# SPILLED ENERGY IN LARGE LANGUAGE MODELS

### Anonymous authors

000

001 002 003

004

005 006 007

008

010

011

012

013

014

015

016

017

018

019

021

023

025

026027028

029

031

032

034

037

040

041

042

043

045

047 048 049

051

052

Paper under double-blind review

#### **ABSTRACT**

We reinterpret the final softmax classifier over the vocabulary of Large Language Models (LLM) as an Energy-based Model (EBM). This allows us to decompose the chain of probabilities used in sequence-to-sequence modeling as multiple EBMs that interact together at inference time. Our decomposition offers a principled approach to measuring where the "energy spills" in LLM decoding, empirically showing that spilled energy correlates well with factual errors, inaccuracies, biases, and failures. Similar to Orgad et al. (2025), we localize the "exact" token associated with the answer, yet, unlike them, who need to train a classifier and ablate which activations to feed to it, we propose a method to detect hallucinations *completely* training-free that naturally generalizes across tasks and LLMs by using the output logits across subsequent generation steps. We propose two ways to detect hallucinations: the first one that measures the difference between two quantities that we call **spilled energy**, measuring the difference between energy values across two generation steps that mathematically should be equal; the other is marginal **energy**, which we can measure at a single step. Unlike prior work, our method is training-free, mathematically principled, and demonstrates strong cross-dataset generalization: we scale our analysis to state-of-the-art LLMs, including LLaMa-3, Mistral, and Qwen-3, evaluating on nine benchmarks and achieving competitive performance with robust results across datasets and different LLMs.

Q/A: ``What is the capital of Italy? Answer:''

Logit

Spilled (Ours)

The capital of Italy is Rome

The capital of Italy is Sydney

X

Reasoning: ''A farmer has 12 chickens. Each chicken lays 2 eggs per day. How many eggs will the farmer collect in 5 days?''

Logit

Spilled (Ours)

12 chickens lay 2 eggs per day . In
5 days , the farmer will collect 12 x
2 x 5 = 120 eggs in 5 days

12 chickens lay 2 eggs per day . In
5 days , the farmer will collect 12 x
2 x 5 = 470 eggs in 5 days

X

12 chickens lay 2 eggs per day . In
5 days , the farmer will collect 12 x
2 x 5 = 470 eggs in 5 days

X

12 chickens lay 2 eggs per day . In
5 days , the farmer will collect 12 x
2 x 5 = 470 eggs in 5 days

X

Figure 1: Color-coded comparison of hallucination detection with LLaMa-3 8B using logit confidence and **our spilled energy**. Our method generalizes well across topics (e.g., Q&A, reasoning) and diverse LLMs. ✓ indicates a correct answer and ✗ an incorrect one. While our approach focuses on the exact answer tokens (e.g. Rome/Sydney and 120/470, see Section 4.2), here we apply min–max normalization to the full answer for visualization, as truthful hallucination.

# 1 Introduction

055 056

057

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076 077

078

079

081

083

084

085

086 087

880

091

092

094

095

096

097

098

100 101

102

103

104

105

106

107

The widespread adoption of Large Language Models (LLMs) across various domains has brought increasing attention to their critical limitation: their tendency to generate incorrect or misleading information—commonly referred to as "hallucinations." This issue supports the idea that LLMs are just stochastic parrots (Bender et al., 2021) answering in a way that is statistically plausible with respect to the input prompt despite not having a real understanding of it. On the other side, recent reasoning capabilities proper to ChatGPT 40 (OpenAI-Team, 2023) or Deepseek (Liu et al., 2024) offer counter evidence to actually support this. Ongoing research seeks to characterize and categorize hallucinations, setting them apart from other error types (Liu et al., 2022; Ji et al., 2023; Huang et al., 2023; Rawte et al., 2023). At the same time, recent discussions have introduced terms such as confabulations (Millidge, 2023) and fabrications (McGowan et al., 2023), sometimes attributing a form of "intention" to LLMs—though the very idea of LLM "intentionality" and other human-like qualities remains contested (Salles et al., 2020; Serapio-García et al., 2023; Harnad, 2024). Research on LLM hallucinations can be categorized into two main branches: the first one is the extrinsic branch, where the hallucinations are measured with respect to the interpretation that humans give to those errors (Bang et al., 2023; Ji et al., 2023; Huang et al., 2023; Rawte et al., 2023). The second branch was started by Kadavath et al. (2022), proposing to study the hallucinations within the model itself. Following Kadavath et al. (2022), the work in Li et al. (2024) proposes Inference-Time Intervention (ITI) as a way to improve the "truthfulness" of LLMs at inference time. ITI functions by altering model activations at inference time, steering them along specific directions within a restricted set of attention heads. Our work is also different from Yin et al. (2023), since we care about detecting errors in LLMs, whereas they introduce an automated methodology to detect when LLMs are aware that they do not know how to answer.

In this work, we follow the definition of hallucinations given by Orgad et al. (2025) as any form of error produced by an LLM—including factual mistakes, biased outputs, breakdowns in common-sense reasoning, and related issues. Like them, we also confirm that the truthfulness signal is concentrated in the "exact answer tokens." Nevertheless, unlike them, we abandon the idea of using a probe classifier (Belinkov, 2022) trained for each task and dataset. Given that LLMs are foundational models, user interactions typically occur *in the wild*, making it difficult to predict which probe classifier is best suited for detecting hallucinations in real-world scenarios. Furthermore, in this setting, classifier weights should not only be updated dynamically for each task, but the optimal token—layer combination is also dataset-dependent, which conflicts with the broad LLM applicability. Indeed, in the work by Orgad et al. (2025), the article reports:

"We find that probing classifiers do not generalize across different tasks."

In our paper, we propose to solve this problem with a training-free method that generalizes better across different tasks and is mathematically principled using the framework of Energy-based Models (EBMs). Fig. 1 reports a qualitative comparison across tasks, comparing to the logit confidence. Additional samples are shown in Appendix D.3.

We reinterpret the final softmax classifier over the vocabulary of LLM as an EBM, taking inspiration from what Grathwohl et al. (2020) did five years ago for classifiers. This perspective enables us to decompose the sequence-to-sequence probability chain into multiple interacting EBMs that operate jointly during inference. Through this decomposition, we introduce the notion of "spilled energy" in LLM decoding and show empirically that such spill strongly correlates with errors. Given that our method is solely based on the mathematics of EBMs and the chain rule of probability, we do not have to train or tune our detector, striking a good generalization across tasks and LLMs. Building on this foundation, our contributions are as follows:

- Training-free, LLM hallucination detection generalizing across tasks using the EBM framework. We introduce a method for detecting hallucinations that requires no additional training, in contrast to prior work that relies on trained classifiers and ablations of model activations. Our approach directly reads values inside the LLM, enabling natural generalization across tasks and performing better than logit-based detection.
- $\diamond$  Two energy-based metrics. We define two complementary measures of energy spills: (i) delta energy  $\Delta E_{\theta}(\mathbf{x}_{i:1})$ , which captures discrepancies between energy values across two time steps that

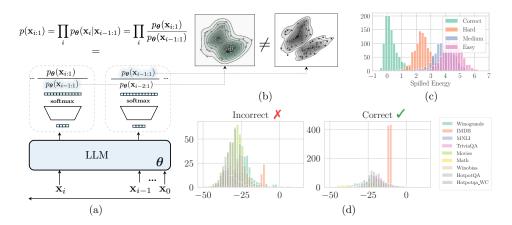


Figure 2: **How energy spills in LLMs**. (a) Language Modeling  $p(\mathbf{x}_{i:1})$  is attained as a decomposition problem following the chain rule of probability, implemented as autoregressive: we recursively apply a discriminative classifier over the vocabulary  $\mathcal{V}$  to attain generative modeling with larger context size i.e.  $p(\mathbf{x}_i|\mathbf{x}_{i-1:1})$ . (b) We reinterpret each discriminative classifier as a generative EBM, finding a connection between two quantities that should be the same across time steps yet are different. We call this difference "the spilled energy"  $\Delta E_{\theta}(\mathbf{x}_{i:1})$  in Eq. (8). (c) Given that we simply read values inside the LLM, our approach is training-free and correlates well with hallucinations on a synthetic math dataset with increasing difficulty; (d) histograms of spilled energy values, for incorrect and correct answers on all nine datasets using min pooling for Llama-3-Instruct. The two distributions are easily separable by using a simple threshold, resulting in a generalization across real-world tasks.

should be mathematically equivalent, and (ii) marginal energy  $E^m_{\theta}(\mathbf{x}_{i:1})$ , which can be evaluated at a single time step.

 Scalable and generalizable analysis. Our framework is mathematically principled, training-free, and exhibits strong cross-dataset generalization. We scale our analysis to state-of-the-art LLMs, including Llama 3-8B-Instruct and Mistral-7B-Instruct, and demonstrate competitive performance across nine benchmarks, showing robustness across datasets and architectures.

Fig. 2(a) illustrates the core idea of our method: rather than using a naïve approach, such as simply recording the logit or training a probe classifier at the activations of the answer token, we first reinterpret the LLM as an autoregressive EBM via the chain rule of probabilities. We then further decompose each conditional probability, incorporating insights from Grathwohl et al. (2020). At the time step of the exact token i-1, we extract the energy, which corresponds to the logit, and compare it with the marginal energy at the next time step i, corresponding to the denominator of the softmax. According to the chain rule, these two quantities should be identical; however, they differ in the LLM implementation—Fig. 2(b). We find that the discrepancy, which we term spilled energy  $\Delta E_{\theta}(\mathbf{x}_{i:1})$ , correlates strongly with instances where the LLM produces an incorrect output—see Fig. 2(c). Moreover, its detection signal separates well correct and incorrect classes across datasets, reflecting the model's confidence, as shown in Fig. 2(d).

### 2 Related Work

EBM applications to Trustworthy AI. EBMs have been applied to improve the reliability and interpretability of Deep Nets. For example, Energy-Based Out-of-Distribution Detection (OOD) (Liu et al., 2020) uses the energy score as a more robust alternative to the softmax confidence. At the same time, Grathwohl et al. (2020) presents how to reinterpret a discriminative classifier as EBM to train models both discriminative and generative. Following this work, Zhu et al. (2021) gives new insights into the role of energy when training EBMs and robust classifiers using adversarial training. Instead, Mirza et al. (2024; 2025) explain adversarial attacks by reinterpreting the softmax classifier as an EBM, showing that these perturbations correspond to shifts in the underlying energy landscape.

