# MEASURING INTENT COMPREHENSION IN LLMS: A VARIANCE DECOMPOSITION FRAMEWORK

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

016

017

018

019

021

023

024

025

026

027

028

029

031

033 034

037

038

040

041

042 043

044

046

047

048

051

052

#### **ABSTRACT**

People judge interactions with large language models (LLMs) as successful when outputs match what they want, not what they type. Yet LLMs are trained to predict the next token solely from text input, not underlying intent. Because written language is an imperfect proxy for intent, and correlations between phrasing and desired outcomes can break down in training data, models that rely too heavily on surface cues may respond inconsistently to semantically equivalent prompts. This makes it essential to evaluate whether LLMs can reliably infer user intent—especially in high-stakes settings where robustness and generalization are critical. We introduce a formal framework for assessing intent comprehension in LLMs: whether a model demonstrates robust understanding of user intent by producing consistent outputs across semantically equivalent prompts while differentiating between prompts with distinct intents. Our evaluation approach is based on a variance decomposition of model responses into three components: variability due to user intent, user articulation, and model uncertainty. Models that understand what users want, and are not overly sensitive to textual cues, should attribute most output variance to intent differences, rather than articulation style. Applying this framework across diverse domains, we find that larger models typically assign a greater share of variance to intent, indicating stronger comprehension of intent, although gains are uneven and often modest with increasing model size. These results motivate moving beyond accuracy-only benchmarks toward semantic diagnostics that directly assess whether models understand what users intend.

#### 1 Introduction

Human communication involves a fundamental socio-cognitive process: we form an intent we want to articulate, then we choose words to express it, hoping the recipient will sympathetically decipher our intended meaning Smith (1759). Consider a frustrated traveler at an airport asking "Is there any way to get to Terminal B faster?" They might equally say "I need to catch my flight—what's the quickest route to B?" or "Terminal B is so far—any shortcuts?" They could even ask indirectly, "How do I get to the gate for flight 718?" Each phrasing differs dramatically, yet all express the same underlying intent: finding the fastest path to their destination. Much communication, therefore, centers on extracting underlying intent from observable signals, with success depending on the recipient's ability to see past surface variation to grasp the purpose beneath.

This challenge of intent extraction becomes particularly critical for large language models (LLMs), where all interaction occurs through text, making the model's ability to understand user intent the first step in effective human-AI interaction. Despite the centrality of intent understanding, most model evaluation approaches focus on whether models can perform specific tasks. These evaluation frameworks assume users express their intentions clearly, but this assumption often fails in practice. Users frequently struggle to articulate their needs and may revise their requests multiple times until they express their intent clearly. This implies that for models to be reliable from a user perspective, they must not simply respond to literal text, but rather infer the underlying intent from limited textual cues, disregard unimportant surface variations, and respond to the user's true purpose.

In this paper, we propose a framework for measuring how well models capture users' underlying intentions. We treat intent as a latent variable that underlies every prompt: while the same purpose can be expressed in many ways, a model that understands intent should produce consistent response

distributions across surface variations, yet shift appropriately when purpose changes. We define intent comprehension as the property that responses remain invariant to phrasing when purpose is fixed, but vary systematically when purpose changes.

To operationalize this idea, we present a diagnostic method that decomposes variation of model outputs into three components: Intent Sensitivity, the share of variation due to changes in purpose or intent; Articulation Sensitivity, the share due to phrasing variations; and Model Uncertainty, the residual variation stemming from the model's inherent uncertainty. The variation in a model that understands user intent should exhibit high Intent Sensitivity and low Articulation Sensitivity, indicating consistent meaning extraction that disregards superficial linguistic cues.

Because intent is unobserved and generating equivalent prompts that differ only in form is non-trivial, we construct semantically equivalent prompts through cross-lingual translation, inspired by universal semantic structures across languages Youn et al. (2016). Starting from a base prompt, we translate it through a sequence of typologically diverse languages and back to English, inducing natural variation in phrasing while preserving intent. Translation serves this purpose well because its core objective aligns with our needs: altering surface form while maintaining meaning. Additionally, LLMs demonstrate strong automated translation capabilities, making this approach both theoretically motivated and practically feasible.

Using this pipeline, we evaluate five language models of varying sizes across tasks in health, logistics, finance, travel, and social planning. Our findings reveal that larger models generally attribute a greater share of output variability to changes in user intent (Intent Sensitivity), indicating stronger alignment with user goals and more coherent internal representations. However, this improvement is not uniform: the 70B model does not consistently outperform smaller models across all domains, and its robustness advantages are often modest. Larger models also demonstrate greater sensitivity to prompt articulation, reflecting a trade-off between semantic generalization and responsiveness to surface cues. Furthermore, model robustness varies significantly across domains. These findings demonstrate our framework's diagnostic value and highlight the need for more targeted evaluations of models' ability to respond to user intent rather than varied textual cues.

Intent in Computational Models. The concept of intent has emerged as a fundamental construct in artificial intelligence to bridge observable behaviors and underlying cognitive states. In natural language processing, intent traditionally represents the underlying purpose or goal behind user utterances Qin et al. (2021); Zhang and Wang (2022), evolving from early slot-filling paradigms to sophisticated neural architectures Louvan and Magnini (2020). This aligns with philosophical foundations in Brentano's intentionality theory—the directedness of mental states toward objects Jacob (2019)—and Gricean pragmatics, which emphasizes the role of communicative intentions in meaning Grice (1957; 1989). Recent work has formalized these intuitions through a Bayesian calculus, most notably Baker, Tenenbaum, and Saxe's inverse planning framework Baker et al. (2007; 2009), which models intent recognition as Bayesian inference over an agent's goals given observations of their actions under an assumption of approximate rationality.

For large language models, the relationship between intent understanding and world models presents both theoretical and empirical challenges. While LLMs demonstrate sophisticated performance on theory of mind tasks Kosinski (2023) and pragmatic reasoning Hu et al. (2023), questions persist about whether they possess genuine intentionality or rely on surface-level statistical associations Bender et al. (2021); Marcus (2022). The distinction between surface form and underlying intent becomes particularly salient when evaluating whether LLMs develop genuine world models—causal representations supporting counterfactual reasoning—versus sophisticated pattern matching Andreas (2022); Li et al. (2023). Recent Bayesian frameworks for intent recognition incorporate multiple observation sources through recursive filtering Javdani et al. (2018); Jiang et al. (2025), providing a principled approach for modeling uncertainty over user goals. This variance decomposition approach, where intent serves as a latent variable influencing response distributions, offers a promising methodology for distinguishing genuine semantic understanding from statistical association in LLMs, with consistent intent-response mappings across surface variations indicating robust world model properties versus brittleness to linguistic changes suggesting surface-level processing Vafa et al. (2024).

**Contributions** This paper makes three key contributions. First, we use these foundations to propose a formal definition of intent comprehension, specifically tailored to language models, grounded in the

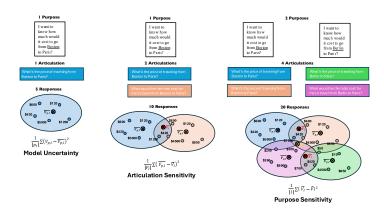


Figure 1: Conceptual Illustration of our decomposition Measures

notion of response invariance to surface variations. Second, we introduce a variance decomposition method that quantifies the relative contributions of intent, phrasing, and uncertainty to model outputs, providing interpretable measures of model behavior. Third, we apply this method to real-world tasks across multiple domains, demonstrating how our framework reveals meaningful differences in the the ability of models to understand users intent.

#### 2 CONCEPTUAL FRAMEWORK

In this section, we provide the conceptual foundation for our proposed measure in section 3. Let  $T \in \mathcal{T}$  denote users' intent. This intent or purpose represents the communicative goal underlying a user's request to the model. Users articulate their intent through prompts  $p_i$ , which may vary widely in surface form even when seeking to convey the same underlying objective. For example, the intent of knowing the capital of France may be expressed through prompts such as "What's the capital of France?" or "Which city is France's capital?". Finally, let  $a \in \mathcal{A}$  be the model response.

We say that a language model understands the user's intent if it produces identical response distributions for requests or prompts that express the same underlying intent, and distinct response distributions for prompts that express different intents. We formalize this as follows:

**Definition 2.1 (Intent Comprehension)** Let  $\tau: \mathcal{P} \to \mathcal{T}$  be the ground-truth mapping that assigns each prompt an intent label, and write  $\tau^{-1}(T) = \{ p \in \mathcal{P} \mid \tau(p) = T \}$  for the set of prompts that express intent T. A language model that generates a response distribution  $\pi$  over  $\mathcal{A}$ , is said to understand the user's intent if, for every pair of prompts  $p_i, p_j \in \mathcal{P}$ , the following properties hold:

1. Consistency: If  $p_i, p_j \in \tau^{-1}(T)$  for some  $T \in \mathcal{T}$ , then

$$\pi(\cdot \mid p_i) = \pi(\cdot \mid p_j).$$

2. Sensitivity: If  $\tau(p_i) \neq \tau(p_j)$ , then

$$\pi(\cdot \mid p_i) \neq \pi(\cdot \mid p_j).$$

Our definition requires a model that understands user intent to distinguish between surface-level variations in articulation and substantive changes in user intent or purpose. In particular, a model that varies its response distribution in reaction to inconsequential articulation changes is likely overfitting to superficial correlations in the training data rather than genuinely understanding the user's underlying purpose.

Remark 1 Our definition of Intent Comprehension (IC) is directly tied to the broader notion of a world model. In robotics and model-based control, agents infer a latent state  $s_t$  and learn dynamics/observation maps (f,g) so that behavior depends on  $s_t$  rather than raw measurements. Performance is typically judged by correctness-centric criteria such as one-step prediction error, trajectory likelihood, or task return under a known oracle. In our text-only setting, the prompt stream serves as observations, and the payoff-relevant latent state is the user's intent. As in robotics, a language model with a usable world model should (i) behave invariantly across paraphrases that preserve this state and (ii) shift its output distribution when the state (intent) changes. Crucially, unlike standard evaluations of world models—which assess both invariance across equivalent states and whether outcomes are correct relative to an external oracle—IC is consistency-centric: it tests whether the LLM can identify and respond to the latent state rather than the surface cues. IC therefore focuses on the first property of a world model: recognizing its stable latent state.

Remark 2 In our framework, intent captures what the user wants the model to say—specifically, the expected distribution over exact responses. Two prompts have the same intent if they reliably elicit the same output distribution from the model, even if they differ in wording or domain. For example, "Will a fair coin land heads or tails? Answer 1 for heads, 0 for tails" and "I drew a number from [-1, 1]. Is it positive or negative? Answer 1 for positive, 0 for negative" differ in content but share the same intent: both aim to produce a 50/50 distribution over identical outputs—'1' and '0'.

Our criteria for whether a model has intent comprehension may be over-exacting. Precisely defining intent is inherently difficult due to its latency and potential ambiguity. Evaluating model responses adds another layer of complexity. This often requires estimation of the distribution of full-text responses, which are inherently high-dimensional (e.g., a detailed description of a driving itinerary from Boston to Miami.) Focusing on the distribution of responses, an evaluator may view distinct answers as carrying the same underlying meaning. For example, in response to the prompt "How much is 1+2?", a model might answer "three" or "3." Although slight variations in the prompt may influence the likelihood of generating "three" versus "3," we should still regard the model as consistent.

To overcome these challenges, we propose a more relaxed notion of intent comprehension: a *sufficient* intent comprehension. This concept softens the standard definition by introducing an evaluator who judges prompts and responses. Specifically, we assume the existence of two functions: one that determines whether two prompts express the same intent, and another that assigns values to different responses. We say that a language model understands an the user intent if, whenever two prompts are judged to share the same intent, the distributions over evaluator-assigned values of responses remain identical. In this sense, we say that it's sufficient to say that a model understands the user's latent intent if the user cannot distinguish between model responses to the same intent.

Formally, let  $\tilde{\tau}: \mathcal{P} \to \mathcal{T}$  map from prompts to intents (i.e., articulations to purposes), defined according to some evaluation criterion. Let  $V: \mathcal{A} \to \mathbb{R}$  be a function that maps responses to values. Finally, let  $\pi_V(\cdot \mid p)$  denote the distribution over response values induced by model responses to prompt p, as evaluated by V. The relaxed definition is as follows:

**Definition 2.2 (Sufficient Intent Comprehension with respect to an evaluator**  $\tilde{\tau}$  **and** V) *Let*  $\tilde{\tau}$  :  $\mathcal{P} \to \mathcal{T}$  *be a fixed intent–evaluator that assigns an intent label for each prompt, and let*  $\pi_V(\cdot \mid p)$  *be the value distribution induced by the model, in response to prompt p. A language model possesses a sufficient intent comprehension with respect to \tilde{\tau} and V if for every pair of prompts p\_i, p\_i \in \mathcal{P}:* 

1. (Sufficient Consistency) If  $\tilde{\tau}(p_i) = \tilde{\tau}(p_j)$ , then

$$\pi_V(\cdot \mid p_i) = \pi_V(\cdot \mid p_j).$$

2. (Sufficient Sensitivity) If  $\tilde{\tau}(p_i) \neq \tilde{\tau}(p_j)$ , then

$$\pi_V(\cdot \mid p_i) \neq \pi_V(\cdot \mid p_j).$$

intent, and not simply their textual variation.

Remark 3 Both Intent Comprehension and sufficient Intent Comprehension criteria focus on measuring consistency, rather than assessing whether a language model's responses are factually correct. This is because our primary goal is to evaluate how well the model understands user intent—not whether it produces the correct answer. This approach differs from traditional benchmarks, which typically rely on a binary notion of correctness, classifying answers as strictly right or wrong. To understand why consistency matters, consider a language model trained only on economic data from 2020. If asked about economic conditions in 2023, it may consistently produce outdated yet internally coherent answers across differently phrased questions, reflecting the information on which it was trained. While these answers are incorrect in light of present-day facts, their internal consistency suggests that the model is responding to the user's underlying

Our definition of a (sufficient) intent comprehension is motivated by LLM training data, which consists of human-generated text, replete with spurious correlations between writing style and user intent. Small stylistic changes—like tone or word choice—often co-occur with shifts in latent attributes of the writer, such as personality, mood, or communicative norms, even when the underlying intent and objective remains constant. This makes prompt phrasing a confounded signal, blending true intent with incidental articulation. As a result, a model may change its responses not because users change their intent, but because it has learned correlations between superficial cues and different responses. Evaluating whether an LLM understands the latent user intent, thus mirrors the causal inference problem: Can it separate variation due to intent from variation due to articulation? In a hypothetical world where prompt phrasing perfectly reveals intent and carries no noise, models would trivially satisfy our definition. But in reality, disentangling intent from correlation is non-trivial—making consistency a meaningful test of deeper understanding.

