

# A Head to Predict and a Head to Question: Pre-trained Uncertainty Quantification Heads for Hallucination Detection in LLM Outputs

Artem Shelmanov<sup>1</sup>  $\diamond$  Ekaterina Fadeeva<sup>2</sup>  $\diamond$  Akim Tsvigun<sup>5</sup> Ivan Tsvigun<sup>6</sup>,  
Zhuohan Xie<sup>1</sup> Igor Kiselev<sup>7</sup> Nico Daheim<sup>2</sup> Caiqi Zhang<sup>3</sup> Artem Vazhentsev<sup>8</sup>  
Mrinmaya Sachan<sup>2</sup> Preslav Nakov<sup>1</sup> Timothy Baldwin<sup>1,4</sup>  
<sup>1</sup>MBZUAI <sup>2</sup>ETH Zürich <sup>3</sup>University of Cambridge <sup>4</sup>The University of Melbourne  
<sup>5</sup>Nebius.AI <sup>6</sup>Behavox <sup>7</sup>Accenture <sup>8</sup>Computational Semantics Group  
{artem.shelmanov,preslav.nakov,timothy.baldwin}@mbzuai.ac.ac  
{efadeeva,msachan}@ethz.ch aktsvigun@nebius.com

## Abstract

Large Language Models (LLMs) have the tendency to hallucinate, i.e., to sporadically generate false or fabricated information. This presents a major challenge, as hallucinations often appear highly convincing and users generally lack the tools to detect them. Uncertainty quantification (UQ) provides a framework for assessing the reliability of model outputs, aiding in the identification of potential hallucinations. In this work, we introduce pre-trained UQ heads: supervised auxiliary modules for LLMs that substantially enhance their ability to capture uncertainty compared to unsupervised UQ methods. Their strong performance stems from the transformer architecture in their design, in the form of informative features derived from LLM attention maps and logits. Our experiments show that these heads are highly robust and achieve state-of-the-art performance in claim-level hallucination detection across both in-domain and out-of-domain prompts. Moreover, these modules demonstrate strong generalization to languages they were not explicitly trained on. We pre-train a collection of UQ heads for popular LLM series, including Mistral, Llama, and Gemma. We publicly release both the code and the pre-trained heads.<sup>1</sup>

## 1 Introduction

Uncertainty quantification (UQ) (Gal and Ghahramani, 2016; Baan et al., 2023; Geng et al., 2024; Zhang et al., 2024a) has become an increasingly important topic in natural language processing (NLP),

particularly for addressing challenges with hallucinations (Huang et al., 2025) and low-quality outputs of large language models (LLMs) (Malinin and Gales, 2021; Kuhn et al., 2023; Fadeeva et al., 2024). UQ offers the potential to improve the safety and reliability of LLM-based applications by flagging highly uncertain generations. Such generations could be discarded or marked as untrustworthy, thus reducing the risk of misleading information reaching users (Zhang et al., 2024a,b; Huang et al., 2024). Contrary to other methods for detecting hallucinations that rely on external knowledge bases or additional LLMs (Ji et al., 2023; Min et al., 2023; Chen et al., 2023), UQ methods assume that LLMs naturally encode information about their own limitations, and this self-knowledge can be efficiently accessed.

There are many existing UQ techniques for well-defined tasks such as classification and regression (Zhang et al., 2019; He et al., 2020; Xin et al., 2021; Wang et al., 2022; Vazhentsev et al., 2023a; He et al., 2024a). However, applying UQ to text generation has unique challenges, including (i) potentially multiple correct answers with different surface forms (Kuhn et al., 2023), (ii) the need to aggregate uncertainties across multiple conditionally dependent predictions (Zhang et al., 2023), (iii) generated tokens not contributing to uncertainty equally (Duan et al., 2024), and (iv) some sources of uncertainty being irrelevant for hallucination detection (Fadeeva et al., 2024). These challenges hinder the performance of classical unsupervised UQ techniques, as they are difficult to address explicitly within a single method. Recently, researchers have proposed learning the aforementioned intricacies from the annotated data and developed supervised methods for UQ and hallucina-

$\diamond$  Equal contribution

<sup>1</sup><http://uncertainty-head.nlpresearch.group>

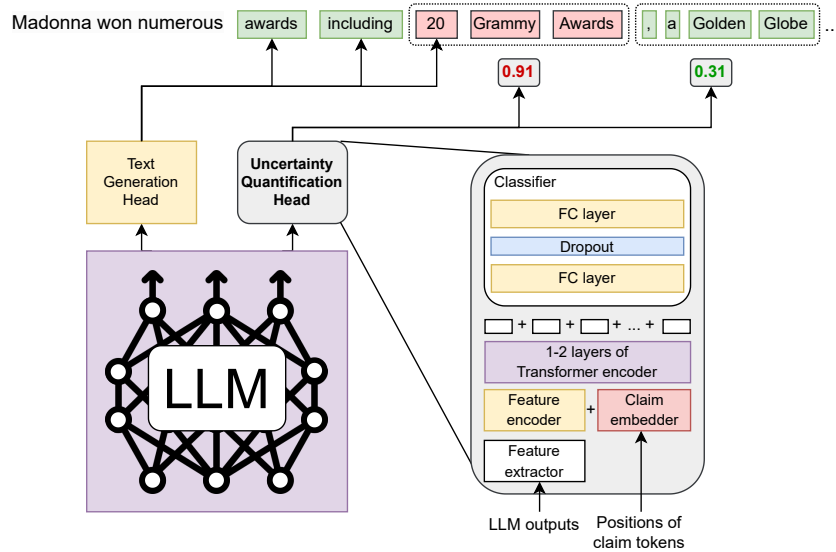


Figure 1: The architecture of uncertainty quantification heads. The example represents a text generated using an LLM, containing the hallucination *20 Grammy Awards* highlighted in red.

tion detection (Azaria and Mitchell, 2023; Li et al., 2024; He et al., 2024b; Chuang et al., 2024).

We continue this line of work by introducing pre-trained UQ heads: supervised auxiliary modules for LLMs that substantially enhance their ability to capture uncertainty compared to unsupervised UQ methods (Figure 1). Their strong performance stems from the transformer architecture, in deriving informative features from LLM attention maps and logits. These heads do not require re-training of the entire LLM and do not alter its outputs. In addition to their high performance, these modules maintain a relatively small memory and computational footprint, ensuring practical usability.

Our experiments show that the proposed uncertainty heads are highly robust and achieve state-of-the-art performance in claim-level hallucination detection across both in-domain and out-of-domain prompts, outperforming other supervised and unsupervised techniques. Moreover, these modules demonstrate strong generalization to languages they were not explicitly trained on.

Training UQ heads requires annotated hallucinations in LLM outputs. To construct training data, we created an automatic annotation pipeline, which allowed us to scale experiments and to pre-train UQ heads for various LLMs. We release a collection of pre-trained UQ heads for popular open-source instruction-following LLMs, including the Llama series (Grattafiori et al., 2024), Gemma 2 (Team et al., 2023), and Mistral (Jiang et al., 2023). The

**contributions** of this work are as follows:

- We design a pre-trained uncertainty quantification head: a supplementary module for an LLM that yields substantially better performance at claim-level hallucination detection than classical unsupervised UQ methods and state-of-the-art supervised techniques.
- We conduct an extensive empirical investigation and find that uncertainty heads show good generalization across various domains and languages. We perform a comprehensive ablation study that compares various feature sets architectures, and approaches to training data generation.
- We build and release a collection of pre-trained UQ heads for popular series of open-source instruction-tuned LLMs. These modules could be seamlessly integrated into text generation code and be used as off-the-shelf hallucination detection tools.

## 2 Related Work

**Unsupervised UQ methods** for LLMs can be broadly categorized into five groups: information-based approaches (Kuhn et al., 2023; Farquhar et al., 2024), density-based scores (Vazhentsev et al., 2022, 2023b; Ren et al., 2023), self-consistency methods (Manakul et al., 2023; Lin et al., 2024; Zhang et al., 2024a; Qiu and Miikkulainen, 2024), methods grounded in mechanistic analysis of LLMs (Yüksekgönül et al., 2024; Qiu

and Miikkulainen, 2024), and verbalized (reflexive) strategies (Kadavath et al., 2022; Tian et al., 2023). Although all of them have demonstrated potential, their effectiveness in hallucination detection remains limited (Vashurin et al., 2025).

**Supervised UQ methods** leverage the internal states of LLMs as features for predicting hallucinations (Azaria and Mitchell, 2023; Slobodkin et al., 2023; Su et al., 2024; CH-Wang et al., 2024; He et al., 2024b; Chuang et al., 2024; Vazhentsev et al., 2025b,a). These recently-developed methods achieve substantial performance gains over unsupervised approaches, especially for in-domain data.