Foundations of Hallucination in LLMs. LLMs are prone to diverse errors—including bias, reasoning failures, and generation of factually incorrect information unsupported by reliable sources. Karpowicz (2025) frames hallucination and imagination as mathematically identical phenomena, both emerging from a necessary violation of information conservation. Also Xu et al. (2025) provides a formal learning-theoretic proof that hallucinations are unavoidable. They define a *formal world* in which both the LLM and the ground-truth are computable functions, showing through classic results in computability theory, that no LLM can learn all such functions. As a consequence, hallucination is not just a practical artifact but a fundamental limitation of LLMs, valid even under idealized conditions. Recently Kalai et al. (2025) showed that hallucinations come from the statistical problem of the pretraining methodology: minimizing the cross entropy naturally causes errors because it does not train the model to express uncertainty and say "I do not know." Kalai et al. (2025) proposes to change the evaluation practices to not reward models for guessing, but rather to mimic the human exams that penalize only wrong answers.

Detecting and Mitigating LLM Hallucinations. Orgad et al. (2025) train classifiers on the internal representations of the LLMs to predict, based on the features, the correctness of the answer. Given an LLM in a white-box setting, an input prompt, and the generated response  $\hat{y}$ , the classifier's task is to predict whether  $\hat{y}$  is a hallucination. Orgad et al. suggested that LLMs may encode more factual knowledge in their latent subspaces than is revealed in their outputs. Gekhman et al. (2025) propposed a framework for studying hidden knowledge. Finally, Santilli et al. (2025) point out that uncertainty quantification in language models is often evaluated using metrics like AuROC. This shares biases between detection methods and correctness functions (e.g., length effects) that systematically distort results. One way to mitigate hallucinations is to act at the decoding stage, where the output generation can be steered Subramani et al. (2022). Steering vectors provide a straightforward way to control a model by adding a fixed vector to its activations (Dunefsky & Cohan, 2025). Fu et al. (2025) introduced DeepConf, a test-time method that leverages model-internal confidence signals to filter out low-quality reasoning traces during or after generation. Kuhn et al. (2023b); Fadeeva et al. (2024); Farquhar et al. (2024), and its follow-up by Kossen et al. (2025) in which they approximate the semantic entropy in a more efficient way. Constrained decoding approaches Li et al. (2023); Peng et al. (2023) modify token selection policies. Similarly, reinforcement learning with fact-based rewards Ouyang et al. (2022) has been used to bias decoding trajectories toward verifiable outcomes. Incorrect answers may also be given due to an ambiguous prompt: Kuhn et al. (2023a)'s CLAM framework uses few-shot prompts to classify a question's ambiguity and then asks the user to clarify.

# 3 BACKGROUND AND PRELIMINARIES

#### 3.1 ENERGY-BASED MODELS

We give an overview of Energy-based Models (EBMs) and their use in discriminative classifiers.

**EBMs.** Energy-Based Models are a class of probabilistic models in which the probability distribution over data points  $\mathbf{x}$  is defined in terms of an energy function  $E_{\theta}(\mathbf{x})$ . The energy function, parameterized by a neural network  $\theta$  (Lecun et al., 2006), assigns a scalar energy to each configuration of  $\mathbf{x}$ , where lower energy values correspond to higher likelihood. The resulting probability distribution is given by  $p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z_{\theta}}$  where  $Z_{\theta}$  denotes the partition function (normalizing constant), defined as  $Z_{\theta} = \sum_{\mathbf{x}} \exp(-E_{\theta}(\mathbf{x}))$  for discrete  $\mathbf{x}$ , or equivalently  $Z_{\theta} = \int \exp(-E_{\theta}(\mathbf{x})) d\mathbf{x}$  for continuous  $\mathbf{x}$ . Standard neural networks are often deterministic function approximators, mapping  $\mathbf{x} \mapsto y$ , EBMs instead define a full probability distribution over data or latent variables.

One of the strengths of EBMs is their flexibility in modeling arbitrary distributions without being tied to a specific parametric form. This flexibility comes from the fact that the energy function  $E(\mathbf{x})$  can be defined in various ways. Training involves learning the parameters of the energy function such that the probability distribution  $p_{\theta}(\mathbf{x})$  matches the empirical distribution of the data. This is typically done using techniques like contrastive divergence, score matching, or maximum likelihood.

**Notation.** Let  $\mathcal{V}$  denote the vocabulary of the LLM, i.e., the set of all tokens that can be processed as input and generated at each decoding step, with size  $|\mathcal{V}| = V$ . We shorten the sequence of tokens  $\{\mathbf{x}_N, \dots, \mathbf{x}_1\}$  as  $\mathcal{X} = \{\mathbf{x}_{N:1}\}$ , and  $\mathbf{x}_i \in \mathcal{V}$  denotes the token in the *i*-th position along the sequence. We model the LLM as a function  $\boldsymbol{\theta} : \mathbb{R}^{N \times V} \to \mathbb{R}^V$ , implemented by a transformer, or any other sequence-to-sequence mechanism. For a sequence  $\{\mathbf{x}_{i:1}\}$  as input, we write  $\boldsymbol{\theta}(\mathbf{x}_{i:1})[k]$  to denote the

predicted logit assigned to the k-th token class in  $\mathcal{V}$  for the i+1 token in the sequence, as is standard in autoregressive LLM training (Ouyang et al., 2022).

#### 3.2 AUTOREGRESSIVE LARGE LANGUAGE MODELS

Generative modeling has been pursued through a variety of approaches beyond autoregression (AR). Variational Autoencoders (VAEs) (Kingma & Welling, 2014) learn a probabilistic latent variable model by encoding inputs into a latent space and decoding samples back to the data domain. Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) frame generation as a min-max game between a generator and a discriminator. The diffusion process has been incorporated into neural nets (Sohl-Dickstein et al., 2015) and, more recently, Diffusion Models (Ho et al., 2020) have emerged as a powerful class of generative models. While these paradigms differ in how they approximate the data distribution, AR models are special in their kind and take a more direct route by factorizing the joint probability of sequences into conditionals, making them especially suitable for language modeling. We now focus on the AR formulation that underlies most LLMs. Textual data is segmented into a sequence of tokens  $\mathcal{X} = \{\mathbf{x}_i, \dots, \mathbf{x}_1\}$ , and a language modeling objective is employed to maximize the likelihood of such data (Radford & Narasimhan, 2018). In other words, we model the joint probability of tokens in the sequence  $\mathcal{X}$ , through a conditional probability parameterized by  $\boldsymbol{\theta}$ :

$$p(\mathbf{x}_{i:1}) = p(\mathbf{x}_i \mid \mathbf{x}_{i-1:1}) \dots p(\mathbf{x}_2 \mid \mathbf{x}_1) \ p(\mathbf{x}_1) = \prod_i \underbrace{p_{\boldsymbol{\theta}}(\mathbf{x}_i \mid \mathbf{x}_{i-1:1})}_{\text{discriminative model}} p_{\boldsymbol{\theta}}(\mathbf{x}_1). \tag{1}$$

What we find interesting about this factorization is that, although it seeks to attain *generative modeling*, i.e.,  $p(\mathbf{x}_{i:1})$ , it actually uses recursively *discriminative classifiers*, parametrized by a transformer network  $\boldsymbol{\theta}$ , that predicts a discrete distribution of the next token  $\mathbf{x}_i$  over the vocabulary  $\mathcal{V}$ , given previous tokens  $\mathbf{x}_{i-1:1}$ . This is used to model each conditional probability.

# 4 HOW ENERGY SPILLS IN LLMS

When predicting the token at position i, the conditional probability modeled by  $\theta$  can be decomposed using the probabilities of the sequences. As a result, the marginal term from step i cancels out with the sequence probability from the decomposition at the previous step i-1, which means we have:

$$p(\mathbf{x}_{i:1}) = \prod_{i} p_{\boldsymbol{\theta}}(\mathbf{x}_{i} | \mathbf{x}_{i-1:1}) = \prod_{i} \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{i:1})}{p_{\boldsymbol{\theta}}(\mathbf{x}_{i-1:1})} \Longrightarrow \dots \underbrace{\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{i:1})}{p_{\boldsymbol{\theta}}(\mathbf{x}_{i-1:1})}}_{\text{step } i} \underbrace{\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{i-1:1})}{p_{\boldsymbol{\theta}}(\mathbf{x}_{i-2:1})}}_{\text{step } i} \dots = p(\mathbf{x}_{i:1}).$$
(2)

This indeed confirms that Eq. (1) results in the correct formulation for language modeling, which is  $p(\mathbf{x}_{i:1})$ . Following the mathematics, these quantities should cancel out along the sequence, but we will now show that, in practice, this constraint is not explicitly optimized for, and we can exploit it for hallucination detection.

# 4.1 Interpreting LLMs as Energy-based models (EBMs)

Let us continue the expansion from Eq. (2). Writing the conditional as the ratio between the joint distribution in the numerator and the marginal distribution in the denominator, we note that this ratio is actually implemented in LLMs as a softmax classifier that digests the embedding of the prior sentence  $\mathbf{x}_{i-1:1}$  and predicts the next token  $\mathbf{x}_i$ , thus this chain of equality holds true. We can then apply the "trick" from Grathwohl et al. (2020) as:

$$p_{\boldsymbol{\theta}}(\mathbf{x}_{i}|\mathbf{x}_{i-1:1}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{i:1})}{p_{\boldsymbol{\theta}}(\mathbf{x}_{i-1:1})} = \frac{\exp \boldsymbol{\theta}(\mathbf{x}_{i-1:1}) \left[ \mathrm{id}(\mathbf{x}_{i}) \right]}{\sum\limits_{k=1}^{V} \exp \boldsymbol{\theta}(\mathbf{x}_{i-1:1}) [k]} \text{ where id} : \{0,1\}^{V} \mapsto [1,\dots,V]. \quad (3)$$

id is the map that takes as input a one-hot encoding vector  $\mathbf{x}_i$  for a word token at position i in the text and outputs its index in the vocabulary. A typical cross-entropy loss only optimizes with the

supervision provided by the ground-truth token, through the vocabulary index  $id(\mathbf{x}_i)$ . This loss ignores all other quantities or constraints related to the complete sequence  $\mathcal{X}$ , i.e., ignores all the time steps higher than i+1.

