient and when to invoke a stronger model. In this section, we review confidence-guided model-selection architectures (§4.1) and the control policies (§4.2) that operationalize confidence for routing, cascading, and their hybrids.

4.1 Confidence-guided selection architectures

Prior work instantiates confidence-guided model selection via (i) sequential cascading, (ii) parallel routing, or (iii) hybrids that combine both.

Cascading. Sequential cascading invoke models in increasing strength and use confidence to decide whether to stop or escalate. This paradigm is analyzed in [Gupta et al. \(2024\)](#), operationalized for cost control in [Chen et al. \(2023b\)](#), accelerated via speculative-decoding-based cascading in [Narasimhan et al. \(2024\)](#), improved by confidence-aware tuning of smaller models in [Rabanser et al. \(2025\)](#), and optimized through probabilistic threshold tuning in [Zellinger and Thomson \(2025\)](#).

Routing. Parallel routing selects one model (or a small subset) without sequential escalation, using confidence to choose among candidates and to manage routing overhead. Examples include routing as query-to-model assignment ([Liu et al., 2024c](#)), confidence-token interfaces for downstream decisions ([Chuang et al., 2024](#)), uncertainty-driven routing policies ([Zhang et al., 2025c](#)), and confidence-aware prediction frameworks for scalable selection ([Barrak et al., 2025](#)).

Hybrids. Hybrid systems combine routing and cascading by using confidence to narrow candidates and then escalating selectively. Representative methods include cascade routing as a unified policy class ([Dekoninck et al., 2024](#)) and select-then-route pipelines that shortlist models before adaptive escalation ([Shah and Shridhar, 2025](#)).

4.2 Confidence-guided control policies

Complementary to architectural instantiations, we summarize policy primitives that map a confidence/quality signal to actions (escalate, select a

model, or allocate additional compute) and can be embedded in any architecture.

Thresholding / abstain–escalate. The canonical policy is confidence-gated deferral: a cheaper model answers when confidence is high and defers otherwise. [Gupta et al. \(2024\)](#) show that naive sequence-level confidence can be misleading and propose token-level aggregation (e.g., quantiles) and learned decision functions for improved cost-quality trade-offs. [Rabanser et al. \(2025\)](#) strengthen deferral by training smaller models to be confident when correct and to defer when uncertain. [Zellinger and Thomson \(2025\)](#) address threshold tuning directly by modeling joint calibrated confidences across cascade stages and optimizing thresholds via differentiable error-cost objectives.

Selection among candidates. A second family uses confidence to select among candidate models. [Chuang et al. \(2024\)](#) introduce confidence tokens to support routing and rejection, [Zhang et al. \(2025c\)](#) design uncertainty-driven routers, and [Barrak et al. \(2025\)](#) scale to large candidate sets by predicting likely winners from model-judged comparisons and using predictor uncertainty to control routing overhead.

Allocate compute adaptively. Finally, some methods allocate compute adaptively during inference: [Chen et al. \(2025b\)](#) switch from a small to a large model when confidence drops during reasoning; [Narasimhan et al. \(2024\)](#) combine cascading with speculative decoding via confidence-driven acceptance or deferral to reduce decoding-time compute.

In summary, confidence-guided routing and cascading implement policies that map queries (and sometimes intermediate generation states) to model-invocation decisions, thereby navigating cost-quality trade-offs ([Chen et al., 2023b](#); [Dekoninck et al., 2024](#)). Recent token- or step-level extensions further suggest that confidence can govern not only which model to use, but also when to hand off computation during generation ([Narasimhan et al., 2024](#); [Chen et al., 2025b](#)).

5 Confidence-Gated RAG Systems

We focus on the use of confidence: given a reliable signal, what mechanisms translate confidence into actionable RAG behavior? We organize existing work into three categories: adaptive retrieval triggering (§5.1), context filtering and selection (§5.2), and knowledge conflict resolution (§5.3).

5.1 Adaptive retrieval triggering

RAG faces an efficiency-accuracy trade-off: unconditional retrieval wastes computation, while rigid refusal limits models to stale parametric knowledge. Confidence resolves this by transforming static retrieval policies into dynamic resource allocation.

Whether to retrieve. Most methods implement confidence-gated retrieval. [Jiang et al. \(2023\)](#) generate a draft sentence and trigger retrieval when it contains low-probability tokens, using those spans as queries. [Wang et al. \(2023a\)](#) elicit self-knowledge to decide retrieval necessity, while [Tao et al. \(2024\)](#) align verbalized confidence with response quality so that confidence reliably gates retrieval. Beyond binary gating, [Jeong et al. \(2024\)](#) and [Zubkova et al. \(2025\)](#) route queries across no retrieval, single-step, or multi-hop retrieval based on predicted complexity or entropy, and [Yao et al. \(2025\)](#) extract self-aware uncertainty from activations to select retrieval strategies.

What to retrieve. Confidence also guides what to retrieve: [Su et al. \(2024b\)](#) detect increasing uncertainty by combining token-level entropy, attention-derived influence, and semantic significance, forming queries from context-wide self-attention rather than only recent tokens. Besides, some methods treat retrieval as expected utility: [Chen et al. \(2025a\)](#) compare the parametric answer to a self-generated contextual probe—convergence bypasses retrieval, divergence triggers it; [Asai et al. \(2024\)](#) learn special tokens whose probabilities encode retrieval necessity and critique judgments; [Lymperopoulos and Sarathy \(2025\)](#) model joint uncertainty over the LLM and external tools to support retrieve-or-abstain decisions.

When to retrieve. Pre-generation methods ([Wang et al., 2023a](#); [Jeong et al., 2024](#); [Tao et al., 2024](#)) assess confidence from the query or a probe before answering. Mid-generation methods ([Jiang et al., 2023](#); [Su et al., 2024b](#)) interleave retrieval when uncertainty spikes during decoding. [Yan et al. \(2024\)](#) operate post-retrieval: a lightweight evaluator scores retrieved evidence confidence and triggers corrective actions—including web search escalation—when retrieval quality is insufficient. [Li et al. \(2025a\)](#) further show that retrieval timing can be determined by exploiting trends (first- and second-order entropy differences) to retrieve preemptively, before errors propagate.

5.2 Confidence-based context filtering

After retrieval is triggered, the challenge becomes context allocation: selecting which evidence to in-

clude within limited context budgets. Confidence serves as a utility signal, estimating whether retrieved content would improve generation rather than merely match the query.

