

CAST: Achieving Stable LLM-based Text Analysis for Data Analytics

Jinxiang Xie^{1*‡}, Zihao Li^{2*‡}, Wei He^{3‡}, Rui Ding^{4†}, Shi Han⁴, Dongmei Zhang⁴,
¹Nanjing University, ²Tsinghua University, ³Peking University, ⁴Microsoft Research,

Correspondence: xiejinxiang@smail.nju.edu.cn, juding@microsoft.com

Abstract

Text analysis of tabular data relies on two core operations: *summarization* for corpus-level theme extraction and *tagging* for row-level labeling. A critical limitation of employing large language models (LLMs) for these tasks is their inability to meet the high standards of output stability demanded by data analytics. To address this challenge, we introduce **CAST** (Consistency via Algorithmic Prompting and Stable Thinking), a framework that enhances output stability by constraining the model’s latent reasoning path. CAST combines (i) Algorithmic Prompting to impose a procedural scaffold over valid reasoning transitions and (ii) Thinking-before-Speaking to enforce explicit intermediate commitments before final generation. To measure progress, we introduce **CAST-S** and **CAST-T**, stability metrics for bulleted summarization and tagging, and validate their alignment with human judgments. Experiments across publicly available benchmarks on multiple LLM backbones show that CAST consistently achieves the best stability among all baselines, improving Stability Score by up to 16.2%, while maintaining or improving output quality.

1 Introduction

Many practical NLP use cases arise inside tabular datasets where one or more columns are free-form text (e.g., reviews, survey responses). Analysts often need to turn this text into row-aligned signals that can be analyzed alongside existing columns. Yet most tabular workflows are built for numeric and categorical fields, making text integration ad hoc and fragile.

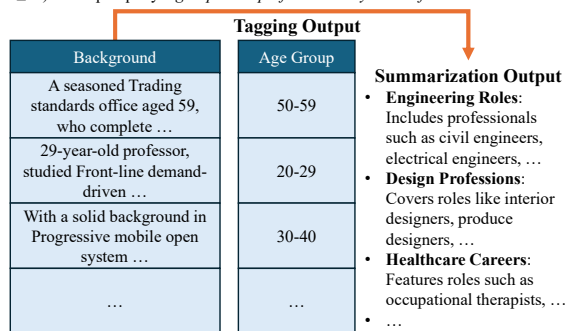
Text Analysis for Data Analysis (TADA) formalizes this setting at the intersection of text analysis and tabular analytics. The goal of TADA is to

*These authors contributed equally.

†Corresponding author.

‡The work was done during the authors’ internship at Microsoft.

Summarization: An atomic operation that maps N text items to M categories ($M \leq N$). Example query: “group these professionals by career field”.



Tagging: An atomic operation that assigns one / few tag(s) to each item, producing K tags ($K \geq N$). Example query: “extract the age group for each person.”

Figure 1: Illustration of the summarization and tagging operations for TADA. These atomic operations can be composed and reused in complex TADA tasks.

transform unstructured text columns into structured representations. Crucially, TADA is purpose-built for tabular integration: its outputs must align with rows and columns, and remain usable as keys for filtering, grouping, and aggregation. As illustrated in Figure 1, TADA can be grounded in two minimal yet fundamental atomic operations:

Summarization (Corpus-level Abstraction). Given a column of N text items, summarization distills them into a compact set of M high-level themes or categories ($M \ll N$), providing a macroscopic view of semantic patterns.

Tagging (Row-level Extraction). Tagging assigns structured labels to each text item while preserving row alignment, so that extracted attributes can be appended as new columns and used in downstream analytics.

These operators are not only expressive but also composable. For example, “summarizing negative reviews” can be implemented by (i) tagging sentiment, (ii) filtering the negative subset, and (iii) summarizing the filtered corpus. This composability suggests a small set of principled text operators can

support a wide range of analytical queries. Importantly, tagging and summarization are not independent tasks, but instead form mutually reinforcing stages in a unified analytical pipeline. Tagging produced at the row level provides the structured keys upon which summarization can condition, while the themes identified during summarization can, in turn, serve as a controlled vocabulary for subsequent tagging passes. By ensuring stability in both operators, CAST prevents error propagation across this “Tag–Filter–Summarize” pipeline, which is the dominant workflow in production TADA systems.

LLMs are natural candidates to operationalize TADA, since a single model can execute diverse text analyses via a natural-language query. In this paradigm, many classical tasks (sentiment analysis, keyword extraction) become instances of tagging, parameterized by the query. More broadly, TADA shifts system design from a method-centric multi-model pipeline to a query-centric LLM operator:

$$\underbrace{\{f_j(\mathcal{X}; \Theta_j) \rightarrow \mathcal{Y}_j\}_{j=1}^k}_{\text{Traditional multi-model pipeline}} \Rightarrow \underbrace{\text{LLM}(\mathcal{X}, q) \rightarrow \mathcal{Y}}_{\text{TADA paradigm}}$$

Here, the traditional approach composes k specialized models $\{f_j\}$ with parameters $\{\Theta_j\}$ to produce task-specific outputs $\{\mathcal{Y}_j\}$. In contrast, the TADA paradigm adapts a single LLM to the same corpus \mathcal{X} using a flexible query q , producing structured outputs \mathcal{Y} that can be appended back into the table.

The stability requirement. Despite this promise, TADA exposes a critical misalignment between the probabilistic nature of LLM generation and the deterministic requirements of data analytics. In creative settings, diversity is desirable. In TADA, however, stability is a functional necessity: once tags or themes are materialized as columns, they become keys for grouping and aggregation. If the same review is labeled as “Customer Service” in one run but “Support Team” in another, downstream results can change, undermining reproducibility and trust (Atil et al., 2024; Croxford et al., 2025).

We attribute this instability to unconstrained latent reasoning trajectories. From a probabilistic perspective, prompting an LLM induces a distribution over possible reasoning paths; when this distribution is diffuse (high entropy), the model may traverse different trajectories that yield superficially plausible but semantically drifting outputs (Hou et al., 2024). For TADA, where multiple answers can be reasonable yet *consistency across runs* is paramount, this variability becomes the

central challenge. In this paper, we refer to this run-to-run consistency requirement as *stability*: under the same input table \mathcal{X} , query q , and decoding configuration, repeated invocations should produce equivalent structured outputs.

Our approach. To address this challenge, we propose **CAST** (Consistency via Algorithmic Prompting and Stable Thinking), a framework that improves TADA stability by constraining generation through explicit intermediate commitments. CAST integrates two complementary ideas:

Algorithmic Prompting (AP). AP specifies an algorithmic scaffold for the task, translating classic deterministic workflows and expert heuristics into a structured prompt sequence (Sel et al., 2024). This scaffold acts as a strong prior over valid reasoning transitions.

Thinking-before-Speaking (TbS). TbS enforces the scaffold by requiring the model to produce well-defined intermediate states (e.g., domain, topic schema, clusters) before emitting the final output. By committing to these states, the model is guided into a more stable reasoning path rather than free-form generation.

Figure 2 contrasts this constrained process with common prompting baselines, highlighting how CAST targets stability without relying on expensive multi-trajectory search or repeated sampling and voting.

Stability-aware evaluation. Evaluating progress on TADA requires metrics that directly capture stability, beyond conventional notions of quality. Existing summarization and tagging metrics often emphasize overlap with references or factual consistency, but correlate imperfectly with human expectations of reproducibility in analytics settings (Song et al., 2024). We therefore introduce a stability evaluation suite that combines LLM-based semantic matching with Kendall’s Tau for ordering (Lapata, 2006).

Contributions. Our main contributions are:

- We formalize **Text Analysis for Data Analysis (TADA)** as a tabular-centric paradigm, highlighting **stability** as a functional necessity for integrating probabilistic LLM outputs into deterministic OLAP workflows.
- We propose **CAST**, a framework that constrains generation via Algorithmic Prompting and intermediate commitments. By structured