# 3 MEASURING (SUFFICIENT) INTENT COMPREHENSION

In this section, we present a simple method, motivated by our conceptual framework, for evaluating whether an LLM understands the user's latent intent. We assume access to a sample of prompts along with corresponding model responses. Some prompts are designed to share the same underlying intent, while others are not. In the experimental section, we provide a detailed description of how we construct such a sample.

Our evaluation approach is based on a variance decomposition of model responses. Specifically, we break down response variance into three components: *Intent Sensitivity* (IS), which captures the model's responsiveness to changes in the intent or purpose of the input prompt; *Articulation Sensitivity* (AS), which reflects the model's sensitivity to variations in how users express the same intent; and Model Uncertainty (MU), which accounts for the inherent ambiguity or variability in the model's responses.

To formalize these components, we consider a domain D consisting of related purposes or intents  $I \in D$ . Each intent can be expressed through multiple prompts  $p_I \in \tau(I)$ , and we denote the model's value of response to a prompt p as v. To enable comparability across tasks, domains, and models, we define the standardized response as  $\tilde{v} = \frac{v}{\operatorname{std}(v|D)}$ . Furthermore, let  $R_I^2$  be the standard coefficient of determination, the proportion of total variance explained by differences in intent, and  $R_{p_I}^2$  be the proportion of variance explained by differences in prompt phrasing, holding intent I fixed. With this notation, we define our three core measures, sketched in Figure 1, as:

$$\begin{split} IS &= Var_I(\mathrm{E}[\tilde{v}|I]) = R_I^2 & \text{(Intent Sensitivity)}, \\ AS &= \mathrm{E}_I\left[Var_{p|I}\left(\mathrm{E}\left[\tilde{v}\mid p_I\right]\right)\right] = \mathrm{E}_I\left[R_{p_I}^2\right] & \text{(Articulation Sensitivity)}, \\ MU &= \mathrm{E}_I\left[\mathrm{E}_{p|I}\left(\mathrm{Var}\left(\tilde{v}\mid p_I\right)\right)\right] = \mathrm{E}_I\left[1-R_{p_I}^2\right] & \text{(Model Uncertainty)}. \end{split}$$

The first term, IS, is the variance of the conditional expectation of responses given intent. It reflects the sufficient Intent Comprehension property in 2.2, quantifying how much the model's mean response shifts with changes in the user's underlying purpose. A model appropriately sensitive to intent should produce distinctly different responses for different tasks. Therefore, we expect IS to be high in models that understand the user's latent intent or purpose. The second term, AS, corresponds to the consistency property. It quantifies how much, on average, the model's mean response varies

with changes in phrasing or wording of the prompt, holding the user's purpose fixed. A model that understands the user's latent intent, rather than being swayed by syntactic or stylistic differences, should exhibit low sensitivity to surface-level changes in articulation.

The third term, MU, captures the residual variance that remains after conditioning on both intent and phrasing. This component conflates at least three sources: (i) sampling stochasticity (e.g., due to temperature or nucleus sampling); (ii) epistemic uncertainty, where the model lacks a confident internal representation and spreads probability mass over multiple plausible answers; and (iii) aleatoric uncertainty, which is inherent to the task itself (e.g., when the prompt is "Tell me a random joke"). Whether a high MU is desirable depends on the context. For deterministic tasks with a clear correct answer (e.g., "What is 1+1?"), high MU reflects unwanted uncertainty. In contrast, for inherently open-ended or subjective tasks—such as those we examine in the experimental section—some degree of response variability is expected and even appropriate.

These three measures are all non-negative and sum to one: IS + AS + MU = 1.. This identity ensures that each term represents the relative contribution of a distinct source of variation in model responses. Together, they provide a principled way to assess the robustness and interpretability of a model's behavior.

An alternative perspective on this decomposition is in terms of  $\mathbb{R}^2$ : how much of the variation in model responses could be predicted using only information about the request. For example, consider asking the model how long it takes to travel from a given city to Paris. If the origin is New York versus London, we would expect the response distribution to change accordingly. A model that understands the user's question should produce responses such that a simple predictor using only the origin can explain nearly all the variation, apart from irreducible noise.

We also define two summary statistics that condense the decomposition into intent-centric diagnostics. The first is the Meaningful Variability Share (MVS),

$$MVS = \frac{IS}{IS + AS},$$

which serves as a signal-to-noise ratio over the *explainable* portion of variability: among the variance attributable to request features—either genuine changes in intent (IS) or superficial changes in articulation (AS)—MVS measures how much is signal rather than surface noise. A high MVS indicates that most of the model's explainable variation reflects meaningful distinctions between user intents, whereas a low MVS suggests the model is overly influenced by superficial differences in wording, indicating a lack of intent understanding. Another interpretation is that, when decoding is effectively deterministic (temperature near zero) so that residual randomness is negligible, MVS approximates the two-way split between intent and articulation, because variability then stems only from intent and phrasing.

To align directly with our definition of sufficient intent comprehension—distributions should move with intent (sensitivity) and remain invariant to non-intent articulation (consistency)—we combine MVS with the overall share of total variance due to intent into a single index,

$$\label{eq:ici} \text{ICI } = \frac{2\,\text{MVS} \cdot IS}{\text{MVS} + IS} \in [0,1],$$

which we call the *Intent Comprehension Index (ICI)*. ICI is the harmonic mean of "intent purity" and "intent coverage." The MVS component captures the sufficient consistency by penalizing articulation-driven variability: if non-intent phrasing moves the distribution, AS rises, MVS falls, and so does ICI. The IS component captures the sufficient sensitivity clause by rewarding models whose outputs shift when intent changes: if the distribution barely moves with intent, IS is small and ICI again falls. Therefore, the ICI summarizes the definition of intent comprehension.

This construction mirrors the classic precision–recall analogy in information retrieval. Precision corresponds to the purity of explainable variance, MVS = IS/(IS + AS), while recall corresponds to the coverage of total variance by intent, IS/(IS + AS + MU) = IS. ICI therefore plays the role of an F1-type summary—an F1 of "precision" (MVS) and "recall" (IS)—balancing the two desiderata in a single scalar.

Remark 4 In practice, when estimating the components of the decomposition, we deliberately restrict attention to domains where requests share a common structure but differ in intent. This controlled setup reflects a conservative and more stringent approach to evaluating whether the model understands the user's intent. By fixing the overall structure of the query and varying only a key element that changes the underlying intent (such as modifying the income level in a tax-related prompt while keeping the rest of the wording unchanged), we eliminate spurious cues and force the model to rely on its understanding. If the model systematically adapts its output to such minimal yet meaningful changes, it suggests the presence of a coherent internal representation capable of capturing user intent and purpose.

# 4 EXPERIMENTS

In this section, we provide a brief description of how we construct our experiments; a detailed description is presented in Section D of the Appendix. We begin by designing an automatic question generator. Our questions focus on open-ended, guesstimation-style tasks. These are chosen because they naturally elicit a wide distribution of plausible responses, which is essential for disentangling intent sensitivity from articulation sensitivity and model uncertainty. We utilize an LLM (GPT-4.1) to construct 24 questions across five domains—transportation, personal finance, health and nutrition, logistics, and social planning—via a two-stage pipeline: first, we generate unit-specified templates with placeholders, and then we populate them with realistic values to define specific intents. For each task, we generate 12 distinct intents. This procedure yields a diverse and structured set of quantitative estimation tasks that reflect naturalistic usage while remaining well controlled.

To evaluate articulation sensitivity, we need to generate sets of prompts that convey the same intent. This is non-trivial, as creating prompts with exactly equivalent definitions is difficult even for humans. To address this, we generate equivalent prompts through cross-lingual translation chains, leveraging structural differences across languages to induce lexical and syntactic diversity while preserving meaning. GPT-based checks filter paraphrases to ensure intent equivalence, and we then select a diverse final set using an embedding-based approach. Each paraphrase is paired with varying input values and repeatedly posed to different LLMs. Specifically, for each of the 24 tasks, we generate 10 paraphrases. We then submit these prompt variations to five different language models—LLaMA 3.2 3B, LLaMA 3.1 8B, LLaMA 3.3 70B, and Google's Gemma 3 12B-IT and 27B-IT—requesting 25 responses for each prompt—intent pair. In total, we obtain 1,772,492 model responses, averaging 354,499 per model. Finally, we perform the variance decomposition described in the section 3. To mitigate finite-sample issues, we apply ANOVA for the variance decomposition, which we discuss further in Section C of the Appendix.

#### 4.1 RESULTS

We now turn to explore our results. Figure 2 presents our key findings across models. The figure shows the average Intent Sensitivity (IS), Articulation Sensitivity (AS), and Model Uncertainty (MU) across all tasks and across models. The figure reveals that the two smallest models, LLaMA 8B and LLaMA 3B, exhibit very high MU across tasks, accompanied by both low IS and low AS. This is accompanied by higher variance in general, as shown in figure 5. In contrast, the three larger models—LLaMA 70B, Gemma 27B, and Gemma 12B—show comparable results across these three components, with the largest share attributed to IS. LLaMA 70B shows slightly higher MU. This suggests that larger models do not necessarily perform uniformly better, and that high IS can be achieved with relatively modest parameter counts, as demonstrated by the 12B model.

Figure 3<sup>2</sup> presents our results for the Intent Comprehension Index (ICI), Meaningful Variability Share (MVS), and IS across models. The blue line illustrates that larger models yield higher IS, indicating that changes in user intention are more effectively reflected in adjustments to the response distribution. Specifically, the smaller LLaMA models account for approximately 20% of the response variation due to changes in intention, while the larger models account for around 50%.

<sup>&</sup>lt;sup>1</sup>We only keep track of valid responses, but sometimes the model fails to reach the desired set of responses within 125 attempts, and in such cases, we move on to the next section.

<sup>&</sup>lt;sup>2</sup>Full distribution of the component for each tasks, by model, is in figure 6 in the Appendix

The second component, MVS, shown by the orange line, demonstrates inverse pattern: the two smaller models exhibit higher MVS than their larger counterparts. This implies that smaller models respond relatively more to changes in intent than changes in prompt articulation, aligning with the general lower responsiveness suggested by their Intent Sensitivity scores. Interestingly, larger models show greater sensitivity to articulation shifts. This suggests that increasing model size may enhance response quality, making

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406 407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425 426

427 428 429

430

431

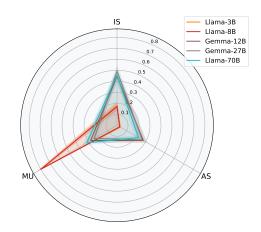


Figure 2: The Average Variance Decomposition Across Models, shaded area is 95% CI

the model more sensitive to prompt text, but at the cost of greater overall susceptibility to surface-level prompt variations, reflecting overfitting rather than an emerging deeper understanding in larger models.

Finally, examining the Intent Comprehension Index, we observe higher values for larger models, but the improvement appears to plateau—the difference between the 70B model and the 12B model is modest. This suggests that intent comprehension is not necessarily an emergent property that scales linearly with model size, and that substantial improvements in understanding user intent can be achieved without requiring the largest available models.

In Figure 4, we examine heterogeneity in model performance across the five topical areas. The results indicate that larger models are not uniformly superior across domains; instead, their sensitivity varies by topic. For example, the largest model, LLaMA 70B, achieves the highest IS scores in Health and Nutrition and Transportation, whereas LLaMA 27B performs best in Logistics, Personal Finance, and Social

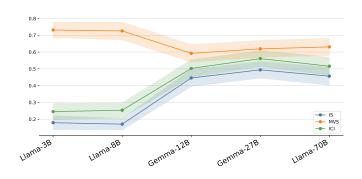


Figure 3: ICI, MVS and IS across models, shaded area is 95% CI

Planning. The figure further highlights differences in the share of AS: questions in Social Planning are consistently more sensitive to writing style across models than those in Health and Nutrition, Transportation, or Personal Finance. Taken together, these findings suggest that model capability is domain-dependent. In this way, intent comprehension is not a uniform property of models but rather varies systematically across topical areas.

#### 5 DISCUSSION AND LIMITATIONS

In this project, we present a formal framework to assess whether LLMs understand the user's latent intent. Instead of focusing solely on whether the model responds correctly to a question, we decompose the response variability of the model into user intent (Intent Sensitivity), irrelevant information (Articulation Sensitivity), and intrinsic randomness (Model Uncertainty). Applying this

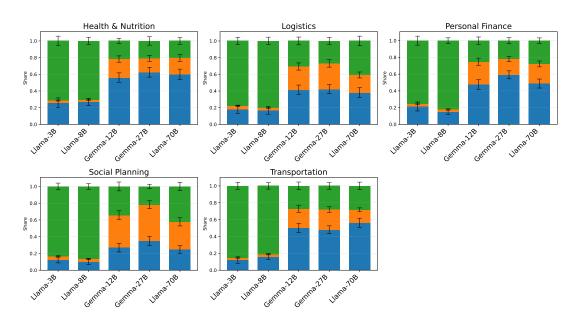


Figure 4: Decomposition across topics, with black bars indicating the 95% confidence intervals.

framework across models of varying sizes and domains, we find that larger models generally exhibit a more robust understanding of users' intent —but an increased number of parameters does not guarantee better intent comprehension. In some cases, we find that smaller models outperform larger ones, and robustness varies across different domains, underscoring the need for evaluation across diverse contexts. Our findings emphasize semantic consistency and generalization as key indicators of model quality, offering a scalable method to distinguish between true understanding and pattern matching.

For practitioners, the IS-AS-MU decomposition offers a comprehensive risk assessment framework suitable for informed decision-making in LLM-enabled products. Accuracy tests alone tell a developer or user whether a model can answer a benchmark question; our metric tells them how the model will behave when real users inevitably rephrase, typo, or embed multiple requests in one prompt. Our metric also informs practitioners on potential fairness issues. Articulation sensitivity often correlates with dialect, accent, or education level; a model whose answers swing with phrasing differences can systematically disadvantage certain groups. Tracking AS gives model producers a concrete, quantitative way to demonstrate equity, complementing broader accuracy and bias checks.

Our measurement approach has several limitations: First, we only consider the first and second moments of response distributions. Although some decision-makers might require higher-order statistics or full distribution analyses, our method strikes a balance between complexity and practicality. Capturing the mean and variance provides essential insights for risk-averse decision-makers and can be easily applied across various models and scenarios. Second, our method is primarily designed for numerical responses. However, some evaluators may be more interested in assessing discrete outputs that lack a natural cardinal or ordinal structure. In the Appendix, we extend our decomposition framework to handle such cases by analyzing the entropy of the discrete response distribution—providing an analogous breakdown to the variance-based approach used in the paper.

