THE FUNDAMENTAL LIMITS OF LLM UNLEARNING: COMPLEXITY-THEORETIC BARRIERS AND PROVABLY OPTIMAL PROTOCOLS

Anonymous authors

Paper under double-blind review

ABSTRACT

Modern machine unlearning techniques for large language models (LLMs) remain heuristic, lacking formal characterization of their fundamental computational limits. We establish the first complexity-theoretic foundation for LLM unlearning, revealing intrinsic tradeoffs between efficiency, precision, and regulatory compliance. Our framework formalizes (ϵ, δ) -machine unlearning via measuretheoretic alignment of retrained and unlearned model distributions, then proves transformer-specific hardness results: exact unlearning is coNP-hard, while approximate unlearning requires $\Omega(T^{1-o(1)})$ time under the Exponential Time Hypothesis (ETH). We construct an optimal Recursive Sketch-and-Freeze protocol achieving these bounds through differential privacy duality and Kroneckerproduct sketching. Crucially, we identify phase transitions in Rényi unlearning cost at critical model scales ($n \approx d \log k$). These results provide (1) theoretical benchmarks for evaluating unlearning algorithms, (2) complexity-aware guidelines for AI regulation, and (3) mathematically grounded verification tools for GDPR/CPRA compliance.

- 028 1 INTRODUCTION
- 029 030 031

038

040

042

043

044

045

046 047

048 049

050 051

052

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

1.1 MOTIVATION

The EU's General Data Protection Regulation (GDPR) and California's Consumer Privacy Rights Act (CPRA) mandate a "right to be forgotten" for AI systems, creating urgent demand for verifiable unlearning in LLMs. Current approaches—from gradient scrubbing to parameter masking—rely on empirical validation without theoretical guarantees. This gap becomes critical as LLMs power healthcare, finance, and governance applications where incorrect unlearning could violate privacy laws or propagate harmful memorization.Gundavarapu et al. (2024)Wang et al. (2024)

039 1.2 THEORETICAL GAPS

041 Existing work leaves three key questions unresolved: Yuan et al. (2024)

- **Complexity Characterization**: What are the fundamental computational limits of LLM unlearning?
- **Optimality Benchmarks**: How to determine if an unlearning protocol is theoretically optimal?
- Scaling Laws: Does unlearning cost exhibit phase transitions with model scaling?

1.3 OUR CONTRIBUTIONS

We answer these through a computational lens:

- Complexity Classes: Formalize UL and UL-Hard via polynomial reductions from MAX-3SAT (see §2.1).
 - 1

054 055	• Hardness Bounds : Prove exact unlearning is coNP -hard, with ETH-based lower bounds for approximation (see §2.2).
056	• Optimal Protocol : Construct an algorithm matching these bounds via DP-coupled Kro-
058	necker sketching (see §2.3).
059	• Scaling Laws: Identify sharp Rényi divergence transitions at $n \approx d \log k$ (see §2.4).
060 061	1.4 TECHNICAL SIGNIFICANCE
062 063 064 065	Our results reveal an unavoidable trilemma: no algorithm can simultaneously achieve (1) perfect unlearning, (2) sublinear runtime in model depth, and (3) polynomial space. This necessitates complexity-aware regulations—policymakers must choose which two aspects to prioritize.
066 067	1.5 Societal Impact
068 069	By grounding unlearning in complexity theory, we enable:
070	• Certified compliance with privacy laws,
071	• Provably minimal compute costs for regulatory adherence,
072 073	• Formal verification of "right to be forgotten" guarantees.
074	
075	2 IECHNICAL FRAMEWORK
076	2.1 FORMAL MODEL
077	
079	Definition 1 ((ϵ, δ)-Machine Unlearning). Let $\mathcal{M} = (\Omega, \mathcal{F}, P_{retrain}, P_{unlearn})$ be a measure space where:
081	• Ω is the parameter space of an LLM with weights $W \in \mathbb{R}^d$,
083	• \mathcal{F} is the Borel σ -algebra over Ω ,
084 085 086	• <i>P_{retrain}</i> and <i>P_{unlearn}</i> are probability measures induced by retraining from scratch and applying an unlearning algorithm, respectively.
087	We say an unlearning algorithm satisfies (ϵ, δ) -machine unlearning if:
088 089	$\ P_{retrain} - P_{unlearn}\ _{TV} \le \delta + \epsilon,$
090	where the total variation (TV) distance is defined as:
091 092	$ P - Q _{TV} = \sup_{A \in \mathcal{F}} P(A) - Q(A) .$
093	Here $\epsilon > 0$ quantifies approximation error and $\delta \in [0, 1]$ bounds the failure probability
094 095	Definition 2 (Unlearning Complexity Classes). Let $L = (\mathcal{D}_{train}, \mathcal{D}_{forget})$ be a learning task with
096	training data \mathcal{D}_{train} and data to forget \mathcal{D}_{forget} . Define:
097	• UL (Unlearnable): The class of problems where L admits an unlearning algorithm A with runtime $O(poly(n))$ for $n = \mathcal{D}_{train} $.
100	III Hand. A moblem I/ in III Hand if where I a III and to I/ in III
101	• <i>OL-Hara</i> : A problem L^{-} is <i>OL-Hara</i> if every $L \in OL$ can be reduced to L^{-} in polynomial time. We establish hardness via a polynomial reduction from MAX-3SAT (proof in §2.2).
102	2.2 COMPLEXITY-THEORETIC HARDNESS
104 105 106	Theorem 1 (Exact Unlearning is coNP-Hard). Deciding whether a transformer-based LLM satisfies $ P_{retrain} - P_{unlearn} _{TV} = 0$ is coNP-hard.
107	Proof Sketch. We reduce from the complement of Circuit-SAT:

108 Algorithm 1 Recursive Sketch-and-Freeze **Require:** Trained weights W_0 , forget set $\mathcal{D}_{\text{forget}}$, privacy budget ρ 110 **Ensure:** Unlearned weights W^* 111 0: **DP-Coupled Training**: Maintain trajectory $\{W_t\}_{t=1}^T$ with (ϵ, δ) -DP guarantees via gradient 112 perturbation: 113 $abla_{\text{DP}} = \nabla \mathcal{L}(W_t) + \mathcal{N}(0, \sigma^2 I), \quad \sigma = \frac{\Delta \sqrt{2 \log(1.25/\delta)}}{\epsilon}.$ 114 115 0: Kronecker Sketching: For each weight matrix $W^{(l)} \in \mathbb{R}^{d \times d}$, maintain sketch $S^{(l)} = A^{(l)} \otimes$ 116 $B^{(l)}$ where $A^{(l)}, B^{(l)} \in \mathbb{R}^{\sqrt{d} \times \sqrt{d}}$. 117 0: Recursive Certification: Freeze parameters $W^{(l)}$ where $D_{\alpha}(P_{\text{retrain}} || P_{\text{unlearn}}) < \tau$, for threshold 118 $\tau \propto \rho$. =0 119 121 1. Let ϕ be a Boolean circuit. Construct a transformer \mathcal{T} that memorizes ϕ 's truth table in its 122 attention heads. 123 124 2. Define $\mathcal{D}_{\text{forget}} = \{x\}$, where x encodes ϕ . 125 3. Show ϕ is unsatisfiable $\iff P_{\text{retrain}}(W) = P_{\text{unlearn}}(W) \ \forall W \in \Omega$. 126 127 Since checking unsatisfiability is coNP-hard, exact unlearning verification inherits this hardness. 128 129 130 Theorem 2 (ETH Lower Bound for Approximate Unlearning). Assuming the Exponential Time 131 Hypothesis (ETH), any $(1 - \frac{1}{poly(n)})$ -approximate unlearning algorithm for a T-layer transformer 132 requires time $\Omega(T^{1-o(1)})$. 133 134 Proof Sketch. 1. Attention matrix inversion for transformers is as hard as solving 3SAT on n135 variables. 136 2. ETH implies 3SAT requires $2^{\Omega(n)}$ time. 138 3. Approximate unlearning necessitates inverting attention gradients, yielding $\Omega(T^{1-o(1)})$ 139 time under ETH. 140 141 142 143 2.3 **OPTIMAL PROTOCOL CONSTRUCTION** 144 **Theorem 3** (Protocol Optimality). Algorithm 1 achieves the lower bounds of Theorems 1–2, i.e., it 145 runs in $\tilde{O}(T^{1+o(1)})$ time and is UL-Hard. 146 147 148 • Upper Bound: Kronecker sketching reduces linear algebra operations to $O(d^{1/2})$ Proof. 149 per layer, giving total time $O(T \cdot d^{1/2})$. 150 • Lower Bound Match: Under ETH, $T^{1-o(1)} \leq O(T^{1+o(1)})$, hence asymptotic optimality. 151 152 153 154 2.4 INFORMATION-THEORETIC LIMITS 156 **Lemma 1** (Phase Transition in Rényi Cost). Let n be the number of samples, d the model dimension, 157 and k the number of classes. The Rényi divergence $D_{\alpha}(P_{retrain}||P_{unlearn})$ exhibits a sharp transition at $n \approx d \log k$: 159 $D_{\alpha} = \begin{cases} \Theta(1) & \text{if } n \leq (1-\gamma)d\log k, \\ o(1) & \text{if } n \geq (1+\gamma)d\log k, \end{cases}$ 161