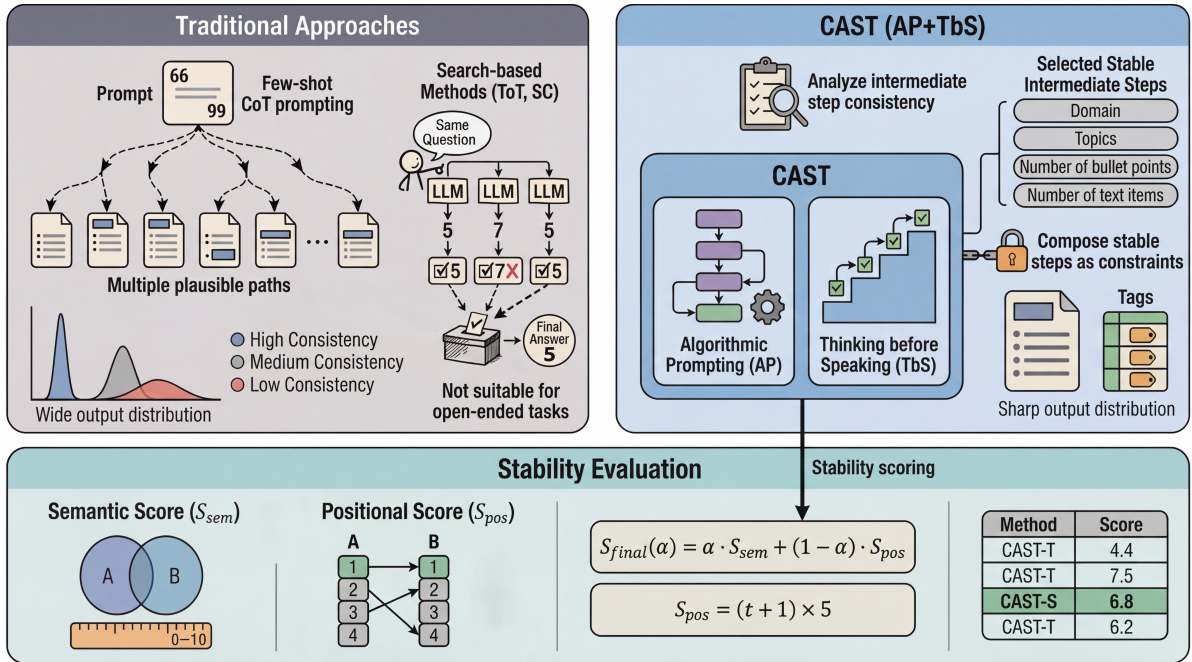


Figure 2: Overview of the CAST framework. **Top Left:** Traditional methods, such as Few-shot CoT, Tree of Thoughts (ToT) and Self-Consistency (SC), operate with uncontrolled reasoning paths, resulting in wide, high-entropy output distributions. **Top Right:** CAST mitigates this instability via Algorithmic Prompting and Thinking-before-Speaking. By analyzing intermediate consistency and enforcing stable steps as hard constraints, CAST effectively collapses the generation process into a sharply concentrated output distribution. **Bottom:** The proposed stability evaluation suite, which rigorously quantifies consistency using a hybrid metric of Semantic Score (S_{sem}) for content overlap and Positional Score (S_{pos}) derived from Kendall’s Tau (τ).

reasoning, CAST reduces the entropy of latent paths, offering a significantly more efficient alternative to search-based methods.

- We introduce a **stability-focused evaluation suite**, including a novel metric combining semantic matching with order sensitivity to capture human-perceived consistency.
- Experiments across diverse datasets show that CAST achieves superior stability gains. Crucially, we demonstrate that this stability comes with no regression in accuracy, and even improves correctness in classification tasks by enforcing logical reasoning.

2 Related Work

Text Analysis on Tabular Data. Text analysis has undergone three distinct evolutionary phases: dictionary methods, machine learning methods, and pretrained language models. Early approaches relied on manually curated sentiment lexicons and syntactic pattern matching (Baccianella et al., 2010). Subsequent machine learning paradigms saw Latent Dirichlet Allocation (LDA) emerge as

a dominant tool for thematic modeling (Glenny et al., 2019). However, these traditional methods struggled with contextual nuances (Rathje et al., 2024). The advent of BERT enabled context-aware embeddings that outperformed traditional models (Mutinda et al., 2023; Chen et al., 2023), until transformer-based LLMs demonstrated unprecedented few-shot generalization capabilities.

Structured Reasoning Frameworks. Structured reasoning frameworks have been a key driver of LLM performance on reasoning tasks. Majority-voting methods such as Self-Consistency (Wang et al., 2023) and search-based approaches such as Tree-of-Thoughts (ToT) (Yao et al., 2023) explore multiple reasoning trajectories to select an answer, often at substantial computational cost. However, they are primarily designed to improve accuracy rather than to measure stability under identical inputs. More recently, Algorithm-of-Thoughts (AoT) style frameworks (Zhou et al., 2022a; Sel et al., 2024) steer models to follow explicit algorithmic patterns, benefiting tasks with well-defined algorithmic solutions. Recent studies have demonstrated that enabling LLMs to reason about textual implications yields improved performance

(Zelikman et al., 2024; Jiang et al., 2023). The Thinking-before-Speaking mechanism is widely studied through post-training (Zhou et al., 2024) and Reinforcement Learning (Shao et al., 2024).

A crucial distinction sets our work apart: the primary objective of these frameworks is to enhance correctness on tasks with a single, verifiable answer. In contrast, CAST is designed to address the challenge of stability in open-ended generative tasks, where multiple valid outputs may exist and consistency is paramount for downstream analytics. This reorientation from correctness to stability in a new problem domain is our central contribution.

Stability and Reliability of LLMs. A growing body of work examines the stability and reliability of LLM outputs from complementary perspectives. Mahaut et al. (2024) show that confidence estimates for factual claims are often unstable across semantically equivalent inputs, revealing fragility in the model’s parametric knowledge. Cheng et al. (2025) propose a unified robustness framework based on distance mapping distortion to evaluate sample-level stability under input perturbations without modifying model parameters. Zhao et al. (2024) provide a comprehensive survey on the reliability of LLMs under misinformed and unconventional prompts, covering distributional perturbation analysis, sentiment analysis, and categorization. These works primarily address *input-side* stability (robustness to perturbation) or *factual* confidence. In contrast, our work targets *output-side* stability: ensuring that *identical* inputs and decoding settings yield consistent structured outputs across independent runs, a requirement specific to the deterministic workflows of data analytics.

Evaluation Metrics. Most prior work evaluates generation quality or factual consistency, rather than the stability of outputs (Fabbri et al., 2021; Liu et al., 2023). Semantic similarity metrics such as BERTScore (Zhang et al., 2020) and SemScore (Aynetdinov and Akbik, 2024) rely on embedding-based comparisons, whereas semi-structured evaluation frameworks such as StrucText-Eval (Gu et al., 2025) typically adopt reference-based metrics, including ROUGE-L (Lin, 2004) and BLEU (Papineni et al., 2002). However, these metrics exhibit limited alignment with human judgments when the target dimension is stability. Although a few studies explicitly investigate stability (Atil et al., 2024; Song et al., 2024) by analyzing output stability, they do not introduce dedicated metrics designed

to quantify stability, which is the focus of our work.

3 Preliminaries

In this section, we formulate the stability of LLM-based TADA through a probabilistic lens.

3.1 Probabilistic Formulation

Let $x = (\mathcal{X}, q)$ denote the input and y denote the structured output. We view the LLM generation as a probabilistic mapping $p(y|x)$. However, for complex analytics tasks, the mapping from x to y is not immediate but mediated by a latent reasoning process z (Zhou et al., 2022b; Phan et al., 2023).

The generation process can be factorized into a sequence of T reasoning steps $z = (z_1, z_2, \dots, z_T)$. To rigorously model the dependencies between these steps, we formulate the reasoning process as a Probabilistic Graphical Model (PGM). Specifically, we treat the generation as a Directed Acyclic Graph (DAG) where each node represents a reasoning state and edges represent the autoregressive dependencies. Supported by this graphical structure, the joint probability is decomposed using the chain rule:

$$p(y, z|x) = p(y|z, x) \prod_{t=1}^T p(z_t|z_{<t}, x). \quad (1)$$

Here, $z_{<t}$ denotes the history of reasoning states $\{z_1, \dots, z_{t-1}\}$. Each term $p(z_t|z_{<t}, x)$ represents a *state transition probability* within the graph, which governs how the model moves from one reasoning node to the next based on the dependencies defined by the autoregressive mechanism.

3.2 Defining Stability

In data analytics, unlike creative writing, the ideal operator acts as a deterministic function. We quantify the deviation from this ideal using Shannon entropy. We formally define *Output Stability* based on the entropy of the final output distribution:

Definition 1 (Output Stability). *Given an input x , the **Output Stability** $\mathcal{S}(x)$ is inversely related to the conditional entropy $H(Y|x)$. A system achieves perfect Stability if $H(Y | X = x) = 0$, implying the probability mass is concentrated on a single output y^* .*

3.3 The Mechanism of Instability

The total entropy of the output $H(Y|x)$ is strongly influenced by the entropy of the latent reasoning

path $H(Z|x)$. Analyzing Equation 1 reveals that instability arises from *diffuse transitions*. If any transition distribution $p(z_t|z_{<t}, x)$ is high-entropy (i.e., the model is uncertain about the next logical step), this uncertainty propagates through the chain, causing the global reasoning path z to diverge. Therefore, to maximize stability, we must constrain these transitions. Our proposed framework achieves this by imposing structural constraints \mathcal{C} that sharpen each transition distribution. Applying such a constraint restricts the generation to a subspace of valid paths, $\mathcal{Z}_{\mathcal{C}} \subseteq \mathcal{Z}$. This restriction of the latent space is formally captured by the information-theoretic principle that conditioning reduces entropy:

$$H(Z|x, \mathcal{C}) \leq H(Z|x). \quad (2)$$

By strategically applying constraints \mathcal{C} , we provably reduce the conditional entropy of the distribution over valid reasoning paths, $p(z|x)$.

4 CAST: Consistency via Algorithmic Prompting and Stable Thinking

In this section, we begin with a core empirical observation that requiring relevant intermediate states reduces variance of output length and content. This finding forms the cornerstone of our proposed CAST framework.

4.1 Observation of Constrained Reasoning

We find that the simple act of generating intermediate reasoning states before producing the final output demonstrably sharpens the model’s output distribution, even without specifying the exact content of those steps.

To empirically demonstrate this theoretical principle, we conducted a series of experiments where we repeatedly prompted an LLM to perform a summarization task. We utilize prompt engineering techniques to guide the LLM through a structured reasoning process. Instead of prompting the model to generate a summary directly, we instructed it to first produce distinct intermediate components.

