
Trilemma of Truth in Large Language Models

Germans Savcisens
Northeastern University
Boston, USA
germans@savcisens.com

Tina Eliassi-Rad
Northeastern University
Boston, USA
tina@eliassi.org

Abstract

The public often attributes human-like qualities to large language models (LLMs) and assumes they “know” certain things. In reality, LLMs encode information retained during training as internal probabilistic knowledge. This study examines existing methods for probing the veracity of that knowledge and identifies several flawed underlying assumptions. To address these flaws, we introduce sAwMIL (Sparse-Aware Multiple-Instance Learning), a multiclass probing framework that combines multiple-instance learning with conformal prediction. sAwMIL leverages internal activations of LLMs to classify statements as true, false, or neither. We evaluate sAwMIL across 16 open-source LLMs, including default and chat-based variants, on three new curated datasets. Our results show that (1) common probing methods fail to provide a reliable and transferable veracity direction and, in some settings, perform worse than zero-shot prompting; (2) truth and falsehood are not encoded symmetrically; and (3) LLMs encode a third type of signal that is distinct from both true and false.

1 Introduction

Can we trust the content generated by large language models (LLMs)? Recent literature suggests that LLMs possess internal probabilistic knowledge [1, 2, 3, 4, 5]. However, our understanding of how LLMs use this internal knowledge (if at all) remains fragmented. It is known that LLMs are indifferent to the veracity of their outputs [6] and often hallucinate [7]. Furthermore, it is often difficult for users to recognize hallucinations because LLMs produce fluent and persuasive text. For example, Church [8] shows that students trust factually incorrect answers from GPT due to their authoritative and confident tone, and Williams et al. [9] demonstrate that users rate disinformation generated by LLMs as equally or even more credible than human-generated content. Thus, we need a method to assess the truthfulness of internal probabilistic knowledge to make user interactions with LLMs more reliable.

Prompt-based evaluations (see Fig. 1A) rely on the idea that we can simply ask an LLM about its knowledge. Abbasi Yadkori et al. [10] introduce an information-theoretic prompt-based evaluation, while Xu et al. [11] propose a training framework to produce prompts with self-reflective rationales, and Farquhar et al. [12] introduce uncertainty estimators to detect inconsistent text generations. However, prompt-based evaluations are sensitive to the input’s phrasing [13] and content [14].

A more direct approach is to examine *how* LLMs represent text internally (see Fig. 1B). Consider a large language model, \mathcal{M} , with a vocabulary \mathcal{V} . The LLM maps the input text \mathbf{x} to a probability distribution over subsequent tokens, denoted $P_{\mathcal{M}}$:

$$\mathcal{M}(\mathbf{x}) = P_{\mathcal{M}}(\tau \mid \mathbf{x}), \text{ where } \tau \in \mathcal{V}. \quad (1)$$

For any token $\tau \in \mathcal{V}$, the distribution $P_{\mathcal{M}}(\tau \mid \mathbf{x})$ denotes the probability that τ is the continuation of the sequence \mathbf{x} . To compute the conditional distribution $P_{\mathcal{M}}$, an LLM transforms \mathbf{x} into *intermediate*

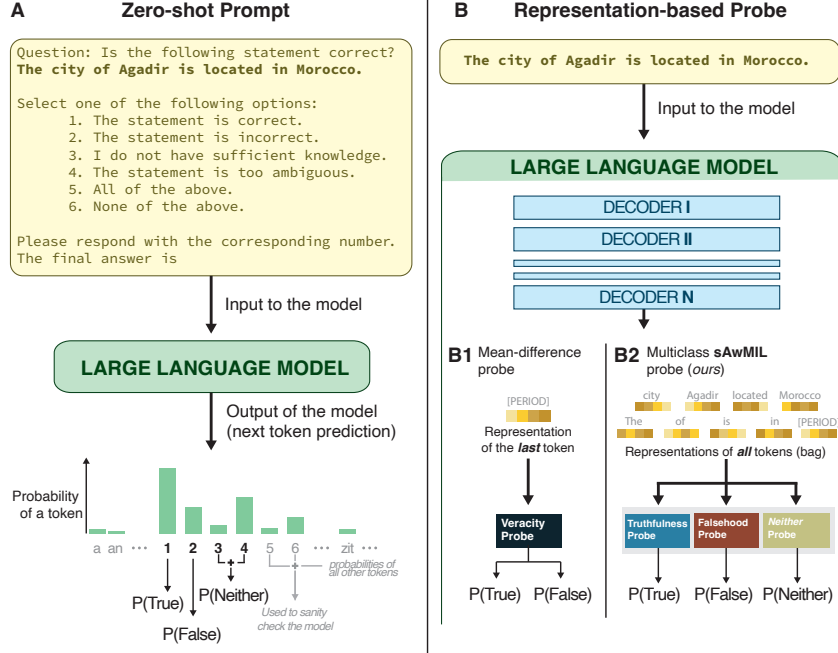


Figure 1: Overview of methods for probing veracity in LLMs. **(A)** In zero-shot prompting, a target statement is inserted into a structured prompt instructing the LLM to select an answer from a specific set of tokens. The LLM’s prediction is based on the probabilities of these tokens. This method treats the model as a black box and examines its (input, output) pairs. **(B)** In representation-based probing, the analysis is done on the internal representations generated by intermediate decoders. **(B1)** The mean-difference probe [15] is a common method for determining the veracity of a statement based on the representation of the last token. This approach outputs probabilities for *true* or *false* statements, but cannot account for statements that lack a definitive truth value. **(B2)** Our probe, multiclass sparse aware MIL (sAwMIL), looks at the representation of every token in a statement and provides probabilities for three classes: *true*, *false*, and *neither*. Multiclass sAwMIL can account for cases when the LLM does not have any knowledge about the statement.

neural activations denoted $h_i(\mathbf{x}) \in \mathbb{R}^{L \times d}$. Here $h_i(\mathbf{x})$ denotes the neural activations of \mathcal{M} after the i th decoder, d stands for the dimensionality of the decoder, and L stands for the length of the sequence \mathbf{x} . We can probe these intermediate activations to identify *veracity signals*, i.e., to isolate activation patterns that correspond to truthful statements.¹ For example, Azaria and Mitchell [16] train a neural network to classify statements as *true* or *false* based on these internal representations. Similarly, Marks and Tegmark [15] use a mean-difference classifier to linearly separate *true* or *false* statements (see Fig. 1.B1). Further examples include the improved linear classifier introduced by Bürger et al. [17] and a semi-supervised method based on the contrastive pairs of statements [18]. Collectively, these works rely on the idea that given a data set $\langle \mathbf{x}, y \rangle \in \mathcal{D}$ with some statements \mathbf{x} and veracity labels $y \in \mathcal{Z}$, we can train a probe g_i that maps neural activations $h_i(\mathbf{x})$ to the distribution $G_{\mathcal{M}}$ over \mathcal{M} ’s veracity labels:

$$g_i(h_i(\mathbf{x})) = G_{\mathcal{M}}(z \mid \mathbf{x}), \text{ where } z \in \{\text{true}, \text{false}\} \quad (2)$$

However, we observe that existing probing methods often rely on flawed assumptions, which limit the reliability of their findings (for an overview, refer to Supplementary Tab. 4). We argue for a three-valued logic approach (as in Fig. 1.B2) as the more appropriate method to model veracity in LLMs. Our method sAwMIL (short for Sparse Aware Multiple-Instance Learning) combines Multiple Instance Learning (MIL) [19] and Conformal Predictions (CP) [20] to allow for a flexible probe that can handle ‘*neither*’ statements and quantify uncertainty.

¹We use the terms ‘pattern’ and ‘signal’ interchangeably. Similarly, ‘neurons’ and ‘features’ are used to refer to individual components of a signal. In this context, a signal or pattern denotes a set of features that operate collectively.

In summary, our contributions include the following.

1. We identify and discuss five flawed assumptions in the current veracity-probing literature.
2. We show that common linear classifiers do not capture reliable veracity directions.
3. We propose a novel multiclass linear probing method sAwMIL based on Multiple Instance Learning (MIL) [19] and Conformal Prediction [20, 21].
4. We present three new data sets containing statements labeled *true*, *false*, and *neither*² to enable more rigorous evaluations of veracity probes.

2 Background and Flawed Assumptions When Probing Veracity in LLMs

An LLM, \mathcal{M} , has **internal probabilistic knowledge** $K_{\mathcal{M}}$, which it acquires during training.³ To determine the veracity of a statement ϕ , the model \mathcal{M} should be able to distinguish between three scenarios:

1. ϕ is **True** if there is sufficient support for ϕ given $K_{\mathcal{M}}$:

$$P(\phi \mid K_{\mathcal{M}}) \geq \zeta, \text{ where } \zeta \in (0, 1] \text{ is a threshold.}$$

2. ϕ is **False** if there is sufficient support for $\neg\phi$ given $K_{\mathcal{M}}$:

$$P(\neg\phi \mid K_{\mathcal{M}}) \geq \zeta, \text{ where } \zeta \in (0, 1] \text{ is a threshold.}$$

3. ϕ is **Neither** if there is not sufficient support for ϕ and $\neg\phi$ given $K_{\mathcal{M}}$:

$$\left[P(\phi \mid K_{\mathcal{M}}) < \zeta \right] \text{ and } \left[P(\neg\phi \mid K_{\mathcal{M}}) < \zeta \right], \text{ where } \zeta \in (0, 1] \text{ is a threshold.}$$

If \mathcal{M} has a mechanism to determine the veracity of a statement ϕ , then \mathcal{M} should encode the signal associated with the veracity in its intermediate activations:

1. **Truthfulness:** \mathcal{M} generates an activation pattern that encodes support for ϕ in $K_{\mathcal{M}}$, reflecting the model’s internal support for the statement ϕ .
2. **Falsehood:** \mathcal{M} produces an activation pattern that reflects a lack of sufficient support for ϕ , instead indicating that the internal knowledge $K_{\mathcal{M}}$ provides stronger support for $\neg\phi$ (e.g., signaling a contradiction or misalignment with known facts).
3. **Neither:** \mathcal{M} should encode the lack of support for ϕ and $\neg\phi$, indicating that the veracity of ϕ is currently undefined. That is, ϕ is neither true nor false.

2.1 Flawed Assumptions When Probing Veracity in LLMs

To train and evaluate a veracity probe g_i , a labeled data set \mathcal{D} is assembled. This data set consists of pairs of neural activations and ground-truth labels, denoted as $\langle h_i(\mathbf{x}), y \rangle$, where h_i is the activations after the i th decoder and labels y specify the veracity label Z . In most cases, $Z \in \{\text{true}, \text{false}\}$. The probe g_i is trained on the train split $\mathcal{D}_{\text{train}} \subseteq \mathcal{D}$ and evaluated on the test split $\mathcal{D}_{\text{test}} \subseteq \mathcal{D}$. The intersection between $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ is empty.

We focus exclusively on linear probes, where the parameters of g_i define a linear direction \vec{v}_i for the veracity signal after the i th decoder of \mathcal{M} :

$$g_i(\mathbf{x}) = \mathbf{x} \boldsymbol{\theta}^T + b, \text{ where } \boldsymbol{\theta} \in \mathbb{R}^{1 \times d}, b \in \mathbb{R} \text{ are parameters learned on } \mathcal{D}_{\text{train}} \text{ and } \mathbf{x} \in \mathbb{R}^{1 \times d}. \quad (3)$$

Next, we provide a detailed overview of the flawed assumptions made in the existing literature. Refer to the Supplementary Tab. 4 for a condensed overview of flawed assumptions.

Flawed Assumption I: Truth and falsehood are bidirectional. To determine the veracity of a statement ϕ , a large language model \mathcal{M} must develop a mechanism to detect ϕ ’s truth or falsehood.⁴

²Throughout the paper, we use the terms *neither*, *neither-valued*, *neither-type*, and *neither-true-nor-false* interchangeably to refer to statements that are neither true nor false. When used as a class label, we italicize *neither* to distinguish it from the regular use of the word. We similarly italicize words such as *true* and *false* when referring to class labels.

³Supplementary Tables 2 and 3, respectively, list the notations and abbreviations used in this paper.

⁴In this example, we assume that the veracity label is $Z \in \{\text{true}, \text{false}\}$.

This mechanism must rely on \mathcal{M} 's neural activations to find support for ϕ by using \mathcal{M} 's internal probabilistic knowledge $K_{\mathcal{M}}$. Existing veracity probes [15, 17, 18, 22] implicitly assume that truth and falsehood are encoded bidirectionally. That is,

$$P(\phi \mid K_{\mathcal{M}}) = 1 - P(\neg\phi \mid K_{\mathcal{M}}) \quad (4)$$

This formulation implies (1) any statement ϕ not confirmed as *true* is considered *false*, and (2) each decoder symmetrically encodes a signal corresponding to falsehood and truthfulness. However, there is little support to justify either scenario. Similarly, Bürger et al. [17] and Marks and Tegmark [15] suggest that veracity exists along more than two directions.

Valid Assumption I (Truthfulness and falsehood have distinct directions). *The representation of truth and falsehood requires more than one direction. This is, $P(\phi \mid K_{\mathcal{M}}) \neq 1 - P(\neg\phi \mid K_{\mathcal{M}})$.*

Flawed Assumption II: LLMs capture and retain everything we know. To train a probe g_i , we use \mathcal{D}_{train} that consists of pairs of factual statements and ground-truth labels $\langle x_i, y_i \rangle$, where (usually) $y \in \{\text{true}, \text{false}\}$. The labels y that we assign to the statements in \mathcal{D} are based on *our* knowledge (i.e., what we know to be true). Thus, the veracity labels in \mathcal{D} are distributed according to $G_{\mathcal{D}}$.

Our goal, however, is to train a veracity probe g_i that classifies *what the LLM deems to be true, false, or 'neither'*. So, the probe g_i should map the statements to the space of the LLM's internal probabilistic knowledge $K_{\mathcal{M}}$. Although most recent studies use open-source models, the precise composition of their training data remains largely unknown [23, 24], and we do not have straightforward methods to verify what has been incorporated into internal probabilistic knowledge $K_{\mathcal{M}}$. Thus, \mathcal{M} may have a different distribution $G_{\mathcal{M}}$ for the veracity labels. That is, $G_{\mathcal{M}}$ may not be equivalent to $G_{\mathcal{D}}$. For example, we know that "The city of Bissau is in Congo" has a ground-truth label $y = \text{false}$, because we can check maps or official sources. On the other hand, we do not know how \mathcal{M} labels it.

Recent probing methods [15, 17, 18, 22, 25] cannot account for the mismatch between the label distributions. Instead, these probes introduce a systemic bias, where g_i captures a signal that reflects *our* labeling choices rather than the model's true internal representations.

Valid Assumption II (LLMs do not capture and retain everything we know). *The distribution of ground-truth labels $G_{\mathcal{D}}$ may not be equivalent to the model's label distribution $G_{\mathcal{M}}$.*

Flawed Assumption III: All veracity probes provide calibrated probabilities. Veracity probes are generally designed to predict discrete labels (e.g., methods introduced by Azaria and Mitchell [16] or Marks and Tegmark [15]). That is, they are classification tasks where the probe assigns one of two labels to a given statement: $g_i : h_i(x) \rightarrow \{\text{true}, \text{false}\}$. However, as Herrmann and Levinstein [26] point out, veracity probes should provide not only discrete labels, but also values that can be interpreted as degrees of belief (or some other alternative that quantifies confidence).

Valid Assumption III (The probabilities generated by veracity probes are not inherently calibrated). *The output of veracity probes g_i may not be calibrated and require additional post-processing to be interpreted as meaningful estimates of confidence.*

Flawed Assumption IV: Every statement is either true or false. There are cases where the LLM lacks definitive evidence to determine if a statement is true or false. Suppose we have a veracity probe g_i , which returns a probability of 0.5 for a given statement ϕ to be true. The question is how to interpret the probability of 0.5. That is, probes that return scores (e.g., distance from the separating hyperplane) rather than probabilities cannot provide uncertainty estimates. For instance, in Supplementary Sec. H, we show an example where the probe assigns high scores to cases when statements do not have any definitive truth-value.⁵

To address the issue, we have to train probes that can account for *neither* cases, so that g_i can reflect the insufficient evidence in $K_{\mathcal{M}}$. For instance, studies with human participants have shown that including options such as "other" or "I do not know" can help with data quality [27].

Valid Assumption IV (Some statements are neither true nor false). *A probe g_i should distinguish between the cases where the model \mathcal{M} lacks sufficient support to assess the truthfulness of the statement ϕ , and the cases where ϕ lacks a veracity value.*

⁵Mean-difference probe consistently mislabeled examples where the statement did not have a fully realized truthfulness, or statements that were neither true, nor false (neither-valued statements).

Flawed Assumption V: We know where the signal for veracity is stored. The majority of veracity probes are trained on the representation of the last token [15, 25]. For example, if the statement is “Boston is in the US.”, they assume the period alone carries the entire veracity signal. Such methods assume that any factual signal appearing n tokens before the end of the statement will be faithfully preserved until the last token. A more reasonable approach is to probe at the exact position where the statement is actualized—e.g., immediately after “in the” in the above example—rather than relying on the LLM to *move* that signal all the way to the end of the statement.

Valid Assumption V (Position of the veracity token is not known a priori). *Probes should include a flexible mechanism for identifying the optimal token positions from which to extract veracity signals, instead of relying on fixed positions such as the final token in the statement.*

A probe that directly addresses these flawed assumptions would better reflect the internal knowledge of the LLM and provide a clearer understanding of (1) the factual information encoded in \mathcal{M} , (2) how \mathcal{M} classifies statements as *true*, *false*, or *neither*, and (3) calibrated measures of \mathcal{M} ’s confidence in its own probabilistic knowledge.

3 Method: sAwMIL

To address the flawed assumptions, we propose a multiclass probe, called sAwMIL (short for *sparse aware multiple-instance learning*). It classifies statements into three classes: *true*, *false*, and *neither*. sAwMIL uses multiple-instance learning (MIL) [28] and conformal prediction (CP) [20].⁶

3.1 Sparse Aware Multiple-Instance Learning

Algorithms such as logistic regression, support vector machines, and mean-difference classifiers belong to the single-instance learning (SIL) family, where each instance in the data set has an individual label. In contrast, multiple-instance learning (MIL) is a type of weakly supervised learning that operates on a set of labeled *bags* [28]. A bag, B , is a set of related instances (e.g., patches extracted from the same image or embeddings of individual words in a sentence). Each bag has an associated binary label,⁷ but the labels for individual instances within the bag remain *unknown*. A positive label ($y = 1$) indicates that at least one instance in the bag B belongs to the positive class. Thus, an MIL algorithm must identify the most influential instances contributing to the bag’s label.

These algorithms must consider the overall structure of the bag and simultaneously suppress irrelevant instances. Bunescu and Mooney [19] introduced sparse balanced MIL (sbMIL), an adaptation of linear support vector machines (SVM). It is designed for cases where bags are sparse and only a few instances within a bag are important. sbMIL has two training stages. In the first stage, it uses the MIL-modified SVM [19] (see Fig. 4 in [19]). During this stage, the objective is to identify the most important instances within positive bags, pushing all other instances and negative bags toward the opposite side of the separating hyperplane. Once the initial model is trained, it computes the scores assigned to each instance in all positive bags. Then, it computes the η -quantile, where η is the *balancing hyperparameter*. Instances scoring above the η threshold are marked as positive. In the second stage, we use a single-instance SVM. It works with individual samples, disregarding their original grouping into bags, and assigns them the labels determined during the first stage. We provide the pseudocode for sAwMIL in the Supplementary Alg. 1.

3.1.1 Workflow

One-vs-all sAwMIL. We modify sbMIL since we have an additional piece of information. Given a statement “The city of Riga is in Latvia,” we know which tokens come from the *actualized* part of the statement (e.g., “Latvia”) and which ones come from the *pre-actualized* part of the statement (e.g., “The city of Riga is in”). Hence, after we apply the η -quantile threshold, we add another round of filtering (see Supplementary Alg. 1, Step 6).

To perform the additional filtering, along with the set of instances in the bag x_i and a binary bag label y , we consider the *intra-bag labels* m_i , where $m_i \in \{0, 1\}^{L_i}$. (L_i is the number of items/tokens

⁶The code is available on GitHub at carlomarxdk/trilemma-of-truth, and the data is available on HuggingFace at carlomarxx/trilemma-of-truth.

⁷For simplicity, we assume a binary label.

in the bag.) These intra-bag labels specify the instances where we expect to find a signal. Given a statement, $\mathbf{x} \leftarrow [\mathbf{x}^p, \mathbf{x}^a]$, all the tokens in the pre-actualized part \mathbf{x}^p have an intra-label of 0 (since the factual statement has not yet been actualized), and all the tokens in the actualized part \mathbf{x}^a have a label of 1. To label a sample, this sample should have a score above η -quantile, and it should be part of an actualized part \mathbf{x}^a .

We train three one-vs-all sAwMIL probes that isolate distinct veracity signals:

- **is-true** probe: separates tokens that carry a true signal from all others.
- **is-false** probe: separates tokens that carry a false signal from all others.
- **is-neither** probe: separates tokens that carry neither (not true or false) signal from all others.

Multiclass sAwMIL Ideally, we want a multiclass probe that assigns probabilities to a statement being *true*, *false*, or *neither*. Thus, we assemble the one-vs-all sAwMIL probes into a multiclass probe via *softmax regression*, which takes the outputs of the one-vs-all probes and transforms them into multiclass probabilities. Formally

$$p_k = \frac{\exp(z_k)}{\sum_j \exp(z_j)}, \quad (5)$$

where $z_k = g_i^k(\mathbf{x}) \cdot \alpha_k + \beta_k$ and $k \in \{\text{is-true}, \text{is-false}, \text{is-neither}\}$. Recall that g_i is the trained probe (Eq. 2). It is given neural activations $h_i(\mathbf{x})$.

3.2 Conformal Predictions

Raw outputs from many models, such as Support Vector Machines (SVMs) – specifically, the distance-to-hyperplane score – are not meaningful as confidence measures. Wrapping SVM scores in a sigmoid function to force them into $[0, 1]$ does *not* create calibrated probabilities. They can underestimate their true confidence unless they are explicitly calibrated. Thus, we introduce conformal learning into our probe.

Conformal learning is a framework [20, 21] that enables us to transform raw scores into prediction sets with *guaranteed coverage*. Hence, it provides a method to account for uncertainty. Conformal prediction methods identify intervals within which the probes’ predictions are correct with a probability of $1 - \alpha$. For a detailed description of the nonconformity scores [29], see Sec. G in the Supplementary Material.

4 Experiments

This section outlines our experimental setup, including data sets, our evaluation procedure, and the set of large language models.

4.1 Data

We introduce three new data sets consisting of factually *true*, factually *false*, and *neither*-valued statements. The *neither* statements are the ones whose truthfulness value cannot be determined at the present moment (due to the lack of information). While several benchmark data sets for veracity and factuality evaluation exist, prior work has shown that some of these may be partially included in the pretraining or fine-tuning stages of LLMs [30]. In contrast, our goal is to minimize the risk of data contamination while also maintaining higher control over data provenance and quality. Hence, we assemble new data sets that involve statements related to specific themes (see Tab. 1 for examples):

- **City Locations** data set contains statements about cities and their corresponding countries extracted from the GeoNames geographical database.
- **Medical Indications** data set consists of statements about the medications and their corresponding indications from the DrugBank 5.1 pharmaceutical knowledge base [31]. Medications include the drug and substance names, while indications specify the symptoms or a disease/disorder.
- **Word Definitions** data set is based on the WordsAPI dictionary. Hence, the statement involves words and their synonyms or relations.

Table 1: Composition of data sets used in this work. Number of *true*, *false*, and *neither*-valued statements per data set. **A** stands for the number of affirmative statements, and **N** stands for the number of negated statements. The last column displays example statements with ground truth labels.

Data Set	True	False	Neither	Examples
City Locations	A: 1392 N: 1376	A: 1358 N: 1374	A: 876 N: 876	(True) The city of Mâcon is located in France. (False) The city of Dharân is located in Ecuador. (Neither) The city of Staakess is located in Marbate.
Word Definitions	A: 1234 N: 1235	A: 1277 N: 1254	A: 1747 N: 1753	(True) Corsage is a synonym of a nosegay. (False) Towner is not a type of a resident. (Neither) Kharter is not a synonym of a greging
Medical Indications	A: 1423 N: 1347	A: 1329 N: 1424	A: 478 N: 522	(True) PR-104 is indicated for the treatment of tumors. (False) Zolpidem is indicated for the treatment of angina. (Neither) Alostast is indicated for the treatment of candigemina.

Every data set consists of negated statements like “The city of Riga **is not** located in Estonia.”⁸ and affirmative ones like “Menadione **is** indicated for the treatment of coughs.”⁹ We provide a detailed description of these data sets in the Supplementary Sec. C.

Neither statements. If a statement ϕ is absent from the LLM’s internal probabilistic knowledge $K_{\mathcal{M}}$, then ϕ is *neither* true nor false. It is difficult to determine which statements are absent from $K_{\mathcal{M}}$ because we generally do not have access to the training data sets used to train the LLMs. However, we can create *neither* statements with *synthetic entities*—i.e., entities that do not exist in the real world or fictional works. Since these objects are specifically generated for our experiments, it is highly unlikely that an LLM has learned anything about them during training. Thus, we can use them as *substitutes* for content that LLMs could not have learned—i.e., from the point of view of an LLM, these should be considered neither *true* nor *false*. For a detailed description of the generation [32] of *neither*-valued statements, see Sec. C.1 in the Supplementary Material.

4.2 Evaluation

Language Models. We evaluate 16 open-source LLMs (ranging from 3 to 14 billion parameters) across four families: Gemma/Gemma-2, Llama-3 (v3.1 and v3.2), Mistral-v0.3, and Qwen-2.5. These models run on consumer-grade hardware and are publicly available through HuggingFace [33]. We provide an overview of these models in Sec. D of the Supplementary Material.

Metric. We evaluate probes based on their ability to correctly classify statements into three classes: *true*, *false*, and *neither*. We use Matthew’s Correlation Coefficient (MCC) to summarize the performance of each probe (on the test sets); refer to Eq. 14 in the Supplementary Sec. I for the definition of the multiclass MCC. An $\text{MCC} = 1$ indicates perfect classification, $\text{MCC} = 0$ corresponds to random prediction, and $\text{MCC} = -1$ indicates inverted predictions with respect to ground-truth labels.

Probing Methods. We compare performance across *six* probing methods, grouped into three categories.

In *zero-shot prompting*, we insert each statement into a prompt formulated as a multiple-choice question (see Fig. 1.A and Supplementary Sec. F). We compute the probabilities over the candidate answers to determine how the statement is classified. We estimate the conformal prediction intervals to improve the performance of the zero-shot prompting.

We evaluate three *binary representation-based probes* trained on the representations of the last token: the mean-difference probe (see Supplementary Sec. H.1) introduced by Marks and Tegmark [15], the TTPD probe (short for Training of Truth and Polarity Direction) introduced by Bürger et al. [17], and a supervised PCA classifier. We refer the reader to Algorithms 2, 3, and 4 in the supplementary material for details.

⁸This statement is a factually true and negated statement.

⁹This statement is a factually false and affirmative statement.

Since MD, TTPD, and sPCA probes are trained to separate *true* and *false* statements, we augment each with conformal prediction intervals. Samples falling outside these intervals are labeled as *neither*. We refer to these methods as MD+CP, TTPD+CP, and sPCA+CP, accordingly.

Finally, we include two *multiclass representation-based probes*: a multiclass SVM trained on the representation of the last token, and (6) the multiclass sAwMIL probe, which operates on the whole sentence (i.e., all token representations within a statement), and has in-built CP intervals.

Criteria. We compare all probes under two complementary criteria: (1) *Correlation* and (2) *Generalization* (see Supplementary Sec. E for more details on the validity criteria).

The **correlation criterion** [25, 34] assesses how well a probe g_i , trained on the training split of \mathcal{D}_{train} , performs on the corresponding test split \mathcal{D}_{test} , assuming that both are drawn from the same distribution. For example, a probe trained on statements about city locations should accurately classify other statements from the same domain. The **generalization criterion** [17, 34] evaluates whether a probe g_i , trained on \mathcal{D}_{train} , successfully generalizes to datasets \mathcal{D}'_{test} containing statements from different domain.

Together, these two criteria measure how well the probe identifies the veracity signal. The higher the performance along both, the more probable it is that the probe captured a robust veracity signal and not mere proxies.

5 Results

We first consider the performance of probes under the *correlation* and *generalization* criteria. Mean performances aggregated across all 16 LLMs and three datasets are shown in Fig. 2 (see Supplementary Tab. 20 for more details). For representation-based probes, we report results for two settings: (1) using only the last token’s representation, and (2) using all token representations (the full bag).¹⁰

Binary representation-based probes. The MD+CP, TTPD+CP, and sPCA+CP probes achieve moderate MCC when evaluated on last-token representation (see Fig. 2.A). However, their performance drops substantially when evaluated on the full bags, suggesting sensitivity to non-actualized (non-informative) tokens and noise. Under the generalization criterion (see Fig. 2.B), these probes exhibit limited transferability and degraded performance on the full bags. Another common failure, illustrated in Fig. 2.D and detailed in Supplementary Tables 25–30, is that these probes frequently confidently misclassify *neither*-valued statements as *true* or *false*. While effective under controlled conditions, these binary probes appear to identify proxy directions that partially reflect veracity but are confounded by spurious correlations.

Zero-shot prompting. Unlike representation-based probes, the prompting does not require information about where a factual statement begins or ends. This makes it more robust when working with unstructured text. However, zero-shot prompting exhibits class imbalance: models tend to overpredict one label. For example, in Fig. 2.C, the prompting of Qwen-2.5-14b-instruct overpredicts false label. We see similar skews in the confusion matrices of other models with zero-shot prompting (see Supplementary Tab. 24).

Multiclass representation-based probes. Both multiclass probes, SVM and sAwMIL, achieve substantially higher performance Fig. 2.A). However, the SVM probe, like its binary counterparts, shows a drop in performance under the bag setting, indicating reliance on non-generalizable cues. In contrast, the sAwMIL probe maintains consistent performance across both settings, and yields superior generalization performance (see Fig. 2.B).

Together, these findings suggest that sAwMIL captures more robust and transferable veracity directions, leveraging distributed information across tokens rather than relying on isolated representations.

In Supplementary Sec. I, we detail results showing sAwMIL’s ability to generalize across datasets, as well as the results of targeted interventions on the identified veracity directions \vec{v}_i . The interventions show that, across the majority of LLMs, one can use sAwMIL directions to elicit truthful or false replies. For brevity, we do not cover these in the manuscript.

¹⁰For single-instance probes, when the full bag is provided, the predicted label corresponds to the class with the highest score among all tokens in the statement.

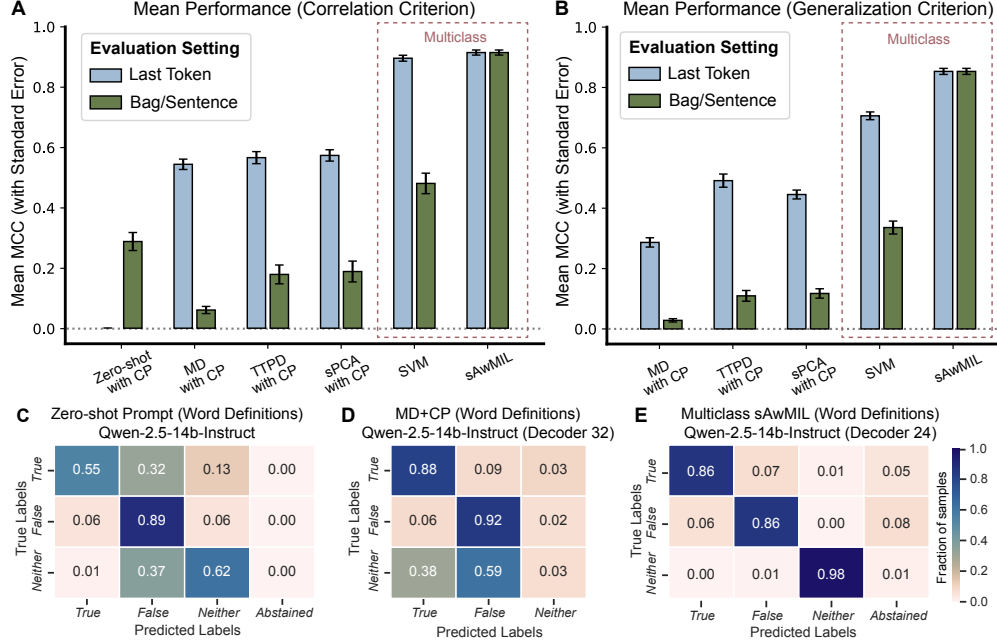


Figure 2: Mean performance of probing methods aggregated across 16 models and three datasets, along with the examples of confusion matrices. Each probe is evaluated under two settings: using only the representation of the last token and using predictions aggregated across the entire bag/sentence. For each probe, we aggregate the statistics using the best-performing layers. (A) *Correlation criterion*. Mean probe performance on the test split of the dataset on which each probe was trained. (B) *Generalization criterion*. Mean probe performance on datasets that were not used for training. Binary probes trained on the last token representation (MD+CP, TPPD+CP, and sPCA+CP) exhibit both lower performance and poorer generalization. The multiclass SVM+CP (also trained on the last token representation) achieves much better performance in A; however, performance degrades when evaluated on the full bag. Importantly, SVM+CP achieves lower generalization performance as compared to the sAwMIL. (C) Confusion matrix for the zero-shot prompting, where the LLM overpredicts the *false* class. (D) Confusion matrix for MD+CP evaluated on the last token representation, where probe incorrectly predicts *neither*-valued statements. (E) Confusion matrix for the multiclass sAwMIL evaluated on the full bag. This probe achieves better accuracy on the *neither*-valued statements.

5.1 Veracity Directions

To show that the truthfulness and falsehood directions are not simple opposites, we examine the *is-true* and *is-false* directions identified by both the multiclass SVM probe and the sAwMIL probe.