Utility-scored filtering. A common pattern is to compute (or learn) a scalar utility score for each evidence unit and retain only high-utility content. [Wang et al. \(2023b\)](#) construct oracle usefulness signals from lexical overlap and conditional cross-mutual information, the confidence increase when a span is provided, then train a sentence-level filter that removes distracting content at test time. [Yan et al. \(2024\)](#) use a lightweight evaluator that scores retrieved-set quality as CORRECT, INCORRECT, or AMBIGUOUS, triggering corrective actions (keep, refine, or escalate to web search) and decomposing documents to focus on key information.

Information-gain filtering. Absolute relevance can be misleading: an on-topic document may add noise rather than signal. [Wang et al. \(2025\)](#) define Document Information Gain as the change in generation confidence with vs. without a document, training a reranker to prioritize evidence that meaningfully shifts confidence. [Isoda \(2025\)](#) filter at sentence level using a self-knowledge confidence score: the model outputs a scalar “know/unknown” probability for each candidate sentence, and an RL policy learns to keep sentences that increase this self-knowledge (and drop the rest). [Zhu et al. \(2024\)](#) assign each document a relevance confidence from control-token judgments, retain only high-score docs, and decode using only their KV caches to reduce noise and latency. [Li et al. \(2024e\)](#) use SNR-based span uncertainty to estimate chunk similarity in long-context retrieval, improving calibration and robustness to chunking artifacts and distribution shift.

Conformal filtering. Conformal prediction calibrates raw confidence scores—retrieval similarities or generation probabilities—into threshold rules with coverage guarantees. [Li et al. \(2024d\)](#) calibrate retrieval scores to form passage sets and generation confidence to form answer sets, composing both for end-to-end guarantees. [Rouzrokh et al. \(2024\)](#) calibrate similarity scores to select context at a user-specified error rate. [Chakraborty et al. \(2025\)](#) show that calibrating any scoring function yields a principled filter that reduces retained context by 2-3 times while preserving coverage. [Feng et al. \(2025\)](#) apply conformal calibration post-hoc to sub-claim factuality scores.

5.3 Confidence for knowledge conflict

Retrieved context may contradict parametric knowledge; confidence quantifies the reliability of each source to arbitrate the conflict.

Confidence signals for conflict. Recent work turns “conflict” into measurable confidence signals rather than answer disagreement. Sun et al. (2024) detect RAG hallucinations via mechanistic interpretability, computing separate scores for parametric-knowledge utilization and external-context integration; hallucinations correlate with FFN dominance coupled with copying-head failure. Fadeeva et al. (2025) separate faithfulness (entailed by retrieved context) from factuality (true in the world) and route uncertainty quantification accordingly, so “unsupported by retrieval” is not conflated with “false.” Asai et al. (2024) formalize support at inference time: learned reflection-token probabilities (relevance, support, usefulness) yield calibrated, inspectable confidence over whether each continuation is grounded in retrieved evidence.

Resolution policies. Once confidence exposes a conflict, it must drive an action. Sun et al. (2025) use confidence to partition queries by whether they fall within the model’s parametric knowledge boundary and the retrieval boundary; DPO then trains quadrant-specific behaviors, answering when at least one source suffices and abstaining when both fail. Goswami and Kurra (2025) apply calibrated NLI-ensemble scores as a post-hoc arbitration layer, abstaining when entailment confidence falls below a threshold. Soudani et al. (2025) diagnose why standard uncertainty estimators fail in RAG and propose an axiomatic calibration that conditions uncertainty jointly on model and evidence.

Confidence turns RAG into a closed-loop policy: it gates retrieval, allocates context by marginal utility, and arbitrates conflicts by triggering resolution, escalation, or abstention strategies.

6 Confidence-Based Risk Management

While previous sections utilize confidence to enhance performance and efficiency, this section focuses on risk control and reliability. We approach risk management through a confidence-based three-stage framework: first diagnosing anomalies such as hallucinations and adversarial attacks (§6.1), then controlling risk through mechanisms like abstention and conformal prediction (§6.2), and finally ensuring the reliability of the confidence signal itself to reduce overconfidence (§6.3).

6.1 Confidence-based risk diagnosis

We categorize diagnostic mechanisms by the source of risk: detecting internal model hallucinations and flagging external input anomalies.

Hallucination detection. Hallucinations represent a misalignment between the model’s generated output and factual reality (Rawte et al., 2023). Confidence scores can be applied to this task as diagnostic signals. A primary application involves checking output consistency: works like Manakul et al. (2023) and Muhammed et al. (2025) use self-consistency scores to measure disagreement across independently sampled responses as a proxy for hallucination. This strategy is further adapted to filter errors in domain-specific applications like autonomous driving (Dona et al., 2024). To address the ambiguity of phrasing, Farquhar et al. (2024) introduce Semantic Entropy over sampled meaning clusters, while Han et al. (2024) efficiently approximate it via linear probes to avoid sampling cost. Beyond sampling, detection mechanisms exploit internal state anomalies, utilizing the EigenScore (Chen et al., 2024) to measure the semantic consistency in the embedding space, or training classifiers on hidden layer activations (Azaria and Mitchell, 2023) and various confidence-derived features (Quevedo et al., 2025). Other approaches employ active diagnosis of hallucination by measuring confidence fluctuations under perturbed queries (Joo et al., 2025) or intervening on learned verbal uncertainty features (Ji et al., 2025).

Adversarial and OOD detection. Confidence signals can serve as indicators for adversarial and out-of-distribution (OOD) anomalies. For instance, (Obadinma et al., 2024) utilizes confidence-aware training to mitigate calibration attacks, while SafeBehavior (Zhao et al., 2025) computes a “safety confidence score” to actively trigger self-revision processes during jailbreak attempts. In the context of OOD detection, confidence patterns are utilized to trace data provenance; for example (Li et al., 2025b) aggregates confidence scores to detect membership inference risks by recognizing previously seen training data.

6.2 Confidence-based risk control

Moving from detection to intervention, we review how confidence thresholds are employed to gate unreliable generations and bound output error rates. **Abstention.** Abstention refers to the refusal of an LLM to provide an answer (Wen et al., 2025). A common confidence-driven approach uses a gate-keeping mechanism to abstain from answering

645 when confidence is below a threshold τ . In post-
646 hoc thresholding, confidence scores are applied
647 directly to frozen models to judge response qual-
648 ity. A good confidence score for abstention should
649 focus on tracking correctness rather than fluency
650 or typicality. Following this principle, [Cole et al.](#)
651 [\(2023\)](#) applies self-consistency for abstention and
652 and shows it yields better performance than log-
653 likelihood, while [Ren et al. \(2023\)](#) utilizes LLM
654 self-evaluation scores to better discriminate low-
655 quality generations. Alternatively, [Huang et al.](#)
656 [\(2025c\)](#) probes internal activations to extract a
657 scalar uncertainty feature specifically optimized
658 for this thresholding decision.