Azaria and Mitchell (2023) proposed one of the first methods of this kind called SAPLMA, where they trained a perceptron with layer activations as features to detect when a LLM “agrees” with false statements. Slobodkin et al. (2023) trained a linear model on hidden states to detect question “answerability”, effectively identifying unanswerable questions that typically lead to hallucinations. CH-Wang et al. (2024) extend this approach to span-level hallucination detection, using manually annotated hallucinations. He et al. (2024b) experiment with activation maps, token ranks, and probabilities from unembedding matrices across layers. Chuang et al. (2024) introduce a feature set derived from LLM attention maps.

**Limitations of previous methods.** Azaria and Mitchell (2023); Slobodkin et al. (2023); Su et al. (2024) focused on sequence-level methods and are not able to detect fragment-level hallucinations. Many models, including Slobodkin et al. (2023); Azaria and Mitchell (2023); Chuang et al. (2024); Su et al. (2024) used non-contextualized architectures such as simple linear probes or multi-layer perceptron. Although He et al. (2024b) integrated a linear model with an attention mechanism and CH-Wang et al. (2024) used a contextualized model combining convolutions, ResNet, and GRU, these architectures are considered outdated and exhibit limitations in quality and computational efficiency. The features of the majority of models included only hidden states (Azaria and Mitchell, 2023; Slobodkin et al., 2023; CH-Wang et al., 2024; Su et al., 2024), which limits their generalization. Only He et al. (2024b) and Chuang et al. (2024) performed more elaborate feature engineering. Finally, synthetic data that is leveraged through enforced decoding is used in some work (Azaria and Mitchell, 2023; Slobodkin et al., 2023). Compared to the

native outputs generated by LLMs, such data may introduce additional biases and adversely affect the performance of hallucination detectors.

In contrast to most prior work, we focus on building UQ heads specifically for detecting hallucinations at the fragment level, i.e., individual atomic claims. Our approach leverages the strengths of previous methods while addressing their key limitations: (i) instead of outdated architectures, we build our solution on the transformer architecture; (ii) we investigate the importance of various feature sets for hallucination detection, finding that the most informative features are derived from attention maps of LLMs; and (iii) we build an automatic pipeline to generate training data using *native* LLM responses. This pipeline allows us to build training data at a larger scale and pre-train UQ heads for a range of popular LLMs.

### 3 Uncertainty Quantification Head

Consider the LLM  $P(t_i | \mathbf{x}, \mathbf{t}_{<i})$  with  $L$  layers receiving a prompt  $\mathbf{x}$  of length  $n$  and generating tokens  $\mathbf{t} = \{t_1, t_2, \dots, t_T\}$ . We also have a set of atomic claims  $C = \{c_1, c_2, \dots, c_K\}$ , each representing a mapping to a subset of tokens in the output. Atomic claims, for example, can be extracted by another lightweight model. In this work, we formalize the claim-level uncertainty quantification task as building a function  $U(c_i | \mathbf{x}, \mathbf{t}) \in [0, 1]$  that determines whether the claim  $c_i \in C$  is a hallucination. A large value of  $U(c_i | \mathbf{x}, \mathbf{t})$  indicates a higher likelihood that the claim  $c_i$  is a hallucination.

#### 3.1 Background on Features for UQ and Hallucination Detection

**Hidden states**  $h(t)$  extracted from LLM layers during the generation of a token  $t$  have been shown to serve as strong indicators of hallucinations in several studies (Azaria and Mitchell, 2023; CH-Wang et al., 2024).

$$F_{\text{hs}}(t) = h(t) \quad (1)$$

**Lookback Lens** (Chuang et al., 2024) derives features from the LLM’s attention maps. The key idea is that when the model attends to the prompt, it attempts to solve the task, whereas attending to generated tokens causes it to disregard the prompt, increasing the likelihood of hallucination. They suggest using the so-called lookback ratio – the ratio of aggregated attention to tokens of the prompt and the generated tokens. Consider each layer of the

LLM contains  $Q$  attention heads,  $q$  is an index of a head, and  $\alpha_{ij}^{q,l}$  represents the softmax-weighted attention score from token  $t_i$  to token  $t_j$ .  $A_{\text{context}}^{q,l}(t_i)$  and  $A_{\text{gen}}^{q,l}(t_i)$  are the average attention weights to the input  $\mathbf{x}$  and to the previously generated output  $\mathbf{t}_{<i}$ , respectively:

$$A_{\text{context}}^{q,l}(t_i) = \frac{1}{|\mathbf{x}|} \sum_{j=1}^{|\mathbf{x}|} \alpha_{ij}^{q,l},$$

$$A_{\text{gen}}^{q,l}(t_i) = \frac{1}{i-1-|\mathbf{x}|} \sum_{j=|\mathbf{x}|+1}^{i-1} \alpha_{ij}^{q,l}.$$

Then the lookback ratio of the model head  $q$  and the layer  $l$  for the token  $t_i$  is defined as follows:

$$LR^{q,l}(t_i) = \frac{A_{\text{context}}^{q,l}(t_i)}{A_{\text{context}}^{q,l}(t_i) + A_{\text{gen}}^{q,l}(t_i)},$$

$$F_{\text{LBLEns}}(t_i) = \{LR^{q,l}(t_i)\}_{q,l}^{Q,L}. \quad (2)$$

**Factoscope** (Min et al., 2023) in addition to model activations, uses a set of features that leverage token probabilities, the similarity of token embeddings across layers, and the evolution of token ranks across layers. Commonly, given a token  $t_i$  at the position  $i$ , the LLM outputs hidden states  $\{h_l(t_i)\}_{l=1}^L$ , where the final hidden state  $h_L(t_i)$  is passed through the unembedding matrix  $E$  to predict token logits. Factoscope applies  $E$  to each LLM layer, obtaining a set of token logits on a specific layer  $l$ :  $z_i^l = E(h_l(t_i))$ . Then, it extracts the logits of the top- $m$  tokens from each layer  $l$ :

$$F_{\text{top-tokens}}(t_i) = \left\{ z_i^l(t) \mid t \in \text{top}_m(z_i^l) \right\}_{l=1}^L. \quad (3)$$

To analyze token evolution across layers, Factoscope computes the cosine similarities between embeddings of top tokens from adjacent layers obtained by applying the unembedding matrix:

$$S^l(t_i) = \{ \cos(E_{w_1}, E_{w_2}) \mid$$

$$w_1 \in \text{top}_m(z_i^l), w_2 \in \text{top}_m(z_i^{l+1}) \}$$

$$F_{\text{tokens-sim}}(t_i) = \{S^l(t_i)\}_{l=1}^{L-1}. \quad (4)$$

Finally, Factoscope tracks token rank evolution across layers:  $R^l(t_i) = \text{rank}[t_i, z_i^l]$ , where rank indicates the position of  $t_i$  in the descending order of  $z_i^l$  values (top-ranked token receives 1). The ranks are further normalized to the range  $[0, 1]$ :

$$F_{\text{rank}}(t_i) = \{R^l(t_i)^{-1}\}_{l=1}^L. \quad (5)$$

### 3.2 Features for Pre-trained UQ Heads

We experimented with all the aforementioned types of features and their combinations. However, we found that all of them exhibited various limitations. Hidden states encode a lot of domain-specific information, increasing the risk of overfitting. Factoscope features incur substantial computational overhead while offering limited additional information beyond what is captured by hidden states. Attention features are quite powerful, but the aggregation suggested in Lookback Lens results in the loss of valuable information. Moreover, they underperform without the addition of logits or probabilities. Therefore, for our pre-trained UQ heads, we use two groups of features.

**Attention maps of the LLM.** Mechanistic analysis of attention weights reveals that attention patterns often reflect the model’s behavior under uncertainty (Yüksekgönül et al., 2024). Moreover, attention encodes the conditional dependency between the generation steps (Zhang et al., 2023). For each token, we obtain the attention maps to  $k$  previous tokens from each attention head and layer and flatten them into a single feature vector without aggregation:

$$F_{\text{att}}(t_i) = \{\alpha_{i,i-j}^{q,l}\}_{j,q,l}^{k,Q,L}. \quad (6)$$

When  $(i-j)$  is negative, we pad the feature vector with zero placeholders. While considering many previous tokens might explode the feature space size, we empirically found that the optimal value of  $k$  is typically very small:  $1 \leq k \leq 5$  (see Figure 5). As a contextualized architecture, the transformer can automatically extract meaningful patterns across the entire generated sequence without requiring explicit features from previous tokens.

**Probability distribution of the LLM** might be misleading, but it still conveys useful information about the model’s conditional confidence at the current generation step. This group of features consists of logarithms of the top- $m$  token probabilities:

$$F_{\text{prob}}(t_i) = \{ \log P(t \mid \mathbf{x}, \mathbf{t}_{<i}) \mid$$

$$t \in \text{top}_m(P(\cdot \mid \mathbf{x}, \mathbf{t}_{<i})) \}. \quad (7)$$

Features from both groups are concatenated into a token-level vector:  $F(t) = F_{\text{att}}(t) \circ F_{\text{prob}}(t)$ .