In Fig. 3.A, we see how different these directions are by computing their cosine similarity. If the two were perfectly opposite, we would expect a cosine similarity of approximately -1 . However, this is not the case. The better-performing and more generalizable sAwMIL probe yields directions that are *less opposed*, suggesting that LLMs encode true and false as related rather than strictly polar concepts. In Fig. 3.B, we compare the predicted scores obtained by projecting statements onto the *is-true* and *is-false* directions.¹¹ If the two directions were perfectly bidirectional, we would expect a strong negative Spearman’s correlation ($\tau \approx -1$), indicating that high scores along *is-true* correspond to low scores along the *is-false*. Again, this is not observed, particularly for the sAwMIL probe, where predictions are less inversely correlated.

Finally, when forming a matrix $V = [\overrightarrow{\text{is-true}}, \overrightarrow{\text{is-false}}]$, we find that its rank is 2 for both SVM and sAwMIL. The effective rank is 1.73 with standard error of 0.012 for SVM and 1.93 ± 0.004 for sAwMIL. If the two directions were linear combinations of one another (i.e., lying on a single axis), the rank would be 1 and the effective rank would be closer to 1. Together with the results in Fig. 3, these findings indicate that the *is-true* and *is-false* directions are not strict opposites but instead

¹¹Here, we do not consider *neither*-valued statements.

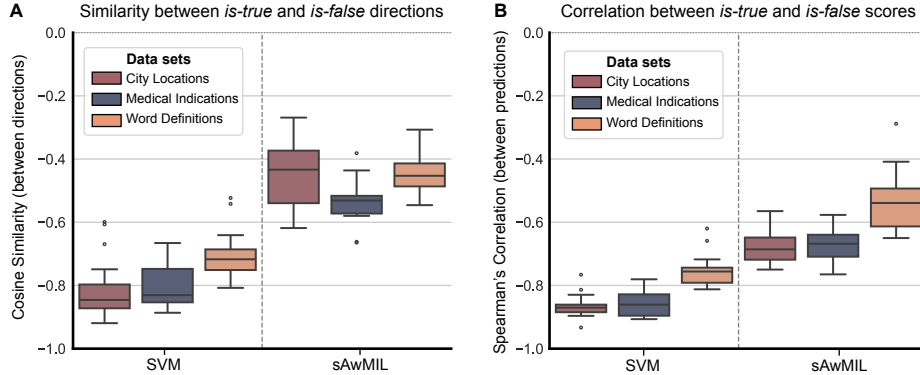


Figure 3: Similarity between *is-true* and *is-false* directions across datasets and probes (values extracted from the best-performing layer and averaged across 16 LLMs). (A) Cosine similarity between the two directions. If the two directions were perfectly opposite, the cosine similarity would be closer to -1 , indicating a single bidirectional axis of truth and falsehood. Instead, all probes exhibit lesser opposition, with the better-performing and more generalizable probe (sAwMIL) showing a smaller angle between *is-true* and *is-false*. (B) Spearman correlation between scores predicted along the two directions (evaluated only on true and false statements). Perfectly bidirectional representations would yield a strong negative correlation: a high score on *is-true* implies a low score in *is-false*. However, the correlation does not approach -1 , indicating partial rather than inverse coupling. Together, the two panels show that *is-true* and *is-false* directions share some structure but are not strict opposites, spanning a multidimensional subspace rather than a single axis.

span a low-dimensional subspace capturing shared yet distinct representational components. This partial overlap may arise because both directions encode aspects of the same underlying concept (e.g., *factuality*) and the relational structure between objects expressed in factual statements (e.g., “The city of Riga is in Latvia”).

6 Conclusion

In this work, we critically examine popular methods for probing the veracity of large language models (LLMs) and show that these probing methods fail to learn reliable and transferable veracity patterns. To address the flaws, we introduce sAwMIL, a multiclass linear probe that combines Multiple Instance Learning with Conformal Prediction Intervals. Unlike prior methods, sAwMIL models veracity using three classes: *true*, *false*, and *neither*. Across sixteen models and three datasets, sAwMIL outperforms existing probes and shows that truthfulness and falsehood are not represented as simple opposites within LLMs, but as directions spanning a subspace. These findings suggest that veracity in language models emerges from distributed and entangled representations rather than a single axis of truthfulness.

Limitations and Future Work. This study focuses on a specific subset of factual statements, namely those that involve the relation between two entities (i.e., city to country, medicine to diseases, noun to noun pairs). As such, it remains unclear how the sAwMIL probe behaves in cases involving multiple relational facts, or in statements describing relations among more than two entities. Extending the evaluation to these higher-order structures represents an important direction for future work.

Furthermore, the *neither*-valued statements used in this study are synthetically generated to simulate cases where LLMs lack knowledge. While these proxies allow for controlled evaluation, they should be further validated to ensure that they accurately capture how models represent unknown and ambiguous information. Future work could explore more natural examples from temporal hold-out datasets or corpora that reflect human uncertainty.

References

- [1] Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17817–17825, 2024.
- [2] Haowen Pan, Xiaozhi Wang, Yixin Cao, Zenglin Shi, Xun Yang, Juanzi Li, and Meng Wang. Precise localization of memories: A fine-grained neuron-level knowledge editing technique for LLMs. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [3] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- [4] David Chanin, Anthony Hunter, and Oana-Maria Camburu. Identifying linear relational concepts in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1524–1535, 2024.
- [5] Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona T Diab, and Marjan Ghazvininejad. A review on language models as knowledge bases. *CoRR*, 2022.
- [6] Michael Townsen Hicks, James Humphries, and Joe Slater. ChatGPT is bullshit. *Ethics and Information Technology*, 26(2):38, 2024.
- [7] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2): 1–55, 2025.
- [8] Kenneth Church. Emerging trends: When can users trust GPT, and when should they intervene? *Natural Language Engineering*, 30(2):417–427, 2024.
- [9] Angus R Williams, Liam Burke-Moore, Ryan Sze-Yin Chan, Florence E Enock, Federico Nanni, Tvesha Sippy, Yi-Ling Chung, Evelina Gabasova, Kobi Hackenburg, and Jonathan Bright. Large language models can consistently generate high-quality content for election disinformation operations. *PloS one*, 20(3): e0317421, 2025.
- [10] Yasin Abbasi Yadkori, Ilja Kuzborskij, András György, and Csaba Szepesvari. To believe or not to believe your LLM: Iterative prompting for estimating epistemic uncertainty. *Advances in Neural Information Processing Systems*, 37:58077–58117, 2024.
- [11] Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. Sayself: Teaching LLMs to express confidence with self-reflective rationales. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5985–5998, 2024.
- [12] Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
- [13] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, 2022.
- [14] Reid McIlroy-Young, Katrina Brown, Conlan Olson, Linjun Zhang, and Cynthia Dwork. Order-independence without fine tuning. *Advances in Neural Information Processing Systems*, 37:72818–72839, 2024.
- [15] Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=aajyHYjjsk>.
- [16] Amos Azaria and Tom Mitchell. The internal state of an LLM knows when it’s lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, 2023.
- [17] Lennart Bürger, Fred A Hamprecht, and Boaz Nadler. Truth is universal: Robust detection of lies in LLMs. *Advances in Neural Information Processing Systems*, 37:138393–138431, 2024.
- [18] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ETKGuby0hcs>.

- [19] Razvan C Bunescu and Raymond J Mooney. Multiple instance learning for sparse positive bags. In *Proceedings of the 24th international conference on Machine learning*, pages 105–112, 2007. doi: 10.1145/1273496.1273510.
- [20] Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33:3581–3591, 2020.
- [21] Anastasios Nikolas Angelopoulos, Stephen Bates, Michael Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=eNdiU_DbM9.
- [22] Benjamin A Levinstein and Daniel A Herrmann. Still no lie detector for language models: Probing empirical and conceptual roadblocks. *Philosophical Studies*, pages 1–27, 2024.
- [23] Abhinav Atla, Rahul Tada, Victor Sheng, and Naveen Singireddy. Sensitivity of different machine learning algorithms to noise. *Journal of Computing Sciences in Colleges*, 26(5):96–103, 2011.
- [24] Shivani Gupta and Atul Gupta. Dealing with noise problem in machine learning data-sets: A systematic review. *Procedia Computer Science*, 161:466–474, 2019.
- [25] Jacqueline Harding. Operationalising representation in natural language processing. *British Journal for the Philosophy of Science*, 2023. doi: 10.1086/728685.
- [26] Daniel A Herrmann and Benjamin A Levinstein. Standards for belief representations in LLMs. *Minds and Machines*, 35(1):1–25, 2025. doi: 10.1007/s11023-024-09709-6.
- [27] Sara Dolnicar and Bettina Grün. Including Don’t know answer options in brand image surveys improves data quality. *International Journal of Market Research*, 56(1):33–50, 2014.
- [28] Marc-André Carboneau, Veronika Cheplygina, Eric Granger, and Ghyslaine Gagnon. Multiple instance learning: A survey of problem characteristics and applications. *Pattern Recognition*, 77:329–353, 2018.
- [29] Ulf Johansson, Henrik Linusson, Tuve Löfström, and Henrik Boström. Model-agnostic nonconformity functions for conformal classification. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 2072–2079. IEEE, 2017.
- [30] Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large language models. *CoRR*, 2024.
- [31] David S Wishart, Yannick D Feunang, An C Guo, Elvis J Lo, Ana Marcu, Jason R Grant, Tanvir Sajed, Daniel Johnson, Carin Li, Zinat Sayeeda, et al. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic acids research*, 46(D1):D1074–D1082, 2018.
- [32] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd edition, 2025. URL <https://web.stanford.edu/~jurafsky/slp3/>. Online manuscript released January 12, 2025.
- [33] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- [34] Daniel A Herrmann and Benjamin A Levinstein. Standards for belief representations in llms. *Minds and Machines*, 35(1):5, 2024.
- [35] John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, 2019.
- [36] Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing Systems*, 37:136037–136083, 2025.
- [37] Yanai Elazar and Yoav Goldberg. Adversarial removal of demographic attributes from text data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21, 2018.

- [38] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A top-down approach to AI transparency. *arXiv preprint arXiv:2310.01405*, 2023. doi: 10.48550/arXiv.2310.01405.
- [39] Anastasios N Angelopoulos, Stephen Bates, et al. Conformal prediction: A gentle introduction. *Foundations and Trends® in Machine Learning*, 16(4):494–591, 2023.
- [40] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge University Press, 2020. ISBN 978-1-108-45514-5.
- [41] Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching. *arXiv preprint arXiv:2404.15255*, 2024.

A Notations and Abbreviations

Table 2: Notations used throughout the paper. Symbols are grouped by category: model definitions, inputs and datasets, internal representations, veracity distributions, and intervention-related symbols.

Symbol	Description	Shape / Notes
\mathcal{M}	Large language model	
$K_{\mathcal{M}}$	Internal probabilistic knowledge of the model \mathcal{M}	
\mathcal{V}	Vocabulary of the model \mathcal{M} , consists of tokens	$[\tau_1, \dots, \tau_{ \mathcal{V} }] \in \mathcal{V}$
$P_{\mathcal{M}}(\tau \mathbf{x})$	Output of \mathcal{M} : a conditional probability distribution on tokens	
Inputs and Datasets		
\mathcal{D}	Dataset of statements	$\mathcal{D}_{train} \cup \mathcal{D}_{test} = \mathcal{D}$
\mathbf{x}	Input token sequence, e.g., “The city of Riga is in Latvia.”	$L = \mathbf{x} $
\mathbf{x}^p	Pre-actualized part of a statement, e.g., “The city of Riga is in”	
\mathbf{x}^a	Actualized part of a statement, e.g., “Latvia.”	
\mathbf{r}	Random sequence with length $ \mathbf{x}^a $	$ \mathbf{r} = \mathbf{x}^a $
y	Veracity label assigned to \mathbf{x}	$y \in Z$
ϕ	A statement evaluated for veracity	
\mathcal{T}_S	Transition matrix for n-gram generation	See Eq. 6
Internal Representations		
d	Size of the hidden representation (of a decoder)	
$h_i(\mathbf{x})$	Activations after the i th decoder	$\mathbb{R}^{L \times d}$
$h_i(\mathbf{x})_{[j]}$	Activation of the token at index j after i th decoder	$\mathbb{R}^{1 \times d}$ and $j \in \{1 \dots L\}$
$h_i(\mathbf{x})_{[n:m]}$	Activations from n th to m th tokens after i th decoder	$\mathbb{R}^{(m-n) \times d}$
Veracity, Probes and Distributions		
Z	Set of veracity labels, e.g., $\{true, false, neither\}$	
$G_{\mathcal{D}}(z \mathbf{x})$	Distribution of veracity labels $z \in Z$ in a dataset \mathcal{D}	
$G_{\mathcal{M}}(z \mathbf{x})$	Distribution of veracity labels $z \in Z$ in the model \mathcal{M}	
g_i	Veracity probe trained on activations of the i th decoder	$g_i : h_i(\mathbf{x}) \mapsto G_{\mathcal{M}}$
\vec{v}_i	Linear direction extracted from the probe g_i	$\mathbb{R}^{1 \times d}$
η	Balancing hyperparameter for sAwMIL	$\eta \in (0, 1)$
\mathbf{m}	sAwMIL’s intra-bag labels (i.e., labels per-token in each \mathbf{x})	$\mathbf{m} \in \{0, 1\}^L$, $L = \mathbf{x} $
Interventions and Effects (Sec. I.1.2)		
I_i^+	Modified representation of $h_i(\mathbf{x}^a)$ after adding $+\vec{v}_i$	
I_i^-	Modified representation of $h_i(\mathbf{x}^a)$ after subtracting $-\vec{v}_i$	
ΔI_i^+	Change in $P_{\mathcal{M}}$ for \mathbf{x}^a after I_i^+ intervention	
ΔI_i^-	Change in $P_{\mathcal{M}}$ for \mathbf{x}^a after I_i^- intervention	
$s_j(\mathbf{x})$	Per-statement success of the intervention; indicator function	See Eq. 20
$\Delta_{correct}$	Total change in $P_{\mathcal{M}}$ for \mathbf{x}^a after both I_i^+ and I_i^- interventions	
Δ_{random}	Total change in $P_{\mathcal{M}}$ for \mathbf{r} random tokens after both interventions	
$\mathbb{E}[\Delta_{correct}]$	Average probability difference across all statements	
\bar{d}_i	Indicator of the dominant direction for the i -th decoder	$\bar{d}_i \in \{-1, 1\}$

Table 3: Abbreviations and naming conventions used throughout this paper.

Abbreviation	Full Form	Description
LLM	Large Language Model	
SIL	Single-Instance Learning	Probes trained on one embedding per example
MIL	Multiple-Instance Learning	Probes trained on multiple embeddings per example
SVM	Support Vector Machine	Type of a classifier
DPO	Direct Preference Optimization	LLM finetuning method
RLHF	Reinforcement Learning from Human Feedback	LLM finetuning method
CP	Conformal Prediction Intervals	Uncertainty calibration method
MD+CP	Mean-Difference with Conformal Prediction Intervals	MD probe with abstention via conformal intervals
sAwMIL	Sparse Aware MIL probe	Multiclass probe handling unknowns
MCC	Matthews Correlation Coefficient	Multiclass performance measure (see Eq. 14)
W-MCC	Weighted-MCC	Using Acceptance Rate as a weight

B Assumptions

Tab. 4 provides an overview of the flawed assumptions in the recent probing methods.

Table 4: **Overview of flawed assumptions in recent methods that probe veracity, their impact on reliability, and our corrective strategies.** Probes that do not account for these issues may lead to biased or unreliable findings.

Flawed Assumptions	Why It Matters	Our Solution/Approach
Truth and falsehood are bidirectional.	There is no conclusive evidence that LLMs treat truth and falsehood as one continuous bidirectional concept. It is more likely that there exist three separate concepts: <code>is-true</code> , <code>is-false</code> , and <code>is-neither</code> ; and they have their own distinct mechanisms.	<code>sAwMIL</code> is a multiclass probe that treats “ <i>true</i> ,” “ <i>false</i> ,” and “ <i>neither</i> ” as separate categories.
LLMs capture and retain everything we know.	We do not know what LLMs have been exposed to during training. Consequently, linear probes that assume every fact in a data set is stored within the LLM are prone to systematic errors in their predictions. We must distinguish between what the LLM <i>actually</i> retains and what <i>we</i> know to be true or false. If a statement \mathbf{x} is unknown to the LLM, it is neither true nor false. In such cases, passing $\langle h_i(\mathbf{x}), \text{true} \rangle$ or $\langle h_i(\mathbf{x}), \text{false} \rangle$ to the probe during training introduces error.	<code>sAwMIL</code> is a linear probe that identifies samples with high support before fitting the linear separator.
All veracity probes provide calibrated probabilities.	Probes such as SVM or mean-group difference classifiers often make a prediction based on the sign (w.r.t. the separation hyperplane). We cannot use these scores to evaluate certainty around the predictions. In other words, these probes are rarely calibrated.	<code>sAwMIL</code> integrates conformal prediction to quantify uncertainty and produce statistically valid prediction regions.
Every token (or statement) is either true or false.	Not every token or sentence expresses a complete factual claim. We should be able to create probes that refrain from making predictions when there is insufficient support.	Instead of training probes to distinguish between true and false statements, <code>sAwMIL</code> is a multiclass classifier that separates statements into “ <i>true</i> ,” “ <i>false</i> ,” and “ <i>neither</i> .”
We know <i>a priori</i> where to look for veracity-related signals.	Most existing probes assume that the last token of a statement has all the information about the veracity.	By using multiple-instance learning, <code>sAwMIL</code> is able to select parts of the input that have the most information about veracity.

C Data Sets

We introduce three new data sets: *City Locations*, *Medical Indications*, and *Word Definitions*. Each dataset consists of statements that are factually *true*, factually *false*, or *neither*. These datasets contain both affirmative and negated statements. An example of a false negated statement is “Guaifenesin is **not** indicated for the treatment of coughs”, and an example of the true affirmative statement is “Shouter is a type of a communicator.”

Data Splits. We split each data set into train, calibration, and test sets using *approximately* 55/20/25 ratios (see Supplementary Tab. 5). We ensure that the objects mentioned in statements are exclusive to the split. For example, if Singapore is mentioned in a statement of the training set, all the statements with Singapore are moved to the training split.

Table 5: Dataset splits. The number of statements per split. In the brackets, we specify the fraction of the total number of statements.

Dataset	Train	Calibration	Test	Total
City Locations	3999 (.55)	1398 (.19)	1855 (.26)	7252 (1.00)
Medical Indications	3849 (.56)	1327 (.19)	1727 (.25)	6903 (1.00)
Definitions	4717 (.55)	1628 (.19)	2155 (.25)	6500 (1.00)

C.1 ‘Neither’ Statements

Since we do not have access to the training data sets of LLMs, we cannot validate whether LLMs retained information about specific facts or entities. That is, we do not know the composition of the internal knowledge $K_{\mathcal{M}}$ of an LLM. Hence, we cannot be certain about what each LLM can (and cannot) verify. To overcome this issue, we create *neither* statements with *synthetic entities*—i.e., entities that do not exist in the real world or fictional works. The *neither* statements are the ones whose value cannot be determined at present (e.g., due to lack of information).

Generation of ‘Neither’ Statements

We use synthetic names to generate *neither*-type statements. For example, “The city of *Staakess* is located in *Soldovadago*” mentions a town and a country that do not exist. From the point of view of an LLM, these statements should be considered *neither-true-nor-false*, as LLMs could not have learned anything about these.

To generate the *neither* statements, we use the Markov-Chain technique [32]. Given a set of existing words $[w_1, w_2 \dots, w_n] \in S$, we break each word w_i into n -grams. For instance, we break “ability” into the following 2-grams: [start]a ab bi il li it ty y[end]. We then compute a transition matrix \mathcal{T}_S , which provides the probability of transitioning from the n -gram i to n -gram j is given by:

$$\mathcal{T}_S(j | i) = \frac{\text{count}(i \rightarrow j)}{\sum_x \text{count}(i \rightarrow x)} \quad (6)$$

In addition, we use \mathcal{T}_S to sample new synthetic words that follow the n -gram distribution of words in S . In our experiments, we use 3-grams for most entities, except for country names, which we generate with 2-grams. We use the `namemaker`¹² package that implements a Markov-Chain word generator.

C.2 Data Selection and Processing

Next, we provide details on the source and processing steps for each dataset.

¹²github.com/Rickmsd/namemaker

City Locations

The *City Locations* dataset is based on the GeoNames¹³ database. GeoNamesCache¹⁴ is a Python package that interacts with the GeoNames API. We use the following criteria to select a $\langle \text{city}, \text{country} \rangle$ pair:

1. The population of the city is at least 30,000.
2. The city has an associated country. If a city name is associated with multiple countries, we include the $\langle \text{city}, \text{country} \rangle$ pair for each country. We exclude all cities that have “Antarctica” as a location or a country.

Since the resulting set of $\langle \text{city}, \text{correct country} \rangle$ pairs is relatively large. We reduce the number of pairs by downsampling. In total, we select 1,400 unique city names: 700 cities with the highest populations, and 700 cities randomly sampled from the rest of the names.

Statement Structure. For each $\langle \text{city}, \text{correct country} \rangle$ pair, we create statements of the form:

The city of [city] is (not) located in [country].

If a city name already contained a word “city” (e.g., “Guatemala City”), we do not start a sentence with “The city of.” We also sample $\langle \text{city}, \text{incorrect country} \rangle$ pairs, and generate statements according to the template above.

Synthetic Entities. We use the technique described in Supplementary Sec. C.1 to generate synthetic city and country names. To generate synthetic *city* names, we collect all the city names in our data set (including those that we did not include) and input them to *namemaker* (with n -gram length of 3). We generate 500 synthetic city names. We validate these synthetic names in two stages:

1. We check whether a synthetic name exists in the GeoNames database by looking for matches in the *name* and *alternative name* fields. We keep 310 cities after this first stage.
2. We use Google Search to validate that each synthetic city name does not exist via the following prompt: “city [city name]”. If the search result returns a city with 1-2 character difference, we remove the synthetic name from the list. We keep 219 cities after this second stage.

For the synthetic country names, we collect all the country names and input them to *namemaker* (with n -gram length of 2). We generate 250 synthetic names and validate them using the workflow described in the previous paragraph. We keep 238 country names after the first stage, and 138 after the second stage. The Google Search prompt is: “country [country name]”. With 25% probability, we add a prefix or suffix to the synthetic country name. The list of prefixes and suffixes include “Island,” “Republic of,” “Kingdom,” “West,” “East,” “North,” “South,” and “Land.” Finally, we randomly match each synthetic name to the name of a synthetic country.

Medical Indications

The *Medical Indications* dataset is based on the DrugBank (version 5.1.12) [31]. We obtain access to the DrugBank on October 4th, 2024, via the academic license (for research purposes only). Our GitHub and Zenodo repositories do not contain the raw data from the DrugBank, but the reader can apply for the academic license.¹⁵ We extract 2 fields from this knowledge base:

1. **Name**, which specifies the official name of the drug or the chemical (e.g., Lepirudin).
2. **Indication**, which is a text field that describes the indication of the drug. If this field consists of multiple sentences, we keep only the first sentence (e.g., “Lepirudin is indicated for anticoagulation in adult patients with acute coronary syndromes (ACS) such as unstable angina and acute myocardial infarction without ST elevation.”)

To extract diseases and conditions from the *Indications* field, we use two named entity recognition (NER) models:

¹³ [geonames.org](https://www.geonames.org)

¹⁴ pypi.org/project/geonamescache/

¹⁵ Here is the link to the DrugBank’s academic license: <https://go.drugbank.com/releases/5-1-12>

1. SciSpacy’s `en_ner_bc5cdr_md` model¹⁶ for the biomedical term annotations
2. BioBERT-based NER¹⁷ for disease annotations

We input the “Indication” text to both models. The disease/condition terms are extracted only if both models mark it as a disease or condition. For example, for “Lepirudin is indicated for anticoagulation in adult patients with acute coronary syndromes (ACS),” the SciSpacy model marks *coronary syndromes* as a disease, but BioBERT does not. Thus, we do not add it to Lepirudin’s disease/condition list. Similarly, we remove the abbreviation if the disease list contains the full name *and* its abbreviation, such as [acute coronary syndromes, ACS].

We further validate the drug names via SciSpacy model, and keep the name only if it is marked as CHEMICAL. Otherwise, we remove the drug from our dataset. Finally, if the disease list (for a given drug) is empty after the preprocessing, we remove the drug from our data set.

Additionally, we use `wordfreq`¹⁸ package to check whether the name of the drug or the name of the indication appears in widely used corpora (e.g., Wikipedia or Books dataset). In other words, we remove the pair if either the drug name or the indication has a Zipf’s frequency of 0 – i.e., the word does not appear in any of the `wordfreq` corporas.

Statement Structure. For each $\langle \text{drug}, \text{correct disease} \rangle$ pair, we create statements of the form:

[drug] is (not) indicated for the treatment of [disease/condition].

We also sample the $\langle \text{drug}, \text{incorrect disease} \rangle$ pairs. We ensure that the “incorrect disease” did not share any words with the diseases in the correct list.

Synthetic Entities. To generate synthetic drug names and disease names, we use the approach described in Supplementary Sec. C.1 (with n -gram length of 3). We generate 500 synthetic drug names. We validate these synthetic names in two stages:

1. We pass each generated name through SciSpacy model and remove the ones marked as CHEMICAL. We keep 315 name after this first stage.
2. We use Google Search to validate that each drug name does not exist via the prompt “medicine [drug name].” If the search result returned a drug with 1-2 character difference, we remove it from the dataset. We keep 243 names after this second stage.

We generate 200 disease names and check whether they exist in our list of diseases. We keep 181 names after this first stage. Next, we use Google Search with the prompt “disease [disease/condition name].” We keep 131 disease names after this second stage. Finally, we randomly match synthetic drug names to synthetic disease names to generate *neither*-type statements.

Word Definitions

The *Word Definitions* dataset is based on the sample data from WordsAPI¹⁹ database. Sample data is publicly available and contains 10% of randomly sampled words from the database.²⁰

For each word in the sample, we keep the ones that satisfy the following criteria:

1. The word is a noun.
2. The word has at least one definition in the *definition* field.
3. The word has at least one of the following fields: *synonym*, *typeOf*, or *instanceOf*.

Statement Structure. Depending on the specified field (i.e., *synonym*, *typeOf*, *instanceOf*), we generate three types of statements:

1. “[word] is (not) [instanceOf].”

¹⁶allenai.github.io/scispacy/

¹⁷[alvaroaalon2/biobert-diseases_ner](https://github.com/alvaroaalon2/biobert-diseases_ner)

¹⁸pypi.org/project/wordfreq

¹⁹[WordsAPI.com](https://wordsapi.com)

²⁰We do not provide a copy of the sample in our GitHub or Zenodo repositories.

2. “[word] is (not) a type of [typeOf].”
3. “[word] is (not) a synonym of [synonym].”

Before inserting a word from *synonym*, *typeOf*, *instanceOf* fields into a corresponding spot, we check which article goes before ‘a’ or ‘an’. When possible, we change words into singular forms. To do so, we use the `inflect` package,²¹

Synthetic Entities. To generate synthetic entities, we use the approach described in Supplementary Sec. C.1 (with *n*-gram length of 3). We generate four categories of synthetic entities:

1. Words that go at the beginning of each statement: We use all the words we have in the dataset.
2. Types: We use all the words from the *typeOf* field for the Markov-Chain generation.
3. Synonyms: We use all the words from the *synonym* field.
4. Instances: We use words from the *instanceOf* field.

We generate 1,000 synthetic words for each of the four categories. We validate the non-existence of words. We use the `english_words` package²² to check whether a word exists in “GNU Collaborative International Dictionary of English 0.53,” or `web2` word list. Furthermore, we check whether there is a word in the *words* list of the `nltk` package.²³ After this stage, we end up with 3,305 words. Finally, we randomly sample pairs of ⟨word, property⟩, where the property is a type, instance, or synonym.

²¹pypi.org/project/inflect/

²²pypi.org/project/english-words/

²³pypi.org/project/nltk/

D Selection of Large Language Models

In this section, we provide an overview of the large language models used in our experiments. Supplementary Tab. 6 provides a list of all the 16 models.

We use default models—i.e., the ones that were pre-trained on general tasks. We also use chat models that have been fine-tuned on instruction- and chat-like interactions. Every default model in our selection has a corresponding chat-based model. We also add two extra chat-tuned Llama models that are specifically fine-tuned on biomedical data. Further, we do not use full official model names but use short names along with a version, such as “chat” or “default”. For example, Llama-3.2 (chat) refers to the Llama-3.2-3b-Instruct model.

Table 6: **List of LLMs used in our experiments.** We provide the official names of the models in the HuggingFace repository. Further, we provide the *type* of the model: default stands for the pre-trained models, and ‘chat’ stands for the chat- or instruction-tuned version of the models. Finally, we provide the number of decoders, the number of parameters, the release date, and the source of the model. These models are publicly available through HuggingFace [33].

Official Model Name	Type	# Decoders	# Parameters	Release Date	Source
Gemma-7b	Default	28	8.54 B	Feb 21, 2024	Google
Gemma-2-9b	Default	26	9.24 B	Jun 27, 2024	Google
Llama-3-8b	Default	32	8.03 B	Jul 23, 2024	Meta
Llama-3.2-3b	Default	28	3.21 B	Sep 25, 2024	Meta
Mistral-7B-v0.3	Default	32	7.25 B	May 22, 2024	Mistral AI
Qwen2.5-7B	Default	28	7.62 B	Sep 19, 2024	Alibaba Cloud
Qwen2.5-14B	Default	38	14.80 B	Sep 19, 2024	Alibaba Cloud
Gemma-7b-it	Chat	28	8.54 B	Feb 21, 2024	Google
Gemma-2-9b-it	Chat	26	9.24 B	Jul 27, 2024	Google
Llama-3.2-3b-Instruct	Chat	28	3.21 B	Sep 25, 2024	Meta
Llama-3.1-8b-Instruct	Chat	32	8.03 B	Jul 23, 2024	Meta
Llama3-Med42-8B	Chat	32	8.03 B	Aug 12, 2024	M42 Health
Bio-Medical-Llama-3-8B	Chat	32	8.03 B	Aug 11, 2024	Contact Doctor
Mistral-7B-Instruct-v0.3	Chat	32	7.25 B	May 22, 2024	Mistral AI
Qwen 2.5-7B-Instruct	Chat	28	7.62 B	Aug 18, 2024	Alibaba Cloud
Qwen 2.5-14B-Instruct	Chat	38	14.80 B	Aug 18, 2024	Alibaba Cloud

E Criteria for Validating Veracity Probe

Table 7: **Validity criteria for representation-based probes.** If satisfied, these criteria serve as validation that g_i indeed captures signals associated with veracity Z . Here, we provide a formal definition of each criterion, along with the implications of satisfying the criterion. Finally, we provide the list of similar criteria and concepts used in the literature.

Criteria	Definition	If Satisfied	Similar Concepts
Correlation	A probe g_i trained on $\langle h_i(\mathbf{x}), y \rangle \in \mathcal{D}_{\text{train}}$ should perform well (i.e., have high predictive accuracy) on $\mathcal{D}_{\text{test}}$, assuming the same input and label distributions.	\mathcal{M} encodes information correlated with veracity.	Information [25], Accuracy [26]
Generalization	A probe g_i trained on $\langle h_i(\mathbf{x}), y \rangle \in \mathcal{D}_{\text{train}}$ should have high predictive accuracy on data from different domains.	\mathcal{M} has a universal activation pattern correlated with veracity (see Fig. 2.B and 13).	Generalization as defined by Bürger et al. [17], Uniformity [26]
Selectivity	A probe g_i trained on $\langle h_i(\mathbf{x}), y \rangle \in \mathcal{D}_{\text{train}}$ should not assign <i>true</i> or <i>false</i> labels to samples where truthfulness is absent or undefined.	\mathcal{M} has a distinct mechanism that correlates exclusively with veracity.	Misrepresentation as defined by Harding [25], Control Task [35]
Manipulation	Modifying $h_i(\mathbf{x})$ along \vec{v}_i should systematically alter $P_{\mathcal{M}}(\tau \mid \mathbf{x})$ for tokens τ related to the veracity property Z .	\mathcal{M} has a <i>linear</i> mechanism to track veracity and uses it to compute the output $P_{\mathcal{M}}(\tau \mid \mathbf{x})$ (see Fig. 14).	Use [25], Addition [36], Intervention [3]
Locality	Modifying $h_i(\mathbf{x})$ along \vec{v}_i should not significantly alter $P_{\mathcal{M}}(r \mid \mathbf{x})$ for random tokens r that are unrelated to Z .	\mathcal{M} maintains a separate mechanism that tracks veracity, without being confused with other concepts.	Misrepresentation [25], Leakage [37]

Researchers have proposed criteria to measure the validity of veracity probes [25, 26, 38]. We aggregate these into five major categories and provide an overview in Supplementary Tab. 7. We propose to evaluate a probe g_i along the following criteria:

- (i) **Correlation.** The probe, trained to predict a veracity property $\{true, false\} \in Z$, should achieve high predictive accuracy on unseen samples that possess this property, i.e., on samples from $\mathcal{D}_{\text{test}}$. When the criterion is satisfied, the i -th decoder embeds the information about Z to some degree. We cannot rule out the fact that it captures proxies associated with Z .
- (ii) **Generalization** extends *Correlation* by requiring that the probe generalizes beyond the data set it was trained on. The probe should have high predictive accuracy on samples that have veracity Z , but have different phrasing or come from different domains. For example, if a probe is trained to identify neural activation patterns associated with veracity on statements related to ecology, this probe should have similar predictive accuracy on statements related to biology.
- (iii) **Selectivity** The probe g_i should avoid classifying statements that are not (or cannot be) *true* or *false*. Hence, the probe g_i should abstain from making predictions on the *neither*-valued statements—i.e., statements that the LLM could not have learned from its training data or that inherently lack any truthfulness or falsehood. Poor selectivity indicates that the probe might capture spurious correlations with unrelated properties.
- (iv) **Manipulation.** We should be able to use the identified direction \vec{v}_i to update $h_i(\mathbf{x})$ and have a predictable change in the distribution of the output tokens $P_{\mathcal{M}}(\tau \mid \mathbf{x})$. Since we focus on *linear* probes, we expect that moving $+c$ units along \vec{v}_i should have an opposite effect on $P_{\mathcal{M}}$ compared to moving $-c$ units.