659 Moving from inference-time filtering to model
660 adaption, recent works focus on teaching the model
661 when to abstain, applying confidence scores as su-
662 pervisory signals. [Chen et al. \(2023a\)](#) fine-tunes
663 LLMs with labeled answers to obtain a confidence
664 score that combines self-evaluation and normalized
665 likelihood for selective prediction. To further re-
666 duce label reliance, [Tjandra et al. \(2024\)](#) employs
667 semantic entropy as a pseudo-label to enable label-
668 free abstention fine-tuning, while [An and Xu \(2025\)](#)
669 uses semantic confidence as a reward signal via RL
670 to optimize the model’s boundary recognition.

671 **Conformal prediction.** Conformal prediction (CP)
672 ([Shafer and Vovk, 2008](#)) addresses the limitation
673 of heuristic thresholds by transforming confidence
674 scores into statistically valid prediction sets, en-
675 suring that on expectation a correct² answer is in-
676 cluded with a user-specified probability, or more
677 broadly, that a specific risk measure is controlled
678 ([Angelopoulos et al., 2024](#)). The application of CP
679 in LLMs can be distinguished by the specific con-
680 fidence signals used to construct these sets. Early
681 approaches relied on logit-based scores: [Kumar](#)
682 [et al. \(2023\)](#) applies CP to multiple choice ques-
683 tions using logit scores of the choices, while [Quach](#)
684 [et al. \(2024\)](#) utilizes normalized likelihoods to tune
685 sampling parameters for open-ended generation.

686 However, recognizing that raw logits often fails
687 to capture semantic correctness, recent methods
688 have shifted toward self-consistency and semantic
689 scores. [Mohri and Hashimoto \(2024\)](#) argues that
690 sets of answers are less useful than factual sub-
691 claims, employing self-consistency and LLM-as-a-
692 judge scores to guarantee all selected sub-claims
693 are factual with user-specified probability. This

²In LLMs, correctness is task-defined: matching ground-
truth labels for classification, or using proxies for open-ended
outputs (e.g. semantic match or verifier pass).

694 approach is further refined by [Cherian et al. \(2024\)](#),
695 which exploits the heterogeneity in score quality
696 across prompts to adaptively adjust thresholds. To
697 ensure coherence in reasoning, [Rubin-Toles et al.](#)
698 [\(2025\)](#) extends these scores to model logical de-
699 pendencies via deductibility graphs. For black-box
700 settings, [Su et al. \(2024a\)](#) demonstrates that API-
701 accessible scores like semantic similarity can serve
702 as effective non-conformity scores. Finally, to im-
703 prove the efficiency of these sets, [Xi et al. \(2025\)](#)
704 demonstrates that optimizing confidence distribu-
705 tion via temperature scaling yields tighter coverage.

696 6.3 Ensuring the reliability of confidence 706

707 The efficacy of risk frameworks depends on the re-
708 liability of confidence signals, yet LLMs frequently
709 exhibit significant overconfidence and miscalibra-
710 tion ([Groot and Valdenegro Toro, 2024](#); [Tian et al.,](#)
711 [2025](#)). To resolve this, elicitation strategies are ap-
712 plied to extract more calibrated scores: [Kadavath](#)
713 [et al. \(2022\)](#) established that LLMs can accurately
714 estimate their own correctness via the logit proba-
715 bility $P(\text{True})$ when specifically prompted, while
716 [Tian et al. \(2023\)](#) further showed that soliciting
717 such explicit probabilities or top-k options yields
718 significantly better calibration than open-ended ver-
719 balization. For existing scores, post-hoc calibration
720 utilizes specialized confidence metrics like the TH-
721 score ([Tian et al., 2025](#)) and Net Calibration Error
722 ([Groot and Valdenegro Toro, 2024](#)) to explicitly
723 penalize the gap between confidence and accuracy.

724 In summary, confidence-based risk management
725 transforms the signal from a passive metric into a
726 robust guardrail that detects errors, bounds risks,
727 and ensures system reliability.

728 7 Conclusion 728

729 This survey presents a unified view of confidence as
730 a control signal that governs LLM behavior across
731 the full system lifecycle. During training, it guides
732 data selection, loss weighting, and preference op-
733 timization; during inference, it supports output se-
734 lection, adaptive computation, and contrastive de-
735 coding; and in deployment, it enables cost-aware
736 routing, confidence-gated RAG, and risk control
737 via abstention and conformal prediction. Across
738 these settings, confidence consistently turns dis-
739 crete heuristics—what to learn, when to retrieve,
740 whether to answer—into more principled allocation
741 of computation and evidence. Reliable confidence
742 is therefore not an end, but a prerequisite for sys-
743 tems that act on uncertainty by deferring, retrieving,
744 or abstaining when warranted.

745 **Limitations**

746 This survey focuses on **confidence utilization** as a
747 control signal rather than confidence estimation or
748 calibration. Readers seeking comprehensive cover-
749 age of uncertainty quantification methods should
750 consult complementary surveys.

751 We primarily survey English-language litera-
752 ture and methods evaluated on English bench-
753 marks; confidence utilization in multilingual and
754 low-resource settings remains comparatively under-
755 explored.

756 Control policies inherit failure modes of the un-
757 derlying confidence signal (miscalibration, bias,
758 brittleness under distribution shift or adversarial
759 prompting), which can lead to unsafe acceptance,
760 premature deferral, or unnecessary escalation.

761 Many methods add latency or compute (sam-
762 pling, verifiers, multi-stage retrieval), and formal
763 guarantees (e.g., conformal prediction) depend on
764 assumptions and proxies for correctness that may
765 not hold under deployment drift.

766 **Implication.** Taken together, these limitations sug-
767 gest that confidence-as-control should be viewed
768 as a systems design paradigm rather than a drop-
769 in solution: robust deployment requires validating
770 confidence reliability under realistic shift, aligning
771 proxy correctness with application risk, and ac-
772 counting for compute and human factors alongside
773 model accuracy.