### 3.3 Architecture of UQ Heads

The architecture of the UQ head is depicted in Figure 1. To ensure flexibility and expressive capacity, we build it on top of a transformer backbone.



It consists of a feature size reduction network, a multi-layer transformer encoder, and a two-layer classification neural network. For each component, we use GELU activation functions and dropout regularization.

Consider that for a prompt  $\mathbf{x}$ , a LLM generates a text  $\mathbf{y}$  that includes an atomic claim  $c$ . The uncertainty score of a claim  $u_c$  is defined as:

$$u_c = P(\text{"}c \text{ is a hallucination"} \mid \mathbf{x}, \mathbf{t}).$$

Below is the step-by-step algorithm for obtaining the uncertainty score. For simplicity, we describe the process for a single claim  $c$ .

**Token-level feature extraction.** For each token  $t$  in the sequence  $y$ , we compute a token-level feature vector  $F(t)$ , and the extracted feature vectors are then passed through a fully connected (FC) projection layer,  $\text{FC}_{\text{proj}}$ , to match the hidden dimension of the transformer encoder:

$$\tilde{f}_t = \text{FC}_{\text{proj}}(F(t)).$$

**Claim-specific contextualization.** We use a transformer encoder to obtain the contextualized representation of the claim. To focus the model on the specific claim  $c$ , for each token  $t \in c$ , we augment its projected feature vector with a single, trainable claim-marking embedding  $E$ :

$$\tilde{f}_{t,c} = \tilde{f}_t + E \cdot 1(t \in c),$$

where  $1(\cdot)$  is an indicator function.

**Transformer encoding.** The contextualized feature sequence, denoted as  $\tilde{F}_c = \{\tilde{f}_{t,c}\}_{t \in \mathbf{x} \circ \mathbf{y}}$ , is processed by a transformer encoder, producing a sequence of contextualized hidden states  $H_c = \{h_{t,c}\}_{t \in \mathbf{x} \circ \mathbf{y}}$ :

$$H_c = \text{Transformer}(\tilde{F}_c).$$

**Pooling and representation.** To derive a single representative vector for the claim  $c$ , we perform masked average pooling over the transformer’s output. We average only the contextualized token representations  $h_{t,c} \in H_c$  corresponding to the tokens within the claim:

$$h_c = \frac{1}{|c|} \sum_{t \in c} h_{t,c}.$$

**Classification.** Finally, the claim representation vector  $h_c$  is passed through a two-layer

classification network regularized with dropout,  $\text{MLP}_{\text{classifier}}$ , followed by a sigmoid activation that produces the uncertainty score  $u_c$ :

$$u_c = \sigma(\text{MLP}_{\text{classifier}}(h_c)).$$

**Training loss.** Consider that we have a dataset  $\mathcal{D} = \{(\mathbf{x}, \mathbf{y}, c, v)\}$ , where for each claim  $c$  in answer  $\mathbf{y}$  for the prompt  $\mathbf{x}$ , we have an annotation of its ground-truth veracity  $v$ . The UQ head is trained using a binary cross-entropy loss function:

$$\mathcal{L} = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}, c, v) \sim \mathcal{D}} [v \log u_c + (1-v) \log (1-u_c)].$$

When we train heads, we freeze the “body” of the LLM, so that its generations stay exactly the same.

## 4 Pipeline for Training Data Generation

The training data generation pipeline is presented in Figure 4 in Appendix A. It starts with prompting the LLM to produce responses for a list of questions such as *Write a biography of person X* or *Write the history of the city Y*. We select relatively famous named entities so the task is not very hard for the model based on its parametric knowledge, while at the same time, it is not trivial, so outputs contain a substantial number of hallucinated claims. We also do not use synthetically-generated hallucinations, as they introduce a bias between what the model actually generates vs. the synthetic data. The prompts for other domains can be found in Table 6 in Appendix B.

We split the obtained responses into atomic claims using GPT-4o based on the prompts from (Fadeeva et al., 2024; Vashurin et al., 2025). Each claim is then automatically classified by GPT-4o as *supported*, *unsupported*, or *unknown*. The last category is intended for general claims, for which estimating the veracity is meaningless. To ensure high annotation quality, the claim labeling process is two-staged: first, we ask the model to provide an elaborated answer via chain-of-thought; then, we ask it to summarize its answer into one word. As shown in Table 4 in Appendix A, the performance of such annotation using GPT-4o is high (around 90%, with the exception of German, where it is slightly lower at 83%). This performance could likely be further improved by using more advanced LLMs or by employing model ensembles.

The pipeline enables the cost-effective construction of large datasets annotated with claim-level hallucinations across various LLMs. The cost of

annotating responses from a single LLM on the training biographies dataset, consisting of 3,300 prompts, was approximately \$100. Statistics about the training dataset, as well as the accuracy of LLM responses, are presented in Table 5.

## 5 Experiments

### 5.1 Experimental Setup

**Evaluation datasets.** We constructed eight test sets of English questions designed to prompt LLMs to generate text across various domains: *person biographies*, *cities*, *movies*, *inventions*, *books*, *artworks*, *landmarks*, and *events*. Each test set contains 100 questions, generated by prompting GPT-4o and Claude-3-Opus to output 100 famous domain items, e.g., 100 famous landmarks. Examples of the prompts are presented in Appendix B.1.<sup>2</sup> The labels for the test sets are obtained using the same annotation pipeline as for the training data.

To assess the cross-lingual generalizability of pre-trained UQ modules, we conducted evaluations on Russian and Chinese prompts from Vashurin et al. (2025), and additionally created a similar test set with German prompts. Test sets for each language consist of 100 biography-related questions. The data statistics are presented in Table 6.

**Metrics.** In the main experiments, we measured the claim-level performance of detecting invalid claims. For this purpose, we used PR-AUC, where “unsupported” claims represent the positive class.

**Models.** We conducted our primary experiments with Mistral 7b Instruct v0.2 (Jiang et al., 2023) and Gemma 2 9b Instruct (Team et al., 2023).

**Training procedure and hyper-parameter optimization.** We trained the uncertainty heads using Adam with a linear learning rate decay and warmup. We selected the values of the hyper-parameters on the validation set of the *biographies* dataset using the claim-level PR-AUC metric and the Bayesian optimization algorithm available in the W&B framework. We observed that among important hyper-parameters are the weight of instances with positive labels, the number of epochs, and the learning rate. The best values of hyper-parameters for each of the tested models are presented in Table 14 in Appendix F.

**Baselines.** We compare our method to several unsupervised baselines: Maximum Claim Probability

(an adaptation of Maximum Sequence Probability for claims), Mean Token Entropy, Perplexity, Claim Conditioned Probability (CCP) (Fadeeva et al., 2024), and Attention Score (Qiu and Mikkulainen, 2024). Furthermore, we evaluated our UQ heads against supervised methods, including SAPLMA, Factoscope, and Lookback Lens. SAPLMA predicts token-level uncertainties using a 3-layer perceptron, and the mean uncertainty is calculated over claim-related tokens during inference. We adapt Lookback Lens and Factoscope to the claim level. Lookback Lens uses a Logistic Regression model trained on lookback ratios. Our implementation of Factoscope uses our transformer-based architecture and the feature set that includes hidden states, top token embeddings with similarities, and token ranks. The values of the hyperparameters for the baselines selected after tuning are given in Appendix F.

### 5.2 Results

**Main results.** Table 1 shows the performance of the unsupervised UQ techniques and the supervised detectors trained on persons’ *biographies* for claim-level hallucination detection with Mistral 7B Instruct v0.2. To evaluate supervised methods, the domain *biographies* represents the in-domain test set and all other domains (Cities, Movies, Inventions, Books, Artworks, Landmarks, and Events) represent out-of-domain (OOD) test sets. Note that in this evaluation, both the questions and the LLM’s responses across all domains are in English.

Among the unsupervised techniques, uncertainty scores based on CCP yield the best performance, confidently outperforming other methods on *biographies*, *cities*, *artworks*, and *landmarks*.

Supervised UQ methods greatly outperform unsupervised techniques on the in-domain test set. Moreover, remarkably, all considered supervised methods demonstrate substantial generalization and the ability to perform well beyond the training domain of people’s *biographies*.