- (v) **Locality** When asking a question such as “Is X true? Answer yes or no,” the manipulation should primarily influence the generation process related to the “yes” or “no” responses. It should minimally affect unrelated tokens. For example, if a manipulation does not increase the likelihood of the LLM generating “no”, but increases the likelihood of generating tokens such as “elephant”, then the manipulation degrades the LLM’s abilities.

Evaluating a probe according to these criteria allows us to determine how well g_i captures the signals associated with veracity Z and how manipulations (a.k.a. interventions) affect LLM’s output $P_{\mathcal{M}}$. Part of our future work includes adding a new criterion on whether the probe can assess if an LLM can “reason” logically. For example, if \mathcal{M} classifies a statement ϕ_1 as true and another statement ϕ_2 as true, then will \mathcal{M} also classify $\phi_1 \wedge \phi_2$ as true?

Finally, we demonstrate the evaluation results for the **Correlation** and **Selectivity** in the Results section (see Sec. 5) of the manuscript; we further provide evaluation results for the **Generalization**, **Manipulation** and **Locality** in the Supplementary Sec. I.

F Zero-Shot Prompting: Instructions, Veracity Labeling, and Abstention

Instructions in Zero-shot Prompting. In zero-shot prompting, each LLM receives an instruction along with a statement x as input. These instructions outline the task and describe the format of the output. We use zero-shot prompts to evaluate how well LLMs can assess the veracity of a given statement and provide a correct reply. Our zero-prompts are designed as a multiple-choice question, where an LLM is asked to pick a corresponding number.

In our (zero-shot prompting) experiments, we use three different templates. All default models share the same template as displayed in Supplementary Fig. 4A–B. We do not use this template for chat-models, since they support “turn-based” conversations. Thus, for most chat models, we use templates that model user-assistant interaction. We use the template in Supplementary Fig. 4D–E for the chat models that support context prompts. We use the template in Supplementary Fig. 4F for Gemma models since they do not support context prompts.

Additionally, we use three different phrasings of the instructions:

1. **Original instructions** are displayed in Supplementary Fig. 4A and D. In the main paper, we report results based on these instructions.
2. **True-False instructions**, where we change the phrasing of the original instructions by replacing the words “correct” and “incorrect” to “true” and “false”. These are displayed in Supplementary Fig. 4B and E.
3. **Shuffled instructions**, where we change the original instructions by shuffling choice identifiers. These are displayed in Supplementary Fig. 4C and F.

Supplementary Figures 5–7 depict the Weighted Matthew’s Correlation Coefficient (W-MCC) performance of zero-shot prompting across 3 datasets and 16 LLMs for each template type and instruction type. In some cases, minor changes to the instructions can have a significant impact on performance. Overall, the Qwen and Gemma-2 models provide the best performance compared to other LLMs. However, Qwen models seem to be more sensitive to the phrasing of the prompt.

A Original Instructions (Default Model)	B True-False Instructions (Default Model)	C Shuffled Instructions (Default Model)
<p>Question: Is the following statement correct? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 1. The statement is correct. 2. The statement is incorrect. 3. I do not have sufficient knowledge. 4. The statement is too ambiguous. 5. All of the above. 6. None of the above. <p>Please respond with the corresponding number. The final answer is</p>	<p>Question: Is the following statement true? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 1. The statement is true. 2. The statement is false. 3. I do not have sufficient knowledge. 4. The statement is too ambiguous. 5. All of the above. 6. None of the above. <p>Please respond with the corresponding number. The final answer is</p>	<p>Question: Is the following statement correct? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 2. The statement is correct. 6. The statement is incorrect. 4. I do not have sufficient knowledge. 5. The statement is too ambiguous. 1. All of the above. 3. None of the above. <p>Please respond with the corresponding number. The final answer is</p>
D Original Instructions with context (Chat Model)	E True-False Instructions with context (Chat Model)	F Shuffled Instructions without context (Chat Model)
<p><i>System (Context Prompt)</i></p> <p>You are an expert in fact-checking. Your task is to assist the user by answering questions based on your comprehensive knowledge. Please respond with the corresponding number.</p> <p><i>User</i></p> <p>Question: Is the following statement correct? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 1. The statement is correct. 2. The statement is incorrect. 3. I do not have sufficient knowledge. 4. The statement is too ambiguous. 5. All of the above. 6. None of the above. <p><i>Assistant</i></p> <p>The final answer is</p>	<p><i>System (Context Prompt)</i></p> <p>You are an expert in fact-checking. Your task is to assist the user by answering questions based on your comprehensive knowledge. Please respond with the corresponding number.</p> <p><i>User</i></p> <p>Question: Is the following statement true? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 1. The statement is true. 2. The statement is false. 3. I do not have sufficient knowledge. 4. The statement is too ambiguous. 5. All of the above. 6. None of the above. <p><i>Assistant</i></p> <p>The final answer is</p>	<p><i>User</i></p> <p>Question: Is the following statement correct? The city of Agadir is located in Morocco.</p> <p>Select one of the following options:</p> <ol style="list-style-type: none"> 2. The statement is correct. 6. The statement is incorrect. 4. I do not have sufficient knowledge. 5. The statement is too ambiguous. 1. All of the above. 3. None of the above. <p><i>Assistant</i></p> <p>The final answer is</p>

Figure 4: **Zero-shot prompt templates.** We use these templates in our experiments. **Panel A–B:** The prompts for the default models. **Panel D–F:** Examples of prompts used for the chat models. Note that chat models like Gemma do not have a context (or system) prompt; hence, we provide instructions in the first message (see Panel F). In the main manuscript, we report the performance over the *original instructions* – i.e., instructions in Panels A and D. For the chat LLMs without the context prompt, we apply the two-message template in Panel F, but use the *original instructions*. (Side note: The statement used in these examples is factually correct.)

From Token Probabilities to Veracity Labels in Zero-shot Prompting. Given the original instructions and the statement x , an LLM outputs token-level probabilities over its vocabulary \mathcal{V} as

$$P_{\mathcal{M}}(\tau \mid \text{instruction} \wedge x) \text{ with } \sum_{\tau \in \mathcal{V}} P_{\mathcal{M}}(\tau \mid \text{instruction} \wedge x) = 1.$$

We are interested in the probabilities of the tokens that correspond to the multiple choices – i.e., numbers 1–6 in any of the panels in Supplementary Fig. 4. We denote tokens associated with these numbers as: [1], [2], [3], etc. We map these token-level probabilities $P_{\mathcal{M}}$ into the veracity-label probabilities $G_{\mathcal{M}}$ as follows:

$$G_{\mathcal{M}}(\text{true} \mid x) = P_{\mathcal{M}}([1] \mid \text{instruction} \wedge x) \quad (7)$$

$$G_{\mathcal{M}}(\text{false} \mid x) = P_{\mathcal{M}}([2] \mid \text{instruction} \wedge x) \quad (8)$$

$$G_{\mathcal{M}}(\text{neither} \mid x) = P_{\mathcal{M}}([3] \mid \text{instruction} \wedge x) + P_{\mathcal{M}}([4] \mid \text{instruction} \wedge x) \quad (9)$$

Abstention in Zero-shot Prompting. We include options [5] and [6] to check the “sanity” of the model \mathcal{M} . For example, option #5 in Supplementary Fig. 4A suggests that a statement x is true, false, and ambiguous – all at the same time. If the model assigns most of the probability mass to these tokens, we assume that the model does not follow the instructions. Similarly, if a model assigns most of its probability mass $P_{\mathcal{M}}$ to other tokens in the vocabulary $\{\tau \in \mathcal{V} : \tau \notin \{[1], [2], [3], [4], [5], [6]\}\}$, we also assume that the model did not follow the instructions. Hence, if instructions are not followed, we assume that the model *abstains* from making a prediction, see Eq. 10.

$$G_{\mathcal{M}}(\text{abstain} \mid x) = \sum_{\tau} P_{\mathcal{M}}(\tau \mid \text{instruction} \wedge x), \text{ where } \{\tau \in \mathcal{V} : \tau \notin \{[1], [2], [3], [4]\}\} \quad (10)$$

Note that the zero-shot prompting relies only on the token-level probabilities, i.e., \mathcal{M} ’s output. It does not look at the intermediate hidden representation of x .

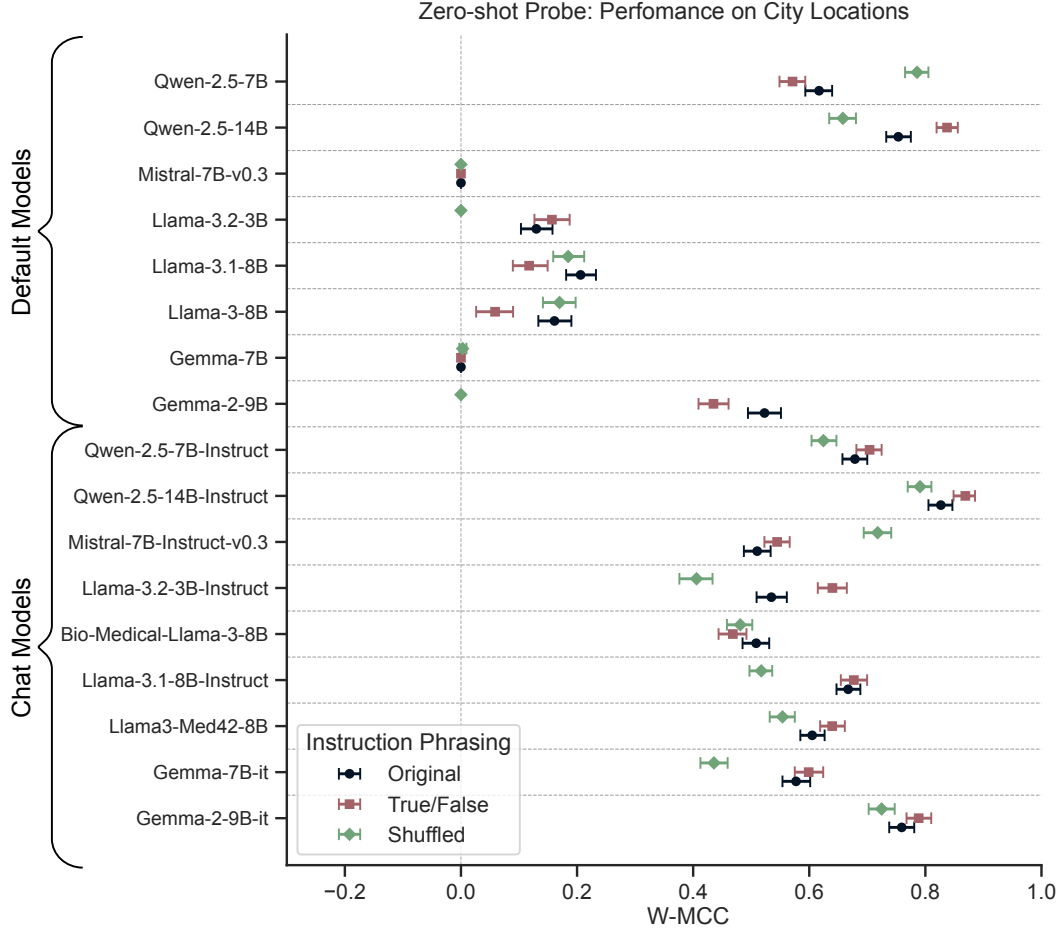


Figure 5: **Performance of zero-shot prompting on the *City Locations* data set across different models and instruction phrasings.** We use the Weighted Matthew's Correlation Coefficient (W-MCC) to quantify the performance. The marker shows the mean value and the error bars show the 95% confidence intervals (based on the bootstrapping with $n = 1,000$ bootstrap samples). Minimal changes to the prompt instructions can skew the performance of zero-shot prompting. Chat models exhibit the highest performance across all instruction phrasings. However, the default Qwen models match the performance of other chat-based models. Shuffled instructions appear to lead to worse performance in chat models. We expected that the phrasings would have only a minor effect on their performance. The default Gemma and Mistral models seem to fail (their performance is around 0).

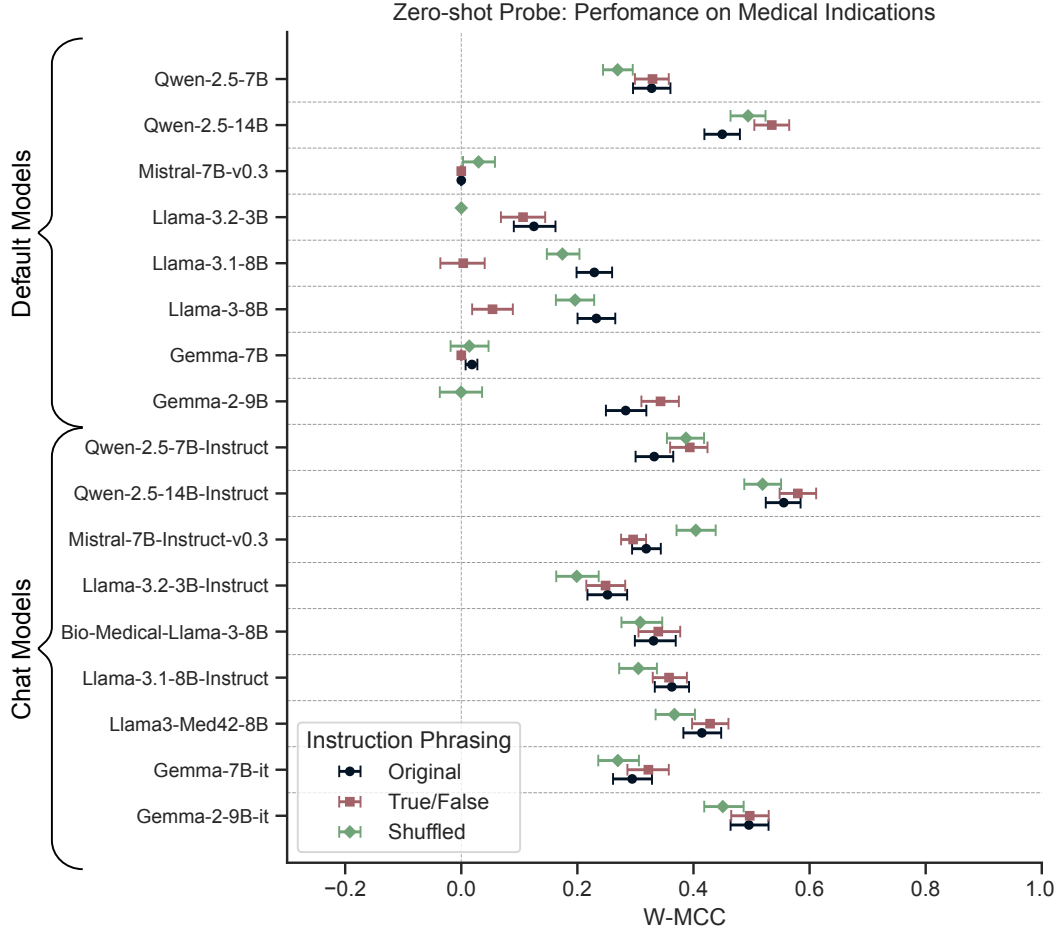


Figure 6: **Performance of zero-shot prompting on the *Medical Indications* data set across different models and instruction phrasings.** We use the Weighted Matthew’s Correlation Coefficient (W-MCC) to quantify performance. The marker shows the mean value and the error bars show the 95% confidence intervals (based on the bootstrapping with $n = 1,000$ bootstrap samples). Minimal changes to the prompt instructions can skew the performance of zero-shot prompting. We observe a slight performance misalignment depending on the instruction phrasing. The best-performing LLMs are the largest chat models: Gemma-2-9b and Qwen-2.5-14b. We expected the biomedical Llama models to outperform on the medical indications dataset.

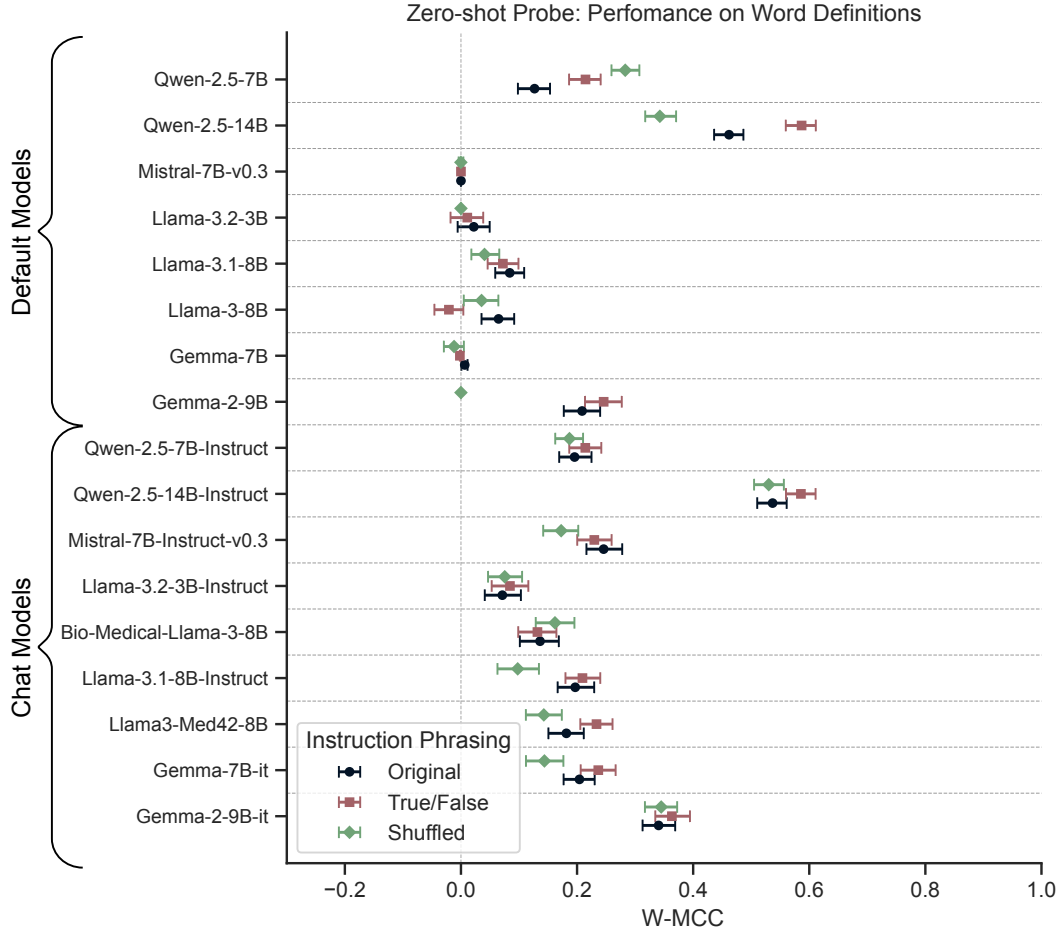


Figure 7: **Performance of the zero-shot prompting on the *Word Definitions* data set across different LLMs and instruction phrasings.** We use the Weighted Matthew's Correlation Coefficient (W-MCC) to quantify the performance: the marker shows the mean value and the error bars show the 95% confidence intervals (based on bootstrapping with $n = 1,000$ bootstrap samples). Minimal changes to the prompt instructions can skew the performance of zero-shot prompting. The overall performance on the *Word Definitions* data set is much lower compared to the performances on the other data sets. Generally, the misalignment in the performance between different instructions is much lower (except for the default Qwen models, where the difference is significant). The largest Qwen-2.5-14b are the top-performing models on this task.

G Conformal Prediction Intervals

In our work, we focus on “split conformal learning” [21], which requires a hold-out (or calibration) data set to compute conformal prediction intervals.

Given a probe g_i , we use a calibration data set \mathcal{D}_{calib} of activation-label pairs $\langle h_i(\mathbf{x}), y \rangle$ to find prediction regions that ensure, for example, that a sample falling within the region is correctly classified 90% of the time. If a prediction falls into the overlapping conformal prediction intervals of two (or more) classes, or if it does not fall within any interval, the probe abstains from making any prediction. We provide pseudocode for the nonconformity functions in Supplementary Alg. 5 and 6.

G.1 Nonconformity score

To identify the conformal intervals, we compute a nonconformity score for each sample in \mathcal{D}_{calib} . For the binary cases (such as mean-difference probe and one-vs-all sAwMIL), we use the binary nonconformity scoring; see Supplementary Alg. 5. It is based on the distance between the prediction and the classifier’s separating hyperplane. In Eq. 11, s is the signed distance of the sample to the separation hyperplane, and y is a ground-truth (or candidate) label:²⁴

$$binaryNC(s, y) = \exp(-y \cdot s), \quad y \in \{-1, 1\} \text{ and } s \in \mathbb{R}. \quad (11)$$

If the sample ends up on the wrong side of the separation hyperplane (e.g., $s > 0$ and $y = -1$), then the nonconformity score in Eq. 11 is high and the candidate label is weakly supported by the model.

For the multiclass sAwMIL, we use the multiclass nonconformity score [29]; see Supplementary Alg. 6. For a given candidate label y , the label is defined in terms of the difference between the predicted probability of the true class and the highest probability among the other classes (with K denoting the total number of classes). Formally, for a candidate label y with predicted probability p_y , we calculate the multiclass nonconformity score with the following function:

$$multiclassNC(\mathbf{p}) = \frac{1 - (p_y - \max_{i \neq y} p_i)}{2} \quad (12)$$

where $\mathbf{p} \in \Delta^{K-1} := \left\{ \mathbf{p} \in \mathbb{R}^K \mid p_i \geq 0, \sum_{i=1}^K p_i = 1 \right\}$ and Δ^{K-1} is a simplex.

In both cases, lower scores in Eq. 11 and Eq. 12 indicate that the candidate label y is strongly supported by the model. In our work, we set $\alpha = 0.1$. Thus, if the nonconformity score of a new sample exceeds the 90th quantile, the probe abstains from prediction (see Supplementary Alg. 6). The addition of conformal intervals enables us to distinguish between cases where the statements originate from different distributions, as compared to those in the calibration data set.

²⁴In this case, labels should be either -1 or 1 . Thus, all samples with label 0 are assigned label -1 .

H More on Representation-based Probing Methods

H.1 Mean-difference Probe with Conformal Prediction Intervals

The mean-difference probe (MD+CP) consists of two components: binary mean-difference classifier (MD) and the conformal prediction intervals (CP).

First, we fit the binary classifier with a linear decision boundary [15]. We use it to separate *true* and *false* statements based on the internal activations h_i . For each pair $\langle x_j, y_j \rangle$, we extract the activation of the last token $h_i(x_j)_{[L]}$ and assemble a set of factually true $\mathcal{X}^+ = \{h_i(x_j)_{[L]} : y_j = \text{true}\}$ and a set of false $\mathcal{X}^- = \{h_i(x_j)_{[L]} : y_j = \text{false}\}$ activations. Here, L is the index of the last token in x . We then compute the means of each set, denoted μ^+ and μ^- , and compute a direction vector:

$$\theta = (\mu^+ - \mu^-) \Sigma^{-1} (\mu^+ - \mu^-)^T \quad (13)$$

In Eq. 13, Σ is a pooled covariance matrix. See Alg. 2 in Supplementary Materials for the detailed pseudo-code.

Second, we augment MD with conformal prediction intervals [39]. Conformal intervals help detect statements that fall outside MD’s high-confidence regions for *true* or *false* classes. We use $\alpha = 0.1$ in our experiments; thus, predictions in the high-confidence regions are guaranteed to be correct at least 90% of the time. Note that we use the *true* and *false* statements from the calibration set to find the conformal prediction intervals. Finally, we test the MD+CP probe using *true*, *false*, and *neither* statements from the test set. *How can the binary MD+CP classifier identify neither statements in addition to true and false statements?* If the MD+CP probe accurately captures the veracity signal, the *neither* statements (from the test set) should fall outside of the conformal prediction intervals. Below, we observe that this is not the case. MD+CP assigns high-confidence scores to the *neither*-valued statements in the *true* or *false* regions.

Supplementary Fig. 8 shows the score distributions²⁵ of MD+CP on the best performing decoder (i.e., 13th) of the default Llama-3-8B model on the *City Locations* data set. There are three distributions: one for *true*, one for *false*, and one for *neither*. The distributions are based on samples from the test set. If MD+CP correctly captures the veracity signal, we expect the distribution of *neither* statements (green bars) to be outside of the conformal prediction interval (i.e., in the gray area). However, this is not the case. Most of the *neither* statements fall within the conformal prediction intervals (i.e., not in the area colored gray) and get labeled as *true* or *false*.

Supplementary Fig. 9A illustrates per-token MD+CP predictions across entire statements. If MD+CP correctly identifies the veracity signal, then (1) it should not assign any labels to the tokens in the pre-actualized parts of statements x^p , and (2) the label should be consistent across the actualized path x^a . Given a statement “The city of Tokyo is in Japan.”, the pre-actualized part is “The city of Tokyo is in” and the actualized part is “Japan.” We observe that MD+CP assigns scores to the pre-actualized tokens (“The city of X”) that fall within the conformal prediction intervals in cases #1–5 (see Supplementary Fig. 9). In case #5, MD+CP assigns a correct prediction at the *period sign* ($p = 0.39$ corresponds to a false label), but the prediction flips at the end of the text, where the *question mark* gets $p = 0.95$ corresponding to the *true* label. Similarly, in the #7 case of Supplementary Fig. 9A, the sentence does not have *any* veracity value (i.e., it is not a factual claim). However, the MD+CP probe assigns high confidence scores to some of its tokens. These findings suggest that MD+CP probe captures proxies or spurious correlations. One cannot use it in real scenarios, where we do not know a priori where the factual claim ends. In contrast, Supplementary Fig. 9 B–D show the per-token predictions for the one-vs-all sAwMIL. These probes correctly identify positions where the veracity is actualized. For example, they do not assign predictions to the non-actualized parts of the statements and only label tokens in the actualized part. Moreover, these probes do not label tokens in cases where veracity is absent (e.g., see case #7 in Supplementary Fig. 9D).

H.2 Multiclass Single-Instance Support Vector Machine

The multiclass sAwMIL probe is a multiple-instance learning (MIL) version of Support Vector Machine (SVM), designed to operate on bags of token representations. To assess whether the MIL formulation

²⁵We use the embedding of the last token to compute the scores in Supplementary Fig. 8: $g_i(h_i(x)) = \theta^T h_i(x)_{[L]} + \beta$. Here, L is the total length of the statement, which is the same as the index of the last token.

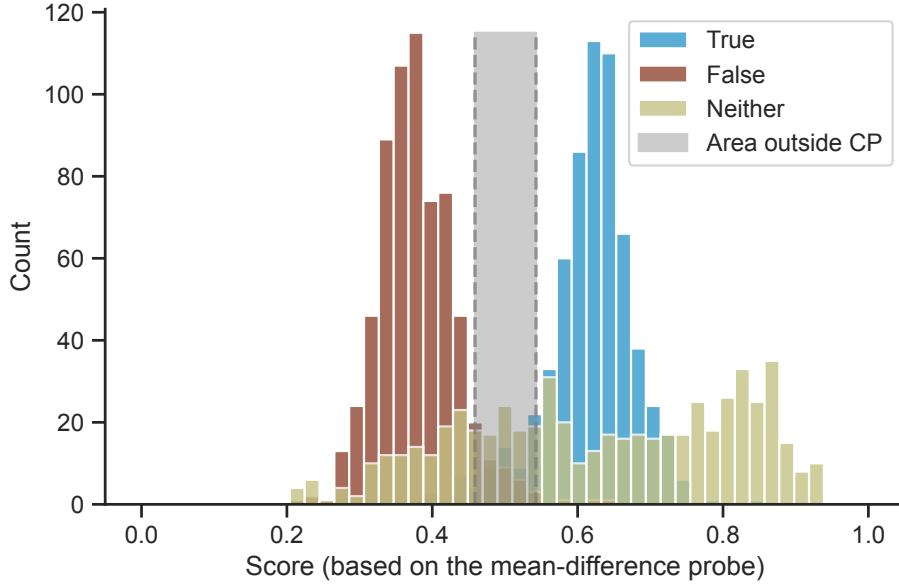


Figure 8: **Score distributions of mean-difference probe with conformal prediction intervals (MD+CP) on the 13th decoder activations of the default Llama-3-8B model for the City Locations dataset.** The probe provides a good separation between *true* and *false* statements. Note, if the MD+CP probe *truly* captures only the veracity signal, we expect the scores for the *neither* statements to fall outside the conformal intervals (i.e., be in the area highlighted with gray color). However, MD+CP assigns high-confidence scores for the *neither*-valued statements, labeling them *true* or *false*. This finding suggests that MD+CP relies on spurious proxies rather than genuine veracity signals.

offers any benefits, we construct a single-instance baseline by training a multiclass SVM on the last token representation only $h_i(\mathbf{x})_{[L]}$. As with multiclass sAwMIL, we first train three one-vs-all probes: *is-true*, *is-false*, and *is-neither*. Then, these one-vs-all classifiers are assembled into a multiclass SVM using the same procedure described in Sec. 3.1.1 and Supplementary Alg. 1. Finally, we augment the multiclass SVM with conformal prediction intervals to provide calibrated estimates, mirroring the multiclass sAwMIL setup.

As before, to evaluate performance, we provide all token representations $h_i(\mathbf{x})$ (not only the last one), where the final prediction is computed based on

$$\hat{g}_i(\mathbf{x}) = \max_{1 \leq j \leq L} g_i(h_i(\mathbf{x})_{[j]}), \text{ where } L = |\mathbf{x}| \text{ (number of tokens in } \mathbf{x}\text{)}.$$

Supplementary Fig. 10 depicts the performance of multiclass sAwMIL vs. multiclass SVM. The performance of multiclass SVM is closer to the performance of multiclass sAwMIL (as compared to MD+CP probe in Fig. 10.A). However, multiclass sAwMIL still outperforms the multiclass single-instance SVM in 46 out of 48 cases (= 16 LLMs \times 3 data sets), and is competitive in the remaining two cases. For more results, we refer the reader to Tables 19 and 10 of the Supplementary materials.

In Supplementary Fig. 10, we also observe that the multiclass sAwMIL performs better on the chat models (see bottom right portion of the plot) than the multiclass SVM. This supports our claim that veracity signals often emerge at positions other than the final token, and that multiple-instance learning can better isolate the veracity signal. Recall that the multiclass sAwMIL probe considers all the tokens in the statement and has additional training stages

Supplementary Fig. 11 visualizes the per-token predictions. The one-vs-all SVM-based probes have better selectivity than the mean-difference probe with conformal prediction intervals (MD+CP). However, in some cases, one-vs-all SVM assigns labels to the tokens in the pre-actualized part of the statement (e.g., see Supplementary Fig. 11B, #5 statement). This suggests that one-vs-all SVM probes are capturing spurious correlations with potential proxies.



Figure 9: **Per-token predictions of mean-difference with conformal prediction intervals (MD+CP) and one-vs-all sAwMIL on the 13th decoder activations of the default Llama-3-8B model.** Statements are from the *City Locations* data set. We show per-token probabilities (printed beneath each word), assigned based on the token’s representation. Words are shaded based on the predicted probability. When MD+CP outputs 0, the statement is labeled false; when it outputs 1, the statement is labeled true. If the per-token score falls outside the conformal intervals, MD+CP assigns a score of 0.5 to that token (which corresponds to the highest uncertainty). The one-vs-all sAwMIL probe for *is-true* outputs 1 when the probe is 100% confident that the statement is true; and it outputs 0 when the per-token score is outside the conformal intervals (i.e., there is an absence of truthfulness signal). Similarly, the one-vs-all sAwMIL probe for *is-false* outputs 1 when the probe is 100% confident that the statement is false; and it outputs 0 when the per-token score is outside the conformal intervals (i.e., there is an absence of falsehood signal). The same logic applies for the one-vs-all sAwMIL probe for *is-neither*. **Panel A** shows the MD+CP predictions. It often assigns high confidence scores to pre-actualization tokens and makes mistakes on the *wrapped* prompts in cases #5 and 6 (e.g., statement #6: “Hey,___ Is this correct?”). Also, MD+CP probe assigns labels to the statement without any veracity value (e.g., case #7). **Panels B–D** display one-vs-all sAwMIL probes (*is-true*, *is-false*, and *is-neither*). Unlike MD+CP, one-vs-all sAwMIL localizes the veracity signal to the actualized token and abstains elsewhere, demonstrating superior selectivity.

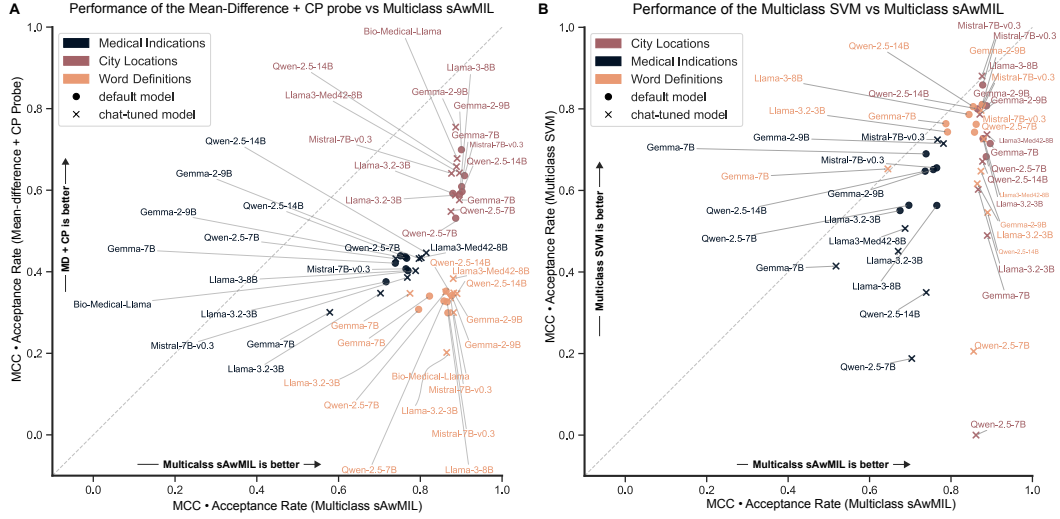


Figure 10: Comparison of performances for MD+CP, a multiclass SVM and multiclass sAwMIL probes. Panels A & B: Each marker shows a probe’s performance for a \langle model, dataset \rangle pair. Default models are shown with circles, while chat models are shown with crosses. The different colors indicate the different data sets. **Panel A** shows the comparison between the multiclass sAwMIL probe on the x-axis and the mean-difference probe with conformal prediction intervals (MD+CP) on the y-axis. **Panel B:** Performance of multiclass sAwMIL vs. multiclass single-instance SVM probes. The performance of the multiclass sAwMIL probe is specified on the x-axis, while the performance of the multiclass SVM is specified on the y-axis. We observe that multiclass sAwMIL outperforms multiclass SVM. The only exceptions are the Gemma-7B chat model and the Mistral-7B-v0.3 chat model on *Word Definitions*, where the performances of the single-instance and multiple-instance are competitive. This experiment shows that multi-instance learning (i.e., training on all the tokens in the statement) is beneficial when tracking the veracity of an LLM. Overall, multiclass sAwMIL probe outperforms MD+CP and the multiclass SVM.