Our research highlights the importance of understanding the intent that drives users to utilize LLMs and other AI models. In this study, we show that LLMs often do not fully understand the user intent and respond to differences in superficial phrasing. Future research could more deeply and directly quantify user intent, processes through which intent emerges from preferences or through interaction with generative models, and the inherent limitations of fully expressing complex and original intent. Improved intent identification will enable to better measure models' ability to capture and respond appropriately.

# **Appendix**

# A EQUIVALENT DECOMPOSITION FOR DISCRETE OUTPUTS

In the main text, we use the variance decomposition for our measure

$$\operatorname{Var}(v) = \underbrace{\operatorname{Var}(\operatorname{E}[v \mid I])}_{PS} + \underbrace{\operatorname{E}_I[\operatorname{Var}(\operatorname{E}[v \mid p_I] \mid I)]}_{AS} + \underbrace{\operatorname{E}_I[\operatorname{E}[\operatorname{Var}(v \mid p_I) \mid I]]}_{MU},$$

which is appropriate when v is continuous (or ordinal), as  $\mathrm{Var}(\cdot)$  is then well-defined. For categorical or otherwise non-ordinal outcomes, however, variance cannot be uniquely specified without imposing an arbitrary numerical scale. For example, consider a travel agency that uses a chatbot to recommend destinations based on user preferences. From the agency's perspective, the recommendations correspond to discrete destinations. To handle such cases, we propose an information-theoretic analogue, derived directly by applying the chain rule of entropy twice.

For a discrete random variable X with distribution P(X) we write  $H(X) = -\sum_x P(X = x) \log P(X = x)$  for its Shannon entropy. We denote the mutual information between random variables X and Y as  $I(X;Y) = H(X) - H(X \mid Y) = H(Y) - H(Y \mid X)$ . Let a denote the one-hot (indicator) representation of the discrete outcome v. The variables I (intent) and  $p_I$  (prompt constructed from intent) are defined in the main text.

With these definitions, we obtain the decomposition

$$H(v) \ = \ \underbrace{I\!\!\left(I;v\right)}_{IS} \ + \ \underbrace{\mathrm{E}_i\!\!\left[I\!\!\left(p_I;v\mid I=i\right)\right]}_{AS} \ + \ \underbrace{\mathrm{E}_i\!\!\left[H\!\!\left(v\mid p_I,I=i\right)\right]}_{MU}.$$

where IS (Intent Sensitivity) is the mutual information between the intent I and the model output quantifies how much output reveals about the intended meaning before the prompt is supplied. AS (Articulation Sensitivity) is now the conditional mutual information that measures the extra information carried by the prompt  $p_I$ , given the intent. Finally, MU (Model uncertainty) is the remaining conditional entropy that captures irreducible uncertainty in the model's prediction once both intent and prompt are known.

We can further divide by the total entropy to get a unit-sum normalization that facilitates comparison across models:

$$1 = \frac{I(I;v)}{H(v)} \ + \ \mathrm{E}_i \bigg[ \frac{I(p_I;v \mid I=i)}{H(v)} \bigg] \ + \ \mathrm{E}_i \bigg[ \frac{H(v \mid p_I,I=i)}{H(v)} \bigg] \,.$$

This expression mirrors the continuous-outcome variance decomposition in structure while remaining well-defined for any discrete output space.

#### B RELATED LITERATURE

LLMs Prompt Robustness Recent papers have investigated the sensitivity of large language models (LLMs) to prompting using a variety of frameworks and metrics. Zhuo et al. (2024) introduce the ProSA framework, which quantifies prompt sensitivity through a novel metric called PromptSensiScore (PSS). PSS measures the average variation in performance across semantically equivalent prompt variants at the instance level, offering insight into how much a model's output changes with different prompt formulations. Their findings show that prompt sensitivity varies across tasks, models, and prompt types, with particularly high sensitivity observed in reasoning and creative tasks. Similarly, Sclar et al. (2023) examine the role of "spurious features" in prompt formatting—such as punctuation, spacing, and capitalization—and propose FORMATSPREAD, a metric that captures the spread in task accuracy across equivalent prompt formats. FORMATSPREAD is defined as the difference in performance between the best and worst format, and their results reveal that such superficial changes can lead to performance swings of up to 76 accuracy points—effects that are not mitigated by model size or few-shot learning. Brucks and Toubia (2025) take a different perspective, framing prompt sensitivity as a methodological artifact inspired by the concept of choice architecture.

Using full-factorial experiments, they show that prompt order, labeling, framing, and justification systematically bias model responses—and that even instructing the model to ignore these features does not eliminate the effect. Their central claim is that no prompt is neutral. To address this, they recommend aggregating responses across multiple prompts, akin to ensemble methods, to counteract individual prompt biases.

Our behavioral framing complements and extends this literature by offering a diagnostic framework that quantifies prompt sensitivity in terms of structured variance. Unlike prior work that focuses on scalar measures of how model responses change, we take seriously the idea that model outputs should be stochastic—but emphasize that the source of this randomness should not arise from arbitrary variations in prompt phrasing. As we know, no other measure of prompt sensitivity linked the importance of robustness to the world model and focused on the variance aspect of responses.

Measuring Semantic Robustness and Fairness. Our decomposition also offers practical diagnostic value. High *articulation sensitivity* suggests that small variations in phrasing disproportionately affect model outputs, raising concerns for *fairness* and *accessibility*. Prior work shows that models may perform differently for users with dialectal, non-standard, or less "typical" phrasing (Bolukbasi et al., 2016; Si et al., 2022; Tan and Celis, 2021; Guo et al., 2024). Our framework allows developers to quantify this fragility and track improvements over time. In this respect, our method serves as a form of semantic reliability auditing, applicable in safety-critical domains such as healthcare, finance, and legal reasoning, where output variability under minor phrasing shifts can lead to unacceptable inconsistency.

**Translation Chains** Back-translation is a data augmentation technique in natural language processing, particularly in neural machine translation (Sennrich et al. (2016)). Back-translation involves translating monolingual target-language text into the source language using a reverse translation model, and then using the resulting synthetic parallel data to train the forward model. Subsequent work has extended the approach to include round-trip translation and multilingual back-translation, where intermediate languages are introduced to further diversify the training data and improve generalization (Fadaee et al., 2017; Edunov et al., 2018; Xie et al., 2020; Youn et al., 2016). These variants exploit linguistic variation introduced during translation to create richer training distributions, which have been beneficial not only for translation tasks but also for classification, question answering, and style transfer. In our setup, we do not use translation chains as a data augmentation tool, but instead think of them as a way to diversify language while preserving meaning.

#### B.1 RELATION TO OTHER WORK ON MEASURING WORLD MODELS

As discussed in Remark 1, our definition of Intent Comprehension, is highly related to world models. Our notion of a world model draws inspiration from research in reinforcement learning (RL) and model-based control, where a world model captures the dynamics of an environment and supports planning. In these contexts, a world model is typically a learned transition function or latent representation that abstracts the environment's relevant state space [Ha and Schmidhuber (2018b); Hafner et al. (2019); Ha and Schmidhuber (2018a)].

Analogously, we treat user intent as a latent variable and evaluate whether a model can infer and represent this variable consistently across diverse inputs. This is conceptually related to state abstraction in RL Li et al. (2006), where different observations map to the same underlying state if they yield equivalent value functions. In our case, different prompts map to the same intent if they elicit the same output distribution (or value-weighted distribution).

Recent work has examined whether generative models develop an internal world model in the context of games. Toshniwal et al. (2022) and Li et al. (2023) helped establish games—such as chess and Othello—as testbeds for evaluating the emergence of world models. A common approach in this literature involves using probes to assess whether a model's internal representations encode latent game states. In contrast, our evaluation metrics are model-agnostic: rather than probing representations or relying on an external notion of state feasibility, we focus on the internal consistency of a model's responses as an indicator of world model quality.

Most closely aligned with our work is Vafa et al. (2024), who test an LLM's world model by asking whether it recovers the deterministic finite automaton (DFA) governing a sequence–generation

task. Leveraging the Myhill–Nerode theorem, they introduce two metrics: (i) sequence compression—do any two prefixes that land in the same DFA state admit the same set of valid continuations? and (ii) sequence distinction—do prefixes that reach different states generate different permissible continuations?

Both their framework and ours impose a common behavioral mandate: (1) inputs sharing a latent condition (DFA state or user intent) must elicit indistinguishable outputs, and (2) inputs differing in that condition must yield measurably different outputs. Crucially, both approaches assess world-model quality based solely on outputs, without inspecting internal parameters.

The divide lies in how "correctness" is anchored. Vafa et al. (2024) use a built-in, binary oracle—the known DFA—to label continuations as right or wrong at the syntax level. Our variance-decomposition framework instead delegates that judgment to an external evaluator V, which maps a  $\langle$ prompt, response $\rangle$  pair to latent intent and task value. V can imitate a hard 0/1 oracle, but it can also apply graded semantic scores. Importantly, our test goes further: when V deems two prompts equivalent, we require the entire distribution of (possibly stochastic) responses to match across those prompts. If V itself enforces strict correctness, this collapses to zero model uncertainty—any variability in outputs must stem from differing intents. Thus, while both methods rely on an evaluator, Keyon's oracle is intrinsic and binary, whereas ours is external, tunable, and distribution-aware.

Finally, our notion of a IC demands distributional equality across equivalent prompts, whereas the DFA test assumes deterministic continuations. Extending Myhill—Nerode to stochastic automata would miss the point we target: open-ended tasks with no unique correct answer. In such settings, the likelihoods assigned to alternative continuations encode the model's assumptions. A genuine world model, therefore, aligns those likelihoods whenever the underlying intent is the same, preserving semantic sufficiency even under uncertainty.

### C ESTIMATING VARIANCE DECOMPOSITION VIA ANOVA

Our goal is to decompose variability in numeric responses into three components: PS, AS, and MU. One approach is to compute the sample analogs of the variance decomposition expression in 3, but this approach may induce finite sample bias when cells are sparse or unbalanced. Specifically, estimation of the variance of means may be exaggerated if the means are estimated with noise. To overcome this, we suggest using a hierarchical (mixed-effects) specification that *borrows strength* across related cells via *partial pooling*: group means are shrunk toward a common mean in proportion to their sampling variance and group size. This increases stability, especially when some tasks or prompt variations have few observations.

Index tasks by  $i=1,\ldots,I$ , the prompts variation (nested within task) by  $p=1,\ldots,P_i$ , and individual observations within a prompt by k. We model the standardized response  $y_{ipk}$  as

$$y_{ipk} = \mu + u_i + v_{ip} + \varepsilon_{ipk}, \qquad u_i \sim \mathcal{N}(0, \sigma_I^2), \quad v_{ip} \sim \mathcal{N}(0, \sigma_A^2), \quad \varepsilon_{ipk} \sim \mathcal{N}(0, \sigma_R^2), \quad (1)$$

with  $v_{ip}$  nested in task i. Here  $\mu$  is a global intercept,  $u_i$  captures task-level heterogeneity,  $v_{ip}$  captures prompt-level heterogeneity within tasks, and  $\varepsilon_{ipk}$  is idiosyncratic noise (which we label MU).

Let  $\sigma_T^2 = \sigma_I^2 + \sigma_A^2 + \sigma_R^2$ . We report variance *shares* 

$$\mathrm{TU} = rac{\sigma_I^2}{\sigma_T^2}, \qquad \mathrm{PU} = rac{\sigma_A^2}{\sigma_T^2}, \qquad \mathrm{MU} = rac{\sigma_R^2}{\sigma_T^2},$$

which sum to 1 and are invariant to scale. If we standardize y before fitting, these coincide with absolute contributions on the standardized scale. We estimate  $(\mu, \sigma_I^2, \sigma_A^2, \sigma_R^2)$  by Restricted Maximum Likelihood (REML), which reduces the small-sample bias of variance component MLEs. Practically, we fit (1) via lme4: lmer (through pymer4). Finally, we report uncertainty for the shares using a bootstrap that respects the nesting.

#### D EXPERIMENTAL DETAILS

In this section, we describe how we construct our experiments. We first describe how we construct our set of tasks and intents, and then discuss how we construct an intent-equivariant set of prompts.

Constructing the set of tasks To construct our evaluation metric, we focus on generating open-ended, guesstimation-type questions, using a two-stage LLM workflow. We focus on these types of questions as our IC metric for two reasons. First, our measure is defined over response distributions. To identify whether variation is driven by intent (IS) rather than articulation (AS) or model uncertainty (MU), we need tasks that produce genuine dispersion in plausible answers<sup>3</sup>. Open-ended questions, with various sets of answers, also match real usage, where users ask for estimates, recommendations, and contextual judgments; Together, this yields a naturalistic yet controlled setting in which IS, AS, and MU are identifiable and informative.

We focus on 5 areas of interest: Transportation, Personal Finance, health and nutrition, logistics, and social planning. and use an automatic procedure to generate 24 questions for each topic. To construct our dataset of guesstimation-style prompts, we implemented an automated pipeline that leverages large language models guided by structured templates and semantic constraints. Each generated question contains exactly one {placeholder} token, which is later substituted with concrete values. The system prompt enforces that questions are self-contained, specify explicit reporting units (e.g., "dollars per year," "kWh per month"), and use the placeholder to denote a general semantic role such as a material, process, or population subgroup. To ensure coverage across domains, we draw from a library of question templates spanning categories such as counts, rates, intensities, costs, and probabilities (see Appendix A). For each placeholder, a second prompt generates a set of plausible replacements, or intents, which are short noun phrases consistent with the semantic role and unit of the question. This two-step process—first generating unit-specified question structures and then populating them with realistic values—produces a diverse set of open-ended, quantitative guesstimation questions (see Appendix B for the exact wording of the prompts).

Creating a Set of Intent-Equivalent Prompts Given the set of categories and tasks, our next challenge is to generate a collection of prompts that express the *same intent* but differ in surface form and semantics. This is a non-trivial challenge. Even for humans, writing multiple prompts that convey exactly the same intent without introducing subtle shifts in meaning is difficult. Language models are highly sensitive to linguistic cues, and small differences in phrasing can unintentionally signal different goals or assumptions. This makes it essential to generate paraphrases that are semantically faithful to the original prompt but lexically and syntactically distinct.

