

ICLR 2026

# Spilled Energy in Large Language Models

*A training-free, energy-based framework for hallucination detection in LLMs*

Adrian R. Minut · Hazem Dewidar · Iacopo Masi

*Sapienza University of Rome · OmnAI Lab · GLADIA*

# The Hallucination Problem

## What are Hallucinations?

**Hallucinations** are any form of error produced by an LLM, including:

- Factual mistakes
- Biased outputs
- Commonsense reasoning failures
- Unverifiable claims

## Limitations of Existing Methods

### Probe Classifiers

Require training per task; don't generalize across datasets

### Logit Confidence

Weak baseline; degrades with instruction tuning

### $p(\text{true})$

Marginally better than random guessing (~51%)

→ We need a training-free, task-agnostic hallucination detector

# Core Idea: LLMs as Energy-Based Models

## Energy-Based Models (EBMs)

EBMs assign a scalar energy to each configuration: lower energy = higher likelihood.  
The probability distribution is:  $p(\mathbf{x}) = \exp(-E(\mathbf{x})) / Z$  where  $Z$  is the partition function.

### 1 Autoregressive LLMs

LLMs model  $p(x_{N:1})$  as a chain of conditional probabilities:  $\prod_{i=1}^N p(x_i | x_{i-1:1})$

Each step is a softmax classifier over the vocabulary  $V$ .



### 2 Reinterpret as EBMs

Each softmax classifier can be recast as an EBM ratio:  $\log p(x_i | x_{i-1:1}) = -E_{i:1}^l + E_{i-1:1}^m$



### 3 Discover the Spill

Across consecutive steps, **these two energy quantities should be equal**, but they differ in practice since there is **no training constraint**. We call this discrepancy *spilled energy*.

# Core Idea: LLMs as Energy-Based Models

## Energy-Based Models (EBMs)

EBMs assign a scalar energy to each configuration: lower energy = higher likelihood.

The probability distribution is:  $p(x) = \exp(-E(x)) / Z$  where  $Z$  is the partition function.

## Language Modeling as Chain of Probabilities

$$p(x_{i:1}) = \exp(-E(x_{i:1})) / Z$$

$$p(x_{N:1}) = \prod_{i=1}^N p(x_i | x_{i-1:1}) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \dots$$

$$\log p(x_3, x_2, x_1) = \log p(x_1) + \log p(x_2 | x_1) + \log p(x_3 | x_2, x_1)$$

$$\log p(x_3, x_2, x_1) = \underbrace{-E(x_1)} - \log Z - \underbrace{E(x_2, x_1)} - \log Z + \underbrace{E(x_1)} + \log Z - \underbrace{E(x_3, x_2, x_1)} - \log Z + \underbrace{E(x_2, x_1)} + \log Z = -E(x_3, x_2, x_1) - \log Z$$

# Core Idea: LLMs as Energy-Based Models

## Energy-Based Models (EBMs)

EBMs assign a scalar energy to each configuration: lower energy = higher likelihood.

The probability distribution is:  $p(\mathbf{x}) = \exp(-E(\mathbf{x})) / Z$  where  $Z$  is the partition function.

## Language Modeling in Practice

---

$$p(x_{i:1}) = \exp(-E(x_{i:1})) / Z$$

$$p(x_i | x_{i-1:1}) = \frac{\exp(\theta(x_{i-1:1})[id(x_i)])}{\sum_k \exp(\theta(x_{i-1:1})) [k]} = [\exp(-E(x_{i-1:1}) + E(x_{i:1}))] Z / Z$$

$$E^l(x_{i:1}) = \theta(x_{i-1:1})[id(x_i)]$$
$$E^m(x_{i-1:1}) = \log \sum_k \exp(\theta(x_{i-1:1})) [k]$$

# What is Spilled Energy?

## Definition 4.1 — Spilled Energy $\Delta E(x_{i:1})$

The discrepancy between two energy values that, in principle, should be equal but are measured at different time steps and from different components of the LLM softmax:

$$\Delta E(x_{i:1}) = -\log \sum \exp(\theta(x_{i:1})) [k] + \theta(x_{i-1:1}) [id(x_i)]$$

## $E^\ell$ — Logit Energy

$$E^\ell = \theta(x_{i-1:1}) [id(x_i)]$$

Measured at step  $i$ :  
The **logit** of the **sampled token**, i.e., the numerator of the softmax.