Our UQ head (UHead) demonstrates the best results in both in-domain and out-of-domain evaluations. For in-domain evaluation, UHead outperforms the best unsupervised method CCP by 16% (absolute) in terms of PR-AUC. The gap is also large for out-of-domain evaluation, e.g., for *books*, UHead outperforms CCP by 23%, for *movies* and *events* by 20%, for *artworks* by 18%. Table 9 in Appendix E also shows that UHead pre-trained on *biographies* generalizes to question an-

<sup>2</sup>All data used for training and testing is available at <https://huggingface.co/llm-uncertainty-head>

Method \ Test Sets	Biographies (in domain)	Cities	Movies	Inventions	Books	Artworks	Landmarks	Events
Random	.291 ± .014	.205 ± .013	.099 ± .008	.163 ± .012	.110 ± .010	.264 ± .014	.117 ± .010	.113 ± .009
MCP	.412 ± .020	.310 ± .023	.205 ± .020	.319 ± .024	.145 ± .014	.317 ± .017	.135 ± .012	.141 ± .017
Perplexity	.361 ± .018	.231 ± .017	.170 ± .016	.232 ± .019	.138 ± .014	.335 ± .018	.128 ± .011	.123 ± .013
Max Token Entropy	.416 ± .019	.289 ± .022	.241 ± .025	.381 ± .031	.171 ± .019	.321 ± .015	.141 ± .017	.161 ± .019
Attention Score	.333 ± .019	.279 ± .018	.114 ± .011	.211 ± .017	.114 ± .011	.202 ± .010	.125 ± .011	.132 ± .012
CCP	.496 ± .019	.368 ± .025	.267 ± .027	.380 ± .028	.167 ± .018	.382 ± .022	.196 ± .019	.171 ± .018
SAPLMA	.536 ± .021	.435 ± .027	.269 ± .030	.350 ± .025	.292 ± .029	.534 ± .020	<b>.350</b> ± .030	.235 ± .025
Factoscope	<b>.611</b> ± .021	<b>.468</b> ± .029	<b>.344</b> ± .028	<b>.424</b> ± .029	<b>.315</b> ± .030	<b>.485</b> ± .019	.279 ± .026	<b>.265</b> ± .025
Lookback Lens	.557 ± .021	.449 ± .025	.254 ± .025	.391 ± .027	.259 ± .028	.464 ± .021	.257 ± .025	<b>.295</b> ± .030
UHead (Ours)	<b>.660</b> ± .020	<b>.487</b> ± .028	<b>.466</b> ± .036	<b>.485</b> ± .027	<b>.395</b> ± .033	<b>.561</b> ± .020	<b>.340</b> ± .024	<b>.369</b> ± .030

Table 1: PR-AUC for various UQ methods for hallucination detection of the Mistral 7B Instruct v0.2 model on English datasets. Biographies represent the in-domain dataset for supervised UQ methods. The standard deviation is estimated using the bootstrap method.

Method \ Language	English (in domain)	Russian	Chinese	German
Random	.133 ± .010	.337 ± .012	.226 ± .012	.152 ± .010
MCP	.180 ± .017	.433 ± .016	.307 ± .017	.203 ± .016
Perplexity	.136 ± .012	.395 ± .015	.287 ± .016	.149 ± .009
Max Token Entropy	.202 ± .020	.437 ± .014	.444 ± .021	.217 ± .017
Attention Score	.146 ± .018	.446 ± .026	.230 ± .017	.229 ± .023
CCP	.307 ± .024	.493 ± .014	.439 ± .023	.306 ± .024
SAPLMA	.342 ± .023	.514 ± .019	.331 ± .019	<b>.391</b> ± .023
Factoscope	.354 ± .026	.532 ± .018	.350 ± .023	.380 ± .023
Lookback Lens	<b>.359</b> ± .025	<b>.576</b> ± .016	<b>.479</b> ± .024	.390 ± .023
UHead (Ours)	<b>.457</b> ± .026	<b>.581</b> ± .017	<b>.556</b> ± .023	<b>.455</b> ± .025

Table 2: PR-AUC of UQ methods on various languages using the Gemma 2 9b Instruct model. Supervised detectors were trained on English-only *biographies* data. The standard deviation is estimated using bootstrap.

Method \ Test Set	Biographies (dev)
UHead (only hidden states)	.582
UHead (att. + probs. + hs.)	.589
UHead (Factoscope)	.588
UHead (LookBack Lens)	.609
UHead (att.)	.617
UHead (att. + probs.) (ours)	<b>.642</b>

Table 3: PR-AUC scores for UQ heads trained with various feature sets on the Mistral 7B Instruct v0.2 model. Performance was evaluated using the validation set of the *biographies* domain after hyperparameter tuning.

swering on the TruthfulQA and SciQ datasets (Lin et al., 2022), outperforming unsupervised baselines.

When evaluated alongside supervised methods, UHead surpasses the closest competitor, Factoscope, by 5% for the in-domain evaluation. In OOD evaluation, it confidently outperforms other supervised methods across all domains, except for landmarks, where it is slightly below the closest competitor by 1%.

Analyzing other supervised methods, the second-best scores are usually demonstrated by Factoscope. We assume that the underperformance of the base-

line based on the Factoscope features compared to UHead lies in the use of layer activations, which limits its generalization. Another module that relies on hidden states is SAPLMA. In addition to the feature limitations, it also has architectural limitations, which further hurt its performance. For *landmarks*, SAPLMA shows good results, but for other test sets, it stays behind Factoscope and UHead. Compared to UHead, it lags by 12% on in-domain evaluation and up to 20% on OOD evaluation. Lookback Lens also usually falls behind UHead and Factoscope; we believe that its main problem is its weak linear architecture.

The similar evaluation for Llama-3.1 (Table 13 in Appendix E) shows a similar pattern; UHead outperforms all other supervised and unsupervised baselines in the majority of domains. Only for *inventions* and *events*, UHead slightly falls behind Factoscope and Lookback Lens.

**Cross-lingual generalization.** Table 2 presents the cross-lingual results for Gemma 2 9b Instruct. In this experiment, we train UQ modules on the English person’s *biographies* as in the previous experiment, but we evaluate the performance on other languages. Surprisingly, UHead achieves strong cross-lingual generalization. For all OOD languages, UHead achieves substantial improvements over the best unsupervised methods. For Chinese, UHead is better than MTE by 10%; for Russian, it is better than CCP by 9%; and for German by 13%. Notably, other supervised methods also demonstrate some level of generalization, but in most cases, they have substantially worse performance. Overall, these results show that UQ heads, even if they are pre-trained on English data, can serve as effective off-the-shelf hallucination detectors for LLM outputs in other languages.

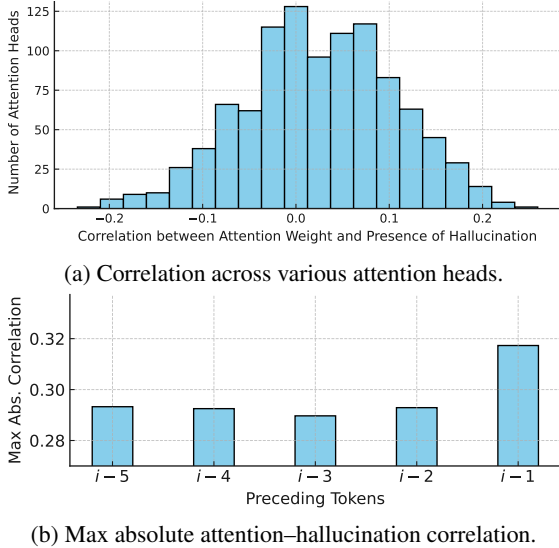


Figure 2: (a) The distribution of correlations between attention on the  $i - 1$ -th token and presence of the  $i$ -th token in hallucinated claim. (b) The maximum absolute correlation across heads and layers for the same phenomenon. All scores were computed using the Mistral model and the *biographies* dataset.

**Analysis of feature sets.** Table 3 presents the comparison of various feature sets in combination with the UHead architecture on the in-domain validation set. For each feature set, we perform an extensive hyper-parameter value search in the same way as for the main results. We can see that all feature sets that leverage hidden states fall substantially behind attention-based features. The analysis of the validation loss dynamics shows that this is probably due to quick overfitting. Models that leverage hidden states start overfitting after 1–3 epochs, while models that leverage attention might not overfit even after 10 epochs. We also note that Lookback Lens features combined with the UHead architecture provide strong performance. However, simple attention maps without feature engineering used in UHead yield even better results. Finally, without probability-based features, UHead loses around 2.5% PR-AUC, which marks their importance.

**Analysis of attention-based features.** We examined attention patterns that may serve as indicators of hallucinations in generated tokens. Figure 2a shows the correlation between the presence of hallucinations and attention weights from the generated token to the immediately preceding token across various attention heads. While most heads show negligible correlation, a subset of heads exhibits moderate positive or negative associations.

Figure 2b further highlights that this correlation is strongest for the token immediately preceding the generated one. Thus, a small subset of attention heads encodes informative signals related to hallucinations and reflects distinct model behavior under uncertainty during generation.

These findings are also confirmed by Figures 5 and 6 in Appendix D. Figure 6 illustrates that attention weights from individual middle layers could serve as relatively strong hallucination detectors. Figure 5 shows that optimal performance is obtained by UHead when using attention weights from only 1–5 preceding tokens.