To tackle this, we adopt a two-stage approach. Our primary method relies on *cross-lingual translation*, a process inherently designed to preserve meaning while altering surface expression (e.g., Youn et al. (2016)). Translation captures the core idea of transferring intent across languages while abstracting away from specific phrasings. LLMs have demonstrated impressive performance in this task, often on par with human translators in consistency and fluency (Karpinska and Iyyer (2023); Yan et al. (2024a;b)). We leverage this ability by using GPT-4.1 to perform the translations: starting with an original English prompt, we translate it sequentially through two randomly selected intermediate languages and then back into English. The result is a paraphrase that retains the original intent but exhibits different syntactic and semantic features due to translation-induced variation. Different languages vary in how they structure, resulting in varying cross-translations as we vary the source languages into and from which we translate the prompts (Lewis et al. (2023)).

To maximize semantic divergence while preserving intent, we randomly sample intermediate languages from a diverse set.<sup>4</sup> Each translation chain is required to include at least one of Chinese, Japanese, or Arabic, given their significant structural and lexical distance from English (Chiswick and Miller (2005); Lewis et al. (2023)), which helps introduce greater variation in the back-translated output.

To ensure the resulting prompts still reflect the original intent, we also use GPT-4.1 as an additional check to verify that the intent remains unchanged. Given the original and translated prompt, the model assesses whether they express the same underlying request. Only those pairs judged equivalent are retained. To further enrich the diversity of the prompt set, we generate 500 paraphrases for each original question, embed all prompts using sentence embeddings, and apply a greedy selection algorithm to identify the 10 most diverse prompts—those with the highest mutual semantic distance.

<sup>&</sup>lt;sup>3</sup>Single-answer tasks cam make the IC estimator ill-conditioned if the model returns just a single response, because there is no variance to decompose; the shares (and ratios like MVS, ICI) become numerically unstable or undefined, and the standardized scale becomes ill-posed.

<sup>&</sup>lt;sup>4</sup>Chinese, Japanese, Arabic, Korean, Portuguese, Spanish, German, Russian, Italian, French, Hindi

This final step ensures that our prompt set not only shares intent but also spans a wide semantic space, enabling robust evaluation of the model's Articulation Sensitivity.

After generating intent-equivalent prompts, we present them to the language model under evaluation. Each prompt is paired with three input values that vary the prompt's main intent — for example, income level in a tax question. Then, for each of the 24 prompts, at each of the 5 topics, and for each of the 12 values, we prompt the candidate LLM and elicit 25 model responses to capture the within-prompt variation.

As discussed in the definition of a sufficient IC, we need to determine how to evaluate the models' responses. We focus on extracting a clear bottom-line value: if the model provides a range, we take the average. Responses without a definitive answer are discarded, and we continue generating responses until we obtain the desired 25 responses. All outputs are generated with a temperature of 1 to accurately reflect the model's response distribution <sup>5</sup>. To extract a usable numeric value from each response, we employ GPT-4.1-mini as a post-processor that identifies and retrieves the relevant quantity from the model's output.

Finally, in our analysis, we evaluate five models: Meta's LLaMA 3.2 3B Instruct, 3.1 8B Instruct and LLaMA 3.3 70B Instruct, as well as Google's Gemma 3 12B-IT and 27B-IT. We use the OpenRouter API to generate multiple responses from each model.

#### E ADDITIONAL FIGURES

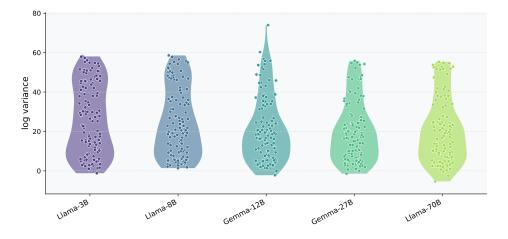


Figure 5: Log variance distribution across the 100 tasks, by model.

<sup>&</sup>lt;sup>5</sup>in the Appendix ?? we show the effect of temperture on the results

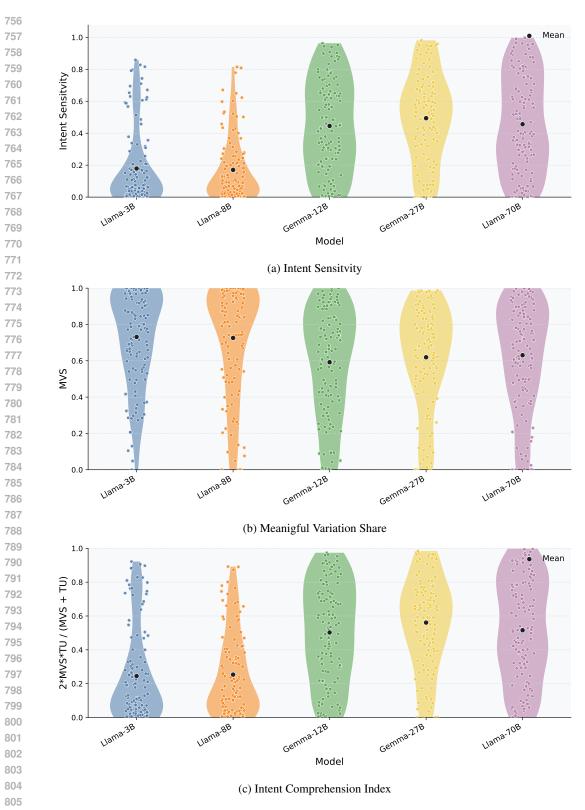


Figure 6: Overall caption for the three stacked figures

# F TEMPERATURE

The main results report model performance under the natural temperature of 1. Figure 2 illustrates this with a spider graph estimated on 5 questions for each of the 5 topics. To explore the effect of sampling temperature, Figures 7 and 8 compare the same tasks estimated with temperature 0.2 versus temperature 1. Reducing the temperature decreases model uncertainty across all models, with the effect particularly pronounced for the smaller models. The figure further suggests that while larger models may overfit to the text, smaller models appear more responsive to meaningful changes—knowing less overall, but perhaps more attuned to user intent.

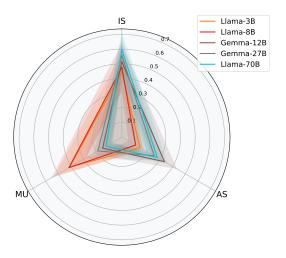


Figure 7: Variance decomposition estimated with temperature 0.2

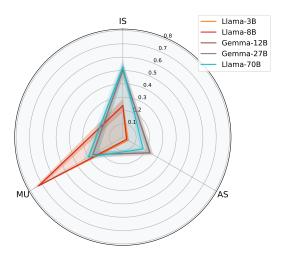


Figure 8: Variance decomposition estimated with temperature 1.

864 865	G	PROMPTS
866	G	.1 Example questions
867	Ů.	LAAMI LE QUESTIONS
868	E	XAMPLE QUESTION TEMPLATES
869		
870 871	Pi	REVALENCE / COUNTS
872	1.	How many {placeholder} are there in a typical city?
873	2.	Approximately how many {placeholder} are produced each farm?
874	3.	What is the total number of {placeholder} registered voters?
875	4.	Roughly how many {placeholder} exist worldwide?
876 877	5.	How many new {placeholder} were added in 2022?
878	٨٠	verages / Means
879		
880		What is the average {placeholder} per high school student?
881 882		What is the mean {placeholder} recorded each year?
883		On average, how much {placeholder} does an individual consume in a day?
884		What is the typical {placeholder} per unit of output?
885	5.	What is the long-run average {placeholder} for a US county?
886 887	R	ATES / SHARES / PERCENTAGES
888	1	By what percentage did {placeholder} change between 2010 and 2020?
889		What fraction of total {scope} is accounted for by {placeholder}?
890		What is the annual growth rate of {placeholder}?
891		What is the aimed growth face of {placeholder}?  What proportion of US population owns at least one {placeholder}?
892 893		What share of all household expenditures goes to {placeholder}?
894	٥.	what share of all household expenditures goes to {placeholder}:
895	To	OTALS / AGGREGATES
896	1.	Estimate the total {placeholder} required over 5 years horizon.
897 898	2.	What is the cumulative {placeholder} expected over the next 5 years?
899	3.	What is the projected lifetime total of {placeholder} for a typical student?
900	4.	What is the overall stock of {placeholder} currently in use?
901	5.	How much {placeholder} will be needed to cover national healthcare system for the next decade?
902 903	L	TENSITIES / PER-UNIT METRICS
904		
905		What is the {placeholder} per unit of output?
906		How much {placeholder} is required per kilometre travelled?
907		What is the {placeholder} per capita in Japan?
908		What is the carbon intensity measured as {placeholder} per kWh?
910	5.	What is the average {placeholder} per square metre?
911	М	AXIMA / MINIMA / RECORDS
912		
913		What is the historical mark (placeholder) and single year?
914 915		What is the leavest recorded (placeholder) observed since 2020?
916		What is the lowest recorded {placeholder}?
917		What is the record-high {placeholder} achieved by a single entity?
	5.	What is the upper bound of {placeholder} under current regulations?

1. What is the average cost of {placeholder} per squre feet?

3. What is the median price paid for {placeholder} in the US?

1. What is the average time needed to complete {placeholder}?

4. What is the expected waiting time until {placeholder} occurs?

5. What is the total expenditure on {placeholder} in 2015?

3. What is the mean lifetime of a {placeholder}?

4. How much investment is required for one unit of {placeholder}?

2. What is the expected budget share spent on {placeholder} each decade?

2. How long does it typically take for {placeholder} to reach completion?

5. What is the typical duration of {placeholder} in the manufacturing sector?

COSTS / MONETARY

**DURATIONS / TIMES** 

918

919

920 921

922

923

924 925

926 927

928

929 930

931

932 933

934

935

936	Densities / Concentrations
937	
938	1. What is the density of {placeholder} per square kilometre?
939 940	2. What is the concentration of {placeholder} per litre in the sample?
941	3. What is the average number of {placeholder} per household?
942	4. What is the typical {placeholder} per lane-kilometre of road?
943 944	5. What is the median {placeholder} per employee in the sector?
945	
946	Probabilistic / Risk
947	1. What is the probability that {placeholder} occurs within a given year?
948 949	2. What is the expected frequency of {placeholder} per decade?
950	3. What is the chance of observing at least one {placeholder} in a week?
951	4. What is the return probability to {placeholder} after 5 years?
952	5. What is the expected failure rate expressed as {placeholder} per 1,000 units?
953 954	
955	RESOURCE / INPUT COEFFICIENTS
956	1. How much {placeholder} is consumed per tonne of output?
957 958	2. What is the marginal {placeholder} required for one additional unit?
959	3. What is the input output coefficient of {placeholder} to gross output?
960	4. What is the elasticity of {placeholder} with respect to price?
961 962	5. What is the shadow cost of one unit of {placeholder}?
963	t and the second
964	VOLUME & CAPACITY
965	1. How many {placeholder} would it take to fill a standard school bus?
966 967	2. What is the total volume of {placeholder} that flows over Niagara Falls in a single day?
968	3. If all the {placeholder} consumed in the United States in a year were put in a single container,
969	how large would it be?
970	4. What is the total annual production of {placeholder} in France, in liters?
971	5. How many bathtubs could you fill with the amount of {placeholder} consumed globally each day?

# 972 WEIGHT & MASS

974 975

978 979

980 981

982

983 984

985

986

987

988 989

990 991

992

993

994 995

996

997

9989991000

1001

1002

1003 1004

1005

1006

1007 1008

1009 1010

1011

1012 1013

1014

1015

1016

1018

1019

1022

1024

1025

- 1. What is the total weight of all the {placeholder} on Earth?
- 2. Estimate the total mass of all {placeholder} currently in the Netherlands.
- 976 977 3. What is the weight of all the {placeholder} produced by New York City each week?
  - 4. What is the total weight of all the {placeholder} in the state of Texas?
    - 5. Estimate the total mass of all {placeholder} currently airborne over the United States.

#### LENGTH, DISTANCE & AREA

- 1. What is the total length, in miles, of all the {placeholder} in Germany?
- 2. If you laid every {placeholder} eaten in America on July 4th end-to-end, how far would the line stretch?
- 3. What is the total surface area of all the {placeholder} in China?
- 4. How many times does an average {placeholder} rotate during its operational lifetime?
- 5. Estimate the total length of all the {placeholder} sold in North America each holiday season.

#### FINANCIAL & ECONOMIC

- 1. How much money, in loose change, is currently in all the {placeholder} in the United States?
- 2. What is the total cost to fuel all the {placeholder} in California for one day?
- 3. What is the total annual revenue of all the {placeholder} operating in the United States?
- 4. How much money is spent on {placeholder} in Canada annually?
  - 5. What is the total market value of all items listed as {placeholder} on eBay worldwide?

#### TIME & DURATION

- 1. How many total hours do all people in the United States spend doing {placeholder} each year?
- 2. How long would it take one person to watch every {placeholder} on YouTube?
  - 3. On average, how many times a day does a person in Japan check their {placeholder}?
  - 4. Estimate the total person-years spent on {placeholder} globally each day.
  - 5. What is the average wait time for a {placeholder} in London during peak hours?

# RATES & FREQUENCY

- 1. Estimate the total number of {placeholder} sent in India every day.
- 2. How many {placeholder} are sold in the United Kingdom each year?
- 3. How many {placeholder} are fixed in Chicago each year?
  - 4. What is the consumption rate of {placeholder}, in units per second, in the United States on a Friday night?
- 5. How many {placeholder} are uploaded to Instagram worldwide every minute?

#### POPULATION & PROFESSION

- 1. How many {placeholder} are there in the state of Illinois?
- 2. Estimate the number of {placeholder} in Brazil.
- 3. What is the total number of {placeholder} on all the people in Japan?
  - 4. Estimate the total number of people currently airborne in {placeholder} around the world.
    - 5. How many {placeholder} work in Paris?

# **EVERYDAY OBJECTS & CONSUMPTION**

- 1. How many pairs of {placeholder} does the average person in America own in their lifetime?
- 2. Estimate the total number of {placeholder} used by all babies in the United States in one year.
- 3. What is the total amount of {placeholder} consumed in the United States annually?
  - 4. How many gallons of {placeholder} are used by the U.S. newspaper industry annually?
  - 5. How many words does the average person read per day on their {placeholder}?