Figure 11: **Per-token predictions of one-vs-all single-instance SVM on the 13th decoder activations of the default Llama-3-8B** . Statements are from the *City Locations* data set. We show per-token probabilities (printed beneath each word), assigned based on the token’s representation. Words are shaded based on their predicted probability. **Panels A–C:** display the one-vs-all SVM probes (is-true, is-false, and is-neither). The one-vs-all SVM probe isolates the signal better than the MD+CP probe. (See Supplementary Fig. 9A for the MD+CP results.) However, in some instances, the one-vs-all SVM probe assigns high-certainty scores to tokens that do not have any veracity signal. For example, in Panel A (case #6, the probe picks up on tokens including ‘this,’ ‘is,’ and ‘?’, which do not have inherent veracity value. Overall, the multiclass sAwMIL in Supplementary Fig. 9B–C has better selectivity.

I Additional Evaluation Details

In the Supplementary Sec. E, we provide the full list of the validity criteria for the. In the manuscript, we only cover the results related to **Correlation** and **Locality** criteria (see Sec. 5). Here, we provide additional details behind the evaluation of the **Correlation** and **Locality** (see below in Supplementary Sec. I.1.1) Further, we describe the experimental setup related to the Generalization (also see below in Supplementary Sec. I.1.1), and setups for the **Manipulation** and **Locality** criteria (see Supplementary Sec. I.1.2).

Finally, we provide the evaluation results for the **Generalization**, **Manipulation** and **Locality** in Supplementary Sec. I.2

I.1 Evaluation Setup

I.1.1 Performance and Validity

Here, we describe a pipeline to evaluate our sAwMIL probe over the validity criteria specified in Sec. 3 and Supplementary Sec. E.

Correlation and Selectivity. We use the test split of each data set to evaluate the performance of the probe. We use Matthew’s Correlation Coefficient (MCC) to summarize the statistical accuracy of probes. The multiclass MCC value is calculated using Eq. 14.

$$\text{MCC} = \frac{c \times s - \sum_k^K p_k \times t_k}{\sqrt{\left(s^2 - \sum_k^K p_k^2\right) \times \left(s^2 - \sum_k^K t_k^2\right)}}, \quad (14)$$

where c is the number of correct predictions, s is the total number of samples, K is the total number of classes, t_k is the number of k -class samples in the data set, and p_k is the number of times k -class was predicted. $\text{MCC} = 1$ indicates that a classifier predicted every instance correctly. $\text{MCC} = 0$ implies that the predictions are random. $\text{MCC} = -1$ indicates that the predictions are inversely correlated with the ground-truth labels.

Since *neither* statements are included in the test data set, MCC provide a sense of how well the probe classifies factually *true* or *false* statements and indicate whether the probes can handle *neither*-type cases.

Generalizability. To test how well a particular probe g_i trained on data set \mathcal{D}_i generalizes, we evaluate its performance using the test split of other data set \mathcal{D}_j —e.g., g_i trained on the city locations data set is evaluated using the test split of the word definitions data set.

I.1.2 Interventions and Validity

In this experiment, we assess whether perturbing the hidden representation $h_i(\mathbf{x})$ along the veracity direction \vec{v}_i affects the model’s outputs, and whether these interventions satisfy the manipulation and locality criteria defined in Supplementary Sec. E. In other words, we use \vec{v}_i to change the distribution of the output tokens $P_{\mathcal{M}}$ and force true or false responses.

We look at each factually true statement $\mathbf{x} \in \mathcal{D}_{test}$ —e.g., “The city of Santo Domingo is in the Dominican Republic.” We split these statements into two segments: a pre-actualized part, \mathbf{x}^p , such as “The city of Santo Domingo is in”; and second, an actualized part, \mathbf{x}^a such as “the Dominican Republic”. Given \mathbf{x}^p , we can compute the probability of the actualized part according to Eq. 15:

$$P_{\mathcal{M}}(\mathbf{x}_{[1:L]}^a \mid \mathbf{x}^p) = \prod_{l=1}^L P_{\mathcal{M}}\left(\mathbf{x}_{[l]}^a \mid \mathbf{x}_{[0:(l-1)]}^a, \mathbf{x}^p\right). \quad (15)$$

In Eq. 15, L is the number of tokens in the actualized part \mathbf{x}^a .²⁶ The subscript $[l]$ specifies the index of a token in \mathbf{x}^a , and $[1 : l]$ specifies the range of tokens in \mathbf{x}^a , while $\mathbf{x}_{[0]}^a$ refers to an empty set. Further, we do not specify a subscript $[1 : L]$ for brevity unless it is necessary for clarity.

²⁶The length of the actualized part depends on the tokenization technique the language model uses. For example, Llama-3 and Mistral models split *Albania* into [A1] [ban] [ia], while Gemma models have one reserved token [Albania].

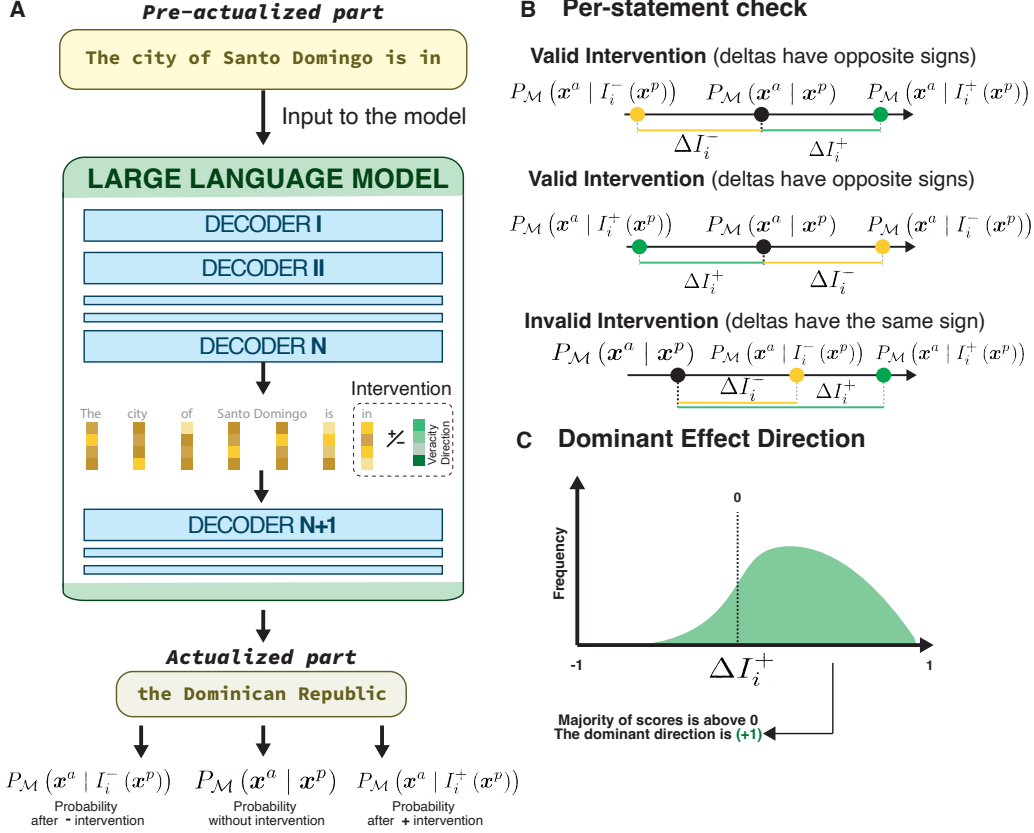


Figure 12: **Workflow for evaluating the success of directional interventions.** (**Panel A**) We provide a pre-actualized part x^p to the LLM, and intervene at the final token embedding after the n -th decoder by adding or subtracting the learned veracity direction \vec{v}_i . The modified representation is then passed to decoder $n + 1$. Further, we compute the conditional probability of the *correct* actualized part x^a . For each statement, we obtain three values: the original conditional probability, the probability after a positive shift $+\vec{v}_i$, and the probability after a negative shift $-\vec{v}_i$. (**Panel B**) We assess the success of each intervention at the statement level, comparing the change in conditional probabilities. If the change is caused by the positive shift (ΔI_i^+) and the negative shift (ΔI_i^-) are of opposite sign, we consider the intervention to have a consistent directional effect and mark it as successful for that statement (see Eq. 19). (**Panel C**) To evaluate the success of the intervention, we also identify the *dominant effect direction*—i.e., the sign of ΔI_i^+ that appears most frequently across all statements. The procedures in panels B and C determine the per-statement success (see Eq. 19). If more than half of the statements show consistent and directionally aligned changes (see Eq. 20), we consider the overall intervention along \vec{v}_i to be successful.

Direction vector \vec{v}_i . We train one-vs-all sAwMIL probes via dual optimization. We use the obtained solution to extract the linear direction that points towards the class of interest. For example, in our *is-true* probe, the class of interest is the true statements.

$$\vec{v}_i = \sum_{j \in \mathcal{S}} \alpha_j y_j h_i(\mathbf{x}_j) \text{ with } \mathcal{S} = \{j \mid \alpha_j > 0\}, \alpha_j \in \mathbb{R}, y_j \in \{-1, 1\}. \quad (16)$$

In Eq. 16, \mathcal{S} is a set of support vectors, α_j is the Lagrangian multiplier [40], $h_i(\mathbf{x}_j)$ is the activation after the i^{th} decoder for statement \mathbf{x}_j , and y_j is the class label for \mathbf{x}_j .

Interventions. Given a pre-actualized part of the statement x^p , model \mathcal{M} , and a hidden representation $h_i(x^p)$, we apply directional interventions by translating the representation of the last token $x_{[L]}^p$ along $\pm \vec{v}_i$ (here, i stands for the index of the decoder):

- **Positive directional shift:**

$$I_i^+ (h_i(\mathbf{x}^p)) = [h_i(\mathbf{x}^p)_{[1]}, \dots, h_i(\mathbf{x}^p)_{[L-1]}, h_i(\mathbf{x}^p)_{[L]} + \vec{v}_i]$$

- **Negative directional shift:**

$$I_i^- (h_i(\mathbf{x}^p)) = [h_i(\mathbf{x}^p)_{[1]}, \dots, h_i(\mathbf{x}^p)_{[L-1]}, h_i(\mathbf{x}^p)_{[L]} - \vec{v}_i]$$

These interventions return modified representations, which we denote as $I_i^+(\mathbf{x}^p)$ and $I_i^-(\mathbf{x}^p)$, respectively. Furthermore, we compute the per-sample effect of the directional interventions. It is defined as a difference between the original probability and the probability we get after the intervention:

$$\Delta I_i^+ (\mathbf{x}^a, \mathbf{x}^p) \leftarrow P_{\mathcal{M}} (\mathbf{x}^a | I_i^+ (\mathbf{x}^p)) - P_{\mathcal{M}} (\mathbf{x}^a | \mathbf{x}^p) \quad (17)$$

$$\Delta I_i^- (\mathbf{x}^a, \mathbf{x}^p) \leftarrow P_{\mathcal{M}} (\mathbf{x}^a | I_i^- (\mathbf{x}^p)) - P_{\mathcal{M}} (\mathbf{x}^a | \mathbf{x}^p) \quad (18)$$

To compare interventions across decoders and models, we look at the success rate of the interventions. The per-statement intervention is successful if ΔI_i^+ and ΔI_i^- have opposing effects on the conditional probability $P_{\mathcal{M}} (\mathbf{x}^a | \mathbf{x}^p)$. In other words, if ΔI_i^+ is positive, then ΔI_i^- must be negative, and vice-versa.

At the same time, we must ensure that the effect of the intervention is consistent in most statements $\mathbf{x} \in \mathcal{D}_{test}$. If half of the statements have positive ΔI_i^+ and another have negative ΔI_i^+ , then our intervention produces a random change in $P_{\mathcal{M}}$. To ensure that the effect is consistent, we look at the *dominant effect direction* of the intervention, $\bar{d}_i \in \{-1, +1\}$. Here, $\bar{d}_i = +1$, if more than half ΔI_i^+ are positive, and $\bar{d}_i = -1$ if more than half is negative.

In summary, a successful directional intervention is one that produces opposing effects when shifting along $\pm \vec{v}_i$, and aligns ΔI_i^+ with the dominant direction. Supplementary Fig. 12 provides an overview of the workflow to determine the per-statement success of the intervention. Formally, we define a per-statement success s as

$$s(\mathbf{x}) = \mathbb{I} \left[[\text{sign} (\Delta I_i^+ (\mathbf{x}^a, \mathbf{x}^p)) \neq \text{sign} (\Delta I_i^- (\mathbf{x}^a, \mathbf{x}^p))] \wedge [\text{sign} (\Delta I_i^+ (\mathbf{x}^a, \mathbf{x}^p)) = \text{sign} (\bar{d}_i)] \right] \quad (19)$$

where

$$\mathbb{I} [\cdot] = \begin{cases} 1, & \text{if the condition holds,} \\ 0, & \text{otherwise.} \end{cases}$$

Why do we only look at the sign? During our initial experiments, we observed that even when \vec{v}_i was trained to separate true and false statements, the effect of the intervention on $P_{\mathcal{M}}(\mathbf{x}^a | \mathbf{x}^p)$ could vary across decoders. Specifically, in certain decoders, shifting along $\pm \vec{v}_i$ consistently increased the probability of \mathbf{x}^a , while in others it decreased it. These effects were consistent in the sense that shifting in the opposite direction produced the opposing effect. This phenomenon can be attributed to the complex interactions within each decoder of the model. As highlighted by Heimersheim and Nanda [41], activation patching experiments have revealed that interventions can have varying directions of effect depending on the decoder. Therefore, observing a sign flip in the effect of an intervention does not invalidate the direction \vec{v}_i , as long as it is consistent.

Given a set of true statements from \mathcal{D}_{test} , we compute the overall success rate at a decoder i as follows:

$$\omega_i = \frac{1}{N} \sum s(\mathbf{x}_j), \quad (20)$$

where N stands for the number of the true statements in \mathcal{D}_{test} and $s(\mathbf{x}_j)$ is success for the j^{th} statement as defined in Eq. 19.

If the overall success rate at decoder i is greater than 50%, we claim that the intervention at decoder i is successful. Hence, the manipulation criterion is fulfilled. We use a one-sided binomial test to confirm whether the overall success rate is significant:

$$H_0 : \omega_i \leq 0.5 \quad (21)$$

$$H_A : \omega_i > 0.5 \quad (22)$$

For example, the overall success rate of 61% with the dominant direction $\bar{d}_i = +1$ tells us that if we have 100 statements, on average, we increase the probability of the correct answer in 61 statements. Similarly, a success rate of 98% with the dominant direction $\bar{d}_i = -1$ indicates that shifting along $+\bar{v}_i$ decreases the probability of the correct answer in approximately 98 out of 100 statements.

Locality. To further determine the quality of the directional intervention, we assess whether changes in probability are concentrated on the actualized part, rather than being diffused across random tokens in the vocabulary \mathcal{V} . In other words, our intervention should change $P_{\mathcal{M}}(\mathbf{x}^a \mid \mathbf{x}^p)$ and should not change the probability of random tokens $P_{\mathcal{M}}(\mathbf{r} \mid \mathbf{x}^p)$.

Specifically, we expect the intervention to primarily affect the likelihood of the correct continuation, $P_{\mathcal{M}}(\mathbf{x}_{[1:L]}^a \mid \mathbf{x}^p)$, while leaving the probability of a randomly sampled continuation, $P_{\mathcal{M}}(\mathbf{r}_{[1:L]} \mid \mathbf{x}^p)$, mostly unchanged. Here, $\mathbf{r}_{[1:L]}$ denotes a random sequence sampled from the vocabulary of the model \mathcal{M} . We quantify these changes as:

$$\Delta_{\text{Correct}} = |P_{\mathcal{M}}(\mathbf{x}^a \mid I_i^+(\mathbf{x}^p)) - P_{\mathcal{M}}(\mathbf{x}^a \mid I_i^-(\mathbf{x}^p))| \quad (23)$$

$$\Delta_{\text{Random}} = |P_{\mathcal{M}}(\mathbf{r} \mid I_i^+(\mathbf{x}^p)) - P_{\mathcal{M}}(\mathbf{r} \mid I_i^-(\mathbf{x}^p))| \quad (24)$$

Further, we say that the intervention satisfies the *locality* criterion if

$$\mathbb{E}[\Delta_{\text{Correct}}] > \mathbb{E}[\Delta_{\text{Random}}]. \quad (25)$$

That is, the expected change in probability for the correct output, \mathbf{x}^a , exceeds the expected change for a randomly sampled output \mathbf{r} .

I.2 Results

I.3 Generalization Across Data Sets

To further support the claim that the multiclass sAwMIL captures veracity signals (and not merely a proxy), we demonstrate generalization performance across data sets. Supplementary Fig. 13 provides results for each data set and LLM. The columns correspond to three test data sets, and the cells specify the multiclass sAwMIL’s performance for a specific LLM. Multiclass sAwMIL provides reasonable generalization performance (see Supplementary Tab. 8). However, it is potentially overfitting to the highly specialized *City Locations* data set. Using more diverse data sets that contain a broader range of entities and cover a larger set of topics isolates the veracity signal better and produces better generalization performance. We refer the reader to Tables 21 through 23 in Supplementary for detailed statistics.

Table 8: **Aggregated generalization performance of the multiclass sAwMIL for each dataset.** Each cell shows a MCC value, which quantifies the performance of the multiclass sAwMIL trained and tested on different combinations of the datasets. The value in the bracket is the standard error. *Word Definitions* provides better generalization performance because it contains statements covering a diverse set of topics, while the *City Locations* provide lower generalization performance.

Training Dataset	Testing Dataset		
	City Locations	Medical Indications	Word Definitions
City Locations	0.963 (0.003)	0.624 (0.030)	0.633 (0.025)
Medical Indications	0.818 (0.033)	0.790 (0.009)	0.698 (0.018)
Word Definitions	0.896 (0.015)	0.723 (0.016)	0.868 (0.008)

We also observe that generalization performance is higher for chat models, where the average MCC score (on the non-training data set) is 77.2% (standard error: 0.2%), compared to 68.2% (standard error: 0.2%) achieved for the default models. This is more noticeable in Supplementary Fig. 13A, where default models have much lower MCC values than their chat model counterparts. For example, the chat model Llama-3.2-3B has 1.6 times higher MCC value than the default Llama-3.2-3B. Over the three panels, Gemma-7B seems to be an outlier, since the generalization performance drops significantly for the chat version of the model.

Multiclass sAwMIL satisfies the generalization criterion defined in Supplementary Sec. E by transferring veracity probes trained on one data set to another while maintaining strong performance. This

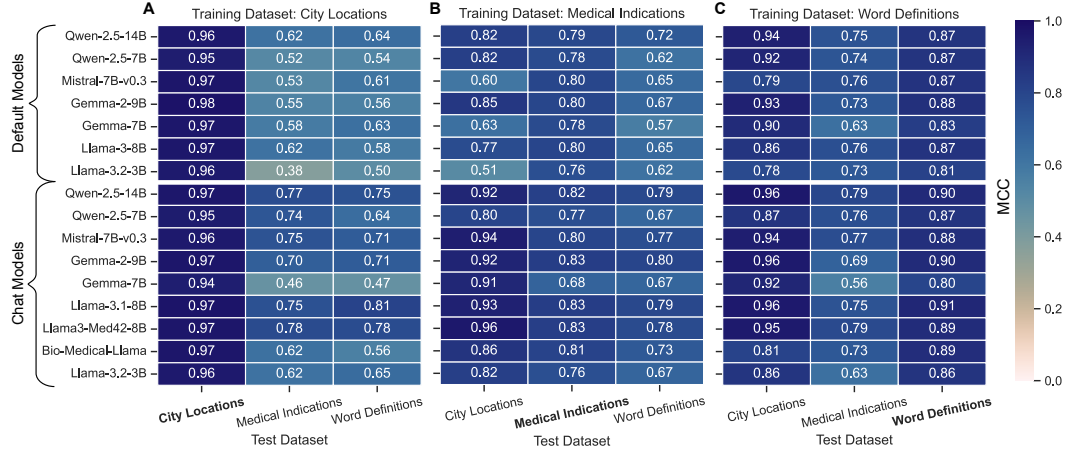


Figure 13: **Generalization performance of the multiclass sAwMIL probe across data sets.** Each panel corresponds to a different *training* data set: *City Locations*, *Medical Indications*, and *Word Definitions*. Each column corresponds to a different *test* data set. Each cell displays MCC values, which quantify how well the probe generalizes to the test data set (a higher value is better). For each model and data set, we report the maximum MCC achieved across all decoders. Generally, probes trained on the chat models have better generalization performance than the default models. **Panel A:** Generalization performance of multiclass sAwMIL when trained on *City Locations*. While the MCC values are significantly higher than random baseline (with $MCC = 0$), the generalization ability is lower than those in Panels B or C. **Panel B:** Generalization performance of multiclass sAwMIL when trained on *Medical Indications*. In Panel A, we observe that training on *City Locations* and testing on *Medical Indications* provides good but not excellent MCC values (average MCC of 0.624 with standard error of 0.030). This is not the case in this panel, where *Medical Indications* is the training data set and *City Locations* is the test data set (average MCC of 0.818 with standard error of 0.033). **Panel C:** Generalization performance of multiclass sAwMIL when trained on *Word Definitions*. This probe has high generalization performance across data sets. When *City Locations* is the test data set, the average MCC is 0.896 with standard error of 0.015; and when *Medical Indications* is the test data set, the average MCC is 0.723 with standard deviation of 0.016. For aggregated statistics, see Supplementary Tab. 8.

provides further evidence that multiclass sAwMIL captures a veracity signal that is not specific to a data set.

I.4 Interventions: Manipulation and Locality

Previous experiments have shown that the multiclass sAwMIL identifies a strong and transferable veracity signal. We further look at how this signal is connected to the output of an LLM, $P_{\mathcal{M}}$. Here, we look at the interventions of one-vs-all sAwMIL probes, which are the building blocks of the multiclass sAwMIL. Specifically, we assess the effectiveness of interventions applied along the is-true and is-false directions. For simplicity, we exclude the one-vs-all probe for is-neither from the intervention analysis.

The overall success rate measures how often directional interventions (adding or subtracting $\pm \vec{v}_i$) produce a consistent change in the model’s output $P_{\mathcal{M}}$. We describe the setup in Supplementary Sec. I.1.2. A higher success rate indicates that shifting along $\pm \vec{v}_i$ has higher chances to skew the conditional probability of the correct answers $P_{\mathcal{M}}(\mathbf{x}^a | \mathbf{x}^p)$. Supplementary Fig. 14 shows the overall success rate for each model and data set. Success rates below .5 suggest that interventions have close to a random effect on $P_{\mathcal{M}}$.

Notably, some models have a success rate of 0. This occurs when:

- Interventions along $\pm \vec{v}_i$ failed to induce opposing changes in $P_{\mathcal{M}}$, e.g., both $+\vec{v}_i$ and $-\vec{v}_i$ increased or decreased probabilities; or

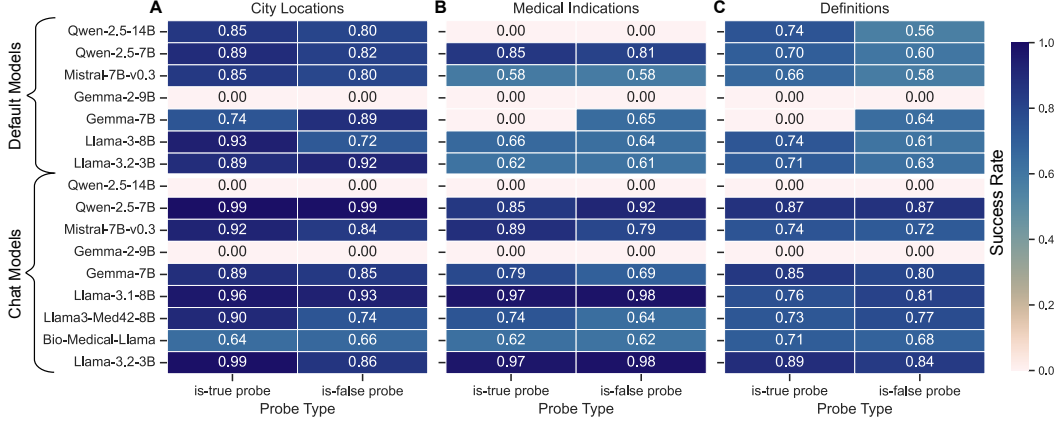


Figure 14: **Overall intervention success rate for the one-vs-all sAwMIL probes.** We report the maximum achievable success rate along the `is-true` and `is-false` directions. Panels A–C show results for probes g_i trained on specific data sets, while columns correspond to the `is-true` or `is-false` probe. Default and chat Gemma-2-9b, as well as, some experiments with the default and chat Qwen-2.5-14B and Gemma-7B models did not pass the consistency check – suggesting that the interaction between directions \vec{v}_i and the conditional probability $P_{\mathcal{M}}$ are not linear (or that the one-vs-all sAwMIL probes failed to identify signals that linearly affect the model output $P_{\mathcal{M}}$).

- The average change in $P_{\mathcal{M}}(x^a | x^p)$ (i.e., Δ_{correct}) matched the change in random sequence continuation probabilities (Δ_{random}), violating the locality criterion in Equations 23 through 25.

The average success rate is 80.1% (standard error: 0.2%) for the `is-true` direction and 76.2% (standard error: 0.2%) for the `is-false` direction. We exclude the models whose intervention success rate was 0 in Supplementary Fig. 14.

This experiment shows that in the majority of cases, we can use the `is-true` and `is-false` directions to manipulate the output of LLMs. The interventions are more successful for the chat models. The average success rate for chat models is 80.1% (standard error: 0.3%), and for default models is 70.9% (standard error: 0.2%).

Anomalies. We observe some anomalous behavior when intervening in the LLMs:

- In the Gemma-2-9B models, interventions along the direction $\pm \vec{v}_i$ consistently increased or decreased the probability $P_{\mathcal{M}}(x^a | x^p)$, regardless of the sign of the intervention. This indicates that the direction \vec{v}_i identified by sAwMIL does not have a clear relationship with the model’s output probabilities. Thus, we cannot use \vec{v}_i to increase and decrease the probability of correct answers. A similar phenomenon is observed in the Qwen-2.5-14B (chat) model.
- For the Gemma-7B (default) model, interventions along the `is-false` direction provide a higher success rate compared to the `is-true` direction.

The exact reasons for these anomalies are unclear. However, we know that these models have additional fine-tuning processes. These additional training procedures may have influenced the internal representations of the models. Except for the Gemma-2-9B models and the Qwen-2.5-14B (chat) model, the one-vs-all sAwMIL probes pass both manipulation and locality criteria. We expect that the non-linear version of sAwMIL will overcome the issues with the models that have additional fine-tuning processes. This is part of our future work.

I.5 Recap: Overall Validity

In this section and Sec. 5 of the manuscript, we established that the multiclass sAwMIL probe satisfies the validity criteria. Specifically, we confirmed that it satisfies the correlation and selectivity criteria, outperforming zero-shot prompting and mean-difference probe with conformal prediction intervals.

We further demonstrated that multiclass sAwMIL satisfies the generalization criterion, indicating that we can successfully apply probes trained on multiclass sAwMIL to statements from other domains. In addition, we showed that the one-vs-all sAwMIL probes satisfy the manipulation and locality criteria. In a majority of cases, we can perform interventions that change the probabilities of correct replies. Together, these findings provide strong evidence for the overall validity of sAwMIL probes. Not all LLMs have a veracity mechanism that has a linear relationship with the output. Exploring this non-linear relationship is part of our future work.

J Algorithms

In this section, we provide pseudo-codes for several procedures described in the main text. In Supplementary Alg. 1, we provide pseudo-code for the Sparse Aware Multiple-Instance Learning (sAwMIL) probe,. In Supplementary Alg. 2, 3 and 4, we show pseudo-codes for the mean-difference (MD), the training with truth and polarity directions (TPPD), and the supervised PCA (sPCA) probes, accordingly. In Supplementary Alg. 5 and 6, we describe the procedure for the binary and multiclass conformal learning (described in Sec. 3.2).

Algorithm 1 Training a one-vs-all sAwMIL classifier

Input: A training data set $\{(\mathbf{x}_i, y_i, \mathbf{m}_i)\}_{i=1}^n$ with binary bag labels $y_i \in \{0, 1\}$, bags $\mathbf{x}_i \in \mathbb{R}^{L_i \times d}$, and intra-bag confidences $\mathbf{m}_i \in \{0, 1\}^{L_i}$, where L_i is the number of items in a bag \mathbf{x}_i ; also, a balancing parameter $\eta \in (0, 1]$.

Output: Parameters $\boldsymbol{\theta} \in \mathbb{R}^{1 \times d}$ and $b \in \mathbb{R}$. These parameters are subsequently given to a function f along with \mathbf{z} to compute $f(\mathbf{z}) = \sigma(\mathbf{z} \boldsymbol{\theta}^T + b)$, where σ is a sigmoid function.

- 1: Partition data into positive and negative sets:

$$\mathcal{X}^+ = \{(\mathbf{x}_i, \mathbf{m}_i) : y_i = 1\}, \quad \mathcal{X}^- = \{(\mathbf{x}_i, \mathbf{m}_i) : y_i = 0\}$$

- 2: Compute the **initial** coefficient vector and the intercept ▷ See Bunescu and Mooney [19]

$$(\hat{\boldsymbol{\theta}}, \hat{b}) \leftarrow \text{solve_SMIL}(\mathcal{X}^+, \mathcal{X}^-), \text{ where } \hat{\boldsymbol{\theta}} \in \mathbb{R}^{1 \times d} \text{ and } \hat{b} \in \mathbb{R}.$$

- 3: Let $\bar{\mathcal{X}}^+$ denote the set of all instances from the positive bags and $\bar{\mathcal{X}}^-$ all instances from the negative bag.
- 4: Compute scores for every instance in a positive set

$$S^+ \leftarrow \bar{\mathcal{X}}^+ \hat{\boldsymbol{\theta}}^T + \hat{b}, \text{ where } \bar{\mathcal{X}}^+ \in \mathbb{R}^{|\bar{\mathcal{X}}^+| \times d}.$$

- 5: Compute the threshold

$$q \leftarrow \text{quantile}(S^+, 1 - \eta).$$

- 6: **for all** positive instances $\langle \bar{\mathbf{x}}_j, \bar{\mathbf{m}}_j, \bar{y}_j \rangle \in \bar{\mathcal{X}}^+$, where $\bar{\mathbf{x}}_j \in \mathbb{R}^{1 \times d}$, $\bar{\mathbf{m}}_j \in \{0, 1\}$ and $\bar{y}_j = \emptyset$ **do**

if $(\bar{\mathbf{x}}_j \hat{\boldsymbol{\theta}}^T + b) \geq q_\eta$ and $\bar{\mathbf{m}}_j = 1$ **then** set $\bar{y}_j = 1$;

else set $\bar{y}_j = 0$.

- 7: **end for**

- 8: Compute the **final** coefficient vector and the intercept ▷ via simple support vector machine

$$(\boldsymbol{\theta}, \beta) \leftarrow \text{solve_SIL}(\bar{\mathcal{X}}^+, \bar{\mathcal{X}}^-).$$

return Coefficient vector $\boldsymbol{\theta} \in \mathbb{R}^{1 \times d}$ and the intercept $\beta \in \mathbb{R}$.

Algorithm 2 Training a mean-difference (MD) probe [15], sometimes referred to as mean-mass/mean-cluster difference classifier or linear discriminant analysis.

Input: A training dataset $\{\langle \mathbf{z}_i, y_i \rangle\}_{i=1}^n$ with binary labels $y_i \in \{0, 1\}$ and $\mathbf{z}_i \in \mathbb{R}^{1 \times d}$. In our experiments, \mathbf{z} is the embedding of the last token (unless otherwise noted).

Output: Parameters $\boldsymbol{\theta} \in \mathbb{R}^{1 \times d}$, $\beta \in \mathbb{R}$, and $\boldsymbol{\Sigma}^{-1} \in \mathbb{R}^{d \times d}$. These parameters are subsequently given to a function f along with \mathbf{z} to compute $f(\mathbf{z}) = \sigma(\mathbf{z}(\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}^T) + \beta)$, where σ is a sigmoid function.

- 1: Partition data into positive and negative sets:

$$\mathcal{X}^+ = \{\mathbf{z}_i : y_i = 1\}, \quad \mathcal{X}^- = \{\mathbf{z}_i : y_i = 0\}$$

- 2: Compute class means $\boldsymbol{\mu}^+$ and $\boldsymbol{\mu}^-$, and covariance matrices $\boldsymbol{\Sigma}^+$ and $\boldsymbol{\Sigma}^-$ for \mathcal{X}^+ and \mathcal{X}^- .
- 3: Compute pooled covariance matrix (where $n^+ = |\mathcal{X}^+|$ and $n^- = |\mathcal{X}^-|$):

$$\boldsymbol{\Sigma} = \frac{(n^+ - 1)\boldsymbol{\Sigma}^+ + (n^- - 1)\boldsymbol{\Sigma}^-}{n^+ + n^- - 2}$$

- 4: Compute the coefficient vector: $\boldsymbol{\theta} = \boldsymbol{\mu}^+ - \boldsymbol{\mu}^-$, where $\boldsymbol{\theta} \in \mathbb{R}^{1 \times d}$.
- 5: Compute scores for positive and negative sets:

$$s^+ \leftarrow \mathcal{X}^+ (\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}^T) \text{ and } s^- \leftarrow \mathcal{X}^- (\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}^T), \text{ where } s^+ \in \mathbb{R}^{n^+} \text{ and } s^- \in \mathbb{R}^{n^-}.$$

- 6: Compute the intercept:

$$b = \frac{1}{2} (\text{mean}(s^+) + \text{mean}(s^-))$$

return Coefficient vector $\boldsymbol{\theta}$, intercept β , and the inverse covariance matrix $\boldsymbol{\Sigma}^{-1}$.