We can write the conditional probability of Eq. (3) as a ratio of two EBMs as:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{i}|\mathbf{x}_{i-1:1}) = \log \frac{\exp(-E_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}_{i:1}))}{\exp(-E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1}))} \frac{\widetilde{Z}(\boldsymbol{\theta})}{Z(\boldsymbol{\theta})} = -E_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}_{i:1}) + E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1}). \tag{4}$$

Following Zhu et al. (2021), the partition functions simplify since  $\log \widetilde{Z}(\theta) = \log Z(\theta)^1$ .

 $E^\ell_{m{ heta}}, \ E^m_{m{ heta}}$  are computed from the output of the model, but with two big differences:  $E^\ell_{m{ heta}}$  as a single logit extracted using the id of the sampled token,  $E^m_{m{ heta}}$  by marginalizing over all ids in the vocabulary.

The two energies can be derived from the softmax of the logits, by connecting Eq. (4) and Eq. (3):

$$-\log p_{\boldsymbol{\theta}}(\mathbf{x}_i \mid \mathbf{x}_{i-1:1}) = -\log \left( \frac{\exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\operatorname{id}(\mathbf{x}_i)])}{\sum\limits_{k} \exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k])} \right) =$$
 (5)

$$= \underbrace{-\boldsymbol{\theta}(\mathbf{x}_{i-1:1})\left[\mathrm{id}(\mathbf{x}_{i})\right]}_{E_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}_{i:1})} + \underbrace{\log \sum_{k=1}^{V} \exp \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]}_{-E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1})}$$
(6)

where  $\theta(\mathbf{x}_{i-1:1})$  produces the logits over the entire vocabulary  $\mathcal{V}$ , and  $id(\mathbf{x}_i)$  allows us to extract the logit of the sampled token at decoding step i.

We can think of  $E^{\ell}_{\theta}(\mathbf{x}_{i:1})$  as the energy of the sampled tokens  $\{\mathbf{x}_{i:1}\}$ , and  $E^{m}_{\theta}(\mathbf{x}_{i-1:1})$  as the energy  $E_{\theta}(\mathbf{x}_{i:1})$ , marginalized over all possible  $\mathbf{x}_{i}$ . Considering the decoding at step i in Eq. (4), we get:

$$E_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}_{i:1}) = -\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\mathrm{id}(\mathbf{x}_i)], \quad E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1}) = -\log \sum_{k=1}^{V} \exp \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]. \tag{7}$$

Using the chain rule and Eq. (6), we can write the negative log-likelihood in terms of energies as:

$$-\log p(\mathbf{x}_{N:1}) = -\log \prod_{i} p_{\theta}(\mathbf{x}_{i}|\mathbf{x}_{i-1:1}) = \sum_{i} E_{\theta}^{\ell}(\mathbf{x}_{i:1}) - E_{\theta}^{m}(\mathbf{x}_{i-1:1})$$

without considering the base case  $p_{\theta}(\mathbf{x}_1)$ . Now, if we develop the above equation as done for Eq. (2), we write the total energy of a sequence of length N as  $E_{\theta}(\mathbf{x}_{N:1})$ . Observe that the two energies, not interacting at the same step but at steps i and i-1, should be equal, but they are measured in the LLM at different generation steps and from different components.

$$E_{\theta}(\mathbf{x}_{N:1}) = \sum_{i=1}^{N-1} E_{\theta}^{\ell}(\mathbf{x}_{i+1:1}) - E_{\theta}^{m}(\mathbf{x}_{i:1}) = \dots \underbrace{E_{\theta}^{\ell}(\mathbf{x}_{i+1:1})}_{\Delta E_{\theta}(\mathbf{x}_{i:1})} \underbrace{-E_{\theta}^{m}(\mathbf{x}_{i:1}) + E_{\theta}^{\ell}(\mathbf{x}_{i:1})}_{\Delta E_{\theta}(\mathbf{x}_{i:1})} - E_{\theta}^{m}(\mathbf{x}_{i-1:1}) \dots$$

At timestep i+1, first  $-E^m_{\theta}(\mathbf{x}_{i:1})$  is measured, taking the denominator in the softmax as in the right part of Eq. (6), whereas at timestep i, the second  $E^{\ell}_{\theta}(\mathbf{x}_{i:1})$  is taken, reading the logit in the softmax, left part of Eq. (6). We thus define the discrepancy between the two quantities as **spilled energy**:

**Definition 4.1** (Spilled Energy  $\Delta E_{\theta}(\mathbf{x}_{i:1})$ ). The spilled energy in an LLM is the difference between two energies that, in principle, should be equal, but given that they are measured i) at different time steps ii) in different components, could be different.

$$\Delta E_{\theta}(\mathbf{x}_{i:1}) \triangleq -E_{\theta}^{m}(\mathbf{x}_{i:1}) + E_{\theta}^{\ell}(\mathbf{x}_{i:1}) = \underbrace{-\log \sum_{k} \exp(\theta(\mathbf{x}_{i:1})[k])}_{\text{timestep } i+1} + \underbrace{\theta(\mathbf{x}_{i-1:1})[\text{id}(\mathbf{x}_{i})]}_{\text{timestep } i}$$
(8)

Since both terms on the right side should be equal to  $E_{\theta}(\mathbf{x}_{i:1})$ , delta values should always be zero when we are correctly modeling the energy at timestep i.

<sup>&</sup>lt;sup>1</sup>For a formal proof, please see Appendix A.1.

### 4.2 DETECTING HALLUCINATIONS WITH SPILLED ENERGY

EBMs have previously been used to assess neural network credibility (Liu et al., 2020), and calibration for LLMs has been explored by the Anthropic team (Kadavath et al., 2022). However, dominant training-free baselines such as logits or "p(true)" remain weak. We likewise adopt a training-free approach, but rely on Eq. (8) and its variants as discriminants.

We feed the prompt  $\{\mathbf{x}_{i-1},\ldots,\mathbf{x}_1\}$  to the LLM  $\boldsymbol{\theta}$  and obtain the completion  $\{\mathbf{x}_N,\ldots,\mathbf{x}_i\}$ . Following Orgad et al. (2025), we focus on the "exact answer" tokens—those in [i+1,N] that contain the precise answer (e.g., Rome in Fig. 1), denoted  $[u,w]\subseteq [i+1,N]$ . For instance, it would be the tokens associated with Rome in the question in Fig. 1. We identify this span by prompting the LLM for a brief answer. When the answer spans multiple tokens, we apply a pooling strategy, which we ablate in Section 5. We propose measuring two values that correlate well with hallucinations:

- 1. the marginal energy  $E_{\theta}^{m}(\mathbf{x}_{i:1})$ ;
- 2. the spilled energy  $\Delta E_{\theta}(\mathbf{x}_{i:1})$  by definition of Eq. (8).

We also attempt to combine the two metrics into scaled spilled energy  $\Delta E_s$ , where the spilled energy is multiplied by the absolute value of the marginal energy as  $\Delta E_s(\mathbf{x}_{i:1}) = |E^m_{\theta}(\mathbf{x}_{i:1})|\Delta E_{\theta}(\mathbf{x}_{i:1})$ . The metrics proposed here are independent, new for LLMs, and can all be tested efficiently. These measures can be computed over the full sequence, but for error detection, as discussed in Section 5.2, we must extract the values in the localized exact interval [u, w] to avoid false positives. Note that  $E^{\ell}_{\theta}(\mathbf{x}_{i:1})$  is the classic baseline which in literature is referred to as "logits" or "logits confidence".

### 5 EXPERIMENTS

To evaluate spilled energy, we consider two complementary settings. First, a controlled synthetic environment, where we generate both correct and incorrect multi-digit arithmetic solutions. Second, established real-world benchmarks, where errors arise naturally across diverse reasoning and comprehension tasks. Together, these experiments test whether insights from the clean synthetic setup transfer to the complexity of open-domain language understanding.

# 5.1 SPILLED ENERGY UNDER SYNTHETIC ARITHMETIC

**Experimental Setting.** We first evaluate spilled energy in a controlled setting: multi-digit arithmetic problems with more than 14 digits. For each instance, we generate both correct and incorrect solutions. We tested three different LLMs: Llama-3 8B (Dubey et al.), Qwen-3 8B (Qwen-Team), and Mistral-7B-Instruct v0.3 (Jiang et al.). Incorrect solutions are obtained by introducing random numerical errors of varying magnitude. Specifically, we define three error ranges that differ in their difficulty of detection:

- ♦ **Easy**: random offset in the range [1000, 10000], which are typically easier to identify.
- $\diamond$  **Medium**: random offset in the range [100, 1000], where detection requires closer inspection.
- $\diamond$  **Hard**: random offset in [1, 10], much harder to detect since they appear plausible at first glance.

This design allows us to systematically probe whether spilled energy can distinguish between correct and incorrect generations across different levels of error subtlety.

**Results.** We observe that spilled energy values separate correct from incorrect solutions with high reliability across all error ranges and across all LLMs. In particular, spilled energy consistently assigns lower values to correct generations and higher values to incorrect ones, producing a clear margin of separation. Compared to standard baselines such as *logits*, spilled energy achieves superior discriminative power, especially for errors in the more challenging range [1, 10], see Fig. 3. We offer more results in Fig. 4. Larger, better-detailed ROC and histograms are in Figs. 5 and 6 respectively.

### 5.2 Cross-dataset Results in Real-World Benchmarks

Experimental Setting. We evaluate our methods on a diverse set of established NLP benchmarks, including Math (Hendrycks et al.), TriviaQA (Joshi et al.), HotpotQA (Yang et al.), Winogrande

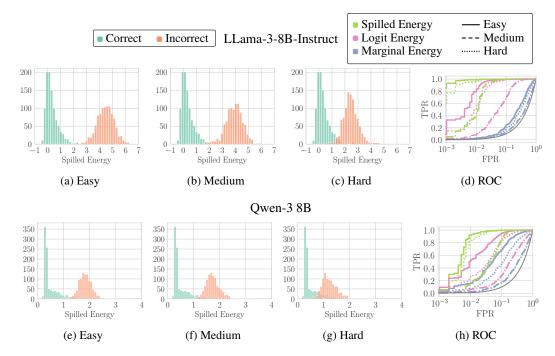


Figure 3: Histograms of Spilled Energy values across models (rows) on Math Sums with different error ranges in the answer (columns, decreasing range left to right, making it harder to detect errors). All sums are performed on 13-digit integers. In the fourth column, we show ROC curves for Hallucination Detection across the error ranges (colors) and methods (line styles).