#### Infrastructure & Urban

1026

1027

1028

1031

1032

1033 1034

1035 1036

1037

1038

1040 1041

1042 1043

1044

1045

1046

1048

1049

1050 1051

1052

1053 1054

1055

1057

1058

1059

1060 1061

1062

1063

1064

1065 1066 1067

1068

1070 1071

1072

1074 1075

1076

1077

1078

1079

- 1. What is the total number of {placeholder} in all the buildings of downtown Manhattan?
- 2. How many {placeholder} are there in the entire city of Tokyo?
  - 3. Estimate the total number of {placeholder} in New York's Central Park?
- 4. What is the total number of {placeholder} in the Empire State Building?
- 5. Estimate the total annual electricity consumption of all the {placeholder} in the world.

#### CREATIVE & ABSTRACT

- 1. How many individual {placeholder} are on a professional soccer field?
- 2. How many {placeholder} could you fit in Grand Central Terminal's main concourse?
- 3. What is the total number of {placeholder} manufactured globally in a single day?
  - 4. What is the total population of {placeholder} in Venice?
  - 5. How many {placeholder} would it take to build a chain to the Moon?

#### PROBABILISTIC & ODDS

- 1. What is the probability that a randomly selected person in the United States has {placeholder}?
- 2. What are the odds of a flight being delayed at {placeholder} International Airport?
- 3. What is the daily probability that the {placeholder} in a major city experiences a major outage?
- 4. What is the chance that a new {placeholder} in the United States fails within its first year?
- 5. If you pick a random word from the New York Times, what is the probability it is the word '{placeholder}'?
- 6. What is the likelihood of experiencing {placeholder} in London on a given day in July?
- 7. What is the probability that a car driving one mile on a US highway will get a {placeholder}?
- 8. Estimate the chance that a randomly selected email in an average inbox is {placeholder}.
- 9. What is the annual probability of a {placeholder} causing significant damage in California?
- 10. What are the odds that a randomly chosen {placeholder} from a large supermarket is expired?

## H PROMPT SPECIFICATIONS AND DATA COLLECTION PROCEDURE

This section describes all prompts and control settings used in data collection.

#### H.1 System Prompt Used for Answer Generation

For every question posed to a model (the *estimation task*), we attach the following system message and then a user message containing the question text:

```
You are a helpful assistant designed to answer users questions that involve estimating real-world quantities. When asked for a numerical value (e.g., average, frequency, duration, or count), always provide your best-guess estimate, even if you lack exact data. Avoid generic refusals like "I don't have that information." If needed,
```

rely on general knowledge, plausible assumptions, or known ranges from similar contexts. The goal is to be useful by offering a grounded and reasoned numerical estimate. Return a single numeric value in your response, followed by the appropriate unit of measurement (e.g., 3 days, 3 kg, 10 mg).

#### H.2 BACK-TRANSLATION PROMPTING AND DIVERSITY PROCEDURE

We generate alternative wordings of a seed question via back-translation. Each back-translation hop uses the same system translator instruction:

```
You are a professional translator. Translate the user text to {target_lang}, preserving {placeholder}.

Preserve meaning exactly. Do not add, remove facts and information.

Keep sentence boundaries and speaker perspective the same.

Return *only* the translation--no commentary.
```

The corresponding user content is the current question snippet (potentially containing a literal  $\{placeholder\}$  token). We first translate English  $\rightarrow$  intermediate language(s), then back to English using the same instruction with  $\{target\_lang\} = English$ . Placeholders enclosed in curly braces are preserved verbatim through all hops.

**Language chain.** We use a chain of H languages (default H=2) selected to encourage script and typological variation. One of {Chinese, Japanese, Arabic} is always included; the remaining hop(s) are sampled from: { Chinese, Japanese, Arabic, Korean, Portuguese, Spanish, German, Russian, Italian, French, Hindi}. (Exact draws are randomized per back–translation.)

EQUIVALENCE FILTER.

```
system: You are an expert at judging whether two English
sentences mean exactly the same thing.
user: Do these two sentences convey the same meaning and
intent?
Sentence A: "..."
Sentence B: "..."
Respond only with "Yes" or "No".
```

When a concrete value for {placeholder} is available, the comparison is done after substituting that value in both sentences (the raw candidate must still contain the literal {placeholder} token).

SPEAKER/PERSPECTIVE FILTER.

```
system: You are an expert in grammar and meaning. Your
job is to assess whether two English sentences use the same
grammatical subject or speaker perspective. That includes
whether both use "I", "we", "you", passive voice, or the
same named subject (e.g., "the city", "Jiji").
user: Do these two sentences use the same speaker or
subject perspective?
Sentence A: "..."
Sentence B: "..."
Respond only with "Yes" or "No". If the speaker changes or
a name is added or removed, respond "No".
```

H.2.1 LIST OF QUESTION PROMPTS AND EXTRACTIONS

#### **HEALTH & NUTRITION**

**Question:** On average, what is the energy expenditure from {placeholder} for an adult during one hour, in kilocalories per hour?

Values: ["walking", "running", "cycling", "swimming", "yoga", "dancing", "gardening", "standing", "typing", "reading", "cleaning", "driving"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing kilocalories per hour. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average amount of dietary fiber contained in one serving of {placeholder}, in grams per serving.

**Values:** ["black beans", "oatmeal", "whole wheat bread", "broccoli", "chia seeds", "raspberries", "almonds", "sweet potato", "green peas", "avocado", "brown rice", "carrots"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing grams per serving. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average number of times a typical household uses a {placeholder} per week, in uses per week.

**Values:** ["microwave", "dishwasher", "vacuum cleaner", "washing machine", "television", "oven", "toaster", "coffee maker", "blender", "refrigerator door", "air conditioner", "hair dryer"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the number of uses per week; unit: 'uses per week'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average annual revenue generated by {placeholder} in the health and nutrition industry in the United States, in dollars per year.

**Values:** ["dietary supplements", "organic food products", "vitamin sales", "weight loss programs", "nutrition consulting services", "sports nutrition drinks", "meal replacement shakes", "functional foods", "natural health stores", "health coaching services", "nutritional app subscriptions", "protein powder brands"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average body mass index (BMI) of {placeholder} in the United States, in kilograms per square meter ( $kg/m\neg$ ).

Values: ["adults", "children", "teenagers", "elderly men", "pregnant women", "college students", "African Americans", "Asian Americans", "male athletes", "female nurses", "Hispanic women", "middle-aged adults"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is kilograms per square meter  $(kg/m\neg)$ . Return "None" if there is no numerical value.

Answer text: {answer\_text}

1188 1189 1190 Question: Estimate the average shelf life of {placeholder} in a typical American grocery store, in 1191 days per item. Values: ["milk", "eggs", "lettuce", "yogurt", "chicken breast", "ground beef", "bread loaf", "apples", 1192 "bananas", "cheese block", "carrots", "spinach"], You extract a single numeric value from an answer 1193 string. Return only the number and ensure it matches the requested unit. 1194 Extraction Prompt: Extract a single numeric value representing average shelf life; unit must be 1195 'days per item'. Return "None" if there is no numerical value. 1196 1197 Answer text: {answer\_text} 1198 1199 **Question:** Estimate the average amount of vitamin C contained in one medium-sized {placeholder}, 1200 in milligrams per item. 1201 Values: ["orange", "kiwi", "strawberry", "broccoli floret", "bell pepper", "tomato", "grapefruit", 1202 "mango", "papaya", "brussels sprout", "cabbage leaf", "pineapple slice"], You extract a single numeric 1203 value from an answer string. Return only the number and ensure it matches the requested unit. 1204 **Extraction Prompt:** Extract a single numeric value and report it in milligrams per item. Return 1205 "None" if there is no numerical value. 1206 Answer text: {answer\_text} 1207 1208 1209 1210 tons per year. 1211 1212 only the number and ensure it matches the requested unit. 1213 1214 'metric tons per year'. Return "None" if there is no numerical value. 1215

Question: Estimate the average annual production of {placeholder} in the United States, in metric

Values: ["corn", "soybeans", "wheat", "cotton", "rice", "sugar beets", "potatoes", "tomatoes", "coal", "steel", "aluminum", "cement"], You extract a single numeric value from an answer string. Return

**Extraction Prompt:** Extract a single numeric value representing annual production. Only report in

Answer text: {answer\_text}

**Question:** Estimate the average maintenance cost of a {placeholder} used in a typical American hospital, in dollars per device per year.

Values: ["MRI machine", "X-ray machine", "ultrasound scanner", "ventilator", "defibrillator", "ECG monitor", "infusion pump", "anesthesia machine", "patient monitor", "CT scanner", "sterilizer", "dialysis machine"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it using the unit 'dollars per device per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average failure rate for {placeholder} in clinical nutrition settings, in percent

Values: ["infusion pumps", "enteral feeding tubes", "peripheral IV catheters", "central venous catheters", "parenteral nutrition bags", "glucometers", "electronic medication carts", "feeding pumps", "blood glucose monitors", "nutrition software systems", "IV fluid warmers", "nutritional supplement dispensers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in percent (%) per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

1239 1240 1241

1216

1217 1218

1219

1220

1221

1222

1223

1224

1225

1226

1227 1228 1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

Question: Estimate the average annual revenue generated by {placeholder} per grocery store in the United States, in dollars per year.

Values: ["fresh produce sales", "dairy products", "bakery items", "meat department", "seafood sales", "beverage section", "frozen foods", "prepared meals", "snack foods", "organic products", "household goods", "health and beauty aids"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value expressed in 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average annual treatment cost, in dollars per year, for a patient with {placeholder}.

**Values:** ["diabetes", "hypertension", "asthma", "rheumatoid arthritis", "multiple sclerosis", "chronic kidney disease", "heart failure", "COPD", "psoriasis", "Parkinson's disease", "HIV infection", "breast cancer"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing cost and ensure the unit is 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the estimated annual economic cost of {placeholder} in the United States, measured in dollars per year?

**Values:** ["Diet-related chronic diseases", "foodborne illnesses", "drug abuse", "workplace injuries", "medical errors", "Childhood obesity", "Adult obesity", "Micronutrient deficiencies", "chronic diseases", "environmental pollution", "Type 2 diabetes", "High sodium consumption"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the estimated average indoor concentration of {placeholder} in residential homes, measured in micrograms per cubic meter?

**Values:** ["particulate matter", "formaldehyde", "benzene", "toluene", "nitrogen dioxide", "carbon monoxide", "ammonia", "radon", "ozone", "acetaldehyde", "volatile organic compounds", "mold spores"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value. The unit must be micrograms per cubic meter. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of meals prepared using {placeholder} per household per month? (unit: meals per month)

**Values:** ["weekdays", "weekends", "holidays", "summer", "winter", "schooldays", "vacation", "festivals", "spring", "autumn", "busy days", "quiet days"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average number of meals prepared using the specified time window/type, reported in meals per month. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the estimated annual revenue generated from {placeholder} in the United States, measured in dollars per year?

Values: ["online advertising", "pharmaceutical sales", "movie ticket sales", "automobile manu-

facturing", "telecommunications services", "video game industry", "fast food chains", "streaming subscriptions", "health insurance premiums", "retail e-commerce", "professional sports leagues", "music industry"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the annual revenue from {placeholder} in dollars per year. Only report a number and 'dollars per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What fraction of total public health spending is allocated to {placeholder}, measured in percent (%)?

**Values:** ["mental health services", "immunization programs", "maternal care", "chronic disease management", "substance abuse prevention", "emergency preparedness", "HIV/AIDS treatment", "elderly care", "primary care initiatives", "health education campaigns", "rural health outreach", "tuberculosis control"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the fraction as a percentage (%). Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the average annual interest paid on {placeholder} by U.S. households, measured in dollars per year?

Values: ["credit cards", "mortgages", "auto loans", "student loans", "personal loans", "home equity lines", "payday loans", "installment loans", "private loans", "business loans", "medical debt", "store credit accounts"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value that answers the question, reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average daily calorie intake for {placeholder} in kilocalories per day?

**Values:** ["adult men", "adult women", "teenagers", "infants", "preschool children", "pregnant women", "lactating mothers", "elderly adults", "athletes", "office workers", "manual laborers", "vegetarians"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing daily calorie intake and report the unit as 'kilocalories per day'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the annual growth rate of cost of {placeholder} in the healthcare sector, measured in percent (%) per year?

**Values:** ["prescription drugs", "medical devices", "hospital services", "physician fees", "nursing care", "diagnostic tests", "surgical procedures", "insurance premiums", "emergency care", "laboratory services", "imaging services", "rehabilitation therapy"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual growth rate, reported in percent (%) per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the number of {placeholder} operating in New York City on an average weekday, measured in vehicles per day.

Values: ["taxis", "buses", "delivery trucks", "rideshare cars", "garbage trucks", "ambulances", "fire

1350 engines", "limousines", "school buses", "police cars", "construction vehicles", "motorcycles"], You 1351 extract a single numeric value from an answer string. Return only the number and ensure it matches 1352 the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it using the unit: vehicles per day. 1354 Return "None" if there is no numerical value.

Answer text: {answer\_text}

1356 1357 1358

1359

1360

1361

1362

1363

1364

1365

1366

1353

1355

Question: What is the median distance traveled annually by a {placeholder} in the United States, measured in kilometers per year?

Values: ["passenger car", "pickup truck", "SUV", "motorcycle", "school bus", "delivery van", "taxicab", "minivan", "city bus", "semi-truck", "ambulance", "fire engine"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

Extraction Prompt: Extract a single numeric value representing the median annual distance traveled by a {placeholder}. The allowed unit is kilometers per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

1367 1368 1369

1370

1371

1372

1373

1374

1375

1376 1377

**Question:** What is the average distance traveled per day by a {placeholder} in urban areas, measured in kilometers per day?

Values: ["taxi", "bus", "bicycle", "electric scooter", "motorcycle", "delivery van", "ride-share car", "ambulance", "garbage truck", "fire engine", "postal vehicle", "private car"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested

Extraction Prompt: Extract a single numeric value for average daily distance and use 'kilometers per day' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

1378 1379 1380

1381

1383

1384

1385

1386

1387

Question: Estimate the total mass of all {placeholder} currently stored in U.S. hospitals, measured in kilograms.

Values: ["MRI machines", "CT scanners", "ultrasound devices", "X-ray tubes", "ventilators", "infusion pumps", "dialysis machines", "defibrillators", "anesthesia workstations", "surgical robots", "incubators", "ECG monitors"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in kilograms. Return "None" if there is no numerical value.

Answer text: {answer\_text}

1388 1389 1390

1391

1392

1393

1394

1395

1396

Question: How many servings of {placeholder} are typically consumed by an adult in the United States per week, measured in servings per week?

