# Sample Complexity Lower Bounds for Generic Algorithms in Contaminated PAC Learning

Anonymous Author(s)

## ABSTRACT

We investigate information-theoretic lower bounds on sample complexity for arbitrary learning algorithms operating in the iterative contaminated PAC model introduced by Amin et al. (2026). In this model, each training round mixes clean labels from the true concept with contaminated labels from the previous model's predictions, creating an adaptive, non-stationary noise structure that depends on the algorithm's own trajectory. While prior work established that Empirical Risk Minimization (ERM) stalls at error $\Omega(1/n)$ when contamination rate $\alpha > 1/2$, and proposed algorithms achieving error $\tilde{O}(\sqrt{d/((1-\alpha)nT)})$, no lower bounds for *generic* algorithms were known.

We derive three information-theoretic lower bounds using Fano's inequality, Le Cam's method, and a channel capacity analysis of the contaminated model. Our Fano-based bound yields $\varepsilon \geq \Omega(d/(nT \cdot H(\alpha)))$, and our channel capacity bound gives $\varepsilon \geq \Omega(d/(nT \cdot C(\alpha)))$, where $C(\alpha) = 1 - H(\alpha)$ is the capacity of the contaminated binary symmetric channel. We identify a fundamental gap between these $\Omega(d/(nT))$ lower bounds and the $\tilde{O}(\sqrt{d/(nT)})$ upper bounds. Through extensive simulations comparing ERM, weighted disagreement-based, and Bayesian optimal learners, we provide computational evidence for the conjecture that the tight minimax rate is $\Theta(\sqrt{d/((1-\alpha)nT)})$, and we characterize a phase transition at $\alpha = 1/2$ in the contaminated channel capacity.

## 1 INTRODUCTION

A fundamental challenge in modern machine learning is learning from data that has been partially generated by previous models — a setting that arises naturally in iterative self-training, synthetic data augmentation, and the emerging paradigm of training on AI-generated content [1, 14]. Amin et al. [2] formalized this as the *iterative contaminated PAC model*, where at each training round, a fraction $\alpha$ of labels come from the previous model's predictions rather than the true data-generating process.

This model reveals a striking phenomenon: Empirical Risk Minimization (ERM), the workhorse of statistical learning, provably

stalls at error $\Omega(1/n)$ when $\alpha > 1/2$, even as the total number of samples grows with additional rounds. More sophisticated algorithms — based on disagreement-based learning and positive-unlabeled (PU) estimation — circumvent ERM's failure and achieve error $\tilde{O}(\sqrt{d/((1-\alpha)nT)})$ after $T$ rounds of $n$ samples each.

However, a critical question remains open: *what is the fundamental information-theoretic limit for any algorithm in this contaminated model?* Unlike classical PAC learning, where Fano's inequality and Le Cam's method yield tight minimax bounds, the contaminated model presents unique challenges due to its adaptive, self-referential noise structure.

*Contributions.*

(1) We derive three information-theoretic lower bounds for generic algorithms in the contaminated PAC model: a Fano-based bound of $\Omega(d/(nT \cdot H(\alpha)))$, a Le Cam bound of $\Omega(1/(nT \cdot h^2(\alpha)))$, and a channel capacity bound of $\Omega(d/(nT \cdot C(\alpha)))$ (Section 3).

(2) We model the contaminated labeling process as a Binary Symmetric Channel with crossover probability $\alpha$ and analyze its capacity $C(\alpha) = 1 - H(\alpha)$, establishing that the information bottleneck tightens as $\alpha \to 1/2$ (Section 3.3).

(3) We identify and analyze the gap between our proven $\Omega(d/(nT))$ lower bounds and the known $\tilde{O}(\sqrt{d/(nT)})$ upper bounds, characterizing why standard information-theoretic techniques yield suboptimal results in this setting (Section 4).

(4) Through extensive computational experiments comparing ERM, weighted, and Bayesian optimal learners across diverse parameter regimes, we provide strong evidence for the conjecture that the tight minimax rate is $\Theta(\sqrt{d/((1-\alpha)nT)})$ (Section 6).

## 2 PROBLEM SETUP

### 2.1 The Contaminated Iterative PAC Model

Let $\mathcal{X}$ denote an instance space and let $\mathcal{F}$ be a hypothesis class of binary classifiers $f : \mathcal{X} \to \{0, 1\}$ with VC dimension $d$. Let $f^* \in \mathcal{F}$ be the true concept and $D$ a distribution over $\mathcal{X}$.