This is the energy assigned to the chosen token sequence.



## $E^m$ — Marginal Energy

$$E^m = \log \sum_k \exp(\theta(x_{i:1})) [k]$$

Measured at step  $i + 1$ :  
The **log-sum-exp** of **ALL logits**, i.e., the denominator of the softmax.

This marginalizes over all possible next tokens.

# Two Training-Free Metrics



## Spilled Energy $\Delta E$

Two consecutive steps

Captures the discrepancy between energy values that should theoretically be equal across consecutive generation steps.

$\Delta E$  is large  $\rightarrow$  model internally inconsistent  $\rightarrow$  likely hallucinating.

$$\Delta E(x_i : i) = \theta(x_{i-1} : i)[\text{id}(x_i)] - \log \sum \exp(\theta(x_i : i)[k])$$



## Marginal Energy $E^m$

Single step

Measurable at a single decoding step.  
Corresponds to the log-partition function: how “spread out” the probability mass is over the entire vocabulary.

Low marginal energy  $\rightarrow$  model is more certain  $\rightarrow$  likely correct.

$$E^m(x_{i-1} : i) = -\log \sum \exp(\theta(x_{i-1} : i)[k])$$

# Key Insight: Exact Answer Token Localization

*Detection is most reliable when restricted to the "exact answer" tokens: those carrying the semantic content of the answer (e.g., "Rome" in "The capital of Italy is Rome").*

Example: "What is the capital of Italy?"



## Answer Extraction Strategies

### Heuristic Matching

For classification & multiple-choice tasks with a closed label set: string matching.

### LLM-based Extraction

For open-ended tasks (TriviaQA, Math): an auxiliary Mistral-7B model extracts the short answer from the full generation.

### Pooling Strategy

When the answer spans multiple tokens, apply min / max / mean pooling over the localized span to get a final score.

★ Using exact answer tokens provides +24% AuROC improvement for spilled energy vs. full sequence evaluation

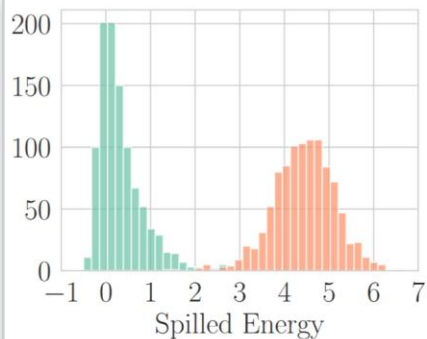
# Experiment 1: Synthetic Arithmetic

Setup: Multi-digit arithmetic (13-digit integers) with correct vs. deliberately wrong answers. Tested on LLaMA-3, Qwen-3, and Mistral-7B.

## Easy

Offset [1,000–10,000]

Large errors — clearly wrong answers. Spilled energy cleanly separates correct vs. incorrect distributions.

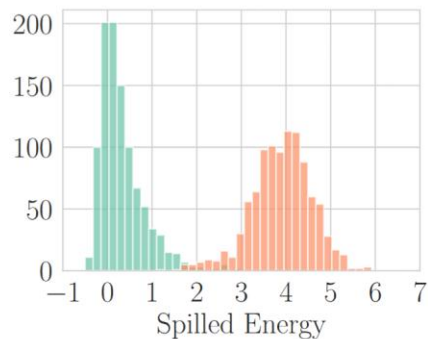


(a) Easy

## Medium

Offset [100–1,000]

Moderate errors. Spilled energy still maintains strong separation even as the gap shrinks.

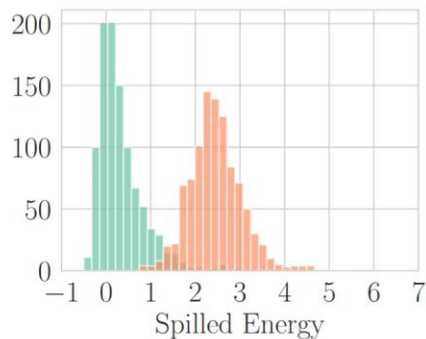


(b) Medium

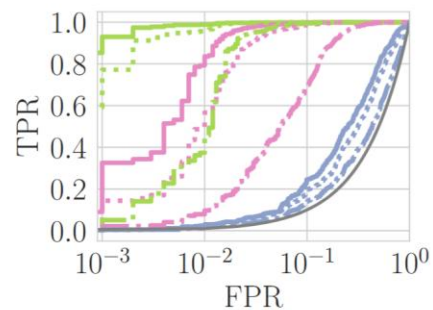
## Hard

Offset [1–10]

Subtle errors — hardest to detect. Spilled energy outperforms logit-based detection at this difficulty.