Values: ["breakfast", "lunch", "dinner", "snack", "dessert", "vegetables", "fruits", "meat", "fish", "dairy products", "grains", "salads"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and specify the unit as 'servings per week'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

1398 1399 1400

**Question:** On average, how many minutes per week does a person spend on {placeholder}? **Values:** ["exercise", "reading", "cooking", "commuting", "watching television", "cleaning", "working", "shopping", "studying", "socializing", "gardening", "sleeping"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

1402 1403

1401

**Extraction Prompt:** Extract a single numeric value for time spent per week. Allowed unit: minutes per week. Return "None" if there is no numerical value.

Answer text: {answer\_text}

#### LOGISTICS

 **Question:** What is the estimated number of {placeholder} employed in warehouse logistics worldwide, reported as individuals?

**Values:** ["forklift operators", "inventory managers", "shipping coordinators", "order pickers", "warehouse supervisors", "logistics analysts", "packers", "receiving clerks", "distribution managers", "material handlers", "customs brokers", "freight forwarders"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing individuals. The allowed unit is 'individuals'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual number of {placeholder} incidents reported in major logistics networks worldwide, measured in cases per year.

**Values:** ["cargo theft", "lost shipment", "delayed delivery", "damaged goods", "customs violation", "fraudulent invoice", "piracy attack", "cyber breach", "hazardous spill", "stolen container", "misrouted package", "supply chain disruption"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated total annual number of {placeholder} incidents, reported in cases per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the average financial loss caused by {placeholder} in global supply chains per year, reported in US dollars per year?

**Values:** ["cyberattacks", "natural disasters", "port congestion", "trade wars", "labor strikes", "piracy", "regulatory changes", "pandemics", "supplier insolvency", "transportation delays", "counterfeit goods", "customs bottlenecks"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure it is reported in US dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the historical peak number of {placeholder} operating in global logistics, measured in vehicles per year?

**Values:** ["container ships", "cargo planes", "delivery trucks", "freight trains", "oil tankers", "bulk carriers", "vans", "electric trucks", "autonomous vehicles", "motorcycles", "reefer trucks", "intermodal trailers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value for the a quantity in the response. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** How much revenue is generated from {placeholder} in the logistics industry worldwide each year, measured in US dollars per year?

**Values:** ["freight forwarding", "warehousing", "customs brokerage", "last-mile delivery", "cold chain logistics", "e-commerce fulfillment", "express shipping", "third-party logistics services", "reverse logistics", "container leasing", "transportation management systems", "supply chain consulting"],

You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual revenue. The allowed unit is US dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the probability that a shipping container in transit will experience {placeholder}, measured in percent per shipment?

**Values:** ["mold growth", "water damage", "infestation", "corrosion", "contamination", "spoilage", "condensation", "bacterial infection", "fungal contamination", "rusting", "odorous emission", "rot"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the probability, reported in percent per shipment. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the consumption rate of {placeholder} in a major logistics hub, measured in units per hour?

**Values:** ["pallets", "shipping containers", "fuel drums", "packaging materials", "barcode labels", "forklift batteries", "loading crates", "delivery vans", "conveyor belts", "storage bins", "handheld scanners", "sorting trays"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the consumption rate and specify the unit as 'units per hour'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the number of {placeholder} operating in Japan at any given moment, measured in vehicles.

**Values:** ["taxis", "buses", "trains", "ambulances", "fire trucks", "police cars", "delivery vans", "motorcycles", "rental cars", "garbage trucks", "private cars", "ride-sharing vehicles"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in vehicles. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** Estimate the total energy consumption attributed to {placeholder} in large logistics centers worldwide each year, measured in megawatt-hours per year.

**Values:** ["lighting", "heating", "cooling", "ventilation", "material handling equipment", "refrigeration", "security systems", "conveyor belts", "automated sorting systems", "charging electric vehicles", "water heating", "data processing centers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual energy consumption, using megawatt-hours per year as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the probability that {placeholder} will experience a major logistics disruption in a given year, measured in percent per year?

**Values:** ["regional warehouse", "supply chain", "distribution center", "retail outlet", "manufacturing facility", "logistics network", "shipping hub", "inventory system", "transport fleet", "customs terminal", "port authority", "fulfillment center"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing probability. The allowed unit is percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average interest rate charged for {placeholder} used by logistics companies in percent per year.

**Values:** ["working capital loans", "equipment leases", "revolving credit lines", "invoice factoring", "asset-based loans", "vehicle financing", "commercial mortgages", "bridge loans", "trade credit facilities", "term loans", "letters of credit", "fleet leasing agreements"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average interest rate, and specify the unit as percent per year (%/year). Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual salary earned by {placeholder} working in the logistics industry, measured in US dollars per year?

**Values:** ["warehouse manager", "forklift operator", "supply chain analyst", "logistics coordinator", "inventory specialist", "transportation manager", "customs broker", "shipping clerk", "freight dispatcher", "delivery driver", "operations supervisor", "procurement officer"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average annual salary, reported in US dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** How much maintenance cost for {placeholder} is incurred per kilometer traveled by a delivery truck, measured in US dollars per kilometer?

**Values:** ["engine repairs", "tire replacement", "oil changes", "brake servicing", "transmission maintenance", "suspension repairs", "coolant system upkeep", "battery replacement", "exhaust system repair", "air conditioning maintenance", "electrical system servicing", "fuel system cleaning"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in US dollars per kilometer. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the total annual value of {placeholder} issued to logistics companies worldwide, measured in US dollars per year?

**Values:** ["trade finance", "invoice factoring", "letters of credit", "equipment leases", "supply chain loans", "working capital loans", "commercial paper", "asset-backed securities", "revolving credit facilities", "warehouse receipts financing", "export credits", "fleet insurance policies"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in 'US dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the annual growth rate of revenue from {placeholder} in the logistics sector, measured in percent per year?

**Values:** ["freight forwarding", "warehousing services", "last-mile delivery", "customs brokerage", "cold chain logistics", "express shipping", "reverse logistics", "e-commerce fulfillment", "intermodal transport", "fleet management", "supply chain consulting", "contract logistics"], You extract a single

numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual growth rate. The allowed unit is percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of gallons of {placeholder} consumed by the global logistics industry each year.

**Values:** ["diesel", "gasoline", "jet fuel", "marine fuel", "biodiesel", "hydrogen", "liquefied natural gas", "ethanol blend", "synthetic fuel", "renewable diesel", "compressed natural gas", "aviation fuel"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is gallons per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Approximately how many {placeholder} are loaded onto a cargo ship per voyage? (units: items per voyage)

**Values:** ["containers", "cranes", "forklifts", "pallets", "vehicles", "generators", "refrigerators", "tractors", "bulldozers", "excavators", "computers", "machinery"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the number of items loaded per voyage, using 'items per voyage' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of crates of {placeholder} delivered to supermarkets in New York City per day? (units: crates per day)

**Values:** ["apples", "oranges", "bananas", "lettuce", "tomatoes", "potatoes", "carrots", "onions", "grapes", "spinach", "broccoli", "peppers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is 'crates per day'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** Estimate the total weight of {placeholder} transported by trucks across Europe in metric tons per year.

**Values:** ["construction materials", "fresh produce", "electronic goods", "automobiles", "industrial machinery", "textiles", "petroleum products", "furniture", "pharmaceuticals", "beverages", "household appliances", "steel"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total weight, and ensure the unit is metric tons per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

Question: How long does it typically take for {placeholder} to reach completion in days per shipment?

**Values:** ["customs clearance", "order processing", "quality inspection", "inventory restocking", "freight consolidation", "packaging preparation", "route planning", "document verification", "payment confirmation", "load scheduling", "cargo unloading", "delivery coordination"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

 **Extraction Prompt:** Extract a single numeric value representing the typical completion time for {placeholder}, reported in days per shipment. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What percentage of warehouse workers experience {placeholder} each year due to repetitive lifting? (units: percent per year)

**Values:** ["lower back pain", "shoulder strain", "tendinitis", "herniated discs", "carpal tunnel syndrome", "muscle fatigue", "sprains", "rotator cuff injuries", "joint inflammation", "elbow pain", "chronic soreness", "ligament tears"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value. Unit must be 'percent per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average amount of {placeholder} transported by a single refrigerated truck in kilograms per trip?

**Values:** ["beef", "chicken", "lettuce", "milk", "cheese", "yogurt", "apples", "broccoli", "carrots", "ice cream", "fish fillets", "tomatoes"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract one numeric value representing the average quantity, using kilograms per trip as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of {placeholder} handled by a medium-sized logistics company in one month (units: shipments per month).

**Values:** ["overnight shipments", "express packages", "international deliveries", "standard parcels", "bulk consignments", "fragile items", "return shipments", "same-day deliveries", "temperature-controlled goods", "e-commerce orders", "seasonal shipments", "high-value packages"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated number of shipments handled per month. The allowed unit is 'shipments per month'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the annual turnover rate for {placeholder} working in warehouse logistics, measured in percent per year?

**Values:** ["female employees", "male employees", "temporary staff", "full-time workers", "part-time workers", "seasonal workers", "shift supervisors", "older workers", "younger employees", "new hires", "contractors", "management staff"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual turnover rate, reported in percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** If all the {placeholder} incidents reported by logistics companies in one year were placed into a single file, how many pages would it contain? (units: pages per year)

**Values:** ["vehicle breakdown", "missed delivery", "cargo theft", "shipment delay", "lost package", "inventory discrepancy", "equipment malfunction", "damaged goods", "routing error", "documentation error", "fuel shortage", "container misplacement"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

 **Extraction Prompt:** Extract the estimated total number of pages and report the value in 'pages per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of cases of {placeholder} reported among long-haul truck drivers in the United States per year (units: cases per year).

**Values:** ["sleep apnea", "hypertension", "diabetes", "depression", "obesity", "back pain", "lung cancer", "hepatitis C", "skin infections", "substance abuse", "cardiovascular disease", "chronic fatigue"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the total annual cases, using 'cases per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total market value of all {placeholder} currently stored in U.S. warehouses, measured in US dollars.

**Values:** ["soybeans", "automobiles", "furniture", "pharmaceuticals", "electronics", "apparel", "petroleum", "coffee beans", "lumber", "steel coils", "corn", "copper wire"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated total market value. The allowed unit is US dollars. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the historical peak number of {placeholder} operating in global logistics, measured in vehicles per year?

**Values:** ["container ships", "cargo planes", "delivery trucks", "freight trains", "oil tankers", "bulk carriers", "vans", "electric trucks", "autonomous vehicles", "motorcycles", "reefer trucks", "intermodal trailers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value followed by 'vehicles per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** How much revenue is generated from {placeholder} in the logistics industry worldwide each year, measured in US dollars per year?

**Values:** ["freight forwarding", "warehousing", "customs brokerage", "last-mile delivery", "cold chain logistics", "e-commerce fulfillment", "express shipping", "third-party logistics services", "reverse logistics", "container leasing", "transportation management systems", "supply chain consulting"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual revenue. The allowed unit is US dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the probability that a shipping container in transit will experience {placeholder}, measured in percent per shipment?

**Values:** ["mold growth", "water damage", "infestation", "corrosion", "contamination", "spoilage", "condensation", "bacterial infection", "fungal contamination", "rusting", "odorous emission", "rot"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

Extraction Prompt: Extract a single numeric value representing the probability, reported in percent per shipment. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the consumption rate of {placeholder} in a major logistics hub, measured in units per hour?

**Values:** ["pallets", "shipping containers", "fuel drums", "packaging materials", "barcode labels", "forklift batteries", "loading crates", "delivery vans", "conveyor belts", "storage bins", "handheld scanners", "sorting trays"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the consumption rate and specify the unit as 'units per hour'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the number of {placeholder} operating in Japan at any given moment, measured in vehicles.

**Values:** ["taxis", "buses", "trains", "ambulances", "fire trucks", "police cars", "delivery vans", "motorcycles", "rental cars", "garbage trucks", "private cars", "ride-sharing vehicles"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in vehicles. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total energy consumption attributed to {placeholder} in large logistics centers worldwide each year, measured in megawatt-hours per year.

**Values:** ["lighting", "heating", "cooling", "ventilation", "material handling equipment", "refrigeration", "security systems", "conveyor belts", "automated sorting systems", "charging electric vehicles", "water heating", "data processing centers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual energy consumption, using megawatt-hours per year as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the probability that {placeholder} will experience a major logistics disruption in a given year, measured in percent per year?

**Values:** ["regional warehouse", "supply chain", "distribution center", "retail outlet", "manufacturing facility", "logistics network", "shipping hub", "inventory system", "transport fleet", "customs terminal", "port authority", "fulfillment center"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing probability. The allowed unit is percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average interest rate charged for {placeholder} used by logistics companies in percent per year.

**Values:** ["working capital loans", "equipment leases", "revolving credit lines", "invoice factoring", "asset-based loans", "vehicle financing", "commercial mortgages", "bridge loans", "trade credit facilities", "term loans", "letters of credit", "fleet leasing agreements"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average interest rate, and specify the unit as percent per year (%/year). Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual salary earned by {placeholder} working in the logistics industry, measured in US dollars per year?

**Values:** ["warehouse manager", "forklift operator", "supply chain analyst", "logistics coordinator", "inventory specialist", "transportation manager", "customs broker", "shipping clerk", "freight dispatcher", "delivery driver", "operations supervisor", "procurement officer"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average annual salary, reported in US dollars per year. Return "None" if there is no numerical value.

Answer text: {answer text}

 **Question:** How much maintenance cost for {placeholder} is incurred per kilometer traveled by a delivery truck, measured in US dollars per kilometer?

**Values:** ["engine repairs", "tire replacement", "oil changes", "brake servicing", "transmission maintenance", "suspension repairs", "coolant system upkeep", "battery replacement", "exhaust system repair", "air conditioning maintenance", "electrical system servicing", "fuel system cleaning"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in US dollars per kilometer. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the total annual value of {placeholder} issued to logistics companies worldwide, measured in US dollars per year?

**Values:** ["trade finance", "invoice factoring", "letters of credit", "equipment leases", "supply chain loans", "working capital loans", "commercial paper", "asset-backed securities", "revolving credit facilities", "warehouse receipts financing", "export credits", "fleet insurance policies"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in 'US dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the annual growth rate of revenue from {placeholder} in the logistics sector, measured in percent per year?

**Values:** ["freight forwarding", "warehousing services", "last-mile delivery", "customs brokerage", "cold chain logistics", "express shipping", "reverse logistics", "e-commerce fulfillment", "intermodal transport", "fleet management", "supply chain consulting", "contract logistics"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual growth rate. The allowed unit is percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of gallons of {placeholder} consumed by the global logistics industry each year.