774 **Acknowledgment on AI Assistance.**

775 GPT 5.2, Claude Opus 4.5, and Gemini 3 Pro
776 were used solely to polish the language of this
777 manuscript (e.g., improving grammar, clarity, and
778 style). These tools did not generate or alter the
779 scientific content, interpretations, results, or con-
780 clusions. All content was reviewed, verified, and
781 approved by the authors, who take full responsibil-
782 ity for the accuracy and integrity of the work.

783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836

References

Hao An and Yang Xu. 2025. Teaching llms to abstain via fine-grained semantic confidence reward. *arXiv preprint arXiv:2510.24020*.

Anastasios Nikolas Angelopoulos, Stephen Bates, Adam Fisch, Lihua Lei, and Tal Schuster. 2024. Conformal risk control. In *The Twelfth International Conference on Learning Representations*.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-rag: Learning to retrieve, generate, and critique through self-reflection.

Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when it’s lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976.

Debangshu Banerjee and Aditya Gopalan. 2024. Towards reliable alignment: Uncertainty-aware rlhf. *arXiv preprint arXiv:2410.23726*.

Amine Barrak, Yosr Fourati, Michael Olchawa, Emna Ksontini, and Khalil Zoghalmi. 2025. Cargo: A framework for confidence-aware routing of large language models. *arXiv preprint arXiv:2509.14899*.

Debashish Chakraborty, Eugene Yang, Daniel Khashabi, Dawn Lawrie, and Kevin Duh. 2025. Principled context engineering for rag: Statistical guarantees via conformal prediction. *arXiv preprint arXiv:2511.17908*.

Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024. INSIDE: LLMs’ internal states retain the power of hallucination detection. In *The Twelfth International Conference on Learning Representations*.

Jiefeng Chen, Jinsung Yoon, Sayna Ebrahimi, Sercan Arik, Tomas Pfister, and Somesh Jha. 2023a. Adaptation with self-evaluation to improve selective prediction in llms. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5190–5213.

Jiuhai Chen and Jonas Mueller. 2024. Automated data curation for robust language model fine-tuning. *arXiv preprint arXiv:2403.12776*.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. Frugalgpt: How to use large language models while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176*.

Wang Chen, Guanqiang Qi, Weikang Li, Yang Li, Deguo Xia, and Jizhou Huang. 2025a. Pairs: Parametric-verified adaptive information retrieval and selection for efficient rag. *arXiv preprint arXiv:2508.04057*.

Zhuokun Chen, Zeren Chen, Jiahao He, Lu Sheng, Mingkui Tan, Jianfei Cai, and Bohan Zhuang. 2025b. R-stitch: Dynamic trajectory stitching for efficient reasoning. *arXiv preprint arXiv:2507.17307*.

John Cherian, Isaac Gibbs, and Emmanuel Candes. 2024. Large language model validity via enhanced conformal prediction methods. *Advances in Neural Information Processing Systems*, 37:114812–114842.

Yu-Neng Chuang, Prathusha Kameswara Sarma, Parikshit Gopalan, John Boccio, Sara Bolouki, Xia Hu, and Helen Zhou. 2024. Learning to route llms with confidence tokens. *arXiv preprint arXiv:2410.13284*.

Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. *arXiv preprint arXiv:2309.03883*.

Jeremy Cole, Michael Zhang, Dan Gillick, Julian Eisen-schlos, Bhuwan Dhingra, and Jacob Eisenstein. 2023. Selectively answering ambiguous questions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 530–543.

Jasper Dekoninck, Maximilian Baader, and Martin Vechev. 2024. A unified approach to routing and cascading for llms. *arXiv preprint arXiv:2410.10347*.

Malsha Ashani Mahawatta Dona, Beatriz Cabrero-Daniel, Yinan Yu, and Christian Berger. 2024. Llm can check their own results to mitigate hallucinations in traffic understanding tasks. In *IFIP International Conference on Testing Software and Systems*, pages 114–130. Springer.

Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. 2023. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*.

He Du, Bowen Li, Chengxing Xie, Chang Gao, Kai Chen, and Dacheng Tao. 2025. Confidence as a reward: Transforming llms into reward models. *arXiv preprint arXiv:2510.13501*.

Ekaterina Fadeeva, Aleksandr Rubashevskii, Roman Vashurin, Shehzaad Dhuliawala, Artem Shelmanov, Timothy Baldwin, Preslav Nakov, Mrinmaya Sachan, and Maxim Panov. 2025. Faithfulness-aware uncertainty quantification for fact-checking the output of retrieval augmented generation. *arXiv preprint arXiv:2505.21072*.

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630.

Naihe Feng, Yi Sui, Shiyi Hou, Jesse C Cresswell, and Ga Wu. 2025. Response quality assessment for retrieval-augmented generation via conditional conformal factuality. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2832–2836.

891	Yichao Fu, Xuwei Wang, Yuandong Tian, and Jiawei Zhao. 2025. Deep think with confidence. <i>arXiv preprint arXiv:2508.15260</i> .	Zhiqi Huang, Vivek Datla, Chenyang Zhu, Alfy Samuel, Daben Liu, Anoop Kumar, and Ritesh Soni. 2025c. Confidence-based response abstinence: Improving llm trustworthiness via activation-based uncertainty estimation. In <i>Proceedings of the 2nd Workshop on Uncertainty-Aware NLP (UncertainNLP 2025)</i> , pages 184–193.	944
892			945
893			946
894	Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. 2024. A survey of confidence estimation and calibration in large language models. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6577–6595.		947
895			948
896			949
897			950
898			951
899			952
900			953
901			954
902	Tilmann Gneiting and Adrian E Raftery. 2007. Strictly proper scoring rules, prediction, and estimation. <i>Journal of the American statistical Association</i> , 102(477):359–378.		955
903			956
904			957
905			958
906	Saumya Goswami and Siddharth Kurra. 2025. Halt-rag: A task-adaptable framework for hallucination detection with calibrated nli ensembles and abstention. <i>arXiv preprint arXiv:2509.07475</i> .		959
907			960
908			961
909			962
910	Tobias Groot and Matias Valdenegro Toro. 2024. Overconfidence is key: Verbalized uncertainty evaluation in large language and vision-language models. In <i>Proceedings of the 4th Workshop on Trustworthy Natural Language Processing (TrustNLP 2024)</i> , pages 145–171, Mexico City, Mexico. Association for Computational Linguistics.		963
911			964
912			965
913			966
914			967
915			968
916			969
917	Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkritum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. 2024. Language model cascades: Token-level uncertainty and beyond. <i>arXiv preprint arXiv:2404.10136</i> .		970
918			971
919			972
920			973
921			974
922	Jiatong Han, Jannik Kossen, Muhammed Razzak, Lisa Schut, Shreshth A Malik, and Yarin Gal. 2024. Semantic entropy probes: Robust and cheap hallucination detection in LLMs. In <i>ICML 2024 Workshop on Foundation Models in the Wild</i> .		975
923			976
924			977
925			978
926			979
927	Jindong Han, Hao Liu, Jun Fang, Naiqiang Tan, and Hui Xiong. 2025. Automatic instruction data selection for large language models via uncertainty-aware influence maximization. In <i>Proceedings of the ACM on Web Conference 2025</i> , pages 4969–4979.		980
928			981
929			982
930			983
931			984
932	Cheng Peng Huang and Hao-Yuan Chen. 2025. Delta-contrastive decoding mitigates text hallucinations in large language models. <i>arXiv preprint arXiv:2502.05825</i> .		985
933			986
934			987
935			988
936	Chengsong Huang, Langlin Huang, Jixuan Leng, Jiacheng Liu, and Jiaxin Huang. 2025a. Efficient test-time scaling via self-calibration. <i>arXiv preprint arXiv:2503.00031</i> .		989
937			990
938			991
939			992
940	Haiduo Huang, Jiangcheng Song, Yadong Zhang, and Pengju Ren. 2025b. Selectkd: Selective token-weighted knowledge distillation for llms. <i>arXiv preprint arXiv:2510.24021</i> .		993
941			994
942			995
943			996
			997
			998
			999
			1000