(Sakaguchi et al.), Winobias (Zhao et al.), Movies (Tapaswi et al.), MNLI (Williams et al.) and IMDB (Maas et al.). These datasets span a wide range of reasoning and error-detection tasks, allowing us to test whether the patterns observed in the synthetic arithmetic setting extend to real-world, open-domain scenarios. Here too, we evaluate multiple LLMs that are either instruction-aligned or not aligned, such as LLaMA-3-Instruct (Dubey et al.), Mistral-Instruct, and Mistral (Jiang et al.). As emphasized by Orgad et al., it is essential to first localize the tokens most relevant to the final answer before applying error detection. Following this principle, we restrict our analysis to the identified exact answer tokens and compare spilled energy against baselines such as the probing classifiers of Orgad et al. and *logit* confidence. Since exact answer tokens may consist of multiple tokens, we further adopt a pooling strategy across the localized span to obtain a final score per sentence.

**Ablation of the exact answer token.** We provide an ablation experiment on the impact of selecting the exact answer tokens. Table 2 reports average AuROC over 9 benchmarks and 3 LLMs with the exact answer, and then another column that offers the improvement provided by using the exact answer. Like prior work, we confirm that searching for the exact answer provides a notable boost: the improvement is very pronounced ( $\sim 25\%$ ) for spilled and marginal energy, while the logit baseline receives a modest increase of 9%.

Cross-dataset results. We next evaluate in the more general setting of cross-dataset transfer, which better reflects real-world usage. For methods requiring training, we report the average performance on each dataset when trained on the remaining datasets (e.g., performance on IMDB is the average of classifiers trained on each of the other nine datasets). Table 1 summarizes results across nine benchmarks. Spilled energy consistently outperforms *logit* confidence, and substantially surpasses the probing classifiers of Orgad et al. (2025). While probing classifiers perform well when trained and tested on the same dataset, their performance drops sharply under cross-dataset evaluation, as reflected in their higher standard deviations. By contrast, spilled energy requires no training and generalizes robustly across diverse benchmarks, providing a lightweight and broadly applicable solution for error detection in LLMs. We also observe that instruction-tuned models tend to amplify the margin by which spilled energy outperforms other methods, whereas on non-aligned models such as Mistral, spilled ranks slightly behind marginal energy. Although this suggests a possible

448 449

450

452 453

454

455

456

457

458

459

460

461

462 463 464

465 466 467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

	Pool	HotpotQA	HotpotQA-WC	IMDB	Math	MNLI	Movies	TriviaQA	Winobias	Winogrande	Average (%)
	LLaMA-3-Instruct (Dubey et al., 2024)										
Logit $E^{\ell}$		72.85±2.12	91.11±1.52		57.81±3.82						54.62±2.60
Orgad et al.	Mean	66.56±9.10	59.00±8.14	<b>69.</b> /8±14.76	66.56±17.04	60.56±12.53	<b>00.44</b> ±8.06	63.22±11.11	67.33±11.97	<b>58.00</b> ±7.79	64.16±11.17
Spilled $\Delta E_s$	Max	$53.65{\scriptstyle\pm1.40}$	$36.28 \pm 2.99$	$55.80{\scriptstyle\pm4.32}$	$35.44 \pm 3.41$	$58.81{\scriptstyle\pm2.58}$	$70.30{\scriptstyle\pm1.49}$	$48.70{\scriptstyle\pm2.44}$	$36.53{\scriptstyle\pm2.98}$	$44.32 \pm 1.70$	$48.87 \pm 2.59$
Marginal $E^m$	Max	$76.72 \pm 1.38$	$30.74 \pm 3.45$	$85.63 \pm 2.39$	$27.08{\scriptstyle\pm5.06}$	$89.90 \pm 1.25$	$96.17 \pm 0.63$	$80.13_{\pm 1.87}$	$57.67 \pm 2.94$	$47.47 \pm 1.83$	$65.72 \pm 2.31$
Marginal $E^m$	Min	$75.91 \pm 1.62$	$97.57 \pm 0.75$	$14.37 \pm 2.39$	$70.55 \pm 2.43$	$61.21 \pm 3.24$	$72.21 \pm 1.60$	$73.38 \pm 1.86$	$47.19_{\pm 2.71}$	$53.98 \pm 2.30$	$62.93 \pm 2.10$
Spilled $\Delta E$	Min	$85.98 \pm 1.09$	$93.00{\scriptstyle\pm1.61}$	$47.66{\scriptstyle\pm4.06}$	$65.58{\scriptstyle\pm3.02}$	$73.95{\scriptstyle\pm1.97}$	$89.34_{\pm 1.04}$	$\pmb{87.07} \scriptstyle{\pm 1.33}$	$60.72{\scriptstyle\pm2.74}$	$55.11_{\pm 2.05}$	<b>73.16</b> ±2.10
		Mistral-Instruct (Jiang et al., 2023)									
Logit $E^{\ell}$	Max	77.24±1.66	83.84±1.66	22.28±2.54	57.67±3.29	78.98±1.58	76.89±1.49	80.35±1.88	45.53±2.60	48.17±1.97	63.44±2.07
Orgad et al.	Mean	$64.78 \scriptstyle{\pm 10.56}$	$56.78{\scriptstyle\pm7.95}$	$\pmb{82.67} \scriptstyle{\pm 11.63}$	$\textbf{68.78} \scriptstyle{\pm 11.43}$	$64.22{\scriptstyle\pm12.12}$	$64.89{\scriptstyle\pm11.55}$	$65.44{\scriptstyle\pm12.10}$	$\textbf{61.00} \scriptstyle{\pm 12.23}$	$\textbf{61.44} \scriptstyle{\pm 11.31}$	$65.56 \scriptstyle{\pm 11.21}$
Spilled $\Delta E_s$	Max	49.13±2.50	36.37±2.40	$46.45{\scriptstyle\pm2.56}$	$29.05{\scriptstyle\pm2.57}$	$53.79{\scriptstyle\pm1.55}$	55.24±2.17	46.73±1.98	53.30±3.66	51.20±1.84	46.81±2.36
Marginal $E^m$	Min	$87.58 \pm 1.35$	$97.94 \pm 0.62$	$18.67{\scriptstyle\pm2.27}$	$67.58 \pm 3.37$	<b>97.96</b> ±0.55	$84.90 \pm 1.37$	$87.75 \pm 1.73$	$49.19 \pm 3.97$	$48.49 \pm 1.86$	71.12±1.90
Marginal $E^m$	Max	$64.63 \pm 1.97$	$33.42 \pm 1.90$	$81.33 \pm 2.32$	$26.52{\scriptstyle\pm2.28}$	$17.62 \pm 1.20$	$86.60_{\pm 1.20}$	$65.46 \pm 2.25$	56.41±4.44	$51.14 \pm 1.71$	$53.68 \pm 2.14$
Spilled $\Delta E$	Min	$91.12 \pm 1.10$	$97.47 \pm 0.78$	$59.77 \scriptstyle{\pm 2.57}$	$66.63{\scriptstyle\pm3.46}$	$95.95{\scriptstyle\pm0.83}$	<b>94.99</b> ±0.93	$91.75 {\scriptstyle\pm1.01}$	$50.74{\scriptstyle\pm3.15}$	$49.00{\scriptstyle\pm1.92}$	<b>77.49</b> ±1.75
		Mistral (Jiang et al., 2023)									
Logit $E^{\ell}$	Max	49.54±1.42	52.47±1.61	32.72±2.89	57.21±3.89	92.49±1.15	30.52±2.00	39.73±2.03	46.53±3.80	44.41±2.42	49.51±2.36
Orgad et al.	Mean	$61.78 \scriptstyle{\pm 9.27}$	$57.44 \pm 6.95$	$\textbf{76.22} \scriptstyle{\pm 12.82}$	$65.78{\scriptstyle\pm15.27}$	$56.67{\scriptstyle\pm11.83}$	$64.22{\scriptstyle\pm8.91}$	$64.33{\scriptstyle\pm10.40}$	$\textbf{58.00} \scriptstyle{\pm 12.29}$	$54.56 \pm 4.36$	$62.11{\scriptstyle\pm10.23}$
Spilled $\Delta E_s$	Max	60.54±1.81	60.18±1.84	43.47±2.76	$71.93{\scriptstyle\pm3.62}$	45.94±2.40	78.84±1.53	$67.92{\scriptstyle\pm1.32}$	57.24±3.72	51.88±1.90	59.77±2.32
Marginal $E^m$	Min	$87.52 \pm 1.31$	$90.91 \pm 1.58$	$54.69{\scriptstyle\pm2.49}$	$86.21 \pm 1.96$	$98.80 \scriptstyle{\pm 0.35}$	$94.41_{\pm 0.62}$	83.66±2.16	$52.15 \pm 1.74$	$46.37{\scriptstyle\pm2.02}$	<b>77.19</b> ±1.58
Marginal $E^m$	Max	$83.57{\scriptstyle\pm1.13}$	$86.83_{\pm 1.70}$	$45.31{\scriptstyle\pm2.49}$	$62.26{\scriptstyle\pm4.29}$	$96.03{\scriptstyle\pm0.83}$	99.27±0.24	$92.26 \pm 1.31$	$51.31 \pm 3.35$	$54.49_{\pm 2.48}$	74.59±1.98
Spilled $\Delta E$	Min	$\textbf{84.24} \scriptstyle{\pm 1.18}$	$83.74{\scriptstyle\pm1.41}$	$57.43_{\pm 2.99}$	$\textit{78.26}\scriptstyle{\pm 2.93}$	$\textbf{96.69} \scriptstyle{\pm 0.62}$	$84.47 \scriptstyle{\pm 1.17}$	$81.27{\scriptstyle\pm1.83}$	$50.62{\scriptstyle\pm1.72}$	$48.72{\scriptstyle\pm1.75}$	$\textbf{73.94} \scriptstyle{\pm 1.74}$

Table 1: Hallucination detection performance, in terms of AuROC, across nine benchmarks and different LLMs. We measure the generalization across all tasks by computing the average.

interaction with alignment, we cannot claim a definitive link. Finally, we compare pooling strategies and find that min pooling yields the best overall performance across methods.