**Values:** ["diesel", "gasoline", "jet fuel", "marine fuel", "biodiesel", "hydrogen", "liquefied natural gas", "ethanol blend", "synthetic fuel", "renewable diesel", "compressed natural gas", "aviation fuel"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is gallons per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Approximately how many {placeholder} are loaded onto a cargo ship per voyage? (units: items per voyage)

**Values:** ["containers", "cranes", "forklifts", "pallets", "vehicles", "generators", "refrigerators", "tractors", "bulldozers", "excavators", "computers", "machinery"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the number of items loaded per voyage, using 'items per voyage' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of crates of {placeholder} delivered to supermarkets in New York City per day? (units: crates per day)

**Values:** ["apples", "oranges", "bananas", "lettuce", "tomatoes", "potatoes", "carrots", "onions", "grapes", "spinach", "broccoli", "peppers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is 'crates per day'. Return "None" if there is no numerical value.

1858 Answer text: {answer\_text}

 **Question:** Estimate the total weight of {placeholder} transported by trucks across Europe in metric tons per year.

**Values:** ["construction materials", "fresh produce", "electronic goods", "automobiles", "industrial machinery", "textiles", "petroleum products", "furniture", "pharmaceuticals", "beverages", "household appliances", "steel"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total weight, and ensure the unit is metric tons per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

Question: How long does it typically take for {placeholder} to reach completion in days per shipment?

**Values:** ["customs clearance", "order processing", "quality inspection", "inventory restocking", "freight consolidation", "packaging preparation", "route planning", "document verification", "payment confirmation", "load scheduling", "cargo unloading", "delivery coordination"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit

**Extraction Prompt:** Extract a single numeric value representing the typical completion time for {placeholder}, reported in days per shipment. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What percentage of warehouse workers experience {placeholder} each year due to repetitive lifting? (units: percent per year)

**Values:** ["lower back pain", "shoulder strain", "tendinitis", "herniated discs", "carpal tunnel syndrome", "muscle fatigue", "sprains", "rotator cuff injuries", "joint inflammation", "elbow pain", "chronic soreness", "ligament tears"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value. Unit must be 'percent per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average amount of {placeholder} transported by a single refrigerated truck in kilograms per trip?

**Values:** ["beef", "chicken", "lettuce", "milk", "cheese", "yogurt", "apples", "broccoli", "carrots", "ice cream", "fish fillets", "tomatoes"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract one numeric value representing the average quantity, using kilograms per trip as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** Estimate the total number of {placeholder} handled by a medium-sized logistics company in one month (units: shipments per month).

**Values:** ["overnight shipments", "express packages", "international deliveries", "standard parcels", "bulk consignments", "fragile items", "return shipments", "same-day deliveries", "temperature-controlled goods", "e-commerce orders", "seasonal shipments", "high-value packages"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated number of shipments handled per month. The allowed unit is 'shipments per month'. Return "None" if there is no numerical value.

1912 Answer text: {answer\_text}

**Question:** What is the annual turnover rate for {placeholder} working in warehouse logistics, measured in percent per year?

**Values:** ["female employees", "male employees", "temporary staff", "full-time workers", "part-time workers", "seasonal workers", "shift supervisors", "older workers", "younger employees", "new hires", "contractors", "management staff"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual turnover rate, reported in percent per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** If all the {placeholder} incidents reported by logistics companies in one year were placed into a single file, how many pages would it contain? (units: pages per year)

**Values:** ["vehicle breakdown", "missed delivery", "cargo theft", "shipment delay", "lost package", "inventory discrepancy", "equipment malfunction", "damaged goods", "routing error", "documentation error", "fuel shortage", "container misplacement"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract the estimated total number of pages and report the value in 'pages per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** Estimate the total number of cases of {placeholder} reported among long-haul truck drivers in the United States per year (units: cases per year).

**Values:** ["sleep apnea", "hypertension", "diabetes", "depression", "obesity", "back pain", "lung cancer", "hepatitis C", "skin infections", "substance abuse", "cardiovascular disease", "chronic fatigue"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the total annual cases, using 'cases per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total market value of all {placeholder} currently stored in U.S. warehouses, measured in US dollars.

**Values:** ["soybeans", "automobiles", "furniture", "pharmaceuticals", "electronics", "apparel", "petroleum", "coffee beans", "lumber", "steel coils", "corn", "copper wire"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated total market value. The allowed unit is US dollars. Return "None" if there is no numerical value.

Answer text: {answer\_text}

## PERSONAL FINANCE

**Question:** On average, how many grams of {placeholder} does a U.S. adult consume per day? **Values:** ["protein", "fiber", "sugar", "fat", "carbohydrate", "sodium", "cholesterol", "calcium", "potassium", "magnesium", "iron", "vitamin C"], You extract a single numeric value from an answer

string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the amount consumed per day, reported in grams. Return "None" if there is no numerical value.

1966 Answer text: {answer\_text}

**Question:** What is the average monthly electricity usage attributed specifically to {placeholder} in a typical U.S. household (kWh per month)?

**Values:** ["refrigerator", "air conditioning", "lighting", "water heater", "clothes dryer", "dishwasher", "television", "microwave oven", "freezer", "space heating", "computer equipment", "washing machine"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and the exact allowed unit: kWh per month. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** On average, how much do households in the United States spend on transportation for {placeholder} per year (dollars per year)?

**Values:** ["retirees", "single adults", "families with children", "urban residents", "rural households", "college students", "low-income families", "high-income households", "immigrants", "senior citizens", "military families", "recent graduates"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing dollars spent per year; the allowed unit is 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average financial recovery time for a household after experiencing {placeholder} in the United States (months per incident)?

**Values:** ["job loss", "medical emergency", "natural disaster", "house fire", "identity theft", "car accident", "flooding", "burglary", "divorce", "eviction", "cyberattack", "bankruptcy"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average number of months required for financial recovery per incident; unit: months per incident. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** Estimate the total number of active {placeholder} accounts in the United States (accounts). **Values:** ["monthly", "daily", "weekly", "annual", "student", "business", "retail", "savings", "checking", "joint", "corporate", "online"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the estimated total number of active {placeholder} accounts in the United States. Use 'accounts' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average annual child care expense for households with {placeholder} in the United States (dollars per year).

**Values:** ["single mothers", "two working parents", "military families", "immigrant families", "foster children", "parents under 25", "rural residents", "urban households", "Latino families", "Asian American parents", "low-income households", "same-sex couples"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average annual child care expense. The allowed unit is 'dollars per year.' Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average monthly income earned from {placeholder} by a typical household in the United States (dollars per month)?

**Values:** ["rental properties", "dividends", "social security", "freelance work", "pension", "stock investments", "side business", "online sales", "royalties", "child support payments", "interest income", "government assistance"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value. The unit must be dollars per month. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average monthly expenditure on {placeholder} for a single adult in the United States (dollars per month)?

**Values:** ["groceries", "clothing", "toiletries", "household supplies", "electronics", "furniture", "pet food", "medications", "personal care products", "cleaning products", "books", "stationery"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average monthly expenditure, reported in dollars per month. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual amount spent on replacement of {placeholder} in the United States (dollars per year).

**Values:** ["automobile tires", "roof shingles", "water heaters", "air conditioners", "refrigerators", "light bulbs", "cell phones", "laptop computers", "washing machines", "televisions", "furnaces", "car batteries"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure it is reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

Question: Estimate the average annual household spending on internet services for {placeholder} in the United States (dollars per year).

**Values:** ["urban families", "rural households", "millennials", "retirees", "college students", "single-parent families", "low-income households", "high-income households", "suburban residents", "tech enthusiasts", "remote workers", "senior citizens"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual household spending. Use 'dollars per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the total number of operational {placeholder} currently in use by households in the United States (units)?

**Values:** ["refrigerators", "microwave ovens", "dishwashers", "washing machines", "televisions", "air conditioners", "water heaters", "clothes dryers", "vacuum cleaners", "computers", "smoke detectors", "space heaters"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single integer value representing the total count. The allowed unit is 'units'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual electricity consumption of all the {placeholder} in the United States (kWh per year).

**Values:** ["refrigerators", "air conditioners", "washing machines", "televisions", "microwave ovens", "dishwashers", "water heaters", "clothes dryers", "computers", "freezers", "electric stoves", "space heaters"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value for total annual electricity consumption. The allowed unit is kWh per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average monthly payment required to maintain a {placeholder} in the United States (dollars per month)?

**Values:** ["mortgage", "car lease", "health insurance plan", "cell phone plan", "student loan", "gym membership", "internet subscription", "rent", "childcare service", "homeowners insurance policy", "streaming service subscription", "utility bill"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in dollars per month. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average financial loss from one incident of {placeholder} for a household in the United States (dollars per incident).

**Values:** ["burglary", "fire", "flooding", "identity theft", "car theft", "vandalism", "cyberattack", "medical emergency", "appliance failure", "water leak", "windstorm damage", "earthquake"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value in dollars per incident, e.g., '700 dollars per incident'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual spending on repairs for {placeholder} by households in the United States (dollars per year).

Values: ["roofing", "plumbing systems", "heating systems", "air conditioning units", "water heaters",
"kitchen appliances", "garage doors", "windows", "electrical wiring", "flooring surfaces", "exterior
siding", "bathroom fixtures"], You extract a single numeric value from an answer string. Return only
the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total annual spending. The allowed unit is 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual spending on {placeholder} for all households in the United States (dollars per year).

**Values:** ["toilet paper", "laundry detergent", "pet food", "paper towels", "coffee beans", "bottled water", "diapers", "cleaning supplies", "milk", "bread", "trash bags", "dish soap"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of miles traveled per year by all {placeholder} in the United States (miles per year).

**Values:** ["cars", "pickup trucks", "motorcycles", "school buses", "delivery vans", "semi trucks", "city buses", "SUVs", "minivans", "ambulances", "fire trucks", "taxis"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total miles per year. The allowed unit is 'miles per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual federal budget allocation for {placeholder} in the United States (dollars per year).

**Values:** ["Medicaid", "defense spending", "education", "infrastructure", "scientific research", "veterans benefits", "environmental protection", "homeland security", "agriculture subsidies", "public transportation", "student loans", "healthcare"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual amount, reported in 'dollars per year.' Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the projected lifetime total spent on {placeholder} for an average adult in the United States (dollars over a lifetime)?

**Values:** ["groceries", "healthcare", "education", "transportation", "housing", "vacations", "clothing", "dining out", "insurance premiums", "pet care", "entertainment", "childcare"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing dollars over a lifetime. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual insurance payout for claims related to {placeholder} in the United States (dollars per year)?

**Values:** ["diabetes", "cancer", "stroke", "heart attack", "asthma", "arthritis", "COPD", "hypertension", "kidney failure", "depression", "HIV/AIDS", "multiple sclerosis"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average annual insurance payout for claims related to {placeholder}, reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average amount of {placeholder} contributed to retirement accounts per working adult in the United States each year (dollars per year)?

**Values:** ["income", "salary", "wages", "bonus", "commission", "overtime pay", "dividends", "interest", "tax refund", "gift money", "inheritance", "profit"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average annual contribution and include the unit 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual revenue generated from {placeholder} by a small business in the United States (dollars per year)?

**Values:** ["online sales", "consulting services", "product subscriptions", "advertising", "affiliate marketing", "retail operations", "event hosting", "franchise fees", "membership dues", "service contracts", "licensing agreements", "training workshops"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual revenue. The allowed unit is dollars per year. Return "None" if there is no numerical value.

2184 Answer text: {answer\_text}

**Question:** What is the average annual depreciation rate of a {placeholder} in the United States (% per year)?

**Values:** ["sedan", "pickup truck", "SUV", "motorcycle", "RV", "commercial van", "luxury car", "electric vehicle", "hybrid car", "sports car", "minivan", "cargo trailer"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing an annual depreciation rate, using percent per year (% per year) as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual household spending on managing {placeholder} in the United States (dollars per year)?

**Values:** ["diabetes", "asthma", "hypertension", "arthritis", "obesity", "depression", "allergies", "cancer", "heart disease", "migraine", "eczema", "osteoporosis"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual household spending. Unit must be 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual household consumption of {placeholder} in the United States (kilograms per year).

**Values:** ["sugar", "rice", "beef", "cheese", "chicken", "potatoes", "wheat flour", "eggs", "milk powder", "pasta", "apples", "fish"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing annual household consumption. Only use 'kilograms per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

Question: Estimate the total number of {placeholder} applications submitted in the United States in one year (applications per year).

Values: ["patent", "trademark", "copyright", "asylum", "immigration", "student visa", "work permit", "green card", "welfare", "unemployment benefit", "Medicaid", "social security"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing applications per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

## SOCIAL PLANNING

**Question:** What is the estimated average salary for {placeholder} working in metropolitan areas of Canada? (dollars per year)

Values: ["software engineers", "accountants", "registered nurses", "civil engineers", "marketing managers", "financial analysts", "primary school teachers", "construction managers", "graphic designers", "pharmacists", "electricians", "lawyers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value in dollars per year, representing the average annual salary for {placeholder} working in Canadian metropolitan areas. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual revenue generated from {placeholder} by a midsize US city (dollars per year)?

**Values:** ["parking fees", "property taxes", "sales taxes", "hotel occupancy taxes", "business licenses", "building permits", "utility services", "transit fares", "recreational facilities", "zoning applications", "waste collection services", "franchise agreements"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average annual revenue, and ensure the unit is 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the expected waiting time until {placeholder} is approved by local government in a typical medium-sized city (days)?

**Values:** ["zoning variance", "building permit", "business license", "environmental impact assessment", "liquor license", "noise ordinance waiver", "parking permit", "signage approval", "health code exception", "short-term rental permit", "street closure request", "public event permit"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the expected waiting time, reported in days. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the cumulative hours volunteered by {placeholder} in public community programs over one year in a typical large city (hours per year)?

**Values:** ["teenagers", "retirees", "college students", "corporate employees", "high school teachers", "medical professionals", "single parents", "immigrants", "disabled adults", "faith group members", "government workers", "young adults"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total hours volunteered per year; report the answer in 'hours per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What fraction of total emergency shelter usage is accounted for by {placeholder} in a typical large metropolitan area (% of total usage)?

Values: ["fire", "flood", "earthquake", "hurricane", "tornado", "winter storm", "heatwave", "pandemic outbreak", "power outage", "chemical spill", "civil unrest", "building collapse"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing a percentage. Only provide the answer as '% of total usage'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the estimated total area covered by {placeholder} in public parks of a typical large city (square meters)?