998	Bhawesh Kumar, Charlie Lu, Gauri Gupta, Anil Palepu, David Bellamy, Ramesh Raskar, and Andrew Beam. 2023. Conformal prediction with large language models for multi-choice question answering. <i>arXiv preprint arXiv:2305.18404</i> .	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.	1054 1055 1056 1057
1003	Jixuan Leng, Chengsong Huang, Banghua Zhu, and Jiaxin Huang. 2024. Taming overconfidence in llms: Reward calibration in rlhf. <i>arXiv preprint arXiv:2410.09724</i> .	Yubo Li, Yidi Miao, Xueying Ding, Ramayya Krishnan, and Rema Padman. 2025e. Firm or fickle? evaluating large language models consistency in sequential interactions. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 6679–6700, Vienna, Austria. Association for Computational Linguistics.	1058 1059 1060 1061 1062 1063 1064
1007	Bo Li, Tian Tian, Zhenghua Xu, Hao Cheng, Shikun Zhang, and Wei Ye. 2025a. Modeling uncertainty trends for timely retrieval in dynamic rag. <i>arXiv preprint arXiv:2511.09980</i> .	Zixuan Li, Jing Xiong, Fanghua Ye, Chuanyang Zheng, Xun Wu, Jianqiao Lu, Zhongwei Wan, Xiaodan Liang, Chengming Li, Zhenan Sun, and 1 others. 2024e. Uncertaintyrag: Span-level uncertainty enhanced long-context modeling for retrieval-augmented generation. <i>arXiv preprint arXiv:2410.02719</i> .	1065 1066 1067 1068 1069 1070 1071
1011	Haodong Li, Jingqi Zhang, Xiao Cheng, Peihua Mai, Haoyu Wang, and Yang Pan. 2025b. As if we've met before: Llms exhibit certainty in recognizing seen files. <i>arXiv preprint arXiv:2511.15192</i> .	Aiwei Liu, Haoping Bai, Zhiyun Lu, Yanchao Sun, Xiang Kong, Simon Wang, Jiulong Shan, Albin Madappally Jose, Xiaojiang Liu, Lijie Wen, and 1 others. 2024a. Tis-dpo: Token-level importance sampling for direct preference optimization with estimated weights. <i>arXiv preprint arXiv:2410.04350</i> .	1072 1073 1074 1075 1076 1077
1015	Jiaqi Li, Yixuan Tang, and Yi Yang. 2025c. Know the unknown: An uncertainty-sensitive method for llm instruction tuning. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 2972–2989.	Haotian Liu, Shuo Wang, and Hongteng Xu. 2025a. c^2 gspg: Confidence-calibrated group sequence policy gradient towards self-aware reasoning. <i>arXiv preprint arXiv:2509.23129</i> .	1078 1079 1080 1081
1020	Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jixiang Gu, and Tianyi Zhou. 2024a. Selective reflection-tuning: Student-selected data recycling for llm instruction-tuning. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 16189–16211.	Liangxin Liu, Xuebo Liu, Derek F Wong, Dongfang Li, Ziyi Wang, Baotian Hu, and Min Zhang. 2024b. Selectit: Selective instruction tuning for large language models via uncertainty-aware self-reflection. <i>arXiv preprint arXiv:2402.16705</i> .	1082 1083 1084 1085 1086
1026	Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024b. Superfiltering: Weak-to-strong data filtering for fast instruction-tuning. <i>arXiv preprint arXiv:2402.00530</i> .	Xiaoou Liu, Tiejin Chen, Longchao Da, Chacha Chen, Zhen Lin, and Hua Wei. 2025b. Uncertainty quantification and confidence calibration in large language models: A survey. <i>arXiv preprint arXiv:2503.15850</i> .	1087 1088 1089 1090
1031	Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024c. From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7602–7635.	Yueyue Liu, Hongyu Zhang, Yuantian Miao, Van-Hoang Le, and Zhiqiang Li. 2024c. Optllm: Optimal assignment of queries to large language models. In <i>2024 IEEE International Conference on Web Services (ICWS)</i> , pages 788–798. IEEE.	1091 1092 1093 1094 1095
1040	Pengyi Li, Matvey Skripkin, Alexander Zubrey, Andrey Kuznetsov, and Ivan Oseledets. 2025d. Confidence is all you need: Few-shot rl fine-tuning of language models. <i>arXiv preprint arXiv:2506.06395</i> .	Junjie Lu, Yuliang Liu, Chaofeng Qu, Wei Shen, Zhouhan Lin, and Min Xu. 2025. Enhancing llm reasoning via non-human-like reasoning path preference optimization. <i>arXiv preprint arXiv:2510.11104</i> .	1096 1097 1098 1099
1044	Shuo Li, Sangdon Park, Insup Lee, and Osbert Bastani. 2024d. Traq: Trustworthy retrieval augmented question answering via conformal prediction. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3799–3821.	Junyu Luo, Xiao Luo, Xiushi Chen, Zhiping Xiao, Wei Ju, and Ming Zhang. 2024. Semievol: semi-supervised fine-tuning for llm adaptation. <i>arXiv e-prints</i> , pages arXiv–2410.	1100 1101 1102 1103
1051	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori B Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. Contrastive decoding:	Panagiotis Lymeropoulos and Vasanth Sarathy. 2025. Tools in the loop: Quantifying uncertainty of llm question answering systems that use tools. <i>arXiv preprint arXiv:2505.16113</i> .	1104 1105 1106 1107