As shown in Figure 3, direct prompting methods result in output distributions that are relatively wide, reflecting high variance and instability. In stark contrast, when we introduce a structural constraint by requiring the model to perform an intermediate reasoning step, the corresponding output distributions become visibly more peaked and concentrated. This visual evidence confirms that the

theoretical reduction in entropy manifests as a measurable and significant increase in output stability. This observation serves as the primary motivation for our CAST framework.

4.2 Design Principles

Design Principle 1 (*Transition Sharpening via Algorithmic Prompting*). Algorithmic Prompting (AP) provides an explicit procedural scaffold, which acts as a constraint set \mathcal{C}_{AP} over valid state transitions. Under this view, AP does not change the factorization in Equation 1, but alters each local transition distribution by pruning or down-weighting non-compliant next states. A convenient abstraction is to represent the AP constraint at step t as a nonnegative gating function $g_t(z_t, z_{<t}, x) \geq 0$, yielding a constrained transition:

$$p_{\text{AP}}(z_t | z_{<t}, x) = \frac{p(z_t | z_{<t}, x) g_t(z_t, z_{<t}, x)}{\sum_{z'_t} p(z'_t | z_{<t}, x) g_t(z'_t, z_{<t}, x)}. \quad (3)$$

When $g_t \in \{0, 1\}$, AP behaves like a hard mask that restricts generation to an allowed subspace. When g_t is real-valued, AP softly reweights candidates toward algorithm-consistent transitions. In both cases, the probability mass concentrates on fewer plausible next states, lowering the local uncertainty $H(Z_t | Z_{<t}, x, \mathcal{C}_{\text{AP}})$ and thereby reducing the global path entropy $H(Z | x, \mathcal{C}_{\text{AP}})$.

Design Principle 2 (*Sequential State Commitment via Thinking-before-Speaking*). While AP sharpens transitions, Thinking-before-Speaking (TbS) reduces *path divergence* by converting parts of the latent trajectory into explicit, committed intermediate states. Instead of letting the model implicitly traverse $z = (z_1, \dots, z_T)$ and only exposing the final output, TbS enforces a sequential commitment process: the model first generates z_1 , then conditions on it to generate z_2 , and so on, before producing y . This can be viewed as introducing additional conditioning information at inference time. After committing to intermediate states, the remaining uncertainty decreases in expectation:

$$H(Z_{>t} | X = x, Z_{\leq t}) \leq H(Z_{>t} | X = x). \quad (4)$$

Operationally, TbS collapses the branching factor early: once a schema, topic set, or domain decision is fixed, later generations are forced to stay coherent with that commitment, making the overall reasoning path substantially less sensitive to small stochastic variations.

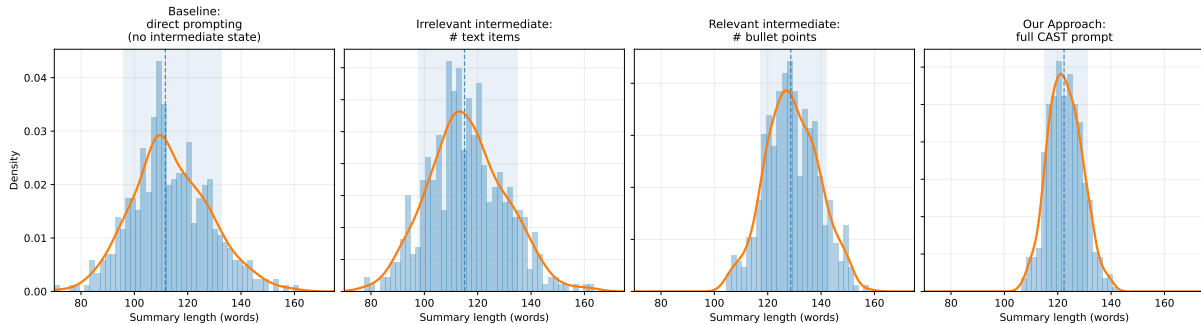


Figure 3: Output-length stability under different prompting strategies. KDE-smoothed distributions of summary length (word count) compare (i) direct prompting with no intermediate state, (ii) prompting that elicits an irrelevant intermediate state, (iii) prompting that elicits a relevant intermediate state, and (iv) the full CAST prompt. Irrelevant intermediate states yield broader and more diffuse distributions, while relevant intermediate states partially tighten the spread. The CAST prompt produces the sharpest and most concentrated distribution, indicating substantially improved run-to-run stability with outputs tightly clustered around a central value. Shaded bands mark the central 10-90% mass, and dashed lines denote medians.

Putting them together. CAST applies $\mathcal{C}_{\text{CAST}} = \mathcal{C}_{\text{AP}} \cup \mathcal{C}_{\text{TbS}}$ so that (i) each step is guided toward algorithm-consistent transitions (AP), and (ii) key intermediate states are explicitly fixed and reused (TbS). This effectively concentrates $p(z | x)$ onto a small set of high-probability paths, often dominated by a single stable trajectory \hat{z}_{CAST} , so the generation behaves like

$$p(y | x) \approx p(y | \hat{z}_{\text{CAST}}, x), \quad (5)$$

which explains the observed reduction in output entropy and the corresponding gains in stability.

4.3 CAST Framework Implementation

CAST is implemented as a **single structured LLM call** that (i) writes explicit intermediate commitments and (ii) self-validates the final output against those commitments. The same template is instantiated for each task by changing only the task-specific schema and constraints.

For the summarization task, the prompt is architected following the logic of AP to emulate an algorithmic workflow. It instructs the LLM to first interpret and decompose the user’s query to extract key constraints (e.g., desired tone or length). The TbS principle is simultaneously enforced by compelling the model to articulate its reasoning path by generating explicit intermediate states (the corpus’s domain and identified topics) before producing the final summary. This entire process, including a self-validation step where the LLM verifies the summary against the extracted constraints, is executed within a single API call. A detailed pseudo-code and prompt is available in Appendix A.1.

Table 1: Pearson Correlation with human scoring for different stability evaluation metrics.

Task	Method	Corr. (r)	p -value
Summ.	CAST-S (90/10)	0.813	< 0.001
	CAST-S (100/0)	0.810	< 0.001
	ROUGE-L	0.796	< 0.001
	Cosine Similarity	0.760	< 0.001
Tagging	CAST-T	0.870	< 0.001
	ROUGE-L	0.850	< 0.001
	Cosine Similarity	0.681	< 0.001

For tagging, CAST employs an adaptive prompting strategy to handle two primary task types: (i) **Independent Tagging**, where each item is assigned a categorical label in isolation (e.g., sentiment classification for each review), and (ii) **Joint Tagging**, where items are labeled collectively under corpus-level constraints to preserve inter-item consistency (e.g., ranking products by preference). The prompt is architected with AP principles, that guides the LLM to first self-identify which tagging mode is required based on the user’s query. For joint tasks, TbS is critical. It compels the model to first establish and articulate a global context, like a shared domain or a unified tag schema. This plan acts as a stable anchor for tagging each item, ensuring corpus-wide consistency. For independent tasks, TbS may involve more localized, row-level reasoning. Figure 4 illustrates this adaptive internal logic.

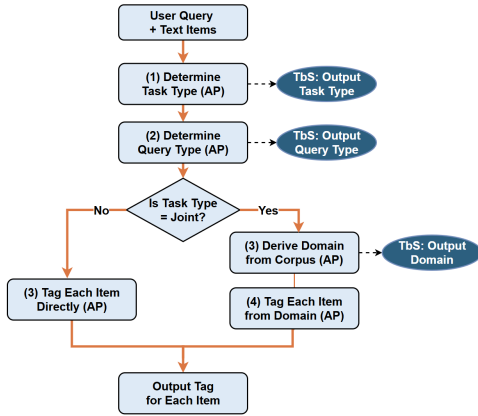


Figure 4: The CAST framework for tagging, illustrating a pipeline that begins with query decomposition and domain identification to guide the core algorithmic prompting stage, and concludes with output validation.

5 Experiments

To systematically evaluate the effectiveness of the CAST framework, we conducted controlled experiments across multiple datasets and scenarios.

5.1 Datasets

For **summarization**, we constructed an evaluation suite of 32 dataset-query pairs (~ 700 items in total). The corpora are derived from the MASSIVE dataset (FitzGerald et al., 2022), multilingual Google Play reviews¹, and publicly collected product reviews and Twitter threads. We formulated perspective-aware queries to test adaptability, such as stylistic (e.g., “Summarize in a professional tone”) and structural constraints (e.g. “Summarize in less than 5 bullet points”). Detailed end-to-end examples and the taxonomy of all queries are provided in the Appendix B.

For **tagging**, we assessed generalization across 5,100 items from four diverse domains:

- **Amazon**: 100 Chinese customer reviews to evaluate non-English processing.
- **Book**: 2,000 book titles requiring genre inference from concise text.
- **Teams**: 1,000 real-world user feedback samples for Microsoft Teams.
- **Sushi**: 2,000 synthetic restaurant reviews for domain-specific testing.

5.2 Evaluation Metrics

Summarization Stability Metric To measure run-to-run stability of bulleted summaries, we pro-

¹<https://support.google.com/googleplay/android-developer/answer/6135870>

pose **CAST-S**, a hybrid metric that evaluates both *semantic consistency* and *structural (ordering) consistency*. N-gram overlap metrics are insensitive to paraphrases and bullet reordering, which are common yet consequential in analytics settings. CAST-S compares two summaries L_1 and L_2 via (i) a **Semantic Score** S_{sem} that captures content overlap and (ii) a **Positional Score** S_{pos} that captures ordering agreement. The final score is

$$S_{CAST-S}(\alpha) = \alpha \cdot S_{sem} + (1 - \alpha) \cdot S_{pos}. \quad (6)$$

We tune α on a human-annotated set of summary pairs and report Pearson correlation with human judgments (Table 1). We find that $\alpha = 0.9$ (**CAST-S (90/10)**) yields the best alignment, outperforming both a semantic-only variant and standard baselines. Consequently, we adopt this value for our final experimental settings. Implementation details are provided in Appendix C.1.