(c) Hard



(d) ROC

# Experiment 2: Cross-Dataset Generalization

*Probing classifiers degrade sharply out-of-distribution, while spilled energy generalizes without any retraining.*

	Pool	HotpotQA	HotpotQA-WC	IMDB	Math	MNLI	Movies	TriviaQA	Winobias	Winogrande	Average
LLaMA-Instruct <a href="#">Dubey et al. (2024)</a>											
$p(\text{true})$	—	58.31 $\pm$ 0.32	51.66 $\pm$ 1.05	50.72 $\pm$ 1.20	49.53 $\pm$ 2.16	52.33 $\pm$ 0.98	59.30 $\pm$ 0.85	45.99 $\pm$ 0.51	45.47 $\pm$ 1.58	48.33 $\pm$ 0.68	51.29 $\pm$ 04.86
<a href="#">Orgad et al.</a> Mean		66.56 $\pm$ 9.10	59.00 $\pm$ 8.14	69.78 $\pm$ 14.76	66.56 $\pm$ 17.04	60.56 $\pm$ 12.53	66.44 $\pm$ 8.06	63.22 $\pm$ 11.11	<b>67.33</b> $\pm$ 11.97	<b>58.00</b> $\pm$ 7.79	64.16 $\pm$ 03.90
Logit $E^\ell$ Max		72.85 $\pm$ 2.12	91.11 $\pm$ 1.52	42.08 $\pm$ 5.07	57.81 $\pm$ 3.82	25.52 $\pm$ 3.00	43.97 $\pm$ 1.38	68.89 $\pm$ 1.96	39.95 $\pm$ 2.41	49.40 $\pm$ 2.16	54.62 $\pm$ 18.97
Marginal $E^m$ Max		76.72 $\pm$ 1.38	30.74 $\pm$ 3.45	<b>85.63</b> $\pm$ 2.39	27.08 $\pm$ 5.06	<b>89.90</b> $\pm$ 1.25	<b>96.17</b> $\pm$ 0.63	80.13 $\pm$ 1.87	57.67 $\pm$ 2.94	47.47 $\pm$ 1.83	65.72 $\pm$ 24.39
Marginal $E^m$ Min		75.91 $\pm$ 1.62	<b>97.57</b> $\pm$ 0.75	14.37 $\pm$ 2.39	<b>70.55</b> $\pm$ 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 $\pm$ 21.89
Spilled $\Delta E_s$ Max		53.65 $\pm$ 1.40	36.28 $\pm$ 2.99	55.80 $\pm$ 4.32	35.44 $\pm$ 3.41	58.81 $\pm$ 2.58	70.30 $\pm$ 1.49	48.70 $\pm$ 2.44	36.53 $\pm$ 2.98	44.32 $\pm$ 1.70	48.87 $\pm$ 11.26
Spilled $\Delta E$ Min		<b>85.98</b> $\pm$ 1.09	93.00 $\pm$ 1.61	47.66 $\pm$ 4.06	65.58 $\pm$ 3.02	73.95 $\pm$ 1.97	89.34 $\pm$ 1.04	<b>87.07</b> $\pm$ 1.33	60.72 $\pm$ 2.74	55.11 $\pm$ 2.05	<b>73.16</b> $\pm$ 15.64
LLaMA <a href="#">Dubey et al. (2024)</a>											
$p(\text{true})$	—	52.83 $\pm$ 0.71	49.33 $\pm$ 0.86	52.30 $\pm$ 0.58	58.63 $\pm$ 1.26	53.78 $\pm$ 0.70	60.76 $\pm$ 0.69	62.94 $\pm$ 0.51	50.02 $\pm$ 1.24	53.47 $\pm$ 0.54	54.90 $\pm$ 04.77
<a href="#">Orgad et al.</a> Mean		61.22 $\pm$ 9.95	56.78 $\pm$ 8.70	<b>72.67</b> $\pm$ 13.91	69.67 $\pm$ 15.07	60.33 $\pm$ 13.77	64.00 $\pm$ 8.40	66.44 $\pm$ 8.20	<b>60.89</b> $\pm$ 12.60	<b>53.56</b> $\pm$ 4.36	62.84 $\pm$ 05.71
Logit $E^\ell$ Max		53.47 $\pm$ 2.13	49.02 $\pm$ 1.79	48.27 $\pm$ 1.32	57.38 $\pm$ 6.09	91.76 $\pm$ 0.91	57.42 $\pm$ 1.43	52.77 $\pm$ 2.58	50.74 $\pm$ 1.51	51.17 $\pm$ 1.83	56.89 $\pm$ 12.70
Marginal $E^m$ Max		78.00 $\pm$ 1.30	76.90 $\pm$ 1.09	48.29 $\pm$ 1.16	68.77 $\pm$ 8.33	10.93 $\pm$ 1.42	<b>80.70</b> $\pm$ 1.98	67.49 $\pm$ 1.69	51.91 $\pm$ 2.32	51.28 $\pm$ 2.47	59.36 $\pm$ 20.69
Marginal $E^m$ Min		58.39 $\pm$ 2.79	59.20 $\pm$ 1.95	51.71 $\pm$ 1.16	34.13 $\pm$ 8.78	97.42 $\pm$ 0.51	50.37 $\pm$ 2.43	69.88 $\pm$ 1.40	49.05 $\pm$ 2.20	49.00 $\pm$ 2.30	57.68 $\pm$ 16.75
Spilled $\Delta E_s$ Min		77.75 $\pm$ 1.52	79.44 $\pm$ 2.05	43.39 $\pm$ 1.82	72.87 $\pm$ 6.10	<b>99.97</b> $\pm$ 0.08	61.56 $\pm$ 2.95	77.55 $\pm$ 1.62	52.34 $\pm$ 2.57	48.17 $\pm$ 1.62	68.12 $\pm$ 17.15
Spilled $\Delta E$ Min		<b>79.04</b> $\pm$ 1.78	<b>80.83</b> $\pm$ 1.87	43.22 $\pm$ 1.67	<b>74.36</b> $\pm$ 5.54	<b>99.97</b> $\pm$ 0.08	61.97 $\pm$ 2.81	<b>78.54</b> $\pm$ 1.57	52.11 $\pm$ 2.58	48.21 $\pm$ 1.62	<b>68.69</b> $\pm$ 17.48