1108	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023.	Victor Quach, Adam Fisch, Tal Schuster, Adam Yala,	1163
1109	Selfcheckgpt: Zero-resource black-box hallucination	Jae Ho Sohn, Tommi S. Jaakkola, and Regina Barzi-	1164
1110	detection for generative large language models. In	lay. 2024. Conformal language modeling. In <i>The</i>	1165
1111	<i>Proceedings of the 2023 conference on empirical</i>	<i>Twelfth International Conference on Learning Repre-</i>	1166
1112	<i>methods in natural language processing</i> , pages 9004–	<i>sentations</i> .	1167
1113	9017.		
1114	Christopher Mohri and Tatsunori Hashimoto. 2024.	Ernesto Quevedo, Jorge Yero Salazar, Rachel Koerner,	1168
1115	Language models with conformal factuality guaran-	Pablo Rivas, and Tomas Cerny. 2025. Detecting hal-	1169
1116	tees. In <i>International Conference on Machine Learn-</i>	lucinations in large language model generation: A	1170
1117	<i>ing</i> , pages 36029–36047. PMLR.	token probability approach. In <i>Artificial Intelligence</i>	1171
1118	Mohammad Mahdi Moradi, Hossam Amer, Sudhir	<i>and Applications</i> , pages 154–173, Cham. Springer	1172
1119	Mudur, Weiwei Zhang, Yang Liu, and Walid	Nature Switzerland.	1173
1120	Ahmed. 2025. Continuous self-improvement of	Stephan Rabanser, Nathalie Rauschmayr, Achin Kul-	1174
1121	large language models by test-time training with	shrestha, Petra Poklukar, Wittawat Jitkrittum, Sean	1175
1122	verifier-driven sample selection. <i>arXiv preprint</i>	Augenstein, Congchao Wang, and Federico Tombari.	1176
1123	<i>arXiv:2505.19475</i> .	2025. Gatekeeper: Improving model cascades	1177
1124	Diyana Muhammed, Gollam Rabby, and Sören Auer.	through confidence tuning. In <i>The Thirty-ninth An-</i>	1178
1125	2025. Selfcheckagent: Zero-resource hallucination	<i>nual Conference on Neural Information Processing</i>	1179
1126	detection in generative large language models. <i>arXiv</i>	<i>Systems</i> .	1180
1127	<i>preprint arXiv:2502.01812</i> .	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	1181
1128	William Muldrew, Peter Hayes, Mingtian Zhang, and	pher D Manning, Stefano Ermon, and Chelsea Finn.	1182
1129	David Barber. 2024. Active preference learn-	2023. Direct preference optimization: Your language	1183
1130	ing for large language models. <i>arXiv preprint</i>	model is secretly a reward model. <i>Advances in neural</i>	1184
1131	<i>arXiv:2402.08114</i> .	<i>information processing systems</i> , 36:53728–53741.	1185
1132	Harikrishna Narasimhan, Wittawat Jitkrittum,	Amir Hossein Rahmati, Sanket Jantre, Weifeng Zhang,	1186
1133	Ankit Singh Rawat, Seungyeon Kim, Neha Gupta,	Yucheng Wang, Byung-Jun Yoon, Nathan M Urban,	1187
1134	Aditya Krishna Menon, and Sanjiv Kumar. 2024.	and Xiaoning Qian. 2025. C-lora: Contextual low-	1188
1135	Faster cascades via speculative decoding. <i>arXiv</i>	rank adaptation for uncertainty estimation in large	1189
1136	<i>preprint arXiv:2405.19261</i> .	language models. <i>arXiv preprint arXiv:2505.17773</i> .	1190
1137	Stephen Obadinma, Xiaodan Zhu, and Hongyu Guo.	Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A	1191
1138	2024. Calibration attacks: A comprehensive study	survey of hallucination in large foundation models.	1192
1139	of adversarial attacks on model confidence. <i>arXiv</i>	<i>arXiv preprint arXiv:2309.05922</i> .	1193
1140	<i>preprint arXiv:2401.02718</i> .	Jie Ren, Yao Zhao, Tu Vu, Peter J. Liu, and Balaji	1194
1141	Yijun Pan, Taiwei Shi, Jieyu Zhao, and Jiaqi W Ma.	Lakshminarayanan. 2023. Self-evaluation improves	1195
1142	2025. Detecting and filtering unsafe training data via	selective generation in large language models. In <i>Pro-</i>	1196
1143	data attribution. <i>arXiv preprint arXiv:2502.11411</i> .	<i>ceedings on "I Can't Believe It's Not Better: Failure</i>	1197
1144	Jinlong Pang, Na Di, Zhaowei Zhu, Jiaheng Wei, Hao	<i>Modes in the Age of Foundation Models" at NeurIPS</i>	1198
1145	Cheng, Chen Qian, and Yang Liu. 2025. Token clean-	<i>2023 Workshops</i> , volume 239 of <i>Proceedings of Ma-</i>	1199
1146	ing: Fine-grained data selection for llm supervised	<i>chine Learning Research</i> , pages 49–64. PMLR.	1200
1147	fine-tuning. <i>arXiv preprint arXiv:2502.01968</i> .	Pouria Rouzrokh, Shahriar Faghani, Cooper U Gam-	1201
1148	Rhitabrat Pokharel, Yufei Tao, and Ameeta Agrawal.	ble, Moein Shariatnia, and Bradley J Erickson. 2024.	1202
1149	2025. Capo: Confidence aware preference optimiza-	Conflare: conformal large language model retrieval.	1203
1150	tion learning for multilingual preferences. <i>arXiv</i>	<i>arXiv preprint arXiv:2404.04287</i> .	1204
1151	<i>preprint arXiv:2511.07691</i> .	Maxon Rubin-Toles, Maya Gambhir, Keshav Ramji,	1205
1152	Mihir Prabhudesai, Lili Chen, Alex Ippoliti, Katerina	Aaron Roth, and Surbhi Goel. 2025. Conformal lan-	1206
1153	Fragkiadaki, Hao Liu, and Deepak Pathak. 2025.	guage model reasoning with coherent factuality. In	1207
1154	Maximizing confidence alone improves reasoning.	<i>The Thirteenth International Conference on Learning</i>	1208
1155	<i>arXiv preprint arXiv:2505.22660</i> .	<i>Representations</i> .	1209
1156	Ziqing Qiao, Yongheng Deng, Jiali Zeng, Dong	Noveen Sachdeva, Benjamin Coleman, Wang-Cheng	1210
1157	Wang, Lai Wei, Guanbo Wang, Fandong Meng, Jie	Kang, Jianmo Ni, Lichan Hong, Ed H Chi, James	1211
1158	Zhou, Ju Ren, and Yaoxue Zhang. 2025. Concise:	Caverlee, Julian McAuley, and Derek Zhiyuan Cheng.	1212
1159	Confidence-guided compression in step-by-step ef-	2024. How to train data-efficient llms. <i>arXiv preprint</i>	1213
1160	ficient reasoning. In <i>Proceedings of the 2025 Con-</i>	<i>arXiv:2402.09668</i> .	1214
1161	<i>ference on Empirical Methods in Natural Language</i>	Glenn Shafer and Vladimir Vovk. 2008. A tutorial on	1215
1162	<i>Processing</i> , pages 8021–8040.	conformal prediction. <i>Journal of Machine Learning</i>	1216
		<i>Research</i> , 9(3).	1217