Algorithm 3 Training with Truth and Polarity Directions, aka TTPD probe [17]

Input: A training dataset $\{\langle \mathbf{z}_i, y_i, p_i \rangle\}_{i=1}^n$ with binary truthfulness labels $y_i \in \{-1, 1\}$, binary polarity labels $p_i \in \{-1, 1\}$ that specify affirmative/negated statements and $\mathbf{z}_i \in \mathbb{R}^{1 \times d}$. In our experiments, \mathbf{z} is the embedding of the last token (unless otherwise noted).

Output: Parameters $\Theta \in \mathbb{R}^{2 \times d}$, $\boldsymbol{\theta} \in \mathbb{R}^{1 \times 2}$ and $\beta \in \mathbb{R}$. These parameters are subsequently given to a function f along with \mathbf{z} to compute $f(\mathbf{z}) = \sigma((\mathbf{z}\Theta^T)\boldsymbol{\theta}^T + \beta)$, where σ is a sigmoid function.

- 1: Given data matrix $X = [\mathbf{z}_0, \dots, \mathbf{z}_n] \in \mathbb{R}^{n \times d}$ compute the centered data matrix

$$\bar{X} = X - \text{mean}(X)$$

- 2: Find the truth direction $\boldsymbol{\theta}_t \in \mathbb{R}^{1 \times d}$ via the Ordinary Least Squares:

$$\boldsymbol{\theta}_t = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \mathbf{y}, \text{ where } \mathbf{y} = [y_0, \dots, y_n]$$

- 3: Find the polarity direction $\boldsymbol{\theta}_p \in \mathbb{R}^{1 \times d}$ via the Logistic Regression:

$$\boldsymbol{\theta}_p \leftarrow \text{LogisticRegression}(X, \mathbf{p}), \text{ where } \mathbf{p} = [p_0, \dots, p_n]$$

- 4: Project X onto $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}_p$:

$$\hat{X} \leftarrow \text{stack}([X\boldsymbol{\theta}_t^T, X\boldsymbol{\theta}_p^T]), \text{ where } \hat{X} \in \mathbb{R}^{n \times 2}$$

- 5: Use the *projected* data matrix \hat{X} to get coefficient vector $\boldsymbol{\theta}$ and the intercept β :

$$\boldsymbol{\theta}, \beta \leftarrow \text{LogisticRegression}(\hat{X}, \mathbf{y})$$

return Projection matrix $\Theta = \text{stack}([\boldsymbol{\theta}_t, \boldsymbol{\theta}_p])$, coefficient vector $\boldsymbol{\theta}$ and intercept β .

Algorithm 4 Training a Supervised Principal Component Analysis (sPCA) probe

Input: A training dataset $\{\langle \mathbf{z}_i, y_i \rangle\}_{i=1}^n$ with binary labels $y_i \in \{0, 1\}$ and $\mathbf{z}_i \in \mathbb{R}^{1 \times d}$, and $k \in \mathbb{R}$ to specify the number of components. In our experiments, \mathbf{z} is the embedding of the last token (unless otherwise noted).

Output: Parameters $\Lambda \in \mathbb{R}^{k \times d}$, $\boldsymbol{\theta} \in \mathbb{R}^{1 \times 2}$ and $\beta \in \mathbb{R}$. These parameters are subsequently given to a function f along with \mathbf{z} to compute $f(\mathbf{z}) = \sigma((\mathbf{z} \Lambda^T) \boldsymbol{\theta}^T + \beta)$, where σ is a sigmoid function.

- 1: Partition data into positive and negative sets:

$$\mathcal{X}^+ = \{\mathbf{z}_i : y_i = 1\}, \quad \mathcal{X}^- = \{\mathbf{z}_i : y_i = 0\}$$

- 2: Compute means and centered matrices:

$$\boldsymbol{\mu}^+ = \text{mean}(\mathcal{X}^+), \quad \boldsymbol{\mu}^- = \text{mean}(\mathcal{X}^-), \quad \boldsymbol{\mu} = \text{mean}(\mathcal{X}), \quad \mathcal{X}_c^+ = \mathcal{X}^+ - \boldsymbol{\mu}^+, \quad \mathcal{X}_c^- = \mathcal{X}^- - \boldsymbol{\mu}^-.$$

- 3: Compute class covariance matrices $\boldsymbol{\Sigma}^+$ and $\boldsymbol{\Sigma}^-$ for \mathcal{X}_c^+ and \mathcal{X}_c^- .
- 4: Compute within-class covariance matrix (where $n^+ = |\mathcal{X}_c^+|$ and $n^- = |\mathcal{X}_c^-|$):

$$\boldsymbol{\Sigma}_w = \frac{(n^+ - 1) \boldsymbol{\Sigma}^+ + (n^- - 1) \boldsymbol{\Sigma}^-}{n^+ + n^- - 2}$$

- 5: Compute between-class covariance matrix:

$$d_p = (\boldsymbol{\mu}^+ - \boldsymbol{\mu}), \quad d_n = (\boldsymbol{\mu}^- - \boldsymbol{\mu}), \quad \boldsymbol{\Sigma}_b = \frac{n^+ d_p d_p^\top + n^- d_n d_n^\top}{n}.$$

- 6: Build symmetric scatter matrix with ridge penalty:

$$\mathbf{M} = \boldsymbol{\Sigma}_b + \boldsymbol{\Sigma}_w + \lambda \mathbf{I}_d, \quad \mathbf{M} \leftarrow \frac{1}{2}(\mathbf{M} + \mathbf{M}^\top).$$

- 7: Compute top- k eigenpairs of \mathbf{M} (largest algebraic), assemble the eigenmatrix $\Lambda \in \mathbb{R}^{k \times d}$ and corresponding eigenvalue vector $\boldsymbol{\nu} \in \mathbb{R}^{1 \times k}$:

$$(\boldsymbol{\nu}_1, v_1), \dots, (\boldsymbol{\nu}_k, v_k) = \text{TopKEigs}(\mathbf{M}, k) \rightarrow \Lambda = [\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_k], \quad \mathbf{v} = [v_1, \dots, v_k]^\top.$$

- 8: Project data matrix $\mathbf{X} = [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathbb{R}^{n \times d}$ via the eigenmatrix Λ , followed by whitening:

$$\hat{\mathbf{X}} = \mathbf{X} \Lambda^T \text{diag}(\mathbf{v})^{-\frac{1}{2}}, \text{ where } \hat{\mathbf{X}} \in \mathbb{R}^{n \times k}$$

- 9: Use the *projected* data matrix $\hat{\mathbf{X}}$ to get coefficient vector $\boldsymbol{\theta}$ and the intercept β :

$$\boldsymbol{\theta}, \beta \leftarrow \text{LogisticRegression}(\hat{\mathbf{X}}, \mathbf{y})$$

return Projection matrix Λ , coefficient vector $\boldsymbol{\theta}$ and intercept β .

Algorithm 5 Inductive Conformal Predictions with binary nonconformity score.

Input: A calibration dataset $\{(z_i, y_i)\}_{i=1}^n \subseteq \mathcal{D}_{cal}$, where $n = |\mathcal{D}_{cal}|$, $y_i \in \{-1, 1\}$ and $z_i \in \mathbb{R}^{L \times d}$ ($L = 1$ in a single instance setup); a confidence level α , a pretrained binary classifier g , and a new sample z_{new} .

Output: Prediction set \mathcal{Y}_{new}

```

1: Initialize an empty score list  $S \leftarrow \emptyset$ 
2: for all samples  $\langle z_i, y_i \rangle \in \mathcal{D}_{cal}$  do
3:    $z_i = g(z_i)$   $\triangleright z_i \in \mathbb{R}$  is a score
4:    $s_i = \exp(-y_i \cdot z_i)$   $\triangleright$  Nonconformity score, see Eq. 11.
5:    $S \leftarrow S \cup \{s_i\}$ 
6: end for
7: Compute a score for the new samples:  $z_{new} = g(z_{new})$ 
8: Initialize empty prediction set  $\mathcal{Y}_{new} \leftarrow \emptyset$ 
9: for all  $y \in \{-1, 1\}$  do
10:   $s_{new} = \exp(-y \cdot z_{new})$ 
11:   $\psi_y = \frac{\mathbb{I}(s_{new} < s_i) : \forall s_i \in S}{|S|}$   $\triangleright \mathbb{I}(\cdot)$  is an indicator function
12:  if  $\psi_y > 1 - \alpha$  then
13:     $\mathcal{Y}_{new} \leftarrow \mathcal{Y}_{new} \cup \{y\}$ 
14:  end if
15: end for
16: return Prediction set  $\mathcal{Y}_{new}$ 

```

Algorithm 6 Inductive Conformal Predictions with multiclass nonconformity score.

Input: A calibration dataset $\{(z_i, y_i)\}_{i=1}^n \subseteq \mathcal{D}_{cal}$, where $n = |\mathcal{D}_{cal}|$, $y_i \in \{1, \dots, K\}$ (K is the number of classes) and $z_i \in \mathbb{R}^{L \times d}$ ($L = 1$ in a single instance setup); a confidence level α , a pretrained multiclass classifier g , and a new sample z_{new} .

Output: Prediction set \mathcal{Y}_{new}

```

1: Initialize an empty score list  $S \leftarrow \emptyset$ 
2: for all samples  $\langle z_i, y_i \rangle \in \mathcal{D}_{cal}$  do
3:    $\mathbf{p} = g(z_i)$   $\triangleright \mathbf{p} \in \Delta^{K-1}$  is a vector of probabilities
4:    $p_z = \max_{j \neq y_i} p_j$   $\triangleright$  Maximum non-target probability, where  $j \in \{1, \dots, K\}$ 
5:    $d_p = p_{y_i} - p_z$   $\triangleright$  Probability margin
6:    $s_i = \frac{1-d_p}{2}$   $\triangleright$  Non-conformity score, see Eq. 12.
7:    $S \leftarrow S \cup \{s_i\}$ 
8: end for

9: Compute probabilities for the new samples:  $\mathbf{p}_{new} = g(z_{new})$ 
10: Initialize empty prediction set  $\mathcal{Y}_{new} \leftarrow \emptyset$ 
11: for all  $y \in \{1, \dots, K\}$  do
12:    $p_z = \max_{j \neq y} p_j$   $\triangleright$  where  $j \in \{1, \dots, K\}$  and  $p_j \in \mathbf{p}_{new}$ 
13:    $d_p = p_y - p_z$ 
14:    $s_{new} = \frac{1-d_p}{2}$ 
15:    $\psi_y = \frac{\mathbb{I}(s_{new} < s_i) : \forall s_i \in S}{|S|}$   $\triangleright \mathbb{I}(\cdot)$  is an indicator function
16:   if  $\psi_y > 1 - \alpha$  then
17:      $\mathcal{Y}_{new} \leftarrow \mathcal{Y}_{new} \cup \{y\}$ 
18:   end if
19: end for
20: return Prediction set  $\mathcal{Y}_{new}$ 

```

K More Tables on Classification Performance, Generalization Performance, and Confusion Matrices

In this section, we report detailed tables of the following results.

- Classification performance for all $\langle \text{model}, \text{dataset} \rangle$ pairs for the setting when only the last token’s representation is available, and the setting when the prediction is made over the full bag (see Supplementary Tables 9– 19),
- Generalization performance of multiclass sAwMIL across all datasets (see Supplementary Tables 21– 23),
- Confusion matrices for all $\langle \text{model}, \text{dataset} \rangle$ pairs (see Supplementary Tables 24– 31).

Table 9: **Classification performance of the multiclass sAwMIL probe across datasets and models** (evaluated on the *last token’s representation*). This probe is trained on the bag representation of the statements and evaluated using the representations of the last tokens. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	instance	sAwMIL	City Locations	0.99	0.99	1.00	16	0.52
Llama-3.2-3B	default	instance	sAwMIL	City Locations	0.98	0.99	1.00	10	0.37
Mistral-7B-v0.3	default	instance	sAwMIL	City Locations	0.99	0.99	1.00	17	0.55
Qwen-2.5-7B	default	instance	sAwMIL	City Locations	0.98	0.99	0.99	19	0.70
Qwen-2.5-14B	default	instance	sAwMIL	City Locations	0.98	0.99	0.99	28	0.60
Gemma-7B	default	instance	sAwMIL	City Locations	0.99	1.00	1.00	20	0.74
Gemma-2-9B	default	instance	sAwMIL	City Locations	0.99	1.00	1.00	21	0.51
Llama-3.1-8B	chat	instance	sAwMIL	City Locations	0.99	1.00	1.00	13	0.42
Llama-3.2-3B	chat	instance	sAwMIL	City Locations	0.98	0.99	0.99	12	0.44
Mistral-7B-v0.3	chat	instance	sAwMIL	City Locations	0.98	0.99	0.99	10	0.32
Qwen-2.5-7B	chat	instance	sAwMIL	City Locations	0.97	0.98	0.99	20	0.74
Qwen-2.5-14B	chat	instance	sAwMIL	City Locations	0.98	0.99	0.99	29	0.62
Gemma-7B	chat	instance	sAwMIL	City Locations	0.97	0.98	0.99	18	0.67
Gemma-2-9B	chat	instance	sAwMIL	City Locations	0.98	0.99	1.00	23	0.56
Bio-Medical-Llama	chat	instance	sAwMIL	City Locations	0.99	0.99	1.00	28	0.90
Llama3-Med42-8B	chat	instance	sAwMIL	City Locations	0.99	0.99	1.00	14	0.45
Llama-3-8B	default	instance	sAwMIL	Medical Indications	0.84	0.86	0.88	13	0.42
Llama-3.2-3B	default	instance	sAwMIL	Medical Indications	0.83	0.85	0.88	10	0.37
Mistral-7B-v0.3	default	instance	sAwMIL	Medical Indications	0.86	0.88	0.90	13	0.42
Qwen-2.5-7B	default	instance	sAwMIL	Medical Indications	0.83	0.86	0.88	16	0.59
Qwen-2.5-14B	default	instance	sAwMIL	Medical Indications	0.84	0.87	0.89	22	0.47
Gemma-7B	default	instance	sAwMIL	Medical Indications	0.84	0.86	0.89	17	0.63
Gemma-2-9B	default	instance	sAwMIL	Medical Indications	0.86	0.88	0.90	18	0.44
Llama-3.1-8B	chat	instance	sAwMIL	Medical Indications	0.85	0.87	0.89	18	0.58
Llama-3.2-3B	chat	instance	sAwMIL	Medical Indications	0.82	0.85	0.87	15	0.56
Mistral-7B-v0.3	chat	instance	sAwMIL	Medical Indications	0.86	0.88	0.90	16	0.52
Qwen-2.5-7B	chat	instance	sAwMIL	Medical Indications	0.83	0.85	0.87	17	0.63
Qwen-2.5-14B	chat	instance	sAwMIL	Medical Indications	0.86	0.88	0.91	23	0.49
Gemma-7B	chat	instance	sAwMIL	Medical Indications	0.80	0.83	0.85	15	0.56
Gemma-2-9B	chat	instance	sAwMIL	Medical Indications	0.86	0.88	0.90	21	0.51
Bio-Medical-Llama	chat	instance	sAwMIL	Medical Indications	0.85	0.87	0.89	11	0.35
Llama3-Med42-8B	chat	instance	sAwMIL	Medical Indications	0.85	0.88	0.90	8	0.26
Llama-3-8B	default	instance	sAwMIL	Word Definitions	0.86	0.87	0.89	13	0.42
Llama-3.2-3B	default	instance	sAwMIL	Word Definitions	0.83	0.85	0.87	10	0.37
Mistral-7B-v0.3	default	instance	sAwMIL	Word Definitions	0.85	0.87	0.89	13	0.42
Qwen-2.5-7B	default	instance	sAwMIL	Word Definitions	0.86	0.88	0.89	16	0.59
Qwen-2.5-14B	default	instance	sAwMIL	Word Definitions	0.86	0.87	0.89	21	0.45
Gemma-7B	default	instance	sAwMIL	Word Definitions	0.83	0.85	0.87	14	0.52
Gemma-2-9B	default	instance	sAwMIL	Word Definitions	0.88	0.90	0.91	17	0.41
Llama-3.1-8B	chat	instance	sAwMIL	Word Definitions	0.92	0.93	0.94	14	0.45
Llama-3.2-3B	chat	instance	sAwMIL	Word Definitions	0.85	0.86	0.88	12	0.44
Mistral-7B-v0.3	chat	instance	sAwMIL	Word Definitions	0.88	0.90	0.91	11	0.35
Qwen-2.5-7B	chat	instance	sAwMIL	Word Definitions	0.87	0.89	0.90	18	0.67
Qwen-2.5-14B	chat	instance	sAwMIL	Word Definitions	0.91	0.93	0.94	24	0.51
Gemma-7B	chat	instance	sAwMIL	Word Definitions	0.81	0.83	0.85	22	0.81
Gemma-2-9B	chat	instance	sAwMIL	Word Definitions	0.91	0.92	0.94	19	0.46
Bio-Medical-Llama	chat	instance	sAwMIL	Word Definitions	0.88	0.90	0.91	13	0.42
Llama3-Med42-8B	chat	instance	sAwMIL	Word Definitions	0.90	0.92	0.93	14	0.45

Table 10: **Classification performance of the multiclass sAwMIL probe across datasets and models** (evaluated on the *full bag*). This probe is trained and evaluated on the bag representation of the statements. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	bag	sAwMIL	City Locations	0.99	0.99	1.00	16	0.52
Llama-3.2-3B	default	bag	sAwMIL	City Locations	0.98	0.99	1.00	10	0.37
Mistral-7B-v0.3	default	bag	sAwMIL	City Locations	0.99	0.99	1.00	17	0.55
Qwen-2.5-7B	default	bag	sAwMIL	City Locations	0.98	0.99	0.99	19	0.70
Qwen-2.5-14B	default	bag	sAwMIL	City Locations	0.98	0.99	0.99	28	0.60
Gemma-7B	default	bag	sAwMIL	City Locations	0.99	1.00	1.00	20	0.74
Gemma-2-9B	default	bag	sAwMIL	City Locations	0.99	1.00	1.00	21	0.51
Llama-3.1-8B	chat	bag	sAwMIL	City Locations	0.99	1.00	1.00	13	0.42
Llama-3.2-3B	chat	bag	sAwMIL	City Locations	0.98	0.99	0.99	12	0.44
Mistral-7B-v0.3	chat	bag	sAwMIL	City Locations	0.98	0.99	0.99	10	0.32
Qwen-2.5-7B	chat	bag	sAwMIL	City Locations	0.97	0.98	0.99	20	0.74
Qwen-2.5-14B	chat	bag	sAwMIL	City Locations	0.98	0.99	0.99	29	0.62
Gemma-7B	chat	bag	sAwMIL	City Locations	0.97	0.98	0.99	18	0.67
Gemma-2-9B	chat	bag	sAwMIL	City Locations	0.98	0.99	1.00	23	0.56
Bio-Medical-Llama	chat	bag	sAwMIL	City Locations	0.99	0.99	1.00	28	0.90
Llama3-Med42-8B	chat	bag	sAwMIL	City Locations	0.99	0.99	1.00	14	0.45
Llama-3-8B	default	bag	sAwMIL	Medical Indications	0.83	0.86	0.88	13	0.42
Llama-3.2-3B	default	bag	sAwMIL	Medical Indications	0.83	0.85	0.88	10	0.37
Mistral-7B-v0.3	default	bag	sAwMIL	Medical Indications	0.86	0.88	0.90	13	0.42
Qwen-2.5-7B	default	bag	sAwMIL	Medical Indications	0.83	0.86	0.88	16	0.59
Qwen-2.5-14B	default	bag	sAwMIL	Medical Indications	0.84	0.87	0.89	22	0.47
Gemma-7B	default	bag	sAwMIL	Medical Indications	0.84	0.86	0.89	17	0.63
Gemma-2-9B	default	bag	sAwMIL	Medical Indications	0.86	0.88	0.90	18	0.44
Llama-3.1-8B	chat	bag	sAwMIL	Medical Indications	0.85	0.87	0.89	18	0.58
Llama-3.2-3B	chat	bag	sAwMIL	Medical Indications	0.83	0.85	0.87	15	0.56
Mistral-7B-v0.3	chat	bag	sAwMIL	Medical Indications	0.86	0.88	0.90	16	0.52
Qwen-2.5-7B	chat	bag	sAwMIL	Medical Indications	0.82	0.85	0.87	17	0.63
Qwen-2.5-14B	chat	bag	sAwMIL	Medical Indications	0.86	0.88	0.90	23	0.49
Gemma-7B	chat	bag	sAwMIL	Medical Indications	0.80	0.83	0.85	15	0.56
Gemma-2-9B	chat	bag	sAwMIL	Medical Indications	0.86	0.88	0.90	21	0.51
Bio-Medical-Llama	chat	bag	sAwMIL	Medical Indications	0.85	0.87	0.89	11	0.35
Llama3-Med42-8B	chat	bag	sAwMIL	Medical Indications	0.85	0.88	0.90	8	0.26
Llama-3-8B	default	bag	sAwMIL	Word Definitions	0.85	0.87	0.89	13	0.42
Llama-3.2-3B	default	bag	sAwMIL	Word Definitions	0.83	0.85	0.87	10	0.37
Mistral-7B-v0.3	default	bag	sAwMIL	Word Definitions	0.85	0.87	0.89	13	0.42
Qwen-2.5-7B	default	bag	sAwMIL	Word Definitions	0.86	0.88	0.89	16	0.59
Qwen-2.5-14B	default	bag	sAwMIL	Word Definitions	0.85	0.87	0.89	21	0.45
Gemma-7B	default	bag	sAwMIL	Word Definitions	0.83	0.85	0.86	14	0.52
Gemma-2-9B	default	bag	sAwMIL	Word Definitions	0.88	0.90	0.91	17	0.41
Llama-3.1-8B	chat	bag	sAwMIL	Word Definitions	0.92	0.93	0.95	14	0.45
Llama-3.2-3B	chat	bag	sAwMIL	Word Definitions	0.85	0.86	0.88	12	0.44
Mistral-7B-v0.3	chat	bag	sAwMIL	Word Definitions	0.88	0.90	0.91	11	0.35
Qwen-2.5-7B	chat	bag	sAwMIL	Word Definitions	0.87	0.89	0.90	18	0.67
Qwen-2.5-14B	chat	bag	sAwMIL	Word Definitions	0.91	0.93	0.94	24	0.51
Gemma-7B	chat	bag	sAwMIL	Word Definitions	0.81	0.83	0.85	22	0.81
Gemma-2-9B	chat	bag	sAwMIL	Word Definitions	0.91	0.92	0.94	19	0.46
Bio-Medical-Llama	chat	bag	sAwMIL	Word Definitions	0.88	0.90	0.92	13	0.42
Llama3-Med42-8B	chat	bag	sAwMIL	Word Definitions	0.90	0.92	0.93	14	0.45

Table 11: **Classification performance of the zero-shot prompting across datasets and models.** We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	-	Zero-Shot	City Locations	0.05	0.09	0.12	-	-
Llama-3.2-3B	default	-	Zero-Shot	City Locations	0.09	0.12	0.15	-	-
Mistral-7B-v0.3	default	-	Zero-Shot	City Locations	0.00	0.00	0.00	-	-
Qwen-2.5-7B	default	-	Zero-Shot	City Locations	0.53	0.56	0.58	-	-
Qwen-2.5-14B	default	-	Zero-Shot	City Locations	0.61	0.63	0.65	-	-
Gemma-7B	default	-	Zero-Shot	City Locations	0.12	0.15	0.18	-	-
Gemma-2-9B	default	-	Zero-Shot	City Locations	0.50	0.53	0.56	-	-
Llama-3.1-8B	chat	-	Zero-Shot	City Locations	0.54	0.56	0.58	-	-
Llama-3.2-3B	chat	-	Zero-Shot	City Locations	0.45	0.47	0.50	-	-
Mistral-7B-v0.3	chat	-	Zero-Shot	City Locations	0.41	0.44	0.47	-	-
Qwen-2.5-7B	chat	-	Zero-Shot	City Locations	0.53	0.55	0.57	-	-
Qwen-2.5-14B	chat	-	Zero-Shot	City Locations	0.80	0.82	0.84	-	-
Gemma-7B	chat	-	Zero-Shot	City Locations	0.44	0.46	0.49	-	-
Gemma-2-9B	chat	-	Zero-Shot	City Locations	0.57	0.59	0.61	-	-
Bio-Medical-Llama	chat	-	Zero-Shot	City Locations	0.48	0.51	0.53	-	-
Llama3-Med42-8B	chat	-	Zero-Shot	City Locations	0.53	0.55	0.57	-	-
Llama-3-8B	default	-	Zero-Shot	Medical Indications	0.18	0.22	0.25	-	-
Llama-3.2-3B	default	-	Zero-Shot	Medical Indications	0.04	0.08	0.12	-	-
Mistral-7B-v0.3	default	-	Zero-Shot	Medical Indications	0.00	0.02	0.04	-	-
Qwen-2.5-7B	default	-	Zero-Shot	Medical Indications	0.31	0.34	0.37	-	-
Qwen-2.5-14B	default	-	Zero-Shot	Medical Indications	0.38	0.41	0.44	-	-
Gemma-7B	default	-	Zero-Shot	Medical Indications	0.13	0.16	0.19	-	-
Gemma-2-9B	default	-	Zero-Shot	Medical Indications	0.28	0.32	0.35	-	-
Llama-3.1-8B	chat	-	Zero-Shot	Medical Indications	0.34	0.37	0.40	-	-
Llama-3.2-3B	chat	-	Zero-Shot	Medical Indications	0.22	0.25	0.28	-	-
Mistral-7B-v0.3	chat	-	Zero-Shot	Medical Indications	0.15	0.19	0.22	-	-
Qwen-2.5-7B	chat	-	Zero-Shot	Medical Indications	0.30	0.33	0.36	-	-
Qwen-2.5-14B	chat	-	Zero-Shot	Medical Indications	0.49	0.52	0.55	-	-
Gemma-7B	chat	-	Zero-Shot	Medical Indications	0.27	0.30	0.34	-	-
Gemma-2-9B	chat	-	Zero-Shot	Medical Indications	0.41	0.44	0.47	-	-
Bio-Medical-Llama	chat	-	Zero-Shot	Medical Indications	0.30	0.33	0.37	-	-
Llama3-Med42-8B	chat	-	Zero-Shot	Medical Indications	0.38	0.41	0.45	-	-
Llama-3-8B	default	-	Zero-Shot	Word Definitions	0.03	0.06	0.09	-	-
Llama-3.2-3B	default	-	Zero-Shot	Word Definitions	-0.01	0.02	0.04	-	-
Mistral-7B-v0.3	default	-	Zero-Shot	Word Definitions	0.00	0.01	0.01	-	-
Qwen-2.5-7B	default	-	Zero-Shot	Word Definitions	0.12	0.14	0.17	-	-
Qwen-2.5-14B	default	-	Zero-Shot	Word Definitions	0.27	0.30	0.33	-	-
Gemma-7B	default	-	Zero-Shot	Word Definitions	-0.03	0.00	0.03	-	-
Gemma-2-9B	default	-	Zero-Shot	Word Definitions	0.05	0.07	0.10	-	-
Llama-3.1-8B	chat	-	Zero-Shot	Word Definitions	0.17	0.20	0.23	-	-
Llama-3.2-3B	chat	-	Zero-Shot	Word Definitions	0.05	0.08	0.11	-	-
Mistral-7B-v0.3	chat	-	Zero-Shot	Word Definitions	-0.09	-0.06	-0.03	-	-
Qwen-2.5-7B	chat	-	Zero-Shot	Word Definitions	0.13	0.16	0.18	-	-
Qwen-2.5-14B	chat	-	Zero-Shot	Word Definitions	0.32	0.35	0.37	-	-
Gemma-7B	chat	-	Zero-Shot	Word Definitions	0.11	0.14	0.17	-	-
Gemma-2-9B	chat	-	Zero-Shot	Word Definitions	0.24	0.27	0.30	-	-
Bio-Medical-Llama	chat	-	Zero-Shot	Word Definitions	0.10	0.14	0.17	-	-
Llama3-Med42-8B	chat	-	Zero-Shot	Word Definitions	0.15	0.17	0.20	-	-

Table 12: **Classification performance of the mean-difference probe with conformal prediction intervals (MD+CP) probe across datasets and models** (evaluated on the *last token’s representation*). This probe is trained and evaluated using the representations of the last tokens. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	instance	MD+CP	City Locations	0.64	0.66	0.69	19	0.61
Llama-3.2-3B	default	instance	MD+CP	City Locations	0.60	0.62	0.64	27	1.00
Mistral-7B-v0.3	default	instance	MD+CP	City Locations	0.65	0.68	0.70	24	0.77
Qwen-2.5-7B	default	instance	MD+CP	City Locations	0.57	0.60	0.62	22	0.81
Qwen-2.5-14B	default	instance	MD+CP	City Locations	0.62	0.65	0.67	32	0.68
Gemma-7B	default	instance	MD+CP	City Locations	0.66	0.68	0.70	21	0.78
Gemma-2-9B	default	instance	MD+CP	City Locations	0.70	0.73	0.75	24	0.59
Llama-3.1-8B	chat	instance	MD+CP	City Locations	0.69	0.71	0.74	27	0.87
Llama-3.2-3B	chat	instance	MD+CP	City Locations	0.62	0.64	0.66	15	0.56
Mistral-7B-v0.3	chat	instance	MD+CP	City Locations	0.70	0.72	0.74	25	0.81
Qwen-2.5-7B	chat	instance	MD+CP	City Locations	0.60	0.63	0.65	19	0.70
Qwen-2.5-14B	chat	instance	MD+CP	City Locations	0.71	0.73	0.76	34	0.72
Gemma-7B	chat	instance	MD+CP	City Locations	0.61	0.64	0.66	19	0.70
Gemma-2-9B	chat	instance	MD+CP	City Locations	0.77	0.80	0.82	24	0.59
Bio-Medical-Llama	chat	instance	MD+CP	City Locations	0.69	0.72	0.74	18	0.58
Llama3-Med42-8B	chat	instance	MD+CP	City Locations	0.67	0.69	0.72	29	0.94
Llama-3-8B	default	instance	MD+CP	Medical Indications	0.50	0.54	0.57	25	0.81
Llama-3.2-3B	default	instance	MD+CP	Medical Indications	0.48	0.51	0.55	13	0.48
Mistral-7B-v0.3	default	instance	MD+CP	Medical Indications	0.51	0.54	0.57	15	0.48
Qwen-2.5-7B	default	instance	MD+CP	Medical Indications	0.54	0.57	0.60	23	0.85
Qwen-2.5-14B	default	instance	MD+CP	Medical Indications	0.56	0.60	0.63	42	0.89
Gemma-7B	default	instance	MD+CP	Medical Indications	0.53	0.56	0.60	21	0.78
Gemma-2-9B	default	instance	MD+CP	Medical Indications	0.52	0.56	0.59	16	0.39
Llama-3.1-8B	chat	instance	MD+CP	Medical Indications	0.52	0.55	0.58	20	0.65
Llama-3.2-3B	chat	instance	MD+CP	Medical Indications	0.48	0.52	0.56	14	0.52
Mistral-7B-v0.3	chat	instance	MD+CP	Medical Indications	0.50	0.54	0.57	24	0.77
Qwen-2.5-7B	chat	instance	MD+CP	Medical Indications	0.52	0.55	0.59	22	0.81
Qwen-2.5-14B	chat	instance	MD+CP	Medical Indications	0.51	0.54	0.57	35	0.74
Gemma-7B	chat	instance	MD+CP	Medical Indications	0.44	0.48	0.52	17	0.63
Gemma-2-9B	chat	instance	MD+CP	Medical Indications	0.54	0.57	0.60	26	0.63
Bio-Medical-Llama	chat	instance	MD+CP	Medical Indications	0.51	0.55	0.58	26	0.84
Llama3-Med42-8B	chat	instance	MD+CP	Medical Indications	0.52	0.56	0.59	17	0.55
Llama-3-8B	default	instance	MD+CP	Word Definitions	0.38	0.41	0.43	17	0.55
Llama-3.2-3B	default	instance	MD+CP	Word Definitions	0.36	0.38	0.41	8	0.30
Mistral-7B-v0.3	default	instance	MD+CP	Word Definitions	0.38	0.40	0.42	14	0.45
Qwen-2.5-7B	default	instance	MD+CP	Word Definitions	0.40	0.42	0.45	14	0.52
Qwen-2.5-14B	default	instance	MD+CP	Word Definitions	0.37	0.40	0.42	29	0.62
Gemma-7B	default	instance	MD+CP	Word Definitions	0.39	0.41	0.44	15	0.56
Gemma-2-9B	default	instance	MD+CP	Word Definitions	0.39	0.41	0.44	14	0.34
Llama-3.1-8B	chat	instance	MD+CP	Word Definitions	0.40	0.42	0.45	26	0.84
Llama-3.2-3B	chat	instance	MD+CP	Word Definitions	0.33	0.36	0.38	10	0.37
Mistral-7B-v0.3	chat	instance	MD+CP	Word Definitions	0.37	0.40	0.43	13	0.42
Qwen-2.5-7B	chat	instance	MD+CP	Word Definitions	0.40	0.42	0.44	18	0.67
Qwen-2.5-14B	chat	instance	MD+CP	Word Definitions	0.39	0.41	0.43	30	0.64
Gemma-7B	chat	instance	MD+CP	Word Definitions	0.38	0.41	0.44	13	0.48
Gemma-2-9B	chat	instance	MD+CP	Word Definitions	0.37	0.40	0.43	17	0.41
Bio-Medical-Llama	chat	instance	MD+CP	Word Definitions	0.38	0.40	0.42	20	0.65
Llama3-Med42-8B	chat	instance	MD+CP	Word Definitions	0.39	0.41	0.43	27	0.87