**Limitations.** A current limitation of spilled energy is that it sometimes produces false positives on tokens that are not semantically informative, as shown in Appendix D.3. We observe this effect most prominently on punctuation tokens (e.g., commas, periods) and on words at the beginning of sentences. In these cases, the probability mass over the next token is naturally spread across many plausible options, leading to inflated spilled energy values even in otherwise correct generations. This highlights the importance of accurately identifying the exact answer tokens, as detection is most reliable when restricted to the parts of the output that carry the semantic content of the answer.

### CONCLUSION

In this work, we reinterpreted the softmax layer of LLMs as an EBM, which allowed us to formalize the notion of spilled energy: the discrepancy between two energy values that should, in principle, be equal across consecutive time steps. We showed both theoretically and empirically that this discrepancy provides a powerful, training-free signal for detecting hallucinations and errors in LLM outputs. Through controlled synthetic arithmetic experiments, we demonstrated that spilled energy separates correct from incorrect generations with high reliability, outperforming standard baselines such as logits and marginal energy, even when errors are subtle and difficult to detect. Extending the analysis to a wide range of real-world NLP benchmarks, we found that spilled energy generalizes robustly across tasks and datasets, achieving superior performance without requiring additional classifiers or task-specific training.

	Pool	Average %	Exact
		w/ exact	answer
		answer	increase
Logit $E^{\ell}$	Max	55.86±02.34	+8.97
Orgad et al.	Mean	$63.94{\scriptstyle\pm10.87}$	_
Spilled $\Delta E_s$	Max	$51.82 \scriptstyle{\pm 02.42}$	+0.43
Marginal $E^m$	Min	$70.41 \pm 01.86$	+23.2
Marginal $E^m$	Max	$64.67{\scriptstyle\pm02.14}$	+4.95
Spilled $\Delta E$	Min	$\textbf{74.86} {\scriptstyle\pm01.86}$	+25.6

Table 2: Improvement in the AuROC provided by searching for the exact answer. Average across 3 LLMs and 9 benchmarks.

Compared to probing approaches, which struggle with cross-dataset transfer, our method provides a lightweight, scalable, and broadly applicable alternative. Overall, our findings highlight spilled energy as a mathematically principled and practically effective framework for error detection in LLMs. Beyond hallucination detection, this perspective opens up promising directions for understanding the internal energy dynamics of autoregressive models, guiding future work on more trustworthy and interpretable AI systems.

# ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. Our study focuses on methodological contributions to error and hallucination detection in Large Language Models. We do not train new models or collect additional data; instead, we rely exclusively on publicly available datasets and widely used benchmark models for evaluation.

We note that part of our evaluation includes the Math dataset, which was publicly accessible at the time of experimentation but has since been taken down following a copyright claim. We emphasize that this dataset was used solely for evaluation purposes of our method, and only prior to the date of the takedown. No redistribution of the dataset was made, and our reported results are limited to demonstrating methodological effectiveness.

Our work does not involve personally identifiable information, sensitive content, or human subjects, and does not raise foreseeable risks of harm. We believe the proposed approach contributes positively to research on trustworthy AI by providing a training-free and generalizable framework for error detection in language models.

# REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our results. All experimental details, including model configurations, evaluation protocols, and datasets used, are described in the main text and Appendix B. Upon acceptance of this work, we will publicly release the code implementing our method, along with instructions to reproduce all reported experiments. This will allow the community to verify our findings and build upon our work.

### REFERENCES

- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*, 2023. 2
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, March 2022. doi: 10.1162/coli\_a\_00422. URL https://aclanthology.org/2022.cl-1.7/. 2
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021. 2
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The LLaMa 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024. 7, 8, 9, 19
- Jacob Dunefsky and Arman Cohan. One-shot optimized steering vectors mediate safety-relevant behaviors in LLMs. In *Second Conference on Language Modeling*, 2025. URL https://openreview.net/forum?id=teW4nIZ1qy. 4
- 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. Fact-checking the output of large language models via token-level uncertainty quantification. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 9367–9385, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.558. URL https://aclanthology.org/2024.findings-acl.558/. 4
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, June 2024. doi: 10.1038/s41586-024-07421-0. URL https://doi.org/10.1038/s41586-024-07421-0. 4

- Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence, 2025. URL https://arxiv.org/abs/2508.15260.4
  - Zorik Gekhman, Eyal Ben-David, Hadas Orgad, Eran Ofek, Yonatan Belinkov, Idan Szpektor, Jonathan Herzig, and Roi Reichart. Inside-out: Hidden factual knowledge in LLMs. In Second Conference on Language Modeling, 2025. URL https://openreview.net/forum?id=f7GG1MbsSM. 4
  - Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 5
  - Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *ICLR*, 2020. 2, 3, 5
  - Stevan Harnad. Language writ large: Llms, chatgpt, grounding, meaning and understanding. *arXiv* preprint arXiv:2402.02243, 2024. 2
  - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021. 7
  - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 5
  - 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. *arXiv* preprint arXiv:2311.05232, 2023. 2
  - Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023. 2
  - Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.7, 8, 9, 19
  - Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611, 2017. 7
  - Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022. 2, 7
  - Adam Tauman Kalai, Ofir Nachum, Santosh S. Vempala, and Edwin Zhang. Why language models hallucinate. Technical report, OpenAI and Georgia Tech, September 2025. Technical Report. 4
  - Michał P. Karpowicz. On the fundamental impossibility of hallucination control in large language models, 2025. URL https://arxiv.org/abs/2506.06382.4
  - Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In ICLR, 2014. 5
  - Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa Schut, Shreshth A Malik, and Yarin Gal. Semantic entropy probes: Robust and cheap hallucination detection in LLMs, 2025. URL https://openreview.net/forum?id=YQvvJjLWX0.4
  - Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Clam: Selective clarification for ambiguous questions with generative language models, 2023a. URL https://arxiv.org/abs/2212.07769.4

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*, 2023b. URL https://openreview.net/forum?id=VD-AYtPOdve. 4

- Yann Lecun, Sumit Chopra, Raia Hadsell, Marc Aurelio Ranzato, and Fu Jie Huang. *A tutorial on energy-based learning*. MIT Press, 2006. 4
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12286–12312, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.687. URL https://aclanthology.org/2023.acl-long.687.4
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 2
- Tianyu Liu, Yizhe Zhang, Chris Brockett, Yi Mao, Zhifang Sui, Weizhu Chen, and Bill Dolan. A token-level reference-free hallucination detection benchmark for free-form text generation. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6723–6737, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022. acl-long.464. URL https://aclanthology.org/2022.acl-long.464. 2
- Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *NeurIPS*, 2020. 3, 7
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL http://www.aclweb.org/anthology/P11-1015. 8
- Alessia McGowan, Yunlai Gui, Matthew Dobbs, Sophia Shuster, Matthew Cotter, Alexandria Selloni, Marianne Goodman, Agrima Srivastava, Guillermo A Cecchi, and Cheryl M Corcoran. Chatgpt and bard exhibit spontaneous citation fabrication during psychiatry literature search. *Psychiatry Research*, 326:115334, 2023. 2
- Beren Millidge. LLMs confabulate not hallucinate. *Beren's Blog*, March 2023. URL https://www.beren.io/2023-03-19-LLMs-confabulate-not-hallucinate/. 2
- Mujtaba Hussain Mirza, Maria Rosaria Briglia, Senad Beadini, and Iacopo Masi. Shedding more light on robust classifiers under the lens of energy-based models. In *ECCV*, 2024. 3
- Mujtaba Hussain Mirza, Maria Rosaria Briglia, Filippo Bartolucci, Senad Beadini, Giuseppe Lisanti, and Iacopo Masi. Understanding adversarial training with energy-based models, 2025. URL https://arxiv.org/abs/2505.22486.3
- OpenAI-Team. Gpt-4 technical report, 2023. URL https://cdn.openai.com/papers/gpt-4.pdf. 2
- Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and Yonatan Belinkov. Llms know more than they show: On the intrinsic representation of llm hallucinations. In *ICLR*, 2025. 1, 2, 4, 7, 8, 9, 15, 19

```
Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL https://arxiv.org/abs/2203.02155.4,5
```

- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. Check your facts and try again: Improving large language models with external knowledge and automated feedback, 2023. URL https://arxiv.org/abs/2302.12813.4
- Qwen-Team. Qwen3: Think deeper, act faster, 2025. URL https://qwen.ai/blog?id=1e3fa5c2d4662af2855586055ad037ed9e555125. Accessed: 2025-09-23. 7
- Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. In *OpenAI Technical Report*, 2018. URL https://cdn.openai.com/research-covers/language-unsupervised/language\_understanding\_paper.pdf. 5
- Vipula Rawte, Swagata Chakraborty, Agnibh Pathak, Anubhav Sarkar, S.M Towhidul Islam Tonmoy, Aman Chadha, Amit Sheth, and Amitava Das. The troubling emergence of hallucination in large language models an extensive definition, quantification, and prescriptive remediations. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2541–2573, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.155. URL https://aclanthology.org/2023.emnlp-main.155/. 2
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- Arleen Salles, Kathinka Evers, and Michele Farisco. Anthropomorphism in ai. *AJOB neuroscience*, 11(2):88–95, 2020. 2
- Andrea Santilli, Adam Golinski, Michael Kirchhof, Federico Danieli, Arno Blaas, Miao Xiong, Luca Zappella, and Sinead Williamson. Revisiting uncertainty quantification evaluation in language models: Spurious interactions with response length bias results. In *ACL*, 2025. 4
- Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. arXiv preprint arXiv:2307.00184, 2023. 2
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 5
- Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from pretrained language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Findings of the Association for Computational Linguistics: ACL 2022, pp. 566–581, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.48. URL https://aclanthology.org/2022.findings-acl.48/. 4
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through question-answering. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4631–4640, 2016. doi: 10.1109/CVPR.2016.501. 8
- Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1112–1122. Association for Computational Linguistics, 2018. URL http://aclweb.org/anthology/N18–1101. 8

Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models, 2025. URL https://arxiv.org/abs/2401.11817. 4

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, 2018. 7

Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuan-Jing Huang. Do large language models know what they don't know? In *ACL*, 2023. 2