1218	Soham Shah and Kumar Shridhar. 2025. Select-then-route : Taxonomy guided routing for LLMs. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 425–441, Suzhou (China). Association for Computational Linguistics.	1274
1219		1275
1220		1276
1221		
1222		
1223		
1224	Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z Ren, and Anirudha Majumdar. 2025. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions. <i>ACM Computing Surveys</i> .	
1225		
1226		
1227		
1228		
1229	Heydar Soudani, Evangelos Kanoulas, and Faegheh Hasebi. 2025. Why uncertainty estimation methods fall short in rag: An axiomatic analysis. <i>arXiv preprint arXiv:2505.07459</i> .	
1230		
1231		
1232		
1233	Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng. 2024a. Api is enough: Conformal prediction for large language models without logit-access. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 979–995.	
1234		
1235		
1236		
1237		
1238	Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024b. Dragin: dynamic retrieval augmented generation based on the information needs of large language models. <i>arXiv preprint arXiv:2403.10081</i> .	
1239		
1240		
1241		
1242		
1243	Xin Sun, Jianan Xie, Zhongqi Chen, Qiang Liu, Shu Wu, Yuehe Chen, Bowen Song, Zilei Wang, Weiqiang Wang, and Liang Wang. 2025. Divide-then-align: Honest alignment based on the knowledge boundary of rag. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11461–11480.	
1244		
1245		
1246		
1247		
1248		
1249		
1250	Zhongxiang Sun, Xiaoxue Zang, Kai Zheng, Yang Song, Jun Xu, Xiao Zhang, Weijie Yu, and Han Li. 2024. Redeeep: Detecting hallucination in retrieval-augmented generation via mechanistic interpretability. <i>arXiv preprint arXiv:2410.11414</i> .	
1251		
1252		
1253		
1254		
1255	Shuchang Tao, Liuyi Yao, Hanxing Ding, Yuexiang Xie, Qi Cao, Fei Sun, Jinyang Gao, Huawei Shen, and Bolin Ding. 2024. When to trust llms: Aligning confidence with response quality. <i>arXiv preprint arXiv:2404.17287</i> .	
1256		
1257		
1258		
1259		
1260	Amir Taubenfeld, Tom Sheffer, Eran Ofek, Amir Feder, Ariel Goldstein, Zorik Gekhman, and Gal Yona. 2025. Confidence improves self-consistency in llms. <i>arXiv preprint arXiv:2502.06233</i> .	
1261		
1262		
1263		
1264	Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5433–5442.	
1265		
1266		
1267		
1268		
1269		
1270		
1271		
1272	Zailong Tian, Zhuoheng Han, Yanzhe Chen, Haozhe Xu, Xi Yang, Richeng Xuan, Houfeng Wang, and Lizi	
1273		
	Liao. 2025. Overconfidence in llm-as-a-judge: Diagnosis and confidence-driven solution. <i>arXiv preprint arXiv:2508.06225</i> .	1274
		1275
		1276
	Benedict Aaron Tjandra, Muhammed Razzak, Jannik Kossen, Kunal Handa, and Yarin Gal. 2024. Fine-tuning large language models to appropriately abstain with semantic entropy. In <i>Neurips Safe Generative AI Workshop 2024</i> .	1277
		1278
		1279
		1280
		1281
	Sahil Tripathi, Md Tabrez Nafis, Imran Hussain, and Jiechao Gao. 2025. The confidence paradox: Can llm know when it’s wrong. <i>arXiv preprint arXiv:2506.23464</i> .	1282
		1283
		1284
		1285
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> .	1286
		1287
		1288
		1289
		1290
	Yile Wang, Peng Li, Maosong Sun, and Yang Liu. 2023a. Self-knowledge guided retrieval augmentation for large language models. <i>arXiv preprint arXiv:2310.05002</i> .	1291
		1292
		1293
		1294
	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023b. Learning to filter context for retrieval-augmented generation. <i>arXiv preprint arXiv:2311.08377</i> .	1295
		1296
		1297
		1298
	Zihan Wang, Zihan Liang, Zhou Shao, Yufei Ma, Huangyu Dai, Ben Chen, Lingtao Mao, Chenyi Lei, Yuqing Ding, and Han Li. 2025. Infogain-rag: Boosting retrieval-augmented generation through document information gain-based reranking and filtering. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 7201–7215.	1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
	Chengwei Wei, Kee Kiat Koo, Amir Tavanaei, and Karim Bouyarmane. 2024. Confidence-aware substructure beam search (cabs): Mitigating hallucination in structured data generation with large language models. <i>arXiv preprint arXiv:2406.00069</i> .	1307
		1308
		1309
		1310
		1311
	Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun Xu, Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang. 2025. Know your limits: A survey of abstention in large language models. <i>Transactions of the Association for Computational Linguistics</i> , 13:529–556.	1312
		1313
		1314
		1315
		1316
	Jiayun Wu, Jiashuo Liu, Zhiyuan Zeng, Tianyang Zhan, and Wenhao Huang. 2025. Mitigating llm hallucination via behaviorally calibrated reinforcement learning. <i>arXiv preprint arXiv:2512.19920</i> .	1317
		1318
		1319
		1320
	HuaJun Xi, Jianguo Huang, Kangdao Liu, Lei Feng, and Hongxin Wei. 2025. Does confidence calibration improve conformal prediction? <i>Transactions on Machine Learning Research</i> .	1321
		1322
		1323
		1324
	Liangru Xie, Hui Liu, Jingying Zeng, Xianfeng Tang, Yan Han, Chen Luo, Jing Huang, Zhen Li, Suhang Wang, and Qi He. 2024. A survey of calibration process for black-box llms. <i>arXiv preprint arXiv:2412.12767</i> .	1325
		1326
		1327
		1328
		1329