**Values:** ["grass", "flowerbeds", "trees", "playgrounds", "ponds", "walking paths", "bushes", "picnic areas", "sports fields", "dog parks", "benches zones", "community gardens"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report it in square meters. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** How many individual {placeholder} are installed in public recreation facilities of a typical mid-sized U.S. city (units)?

**Values:** ["basketball hoops", "treadmills", "picnic tables", "water fountains", "playground swings", "benches", "soccer goals", "volleyball nets", "bike racks", "trash cans", "lighting fixtures", "tennis courts"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value for the estimated number, reported in units. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of hours per week that a typical community center in a mid-sized city allocates to {placeholder} (hours per week)?

**Values:** ["youth programs", "fitness classes", "arts workshops", "senior activities", "sports leagues", "volunteer events", "after-school tutoring", "language courses", "music lessons", "computer training", "parent meetings", "health seminars"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average number of hours per week allocated to {placeholder} in community centers. Only report using 'hours per week'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual energy consumption attributable to {placeholder} in public facilities of a typical large city (kWh per year).

**Values:** ["lighting", "heating", "air conditioning", "ventilation systems", "water heating", "elevators", "computers", "security systems", "kitchen appliances", "pumping stations", "outdoor lighting", "refrigeration"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing energy consumption, using kWh per year as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total person-hours spent responding to {placeholder} by municipal emergency services per year in a typical large metropolitan area (person-hours per year).

Values: ["structure fires", "medical emergencies", "traffic accidents", "hazardous material spills", "water rescues", "wildlife incidents", "natural disasters", "false alarms", "gas leaks", "power outages", "missing persons cases", "active shooter situations"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing total person-hours, and ensure the unit is 'person-hours per year'. Return "None" if there is no numerical value.

2331 Answer text: {answer text}

**Question:** Estimate the total mass of all {placeholder} currently deployed in public transportation systems of a typical large metropolitan area (tons).

**Values:** ["electric buses", "diesel buses", "tram cars", "subway trains", "hybrid buses", "trolleybuses", "light rail vehicles", "autonomous shuttles", "double-decker buses", "articulated buses", "compressed natural gas buses", "monorail cars"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in tons. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of {placeholder} provided by local governments per year in a mid-sized European city (services per year)?

**Values:** ["library books", "building permits", "waste collections", "public events", "health inspections", "recycling pickups", "housing grants", "school meals", "parking tickets", "bus routes", "street repairs", "water quality tests"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in services per year. Return "None" if there is no numerical value.

2353 Answer text: {answer\_text}

**Question:** What is the estimated annual growth rate of {placeholder} implemented for urban safety in a typical large city (% per year)?

**Values:** ["surveillance cameras", "facial recognition systems", "emergency alert apps", "crime prediction algorithms", "traffic monitoring sensors", "smart street lighting", "body-worn cameras", "drones for patrols", "gunshot detection systems", "license plate readers", "panic button networks", "public Wi-Fi hotspots"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the annual growth rate. The allowed unit is % per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the total number of {placeholder} currently accessible in all municipal libraries of a typical large city (units)?

**Values:** ["books", "magazines", "newspapers", "audiobooks", "DVDs", "manuscripts", "maps", "journals", "ebooks", "reference guides", "music CDs", "archives"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value for quantity and use 'units' as the reporting unit. Return "None" if there is no numerical value.

2373
2374
Answer text: {answer\_text}

Question: What is the elasticity of demand for {placeholder} with respect to price in a typical mid-sized U.S. city (unitless)?

**Values:** ["gasoline", "electricity", "apples", "coffee", "bread", "public transportation", "internet service", "movie tickets", "bottled water", "milk", "restaurant meals", "cigarettes"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing elasticity, which must be unitless (no units). Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the estimated annual quantity of {placeholder} required for municipal road maintenance in a typical large city (tons per year)?

**Values:** ["asphalt", "gravel", "salt", "sand", "concrete", "bitumen", "crushed stone", "recycled asphalt pavement", "topsoil", "cement", "road base material", "aggregate"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and report only in 'tons per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average number of hours spent on {placeholder} per month in a typical mid-sized city's social planning department (hours per month)?

**Values:** ["community outreach", "data analysis", "policy drafting", "stakeholder meetings", "public consultations", "report writing", "budget planning", "event coordination", "staff training", "grant applications", "project evaluation", "interdepartmental collaboration"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value only. The allowed unit is 'hours per month'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the typical number of hours per week allocated to {placeholder} in municipal youth programs in a large metropolitan area (hours per week)?

Values: ["physical education", "arts instruction", "STEM activities", "community service", "leadership training", "sports practice", "mentoring sessions", "health education", "language classes", "career exploration", "environmental projects", "technology workshops"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit reported is 'hours per week'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

 **Question:** What is the average number of public health interventions specifically targeting {placeholder} launched by municipal governments per year in a typical urban area (interventions per year)? **Values:** ["diabetes", "asthma", "influenza", "obesity", "hypertension", "tuberculosis", "depression", "HIV/AIDS", "malaria", "measles", "dengue fever", "hepatitis"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the annual count of interventions, using 'interventions per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average response time for a {placeholder} reported to social services in a typical metropolitan area (minutes)?

**Values:** ["child abuse case", "domestic violence incident", "elder neglect report", "sexual assault allegation", "runaway youth report", "human trafficking tip", "mental health crisis", "drug overdose

call", "missing person report", "animal cruelty complaint", "suicide threat alert", "youth truancy notification"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing average response time. The only allowed unit is minutes. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average annual salary paid to {placeholder} working in municipal planning departments in a typical U.S. city (dollars per year)?

**Values:** ["urban planners", "civil engineers", "GIS specialists", "zoning inspectors", "transportation analysts", "environmental planners", "land use planners", "city managers", "planning technicians", "historic preservationists", "community development coordinators", "housing analysts"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average annual salary and ensure the unit is dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the total amount of {placeholder} distributed by social welfare agencies per year in a typical large city (kilograms per year)?

**Values:** ["food", "rice", "flour", "sugar", "vegetables", "meat", "bread", "milk powder", "canned goods", "lentils", "clothing", "diapers"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in kilograms per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total annual revenue generated from {placeholder} offered by local governments in a typical large metropolitan area (dollars per year).

**Values:** ["municipal bonds", "tax-exempt loans", "public pension funds", "lottery tickets", "parking permits", "business licenses", "building permits", "property tax collections", "transit passes", "utility bills", "court fines", "development fees"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the total number of {placeholder} processed by city housing departments per year in a typical medium-sized U.S. city (applications per year).

Values: ["rental applications", "permit applications", "eviction filings", "subsidy requests", "complaint forms", "inspection requests", "appeals", "zoning applications", "lease renewals", "housing vouchers", "building code violations", "affordable housing applications"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit. Extraction Prompt: Extract a single numeric value for the total annual applications, using 'applications per year' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** What is the average maintenance cost of {placeholder} for municipal infrastructure per year in a typical large city (dollars per year)?

Values: ["road resurfacing", "stormwater management", "waste collection", "street lighting", "bridge repairs", "sidewalk upkeep", "traffic signal maintenance", "park landscaping", "public transit facilities", "sewer system repairs", "snow removal operations", "drainage system cleaning"], You extract

a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average maintenance cost, reported in dollars per year. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** On average, how many trips by {placeholder} are made for community events per month in a typical mid-sized city (trips per month)?

**Values:** ["bus", "taxi", "rideshare", "bicycle", "scooter", "carpool", "minivan", "shuttle", "motorcycle", "vanpool", "tram", "light rail"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the number of trips by the specified vehicle/mode per month, using 'trips per month' as the unit. Return "None" if there is no numerical value.

Answer text: {answer\_text}

## TRANSPORTATION

**Question:** What is the average number of incidents of {placeholder} reported per 1,000 train journeys worldwide each year?

**Values:** ["signal failure", "track obstruction", "door malfunction", "brake failure", "engine overheating", "derailment", "power outage", "passenger injury", "collision", "vandalism", "fire outbreak", "communication breakdown"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing the average number of incidents per 1,000 train journeys per year. The allowed unit is 'incidents per 1,000 train journeys per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** How much is spent on {placeholder} for public transportation in the United States per year (in dollars per year)?

**Values:** ["fuel", "maintenance", "labor", "insurance", "vehicles", "infrastructure", "security", "cleaning", "technology upgrades", "administration", "marketing", "energy"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value and ensure the unit is 'dollars per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Approximately what percentage of annual global freight is transported using {placeholder} each year (% per year)?

**Values:** ["container ships", "railroads", "air cargo", "trucks", "bulk carriers", "pipelines", "coastal shipping", "river barges", "automated guided vehicles", "drones", "roll-on/roll-off vessels", "intermodal transport"], You extract a single numeric value from an answer string. Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing a percentage, with unit '% per year'. Return "None" if there is no numerical value.

Answer text: {answer\_text}

**Question:** Estimate the average fuel consumption for {placeholder} in city traffic, measured in liters per 100 kilometers.

Values: ["rush hour", "weekday mornings", "weekend evenings", "holiday season", "summer months",

"winter conditions", "rainy days", "snowy periods", "peak traffic hours", "nighttime driving", "school drop-off hours", "festive weekends"], You extract a single numeric value from an answer string.

Return only the number and ensure it matches the requested unit.

**Extraction Prompt:** Extract a single numeric value representing fuel consumption, reported in liters per 100 kilometers. Return "None" if there is no numerical value.

Answer text: {answer\_text}

## REFERENCES

- Jacob Andreas. Language models as agent models. arXiv preprint arXiv:2212.01681, 2022.
- Chris L. Baker, Joshua B. Tenenbaum, and Rebecca R. Saxe. Goal inference as inverse planning. In *Proceedings of the 29th Annual Conference of the Cognitive Science Society*, pages 779–784, 2007.
- Chris L. Baker, Rebecca Saxe, and Joshua B. Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 2021.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- Melanie Brucks and Olivier Toubia. Prompt architecture induces methodological artifacts in large language models. *PloS one*, 20(4):e0319159, 2025.
- Barry R Chiswick and Paul W Miller. Linguistic distance: A quantitative measure of the distance between english and other languages. *Journal of multilingual and multicultural development*, 26 (1):1–11, 2005.
- Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 489–500, 2018.
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz. Data augmentation for low-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 567–573, 2017.
- H. Paul Grice. Meaning. The Philosophical Review, 66(3):377–388, 1957.
- H. Paul Grice. Studies in the Way of Words. Harvard University Press, 1989.
- Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and Shuo Shuo Liu. Bias in large language models: Origin, evaluation, and mitigation. *arXiv* preprint *arXiv*:2411.10915, 2024.
- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018a.
- David Ha and Jürgen Schmidhuber. World models. arXiv preprint arXiv:1803.10122, 2018b.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv* preprint arXiv:1912.01603, 2019.
  - Jennifer Hu, Roger Levy, Judith Degen, and Sebastian Schuster. In-context learning enables models to learn pragmatics. *arXiv* preprint arXiv:2311.08742, 2023.
  - Pierre Jacob. Intentionality. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Stanford University, 2019.

- Shervin Javdani, Siddhartha S. Srinivasa, and J. Andrew Bagnell. Shared autonomy via deep reinforcement learning. *Robotics: Science and Systems*, 2018.
- Jiaxuan Jiang, Jiapeng Liu, Miłosz Kadziński, and Xiuwu Liao. A bayesian network approach for dynamic behavior analysis: Real-time intention recognition. *Information Fusion*, 118:102873, 2025.
  - Marzena Karpinska and Mohit Iyyer. Large language models effectively leverage document-level context for literary translation, but critical errors persist. *arXiv preprint arXiv:2304.03245*, 2023.
  - Michal Kosinski. Theory of mind might have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*, 2023.
  - Molly Lewis, Aoife Cahill, Nitin Madnani, and James Evans. Local similarity and global variability characterize the semantic space of human languages. *Proceedings of the National Academy of Sciences*, 120(51):e2300986120, 2023.
  - Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. *Proceedings of the International Conference on Learning Representations*, 2023.
  - Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction for mdps. *AI&M*, 1(2):3, 2006.
  - Samuel Louvan and Bernardo Magnini. Recent neural methods on slot filling and intent classification for task-oriented dialogue systems: A survey. *arXiv preprint arXiv:2011.00564*, 2020.
  - Gary Marcus. Deep learning is hitting a wall. Nautilus, 2022.
  - Libo Qin, Qiguang Chen, Wanxiang Che, Yangming Li, Minheng Ni, Yue Li, Min Liu, Weiwei Deng, and Ting Liu. A survey on spoken language understanding: Recent advances and new frontiers. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3032–3045, 2021.
  - Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. arXiv preprint arXiv:2310.11324, 2023.
  - Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, 2016.
  - Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150*, 2022.
  - Adam Smith. The Theory of Moral Sentiments. Penguin Books, 1759. Republished 2010.
  - Yu Tan and L Elisa Celis. Assessing social and intersectional biases in contextualized word representations. In *NeurIPS*, 2021.
  - Shubham Toshniwal, Sam Wiseman, Karen Livescu, and Kevin Gimpel. Chess as a testbed for language model state tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11385–11393, 2022.
  - Keyon Vafa, Justin Y. Chen, Jon Kleinberg, Sendhil Mullainathan, and Ashesh Rambachan. Evaluating the world model implicit in a generative model. *Advances in Neural Information Processing Systems*, 37, 2024.
  - Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Unsupervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, volume 33, pages 6256–6268, 2020.
  - Jianhao Yan, Pingchuan Yan, Yulong Chen, Jing Li, Xianchao Zhu, and Yue Zhang. Benchmarking gpt-4 against human translators: A comprehensive evaluation across languages, domains, and expertise levels. *arXiv* preprint arXiv:2411.13775, 2024a.

Jianhao Yan, Pingchuan Yan, Yulong Chen, Judy Li, Xianchao Zhu, and Yue Zhang. Gpt-4 vs. human translators: A comprehensive evaluation of translation quality across languages, domains, and expertise levels. arXiv preprint arXiv:2407.03658, 2024b. Hyejin Youn, Logan Sutton, Eric Smith, Cristopher Moore, Jon F Wilkins, Ian Maddieson, William Croft, and Tanmoy Bhattacharya. On the universal structure of human lexical semantics. Proceed-ings of the National Academy of Sciences, 113(7):1766–1771, 2016. Xiaodong Zhang and Houfeng Wang. A survey of joint intent detection and slot filling models in natural language understanding. ACM Computing Surveys, 55(8):1–38, 2022. Jingming Zhuo, Songyang Zhang, Xinyu Fang, Haodong Duan, Dahua Lin, and Kai Chen. Prosa: Assessing and understanding the prompt sensitivity of llms. arXiv preprint arXiv:2410.12405,