Since only a small subset of attention heads typically correlates with presence hallucinations, we explore whether the feature set for UHead can be reduced. Table 8 in Appendix D reports PR-AUC results on the evaluation subset when training UHead using only the top- $N$  (layer, head) pairs, ranked by their absolute correlation with the training labels. The results indicate that the feature set can be reduced by roughly eightfold (to 128 heads) without a significant drop in performance. However, further reductions in the number of (layer, head) pairs result in a noticeable decline in performance.

**Analysis of detector architectures.** Table 10 in Appendix E reports the performance of detectors with various architectures trained on our best feature set, consisting of attention maps and top token probabilities. We compare the transformer-based architecture used in UHead against simpler alternatives: MLP and a linear model. Although both simpler models yield notable improvements over the best unsupervised baselines, UHead based on transformer achieves the highest performance.

**Introducing more diverse training data for UHead.** Table 12 in Appendix E presents the results when we train uncertainty heads on *biographies* plus the data from all domains except one, which is used for OOD evaluation. In this scenario, uncertainty heads get access to bigger and more diverse training data. As we can see, expanding the dataset provides slight improvements for certain domains. These results indicate that expanding the training data and enhancing its diversity could further increase the UQ performance, particularly in the OOD setting.

**Using “non-native” training data.** We also analyzed the possibility of using the training data generated for one LLM for training a detector for another LLM. We take the annotated dataset



```

from transformers import AutoModelForCausalLM,
AutoTokenizer
from luh import AutoUncertaintyHead,
CausalLMWithUncertainty

llm = AutoModelForCausalLM.from_pretrained(
    model_name)
tokenizer = AutoTokenizer.from_pretrained(
    model_name)
uhead = AutoUncertaintyHead.from_pretrained(
    uhead_name, base_model=llm)
llm_adapter = CausalLMWithUncertainty(llm, uhead,
    tokenizer=tokenizer)

# tokenize text and prepare inputs ...
output = llm_adapter.generate(inputs)

```

Figure 3: Code example for using uncertainty heads.

generated by Mistral and performed inference of Gemma 2 via forced decoding to generate features. Table 11 in Appendix E compares the results of hallucination detectors for Gemma 2 trained using “native” and “non-native” data. We can see that “non-native” data still yields better results than unsupervised methods, but substantially decreases the performance of the hallucination detector due to distribution shift. Therefore, for each new LLM, we recommend generating a new training dataset.

**Computational efficiency.** We evaluated the computational overhead of various UQ methods. To ensure a fair comparison, we focused only on the time required to generate texts and to compute uncertainty scores, excluding the time spent on claim extraction. Claim extraction could be performed by a small model specifically fine-tuned for this task, and its overhead is negligible compared to LLM inference. Table 7 summarizes the results and provides the memory footprint of methods. MCP and Perplexity incur no additional overhead, serving as baselines for comparison. Our UHead introduces only 5% overhead, which is even better than the best unsupervised method CCP (8.6%). UHead also has a minimal GPU memory footprint (40 MB). Thus, UHead is a very lightweight addition to multi-billion-parameter LLMs and is practical for real-world deployment.

## 6 Collection of Pre-trained Uncertainty Heads for Popular LLMs

We pre-trained a collection of UQ heads for a range of popular 7B–9B LLMs, including Mistral, LLaMA series, and Gemma 2. In addition to model-level UQ, we release token-level UQ heads that can provide uncertainty scores directly for tokens without explicit claim annotation, which enables broader applicability. Our UQ heads are designed for use as an off-the-shelf tool for confi-

dence estimation in LLMs. They could be loaded from the hub using a procedure similar to the `from_pretrained` API in the HuggingFace transformers library and integrated into the LLM generation procedure with an adapter. A code example is provided in Figure 3. Thus, UQ heads could be integrated into third-party code with minimal modifications and could be used as a plug-and-play solution for researchers and practitioners. Examples of UQ head predictions are in Appendix G.

## 7 Conclusion and Future Work

We presented pre-trained UQ heads – supplementary supervised modules for LLMs that help to capture their uncertainty much more effectively than unsupervised UQ methods. We demonstrated that they are quite robust and deliver state-of-the-art results for both in-domain and out-of-domain prompts. They also show remarkable generalization to other languages. Inspired by their good performance, we pre-trained a collection of UQ heads for a series of popular LLMs, including Mistral, Gemma 2, and LLaMA series. We release the code and the pre-trained uncertainty heads so they could be used as off-the-shelf hallucination detectors for other researchers and practitioners. In future work, we plan to scale up the training data and explore the limits of the supervised approach to UQ.

### Limitations

Uncertainty heads cannot solve the problem when LLMs are trained to provide misinformation. In this situation, models are confident in their deceptive answers. Uncertainty heads cannot provide ideal annotation of hallucinations, as some LLMs do not have enough capacity to provide information about what they know and what they do not know. While we see generalization in uncertainty heads, we should acknowledge that, as with any other supervised method, they work best for “in-domain” data and “in-domain” tasks. Further investigation is needed to assess their transferability to other tasks, such as machine translation and summarization. The bias present in LLMs could also be transferred into uncertainty heads.

### Ethical Considerations

In our work, we considered open-weight LLMs and datasets not aimed at harmful content. However, LLMs may generate potentially damaging texts for

various groups of people. Uncertainty quantification techniques can help create a more reliable use of neural networks.

Despite our proposed method demonstrating sizable performance improvements, it can still mistakenly highlight correctly generated text with high uncertainty in some cases. Thus, as with other uncertainty quantification methods, it is an imperfect technology and users should be aware of the limitations of this technology.

We release our source code under the MIT license for broader adoption. We used writing assistants to ensure grammatical correctness throughout the text.

## Acknowledgements

We sincerely thank the reviewers for their constructive and insightful feedback. We are also grateful to Gleb Kuzmin for his valuable contribution to the additional evaluations.

## References

- Amos Azaria and Tom Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, Singapore. Association for Computational Linguistics.
- Joris Baan, Nico Daheim, Evgenia Ilia, Dennis Ulmer, Haau-Sing Li, Raquel Fernández, Barbara Plank, Rico Sennrich, Chrysoula Zerva, and Wilker Aziz. 2023. [Uncertainty in natural language generation: From theory to applications](#). *arXiv preprint arXiv:2307.15703*.
- Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and Chris Kedzie. 2024. [Do androids know they’re only dreaming of electric sheep?](#) In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4401–4420, Bangkok, Thailand. Association for Computational Linguistics.
- Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li, and Yanghua Xiao. 2023. [Hallucination detection: Robustly discerning reliable answers in large language models](#). In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023*, pages 245–255. ACM.
- Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ranjay Krishna, Yoon Kim, and James R. Glass. 2024. [Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1419–1436, Miami, Florida, USA. Association for Computational Linguistics.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2024. [Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5050–5063, Bangkok, Thailand. Association for Computational Linguistics.
- Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. [Fact-checking the output of large language models via token-level uncertainty quantification](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9367–9385, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. [Detecting hallucinations in large language models using semantic entropy](#). *Nature*, 630(8017):625–630.
- Yarin Gal and Zoubin Ghahramani. 2016. [Dropout as a Bayesian approximation: Representing model uncertainty in deep learning](#). In *Proceedings of The 33rd International Conference on Machine Learning Research*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA. PMLR.
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppel, Preslav Nakov, and Iryna Gurevych. 2024. [A survey of confidence estimation and calibration 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 6577–6595, Mexico City, Mexico. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, et al. 2024. [The Llama 3 Herd of Models](#). *arXiv preprint arXiv:2407.21783*.
- Jianfeng He, Linlin Yu, Shuo Lei, Chang-Tien Lu, and Feng Chen. 2024a. [Uncertainty estimation on sequential labeling via uncertainty transmission](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2823–2835, Mexico City, Mexico. Association for Computational Linguistics.