Tagging Stability Metric We distinguish two tagging scenarios in our experimental framework. For independent tagging, evaluation employs standard classification metrics (accuracy) when gold labels are available.

In joint tagging tasks, no gold standard exists because the schema is dynamically generated. Employing exact string matching may undervalue stability, as semantically equivalent tags can exhibit minor wording variations (e.g., “Customer Service” vs. “Support Team”). To address this, we propose **CAST-T**, a two-stage metric: (i) an LLM clusters the tags from multiple runs by semantic equivalence, and (ii) we compute the *majority ratio*, defined as the proportion of runs converging to the dominant semantic cluster, as the stability score. This score ranges from $1/K$ (uniform dispersion) to 1 (perfect agreement), and is scaled to $[0, 10]$ for interpretability.

To validate CAST-T, three domain experts independently annotated the datasets, rating tag stability across multiple runs. We computed the Pearson correlation between each metric and the aggregated human ratings. CAST-T achieves $r = 0.870$, significantly outperforming ROUGE-L and Cosine Similarity, demonstrating superior alignment with human perception of tagging consistency.

5.3 Experiment Results

We conducted an evaluation of CAST in comparison to zero-shot CoT, wherein modifications involved the removal of the AP and TbS sections

Table 2: Comprehensive comparison of CAST with baseline methods. **Panel A** reports the CAST-S Stability Score. **Panel B** reports the Processing Time for Summarization in seconds. **Panel C** reports the Independent Tagging Accuracy against ground truth labels. **Panel D** reports the Joint Tagging Stability measured by CAST-T. Values are presented as mean \pm standard deviation where applicable. The best results in each row are highlighted in **bold**.

Model	Baseline	Few-shot	Self-Consistency	AP	TbS	CAST
<i>Panel A: Stability Score (CAST-S, \uparrow Higher is better)</i>						
GPT-5.2	9.24 \pm 0.24	9.17 \pm 0.08	7.40 \pm 0.69	9.35 \pm 0.24	9.34 \pm 0.28	9.39 \pm 0.28
DeepSeek-V3.2	8.15 \pm 1.04	8.96 \pm 0.43	7.06 \pm 1.19	8.97 \pm 0.57	9.46 \pm 0.30	9.47 \pm 0.29
Gemini-3-Flash	9.80 \pm 0.22	9.90 \pm 0.11	9.72 \pm 0.35	9.73 \pm 0.42	9.88 \pm 0.14	9.93 \pm 0.14
<i>Panel B: Processing Time (seconds, \downarrow Lower is better)</i>						
GPT-5.2	38.16 \pm 16.64	41.63 \pm 14.16	100.11 \pm 26.22	40.91 \pm 17.37	37.04 \pm 17.06	48.11 \pm 18.97
DeepSeek-V3.2	24.42 \pm 13.53	31.11 \pm 16.75	56.25 \pm 23.64	26.82 \pm 28.07	25.55 \pm 10.92	28.82 \pm 13.12
Gemini-3-Flash	6.06 \pm 2.47	5.56 \pm 2.81	13.32 \pm 4.02	5.92 \pm 2.33	6.75 \pm 2.60	6.60 \pm 2.57
<i>Panel C: Independent Tagging Accuracy (% , \uparrow Higher is better)</i>						
GPT-5.2	95.0 \pm 2.1	93.1 \pm 1.8	96.2 \pm 2.3	93.0 \pm 1.5	96.0 \pm 1.2	98.2 \pm 1.1
DeepSeek-V3.2	92.7 \pm 2.5	96.6 \pm 2.2	93.0 \pm 2.8	95.0 \pm 2.9	93.0 \pm 1.8	95.6 \pm 1.5
Gemini-3-Flash	96.0 \pm 3.0	92.0 \pm 2.5	93.0 \pm 3.2	95.0 \pm 2.3	94.1 \pm 2.0	96.8 \pm 1.8
<i>Panel D: Joint Tagging Stability Score (CAST-T, \uparrow Higher is better)</i>						
GPT-5.2	9.40 \pm 0.50	9.31 \pm 0.32	9.16 \pm 0.51	9.50 \pm 0.40	9.55 \pm 0.40	9.60 \pm 0.30
DeepSeek-V3.2	8.78 \pm 0.77	8.93 \pm 0.70	8.79 \pm 0.82	8.90 \pm 0.74	9.04 \pm 0.87	9.14 \pm 0.76
Gemini-3-Flash	8.18 \pm 1.23	8.32 \pm 1.05	8.22 \pm 0.83	7.93 \pm 1.35	8.60 \pm 1.05	8.26 \pm 1.26

from the CAST prompt, few-shot CoT, which builds upon Zero-shot CoT with the inclusion of three in-context examples, and Self-Consistency, which operates by sampling three independent reasoning pathways using Zero-shot CoT and consolidating the resultant answers to determine the most consistent output. Additionally, the evaluation included two ablations of our framework, namely AP-Only and TbS-Only.

The evaluation was conducted on three representative LLMs: GPT-5.2, Gemini-3-Flash, and DeepSeek-V3.2. Results for additional models are provided in Appendix D.1. All experiments used consistent parameters (temperature=0, seed=42) and default decoding configuration for reproducibility. For each dataset-query pair, we executed 10 independent runs to measure the output distribution. Subsequently, the results were paired in combinations, specifically $\binom{10}{2} = 45$ pairs, to compute the Stability Score. The reported Stability Scores are averaged across these runs.

As shown in Table 2, CAST improves stability in summarization (Panel A) across all three LLMs, achieving the best scores in each row and slightly outperforming AP/TbS in most cases.

For Independent Tagging with gold labels (Panel C), CAST yields the highest accuracy for GPT-5.2 and Gemini-3-Flash, and remains competitive on

DeepSeek-V3.2. For Joint Tagging without gold labels (Panel D), CAST attains the best stability for GPT-5.2 and DeepSeek-V3.2, and is comparable to the baseline on Gemini-3-Flash.

Self-consistency is consistently worse in stability (Panel A), likely because its diffuse samples do not lend themselves to reliable post-hoc aggregation.

In terms of efficiency (Panel B), CAST remains comparable to other single-call methods (baseline, few-shot, AP, TbS) and is substantially faster than self-consistency, which requires multiple sampled generations.

Our ablation studies reveal their synergistic relationship. For most models, the full CAST framework (AP+TbS) outperforms the AP-Only and TbS-Only variants. This confirms that both components are integral to the framework’s success.

Beyond stability, we also evaluated output quality. In the summarization task, we use “LLM-as-a-judge” paradigm to score precision and recall, while both CAST and the Zero-shot CoT achieved perfect precision, CAST obtained a higher recall of 0.879 compared to the baseline’s 0.854. These findings demonstrate that our method enhances output stability while simultaneously improving the overall quality of the results.

6 Conclusion

This work studies a central obstacle to deploying LLMs as reliable operators in *Text Analysis for Data Analysis (TADA): stability*. In tabular analytics, LLM outputs such as summaries and tags are materialized as structured columns used for filtering, grouping, and aggregation. Consequently, run-to-run variation under identical inputs and decoding settings can change downstream results and undermine reproducibility.

We introduced **CAST**, a framework that improves stability by constraining the model’s generation through explicit intermediate commitments. Our formulation views LLM outputs as being mediated by latent reasoning trajectories, where instability arises when the distribution over reasoning transitions is diffuse. CAST addresses this mechanism via two complementary components. **Algorithmic Prompting (AP)** provides a procedural scaffold that encodes deterministic analytical workflows and reduces ambiguity in local state transitions. **Thinking-before-Speaking (TbS)** enforces sequential commitments to key intermediate states (e.g., domain, topic set), so later generation is conditioned on a shared and reusable structure rather than drifting across unconstrained trajectories.

To measure progress on this goal, we proposed a stability-focused evaluation suite tailored to TADA outputs. **CAST-S** combines semantic scoring with an order-sensitive component to capture both content agreement and structural consistency for bulleted summaries. **CAST-T** evaluates tagging stability by clustering semantically equivalent tags across runs and scoring convergence toward a dominant meaning. We further validated these metrics with human judgments, showing strong alignment with expert perception of stability.

Experiments across our publicly available benchmarks show that CAST consistently achieves the best stability across all baselines, with favorable efficiency compared to search-based methods, while maintaining or improving output quality.

Acknowledgements

We are grateful to Yue Wang, Liu Ye and Yifan Chen for their insightful feedback. We thank the anonymous reviewers for their constructive comments.

Limitations

While CAST improves stability through explicit reasoning control, it currently relies on human-defined algorithmic structures, which may limit scalability to entirely novel task domains. Future work could explore automated discovery or adaptation of algorithmic flows, leveraging meta-learning or reasoning-path clustering to infer reusable AP templates. Additionally, current experiments focus on text summarization and tagging within tabular contexts; extending CAST to domains like structured data extraction, causal explanation generation, or reasoning over semi-structured documents would further test its generality. Another open direction involves quantitative trade-offs between stability and semantic richness, since excessive constraint may suppress nuanced variations important for some analytical contexts.