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 15–20, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2003. URL https://aclanthology.org/N18-2003/. 8

Yao Zhu, Jiacheng Ma, Jiacheng Sun, Zewei Chen, Rongxin Jiang, Yaowu Chen, and Zhenguo Li. Towards understanding the generative capability of adversarially robust classifiers. In *ICCV*, pp. 7708–7717, 2021. 3, 6, 14

#### A APPENDIX

# A.1 PARTITION FUNCTIONS PROOF USED IN Eq. (4)

We extend the proof of Zhu et al. to the sequence-to-sequence setting by treating next-token prediction as a multi-class classification problem. At step i, the input is the prefix  $\{\mathbf{x}_{i-1:1}\}$ , and the model outputs logits over the vocabulary  $\mathcal V$  of size V. For notational consistency, we define the following energy terms:

$$\begin{cases}
E_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}_{i:1}) = -\log\left(\exp\left(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\operatorname{id}(\mathbf{x}_{i})]\right)\right), \\
E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1}) = -\log\left(\sum_{k=1}^{V}\exp\left(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]\right)\right).
\end{cases} (9)$$

The probability of the sequence up to position i can be expressed as

$$p_{\theta}(\mathbf{x}_{i:1}) = \frac{\exp(-E_{\theta}^{\ell}(\mathbf{x}_{i:1}))}{Z_{\theta}},$$
(10)

where  $Z_{\theta}$  is the global partition function (normalizing constant), defined over all possible continuations of all prefixes:

$$Z_{\boldsymbol{\theta}} = \sum_{\mathbf{x}_{i-1:1}} \sum_{\mathbf{x}_i} \exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\operatorname{id}(\mathbf{x}_i)]) = \sum_{\mathbf{x}_{i-1:1}} \sum_{k=1}^{V} \exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]). \tag{11}$$

Similarly, the probability of the prefix  $x_{i-1:1}$  can be written using the marginal energy:

$$p_{\theta}(\mathbf{x}_{i-1:1}) = \frac{\exp(-E_{\theta}^{m}(\mathbf{x}_{i-1:1}))}{\widetilde{Z}_{\theta}},$$
(12)

where  $\widetilde{Z}_{\theta}$  is the corresponding normalizing constant:

$$\widetilde{Z}_{\boldsymbol{\theta}} = \sum_{\mathbf{x}_{i-1:1}} \exp(-E_{\boldsymbol{\theta}}^{m}(\mathbf{x}_{i-1:1})) = \sum_{\mathbf{x}_{i-1:1}} \exp\left(\log \sum_{k=1}^{V} \exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k])\right).$$
(13)

By expanding the logarithm in Eq. (13), we obtain

$$\widetilde{Z}_{\boldsymbol{\theta}} = \sum_{\mathbf{x}_{i-1:1}} \sum_{k=1}^{V} \exp(\boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]), \qquad (14)$$

which is identical to Eq. (11). Hence, the two partition functions coincide:

$$Z_{\theta} = \widetilde{Z}_{\theta}. \tag{15}$$

### A.2 THE ROLE OF TEMPERATURE IN SPILLED ENERGY

We now analyze how the temperature parameter  $\tau$  affects the definition of spilled energy. Starting from Eq. (3), the probability of the next token under temperature scaling is

$$\log p_{\theta}(\mathbf{x}_i \mid \mathbf{x}_{i-1:1}) = \log \frac{\exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\mathrm{Id}(\mathbf{x}_i)])}{\sum_{k} \exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k])}$$
(16)

$$= \frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[\operatorname{Id}(\mathbf{x}_i)] - \log \sum_{k} \exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]). \tag{17}$$

Accordingly, the spilled energy becomes

$$\Delta E_{\theta}(\mathbf{x}_{i:1}) = \frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1}) [\mathrm{Id}(\mathbf{x}_i)] - \log \sum_{k=1}^{|V|} \exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_i, \dots, \mathbf{x}_1)[k]).$$
(18)

**Limit case**  $\tau \to \infty$ . When the temperature tends to infinity, the logits are scaled down towards zero, making all tokens equally likely:

$$\lim_{\tau \to +\infty} \Delta E_{\theta}(\mathbf{x}_{i:1}) = \lim_{\tau \to \infty} \frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1}) [\mathrm{Id}(\mathbf{x}_{i})] - \log \sum_{k=1}^{|V|} \exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k])$$
(19)

$$= 0 - \log \sum_{k=1}^{|V|} \exp(0)$$
 (20)

$$= -\log|V|. \tag{21}$$

Thus, for  $\tau \to \infty$  the model degenerates into a uniform random classifier over the vocabulary.

**Interpretation.** Varying  $\tau$  perturbs the balance between the two energy terms, introducing a systematic error in  $\Delta E_{\theta}$ . From the perspective of the Boltzmann distribution, scaling by  $\frac{1}{\tau}$  corresponds to injecting or removing energy from the system. At high temperatures  $(\tau \to \infty)$ , the system approaches maximum entropy, where all tokens have equal probability. At low temperatures  $(\tau \to 0^+)$ , the distribution collapses onto the maximum logit token, making the model highly deterministic.

**Error accumulation.** As we generate tokens sequentially, we accumulate deviations in  $\Delta E_{\theta}$ :

$$\log p_{\theta}(\mathbf{x}_{i-1:1}) = \frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1}) [\mathrm{Id}(\mathbf{x}_i)] - \log \sum_{k} \exp(\frac{1}{\tau} \boldsymbol{\theta}(\mathbf{x}_{i-1:1})[k]) + \sum_{i=1}^{t} \Delta E_{\theta}(\mathbf{x}_{j:1}).$$
 (22)

Hence, temperature scaling not only modifies the probabilities but also reshapes the cumulative error landscape traced by spilled energy.

### B REPRODUCIBILITY

For comparability, we adopt the same experimental setting as Orgad et al. (2025), whose implementation is publicly available at https://github.com/technion-cs-nlp/LLMsKnow. This ensures that our baselines and evaluation procedures follow an established and validated protocol.

In addition, we will release our own codebase, which includes:

- computation of the proposed energy-based measures;
- scripts for reproducing the synthetic arithmetic preliminary experiments.

The code and instructions will be made available upon acceptance of this work to facilitate full reproducibility of our results.

# C LLM USAGE

Large language models were used exclusively for text polishing and minor exposition refinements. All substantive research content, methodology, and scientific conclusions were developed entirely by the authors.

# D SUPPLEMENTARY MATERIAL

This supplementary material is intended to complement the main paper by providing further motivation for our assumptions and design choices, as well as additional ablation studies or additional plots, such as ROCs and histograms, that could not fit in the main paper.

### D.1 ADDITIONAL RESULTS FOR SYNTHETIC ARITHMETIC

In Fig. 4 we augmented Fig. 3 in the main paper, adding also the results for **Mistral-7B-Instruct v0.3** and **LLaMa-3-8B**. The same findings of the figure in the paper also translate to this LLM, meaning that our method generalizes across LLMs.

Fig. 5 and Fig. 6 also extend and provide more details of Fig. 3 in the main paper by showing, respectively, the histograms and the ROC at a better resolution and displayed in different frames. Also, we have added results for Mistral-7B-Instruct v0.3 and LLaMa-3-8B.

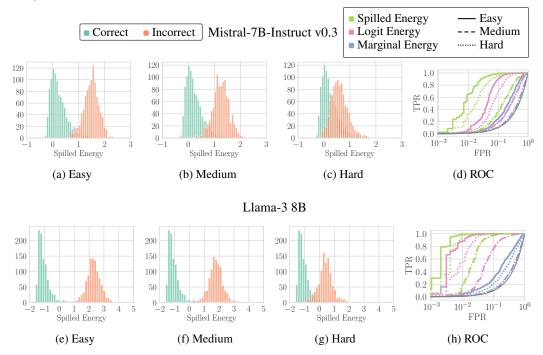


Figure 4: Histograms of Spilled Energy values across models (rows) on Math Sums with different error ranges in the answer (columns, decreasing range left to right, making it harder to detect errors), as described in Section 5.1. In the fourth column, we show ROC curves for Hallucination Detection across the error ranges (colors) and methods (line styles).

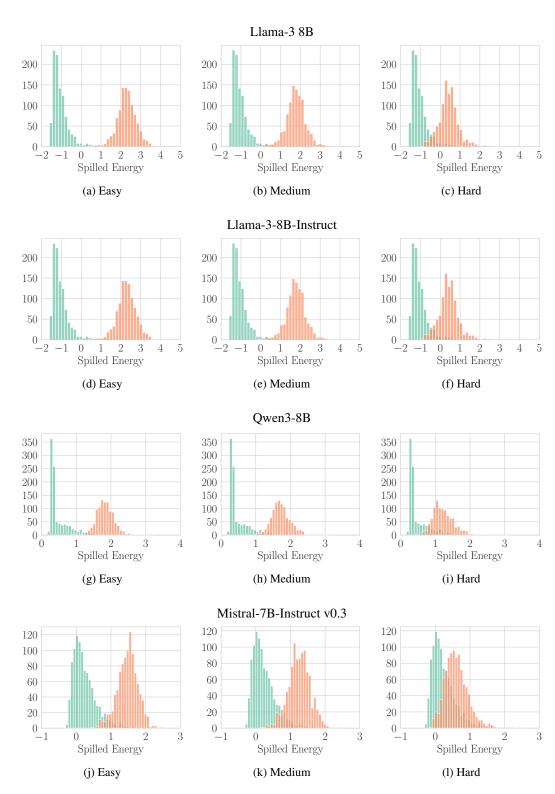


Figure 5: Histograms of Spilled Energy values for Correct and Incorrect answers across models on Math Sums, increasing difficulty from left to right. We compute sums on 13-digit integers, for incorrect answers we add a random offset sampled uniformly from the error interval: Easy  $\sim \mathcal{U}(1e3, 1e4)$  - Medium  $\sim \mathcal{U}(1e2, 1e3)$  - Hard  $\sim \mathcal{U}(1, 10)$ ; for more details see Section 5.1.