- Jianfeng He, Xuchao Zhang, Shuo Lei, Zhiqian Chen, Fanglan Chen, Abdulaziz Alhamadani, Bei Xiao, and ChangTien Lu. 2020. [Towards more accurate uncertainty estimation in text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8362–8372, Online. Association for Computational Linguistics.
- Jinwen He, Yujia Gong, Zijin Lin, Cheng’an Wei, Yue Zhao, and Kai Chen. 2024b. [LLM factoscope: Uncovering LLMs’ factual discernment through measuring inner states](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10218–10230, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Transactions on Information Systems*, 43(2):1–55.
- Yukun Huang, Yixin Liu, Raghuv eer Thirukovalluru, Arman Cohan, and Bhuwan Dhingra. 2024. [Calibrating long-form generations from large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13441–13460, Miami, Florida, USA. Association for Computational Linguistics.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM computing surveys*, 55(12):1–38.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *arXiv preprint arXiv:2310.06825*.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, et al. 2022. [Language models \(mostly\) know what they know](#). *arXiv preprint arXiv:2207.05221*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. [Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- Qing Li, Jiahui Geng, Chenyang Lyu, Derui Zhu, Maxim Panov, and Fakhri Karray. 2024. [Reference-free hallucination detection for large vision-language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4542–4551, Miami, Florida, USA. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. [Generating with confidence: Uncertainty quantification for black-box large language models](#). *Transactions on Machine Learning Research*.
- Andrey Malinin and Mark J. F. Gales. 2021. [Uncertainty estimation in autoregressive structured prediction](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.
- Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. [SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, Singapore. Association for Computational Linguistics.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FActScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100.
- Xin Qiu and Risto Miikkulainen. 2024. [Semantic density: Uncertainty quantification for large language models through confidence measurement in semantic space](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 134507–134533. Curran Associates, Inc.
- Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J Liu. 2023. [Out-of-distribution detection and selective generation for conditional language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. 2023. [The curious case of hallucinatory \(un\)answerability: Finding truths in the hidden states of over-confident large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3607–3625, Singapore. Association for Computational Linguistics.



- Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024. [Unsupervised real-time hallucination detection based on the internal states of large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14379–14391, Bangkok, Thailand. Association for Computational Linguistics.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, et al. 2023. [Gemini: a family of highly capable multimodal models](#). *arXiv preprint arXiv:2312.11805*.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Daniil Vasilev, Akim Tsvigun, Sergey Petrakov, Rui Xing, Abdelrahman Sadallah, Kirill Grishchenkov, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, Maxim Panov, and Artem Shelmanov. 2025. [Benchmarking uncertainty quantification methods for large language models with LM-polygraph](#). *Transactions of the Association for Computational Linguistics*, 13:220–248.
- Artem Vazhentsev, Ekaterina Fadeeva, Rui Xing, Gleb Kuzmin, Ivan Lazichny, Alexander Panchenko, Preslav Nakov, Timothy Baldwin, Maxim Panov, and Artem Shelmanov. 2025a. [Unconditional truthfulness: Learning unconditional uncertainty of large language models](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, Suzhou, China. Association for Computational Linguistics.
- Artem Vazhentsev, Gleb Kuzmin, Artem Shelmanov, Akim Tsvigun, Evgenii Tsymbalov, Kirill Fedyanin, Maxim Panov, Alexander Panchenko, Gleb Gusev, Mikhail Burtsev, Manvel Avetisian, and Leonid Zhukov. 2022. [Uncertainty estimation of transformer predictions for misclassification detection](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8237–8252, Dublin, Ireland. Association for Computational Linguistics.
- Artem Vazhentsev, Gleb Kuzmin, Akim Tsvigun, Alexander Panchenko, Maxim Panov, Mikhail Burtsev, and Artem Shelmanov. 2023a. [Hybrid uncertainty quantification for selective text classification in ambiguous tasks](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11659–11681, Toronto, Canada. Association for Computational Linguistics.
- Artem Vazhentsev, Lyudmila Rvanova, Ivan Lazichny, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2025b. [Token-level density-based uncertainty quantification methods for eliciting truthfulness of large language models](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2246–2262, Albuquerque, New Mexico. Association for Computational Linguistics.
- Artem Vazhentsev, Akim Tsvigun, Roman Vashurin, Sergey Petrakov, Daniil Vasilev, Maxim Panov, Alexander Panchenko, and Artem Shelmanov. 2023b. [Efficient out-of-domain detection for sequence to sequence models](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1430–1454, Toronto, Canada. Association for Computational Linguistics.
- Yuxia Wang, Daniel Beck, Timothy Baldwin, and Karin Verspoor. 2022. [Uncertainty estimation and reduction of pre-trained models for text regression](#). *Transactions of the Association for Computational Linguistics*, 10:680–696.
- Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. [The art of abstention: Selective prediction and error regularization for natural language processing](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1040–1051, Online. Association for Computational Linguistics.
- Mert Yüsekçönül, Varun Chandrasekaran, Erik Jones, Suriya Gunasekar, Ranjita Naik, Hamid Palangi, Ece Kamar, and Besmira Nushi. 2024. [Attention satisfies: A constraint-satisfaction lens on factual errors of language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
- Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024a. [LUQ: Long-text uncertainty quantification for LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5244–5262, Miami, Florida, USA. Association for Computational Linguistics.
- Caiqi Zhang, Ruihan Yang, Zhisong Zhang, Xinting Huang, Sen Yang, Dong Yu, and Nigel Collier. 2024b. [Atomic calibration of LLMs in long-form generations](#). *Preprint*, arXiv:2410.13246.
- Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. 2023. [Enhancing uncertainty-based hallucination detection with stronger focus](#). In



*Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 915–932.

Xuchao Zhang, Fanglan Chen, Chang-Tien Lu, and Naren Ramakrishnan. 2019. [Mitigating uncertainty in document classification](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3126–3136, Minneapolis, Minnesota. Association for Computational Linguistics.

## A Training Data Generation Pipeline

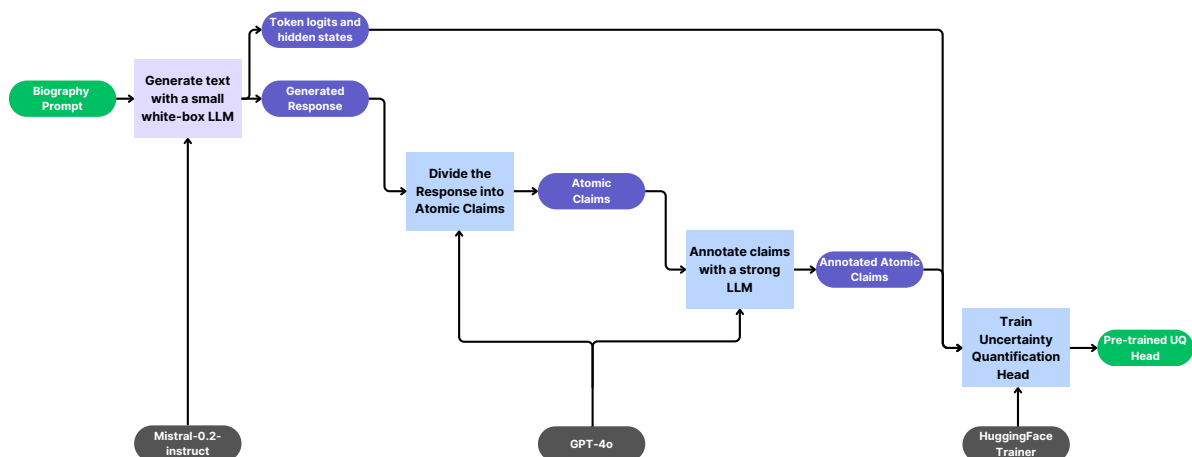


Figure 4: The training data generation pipeline.

Language	Acc.	# claims	% of false claims
English	0.97	97	16.5%
Russian	0.89	275	15.6%
Chinese	0.91	100	35.0%
German	0.83	98	19.4%

Table 4: Performance of GPT-4o annotation pipeline against manual annotation for Mistral-7B-v0.1 model (unsupported claims represent a positive class).

## B Dataset Details

### B.1 Dataset Construction

We used few-shot learning to better guide the LLM to generate the items for the desired domain. The structure of the prompts looks as follows:

```
Continue the list of 100 most famous {domain items}:  
1. <domain-item-1>  
2. <domain-item-2>  
3. <domain-item-3>
```

Example for the “cities” domain:

```
Continue the list of 100 most famous cities:  
1. Paris, France  
2. Amsterdam, Netherlands  
3. Osaka, Japan
```

For claim extraction and their annotation, we use GPT-4o with prompts from (Fadееva et al., 2024). Overall expenses for LLM API calls are approximately \$4000.

### B.2 Dataset Statistics

Model	Dataset	# of texts	# of claims	Claim accuracy, %
Mistral 7b Instruct v0.2	biographies	3,300	68,241	73.7
	multi-domain	700	14,554	86.0
Gemma 2 9b Instruct	biographies	3,300	83,716	88.6

Table 5: Statistics about the training datasets used in our experiments.