Table 13: **Classification performance of the mean-difference probe with conformal prediction intervals (MD+CP) probe across datasets and models** (evaluated on the *full bag*). This probe is trained on the representation of the last tokens and evaluated on the bag representation of the statements. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	default	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	default	bag	MD+CP	City Locations	0.08	0.11	0.14	9	0.29
Qwen-2.5-7B	default	bag	MD+CP	City Locations	0.06	0.08	0.11	12	0.44
Qwen-2.5-14B	default	bag	MD+CP	City Locations	0.07	0.10	0.11	37	0.79
Gemma-7B	default	bag	MD+CP	City Locations	0.29	0.31	0.34	7	0.26
Gemma-2-9B	default	bag	MD+CP	City Locations	0.01	0.04	0.06	33	0.80
Llama-3.1-8B	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	MD+CP	City Locations	0.06	0.08	0.10	16	0.52
Qwen-2.5-7B	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	MD+CP	City Locations	0.18	0.21	0.24	8	0.30
Gemma-2-9B	chat	bag	MD+CP	City Locations	0.26	0.30	0.34	39	0.95
Bio-Medical-Llama	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	MD+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama-3-8B	default	bag	MD+CP	Medical Indications	0.05	0.09	0.12	20	0.65
Llama-3.2-3B	default	bag	MD+CP	Medical Indications	0.21	0.25	0.28	11	0.41
Mistral-7B-v0.3	default	bag	MD+CP	Medical Indications	0.00	0.02	0.04	4	0.13
Qwen-2.5-7B	default	bag	MD+CP	Medical Indications	0.10	0.14	0.17	24	0.89
Qwen-2.5-14B	default	bag	MD+CP	Medical Indications	0.03	0.07	0.10	18	0.38
Gemma-7B	default	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Gemma-2-9B	default	bag	MD+CP	Medical Indications	0.04	0.08	0.11	12	0.29
Llama-3.1-8B	chat	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	MD+CP	Medical Indications	0.04	0.07	0.11	31	1.00
Qwen-2.5-7B	chat	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	MD+CP	Medical Indications	0.14	0.18	0.22	47	1.00
Gemma-7B	chat	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Gemma-2-9B	chat	bag	MD+CP	Medical Indications	0.00	0.04	0.06	23	0.56
Bio-Medical-Llama	chat	bag	MD+CP	Medical Indications	0.05	0.08	0.10	18	0.58
Llama3-Med42-8B	chat	bag	MD+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama-3-8B	default	bag	MD+CP	Word Definitions	0.01	0.03	0.06	22	0.71
Llama-3.2-3B	default	bag	MD+CP	Word Definitions	0.09	0.12	0.14	13	0.48
Mistral-7B-v0.3	default	bag	MD+CP	Word Definitions	0.02	0.04	0.07	20	0.65
Qwen-2.5-7B	default	bag	MD+CP	Word Definitions	-0.01	0.01	0.04	6	0.22
Qwen-2.5-14B	default	bag	MD+CP	Word Definitions	0.04	0.06	0.09	27	0.57
Gemma-7B	default	bag	MD+CP	Word Definitions	-0.01	0.02	0.05	24	0.89
Gemma-2-9B	default	bag	MD+CP	Word Definitions	0.05	0.07	0.10	19	0.46
Llama-3.1-8B	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	MD+CP	Word Definitions	-0.00	0.03	0.06	14	0.45
Qwen-2.5-7B	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	MD+CP	Word Definitions	-0.01	0.03	0.06	23	0.85
Gemma-2-9B	chat	bag	MD+CP	Word Definitions	0.18	0.22	0.27	38	0.93
Bio-Medical-Llama	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	MD+CP	Word Definitions	0.00	0.00	0.00	2	0.06

Table 14: **Classification performance of the TTPD with conformal prediction intervals (TTPD+CP) probe across datasets and models** (evaluated on the *last token’s representation*). This probe is trained and evaluated using the representations of the last tokens. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	instance	TTPD+CP	City Locations	0.53	0.55	0.57	13	0.42
Llama-3.2-3B	default	instance	TTPD+CP	City Locations	0.51	0.54	0.56	12	0.44
Mistral-7B-v0.3	default	instance	TTPD+CP	City Locations	0.65	0.67	0.69	11	0.35
Qwen-2.5-7B	default	instance	TTPD+CP	City Locations	0.54	0.57	0.59	18	0.67
Qwen-2.5-14B	default	instance	TTPD+CP	City Locations	0.55	0.58	0.61	28	0.60
Gemma-7B	default	instance	TTPD+CP	City Locations	0.56	0.58	0.59	12	0.44
Gemma-2-9B	default	instance	TTPD+CP	City Locations	0.59	0.61	0.63	17	0.41
Llama-3.1-8B	chat	instance	TTPD+CP	City Locations	0.83	0.85	0.88	17	0.55
Llama-3.2-3B	chat	instance	TTPD+CP	City Locations	0.45	0.47	0.50	10	0.37
Mistral-7B-v0.3	chat	instance	TTPD+CP	City Locations	0.77	0.80	0.83	14	0.45
Qwen-2.5-7B	chat	instance	TTPD+CP	City Locations	0.69	0.71	0.74	15	0.56
Qwen-2.5-14B	chat	instance	TTPD+CP	City Locations	0.75	0.78	0.80	36	0.77
Gemma-7B	chat	instance	TTPD+CP	City Locations	0.62	0.65	0.67	17	0.63
Gemma-2-9B	chat	instance	TTPD+CP	City Locations	0.58	0.64	0.69	13	0.32
Bio-Medical-Llama	chat	instance	TTPD+CP	City Locations	0.67	0.69	0.72	13	0.42
Llama3-Med42-8B	chat	instance	TTPD+CP	City Locations	0.93	0.95	0.96	16	0.52
Llama-3-8B	default	instance	TTPD+CP	Medical Indications	0.55	0.59	0.63	11	0.35
Llama-3.2-3B	default	instance	TTPD+CP	Medical Indications	0.51	0.54	0.57	11	0.41
Mistral-7B-v0.3	default	instance	TTPD+CP	Medical Indications	0.62	0.65	0.69	13	0.42
Qwen-2.5-7B	default	instance	TTPD+CP	Medical Indications	0.57	0.60	0.64	14	0.52
Qwen-2.5-14B	default	instance	TTPD+CP	Medical Indications	0.63	0.66	0.70	23	0.49
Gemma-7B	default	instance	TTPD+CP	Medical Indications	0.57	0.60	0.64	13	0.48
Gemma-2-9B	default	instance	TTPD+CP	Medical Indications	0.60	0.64	0.67	17	0.41
Llama-3.1-8B	chat	instance	TTPD+CP	Medical Indications	0.69	0.73	0.77	16	0.52
Llama-3.2-3B	chat	instance	TTPD+CP	Medical Indications	0.54	0.58	0.62	11	0.41
Mistral-7B-v0.3	chat	instance	TTPD+CP	Medical Indications	0.62	0.65	0.68	15	0.48
Qwen-2.5-7B	chat	instance	TTPD+CP	Medical Indications	0.50	0.54	0.59	14	0.52
Qwen-2.5-14B	chat	instance	TTPD+CP	Medical Indications	0.63	0.67	0.70	28	0.60
Gemma-7B	chat	instance	TTPD+CP	Medical Indications	0.48	0.52	0.56	17	0.63
Gemma-2-9B	chat	instance	TTPD+CP	Medical Indications	0.58	0.62	0.66	24	0.59
Bio-Medical-Llama	chat	instance	TTPD+CP	Medical Indications	0.54	0.58	0.62	16	0.52
Llama3-Med42-8B	chat	instance	TTPD+CP	Medical Indications	0.68	0.71	0.75	16	0.52
Llama-3-8B	default	instance	TTPD+CP	Word Definitions	0.36	0.39	0.42	9	0.29
Llama-3.2-3B	default	instance	TTPD+CP	Word Definitions	0.38	0.41	0.43	11	0.41
Mistral-7B-v0.3	default	instance	TTPD+CP	Word Definitions	0.40	0.43	0.46	15	0.48
Qwen-2.5-7B	default	instance	TTPD+CP	Word Definitions	0.34	0.37	0.41	19	0.70
Qwen-2.5-14B	default	instance	TTPD+CP	Word Definitions	0.38	0.42	0.46	41	0.87
Gemma-7B	default	instance	TTPD+CP	Word Definitions	0.41	0.43	0.46	15	0.56
Gemma-2-9B	default	instance	TTPD+CP	Word Definitions	0.44	0.46	0.48	17	0.41
Llama-3.1-8B	chat	instance	TTPD+CP	Word Definitions	0.44	0.47	0.50	19	0.61
Llama-3.2-3B	chat	instance	TTPD+CP	Word Definitions	0.32	0.35	0.38	12	0.44
Mistral-7B-v0.3	chat	instance	TTPD+CP	Word Definitions	0.39	0.42	0.45	13	0.42
Qwen-2.5-7B	chat	instance	TTPD+CP	Word Definitions	0.38	0.41	0.44	15	0.56
Qwen-2.5-14B	chat	instance	TTPD+CP	Word Definitions	0.42	0.45	0.47	46	0.98
Gemma-7B	chat	instance	TTPD+CP	Word Definitions	0.36	0.39	0.42	16	0.59
Gemma-2-9B	chat	instance	TTPD+CP	Word Definitions	0.36	0.40	0.43	22	0.54
Bio-Medical-Llama	chat	instance	TTPD+CP	Word Definitions	0.31	0.35	0.38	13	0.42
Llama3-Med42-8B	chat	instance	TTPD+CP	Word Definitions	0.41	0.44	0.47	16	0.52

Table 15: **Classification performance of the TTPD with conformal prediction intervals (TTPD+CP) probe across datasets and models** (evaluated on the *full bag*). This probe is trained on the representation of the last tokens and evaluated on the bag representation of the statements. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	bag	TTPD+CP	City Locations	0.52	0.54	0.56	13	0.42
Llama-3.2-3B	default	bag	TTPD+CP	City Locations	0.47	0.54	0.62	24	0.89
Mistral-7B-v0.3	default	bag	TTPD+CP	City Locations	0.49	0.51	0.52	15	0.48
Qwen-2.5-7B	default	bag	TTPD+CP	City Locations	0.49	0.53	0.56	17	0.63
Qwen-2.5-14B	default	bag	TTPD+CP	City Locations	0.53	0.55	0.58	20	0.43
Gemma-7B	default	bag	TTPD+CP	City Locations	0.56	0.58	0.60	16	0.59
Gemma-2-9B	default	bag	TTPD+CP	City Locations	0.50	0.53	0.55	21	0.51
Llama-3.1-8B	chat	bag	TTPD+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	TTPD+CP	City Locations	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	TTPD+CP	City Locations	0.00	0.00	0.00	2	0.06
Qwen-2.5-7B	chat	bag	TTPD+CP	City Locations	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	TTPD+CP	City Locations	0.55	0.59	0.63	35	0.74
Gemma-7B	chat	bag	TTPD+CP	City Locations	0.04	0.09	0.15	6	0.22
Gemma-2-9B	chat	bag	TTPD+CP	City Locations	0.00	0.00	0.00	3	0.07
Bio-Medical-Llama	chat	bag	TTPD+CP	City Locations	-0.01	0.04	0.09	4	0.13
Llama3-Med42-8B	chat	bag	TTPD+CP	City Locations	0.11	0.13	0.15	4	0.13
Llama-3-8B	default	bag	TTPD+CP	Medical Indications	0.55	0.61	0.67	24	0.77
Llama-3.2-3B	default	bag	TTPD+CP	Medical Indications	0.40	0.44	0.48	13	0.48
Mistral-7B-v0.3	default	bag	TTPD+CP	Medical Indications	0.07	0.22	0.33	17	0.55
Qwen-2.5-7B	default	bag	TTPD+CP	Medical Indications	0.09	0.17	0.26	7	0.26
Qwen-2.5-14B	default	bag	TTPD+CP	Medical Indications	-0.03	0.08	0.18	7	0.15
Gemma-7B	default	bag	TTPD+CP	Medical Indications	0.05	0.13	0.20	5	0.19
Gemma-2-9B	default	bag	TTPD+CP	Medical Indications	-0.04	0.09	0.22	4	0.10
Llama-3.1-8B	chat	bag	TTPD+CP	Medical Indications	0.10	0.16	0.22	5	0.16
Llama-3.2-3B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	TTPD+CP	Medical Indications	0.09	0.13	0.18	8	0.26
Qwen-2.5-7B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Gemma-2-9B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	3	0.07
Bio-Medical-Llama	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	TTPD+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama-3-8B	default	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	default	bag	TTPD+CP	Word Definitions	0.30	0.36	0.41	19	0.70
Mistral-7B-v0.3	default	bag	TTPD+CP	Word Definitions	0.14	0.19	0.25	17	0.55
Qwen-2.5-7B	default	bag	TTPD+CP	Word Definitions	0.17	0.21	0.26	11	0.41
Qwen-2.5-14B	default	bag	TTPD+CP	Word Definitions	0.05	0.07	0.10	20	0.43
Gemma-7B	default	bag	TTPD+CP	Word Definitions	0.30	0.34	0.38	20	0.74
Gemma-2-9B	default	bag	TTPD+CP	Word Definitions	0.31	0.35	0.39	26	0.63
Llama-3.1-8B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Qwen-2.5-7B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Gemma-2-9B	chat	bag	TTPD+CP	Word Definitions	0.33	0.37	0.40	40	0.98
Bio-Medical-Llama	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	TTPD+CP	Word Definitions	0.00	0.00	0.00	2	0.06

Table 16: **Classification performance of the supervised PCA with conformal prediction intervals (sPCA+CP) probe across datasets and models** (evaluated on the *last token's representation*). This probe is trained and evaluated using the representations of the last tokens. We report Matthew's Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The 'Setting' column specifies the evaluation setting. The 'Layer' column specifies the layer at which the probe achieved the highest MCC, and 'Depth' denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	instance	sPCA+CP	City Locations	0.72	0.74	0.77	26	0.84
Llama-3.2-3B	default	instance	sPCA+CP	City Locations	0.65	0.67	0.70	25	0.93
Mistral-7B-v0.3	default	instance	sPCA+CP	City Locations	0.70	0.73	0.75	19	0.61
Qwen-2.5-7B	default	instance	sPCA+CP	City Locations	0.67	0.69	0.72	21	0.78
Qwen-2.5-14B	default	instance	sPCA+CP	City Locations	0.71	0.73	0.76	25	0.53
Gemma-7B	default	instance	sPCA+CP	City Locations	0.70	0.73	0.75	18	0.67
Gemma-2-9B	default	instance	sPCA+CP	City Locations	0.68	0.71	0.73	21	0.51
Llama-3.1-8B	chat	instance	sPCA+CP	City Locations	0.77	0.80	0.82	22	0.71
Llama-3.2-3B	chat	instance	sPCA+CP	City Locations	0.60	0.62	0.64	12	0.44
Mistral-7B-v0.3	chat	instance	sPCA+CP	City Locations	0.71	0.73	0.76	20	0.65
Qwen-2.5-7B	chat	instance	sPCA+CP	City Locations	0.64	0.66	0.68	17	0.63
Qwen-2.5-14B	chat	instance	sPCA+CP	City Locations	0.69	0.71	0.73	39	0.83
Gemma-7B	chat	instance	sPCA+CP	City Locations	0.63	0.65	0.67	19	0.70
Gemma-2-9B	chat	instance	sPCA+CP	City Locations	0.74	0.77	0.79	21	0.51
Bio-Medical-Llama	chat	instance	sPCA+CP	City Locations	0.66	0.69	0.71	13	0.42
Llama3-Med42-8B	chat	instance	sPCA+CP	City Locations	0.81	0.84	0.86	31	1.00
Llama-3-8B	default	instance	sPCA+CP	Medical Indications	0.55	0.58	0.61	8	0.26
Llama-3.2-3B	default	instance	sPCA+CP	Medical Indications	0.54	0.58	0.61	15	0.56
Mistral-7B-v0.3	default	instance	sPCA+CP	Medical Indications	0.57	0.60	0.63	17	0.55
Qwen-2.5-7B	default	instance	sPCA+CP	Medical Indications	0.56	0.59	0.62	21	0.78
Qwen-2.5-14B	default	instance	sPCA+CP	Medical Indications	0.59	0.63	0.67	42	0.89
Gemma-7B	default	instance	sPCA+CP	Medical Indications	0.58	0.61	0.64	15	0.56
Gemma-2-9B	default	instance	sPCA+CP	Medical Indications	0.56	0.58	0.61	18	0.44
Llama-3.1-8B	chat	instance	sPCA+CP	Medical Indications	0.55	0.58	0.61	16	0.52
Llama-3.2-3B	chat	instance	sPCA+CP	Medical Indications	0.54	0.58	0.61	14	0.52
Mistral-7B-v0.3	chat	instance	sPCA+CP	Medical Indications	0.57	0.60	0.64	31	1.00
Qwen-2.5-7B	chat	instance	sPCA+CP	Medical Indications	0.52	0.55	0.58	26	0.96
Qwen-2.5-14B	chat	instance	sPCA+CP	Medical Indications	0.56	0.59	0.61	28	0.60
Gemma-7B	chat	instance	sPCA+CP	Medical Indications	0.47	0.51	0.55	23	0.85
Gemma-2-9B	chat	instance	sPCA+CP	Medical Indications	0.60	0.63	0.67	27	0.66
Bio-Medical-Llama	chat	instance	sPCA+CP	Medical Indications	0.54	0.57	0.60	31	1.00
Llama3-Med42-8B	chat	instance	sPCA+CP	Medical Indications	0.57	0.60	0.63	28	0.90
Llama-3-8B	default	instance	sPCA+CP	Word Definitions	0.39	0.41	0.43	12	0.39
Llama-3.2-3B	default	instance	sPCA+CP	Word Definitions	0.40	0.42	0.45	10	0.37
Mistral-7B-v0.3	default	instance	sPCA+CP	Word Definitions	0.39	0.41	0.43	13	0.42
Qwen-2.5-7B	default	instance	sPCA+CP	Word Definitions	0.41	0.43	0.46	22	0.81
Qwen-2.5-14B	default	instance	sPCA+CP	Word Definitions	0.39	0.42	0.44	39	0.83
Gemma-7B	default	instance	sPCA+CP	Word Definitions	0.38	0.41	0.43	11	0.41
Gemma-2-9B	default	instance	sPCA+CP	Word Definitions	0.41	0.43	0.46	18	0.44
Llama-3.1-8B	chat	instance	sPCA+CP	Word Definitions	0.40	0.42	0.44	13	0.42
Llama-3.2-3B	chat	instance	sPCA+CP	Word Definitions	0.36	0.38	0.41	9	0.33
Mistral-7B-v0.3	chat	instance	sPCA+CP	Word Definitions	0.38	0.40	0.42	17	0.55
Qwen-2.5-7B	chat	instance	sPCA+CP	Word Definitions	0.39	0.41	0.44	11	0.41
Qwen-2.5-14B	chat	instance	sPCA+CP	Word Definitions	0.39	0.41	0.43	39	0.83
Gemma-7B	chat	instance	sPCA+CP	Word Definitions	0.42	0.44	0.46	16	0.59
Gemma-2-9B	chat	instance	sPCA+CP	Word Definitions	0.40	0.42	0.45	19	0.46
Bio-Medical-Llama	chat	instance	sPCA+CP	Word Definitions	0.37	0.40	0.42	19	0.61
Llama3-Med42-8B	chat	instance	sPCA+CP	Word Definitions	0.39	0.41	0.42	17	0.55

Table 17: **Classification performance of the TTPD with conformal prediction intervals (TTPD+CP) probe across datasets and models** (evaluated on the *full bag*). This probe is trained on the representation of the last tokens and evaluated on the bag representation of the statements. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{0.025}	MCC	CI _{0.975}	Layer	Depth
Llama-3-8B	default	bag	sPCA+CP	City Locations	0.38	0.40	0.43	8	0.26
Llama-3.2-3B	default	bag	sPCA+CP	City Locations	0.55	0.57	0.60	27	1.00
Mistral-7B-v0.3	default	bag	sPCA+CP	City Locations	0.00	0.02	0.04	6	0.19
Qwen-2.5-7B	default	bag	sPCA+CP	City Locations	0.44	0.47	0.49	15	0.56
Qwen-2.5-14B	default	bag	sPCA+CP	City Locations	0.47	0.51	0.55	37	0.79
Gemma-7B	default	bag	sPCA+CP	City Locations	0.01	0.04	0.07	6	0.22
Gemma-2-9B	default	bag	sPCA+CP	City Locations	0.11	0.14	0.17	32	0.78
Llama-3.1-8B	chat	bag	sPCA+CP	City Locations	0.74	0.78	0.82	16	0.52
Llama-3.2-3B	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	sPCA+CP	City Locations	0.71	0.73	0.76	20	0.65
Qwen-2.5-7B	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	sPCA+CP	City Locations	0.62	0.64	0.66	17	0.63
Gemma-2-9B	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	3	0.07
Bio-Medical-Llama	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	sPCA+CP	City Locations	0.00	0.00	0.00	2	0.06
Llama-3-8B	default	bag	sPCA+CP	Medical Indications	0.54	0.57	0.60	11	0.35
Llama-3.2-3B	default	bag	sPCA+CP	Medical Indications	0.47	0.51	0.55	13	0.48
Mistral-7B-v0.3	default	bag	sPCA+CP	Medical Indications	0.54	0.58	0.61	15	0.48
Qwen-2.5-7B	default	bag	sPCA+CP	Medical Indications	0.13	0.16	0.20	19	0.70
Qwen-2.5-14B	default	bag	sPCA+CP	Medical Indications	0.37	0.40	0.44	34	0.72
Gemma-7B	default	bag	sPCA+CP	Medical Indications	0.05	0.11	0.18	6	0.22
Gemma-2-9B	default	bag	sPCA+CP	Medical Indications	0.26	0.30	0.33	15	0.37
Llama-3.1-8B	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Qwen-2.5-7B	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	sPCA+CP	Medical Indications	0.55	0.58	0.61	45	0.96
Gemma-7B	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.07
Gemma-2-9B	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	3	0.07
Bio-Medical-Llama	chat	bag	sPCA+CP	Medical Indications	0.00	0.00	0.00	2	0.06
Llama3-Med42-8B	chat	bag	sPCA+CP	Medical Indications	0.20	0.23	0.26	22	0.71
Llama-3-8B	default	bag	sPCA+CP	Word Definitions	0.01	0.03	0.04	9	0.29
Llama-3.2-3B	default	bag	sPCA+CP	Word Definitions	0.03	0.06	0.08	16	0.59
Mistral-7B-v0.3	default	bag	sPCA+CP	Word Definitions	0.14	0.17	0.19	19	0.61
Qwen-2.5-7B	default	bag	sPCA+CP	Word Definitions	0.01	0.04	0.06	6	0.22
Qwen-2.5-14B	default	bag	sPCA+CP	Word Definitions	0.02	0.05	0.07	14	0.30
Gemma-7B	default	bag	sPCA+CP	Word Definitions	0.06	0.08	0.11	20	0.74
Gemma-2-9B	default	bag	sPCA+CP	Word Definitions	0.05	0.08	0.10	28	0.68
Llama-3.1-8B	chat	bag	sPCA+CP	Word Definitions	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	sPCA+CP	Word Definitions	0.00	0.00	0.00	2	0.07
Mistral-7B-v0.3	chat	bag	sPCA+CP	Word Definitions	0.00	0.03	0.06	7	0.23
Qwen-2.5-7B	chat	bag	sPCA+CP	Word Definitions	0.34	0.37	0.39	25	0.93
Qwen-2.5-14B	chat	bag	sPCA+CP	Word Definitions	0.00	0.00	0.00	4	0.09
Gemma-7B	chat	bag	sPCA+CP	Word Definitions	0.17	0.19	0.22	20	0.74
Gemma-2-9B	chat	bag	sPCA+CP	Word Definitions	0.09	0.12	0.14	24	0.59
Bio-Medical-Llama	chat	bag	sPCA+CP	Word Definitions	0.03	0.06	0.09	7	0.23
Llama3-Med42-8B	chat	bag	sPCA+CP	Word Definitions	0.00	0.00	0.00	2	0.06

Table 18: **Classification performance of the multiclass SVM probe across datasets and models** (evaluated on the *last token’s representation*). This probe is trained and evaluated using the representations of the last tokens. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3.2-3B	default	instance	SVM	City Locations	0.95	0.96	0.97	9	0.33
Mistral-7B-v0.3	default	instance	SVM	City Locations	0.96	0.97	0.98	12	0.39
Qwen-2.5-7B	default	instance	SVM	City Locations	0.96	0.97	0.98	18	0.67
Qwen-2.5-14B	default	instance	SVM	City Locations	0.97	0.97	0.98	24	0.51
Gemma-7B	default	instance	SVM	City Locations	0.97	0.98	0.99	15	0.56
Gemma-2-9B	default	instance	SVM	City Locations	0.97	0.98	0.99	21	0.51
Llama-3.1-8B	chat	instance	SVM	City Locations	0.97	0.98	0.99	12	0.39
Llama-3.2-3B	chat	instance	SVM	City Locations	0.94	0.95	0.97	13	0.48
Mistral-7B-v0.3	chat	instance	SVM	City Locations	0.97	0.98	0.98	21	0.68
Qwen-2.5-7B	chat	instance	SVM	City Locations	0.95	0.96	0.97	20	0.74
Qwen-2.5-14B	chat	instance	SVM	City Locations	0.97	0.98	0.98	46	0.98
Gemma-7B	chat	instance	SVM	City Locations	0.95	0.96	0.97	15	0.56
Gemma-2-9B	chat	instance	SVM	City Locations	0.98	0.99	0.99	21	0.51
Bio-Medical-Llama	chat	instance	SVM	City Locations	0.96	0.97	0.98	14	0.45
Llama3-Med42-8B	chat	instance	SVM	City Locations	0.97	0.98	0.99	27	0.87
Llama-3-8B	default	instance	SVM	Medical Indications	0.82	0.84	0.86	14	0.45
Llama-3.2-3B	default	instance	SVM	Medical Indications	0.75	0.78	0.81	12	0.44
Mistral-7B-v0.3	default	instance	SVM	Medical Indications	0.81	0.84	0.86	11	0.35
Qwen-2.5-7B	default	instance	SVM	Medical Indications	0.80	0.82	0.84	17	0.63
Qwen-2.5-14B	default	instance	SVM	Medical Indications	0.83	0.85	0.87	24	0.51
Gemma-7B	default	instance	SVM	Medical Indications	0.80	0.82	0.85	16	0.59
Gemma-2-9B	default	instance	SVM	Medical Indications	0.82	0.84	0.86	21	0.51
Llama-3.1-8B	chat	instance	SVM	Medical Indications	0.81	0.83	0.86	16	0.52
Llama-3.2-3B	chat	instance	SVM	Medical Indications	0.74	0.76	0.79	12	0.44
Mistral-7B-v0.3	chat	instance	SVM	Medical Indications	0.79	0.82	0.84	14	0.45
Qwen-2.5-7B	chat	instance	SVM	Medical Indications	0.78	0.81	0.83	19	0.70
Qwen-2.5-14B	chat	instance	SVM	Medical Indications	0.83	0.85	0.88	30	0.64
Gemma-7B	chat	instance	SVM	Medical Indications	0.71	0.74	0.77	16	0.59
Gemma-2-9B	chat	instance	SVM	Medical Indications	0.80	0.82	0.84	20	0.49
Bio-Medical-Llama	chat	instance	SVM	Medical Indications	0.78	0.81	0.83	14	0.45
Llama3-Med42-8B	chat	instance	SVM	Medical Indications	0.81	0.83	0.85	25	0.81
Llama-3-8B	default	instance	SVM	Word Definitions	0.89	0.90	0.92	12	0.39
Llama-3.2-3B	default	instance	SVM	Word Definitions	0.86	0.88	0.90	12	0.44
Mistral-7B-v0.3	default	instance	SVM	Word Definitions	0.89	0.90	0.92	13	0.42
Qwen-2.5-7B	default	instance	SVM	Word Definitions	0.89	0.90	0.92	16	0.59
Qwen-2.5-14B	default	instance	SVM	Word Definitions	0.90	0.92	0.93	23	0.49
Gemma-7B	default	instance	SVM	Word Definitions	0.87	0.89	0.90	15	0.56
Gemma-2-9B	default	instance	SVM	Word Definitions	0.89	0.91	0.92	18	0.44
Llama-3.1-8B	chat	instance	SVM	Word Definitions	0.90	0.91	0.93	13	0.42
Llama-3.2-3B	chat	instance	SVM	Word Definitions	0.82	0.84	0.85	11	0.41
Mistral-7B-v0.3	chat	instance	SVM	Word Definitions	0.88	0.89	0.91	12	0.39
Qwen-2.5-7B	chat	instance	SVM	Word Definitions	0.89	0.90	0.92	18	0.67
Qwen-2.5-14B	chat	instance	SVM	Word Definitions	0.90	0.91	0.93	23	0.49
Gemma-7B	chat	instance	SVM	Word Definitions	0.87	0.88	0.90	17	0.63
Gemma-2-9B	chat	instance	SVM	Word Definitions	0.87	0.89	0.90	24	0.59
Bio-Medical-Llama	chat	instance	SVM	Word Definitions	0.86	0.88	0.90	13	0.42
Llama3-Med42-8B	chat	instance	SVM	Word Definitions	0.90	0.91	0.93	19	0.61

Table 19: **Classification performance of the multiclass SVM probe across datasets and models** (evaluated on the full bag). This probe is trained on the representation of the last tokens and evaluated on the bag representation of the statements. We report Matthew’s Correlation Coefficient (MCC) with the 95% confidence intervals. Confidence intervals are based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Setting’ column specifies the evaluation setting. The ‘Layer’ column specifies the layer at which the probe achieved the highest MCC, and ‘Depth’ denotes the relative depth of that layer within the model.

Model Name	Type	Setting	Probe	Dataset	CI _{.025}	MCC	CI _{.975}	Layer	Depth
Llama-3-8B	default	bag	SVM	City Locations	0.52	0.54	0.57	8	0.26
Llama-3.2-3B	default	bag	SVM	City Locations	0.71	0.73	0.75	27	1.00
Mistral-7B-v0.3	default	bag	SVM	City Locations	0.86	0.88	0.89	10	0.32
Qwen-2.5-7B	default	bag	SVM	City Locations	0.50	0.52	0.54	26	0.96
Qwen-2.5-14B	default	bag	SVM	City Locations	0.38	0.41	0.43	40	0.85
Gemma-7B	default	bag	SVM	City Locations	0.38	0.40	0.42	11	0.41
Gemma-2-9B	default	bag	SVM	City Locations	0.84	0.86	0.88	17	0.41
Llama-3.1-8B	chat	bag	SVM	City Locations	0.00	0.00	0.00	2	0.06
Llama-3.2-3B	chat	bag	SVM	City Locations	0.49	0.52	0.54	25	0.93
Mistral-7B-v0.3	chat	bag	SVM	City Locations	0.76	0.79	0.81	21	0.68
Qwen-2.5-7B	chat	bag	SVM	City Locations	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	SVM	City Locations	0.36	0.38	0.40	22	0.47
Gemma-7B	chat	bag	SVM	City Locations	0.42	0.44	0.46	16	0.59
Gemma-2-9B	chat	bag	SVM	City Locations	0.91	0.92	0.94	26	0.63
Bio-Medical-Llama	chat	bag	SVM	City Locations	0.06	0.09	0.12	6	0.19
Llama3-Med42-8B	chat	bag	SVM	City Locations	0.34	0.36	0.38	31	1.00
Llama-3-8B	default	bag	SVM	Medical Indications	0.44	0.48	0.51	20	0.65
Llama-3.2-3B	default	bag	SVM	Medical Indications	0.32	0.34	0.36	4	0.15
Mistral-7B-v0.3	default	bag	SVM	Medical Indications	0.65	0.68	0.71	11	0.35
Qwen-2.5-7B	default	bag	SVM	Medical Indications	0.18	0.21	0.25	27	1.00
Qwen-2.5-14B	default	bag	SVM	Medical Indications	0.23	0.26	0.28	12	0.26
Gemma-7B	default	bag	SVM	Medical Indications	0.38	0.40	0.42	23	0.85
Gemma-2-9B	default	bag	SVM	Medical Indications	0.56	0.58	0.61	18	0.44
Llama-3.1-8B	chat	bag	SVM	Medical Indications	0.29	0.32	0.35	18	0.58
Llama-3.2-3B	chat	bag	SVM	Medical Indications	0.35	0.38	0.41	6	0.22
Mistral-7B-v0.3	chat	bag	SVM	Medical Indications	0.59	0.62	0.66	21	0.68
Qwen-2.5-7B	chat	bag	SVM	Medical Indications	0.00	0.00	0.00	2	0.07
Qwen-2.5-14B	chat	bag	SVM	Medical Indications	0.40	0.43	0.45	46	0.98
Gemma-7B	chat	bag	SVM	Medical Indications	0.37	0.40	0.43	9	0.33
Gemma-2-9B	chat	bag	SVM	Medical Indications	0.60	0.63	0.66	19	0.46
Bio-Medical-Llama	chat	bag	SVM	Medical Indications	0.27	0.30	0.32	4	0.13
Llama3-Med42-8B	chat	bag	SVM	Medical Indications	0.52	0.54	0.56	31	1.00
Llama-3-8B	default	bag	SVM	Word Definitions	0.63	0.65	0.66	6	0.19
Llama-3.2-3B	default	bag	SVM	Word Definitions	0.50	0.52	0.54	9	0.33
Mistral-7B-v0.3	default	bag	SVM	Word Definitions	0.62	0.63	0.65	17	0.55
Qwen-2.5-7B	default	bag	SVM	Word Definitions	0.54	0.56	0.58	5	0.19
Qwen-2.5-14B	default	bag	SVM	Word Definitions	0.63	0.65	0.66	22	0.47
Gemma-7B	default	bag	SVM	Word Definitions	0.60	0.62	0.64	8	0.30
Gemma-2-9B	default	bag	SVM	Word Definitions	0.77	0.79	0.82	27	0.66
Llama-3.1-8B	chat	bag	SVM	Word Definitions	0.65	0.67	0.69	12	0.39
Llama-3.2-3B	chat	bag	SVM	Word Definitions	0.37	0.39	0.42	12	0.44
Mistral-7B-v0.3	chat	bag	SVM	Word Definitions	0.74	0.76	0.78	23	0.74
Qwen-2.5-7B	chat	bag	SVM	Word Definitions	-0.00	0.03	0.07	26	0.96
Qwen-2.5-14B	chat	bag	SVM	Word Definitions	0.26	0.29	0.31	35	0.74
Gemma-7B	chat	bag	SVM	Word Definitions	0.56	0.58	0.60	25	0.93
Gemma-2-9B	chat	bag	SVM	Word Definitions	0.81	0.83	0.85	26	0.63
Bio-Medical-Llama	chat	bag	SVM	Word Definitions	0.34	0.36	0.38	10	0.32
Llama3-Med42-8B	chat	bag	SVM	Word Definitions	0.28	0.30	0.32	22	0.71

Table 20: **Aggregated generalization performance.** Each probe is trained on one dataset and evaluated on the remaining two (e.g., a probe trained on *City Locations* is evaluated on *Medical Indications* and *Word Definitions*). The performances are averaged over all large language models (we pick only the performance from the best performing layers). The performance is measured by the Matthew’s Correlation Coefficient (MCC) with the standard error. The ‘Evaluation Setting’ column indicates whether the probe was evaluated using only the representation of the last token in the statement (*Instance-Level*) or the full bag of tokens (*Bag-Level*). Multiclass probes (i.e., SVM and sAwMIL) provide higher generalization performance in both settings, and the multiclass sAwMIL achieves the highest generalization performance.