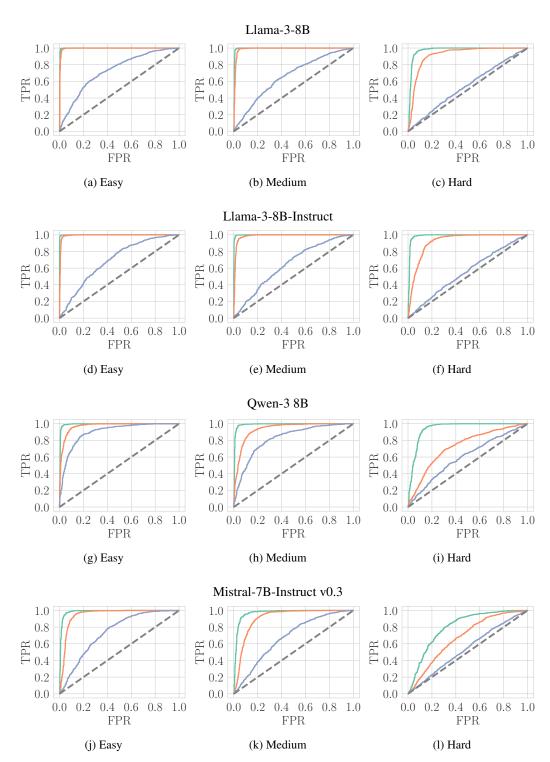


Figure 6: ROC curves for Hallucination Detection across models (rows) on Math Sums with different error ranges in the answer (columns, decreasing range left to right). All sums are performed on 13-digit integers. Legend: **Spilled (ours)** Spilled  $\Delta E$  Logit  $E^{\ell}$  Marginal  $E^{m}$ 

	Pool	HotpotQA	HotpotQA-WC	IMDB	Math	MNLI	Movies	TriviaQA	Winobias	Winogrande	Average
	LLaMA-3-Instruct (Dubey et al., 2024)										
Orgad et al.	Mean	$66.56{\scriptstyle\pm9.10}$	$59.00{\scriptstyle\pm8.14}$	$69.78 \scriptstyle{\pm 14.76}$	$66.56{\scriptstyle\pm17.04}$	$60.56 \scriptstyle{\pm 12.53}$	$66.44{\scriptstyle\pm8.06}$	$63.22{\scriptstyle\pm11.11}$	$\textbf{67.33} \scriptstyle{\pm 11.97}$	$\textbf{58.00} {\scriptstyle\pm7.79}$	$64.16{\scriptstyle\pm11.17}$
Spilled $\Delta E$	Min	<b>85.98</b> ±1.09	93.00±1.61	$47.66{\scriptstyle\pm4.06}$				<b>87.07</b> ±1.33		55.11±2.05	73.16±2.10
Marginal $E^m$	Max	$76.72 \pm 1.38$	$30.74 \pm 3.45$	$85.63 \pm 2.39$					$57.67 \pm 2.94$	$47.47 \pm 1.83$	$65.72 \pm 2.31$
Marginal $E^m$	Min	$75.91 \pm 1.62$	$97.57 \pm 0.75$	$14.37 \pm 2.39$				$73.38{\scriptstyle\pm1.86}$		$53.98 \pm 2.30$	$62.93 \pm 2.10$
Logit $E^{\ell}$	Max	$72.85{\scriptstyle\pm2.12}$	$91.11 \pm 1.52$	$42.08{\scriptstyle\pm5.07}$				$68.89{\scriptstyle\pm1.96}$		$49.40 \pm 2.16$	54.62±2.60
Spilled $\Delta E$	Max	$54.34 \pm 1.58$	$47.68{\scriptstyle\pm2.81}$	$52.34 \pm 4.06$					$38.40 \pm 2.61$		$50.07 \pm 2.52$
Spilled $\Delta E_s$	Max	$53.65 \pm 1.40$	$36.28 \pm 2.99$	$55.80{\scriptstyle\pm4.32}$				$48.70{\scriptstyle\pm2.44}$		$44.32 \pm 1.70$	48.87±2.59
Marginal $E^m$	Mean	$44.64 \pm 1.70$	$4.82 \pm 0.89$	$85.63 \pm 2.39$				$46.12{\scriptstyle\pm2.60}$		$46.41 \pm 2.00$	$46.99 \pm 2.12$
	Last Token		$71.92 \pm 2.97$					$39.97 \pm 2.80$		$49.47 \pm 2.28$	43.35±3.12
Logit $E^{\ell}$	ALT	43.04±1.87	$71.92 \pm 2.97$	$42.08{\scriptstyle\pm5.07}$	$59.89{\scriptstyle\pm5.83}$	19.74±2.28	$21.57{\scriptstyle\pm2.02}$	$39.97 \pm 2.80$	$42.45 \pm 2.94$	$49.47{\scriptstyle \pm 2.28}$	43.35±3.12
		Mistral-Instruct (Jiang et al., 2023)									
Orgad et al.	Mean	$64.78{\scriptstyle\pm10.56}$	56.78±7.95	$\pmb{82.67} \scriptstyle{\pm 11.63}$	$\textit{68.78}\scriptstyle\pm11.43$	$64.22{\scriptstyle\pm12.12}$	$64.89{\scriptstyle\pm11.55}$	$65.44{\scriptstyle\pm12.10}$	$\textbf{61.00} \scriptstyle{\pm 12.23}$	$\textbf{61.44} \scriptstyle{\pm 11.31}$	$65.56 \scriptstyle{\pm 11.21}$
Spilled $\Delta E$	Min	91.12±1.10	97.47±0.78	59.77±2.57	66.63±3.46	95.95±0.83	<b>94.99</b> ±0.93	91.75±1.01	50.74±3.15	$49.00 \pm 1.92$	77.49±1.75
Marginal $E^m$	Min	87.58±1.35	$97.94 \pm 0.62$	$18.67{\scriptstyle\pm2.27}$	$67.58{\scriptstyle\pm3.37}$	<b>97.96</b> ±0.55	$84.90 \pm 1.37$	$87.75 \pm 1.73$	$49.19 \pm 3.97$	$48.49 \pm 1.86$	71.12±1.90
Logit $E^{\ell}$	Max	$77.24 \pm 1.66$	$83.84 \pm 1.66$	$22.28 \pm 2.54$	$57.67 \pm 3.29$	$78.98{\scriptstyle\pm1.58}$	$76.89 \pm 1.49$	$80.35{\scriptstyle\pm1.88}$	$45.53 \pm 2.60$	$48.17 \pm 1.97$	63.44±2.07
Marginal $E^m$	Max	$64.63 \pm 1.97$	$33.42 \pm 1.90$	$81.33 \pm 2.32$	$26.52{\scriptstyle\pm2.28}$	$17.62 \pm 1.20$	86.60±1.20	$65.46 \pm 2.25$	56.41±4.44	$51.14 \pm 1.71$	53.68±2.14
Logit $E^{\ell}$	Last Token	$55.77 \pm 2.38$	$71.26 \pm 2.28$	$22.28 \pm 2.54$	$71.21 \pm 2.42$	$47.78 \pm 2.26$	$42.93{\scriptstyle\pm1.96}$	$58.36{\scriptstyle\pm3.52}$	$45.65 \pm 2.94$	$48.30 \pm 2.04$	51.50±2.48
Logit $E^{\ell}$	ALT	55.77±2.38	$71.26 \pm 2.28$	$22.28 \pm 2.54$	71.21±2.42	$47.78 \pm 2.26$	$42.93 \pm 1.96$	$58.36 \pm 3.52$	$45.65 \pm 2.94$	$48.30 \pm 2.04$	51.50±2.48
Spilled $\Delta E_s$	Max	$49.13 \pm 2.50$	$36.37 \pm 2.40$	$46.45{\scriptstyle\pm2.56}$	$29.05 \pm 2.57$	$53.79{\scriptstyle\pm1.55}$	$55.24 \pm 2.17$	$46.73{\scriptstyle\pm1.98}$	53.30±3.66	$51.20 \pm 1.84$	46.81±2.36
Spilled $\Delta E$	Max	$49.49 \pm 2.38$	$41.07 \pm 2.26$	$40.25{\scriptstyle\pm2.56}$	$30.82{\scriptstyle\pm2.53}$	$49.28{\scriptstyle\pm1.61}$	$53.56{\scriptstyle\pm2.02}$	$47.55{\scriptstyle\pm1.90}$	$53.20 \pm 3.77$	$51.38 \pm 1.97$	46.29±2.33
Logit $E^{\ell}$	Mean	$45.56{\scriptstyle\pm1.78}$	$54.88{\scriptstyle\pm2.06}$	$22.27{\scriptstyle\pm2.53}$	$58.59 \scriptstyle{\pm 2.34}$	$53.85{\scriptstyle\pm2.13}$	$31.88{\scriptstyle\pm2.05}$	$51.26{\scriptstyle\pm2.58}$	$44.69{\scriptstyle\pm2.98}$	$48.40{\scriptstyle\pm1.99}$	45.71±2.27
	Mistral (Jiang et al., 2023)										
Orgad et al.	Mean	$61.78 \scriptstyle{\pm 9.27}$	$57.44{\scriptstyle\pm6.95}$	$\textbf{76.22} {\scriptstyle\pm 12.82}$	$65.78 \scriptstyle{\pm 15.27}$	$56.67 \scriptstyle{\pm 11.83}$	$64.22{\scriptstyle\pm8.91}$	<b>64.33</b> ±10.40	$58.00{\scriptstyle\pm12.29}$	<b>54.56</b> ±4.36	$62.11{\scriptstyle\pm10.23}$
Marginal $E^m$	Min	87.52±1.31	90.91±1.58	54.69±2.49	<b>86.21</b> ±1.96	<b>98.80</b> ±0.35	94.41±0.62	83.66±2.16	52.15±1.74	46.37±2.02	77.19±1.58
Marginal $E^m$	Max	$83.57 \pm 1.13$	$86.83 \pm 1.70$	$45.31 {\scriptstyle \pm 2.49}$	$62.26 \pm 4.29$	$96.03{\scriptstyle\pm0.83}$	99.27±0.24	$92.26 \pm 1.31$	$51.31{\scriptstyle\pm3.35}$	$54.49_{\pm 2.48}$	74.59±1.98
Spilled $\Delta E$	Min	84.24±1.18	$83.74 \pm 1.41$	$57.43 \pm 2.99$	$78.26 \pm 2.93$	$96.69 \pm 0.62$	$84.47 \pm 1.17$	$81.27{\scriptstyle\pm1.83}$	$50.62{\scriptstyle\pm1.72}$	$48.72{\scriptstyle\pm1.75}$	$73.94 \pm 1.74$
Spilled $\Delta E$	Max	$61.50{\scriptstyle\pm1.88}$	$63.60 \pm 1.68$	$42.57{\scriptstyle \pm 2.99}$	$76.27{\scriptstyle\pm3.42}$	$47.01{\scriptstyle\pm2.48}$	$81.84{\scriptstyle\pm1.60}$	$68.07{\scriptstyle\pm1.30}$	<b>58.71</b> ±3.69	$51.13 \pm 1.87$	61.19±2.32
Spilled $\Delta E_s$	Max	$60.54 \pm 1.81$	$60.18 \pm 1.84$	$43.47 \pm 2.76$	$71.93 \pm 3.62$	$45.94 \pm 2.40$	$78.84 \pm 1.53$	$67.92 \pm 1.32$	$57.24 \pm 3.72$	$51.88 \pm 1.90$	59.77±2.32
Logit $E^{\ell}$	Max	$49.54 \pm 1.42$	$52.47 \pm 1.61$	$32.72{\scriptstyle\pm2.89}$				$39.73{\scriptstyle\pm2.03}$		$44.41 \pm 2.42$	49.51±2.36
Spilled $\Delta E$	Last Token		$50.34 \pm 2.13$	$42.57{\scriptstyle\pm2.99}$	$55.15 \pm 3.40$	$46.47{\scriptstyle\pm2.47}$	$49.75{\scriptstyle\pm2.04}$	$45.25{\scriptstyle\pm1.71}$	$58.46 \pm 3.76$	$51.07 \pm 1.86$	48.87±2.47
Spilled $\Delta E$	ALT	$40.79{\scriptstyle\pm1.85}$	$50.34{\scriptstyle\pm2.13}$	$42.57{\scriptstyle\pm2.99}$					$58.46 \pm 3.76$		48.87±2.47
Spilled $\Delta E_s$	ALT	$41.75 \pm 1.80$	$50.64 \pm 2.29$						$57.14{\scriptstyle\pm3.75}$		48.68±2.45
Spilled $\Delta E_s$			$50.64 \pm 2.29$	$43.47{\scriptstyle\pm2.76}$	$51.23 \pm 3.49$	$45.49{\scriptstyle\pm2.37}$	$49.74 \pm 1.96$	$46.98{\scriptstyle\pm1.81}$	$57.14{\scriptstyle\pm3.75}$	$51.67{\scriptstyle\pm1.81}$	48.68±2.45
Logit $E^{\ell}$	ALT	$39.31{\scriptstyle\pm1.26}$	$44.80 \pm 2.05$	$32.72{\scriptstyle\pm2.89}$	$54.35{\scriptstyle\pm4.27}$	$92.29{\scriptstyle\pm1.25}$	$23.69{\scriptstyle\pm1.44}$	$32.58{\scriptstyle\pm1.89}$	$46.53{\scriptstyle\pm3.80}$	$44.42{\scriptstyle\pm2.42}$	45.63±2.37
Logit $E^{\ell}$	Last Token	$39.31 \pm 1.26$	$44.80 \pm 2.05$	$32.72{\scriptstyle\pm2.89}$	$54.35{\scriptstyle\pm4.27}$	$92.29{\scriptstyle\pm1.25}$	$23.69{\scriptstyle\pm1.44}$	$32.58{\scriptstyle\pm1.89}$	$46.53{\scriptstyle\pm3.80}$	$44.42{\scriptstyle\pm2.42}$	45.63±2.37