Split	# of prompts	ChatGPT prompt used to generate questions	Testing prompt	# of claims		Claim accuracy, %	
				Mistral	Gemma	Mistral	Gemma
persons	100	Tell me a list of 100 most famous persons.	Tell me a bio of a <person>	2,234	2,857	72.9	87.4
cities	100	Tell me a list of 100 most famous cities.	Tell me a history of a <city>	2,128	2,684	79.8	87.1
movies	100	Tell me a list of 100 most famous movies.	Tell me about the movie <movie> and its cast.	2,568	3,121	89.7	94.8
inventions	100	Tell me a list of 100 most important inventions.	Tell me about the invention of <invention> and its inventor.	2,269	2,626	84.3	92.1
books	100	Tell me a list of 100 most famous books.	Tell me about the book <book> and its author.	2,530	3,070	89.9	95.9
artworks	100	Tell me a list of 100 most famous artworks.	Tell me about the artwork <artwork> and its artist.	2,464	2,873	75.9	85.1
landmarks	100	Tell me a list of 100 most famous landmarks.	Tell me about the landmark <landmark>.	2,365	2,566	88.5	93.7
events	100	Tell me a list of 100 most significant historical events.	Tell me about <event> event.	2,294	2,665	88.9	94.8
Russian	100	—	Р а с с к а ж и б и о г р а ф и ю <person>	—	3,572	—	66.7
Chinese	100	—	介绍一下 <person>	—	2,248	—	77.8
German	100	—	Erzählen Sie mir eine Biografie von <person>	—	2,815	—	85.1

Table 6: The statistics of the multi-domain test dataset and number of claims generated by Mistral 7B Instruct v0.2 and Gemma 2 9b Instruct models.

## C Hardware and Computational Efficiency

All experiments were conducted on 8 NVIDIA RTX 5880 Ada GPUs. On average, training a single model with hyperparameter search takes around 150 GPU hours.

Method	Computational Overhead	GPU Memory Footprint
MCP	0.0 %	-
Perplexity	0.0 %	-
Max Token Entropy	0.2 %	-
CCP	8.6 %	1,546 MB
SAPLMA	4.7 %	4 MB
Factoscope	6.1 %	32 MB
Lookback Lens	5.5 %	<1 MB
UHead (only hidden states)	8.7 %	73 MB
UHead (att. + prob. + hs.)	9.9 %	82 MB
UHead (Ours)	4.9 %	40 MB

Table 7: Computational overhead of UQ methods evaluated with the Mistral 7B Instruct v0.2 model. Overhead is measured relative to the fastest method, MCP. For CCP, the size of the auxiliary NLI model is reported. The results were obtained using a multi-domain dataset containing 800 texts and a total of 18,852 claims.



## D Analysis of Attention-Based Features

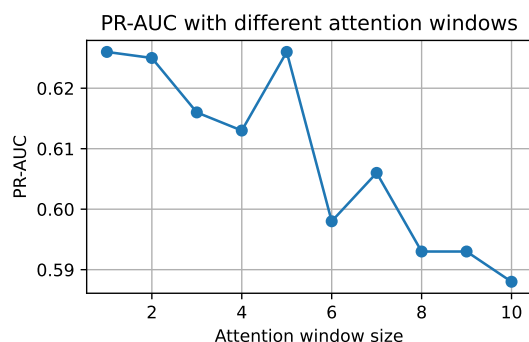


Figure 5: PR-AUC for different attention window sizes using UHead for the Mistral 7B Instruct v0.2 model.

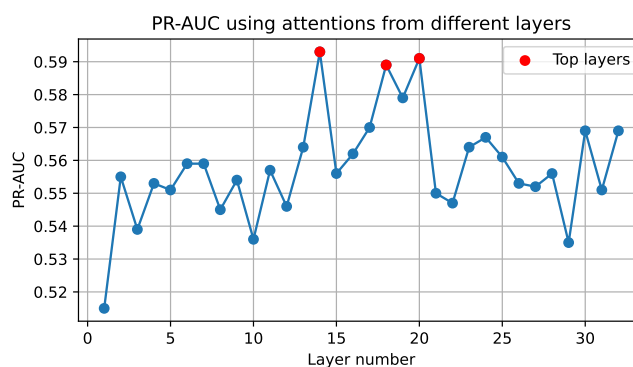


Figure 6: PR-AUC as a function of layer number used for attention features in UHead for the Mistral 7B Instruct v0.2 model. Highlighted points mark layers with highest PR-AUC (layers 14, 18 and 20).

Method	Test Set	Biographies (dev)
UHead, all 1024 heads		.641
UHead, 512 heads		.631
UHead, 256 heads		.638
UHead, 128 heads		<b>.646</b>
UHead, 64 heads		.632
UHead, 32 heads		.614
UHead, 16 heads		.597
UHead, 8 heads		.585
UHead, 4 heads		.492
UHead, 2 heads		.404
UHead, 1 head		.375

Table 8: PR-AUC scores on the *biographies* development set when training UHead using only the top- $N$  (layer, head) pairs ranked by their absolute correlation with training labels. Reducing the number of heads by about eightfold (from 1024 to 128) maintains performance, while further reductions lead to performance degradation.

## E Additional Experimental Results

Method	SciQ	TruthfulQA	MedQUAD
MCP	.259	.405	<b>.561</b>
Perplexity	.242	.408	.549
Max Token Entropy	.273	.379	.551
CCP	.314	.386	.554
UHead	<b>.359</b>	<b>.420</b>	.555

Table 9: PR-AUC for various UQ methods on QA datasets. The results show generalization of UHead trained on *biographies* to the QA task.

Architecture	PR-AUC
Linear	.556
MLP	.626
Transformer (UHead, ours)	<b>.642</b>

Table 10: PR-AUC for different UQ head architectures for the Mistral 7B Instruct v0.2 model on the dev set of *biographies* dataset. The hyperparameters of all detectors are optimized. The results demonstrate the superiority of the transformer architecture.

Method	PR-AUC
MCP	.180
CCP	.307
UHead trained on native dataset (Gemma)	<b>.461</b>
UHead trained on non-native dataset (Mistral)	.435

Table 11: PR-AUC of the hallucination detector for Gemma 2 trained on the “native” data (generated by Gemma 2) in comparison to training on “non-native” data (generated by Mistral). PR-AUC is reported on the test set of English *biographies* dataset. The results show that using “non-native” training data substantially reduces the performance.

Method	Test Sets	Test Sets						
		Cities	Movies	Inventions	Books	Artworks	Landmarks	Events
UHead, bio		.487	.466	.485	.395	.561	.340	.369
UHead, bio + all - 1		.489	.479	.482	.404	.572	.338	.387

Table 12: Introducing more diverse training data. UHead results are shown for two scenarios: when the UQ head is trained solely on the English biographies dataset, and when it is trained on the biographies dataset along with all other domains, excluding the test domain. Adding more data slightly improves the performance in the OOD setting.

Method	Test Sets	Biographies (in domain)	Test Sets						
			Cities	Movies	Inventions	Books	Artworks	Landmarks	Events
Random		.212 ± .012	.207 ± .011	.077 ± .010	.156 ± .011	.080 ± .008	.189 ± .010	.130 ± .010	.092 ± .008
MCP		.354 ± .018	.322 ± .015	.135 ± .009	.254 ± .020	.156 ± .023	.308 ± .017	.185 ± .017	.123 ± .014
Perplexity		.319 ± .016	.269 ± .014	.126 ± .008	.224 ± .016	.115 ± .013	.330 ± .019	.155 ± .013	.102 ± .010
Max Token Entropy		.378 ± .020	.326 ± .015	.182 ± .010	.295 ± .024	.146 ± .019	.368 ± .019	.185 ± .017	.129 ± .014
CCP		.429 ± .018	.440 ± .022	.159 ± .020	.317 ± .022	.128 ± .018	<u>.384</u> ± .011	.220 ± .022	.140 ± .017
SAPLMA		.494 ± .024	.387 ± .020	.211 ± .027	.292 ± .022	.181 ± .025	.364 ± .011	.234 ± .020	.139 ± .014
Factoscope		.499 ± .025	.423 ± .021	.249 ± .026	<b>.393</b> ± .025	.192 ± .023	.382 ± .021	.273 ± .025	.150 ± .016
Lookback lens		.459 ± .022	.441 ± .021	.214 ± .025	.333 ± .020	.199 ± .022	.355 ± .022	<u>.296</u> ± .023	<b>.177</b> ± .019
UHead (Ours)		<b>.553</b> ± .021	<b>.523</b> ± .017	<b>.302</b> ± .008	<u>.379</u> ± .025	<b>.242</b> ± .032	<b>.471</b> ± .024	<b>.319</b> ± .025	<u>.163</u> ± .017

Table 13: Performance comparison of the UQ head using the Llama 3.1 8b Instruct model against various UQ baselines.

## F Hyperparameters

Method	Model	Learning Rate	Num. Epochs	Att. Window Size
SAPLMA	Gemma 2 9b Instruct	1e-4	10	–
	Mistral 7b Instruct v0.2	1e-4	10	–
	Llama 3.1 8B Instruct	1e-4	10	–
Lookback Lens	Gemma 2 9b Instruct	1e-2	13	–
	Mistral 7b Instruct v0.2	1e-2	13	–
	Llama 3.1 8B Instruct	1e-2	13	–
UHead (Factoscope)	Gemma 2 9b Instruct	2e-4	3	–
	Mistral 7b Instruct v0.2	2e-4	5	–
	Llama 3.1 8B Instruct	5e-5	5	–
UHead	Gemma 2 9b Instruct	2e-4	6	2
	Mistral 7b Instruct v0.2	2e-4	7	2
	Llama 3.1 8B Instruct	1e-4	6	5

Table 14: Optimal key hyperparameters for each method and LLM.