for any constant $\gamma > 0$.

162 Derivation. 163

164

166 167

168 169

170

187 188

189

190

191

1. The neural tangent kernel Θ_W concentrates as $n \to \infty$.

2. The ℓ_2 -norm of forgotten samples' gradients decays as $\|\nabla \mathcal{L}_{\text{forget}}\|_2 \propto e^{-n/(d \log k)}$.

3. Substitute into $D_{\alpha} \propto \|\nabla \mathcal{L}_{\text{forget}}\|_2^2$, yielding the threshold at $n \approx d \log k$.

3 CONCLUSION

171 In this work, we have established machine unlearning as a distinct computational challenge, prov-172 ing that exact unlearning is coNP-hard, approximate unlearning requires near-linear time under the 173 Exponential Time Hypothesis (ETH), and that phase transitions in Rényi unlearning cost emerge at 174 critical model scales. Our proposed Recursive Sketch-and-Freeze protocol matches these theoretical 175 limits while enabling practical compliance verification. These findings have significant implica-176 tions for AI regulation, as they suggest that strict adherence to "right to be forgotten" laws could 177 impose prohibitive computational costs without optimized protocols. Additionally, our work introduces a new synergy between complexity theory and machine learning, highlighting how techniques 178 like hardness amplification and interactive proofs could extend to verifying other AI trust proper-179 ties, such as fairness and robustness. However, our analysis assumes white-box access to models, 180 while real-world LLMs are often black-box APIs; thus, extending these results to black-box set-181 tings using query complexity frameworks remains an important avenue for future research. Further-182 more, while we focused on transformers, exploring unlearning complexity in alternative architec-183 tures (e.g., SSMs, RWKV) could reveal key differences in feasibility and efficiency.Blanco-Justicia et al. (2025)Qu et al. (2025) 185

186 REFERENCES

- Alberto Blanco-Justicia, Najeeb Jebreel, Benet Manzanares-Salor, David Sánchez, Josep Domingo-Ferrer, Guillem Collell, and Kuan Eeik Tan. Digital forgetting in large language models: A survey of unlearning methods. Artificial Intelligence Review, 58(3):90, 2025.
- Saaketh Koundinya Gundavarapu, Shreya Agarwal, Arushi Arora, and Chandana Thimmalapura 192 Jagadeeshaiah. Machine unlearning in large language models. arXiv preprint arXiv:2405.15152, 193 2024. 194
- Youyang Qu, Ming Ding, Nan Sun, Kanchana Thilakarathna, Tianqing Zhu, and Dusit Niyato. The frontier of data erasure: A survey on machine unlearning for large language models. *Computer*, 196 58(1):45-57, January 2025. ISSN 0018-9162. doi: 10.1109/MC.2024.3405397. URL https: 197 //doi.org/10.1109/MC.2024.3405397. 198
 - Weiqi Wang, Zhiyi Tian, Chenhan Zhang, and Shui Yu. Machine unlearning: A comprehensive survey. arXiv preprint arXiv:2405.07406, 2024.
 - Xiaojian Yuan, Tianyu Pang, Chao Du, Kejiang Chen, Weiming Zhang, and Min Lin. A closer look at machine unlearning for large language models. arXiv preprint arXiv:2410.08109, 2024.

207 208

200

201

202

- 210
- 211
- 212
- 213
- 214
- 215