Evaluation Setting	Probe	Training Dataset	MCC	Stand. Err.
Bag-Level	MD+CP	City Locations	0.02	0.00
Bag-Level	MD+CP	Medical Indications	0.03	0.01
Bag-Level	MD+CP	Word Definitions	0.04	0.02
Bag-Level	SVM	City Locations	0.24	0.03
Bag-Level	SVM	Medical Indications	0.37	0.04
Bag-Level	SVM	Word Definitions	0.40	0.04
Bag-Level	TTPD+CP	City Locations	0.11	0.03
Bag-Level	TTPD+CP	Medical Indications	0.09	0.03
Bag-Level	TTPD+CP	Word Definitions	0.13	0.04
Bag-Level	sAwMIL	City Locations	0.82	0.02
Bag-Level	sAwMIL	Medical Indications	0.88	0.01
Bag-Level	sAwMIL	Word Definitions	0.86	0.01
Bag-Level	sPCA+CP	City Locations	0.11	0.02
Bag-Level	sPCA+CP	Medical Indications	0.13	0.03
Bag-Level	sPCA+CP	Word Definitions	0.11	0.03
Instance-Level	MD+CP	City Locations	0.22	0.02
Instance-Level	MD+CP	Medical Indications	0.21	0.02
Instance-Level	MD+CP	Word Definitions	0.43	0.02
Instance-Level	SVM	City Locations	0.61	0.02
Instance-Level	SVM	Medical Indications	0.74	0.02
Instance-Level	SVM	Word Definitions	0.77	0.02
Instance-Level	TTPD+CP	City Locations	0.36	0.03
Instance-Level	TTPD+CP	Medical Indications	0.50	0.04
Instance-Level	TTPD+CP	Word Definitions	0.62	0.03
Instance-Level	sAwMIL	City Locations	0.82	0.02
Instance-Level	sAwMIL	Medical Indications	0.88	0.01
Instance-Level	sAwMIL	Word Definitions	0.86	0.01
Instance-Level	sPCA+CP	City Locations	0.38	0.03
Instance-Level	sPCA+CP	Medical Indications	0.42	0.03
Instance-Level	sPCA+CP	Word Definitions	0.54	0.01

Table 21: **Generalization performance of the multiclass sAwMIL trained on the City Locations dataset.** The performance is measured by the Matthew’s Correlation Coefficient (MCC) with 95% confidence intervals, based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Rel. Depth’ column specifies the relative depth of the layer where the multiclass sAwMIL probe achieves the best MCC score.

Model Name	Training Dataset	Test Dataset	CI _{.025}	MCC	CI _{.975}	Rel. Depth
Gemma-7B-it	City Locations	City Locations	0.93	0.94	0.95	0.59
Gemma-7B-it	City Locations	Medical Indications	0.43	0.46	0.49	0.63
Gemma-7B-it	City Locations	Word Definitions	0.44	0.47	0.50	0.63
Gemma-2-9B-it	City Locations	City Locations	0.96	0.97	0.98	0.56
Gemma-2-9B-it	City Locations	Medical Indications	0.67	0.70	0.73	0.44
Gemma-2-9B-it	City Locations	Word Definitions	0.69	0.71	0.74	0.49
Llama-3.2-3B-Instruct	City Locations	City Locations	0.95	0.96	0.97	0.52
Llama-3.2-3B-Instruct	City Locations	Medical Indications	0.59	0.62	0.65	0.56
Llama-3.2-3B-Instruct	City Locations	Word Definitions	0.62	0.65	0.67	0.48
Llama3-Med42-8B	City Locations	City Locations	0.96	0.97	0.98	0.42
Llama3-Med42-8B	City Locations	Medical Indications	0.75	0.78	0.81	0.90
Llama3-Med42-8B	City Locations	Word Definitions	0.75	0.78	0.80	0.45
Llama-3.1-8B-Instruct	City Locations	City Locations	0.96	0.97	0.98	0.48
Llama-3.1-8B-Instruct	City Locations	Medical Indications	0.73	0.75	0.78	0.52
Llama-3.1-8B-Instruct	City Locations	Word Definitions	0.79	0.81	0.83	0.42
Bio-Medical-Llama-3-8B	City Locations	City Locations	0.96	0.97	0.98	0.97
Bio-Medical-Llama-3-8B	City Locations	Medical Indications	0.59	0.62	0.65	0.42
Bio-Medical-Llama-3-8B	City Locations	Word Definitions	0.53	0.56	0.59	0.26
Mistral-7B-Instruct-v0.3	City Locations	City Locations	0.95	0.96	0.97	0.48
Mistral-7B-Instruct-v0.3	City Locations	Medical Indications	0.72	0.75	0.78	0.55
Mistral-7B-Instruct-v0.3	City Locations	Word Definitions	0.69	0.71	0.74	0.35
Qwen-2.5-7B-Instruct	City Locations	City Locations	0.94	0.95	0.96	0.67
Qwen-2.5-7B-Instruct	City Locations	Medical Indications	0.71	0.74	0.76	0.70
Qwen-2.5-7B-Instruct	City Locations	Word Definitions	0.62	0.64	0.67	0.70
Qwen-2.5-14B-Instruct	City Locations	City Locations	0.96	0.97	0.98	0.62
Qwen-2.5-14B-Instruct	City Locations	Medical Indications	0.74	0.77	0.80	0.64
Qwen-2.5-14B-Instruct	City Locations	Word Definitions	0.73	0.75	0.77	0.60
Gemma-7B	City Locations	City Locations	0.96	0.97	0.98	0.74
Gemma-7B	City Locations	Medical Indications	0.55	0.58	0.60	0.41
Gemma-7B	City Locations	Word Definitions	0.61	0.63	0.66	0.52
Gemma-2-9B	City Locations	City Locations	0.97	0.98	0.99	0.63
Gemma-2-9B	City Locations	Medical Indications	0.52	0.55	0.58	0.44
Gemma-2-9B	City Locations	Word Definitions	0.54	0.56	0.58	0.27
Llama-3.2-3B	City Locations	City Locations	0.95	0.96	0.97	0.37
Llama-3.2-3B	City Locations	Medical Indications	0.35	0.38	0.41	0.33
Llama-3.2-3B	City Locations	Word Definitions	0.48	0.50	0.52	0.30
Llama-3-8B	City Locations	City Locations	0.96	0.97	0.98	0.32
Llama-3-8B	City Locations	Medical Indications	0.59	0.62	0.65	0.35
Llama-3-8B	City Locations	Word Definitions	0.56	0.58	0.61	0.26
Mistral-7B-v0.3	City Locations	City Locations	0.96	0.97	0.98	0.42
Mistral-7B-v0.3	City Locations	Medical Indications	0.50	0.53	0.55	0.39
Mistral-7B-v0.3	City Locations	Word Definitions	0.58	0.61	0.63	0.39
Qwen-2.5-7B	City Locations	City Locations	0.94	0.95	0.96	0.70
Qwen-2.5-7B	City Locations	Medical Indications	0.48	0.52	0.55	0.59
Qwen-2.5-7B	City Locations	Word Definitions	0.51	0.54	0.57	0.59
Qwen-2.5-14B	City Locations	City Locations	0.95	0.96	0.97	0.79
Qwen-2.5-14B	City Locations	Medical Indications	0.59	0.62	0.65	0.45
Qwen-2.5-14B	City Locations	Word Definitions	0.62	0.64	0.67	0.43

Table 22: **Generalization performance of the multiclass sAwMIL trained on the *Medical Indications* dataset.** The performance is measured by the Matthew’s Correlation Coefficient (MCC) with 95% confidence intervals, based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC significant confidence intervals. The ‘Rel. Depth’ column specifies the relative depth of the layer where a multiclass sAwMIL probe achieves the best MCC score.

Model Name	Training Dataset	Test Dataset	CI _{.025}	MCC	CI _{.975}	Rel. Depth
Gemma-7B-it	Medical Indications	City Locations	0.90	0.91	0.93	0.70
Gemma-7B-it	Medical Indications	Medical Indications	0.65	0.68	0.72	0.59
Gemma-7B-it	Medical Indications	Word Definitions	0.64	0.67	0.69	0.59
Gemma-2-9B-it	Medical Indications	City Locations	0.91	0.92	0.94	0.63
Gemma-2-9B-it	Medical Indications	Medical Indications	0.81	0.83	0.85	0.49
Gemma-2-9B-it	Medical Indications	Word Definitions	0.78	0.80	0.82	0.61
Llama-3.2-3B-Instruct	Medical Indications	City Locations	0.80	0.82	0.85	0.41
Llama-3.2-3B-Instruct	Medical Indications	Medical Indications	0.73	0.76	0.79	0.48
Llama-3.2-3B-Instruct	Medical Indications	Word Definitions	0.65	0.67	0.70	0.48
Llama3-Med42-8B	Medical Indications	City Locations	0.95	0.96	0.97	0.52
Llama3-Med42-8B	Medical Indications	Medical Indications	0.80	0.83	0.85	0.45
Llama3-Med42-8B	Medical Indications	Word Definitions	0.76	0.78	0.80	0.45
Llama-3.1-8B-Instruct	Medical Indications	City Locations	0.92	0.93	0.94	0.55
Llama-3.1-8B-Instruct	Medical Indications	Medical Indications	0.81	0.83	0.85	0.55
Llama-3.1-8B-Instruct	Medical Indications	Word Definitions	0.77	0.79	0.81	0.42
Bio-Medical-Llama-3-8B	Medical Indications	City Locations	0.85	0.86	0.88	0.39
Bio-Medical-Llama-3-8B	Medical Indications	Medical Indications	0.78	0.81	0.83	0.81
Bio-Medical-Llama-3-8B	Medical Indications	Word Definitions	0.70	0.73	0.75	0.32
Mistral-7B-Instruct-v0.3	Medical Indications	City Locations	0.93	0.94	0.95	0.39
Mistral-7B-Instruct-v0.3	Medical Indications	Medical Indications	0.78	0.80	0.83	0.45
Mistral-7B-Instruct-v0.3	Medical Indications	Word Definitions	0.74	0.77	0.79	0.45
Qwen-2.5-7B-Instruct	Medical Indications	City Locations	0.78	0.80	0.83	0.52
Qwen-2.5-7B-Instruct	Medical Indications	Medical Indications	0.74	0.77	0.79	0.67
Qwen-2.5-7B-Instruct	Medical Indications	Word Definitions	0.64	0.67	0.69	0.63
Qwen-2.5-14B-Instruct	Medical Indications	City Locations	0.90	0.92	0.93	0.49
Qwen-2.5-14B-Instruct	Medical Indications	Medical Indications	0.79	0.82	0.84	0.57
Qwen-2.5-14B-Instruct	Medical Indications	Word Definitions	0.76	0.79	0.81	0.57
Gemma-7B	Medical Indications	City Locations	0.60	0.63	0.66	0.56
Gemma-7B	Medical Indications	Medical Indications	0.75	0.78	0.80	0.63
Gemma-7B	Medical Indications	Word Definitions	0.55	0.57	0.59	0.70
Gemma-2-9B	Medical Indications	City Locations	0.83	0.85	0.87	0.56
Gemma-2-9B	Medical Indications	Medical Indications	0.77	0.80	0.82	0.44
Gemma-2-9B	Medical Indications	Word Definitions	0.64	0.67	0.69	0.39
Llama-3.2-3B	Medical Indications	City Locations	0.49	0.51	0.53	0.41
Llama-3.2-3B	Medical Indications	Medical Indications	0.74	0.76	0.79	0.44
Llama-3.2-3B	Medical Indications	Word Definitions	0.59	0.62	0.64	0.52
Llama-3-8B	Medical Indications	City Locations	0.75	0.77	0.80	0.45
Llama-3-8B	Medical Indications	Medical Indications	0.77	0.80	0.83	0.39
Llama-3-8B	Medical Indications	Word Definitions	0.63	0.65	0.68	0.26
Mistral-7B-v0.3	Medical Indications	City Locations	0.58	0.60	0.63	0.42
Mistral-7B-v0.3	Medical Indications	Medical Indications	0.77	0.80	0.82	0.42
Mistral-7B-v0.3	Medical Indications	Word Definitions	0.63	0.65	0.68	0.42
Qwen-2.5-7B	Medical Indications	City Locations	0.80	0.82	0.84	0.67
Qwen-2.5-7B	Medical Indications	Medical Indications	0.76	0.78	0.81	0.63
Qwen-2.5-7B	Medical Indications	Word Definitions	0.59	0.62	0.65	0.74
Qwen-2.5-14B	Medical Indications	City Locations	0.80	0.82	0.84	0.70
Qwen-2.5-14B	Medical Indications	Medical Indications	0.76	0.79	0.82	0.60
Qwen-2.5-14B	Medical Indications	Word Definitions	0.70	0.72	0.75	0.60

Table 23: **Generalization performance of the multiclass sAwMIL trained on the Word Definitions dataset.** The performance is measured by the Matthew’s Correlation Coefficient (MCC) with 95% confidence intervals, based on bootstrapping with $n = 1,000$ samples. The **bold** values mark MCC with significant confidence intervals. The ‘Rel. Depth’ column specifies the relative depth of the layer where a multiclass sAwMIL probe achieves the best MCC score.

Model Name	Training Dataset	Test Dataset	CI _{.025}	MCC	CI _{.975}	Rel. Depth
Gemma-7B-it	Word Definitions	City Locations	0.90	0.92	0.93	0.70
Gemma-7B-it	Word Definitions	Medical Indications	0.53	0.56	0.60	0.67
Gemma-7B-it	Word Definitions	Word Definitions	0.78	0.80	0.82	0.67
Gemma-2-9B-it	Word Definitions	City Locations	0.94	0.96	0.97	0.56
Gemma-2-9B-it	Word Definitions	Medical Indications	0.65	0.69	0.71	0.41
Gemma-2-9B-it	Word Definitions	Word Definitions	0.88	0.90	0.91	0.54
Llama-3.2-3B-Instruct	Word Definitions	City Locations	0.84	0.86	0.88	0.48
Llama-3.2-3B-Instruct	Word Definitions	Medical Indications	0.60	0.63	0.66	0.44
Llama-3.2-3B-Instruct	Word Definitions	Word Definitions	0.85	0.86	0.88	0.44
Llama3-Med42-8B	Word Definitions	City Locations	0.94	0.95	0.97	0.35
Llama3-Med42-8B	Word Definitions	Medical Indications	0.76	0.79	0.81	0.45
Llama3-Med42-8B	Word Definitions	Word Definitions	0.87	0.89	0.91	0.45
Llama-3.1-8B-Instruct	Word Definitions	City Locations	0.95	0.96	0.97	0.45
Llama-3.1-8B-Instruct	Word Definitions	Medical Indications	0.72	0.75	0.77	0.32
Llama-3.1-8B-Instruct	Word Definitions	Word Definitions	0.90	0.91	0.93	0.45
Bio-Medical-Llama-3-8B	Word Definitions	City Locations	0.78	0.81	0.83	0.35
Bio-Medical-Llama-3-8B	Word Definitions	Medical Indications	0.70	0.73	0.76	0.32
Bio-Medical-Llama-3-8B	Word Definitions	Word Definitions	0.87	0.89	0.90	0.39
Mistral-7B-Instruct-v0.3	Word Definitions	City Locations	0.93	0.94	0.96	0.48
Mistral-7B-Instruct-v0.3	Word Definitions	Medical Indications	0.74	0.77	0.80	0.52
Mistral-7B-Instruct-v0.3	Word Definitions	Word Definitions	0.87	0.88	0.90	0.35
Qwen-2.5-7B-Instruct	Word Definitions	City Locations	0.85	0.87	0.88	0.63
Qwen-2.5-7B-Instruct	Word Definitions	Medical Indications	0.74	0.76	0.79	0.67
Qwen-2.5-7B-Instruct	Word Definitions	Word Definitions	0.85	0.87	0.88	0.63
Qwen-2.5-14B-Instruct	Word Definitions	City Locations	0.95	0.96	0.97	0.66
Qwen-2.5-14B-Instruct	Word Definitions	Medical Indications	0.77	0.79	0.82	0.66
Qwen-2.5-14B-Instruct	Word Definitions	Word Definitions	0.88	0.90	0.92	0.51
Gemma-7B	Word Definitions	City Locations	0.88	0.90	0.91	0.56
Gemma-7B	Word Definitions	Medical Indications	0.59	0.63	0.66	0.48
Gemma-7B	Word Definitions	Word Definitions	0.81	0.83	0.85	0.56
Gemma-2-9B	Word Definitions	City Locations	0.92	0.93	0.94	0.56
Gemma-2-9B	Word Definitions	Medical Indications	0.70	0.73	0.75	0.46
Gemma-2-9B	Word Definitions	Word Definitions	0.86	0.88	0.90	0.41
Llama-3.2-3B	Word Definitions	City Locations	0.76	0.78	0.81	0.41
Llama-3.2-3B	Word Definitions	Medical Indications	0.70	0.73	0.76	0.41
Llama-3.2-3B	Word Definitions	Word Definitions	0.79	0.81	0.83	0.41
Llama-3-8B	Word Definitions	City Locations	0.84	0.86	0.88	0.39
Llama-3-8B	Word Definitions	Medical Indications	0.74	0.76	0.79	0.39
Llama-3-8B	Word Definitions	Word Definitions	0.85	0.87	0.89	0.39
Mistral-7B-v0.3	Word Definitions	City Locations	0.77	0.79	0.81	0.52
Mistral-7B-v0.3	Word Definitions	Medical Indications	0.73	0.76	0.78	0.39
Mistral-7B-v0.3	Word Definitions	Word Definitions	0.85	0.87	0.89	0.45
Qwen-2.5-7B	Word Definitions	City Locations	0.90	0.92	0.93	0.67
Qwen-2.5-7B	Word Definitions	Medical Indications	0.72	0.74	0.77	0.59
Qwen-2.5-7B	Word Definitions	Word Definitions	0.85	0.87	0.88	0.59
Qwen-2.5-14B	Word Definitions	City Locations	0.92	0.94	0.95	0.45
Qwen-2.5-14B	Word Definitions	Medical Indications	0.72	0.75	0.78	0.66
Qwen-2.5-14B	Word Definitions	Word Definitions	0.85	0.87	0.88	0.47

Table 24: **Row-wise confusion matrices for zero-shot prompting across all (model, dataset) pairs** (evaluated on the full bag). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by zero-shot prompting). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.89 – that means that 89% of all the true statements are classified as true.

Model ↓	Ground-truth label → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.89	0.11	0.00	0.00	0.06	0.94	0.00	0.00	0.18	0.82	0.00	0.00
	Medical Indications	0.79	0.21	0.00	0.00	0.32	0.68	0.00	0.00	0.27	0.73	0.00	0.00
	Word Definitions	0.61	0.39	0.00	0.00	0.37	0.63	0.00	0.00	0.05	0.95	0.00	0.00
Gemma-2-9B	City Locations	0.50	0.10	0.39	0.01	0.03	0.53	0.43	0.01	0.05	0.00	0.93	0.02
	Medical Indications	0.70	0.12	0.18	0.00	0.34	0.51	0.14	0.00	0.57	0.19	0.24	0.00
	Word Definitions	0.36	0.20	0.40	0.05	0.17	0.26	0.53	0.04	0.14	0.12	0.72	0.01
Gemma-2-9B-it	City Locations	0.98	0.02	0.00	0.00	0.03	0.97	0.00	0.00	0.06	0.45	0.49	0.00
	Medical Indications	0.87	0.12	0.01	0.00	0.25	0.75	0.00	0.00	0.23	0.59	0.18	0.00
	Word Definitions	0.76	0.16	0.09	0.00	0.25	0.70	0.06	0.00	0.10	0.65	0.25	0.00
Gemma-7B	City Locations	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00
	Medical Indications	0.16	0.00	0.00	0.84	0.09	0.00	0.00	0.91	0.03	0.00	0.00	0.96
	Word Definitions	0.03	0.00	0.00	0.97	0.03	0.01	0.00	0.96	0.00	0.00	0.00	1.00
Gemma-7B-it	City Locations	0.76	0.23	0.01	0.00	0.05	0.95	0.00	0.00	0.04	0.64	0.32	0.00
	Medical Indications	0.69	0.30	0.01	0.00	0.27	0.73	0.00	0.00	0.37	0.61	0.02	0.00
	Word Definitions	0.27	0.63	0.09	0.01	0.09	0.89	0.02	0.00	0.04	0.77	0.15	0.04
Llama-3-8B	City Locations	0.35	0.65	0.00	0.00	0.22	0.78	0.00	0.00	0.47	0.52	0.00	0.01
	Medical Indications	0.33	0.67	0.00	0.00	0.08	0.92	0.00	0.00	0.19	0.81	0.00	0.00
	Word Definitions	0.45	0.55	0.00	0.00	0.33	0.67	0.00	0.00	0.37	0.63	0.00	0.00
Llama-3.1-8B-Instruct	City Locations	0.95	0.05	0.00	0.00	0.03	0.97	0.00	0.00	0.08	0.65	0.27	0.00
	Medical Indications	0.54	0.46	0.00	0.00	0.07	0.93	0.00	0.00	0.13	0.86	0.01	0.00
	Word Definitions	0.54	0.46	0.00	0.00	0.21	0.79	0.00	0.00	0.06	0.94	0.00	0.00
Llama-3.2-3B	City Locations	0.29	0.71	0.00	0.00	0.11	0.89	0.00	0.00	0.43	0.57	0.00	0.00
	Medical Indications	0.46	0.54	0.00	0.00	0.34	0.66	0.00	0.00	0.50	0.50	0.00	0.00
	Word Definitions	0.48	0.52	0.00	0.00	0.44	0.56	0.00	0.00	0.46	0.54	0.00	0.00
Llama-3.2-3B-Instruct	City Locations	0.93	0.07	0.00	0.00	0.15	0.85	0.00	0.00	0.04	0.84	0.13	0.00
	Medical Indications	0.35	0.65	0.00	0.00	0.07	0.93	0.00	0.00	0.15	0.85	0.00	0.00
	Word Definitions	0.60	0.40	0.00	0.00	0.46	0.54	0.00	0.00	0.06	0.93	0.01	0.00
Llama3-Med42-8B	City Locations	0.96	0.04	0.00	0.00	0.04	0.96	0.00	0.00	0.22	0.64	0.14	0.00
	Medical Indications	0.69	0.31	0.00	0.00	0.13	0.87	0.00	0.00	0.17	0.82	0.00	0.00
	Word Definitions	0.47	0.52	0.01	0.00	0.18	0.81	0.01	0.00	0.08	0.89	0.03	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.88	0.00	0.09	0.03	0.04	0.09	0.50	0.37	0.08	0.00	0.90	0.02
	Medical Indications	0.47	0.04	0.49	0.00	0.07	0.09	0.85	0.00	0.03	0.00	0.97	0.00
	Word Definitions	0.63	0.02	0.33	0.02	0.40	0.01	0.55	0.04	0.25	0.00	0.73	0.02
Mistral-7B-v0.3	City Locations	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00
	Medical Indications	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00
	Word Definitions	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00
Qwen-2.5-14B	City Locations	0.95	0.05	0.00	0.00	0.02	0.98	0.00	0.00	0.04	0.47	0.49	0.00
	Medical Indications	0.53	0.46	0.01	0.00	0.04	0.95	0.00	0.00	0.00	0.77	0.23	0.00
	Word Definitions	0.50	0.41	0.09	0.00	0.05	0.90	0.05	0.00	0.00	0.53	0.47	0.00
Qwen-2.5-14B-Instruct	City Locations	0.93	0.07	0.00	0.00	0.02	0.98	0.00	0.00	0.00	0.28	0.71	0.00
	Medical Indications	0.63	0.23	0.14	0.00	0.04	0.87	0.09	0.00	0.00	0.46	0.54	0.00
	Word Definitions	0.55	0.32	0.13	0.00	0.06	0.89	0.06	0.00	0.01	0.37	0.62	0.00
Qwen-2.5-7B	City Locations	0.92	0.07	0.01	0.00	0.03	0.96	0.01	0.00	0.09	0.68	0.23	0.00
	Medical Indications	0.56	0.44	0.00	0.00	0.11	0.89	0.00	0.00	0.25	0.75	0.00	0.00
	Word Definitions	0.60	0.39	0.01	0.00	0.31	0.67	0.02	0.00	0.35	0.64	0.01	0.00
Qwen-2.5-7B-Instruct	City Locations	0.89	0.11	0.00	0.00	0.02	0.98	0.00	0.00	0.02	0.57	0.37	0.03
	Medical Indications	0.41	0.56	0.02	0.00	0.04	0.94	0.02	0.00	0.01	0.92	0.06	0.01
	Word Definitions	0.56	0.39	0.04	0.01	0.27	0.70	0.03	0.01	0.25	0.61	0.12	0.03

Table 25: **Row-wise confusion matrices for mean-difference probe with conformal prediction intervals (MD+CP) across all (model-dataset pairs)** (evaluated on the *last token representation*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True, False, Neither*) with groups of subcolumns that specify the distribution of predictions (*true, false, neither, abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.89 – that means that 89% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.89	0.01	0.10	0.00	0.01	0.92	0.07	0.00	0.34	0.23	0.42	0.00
	Medical Indications	0.68	0.13	0.19	0.00	0.09	0.75	0.16	0.00	0.21	0.58	0.20	0.00
	Word Definitions	0.73	0.09	0.18	0.00	0.09	0.72	0.18	0.00	0.42	0.36	0.22	0.00
Gemma-2-9B	City Locations	0.90	0.02	0.09	0.00	0.02	0.91	0.07	0.00	0.07	0.66	0.27	0.00
	Medical Indications	0.78	0.12	0.10	0.00	0.13	0.77	0.10	0.00	0.15	0.77	0.09	0.00
	Word Definitions	0.80	0.11	0.10	0.00	0.10	0.81	0.09	0.00	0.47	0.39	0.14	0.00
Gemma-2-9B-it	City Locations	0.89	0.09	0.02	0.00	0.05	0.93	0.02	0.00	0.79	0.19	0.02	0.00
	Medical Indications	0.70	0.11	0.19	0.00	0.11	0.73	0.17	0.00	0.42	0.38	0.20	0.00
	Word Definitions	0.75	0.11	0.14	0.00	0.11	0.75	0.14	0.00	0.36	0.42	0.23	0.00
Gemma-7B	City Locations	0.90	0.01	0.08	0.00	0.02	0.92	0.06	0.00	0.04	0.79	0.16	0.00
	Medical Indications	0.71	0.10	0.18	0.00	0.09	0.72	0.19	0.00	0.19	0.53	0.28	0.00
	Word Definitions	0.70	0.10	0.20	0.00	0.08	0.75	0.18	0.00	0.37	0.36	0.28	0.00
Gemma-7B-it	City Locations	0.94	0.03	0.03	0.00	0.06	0.90	0.04	0.00	0.32	0.54	0.14	0.00
	Medical Indications	0.53	0.11	0.36	0.00	0.12	0.61	0.27	0.00	0.19	0.47	0.34	0.00
	Word Definitions	0.75	0.11	0.13	0.00	0.11	0.78	0.11	0.00	0.31	0.48	0.22	0.00
Llama-3-8B	City Locations	0.90	0.01	0.08	0.00	0.00	0.92	0.08	0.00	0.66	0.17	0.17	0.00
	Medical Indications	0.71	0.12	0.17	0.00	0.11	0.75	0.14	0.00	0.40	0.39	0.21	0.00
	Word Definitions	0.73	0.09	0.17	0.00	0.07	0.74	0.19	0.00	0.48	0.31	0.22	0.00
Llama-3.1-8B-Instruct	City Locations	0.89	0.03	0.08	0.00	0.04	0.89	0.06	0.00	0.68	0.27	0.05	0.00
	Medical Indications	0.74	0.14	0.13	0.00	0.10	0.79	0.12	0.00	0.22	0.67	0.11	0.00
	Word Definitions	0.80	0.10	0.10	0.00	0.07	0.83	0.09	0.00	0.45	0.46	0.09	0.00
Llama-3.2-3B	City Locations	0.89	0.03	0.08	0.00	0.03	0.90	0.07	0.00	0.25	0.39	0.35	0.00
	Medical Indications	0.72	0.12	0.16	0.00	0.12	0.71	0.18	0.00	0.19	0.64	0.17	0.00
	Word Definitions	0.72	0.09	0.18	0.00	0.11	0.71	0.18	0.00	0.31	0.44	0.26	0.00
Llama-3.2-3B-Instruct	City Locations	0.87	0.05	0.08	0.00	0.04	0.90	0.06	0.00	0.85	0.06	0.09	0.00
	Medical Indications	0.59	0.11	0.30	0.00	0.11	0.62	0.27	0.00	0.33	0.38	0.30	0.00
	Word Definitions	0.59	0.10	0.31	0.00	0.09	0.62	0.29	0.00	0.37	0.40	0.24	0.00
Llama3-Med42-8B	City Locations	0.91	0.01	0.08	0.00	0.01	0.91	0.09	0.00	0.04	0.73	0.23	0.00
	Medical Indications	0.77	0.14	0.09	0.00	0.10	0.79	0.11	0.00	0.21	0.70	0.09	0.00
	Word Definitions	0.83	0.09	0.08	0.00	0.08	0.88	0.05	0.00	0.59	0.32	0.09	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.88	0.01	0.11	0.00	0.01	0.90	0.09	0.00	0.45	0.40	0.16	0.00
	Medical Indications	0.68	0.11	0.21	0.00	0.13	0.71	0.16	0.00	0.29	0.50	0.20	0.00
	Word Definitions	0.77	0.13	0.10	0.00	0.11	0.81	0.08	0.00	0.24	0.60	0.15	0.00
Mistral-7B-v0.3	City Locations	0.88	0.02	0.10	0.00	0.02	0.91	0.07	0.00	0.31	0.45	0.24	0.00
	Medical Indications	0.69	0.12	0.19	0.00	0.10	0.73	0.17	0.00	0.23	0.54	0.22	0.00
	Word Definitions	0.72	0.08	0.20	0.00	0.08	0.77	0.15	0.00	0.42	0.35	0.24	0.00
Qwen-2.5-14B	City Locations	0.89	0.04	0.07	0.00	0.02	0.93	0.06	0.00	0.21	0.48	0.30	0.00
	Medical Indications	0.72	0.15	0.13	0.00	0.11	0.76	0.14	0.00	0.54	0.25	0.21	0.00
	Word Definitions	0.80	0.09	0.11	0.00	0.11	0.79	0.10	0.00	0.39	0.46	0.15	0.00
Qwen-2.5-14B-Instruct	City Locations	0.88	0.09	0.03	0.00	0.06	0.92	0.02	0.00	0.38	0.59	0.03	0.00
	Medical Indications	0.77	0.12	0.11	0.00	0.11	0.79	0.10	0.00	0.46	0.44	0.10	0.00
	Word Definitions	0.82	0.08	0.10	0.00	0.08	0.85	0.07	0.00	0.30	0.62	0.08	0.00
Qwen-2.5-7B	City Locations	0.91	0.03	0.06	0.00	0.03	0.91	0.06	0.00	0.34	0.55	0.12	0.00
	Medical Indications	0.71	0.11	0.17	0.00	0.10	0.77	0.14	0.00	0.34	0.43	0.23	0.00
	Word Definitions	0.74	0.09	0.16	0.00	0.10	0.74	0.16	0.00	0.32	0.43	0.25	0.00
Qwen-2.5-7B-Instruct	City Locations	0.90	0.03	0.07	0.00	0.02	0.92	0.06	0.00	0.62	0.24	0.14	0.00
	Medical Indications	0.75	0.09	0.15	0.00	0.10	0.78	0.12	0.00	0.67	0.20	0.13	0.00
	Word Definitions	0.78	0.10	0.12	0.00	0.09	0.83	0.08	0.00	0.43	0.44	0.13	0.00