For each tested model, we selected hyperparameters by optimizing the PR-AUC metric on the validation set of the “biographies” dataset. In training, we optimized the learning rate, warmup ratio, number of epochs, and the weight of positive examples in the cross-entropy loss. For the model architecture, we optimized the number of uncertainty layers, the number of heads, and the intermediate dimension. For feature extraction, we optimized the number of layers used to obtain hidden states, token probabilities, and attention weights, as well as the number of preceding tokens considered for attention. The optimal key hyperparameters are summarized in Table 14, the full hyperparameter set is available in the repository.<sup>3</sup> The hyperparameter grid is the following:

**Weight of the positive class:** {1, 3, 4, 5}

**Learning rate:** {1e-5, 3e-5, 5e-5, 1e-4, 2e-4, 5e-4, 1e-2};

**Num. of epochs:**  $\{n \in \mathbb{N} \mid 2 \leq n \leq 15\}$ ;

**Warmup:** {0., 0.05, 0.1};

**Attention window size:** {1, 2, 3, 4, 5, 7, 10};

**Dropout rate:** {0., 0.05, 0.1, 0.2};

**Weight decay:** {0, 1e-2, 1e-1};

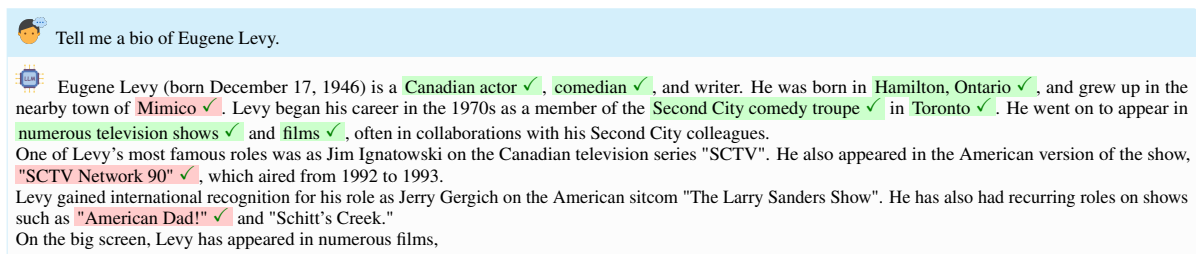
**(For UHead) Number of Layers:** {1, 2};

**(For UHead) Attention head dimension:** {128, 256, 512, 768, 1024};

**(For UHead) Number of attention heads:** {8, 16}.

<sup>3</sup><http://uncertainty-head.nlpresearch.group>

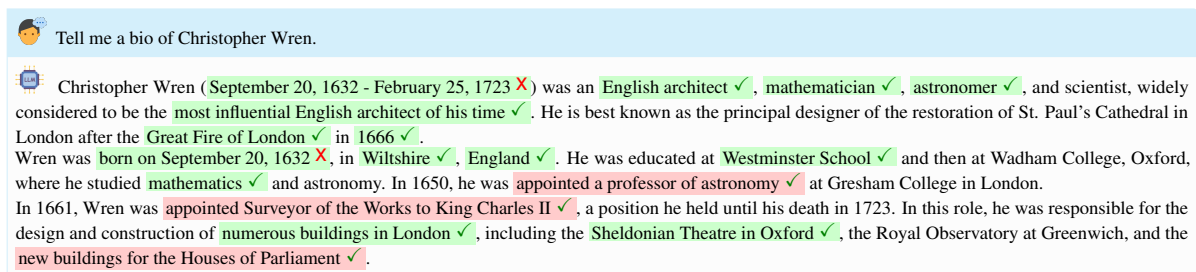
## G Examples of Hallucinations and UQ Head Predictions



Tell me a bio of Eugene Levy.

Eugene Levy (born December 17, 1946) is a Canadian actor ✓, comedian ✓, and writer. He was born in Hamilton, Ontario ✓, and grew up in the nearby town of Mimico ✓. Levy began his career in the 1970s as a member of the Second City comedy troupe ✓ in Toronto ✓. He went on to appear in numerous television shows ✓ and films ✓, often in collaborations with his Second City colleagues. One of Levy's most famous roles was as Jim Ignatowski on the Canadian television series "SCTV". He also appeared in the American version of the show, "SCTV Network 90" ✓, which aired from 1992 to 1993. Levy gained international recognition for his role as Jerry Gergich on the American sitcom "The Larry Sanders Show". He has also had recurring roles on shows such as "American Dad!" ✓ and "Schitt's Creek." On the big screen, Levy has appeared in numerous films,

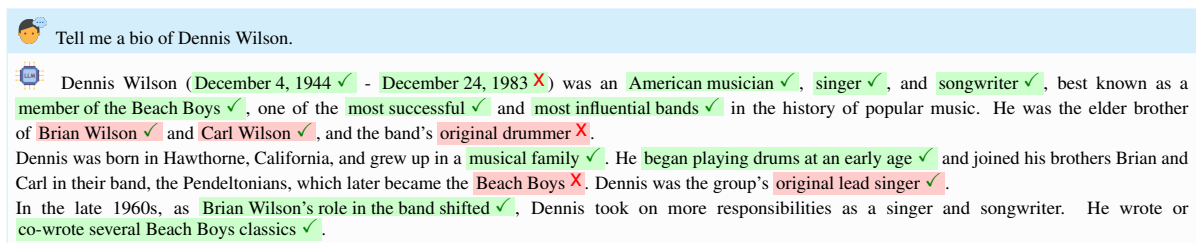
Figure 7: Mistral's response to the question "Tell me a bio of Eugene Levy" and UHead predictions. Claim-relevant phrases are highlighted only when the model shows high confidence (green, uncertainty < 0.2) or high uncertainty (red, uncertainty  $\geq 0.6$ ). Medium-uncertainty phrases (0.2–0.6) are left unmarked. Ticks (✓) indicate confidence matches ground-truth label; crosses (✗) indicate mismatch.



Tell me a bio of Christopher Wren.

Christopher Wren (September 20, 1632 - February 25, 1723 ✗) was an English architect ✓, mathematician ✓, astronomer ✓, and scientist, widely considered to be the most influential English architect of his time ✓. He is best known as the principal designer of the restoration of St. Paul's Cathedral in London after the Great Fire of London ✓ in 1666 ✓. Wren was born on September 20, 1632 ✗, in Wiltshire ✓, England ✓. He was educated at Westminster School ✓ and then at Wadham College, Oxford, where he studied mathematics ✓ and astronomy. In 1650, he was appointed a professor of astronomy ✓ at Gresham College in London. In 1661, Wren was appointed Surveyor of the Works to King Charles II ✓, a position he held until his death in 1723. In this role, he was responsible for the design and construction of numerous buildings in London ✓, including the Sheldonian Theatre in Oxford ✓, the Royal Observatory at Greenwich, and the new buildings for the Houses of Parliament ✓.

Figure 8: Mistral's response to the question "Tell me a bio of Christopher Wren" and UHead prediction. Claim-relevant phrases are highlighted only when the model shows high confidence (green, uncertainty < 0.2) or high uncertainty (red, uncertainty  $\geq 0.6$ ). Medium-uncertainty phrases (0.2–0.6) are left unmarked. Ticks (✓) indicate confidence matches ground-truth label; crosses (✗) indicate mismatch.



Tell me a bio of Dennis Wilson.

Dennis Wilson (December 4, 1944 ✓ - December 24, 1983 ✗) was an American musician ✓, singer ✓, and songwriter ✓, best known as a member of the Beach Boys ✓, one of the most successful ✓ and most influential bands ✓ in the history of popular music. He was the elder brother of Brian Wilson ✓ and Carl Wilson ✓, and the band's original drummer ✗. Dennis was born in Hawthorne, California, and grew up in a musical family ✓. He began playing drums at an early age ✓ and joined his brothers Brian and Carl in their band, the Pendletonians, which later became the Beach Boys ✗. Dennis was the group's original lead singer ✓. In the late 1960s, as Brian Wilson's role in the band shifted ✓, Dennis took on more responsibilities as a singer and songwriter. He wrote or co-wrote several Beach Boys classics ✓.

Figure 9: Mistral's response to the question "Tell me a bio of Dennis Wilson" and UHead predictions. Claim-relevant phrases are highlighted only when the model shows high confidence (green, uncertainty < 0.2) or high uncertainty (red, uncertainty  $\geq 0.6$ ). Medium-uncertainty phrases (0.2–0.6) are left unmarked. Ticks (✓) indicate confidence matches ground-truth label; crosses (✗) indicate mismatch.