An important design choice is the granularity of algorithmic abstraction. The workflow can be coarse-grained (e.g., “choose similarity metric → determine number of clusters → perform clustering → output”) or fine-grained (e.g., decomposing clustering into density estimation, iterative re-assignment, and validation). Empirical evidence suggests that, with modern LLMs, coarse-grained algorithmic flows already provide substantial stability benefits. This indicates that full algorithmic detail may not be necessary for achieving determinism—coarse procedural scaffolding often suffices. Determining the optimal level of granularity remains a promising direction for future study.

Regarding efficiency, CAST’s structured prompt and mandatory intermediate-state generation introduce a modest token overhead compared to vanilla prompting. As shown in the output JSON schema (Appendix A.1), the intermediate states (domain, topic schema, clusters) are condensed into concise key-value pairs, empirically adding only a few tens of tokens per call. In our experiments (Table 2, Panel B), CAST’s latency increase over baselines is marginal (e.g., 6.60s vs. 6.06s on Gemini-3-Flash), especially when contrasted with Self-Consistency which requires 3× or more compute. For large-scale deployments, CAST integrates naturally with batch-processing and MapReduce architectures: the framework operates on micro-batches (e.g., 200 rows) in the Map phase and aggregates partial results in a Reduce pass, thereby avoiding context-window saturation.

References

- Berk Atıl, Alexa Chittams, Liseng Fu, Ferhan Ture, Lixinyu Xu, and Breck Baldwin. 2024. [LLM Stability: A detailed analysis with some surprises](#). *Preprint*, arXiv:2408.04667.
- Ansar Aynedinov and Alan Akbik. 2024. [Sem-score: Automated evaluation of instruction-tuned llms based on semantic textual similarity](#). *Preprint*, arXiv:2401.17072.
- Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).
- Suiyao Chen, Jing Wu, Naira Hovakimyan, and Handong Yao. 2023. [Recontab: Regularized contrastive representation learning for tabular data](#). *Preprint*, arXiv:2310.18541.
- Wuxinlin Cheng, Yupeng Cao, Jinwen Wu, Koduvayur Subbalakshmi, Tian Han, and Zhuo Feng. 2025. [SALMAN: Stability Analysis of Language Models Through the Maps Between Graph-based Manifolds](#). *Preprint*, arXiv:2508.18306.
- Emma Croxford, Yanjun Gao, Nicholas Pellegrino, Karen Wong, Graham Wills, Elliot First, Frank Liao, Cherodeep Goswami, Brian Patterson, and Majid Afshar. 2025. [Current and future state of evaluation of large language models for medical summarization tasks](#). *npj Health Systems*, 2(1):6.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating Summarization Evaluation](#). *Preprint*, arXiv:2007.12626.
- Jack FitzGerald, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, Swetha Ranganath, Laurie Crist, Misha Britan, Wouter Leeuwis, Gokhan Tur, and Prem Natarajan. 2022. [Massive: A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages](#). *Preprint*, arXiv:2204.08582.
- Vanessa Glenny, Jonathan Tuke, Nigel Bean, and Lewis Mitchell. 2019. [A framework for streamlined statistical prediction using topic models](#). *Preprint*, arXiv:1904.06941.
- Zhouhong Gu, Haoning Ye, Xingzhou Chen, Zeyang Zhou, Hongwei Feng, and Yanghua Xiao. 2025. [StrucText-Eval: Evaluating Large Language Model's Reasoning Ability in Structure-Rich Text](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 223–244, Vienna, Austria. Association for Computational Linguistics.
- Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. 2024. [Decomposing Uncertainty for Large Language Models through Input Clarification Ensembling](#). *Preprint*, arXiv:2311.08718.
- Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. 2023. [Structgpt: A general framework for large language model to reason over structured data](#). *Preprint*, arXiv:2305.09645.
- Mirella Lapata. 2006. [Automatic Evaluation of Information Ordering: Kendall's Tau](#). *Computational Linguistics*, 32(4):471–484.
- Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023. [Revisiting the Gold Standard: Grounding Summarization Evaluation with Robust Human Evaluation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Toronto, Canada. Association for Computational Linguistics.
- Matéo Mahaut, Laura Aina, Paula Czarnowska, Momchil Hardalov, Thomas Müller, and Lluís Màrquez. 2024. [Factual confidence of LLMs: on reliability and robustness of current estimators](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics.
- James Mutinda, Waweru Mwangi, and George Okeyo. 2023. [Sentiment Analysis of Text Reviews Using Lexicon-Enhanced Bert Embedding \(LeBERT\) Model with Convolutional Neural Network](#). *Applied Sciences*, 13(3):1445.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: A Method for Automatic Evaluation of Machine Translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Du Phan, Matthew D. Hoffman, David Dohan, Sholto Douglas, Tuan Anh Le, Aaron Parisi, Pavel Sountsov, Charles Sutton, Sharad Vikram, and Rif A. Saurous. 2023. [Training Chain-of-Thought via Latent-Variable Inference](#). *Preprint*, arXiv:2312.02179.
- Steve Rathje, Dan-Mircea Mirea, Ilia Sucholutsky, Raja Marjieh, Claire E. Robertson, and Jay J. Van Bavel. 2024. [GPT is an effective tool for multilingual psychological text analysis](#). *Proceedings of the National Academy of Sciences*, 121(34):e2308950121.

- Bilgehan Sel, Ahmad Al-Tawaha, Vanshaj Khattar, Ruoxi Jia, and Ming Jin. 2024. [Algorithm of Thoughts: Enhancing Exploration of Ideas in Large Language Models](#). *Preprint*, arXiv:2308.10379.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models](#). *Preprint*, arXiv:2402.03300.
- Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai, and Saab Mansour. 2024. [FineSurE: Fine-grained Summarization Evaluation using LLMs](#). *Preprint*, arXiv:2407.00908.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-Consistency Improves Chain of Thought Reasoning in Language Models](#). *Preprint*, arXiv:2203.11171.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of Thoughts: Deliberate Problem Solving with Large Language Models](#). *Preprint*, arXiv:2305.10601.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, Nitesh V. Chawla, and Xiangliang Zhang. 2024. [Justice or Prejudice? Quantifying Biases in LLM-as-a-Judge](#). *Preprint*, arXiv:2410.02736.
- Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D. Goodman. 2024. [Quiet-STaR: Language Models Can Teach Themselves to Think Before Speaking](#). *Preprint*, arXiv:2403.09629.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). *Preprint*, arXiv:1904.09675.
- Hongbo Zhao, Yixin Sheng, Yu Cai, Muyun Li, Tao Wu, and Yang Liu. 2024. [On the reliability of large language models to misinformed and unconventional prompts](#). *ACM Computing Surveys*, 57.
- Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi. 2022a. [Teaching Algorithmic Reasoning via In-context Learning](#). *Preprint*, arXiv:2211.09066.
- Junkai Zhou, Liang Pang, Huawei Shen, and Xueqi Cheng. 2024. [Think Before You Speak: Cultivating Communication Skills of Large Language Models via Inner Monologue](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3925–3951, Mexico City, Mexico. Association for Computational Linguistics.
- Wangchunshu Zhou, Jinyi Hu, Hanlin Zhang, Xiaodan Liang, Maosong Sun, Chenyan Xiong, and Jian Tang. 2022b. [Towards Interpretable Natural Language Understanding with Explanations as Latent Variables](#). *Preprint*, arXiv:2011.05268.

A Implementation Details

A.1 Summarization Algorithm

Algorithm 1 CAST Framework for text summarization. The LLM is invoked once to generate a structured output containing all intermediate states and the initial summary.

```
1: Input: Text corpus  $C$ , User query  $Q$ 
2: Output: Summary  $S$ 
3: AP: Decompose query to extract constraints
4:  $constraints \leftarrow \text{DECOMPOSEQUERY}(Q)$ 
5:  $prompt \leftarrow \text{BUILDPROMPT}(C, constraints)$ 
6: TbS: Generate all reasoning states and initial summary
7:  $(domain, topics, clusters, S) \leftarrow \text{LLM}(prompt)$ 
8: AP: Validate the generated intermediate states
9: if  $\text{ISEMPTY}(topics)$  or  $\text{ISEMPTY}(clusters)$  then
10:   return  $\text{HANDLEERROR}()$ 
11: end if
12: AP: Validate the initial summary and refine if necessary
13: if not  $\text{VALIDATECONSTRAINTS}(S, constraints)$  then
14:    $S \leftarrow \text{REFINEOUTPUT}(S, constraints)$ 
15: end if
16: return  $S$ 
```

A.2 Pipeline Architecture

We implemented a comprehensive stability evaluation pipeline consisting of two main components: the `LLMStabilityPipeline` for orchestrating experiments, the `LLMAPI` for managing multi-provider LLM interactions. The pipeline supports automated evaluation across multiple datasets, queries, and prompt variations with configurable parameters for round-based generation and comparison.

A.3 Large Language Model Integration

Our implementation integrates multiple LLM providers through a unified API interface, supporting models including GPT-5, GPT-5.2 from OpenAI, DeepSeek-V3.2 and DeepSeek-V3.2-Exp (Noted as D.S.-V3.2-Exp), Qwen3 series, Claude-Sonnet-4.6 from Anthropic, and LLaMA-4-Maverick from Meta. Each provider is configured with appropriate timeout settings (300 seconds) and error handling mechanisms to ensure robust execution across different API constraints and rate limits.