1330	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. <i>arXiv preprint arXiv:2306.13063</i> .		
1331		Tuo Zhang, Asal Mehradfar, Dimitrios Dimitriadis, and Salman Avestimehr. 2025c. Leveraging uncertainty estimation for efficient llm routing. <i>arXiv preprint arXiv:2502.11021</i> .	1383
1332			1384
1333			1385
1334			1386
1335	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation.		
1336		Qinjian Zhao, Jiaqi Wang, Zhiqiang Gao, Zhihao Dou, Belal Abuhajja, and Kaizhu Huang. 2025. Safebehavior: Simulating human-like multistage reasoning to mitigate jailbreak attacks in large language models. <i>arXiv preprint arXiv:2509.26345</i> .	1387
1337			1388
1338	Binwei Yao, Zefan Cai, Yun-Shiuan Chuang, Shanglin Yang, Ming Jiang, Diyi Yang, and Junjie Hu. 2024. No preference left behind: Group distributional preference optimization. <i>arXiv preprint arXiv:2412.20299</i> .		1389
1339			1390
1340		Qihuang Zhong, Liang Ding, Li Shen, Juhua Liu, Bo Du, and Dacheng Tao. 2024. Revisiting knowledge distillation for autoregressive language models. <i>arXiv preprint arXiv:2402.11890</i> .	1391
1341			1392
1342	Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao, Linmei Hu, Liu Weichuan, Lei Hou, and Juanzi Li. 2025. Seakr: Self-aware knowledge retrieval for adaptive retrieval augmented generation. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 27022–27043.		1393
1343			1394
1344		Zhanke Zhou, Xiangyu Lu, Chentao Cao, Brando Miranda, Tongliang Liu, Bo Han, and Sanmi Koyejo. 2025a. Codapo: Confidence and difficulty-adaptive policy optimization for post-training language models. In <i>2nd AI for Math Workshop@ ICML 2025</i> .	1395
1345			1396
1346			1397
1347			1398
1348			1399
1349	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuan-Jing Huang. 2023. Do large language models know what they don't know? In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8653–8665.	Ziang Zhou, Tianyuan Jin, Jieming Shi, and Qing Li. 2025b. Steerconf: Steering llms for confidence elicitation. <i>arXiv preprint arXiv:2503.02863</i> .	1401
1350			1402
1351			1403
1352		Yun Zhu, Jia-Chen Gu, Caitlin Sikora, Ho Ko, Yinxiao Liu, Chu-Cheng Lin, Lei Shu, Liangchen Luo, Lei Meng, Bang Liu, and 1 others. 2024. Accelerating inference of retrieval-augmented generation via sparse context selection. <i>arXiv preprint arXiv:2405.16178</i> .	1404
1353			1405
1354	Hee Suk Yoon, Eunseop Yoon, Mark A Hasegawa-Johnson, Sungwoong Kim, and Chang D Yoo. Confpo: Exploiting policy model confidence for critical token selection in preference optimization. In <i>Forty-second International Conference on Machine Learning</i> .		1406
1355			1407
1356			1408
1357			1409
1358		Hanna Zubkova, Ji-Hoon Park, and Seong-Whan Lee. 2025. Sugar: Leveraging contextual confidence for smarter retrieval. In <i>ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	1410
1359			1411
1360	Michael J. Zellinger and Matt Thomson. 2025. Rational tuning of LLM cascades via probabilistic modeling. <i>Transactions on Machine Learning Research</i> .		1412
1361			1413
1362			
1363	Yuanzhao Zhai, Han Zhang, Yu Lei, Yue Yu, Kele Xu, Dawei Feng, Bo Ding, and Huaimin Wang. 2023. Uncertainty-penalized reinforcement learning from human feedback with diverse reward lora ensembles. <i>arXiv preprint arXiv:2401.00243</i> .		
1364			
1365			
1366			
1367			
1368	Hailin Zhang, Defang Chen, and Can Wang. 2022. Confidence-aware multi-teacher knowledge distillation. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 4498–4502. IEEE.		
1369			
1370			
1371			
1372			
1373	Hongxiang Zhang, Hao Chen, Muhao Chen, and Tianyi Zhang. 2025a. Active layer-contrastive decoding reduces hallucination in large language model generation. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 3028–3046.		
1374			
1375			
1376			
1377			
1378			
1379	Ling Zhang, Xianliang Yang, Juwon Yu, Park Cheonyoung, Lei Song, and Jiang Bian. 2025b. Holdout-loss-based data selection for llm finetuning via in-context learning. <i>arXiv preprint arXiv:2510.14459</i> .		
1380			
1381			
1382			

1414
1415
1416

A Notation List

Table 1 summarizes the core notation used throughout this survey; less common symbols are defined where they first appear.

Notation	Name	Description
x_i	query	the i -th query
y_i	LLM's response	response to the i -th query
c_j	candidate	the j -th candidate response of a query or a candidate reasoning trace
\mathcal{C}	candidate set	a set of candidate responses of a query or candidate reasoning traces
z_t	inference state	the model's inference state at step t , such as a partial reasoning trace, an intermediate generation state, or a dialogue state in multi-turn interaction
s_j, s_t	confidence score	confidence score for a candidate c_j or a reasoning step z_t
τ	threshold	confidence score threshold for inference termination or further refinement
ϵ	confidence tolerance	tolerance for changes in confidence used to decide reasoning termination
λ	contrastive weight	controls the strength of contrast between predictive views
π	inference policy	conditional distribution governing inference transitions, such as reasoning-step evolution or token generation

Table 1: Notations list