Table 3: Hallucination detection performance, in terms of AuROC, across nine benchmarks and different LLMs. We measure the generalization across all tasks by computing the average.

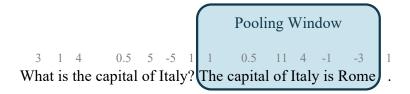


Figure 7: Example of the Pooling Window

### D.2 ADDITIONAL RESULTS FOR CROSS-TESTING CON REAL WORLD BENCHMARKS

Table 3 shows how our method compares with the baselines methods, Orgad et al. (2025) and Logit  $E^{\ell}$ . This table was obtained by using various pooling methods in the pooling frame from which we measure the hallucination. More details below alongside the examples based on Fig. 7:

- $\diamond$  **Min**: minimum energy value in the pooling frame. Energy Measured: -3
- ♦ Max: maximum energy value in the pooling frame. Energy Measured: 11
- ♦ **Mean**: mean among all the energies in the pooling frame. Energy Measured: 2.08
- ♦ **Last Token**: energy on the last token of the pooling frame. Energy Measured: −3
- ♦ After Last Token: energy of the first token after the pooling method. Energy Measured: 1

# D.3 ADDITIONAL QUALITATIVE RESULTS

In this section, we offer additional results of the detection performance following what is shown in Fig. 1. We report both success cases and failure cases. While it is difficult to draw conclusions and predict when, why, and on which topics spilled energy may work or not, we noticed that it appears to perform reliably on knowledge-based factual content but exhibits difficulties with reasoning tasks and numerical information, despite working well on math questions as demonstrated in Section 5.1. Further investigation is required to better understand and validate these patterns.

### D.3.1 Success Cases

 Question: '`Which planet is known as the Red Planet ?''

```
Logits: The Red Planet is Mars . 

Ours: The Red Planet is Mars .
```

```
Logits: The Red Planet is Jupiter . X

Ours: The Red Planet is Jupiter . X
```

Question: ``What is the largest mamm al in the world ?''

```
Logits: The largest mamm al in the world is the Blue Whale 
Ours: The largest mamm al in the world is the Blue Whale
```

```
Logits: The largest mamm al in the world is the House Cat. X

Ours: The largest mamm al in the world is the House Cat. X
```

Question: 'Who painted the Mona Lisa?''

```
Logits: The Mona Lisa was painted by Leonardo da Vinci . ✓
Ours: The Mona Lisa was painted by Leonardo da Vinci . ✓
```

```
Logits: The Mona Lisa was painted by Pablo Esc obar . X

Ours: The Mona Lisa was painted by Pablo Esc obar . X
```

Question: 'What gas do plants breathe in for photosyintesis ?''

```
Logits: They breathe in carbon dioxide ✓
Ours: They breathe in carbon dioxide ✓
```

```
Logits: They breathe in oxygen X

Ours: They breathe in oxygen X
```

```
1080
       Question: ''In which continent is Egypt Located ?''
1081
1082
           Logits: Egypt is located in Africa /
1083
           Ours: Egypt is located in Africa /
1084
1085
1086
           Logits: Egypt is located in Europe X
1087
1088
           Ours: Egypt is located in Europe X
1089
1090
        Question:
                     "What is the fastest land animal ?"
1091
1092
1093
           Logits: The fastest land animal is the che et ah 🗸
1094
           Ours: The fastest land animal is the che et ah \( \sqrt{} \)
1095
1096
1097
           Logits: The fastest land animal is the lion X
1098
           Ours: The fastest land animal is the lion X
1099
1100
1101
                      "What is the hardest natural substance on Earth ?"
       Ouestion:
1102
1103
           Logits: The hardest natural substance is diamond ✓
1104
           Ours: The hardest natural substance is diamond \( \sqrt{} \)
1105
1106
1107
           Logits: The hardest natural substance is gold X
1108
           Ours: The hardest natural substance is gold X
1109
1110
1111
                      "Which ocean is the largest ?"
        Question:
1112
1113
1114
           Logits: The largest ocean is the Pacific Ocean ✓
1115
           Ours: The largest ocean is the Pacific Ocean \checkmark
1116
1117
1118
           Logits: The largest ocean is the Indian Ocean X
1119
           Ours: The largest ocean is the Indian Ocean X
1120
1121
1122
       D.3.2 FAILURE CASES
1123
1124
                      ``Who was the first person to walk on the mood ?^{\prime\prime}
       Question:
1125
1126
           Logits: Neil Armstrong 

✓
1127
           Ours: Neil Armstrong ✓
1128
1129
1130
           Logits: Buzz Ald rin X
1131
           Ours: Buzz Ald rin X
1132
1133
```

```
1134
       Reasoning: 'Marta has 2 sisters and 2 brothers. How many sisters has
1135
        one of the brothers of Marta?"
1136
1137
           Logits: Mart a has 2 sisters and 2 brothers. Each brother has the same 2 sisters
1138
           . So, one of Mart a 's brothers has 2 sisters . \( \square$
1139
1140
           Ours: Mart a has 2 sisters and 2 brothers. Each brother has the same 2 sisters.
1141
           So, one of Mart a 's brothers has 2 sisters.
1142
1143
1144
           Logits: Mart a has 2 sisters and 2 brothers. Each brother has the same 3 sisters
1145
           . So, one of Mart a 's brothers has 3 sisters . X
1146
1147
           Ours: Mart a has 2 sisters and 2 brothers. Each brother has the same 3 sisters.
1148
           So, one of Mart a 's brothers has 3 sisters. X
1149
1150
1151
        Reasoning: ''If there are 3 cars and each car has 4 wheels , how many
1152
       wheels are there in total ? ''
1153
1154
           Logits: Each car has 4 wheels . So , for 3 cars , the total number of wheels is 3
1155
           x = 4 = 12 wheels . \checkmark
1156
1157
           Ours: Each car has 4 wheels. So, for 3 cars, the total number of wheels is 3
1158
           x \mid 4 = 12 wheels . \checkmark
1159
1160
1161
           Logits: Each car has 8 wheels . So , for 3 cars , the total number of wheels is 3
1162
           x = 8 = 14 wheels . X
1163
1164
           Ours: Each car has 8 wheels. So, for 3 cars, the total number of wheels is 3
1165
           x \mid 8 = 14 wheels . X
1166
1167
        Reasoning: 'What is the square root of 64?''
1168
1169
1170
           Logits: The square root of 64 is 8 \checkmark
1171
           Ours: The square root of 64 is 8 \( \square \)
1172
1173
1174
           Logits: The square root of 64 is 10 X
1175
           Ours: The square root of 64 is 10 X
1176
1177
1178
        Question:
                     ''What blood type is known as the universal donor ?''
1179
1180
           Logits: O negative ✓
1181
           Ours: O negative ✓
1182
1183
1184
           Logits: AB positive X
1185
           Ours: AB positive X
1186
1187
```