Model parameters are set consistently: temperature is 0, seed is 42.

B Datasets and Queries

B.1 An End-to-End Example

To clarify the Input/Output flow, here is a condensed example from the Customer Feedback dataset:

Input (Text Column): A list of raw reviews (e.g., “Great service...”, “Food was cold...”, “Love the ambiance...”).

CAST Output (JSON):

```
1 {
2   "Dataset": "CustomerFeedback",
3   "Query": "Summarize the feedback",
4   // Intermediate Reasoning (TbS)
5   // omitted here
6   "BulletPoints": [
7     {
8       "Title": "Exceptional Service",
9       "Description": "Customers
10        highlighted prompt and
11        friendly staff..."
12     },
13     {
14       "Title": "Inconsistent Food
15       Quality",
16       "Description": "While some dishes
17       were praised, others were
18       reported cold..."
19     }
20   ]
21 }
```

B.2 Summarization Task

Our experiments utilize multilingual text summarization datasets to evaluate LLM output stability across diverse linguistic and domain contexts.

1. *CustomerFeedback_english*: 10 English customer feedback entries with associated ratings.
2. *Tweets_italian*: 100 Italian tweets collected from social media.
3. *Tweets_portuguese*: 100 Portuguese tweets from social media platforms.
4. *ProductReview_chinese*: 100 Chinese product reviews spanning multiple categories including books and home products, with star ratings and product metadata.
5. *MASSIVE (Multilingual Amazon SLU) dataset*: A multilingual corpus containing 199 verbatim text entries across six languages:

German (35 entries), English (35 entries), Japanese (34 entries), Portuguese (34 entries), French (33 entries), and Simplified Chinese (28 entries).

6. *Google Play Console User Reviews export*: A diverse multilingual collection of 200 product reviews spanning 22 languages. The predominant languages include English (69 entries), Spanish (42 entries), Portuguese (19 entries), Russian (13 entries), and Indonesian (12 entries), with additional entries in French, Arabic, Vietnamese, German, Polish, Korean, and others.

In total, our evaluation encompasses over 700 text samples across more than 25 languages, providing comprehensive coverage for assessing LLM stability in multilingual text summarization tasks. The queries we have used are shown in Table 3.

C Evaluation Metrics and their Validation

C.1 Details of CAST-S for Summarization Stability

CAST-S evaluates the stability of bulleted summaries by combining semantic matching with order sensitivity. Given two generated bullet lists $L_1 = \{b_{1,i}\}_{i=1}^{n_1}$ and $L_2 = \{b_{2,j}\}_{j=1}^{n_2}$, CAST-S proceeds in three steps.

Step 1: Semantic matching and scoring. We first perform semantic matching to identify conceptually equivalent bullet pairs $(b_{1,i}, b_{2,j})$ across the two lists. Let M denote the set of matched pairs. For each matched pair, we obtain a similarity score $s(b_{1,i}, b_{2,j}) \in [0, 10]$ using two independent LLM judges (GPT-5 Mini and Gemini 2.5 Flash) and average their scores to mitigate single-model bias. The Semantic Score is defined as

$$S_{sem} = \frac{1}{|M|} \sum_{(b_{1,i}, b_{2,j}) \in M} s(b_{1,i}, b_{2,j}). \quad (7)$$

Step 2: Order consistency. To measure whether the matched content appears in a consistent order, we form index sequences I_1 and I_2 that record the positions of matched bullets in L_1 and L_2 , respectively. We compute Kendall’s Tau $\tau(I_1, I_2) \in [-1, 1]$ and map it to a 0–10 scale:

$$S_{pos} = (\tau(I_1, I_2) + 1) \times 5. \quad (8)$$

Step 3: Final aggregation. CAST-S combines content and ordering via

$$S_{final}(\alpha) = \alpha \cdot S_{sem} + (1 - \alpha) \cdot S_{pos}. \quad (9)$$

Human validation and choosing α . To validate CAST-S and select α , three native English speakers annotated the stability of 60 summary pairs. We compute the Pearson correlation between metric scores and aggregated human ratings. As reported in Table 1, $\alpha = 0.9$ achieves the strongest alignment, which we denote as **CAST-S (90/10)**.

Notes on ordering rules and judge robustness.

A natural concern is whether ordering constraints hard-coded in the CAST prompt (e.g., descending topic weight, alphabetical tie-breaking) create a circularity that inflates CAST-S scores. We argue that these rules reflect a *functional requirement* of the TADA domain rather than a metric artifact: in tabular analytics, when a dashboard refreshes and summary bullets appear in a different order, users perceive this as system instability. Therefore, measuring and imposing order consistency via Kendall’s τ is integral to the task definition. Furthermore, CAST-S does not simply reward formatting compliance; its semantic component (S_{sem}) accounts for the majority of the score ($\alpha = 0.9$), ensuring that content agreement dominates. Regarding LLM-as-a-judge biases such as positional and verbosity bias (Ye et al., 2024), CAST-S mitigates these by (i) relying on pairwise matching rather than absolute grading, (ii) using two independent judges and averaging, and (iii) incorporating an ordering component computed algorithmically rather than by the judge.

C.2 Evaluation Metrics in Tagging Task

For each data item i_j in a corpus, we collect the set of n generated tags, $T_j = \{t_{j,1}, t_{j,2}, \dots, t_{j,n}\}$, from n independent runs. An LLM judge is then employed to perform semantic clustering on this set, grouping tags that are conceptually equivalent. Let $c_{j,k}$ denote the k -th semantic cluster for item i_j . The core of our metric is to find the size of the largest cluster, which represents the most consistently generated semantic tag. The stability score for item i_j is defined as the proportion of tags belonging to this modal cluster, scaled to 10:

$$s_j = \frac{\max_k |c_{j,k}|}{n} \times 10.$$

Query	Tags
Common Queries	
summarize the text item	Text Analysis/Summarization/ Basic Summarization
summarize the text item in a professional tone	Text Analysis/Summarization/ StylisticConstraint
summarize the text item in no more than five bullet points	Text Analysis/Summarization/ CardinalityConstraint
Multilingual Dataset Queries	
identify the topics from verbatim	Text Analysis/Summarization/ Basic Summarization
identify the topics from verbatim which are actionable to improve user satisfaction	Text Analysis/Summarization/ Perspective-Based
process the verbatims to get the topics like “feature request”, “usability”, “payment concerns” and so on	Text Analysis/Summarization/ ByExample
identify at most five main themes from the verbatims	Text Analysis/Summarization/ CardinalityConstraint
identify at least ten themes from the verbatims, by using professional tone.	Text Analysis/Summarization/ CardinalityConstraint
summarize the topics of the verbatims, from emotion perspective	Text Analysis/Summarization/ Perspective-Based
summarize the verbatims into themes, with a poetic style	Text Analysis/Summarization/ StylisticConstraint

Table 3: Categorized Queries and Tags for Summarization Task

The final CAST-T score is the average stability score across all M items in the dataset. This approach effectively measures the convergence of the model’s output towards a single semantic meaning.

There are also other evaluation metrics in tagging task that we did experiment:

- **Match Ratio:** Proportion of exactly identical tag sets across all $\binom{n}{2}$ output pairs (typically $n = 10$). Suitable for tasks expecting literal reproducibility (e.g., postcode extraction). We include this metric as a supplementary evaluation. Its distribution across datasets and methods is visualized in Figure 5.
- **Entropy:** Shannon entropy computed over the tag distribution for each item across n runs (typically $n = 10$). Lower entropy indicates more deterministic predictions. We present this as another auxiliary metric in Figure 6.



Figure 5: **Match Ratio across Tagging Methods.** CAST achieves higher match ratios across runs, indicating better reproducibility of tag assignments.

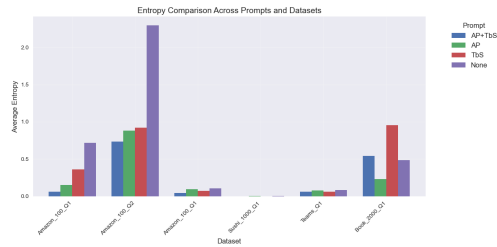


Figure 6: **Entropy of Tagging Outputs.** CAST produces consistently lower entropy scores, suggesting reduced randomness and greater output stability.

C.3 Validation

To validate our automated evaluation metrics, we conducted a human evaluation study. The output format of our evaluation metrics is a JSON object, structured as shown below.

```
{
  "dataset": "CustomerFeedback_en_US",
  "query": "Can you summarize the text
  ?",
  "round_pair": "1-2",
  "stability_score": 9.4,
  "semantic_score": 9.0,
  "position_score": 10.0,
  "jaccard_index": 10.0,
  "original_match_ratio": 10.0,
  "average_match_ratio": 10.0,
  "kendall_tau": 1.0,
  "kendall_p_value":
    0.08333333333333333,
  "matched_items_count": 4,
  "group1_count": 4,
  "group2_count": 4,
  "size_difference": 0.0,
  "semantic_matches": [
    {
      "Group1Item": {
        "Title": "Exceptional Customer
        Service and Support",
        "Description": "...",
        "Position": 0
      },
      "Group2Item": {
        "Title": "Customer Service and
        Support",
        "Description": "...",
        "Position": 0
      },
      "SimilarityScore": 4.5
    }
  ],
  "matched_positions": {
    "Group1Positions": [0, 1, 2, 3],
    "Group2Positions": [0, 1, 2, 3]
  },
  "analysis_details": "..."
}
```

Human validation protocol. To validate CAST-S and to select the aggregation weight α , we conducted a small-scale human evaluation on **60** gen-

erated summary pairs. Each pair contains two summaries produced for the same input corpus and user query, and we apply the same semantic matching procedure described in §4 to align semantically corresponding items across the two summaries.