Table 26: **Row-wise confusion matrices for mean-difference probe with conformal prediction intervals (MD+CP) across all (model-dataset pairs)** (evaluated on the *bag*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 1.00 – that means that 100% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.99	0.00	0.01	0.00	0.95	0.03	0.02	0.00	0.98	0.01	0.01	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Gemma-2-9B	City Locations	1.00	0.00	0.00	0.00	0.97	0.01	0.02	0.00	0.99	0.00	0.01	0.00
	Medical Indications	0.95	0.01	0.04	0.00	0.88	0.04	0.09	0.00	0.98	0.01	0.01	0.00
	Word Definitions	0.80	0.02	0.18	0.00	0.65	0.05	0.29	0.00	0.63	0.06	0.31	0.00
Gemma-2-9B-it	City Locations	0.87	0.01	0.12	0.00	0.30	0.15	0.55	0.00	0.42	0.07	0.51	0.00
	Medical Indications	1.00	0.00	0.00	0.00	0.98	0.00	0.01	0.00	0.99	0.00	0.01	0.00
	Word Definitions	0.57	0.02	0.41	0.00	0.32	0.13	0.55	0.00	0.20	0.12	0.68	0.00
Gemma-7B	City Locations	0.95	0.03	0.02	0.00	0.48	0.46	0.06	0.00	0.76	0.21	0.02	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.93	0.00	0.06	0.00	0.90	0.01	0.09	0.00	0.94	0.00	0.05	0.00
Gemma-7B-it	City Locations	0.83	0.06	0.11	0.00	0.59	0.31	0.11	0.00	0.81	0.12	0.07	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.93	0.00	0.07	0.00	0.88	0.01	0.11	0.00	0.96	0.00	0.04	0.00
Llama-3-8B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.68	0.02	0.30	0.00	0.59	0.05	0.36	0.00	0.69	0.04	0.27	0.00
	Word Definitions	0.96	0.00	0.03	0.00	0.94	0.01	0.05	0.00	0.91	0.01	0.08	0.00
Llama-3.1-8B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama-3.2-3B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.83	0.03	0.14	0.00	0.56	0.22	0.23	0.00	0.95	0.02	0.03	0.00
	Word Definitions	0.92	0.01	0.07	0.00	0.75	0.09	0.16	0.00	0.92	0.02	0.07	0.00
Llama-3.2-3B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama3-Med42-8B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.93	0.00	0.07	0.00	0.93	0.00	0.06	0.00	0.94	0.00	0.06	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.97	0.00	0.03	0.00	0.88	0.02	0.10	0.00	0.97	0.00	0.03	0.00
	Medical Indications	0.96	0.01	0.03	0.00	0.91	0.04	0.05	0.00	0.90	0.04	0.06	0.00
	Word Definitions	1.00	0.00	0.00	0.00	0.97	0.01	0.02	0.00	1.00	0.00	0.00	0.00
Mistral-7B-v0.3	City Locations	0.97	0.01	0.02	0.00	0.86	0.08	0.06	0.00	0.90	0.05	0.06	0.00
	Medical Indications	0.97	0.00	0.03	0.00	0.96	0.00	0.04	0.00	0.97	0.00	0.03	0.00
	Word Definitions	0.94	0.00	0.06	0.00	0.86	0.02	0.12	0.00	0.92	0.02	0.06	0.00
Qwen-2.5-14B	City Locations	1.00	0.00	0.00	0.00	0.89	0.03	0.07	0.00	0.99	0.00	0.01	0.00
	Medical Indications	0.91	0.01	0.07	0.00	0.88	0.04	0.08	0.00	0.96	0.01	0.03	0.00
	Word Definitions	0.98	0.00	0.01	0.00	0.93	0.03	0.04	0.00	0.97	0.01	0.02	0.00
Qwen-2.5-14B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.90	0.01	0.09	0.00	0.66	0.10	0.24	0.00	0.60	0.04	0.36	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Qwen-2.5-7B	City Locations	0.99	0.00	0.01	0.00	0.95	0.03	0.02	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.93	0.01	0.06	0.00	0.76	0.08	0.16	0.00	0.87	0.05	0.09	0.00
	Word Definitions	0.92	0.01	0.07	0.00	0.90	0.01	0.09	0.00	0.90	0.01	0.08	0.00
Qwen-2.5-7B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00

Table 27: **Row-wise confusion matrices for TTPD probe with conformal prediction intervals (TTPD+CP) across all (model-dataset pairs)** (evaluated on the *last token representation*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.87 – that means that 87% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.87	0.08	0.05	0.00	0.04	0.90	0.05	0.00	0.09	0.78	0.14	0.00
	Medical Indications	0.59	0.08	0.33	0.00	0.09	0.64	0.27	0.00	0.03	0.61	0.36	0.00
	Word Definitions	0.44	0.11	0.46	0.00	0.08	0.50	0.41	0.00	0.00	0.69	0.31	0.00
Gemma-2-9B	City Locations	0.91	0.09	0.00	0.00	0.14	0.85	0.00	0.00	0.00	1.00	0.00	0.00
	Medical Indications	0.71	0.09	0.20	0.00	0.11	0.70	0.19	0.00	0.11	0.73	0.16	0.00
	Word Definitions	0.68	0.08	0.23	0.00	0.08	0.74	0.18	0.00	0.08	0.55	0.36	0.00
Gemma-2-9B-it	City Locations	0.86	0.12	0.02	0.00	0.04	0.94	0.01	0.00	0.98	0.02	0.01	0.00
	Medical Indications	0.70	0.07	0.23	0.00	0.11	0.70	0.19	0.00	0.14	0.34	0.51	0.00
	Word Definitions	0.53	0.11	0.36	0.00	0.07	0.59	0.34	0.00	0.03	0.58	0.39	0.00
Gemma-7B	City Locations	0.90	0.05	0.06	0.00	0.08	0.88	0.03	0.00	0.03	0.95	0.02	0.00
	Medical Indications	0.68	0.08	0.25	0.00	0.09	0.68	0.23	0.00	0.10	0.52	0.38	0.00
	Word Definitions	0.75	0.09	0.16	0.00	0.09	0.79	0.12	0.00	0.13	0.66	0.21	0.00
Gemma-7B-it	City Locations	0.88	0.04	0.08	0.00	0.05	0.91	0.03	0.00	0.90	0.01	0.09	0.00
	Medical Indications	0.54	0.07	0.39	0.00	0.13	0.52	0.35	0.00	0.45	0.13	0.41	0.00
	Word Definitions	0.60	0.12	0.28	0.00	0.09	0.64	0.26	0.00	0.17	0.37	0.45	0.00
Llama-3-8B	City Locations	0.88	0.10	0.02	0.00	0.07	0.89	0.04	0.00	0.02	0.98	0.01	0.00
	Medical Indications	0.65	0.10	0.25	0.00	0.11	0.70	0.19	0.00	0.02	0.62	0.36	0.00
	Word Definitions	0.72	0.11	0.17	0.00	0.09	0.79	0.13	0.00	0.41	0.42	0.17	0.00
Llama-3.1-8B-Instruct	City Locations	0.84	0.03	0.13	0.00	0.03	0.92	0.05	0.00	0.00	0.84	0.16	0.00
	Medical Indications	0.68	0.06	0.26	0.00	0.10	0.65	0.25	0.00	0.00	0.74	0.26	0.00
	Word Definitions	0.80	0.08	0.11	0.00	0.09	0.81	0.11	0.00	0.15	0.55	0.30	0.00
Llama-3.2-3B	City Locations	0.88	0.10	0.02	0.00	0.07	0.90	0.02	0.00	0.02	0.98	0.01	0.00
	Medical Indications	0.71	0.08	0.21	0.00	0.10	0.69	0.21	0.00	0.01	0.86	0.12	0.00
	Word Definitions	0.64	0.09	0.27	0.00	0.08	0.66	0.25	0.00	0.10	0.55	0.34	0.00
Llama-3.2-3B-Instruct	City Locations	0.85	0.11	0.04	0.00	0.10	0.87	0.02	0.00	0.73	0.23	0.04	0.00
	Medical Indications	0.51	0.13	0.36	0.00	0.08	0.60	0.32	0.00	0.06	0.41	0.53	0.00
	Word Definitions	0.44	0.12	0.43	0.00	0.07	0.49	0.44	0.00	0.13	0.48	0.40	0.00
Llama3-Med42-8B	City Locations	0.89	0.04	0.07	0.00	0.05	0.92	0.02	0.00	0.22	0.66	0.12	0.00
	Medical Indications	0.69	0.09	0.22	0.00	0.10	0.74	0.16	0.00	0.06	0.07	0.87	0.00
	Word Definitions	0.71	0.10	0.19	0.00	0.07	0.78	0.15	0.00	0.02	0.82	0.17	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.88	0.11	0.00	0.00	0.07	0.93	0.00	0.00	0.30	0.70	0.00	0.00
	Medical Indications	0.73	0.08	0.19	0.00	0.10	0.72	0.19	0.00	0.06	0.54	0.40	0.00
	Word Definitions	0.71	0.10	0.19	0.00	0.09	0.72	0.19	0.00	0.16	0.57	0.27	0.00
Mistral-7B-v0.3	City Locations	0.88	0.08	0.04	0.00	0.06	0.92	0.03	0.00	0.31	0.59	0.10	0.00
	Medical Indications	0.72	0.07	0.21	0.00	0.10	0.73	0.17	0.00	0.15	0.36	0.49	0.00
	Word Definitions	0.72	0.10	0.17	0.00	0.11	0.75	0.15	0.00	0.20	0.50	0.31	0.00
Qwen-2.5-14B	City Locations	0.91	0.07	0.01	0.00	0.08	0.91	0.01	0.00	0.10	0.88	0.02	0.00
	Medical Indications	0.69	0.06	0.25	0.00	0.15	0.65	0.20	0.00	0.36	0.14	0.50	0.00
	Word Definitions	0.60	0.12	0.28	0.00	0.09	0.61	0.30	0.00	0.13	0.49	0.37	0.00
Qwen-2.5-14B-Instruct	City Locations	0.88	0.10	0.03	0.00	0.07	0.91	0.02	0.00	0.39	0.54	0.07	0.00
	Medical Indications	0.80	0.08	0.13	0.00	0.13	0.75	0.11	0.00	0.03	0.55	0.42	0.00
	Word Definitions	0.72	0.10	0.19	0.00	0.11	0.75	0.15	0.00	0.04	0.71	0.25	0.00
Qwen-2.5-7B	City Locations	0.92	0.03	0.05	0.00	0.06	0.89	0.05	0.00	0.01	0.90	0.09	0.00
	Medical Indications	0.67	0.07	0.25	0.00	0.10	0.70	0.20	0.00	0.33	0.19	0.48	0.00
	Word Definitions	0.52	0.11	0.37	0.00	0.07	0.56	0.37	0.00	0.11	0.52	0.37	0.00
Qwen-2.5-7B-Instruct	City Locations	0.90	0.03	0.07	0.00	0.07	0.87	0.06	0.00	0.33	0.10	0.57	0.00
	Medical Indications	0.73	0.09	0.17	0.00	0.15	0.72	0.13	0.00	0.02	0.88	0.10	0.00
	Word Definitions	0.71	0.11	0.18	0.00	0.10	0.71	0.19	0.00	0.11	0.60	0.30	0.00

Table 28: **Row-wise confusion matrices for TTPD probe with conformal prediction intervals (TTPD+CP) across all (model-dataset pairs)** (evaluated on the *bag*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.19 – that means that 19% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.19	0.00	0.81	0.00	0.16	0.01	0.83	0.00	0.03	0.24	0.72	0.00
	Medical Indications	0.08	0.08	0.84	0.00	0.09	0.05	0.86	0.00	0.08	0.02	0.91	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Gemma-2-9B	City Locations	0.92	0.08	0.00	0.00	0.14	0.85	0.00	0.00	0.00	1.00	0.00	0.00
	Medical Indications	0.12	0.03	0.85	0.00	0.08	0.03	0.89	0.00	0.00	0.02	0.98	0.00
	Word Definitions	0.30	0.08	0.62	0.00	0.11	0.16	0.72	0.00	0.00	0.79	0.21	0.00
Gemma-2-9B-it	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.48	0.01	0.51	0.00	0.11	0.20	0.69	0.00	0.01	0.40	0.59	0.00
Gemma-7B	City Locations	0.93	0.04	0.04	0.00	0.05	0.88	0.07	0.00	0.01	0.98	0.01	0.00
	Medical Indications	0.10	0.08	0.82	0.00	0.07	0.13	0.80	0.00	0.00	0.26	0.74	0.00
	Word Definitions	0.38	0.10	0.53	0.00	0.09	0.27	0.64	0.00	0.01	0.59	0.40	0.00
Gemma-7B-it	City Locations	0.23	0.05	0.72	0.00	0.19	0.13	0.67	0.00	0.01	0.51	0.48	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama-3-8B	City Locations	0.88	0.11	0.00	0.00	0.11	0.88	0.01	0.00	0.02	0.98	0.00	0.00
	Medical Indications	0.54	0.08	0.37	0.00	0.10	0.54	0.36	0.00	0.01	0.66	0.32	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama-3.1-8B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.18	0.10	0.73	0.00	0.08	0.14	0.78	0.00	0.20	0.32	0.48	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama-3.2-3B	City Locations	0.91	0.07	0.02	0.00	0.07	0.90	0.03	0.00	0.02	0.98	0.01	0.00
	Medical Indications	0.64	0.05	0.31	0.00	0.10	0.18	0.72	0.00	0.01	0.56	0.43	0.00
	Word Definitions	0.38	0.03	0.59	0.00	0.09	0.08	0.82	0.00	0.01	0.16	0.83	0.00
Llama-3.2-3B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama3-Med42-8B	City Locations	0.82	0.01	0.17	0.00	0.41	0.07	0.52	0.00	0.51	0.47	0.02	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Mistral-7B-Instruct-v0.3	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.35	0.05	0.60	0.00	0.24	0.11	0.65	0.00	0.23	0.32	0.45	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Mistral-7B-v0.3	City Locations	0.88	0.06	0.06	0.00	0.24	0.56	0.20	0.00	0.34	0.50	0.16	0.00
	Medical Indications	0.11	0.14	0.75	0.00	0.06	0.13	0.81	0.00	0.00	0.39	0.61	0.00
	Word Definitions	0.76	0.00	0.24	0.00	0.31	0.01	0.68	0.00	0.34	0.08	0.58	0.00
Qwen-2.5-14B	City Locations	0.90	0.04	0.06	0.00	0.05	0.87	0.08	0.00	0.05	0.88	0.08	0.00
	Medical Indications	0.11	0.02	0.87	0.00	0.09	0.02	0.89	0.00	0.01	0.07	0.93	0.00
	Word Definitions	0.52	0.00	0.47	0.00	0.29	0.01	0.70	0.00	0.29	0.12	0.59	0.00
Qwen-2.5-14B-Instruct	City Locations	0.87	0.00	0.13	0.00	0.02	0.17	0.81	0.00	0.26	0.01	0.73	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Qwen-2.5-7B	City Locations	0.91	0.06	0.03	0.00	0.11	0.85	0.05	0.00	0.28	0.66	0.06	0.00
	Medical Indications	0.10	0.09	0.81	0.00	0.07	0.08	0.85	0.00	0.00	0.13	0.87	0.00
	Word Definitions	0.39	0.01	0.60	0.00	0.11	0.03	0.86	0.00	0.09	0.11	0.81	0.00
Qwen-2.5-7B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00

Table 29: **Row-wise confusion matrices for supervised PCA probe with conformal prediction intervals (sPCA+CP) across all \langle model-dataset pairs \rangle** (evaluated on the *last token representation*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True, False, Neither*) with groups of subcolumns that specify the distribution of predictions (*true, false, neither, abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.90 – that means that 90% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.90	0.02	0.09	0.00	0.01	0.92	0.07	0.00	0.04	0.75	0.21	0.00
	Medical Indications	0.74	0.10	0.16	0.00	0.09	0.79	0.13	0.00	0.06	0.79	0.15	0.00
	Word Definitions	0.80	0.11	0.09	0.00	0.10	0.83	0.07	0.00	0.39	0.54	0.07	0.00
Gemma-2-9B	City Locations	0.89	0.03	0.08	0.00	0.01	0.93	0.06	0.00	0.00	0.98	0.02	0.00
	Medical Indications	0.83	0.11	0.06	0.00	0.09	0.84	0.07	0.00	0.03	0.86	0.11	0.00
	Word Definitions	0.87	0.10	0.03	0.00	0.08	0.89	0.04	0.00	0.42	0.50	0.08	0.00
Gemma-2-9B-it	City Locations	0.86	0.10	0.03	0.00	0.05	0.93	0.02	0.00	0.47	0.49	0.04	0.00
	Medical Indications	0.79	0.10	0.10	0.00	0.07	0.81	0.12	0.00	0.41	0.28	0.31	0.00
	Word Definitions	0.79	0.13	0.08	0.00	0.09	0.86	0.05	0.00	0.38	0.52	0.10	0.00
Gemma-7B	City Locations	0.89	0.00	0.11	0.00	0.00	0.91	0.09	0.00	0.14	0.29	0.57	0.00
	Medical Indications	0.78	0.09	0.13	0.00	0.08	0.81	0.11	0.00	0.23	0.49	0.27	0.00
	Word Definitions	0.84	0.08	0.08	0.00	0.09	0.83	0.07	0.00	0.42	0.42	0.16	0.00
Gemma-7B-it	City Locations	0.86	0.04	0.10	0.00	0.02	0.92	0.07	0.00	0.10	0.60	0.30	0.00
	Medical Indications	0.62	0.09	0.28	0.00	0.10	0.66	0.24	0.00	0.46	0.46	0.07	0.00
	Word Definitions	0.80	0.08	0.12	0.00	0.07	0.84	0.10	0.00	0.33	0.49	0.18	0.00
Llama-3-8B	City Locations	0.88	0.02	0.10	0.00	0.01	0.91	0.08	0.00	0.70	0.12	0.19	0.00
	Medical Indications	0.84	0.12	0.05	0.00	0.11	0.85	0.04	0.00	0.35	0.56	0.09	0.00
	Word Definitions	0.85	0.09	0.07	0.00	0.08	0.86	0.07	0.00	0.49	0.43	0.08	0.00
Llama-3.1-8B-Instruct	City Locations	0.88	0.12	0.01	0.00	0.07	0.93	0.00	0.00	0.52	0.48	0.00	0.00
	Medical Indications	0.83	0.11	0.06	0.00	0.08	0.84	0.08	0.00	0.25	0.59	0.15	0.00
	Word Definitions	0.89	0.07	0.03	0.00	0.07	0.91	0.02	0.00	0.48	0.48	0.04	0.00
Llama-3.2-3B	City Locations	0.92	0.04	0.04	0.00	0.03	0.93	0.04	0.00	0.39	0.58	0.03	0.00
	Medical Indications	0.73	0.11	0.16	0.00	0.09	0.79	0.13	0.00	0.07	0.72	0.21	0.00
	Word Definitions	0.77	0.07	0.16	0.00	0.09	0.81	0.10	0.00	0.28	0.49	0.22	0.00
Llama-3.2-3B-Instruct	City Locations	0.88	0.12	0.00	0.00	0.09	0.91	0.00	0.00	0.89	0.11	0.00	0.00
	Medical Indications	0.67	0.10	0.23	0.00	0.09	0.70	0.20	0.00	0.26	0.49	0.25	0.00
	Word Definitions	0.69	0.11	0.20	0.00	0.07	0.76	0.17	0.00	0.37	0.43	0.21	0.00
Llama3-Med42-8B	City Locations	0.91	0.01	0.08	0.00	0.00	0.92	0.08	0.00	0.26	0.16	0.57	0.00
	Medical Indications	0.80	0.12	0.07	0.00	0.08	0.87	0.05	0.00	0.24	0.65	0.11	0.00
	Word Definitions	0.91	0.09	0.00	0.00	0.07	0.92	0.00	0.00	0.51	0.49	0.00	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.90	0.04	0.06	0.00	0.04	0.89	0.06	0.00	0.46	0.50	0.03	0.00
	Medical Indications	0.79	0.11	0.10	0.00	0.09	0.81	0.10	0.00	0.48	0.46	0.07	0.00
	Word Definitions	0.86	0.13	0.01	0.00	0.11	0.88	0.01	0.00	0.43	0.56	0.02	0.00
Mistral-7B-v0.3	City Locations	0.86	0.00	0.14	0.00	0.00	0.91	0.08	0.00	0.17	0.45	0.38	0.00
	Medical Indications	0.80	0.10	0.10	0.00	0.10	0.84	0.07	0.00	0.42	0.46	0.12	0.00
	Word Definitions	0.86	0.09	0.05	0.00	0.09	0.87	0.03	0.00	0.43	0.48	0.09	0.00
Qwen-2.5-14B	City Locations	0.92	0.01	0.07	0.00	0.01	0.93	0.06	0.00	0.13	0.47	0.40	0.00
	Medical Indications	0.85	0.11	0.04	0.00	0.08	0.88	0.05	0.00	0.37	0.56	0.07	0.00
	Word Definitions	0.89	0.10	0.01	0.00	0.10	0.89	0.01	0.00	0.48	0.49	0.03	0.00
Qwen-2.5-14B-Instruct	City Locations	0.93	0.04	0.03	0.00	0.05	0.92	0.03	0.00	0.51	0.49	0.01	0.00
	Medical Indications	0.87	0.11	0.02	0.00	0.09	0.89	0.02	0.00	0.23	0.71	0.06	0.00
	Word Definitions	0.90	0.10	0.00	0.00	0.09	0.91	0.00	0.00	0.35	0.65	0.00	0.00
Qwen-2.5-7B	City Locations	0.89	0.06	0.06	0.00	0.03	0.92	0.04	0.00	0.52	0.46	0.02	0.00
	Medical Indications	0.80	0.11	0.10	0.00	0.08	0.83	0.09	0.00	0.39	0.47	0.14	0.00
	Word Definitions	0.87	0.10	0.03	0.00	0.08	0.89	0.03	0.00	0.40	0.53	0.07	0.00
Qwen-2.5-7B-Instruct	City Locations	0.89	0.02	0.09	0.00	0.02	0.91	0.07	0.00	0.50	0.43	0.06	0.00
	Medical Indications	0.80	0.11	0.08	0.00	0.12	0.82	0.06	0.00	0.44	0.49	0.06	0.00
	Word Definitions	0.85	0.10	0.05	0.00	0.10	0.88	0.03	0.00	0.45	0.51	0.04	0.00

Table 30: **Row-wise confusion matrices for supervised PCA probe with conformal prediction intervals (sPCA+CP) across all (model-dataset pairs)** (evaluated on the *bag*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by MD+CP). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 1.00 – that means that 100% of all the true statements are classified as true.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.81	0.03	0.17	0.00	0.72	0.07	0.21	0.00	0.68	0.07	0.25	0.00
Gemma-2-9B	City Locations	0.98	0.00	0.02	0.00	0.54	0.04	0.41	0.00	0.06	0.24	0.70	0.00
	Medical Indications	0.86	0.03	0.11	0.00	0.44	0.33	0.23	0.00	0.00	0.84	0.15	0.00
	Word Definitions	0.96	0.01	0.03	0.00	0.88	0.06	0.07	0.00	0.71	0.16	0.13	0.00
Gemma-2-9B-it	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.95	0.00	0.05	0.00	0.61	0.09	0.31	0.00	0.55	0.15	0.30	0.00
Gemma-7B	City Locations	1.00	0.00	0.00	0.00	0.97	0.01	0.02	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.27	0.01	0.72	0.00	0.18	0.03	0.79	0.00	0.40	0.05	0.55	0.00
	Word Definitions	0.82	0.00	0.17	0.00	0.53	0.03	0.44	0.00	0.53	0.03	0.44	0.00
Gemma-7B-it	City Locations	0.90	0.01	0.09	0.00	0.02	0.82	0.16	0.00	0.11	0.51	0.37	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.88	0.01	0.11	0.00	0.37	0.15	0.48	0.00	0.52	0.14	0.35	0.00
Llama-3-8B	City Locations	0.96	0.03	0.02	0.00	0.42	0.49	0.09	0.00	0.54	0.40	0.06	0.00
	Medical Indications	0.82	0.10	0.08	0.00	0.12	0.78	0.10	0.00	0.44	0.45	0.11	0.00
	Word Definitions	0.95	0.01	0.04	0.00	0.86	0.02	0.12	0.00	0.78	0.06	0.16	0.00
Llama-3.1-8B-Instruct	City Locations	0.91	0.00	0.09	0.00	0.00	0.76	0.24	0.00	0.05	0.35	0.59	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama-3.2-3B	City Locations	0.91	0.04	0.05	0.00	0.05	0.87	0.08	0.00	0.24	0.53	0.23	0.00
	Medical Indications	0.75	0.05	0.20	0.00	0.12	0.41	0.47	0.00	0.18	0.48	0.34	0.00
	Word Definitions	0.75	0.00	0.24	0.00	0.43	0.02	0.55	0.00	0.41	0.03	0.56	0.00
Llama-3.2-3B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Llama3-Med42-8B	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.94	0.01	0.05	0.00	0.67	0.17	0.16	0.00	0.89	0.04	0.07	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Mistral-7B-Instruct-v0.3	City Locations	0.88	0.01	0.11	0.00	0.01	0.86	0.13	0.00	0.26	0.16	0.59	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.84	0.01	0.15	0.00	0.76	0.03	0.21	0.00	0.79	0.02	0.19	0.00
Mistral-7B-v0.3	City Locations	1.00	0.00	0.00	0.00	0.98	0.00	0.02	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.80	0.06	0.15	0.00	0.11	0.66	0.23	0.00	0.24	0.49	0.27	0.00
	Word Definitions	0.87	0.02	0.10	0.00	0.49	0.19	0.32	0.00	0.42	0.30	0.28	0.00
Qwen-2.5-14B	City Locations	0.98	0.00	0.02	0.00	0.45	0.32	0.24	0.00	0.53	0.37	0.10	0.00
	Medical Indications	0.83	0.06	0.11	0.00	0.26	0.37	0.37	0.00	0.17	0.49	0.34	0.00
	Word Definitions	0.96	0.00	0.03	0.00	0.87	0.03	0.10	0.00	0.89	0.02	0.09	0.00
Qwen-2.5-14B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	0.84	0.08	0.09	0.00	0.10	0.82	0.09	0.00	0.15	0.71	0.14	0.00
	Word Definitions	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
Qwen-2.5-7B	City Locations	0.93	0.01	0.06	0.00	0.19	0.54	0.28	0.00	0.22	0.51	0.27	0.00
	Medical Indications	0.36	0.05	0.59	0.00	0.32	0.08	0.59	0.00	0.23	0.17	0.60	0.00
	Word Definitions	0.97	0.00	0.03	0.00	0.92	0.01	0.06	0.00	0.97	0.01	0.03	0.00
Qwen-2.5-7B-Instruct	City Locations	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Medical Indications	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
	Word Definitions	0.85	0.09	0.06	0.00	0.16	0.77	0.07	0.00	0.36	0.53	0.11	0.00

Table 31: **Row-wise confusion matrices for the multiclass sAwMIL across all (model-dataset pairs)** (evaluated on the *bag*). Each row corresponds to a specific model and a dataset. Columns are grouped by the ground-truth labels (*True*, *False*, *Neither*) with groups of subcolumns that specify the distribution of predictions (*true*, *false*, *neither*, *abstain*). For each statement in a dataset, the predicted class is the class with the highest probability (as estimated by multiclass sAwMIL). The values in each group of four subcolumns sum to 1 because they are normalized counts. For example, in the first row under the *True* ground-truth label, we see that *true* predictions have the value of 0.80 – that means that 80% of all the true statements are classified as true. In other words, each row is a flattened (and normalized) confusion matrix.

Model ↓	True Labels → Predicted → Data Set ↓	True				False				Neither			
		True	False	Neither	Abstain	True	False	Neither	Abstain	True	False	Neither	Abstain
Bio-Medical-Llama-3-8B	City Locations	0.80	0.01	0.00	0.18	0.00	0.88	0.00	0.11	0.00	0.00	1.00	0.00
	Medical Indications	0.77	0.10	0.00	0.12	0.10	0.81	0.01	0.08	0.00	0.00	1.00	0.00
	Word Definitions	0.86	0.09	0.01	0.03	0.11	0.87	0.01	0.02	0.00	0.01	0.98	0.01
Gemma-2-9B	City Locations	0.87	0.00	0.00	0.13	0.00	0.85	0.00	0.15	0.00	0.00	1.00	0.00
	Medical Indications	0.81	0.08	0.01	0.11	0.10	0.73	0.01	0.16	0.01	0.00	0.97	0.02
	Word Definitions	0.86	0.10	0.00	0.04	0.11	0.85	0.01	0.03	0.01	0.01	0.98	0.00
Gemma-2-9B-it	City Locations	0.85	0.02	0.02	0.12	0.02	0.87	0.00	0.11	0.00	0.00	0.98	0.02
	Medical Indications	0.76	0.10	0.01	0.13	0.07	0.84	0.01	0.09	0.01	0.00	0.99	0.01
	Word Definitions	0.83	0.08	0.00	0.08	0.07	0.88	0.01	0.05	0.00	0.01	0.98	0.01
Gemma-7B	City Locations	0.86	0.01	0.00	0.13	0.01	0.89	0.00	0.11	0.00	0.00	1.00	0.00
	Medical Indications	0.75	0.09	0.00	0.17	0.08	0.73	0.01	0.18	0.01	0.01	0.95	0.04
	Word Definitions	0.82	0.11	0.04	0.02	0.15	0.81	0.01	0.03	0.02	0.01	0.95	0.01
Gemma-7B-it	City Locations	0.86	0.02	0.01	0.11	0.02	0.87	0.00	0.11	0.00	0.00	0.99	0.01
	Medical Indications	0.58	0.05	0.01	0.35	0.15	0.50	0.00	0.35	0.00	0.01	0.94	0.05
	Word Definitions	0.69	0.16	0.04	0.11	0.12	0.77	0.03	0.07	0.01	0.01	0.96	0.01
Llama-3-8B	City Locations	0.87	0.01	0.00	0.12	0.01	0.85	0.00	0.14	0.00	0.00	1.00	0.00
	Medical Indications	0.78	0.09	0.01	0.12	0.11	0.77	0.01	0.11	0.00	0.01	0.98	0.01
	Word Definitions	0.87	0.12	0.00	0.00	0.13	0.85	0.01	0.00	0.01	0.01	0.97	0.00
Llama-3.1-8B-Instruct	City Locations	0.87	0.02	0.00	0.11	0.02	0.86	0.00	0.13	0.00	0.00	1.00	0.00
	Medical Indications	0.80	0.12	0.00	0.09	0.09	0.85	0.00	0.06	0.00	0.00	1.00	0.00
	Word Definitions	0.87	0.06	0.00	0.07	0.07	0.86	0.00	0.06	0.00	0.01	0.98	0.01
Llama-3.2-3B	City Locations	0.84	0.03	0.00	0.12	0.04	0.86	0.01	0.09	0.00	0.01	0.99	0.01
	Medical Indications	0.71	0.10	0.00	0.19	0.09	0.70	0.01	0.20	0.00	0.00	0.98	0.02
	Word Definitions	0.74	0.14	0.02	0.10	0.11	0.77	0.02	0.10	0.01	0.01	0.96	0.02
Llama-3.2-3B-Instruct	City Locations	0.86	0.02	0.00	0.11	0.02	0.82	0.00	0.16	0.00	0.00	1.00	0.00
	Medical Indications	0.67	0.12	0.00	0.21	0.11	0.72	0.00	0.18	0.00	0.00	0.99	0.01
	Word Definitions	0.86	0.11	0.03	0.00	0.14	0.85	0.01	0.00	0.01	0.01	0.98	0.00
Llama3-Med42-8B	City Locations	0.87	0.01	0.00	0.12	0.01	0.84	0.00	0.15	0.00	0.00	0.99	0.01
	Medical Indications	0.81	0.10	0.00	0.10	0.11	0.82	0.00	0.07	0.00	0.00	1.00	0.00
	Word Definitions	0.86	0.06	0.01	0.06	0.08	0.85	0.00	0.07	0.01	0.01	0.97	0.02
Mistral-7B-Instruct-v0.3	City Locations	0.86	0.01	0.00	0.13	0.01	0.79	0.00	0.19	0.00	0.00	1.00	0.00
	Medical Indications	0.77	0.09	0.00	0.14	0.10	0.76	0.00	0.14	0.00	0.01	0.99	0.01
	Word Definitions	0.84	0.08	0.03	0.05	0.09	0.87	0.01	0.03	0.01	0.01	0.98	0.01
Mistral-7B-v0.3	City Locations	0.86	0.01	0.00	0.13	0.01	0.86	0.00	0.13	0.00	0.00	1.00	0.00
	Medical Indications	0.75	0.08	0.00	0.17	0.09	0.76	0.00	0.15	0.00	0.00	0.99	0.01
	Word Definitions	0.87	0.11	0.02	0.00	0.14	0.85	0.01	0.00	0.01	0.01	0.98	0.00
Qwen-2.5-14B	City Locations	0.87	0.01	0.00	0.12	0.01	0.86	0.00	0.13	0.00	0.00	0.99	0.01
	Medical Indications	0.78	0.12	0.00	0.10	0.11	0.76	0.00	0.13	0.00	0.00	0.99	0.00
	Word Definitions	0.86	0.12	0.02	0.00	0.13	0.86	0.01	0.00	0.01	0.02	0.98	0.00
Qwen-2.5-14B-Instruct	City Locations	0.82	0.01	0.00	0.17	0.01	0.87	0.00	0.11	0.00	0.00	1.00	0.00
	Medical Indications	0.80	0.12	0.01	0.07	0.07	0.86	0.01	0.07	0.00	0.00	0.99	0.00
	Word Definitions	0.86	0.07	0.01	0.05	0.06	0.86	0.00	0.08	0.00	0.01	0.98	0.01
Qwen-2.5-7B	City Locations	0.84	0.01	0.00	0.16	0.01	0.84	0.00	0.14	0.00	0.00	1.00	0.00
	Medical Indications	0.74	0.10	0.01	0.15	0.10	0.77	0.01	0.13	0.00	0.00	0.98	0.01
	Word Definitions	0.83	0.12	0.02	0.03	0.10	0.87	0.01	0.03	0.01	0.01	0.97	0.01
Qwen-2.5-7B-Instruct	City Locations	0.82	0.02	0.01	0.15	0.02	0.85	0.00	0.13	0.00	0.00	0.99	0.01
	Medical Indications	0.78	0.08	0.01	0.13	0.14	0.72	0.02	0.12	0.01	0.01	0.98	0.01
	Word Definitions	0.83	0.13	0.02	0.02	0.09	0.87	0.02	0.02	0.01	0.01	0.97	0.01

L Related Resources

Code

The code associated with this manuscript is publicly available on GitHub at `carlomarxdk/trilemma-of-truth` (release version 0.7).

Data

The data used in this study are publicly available on Hugging Face at `carlomarxdk/trilemma-of-truth` (DOI: 10.57967/hf/5900).

Extended Manuscript

The manuscript with additional results is available on ArXiv: 10.48550/arXiv.2506.23921 (Version 2 from July 8, 2025).