# Key Results & Findings

73.2%

Avg AuROC  
LLaMA-3-Instruct

77.5%

Avg AuROC  
Mistral-Instruct

+24%

Boost from  
exact answer tokens

0

Additional  
training needed



## Instruction Tuning Helps

While fine-tuning degrades logit confidence (54.6%), spilled energy improves with instruction tuning (68.7% → 73.2% for LLaMA-3).



## Known Limitation

False positives occur on semantically uninformative tokens (punctuation, sentence-initial words) — confirming the need for exact answer localization.



## Generalizes Across Models

Results hold for LLaMA-3, Mistral, Gemma 1B & 4B — both pretrained and instruction-tuned — without any architecture-specific modifications.



## Min Pooling is Best

Across all methods and models, min pooling over the exact answer span consistently outperforms max, mean, and last-token strategies.

# Conclusion

*Spilled Energy in Large Language Models — ICLR 2026*

## 01 New EBM Perspective

---

The LLM softmax is reinterpreted as an Energy-Based Model, revealing a principled connection between autoregressive generation and energy dynamics.

## 02 Training-Free Detection

---

Spilled energy requires no additional training or probing classifiers — it reads directly from output logits, making it instantly deployable.

## 03 Superior Generalization

---

Consistently outperforms probe classifiers under cross-dataset evaluation (73% vs 62%) and surpasses logit confidence across all tested LLMs.

## 04 Future Direction

---

Open questions: real-time deployment at scale, deeper understanding of false positives, and extension to multi-modal and longer-context generation.