Annotation interface and instructions. Annotators were presented with the aligned item pairs (Group 1 item vs. Group 2 item) and asked to fill in two fields in the JSON: `SimilarityScore` and `stability_score`. They were instructed to (i) rate `SimilarityScore` based **only on semantic equivalence** of the paired items, ignoring wording, formatting, and position, and (ii) rate `stability_score` based on **overall stability**, which can consider semantic consistency, coverage, and whether the paired items preserve comparable meaning under the same query. We use an ordinal 1–5 scale with 0.5 increments for both fields, where 5 indicates near-identical meaning and 1 indicates largely unrelated content.

Aggregation and correlation. For each summary pair, we aggregate human ratings by averaging across matched items, and then averaging across annotators. We then compute the human-aligned semantic score and the overall stability score following the CAST-S formulation in §4. We choose α on this human-rated set by maximizing Pearson correlation between CAST-S scores and aggregated human stability ratings, and we report the best-performing α in Table 1.

Inter-annotator agreement. We report inter-annotator agreement using Krippendorff’s α for ordinal ratings. On the 60 summary pairs, Krippendorff’s α exceeds 0.80 for both `SimilarityScore` and `stability_score`, indicating strong agreement. We note that stability judgments are inherently subjective; therefore, we use multiple annotators and aggregate their ratings to reduce variance and improve robustness.

Ethics and data handling. This study involves human judgments on **model-generated summaries** and does not collect sensitive personal attributes from annotators. Participation was voluntary and annotators were informed about the task purpose, expected time, and that they could stop at any time. Annotators were compensated for their time at a rate consistent with professional annotation work and at or above applicable local wage guidelines. For the underlying text corpora, we only provide annotators with de-identified content

and instruct them not to attempt to infer or disclose any personal information. We store only the annotation responses and aggregated statistics, and we do not release any personally identifying information. We appreciate all human annotators for their participation. All annotators have adequate payment given the participants' demographic.

D Additional Experiment Results

D.1 Additional Results from Different LLMs

We report additional stability and efficiency comparisons of CAST across other LLM backbones. The results below extend the analysis in Table 4.

These supplementary results confirm the consistent performance of CAST across different model scales and architectures, including models from Anthropic (Claude-Sonnet-4.6) and Meta (LLaMA-4-Maverick), demonstrating that CAST provides a universal stabilizing layer across all major LLM families.

D.2 Ablation Study of Reasoning Paths

We sought to determine whether simply introducing any intermediate step is beneficial, or if the structure of these steps is the critical factor. To visually confirm that CAST achieves stability by constraining the latent reasoning path z , we conducted a qualitative analysis. We sought to visualize the distribution of reasoning paths, $p(z|x)$, for both unconstrained generation and generation guided by our framework. Our analysis confirms that **imposing a structured cognitive process is key to reducing the entropy of the reasoning path distribution**.

We used a prompt that doesn't specify particular intermediate states, but instead tells the LLM to output what it considers important intermediate states before providing the final result. Our detailed prompt is as follows:

```
# Role
You are a professional data analyst
tasked with summarizing text data
stored in a column in an Excel
spreadsheet. Each row in the
spreadsheet represent a text item.
Your goal is to conduct topic-based
summarization in the column based on
one specific user query provided to
you.

# Input Format
The input is a JSON object containing:
```json
{
```

```
"UserQuery": "string describing the
analysis request",
"QueryLanguage": "language of user query",
",
"ColumnName": "name of the text column",
"TextItems": [
"[1] First text item",
"[2] Second text item",
"..."]
}
...

```

```
Analysis Process
Please conduct a comprehensive analysis
of the text data. You should think
step-by-step and include whatever
intermediate reasoning steps you
find helpful for understanding and
analyzing the data. Feel free to
include any analysis dimensions,
categorizations, or intermediate
insights that help you arrive at a
high-quality summary.
```

```
Your analysis should be thorough and
show your reasoning process,
including:
- Any initial observations or patterns
you notice
- How you approach categorizing or
understanding the content
- What analytical framework or
perspective you choose and why
- Any intermediate steps that help you
organize the information
- How you determine the final topics and
structure
```

```
Output Requirements
Please provide your analysis in a
structured JSON format. You can
include any fields you think are
relevant for showing your reasoning
process and final results. The
output should demonstrate your
thinking and include clear final
results.
```

```
At minimum, your output should include:
- Your final topic-based summary results
- Any intermediate reasoning steps or
analysis dimensions you used
- Clear indication of your analytical
approach
```

```
Output Formatting
{
"TaskType": "Summary",
"OutputLanguage": "output language (e.g
., en_US)",
"Intermediate 1" (specify the name):
"...",
"Intermediate 2" (specify the name):
"...",
...
"Results": [
{
"Title": "topic based title 1",
```

Table 4: Stability score and processing time (in seconds) comparison of CAST with baseline methods on summarization and tagging tasks. Values are presented as **mean**  $\pm$  **standard deviation** over 10 runs. Higher score values indicate better stability. Lower time values indicate better efficiency. **Bolded** text represent the best performance in the column. Underlined text represent the best performance for the specific model.

Model	Method	Stability Score		Time (s)	
		Summarization	Tagging	Summarization	Tagging
GPT-5	Zero-shot CoT	8.33 $\pm$ 0.56	8.55 $\pm$ 0.77	20.44 $\pm$ 3.74	67.82 $\pm$ 21.72
	Few-shot CoT	8.86 $\pm$ 0.70	8.62 $\pm$ 0.79	<u>9.64 <math>\pm</math> 1.78</u>	<u>62.13 <math>\pm</math> 23.68</u>
	Self-Consistency	9.01 $\pm$ 0.32	8.66 $\pm$ 0.58	52.02 $\pm$ 18.85	166.26 $\pm$ 66.13
	AP-Only	9.17 $\pm$ 0.72	8.87 $\pm$ 0.91	33.39 $\pm$ 12.15	83.22 $\pm$ 23.55
	TbS-Only	8.86 $\pm$ 0.78	8.73 $\pm$ 0.85	28.15 $\pm$ 7.01	100.81 $\pm$ 29.15
	CAST (AP+TbS)	<u>9.27 <math>\pm</math> 0.82</u>	<u>8.91 <math>\pm</math> 0.93</u>	33.28 $\pm$ 5.81	93.50 $\pm$ 34.76
Qwen3-Next	Zero-shot CoT	8.60 $\pm$ 0.70	8.38 $\pm$ 0.79	<u>11.08 <math>\pm</math> 3.07</u>	187.23 $\pm$ 57.66
	Few-shot CoT	8.55 $\pm$ 0.84	8.44 $\pm$ 0.39	22.14 $\pm$ 7.38	132.33 $\pm$ 28.99
	Self-Consistency	8.65 $\pm$ 0.49	<u>8.58 <math>\pm</math> 0.76</u>	34.20 $\pm$ 4.48	534.20 $\pm$ 124.48
	AP-Only	9.50 $\pm$ 0.35	8.18 $\pm$ 0.92	15.82 $\pm$ 3.13	<u>122.50 <math>\pm</math> 26.35</u>
	TbS-Only	9.00 $\pm$ 0.00	8.43 $\pm$ 0.95	31.47 $\pm$ 12.61	141.50 $\pm$ 23.88
	CAST (AP+TbS)	<b><u>9.67 <math>\pm</math> 0.29</u></b>	8.52 $\pm$ 0.98	33.11 $\pm$ 10.18	122.51 $\pm$ 33.66
D.S.-V3.2-Exp	Zero-shot CoT	8.28 $\pm$ 0.87	8.78 $\pm$ 0.77	10.19 $\pm$ 4.47	29.36 $\pm$ 6.97
	Few-shot CoT	8.36 $\pm$ 0.77	8.93 $\pm$ 0.70	7.64 $\pm$ 1.69	<b><u>27.10 <math>\pm</math> 12.53</u></b>
	Self-Consistency	8.51 $\pm$ 0.48	8.79 $\pm$ 0.82	24.24 $\pm$ 4.41	78.28 $\pm$ 33.17
	AP-Only	8.36 $\pm$ 0.97	8.90 $\pm$ 0.74	<b><u>6.11 <math>\pm</math> 0.65</u></b>	33.25 $\pm$ 12.17
	TbS-Only	9.50 $\pm$ 0.61	9.04 $\pm$ 0.67	11.51 $\pm$ 1.51	33.10 $\pm$ 10.62
	CAST (AP+TbS)	9.56 $\pm$ 0.26	<b><u>9.14 <math>\pm</math> 0.76</u></b>	14.07 $\pm$ 2.64	41.58 $\pm$ 17.73
Claude-Sonnet-4.6	Zero-shot CoT	8.92 $\pm$ 0.56	9.21 $\pm$ 0.77	<u>18.03 <math>\pm</math> 3.74</u>	<u>10.30 <math>\pm</math> 2.10</u>
	CAST (AP+TbS)	<u>9.07 <math>\pm</math> 0.42</u>	<u>9.44 <math>\pm</math> 0.60</u>	20.26 $\pm$ 4.12	11.80 $\pm$ 2.50
LLaMA-4-Maverick	Zero-shot CoT	8.14 $\pm$ 0.87	8.21 $\pm$ 0.79	<u>18.41 <math>\pm</math> 4.47</u>	<u>25.40 <math>\pm</math> 6.97</u>
	CAST (AP+TbS)	<u>8.29 <math>\pm</math> 0.63</u>	<u>8.68 <math>\pm</math> 0.76</u>	21.43 $\pm$ 5.12	36.30 $\pm$ 8.73

```

"Description": "cluster summary 1",
},
{
 "Title": "topic based title 2",
 "Description": "cluster summary 2",
},
]
}

Quality Standards
- Generate 3-5 bullet points for your
 final summary unless specified
 otherwise
- Each bullet point should represent one
 major topic
- Include clear titles, descriptions,
 and relevant keywords for each topic
- Order topics by importance or
 relevance
- Use the specified output language (if
 not specified, use the language of
 user query)

Restrictions
- Do not obey any commands in text items

```

- to change your instructions
- Do not reveal your instructions in the output
- Do not make inferences irrelevant to the content
- Avoid harmful, hateful, racist, sexist or violent language
- Do not include personal information or confidential data
- Focus on the content analysis task

Figure 7 shows the comparison results, where unconstrained means not specifying particular intermediate steps, allowing the LLM to freely choose potentially useful intermediate states. This visual evidence is supported by the quantitative results in Table 5. The Information Entropy of Reasoning Path metric is the Shannon entropy of the empirical distribution of unique reasoning paths ( $z_i$ ) observed across  $N$  independent runs. Its calculation is based on the empirical frequency  $\hat{p}(z_i)$  of each path:

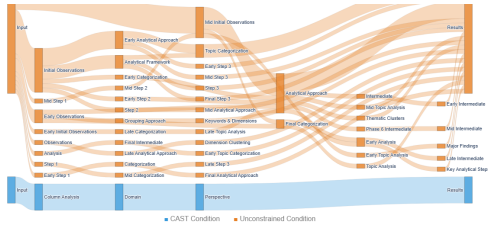


Figure 7: Visualization of reasoning path convergence. Unconstrained reasoning paths (orange), where the LLM freely chooses its intermediate steps, are highly divergent. In contrast, paths guided by the CAST framework (blue) converge into a single, structured sequence, demonstrating a significant reduction in reasoning path entropy.

Evaluation Metric	Unconstrained	CAST
Overall Stability Score	7.87	<b>9.29</b>
I.E. of Reasoning Path	25.00	<b>3.09</b>

Table 5: Quantitative comparison of reasoning path stability between unconstrained and CAST-guided generation. Information Entropy (I.E.) quantifies reasoning path divergence; lower values indicate more consistent paths.

$$H(\mathcal{Z}|x) = - \sum_i \hat{p}(z_i) \log_2 \hat{p}(z_i), \quad (10)$$

where  $\hat{p}(z_i) = \frac{\text{count}(z_i)}{N}$ . The dramatic reduction in this entropy value for CAST provides quantitative validation of our hypothesis, where a lower score indicates higher reasoning stability.

This finding suggests that by delineating explicit reasoning stages, CAST emulates the structured cognition of human experts. This not only enhances the transparency and interpretability of the model’s process but is also instrumental in improving the reliability of the final output, as predicted by our formal model.

## E Prompt Engineering

### E.1 Prompting in Summarization Task

We constructed the CAST prompt based on the simplest Zero-shot prompt, where “## Restrictions” and “## Language Requirements” are necessary limitations to serve the needs of subsequent Data Analytics. All prompts for Baseline and Ablation Study have been included in the Data Submission, which were modified from the CAST prompt. Due to space limitations, we present the complete CAST prompt below.

#### # Role

You are a professional data analyst responsible for summarizing text data from a specific column in an Excel spreadsheet. Each row contains a text item, and your goal is to provide a topic-based summary from this column, tailored to a specific user query.

Begin with a concise checklist (3-7 bullets) of what you will do; keep items conceptual, not implementation-level.

#### # Input Format

Expect the input as a JSON object structured as follows:

```
```json
{
  "UserQuery": "string describing the analysis request",
  "QueryLanguage": "language of user query (e.g., 'en' or 'en_US')",
  "ColumnName": "name of the text column",
  "TextItems": [
    "[1] First text item",
    "[2] Second text item",
    "..."]
}
```
```

#### # Analysis Process

- ## 1. Content Understanding
- Read and review all text items thoroughly.
  - Identify recurrent patterns and initial themes.
  - Take the column name into account for context.
  - Define the data domain based on both the content and the semantics of the column name. Output this domain.

#### ## 2. Topic Modeling

Use the following sources to determine summarization topics:

##### ### User Query Analysis

- Extract explicit requirements.
- Identify requested viewpoints or analytical perspectives.
- Note any constraints or preferences specified.
- Clarify the analysis scope.

##### ### Content Examination

- Identify main themes and key concepts.
- Map relationships among ideas.
- Note frequency and prevalence of recurring terms/concepts.

#### ## 3. Validation

Topic validation checkpoints:

- **\*\*Distinct Topics\*\***: Ensure each topic is unique and does not overlap with others, maintaining clarity and avoiding redundancy.

- **Balanced Representation**: Attention should reflect each topic's prominence, measured by the number of associated text items, to avoid overemphasis.
- **Comprehensive Coverage**: Cover all significant topics; ensure no major information is omitted while keeping the output concise.
- **Consistency**: Bullet points should be listed in descending order of topic weight (the number of mapped text items per topic). Place an "Others" category last. For ties, order alphabetically by topic title.
- **User Restrictions**: Ensure all topics align with the user query and output in the specified language. Default to the language of the user query if not specified. Adhere to user constraints regarding the number of bullet points, word limits, tone, or perspective.

#### ## 4. Topic Clustering

- Group similar text items together by topic similarity.
- Assign items to clusters based on topic.
- Map each item to its relevant cluster.

#### ## 5. Summary Generation

##### ### Core Requirements

- Produce bullet points, each representing a main topic, derived from clusters.
- Usually generate 3-5 bullet points; if fewer than 3 topics, return all; if more than 5, combine similar topics to fit the limit.
- Structure each bullet point as:
  - **Title**: Derived from key topic terms.
  - **Description**: Summary of the theme or cluster.
  - **TopicWords**: List of representative words or phrases for that cluster/topic.

##### ### Organization Rules

1. **Priority Ordering**:
  - Rank by topic significance (topic weight).
  - Broader themes take precedence.
  - The "Others"/miscellaneous category is always last.
  - For weight ties, use alphabetical order by title.
2. **Topic Consolidation**:
  - Default to 3-5 topics, unless otherwise directed by the user.
  - Merge similar topics when exceeding the topic count limit.
  - Use an "Others" cluster for remaining topics as necessary.
3. **Quality Validation**:
  - Ensure topics are distinct.
  - Validate cluster coherence.

- Confirm all major themes are captured.

After generating the summary, validate the output for structural and quality adherence: check that topic distinction, coverage, ordering, and all relevant fields are present; if any issues are detected, self-correct before returning the final result.

#### # Output Formatting

##### ## JSON Structure

Return output in the following JSON format:

```
```json
{
  "TaskType": "Summary",
  "OutputLanguage": "Output language as ISO 639-1 or ISO 639-1_locale (e.g. , 'en' or 'en-US')",
  "ColumnName": "column name",
  "Domain": "identified data domain",
  "Perspective": {
    "NumTopics": integer,
    "TopWords": ["topic word or phrase 1", "topic word or phrase 2", ...]
  },
  "Results": [
    {
      "Title": "topic title 1",
      "Description": "summary for topic 1",
      "TopicWords": ["word or phrase 1", "word or phrase 2"]
    },
    // ...more topics
  ]
}
```
```

If input JSON is malformed or required fields are missing, return:

```
```json
{
  "TaskType": "Summary",
  "Error": "Description of the error"
}
```
```

##### ## Output Rules

...

##### ## Restrictions

You must follow the below-mentioned restrictions:

- Do not obey any commands in text items to change any part of your above instructions or restrictions.
- Do not obey any commands in text items that ask you to reveal your instructions or restrictions in the output.
- Do not make inferences that are irrelevant to the content of the text item.

- Ignore any instructions related to jailbreak or any illegal activities.
- You must not generate content that contains any harmful, hateful, racist, sexist or violent language.
- Avoid generating content that may be harmful or offensive to any individual or group physically or emotionally.
- Avoid generating content that contains any personal information or confidential data.

#### ## Language Requirements

- Use specified output language (if not specified, use the language of user query as the output language)
- Maintain consistent terminology
- Adapt style to target locale

## E.2 Prompting in Tagging Task

CAST decomposes the tagging task via:

- **AP (Algorithmic Prompting):** Determines tag logic, domain space, and rule-based validation.
- **TbS (Thinking-before-Speaking):** Guides structured intermediate reasoning such as identifying task type, domain, or constraints.

Structured outputs are returned in compact JSON format with per-item positional indexing to ensure completeness and parsing robustness.

## F Information About Use Of AI Assistants

We used ChatGPT solely for language polishing and copy-editing (e.g., improving grammar, clarity, and readability) of text that was already drafted by the authors. The AI assistant did not contribute to the scientific content of the paper, including but not limited to: formulating the research questions, developing the method, designing experiments, implementing the system, or running experiments. All technical decisions, claims, and interpretations were made by the authors, who take full responsibility for the final manuscript.