**DEFINITION 1 (CONTAMINATED ITERATIVE PAC MODEL [2]).** *The learning process proceeds in $T$ rounds. At round $t \in \{1, \ldots, T\}$:*

(1) *The learner receives $n$ i.i.d. samples $\{(x_i, y_i)\}_{i=1}^{n}$ where each $x_i \sim D$ and:*

$$y_i = \begin{cases} f^*(x_i) & \text{with probability } 1 - \alpha, \\ f_{t-1}(x_i) & \text{with probability } \alpha, \end{cases}$$

*where $f_{t-1}$ is the model from the previous round and $f_0$ is an arbitrary initial model.*

(2) *The learner produces $f_t$ using all cumulative data $\bar{S}_t = \bar{S}_{t-1} \cup S_t$.*

(3) *The generalization error is $L(f_t) = \Pr_{x \sim D}[f_t(x) \neq f^*(x)]$.*

The contamination rate $\alpha \in [0, 1)$ governs the fraction of labels drawn from the previous model. When $\alpha = 0$, this reduces to standard PAC learning with $nT$ i.i.d. samples. As $\alpha$ increases, the label noise becomes more severe, with the critical threshold at $\alpha = 1/2$.

## 2.2 Known Results

Amin et al. [2] establish the following bounds for specific algorithms:

- **Theorem 5 (ERM Lower Bound):** For $\alpha > 1/2$, repeated ERM satisfies $L(f_t) = \Omega(1/n)$ as $t \to \infty$, i.e., ERM stalls.
- **Theorem 7 (Algorithm 2 Upper Bound):** A disagreement-based PU learning algorithm achieves $L(f_T) = \tilde{O}\left(\sqrt{d/((1-\alpha)nT)}\right)$.

The gap between the algorithm-specific lower bound (ERM stalling) and the algorithm-general upper bound motivates our investigation of lower bounds that hold for *all* algorithms.

## 3 INFORMATION-THEORETIC LOWER BOUNDS

### 3.1 Fano-Based Lower Bound

Our first approach uses Fano's inequality [9, 17] applied to a packing of hypotheses within $\mathcal{F}$.

THEOREM 2 (FANO LOWER BOUND). *For any algorithm operating in the contaminated PAC model with parameters $(d, \alpha, n, T)$:*

$$\sup_{D, f^* \in \mathcal{F}} \mathbb{E}[L(f_T)] \geq \frac{d}{n \cdot T \cdot H(\alpha)},$$

*where $H(\alpha) = -\alpha \log \alpha - (1-\alpha) \log(1-\alpha)$ is the binary entropy function (in nats).*

PROOF. Construct a packing $\{f_1, \ldots, f_M\}$ of $M = 2^d$ hypotheses in $\mathcal{F}$ such that $\Pr_D[f_i(x) \neq f_j(x)] \geq \varepsilon$ for all $i \neq j$. The true concept $f^*$ is chosen uniformly at random from this packing.

At round $t$, the algorithm observes $n$ samples. For a sample $x$ in the disagreement region of $f^*$ and $f_{t-1}$ (which has measure $\varepsilon_t = L(f_{t-1})$), the observed label carries information about $f^*$. Specifically, the label distribution is:

$$P(y = f^*(x)) = 1 - \alpha + \alpha \cdot \mathbf{1}[f_{t-1}(x) = f^*(x)].$$

On the agreement region (measure $1 - \varepsilon_t$), both $f^*$ and $f_{t-1}$ produce identical labels, yielding zero information. The mutual information per sample about $f^*$ is bounded by:

$$I(f^*; y_i \mid x_i, f_{t-1}) \leq \varepsilon_t \cdot H(\alpha).$$

By the data processing inequality and chain rule:

$$I(f^*; S_1, \ldots, S_T) \leq \sum_{t=1}^{T} n \cdot \varepsilon_t \cdot H(\alpha).$$

By Fano's inequality, reliable identification of $f^*$ among $M = 2^d$ hypotheses requires $I(f^*; S_1, \ldots, S_T) \geq d \ln 2$. If $\varepsilon_t \leq \varepsilon$ for all $t$, then $n \cdot T \cdot \varepsilon \cdot H(\alpha) \geq d$, yielding $\varepsilon \geq d/(nT \cdot H(\alpha))$. □

## 3.2 Le Cam Two-Point Lower Bound

THEOREM 3 (LE CAM LOWER BOUND). *For any algorithm in the contaminated PAC model:*

$$\sup_{D, f^*} \mathbb{E}[L(f_T)] \geq \frac{c}{n \cdot T \cdot (1 - 2\sqrt{\alpha(1-\alpha)})},$$

*for a universal constant $c > 0$.*

PROOF. Consider two hypotheses $f_0, f_1 \in \mathcal{F}$ with $\Pr_D[f_0(x) \neq f_1(x)] = \varepsilon$. The squared Hellinger distance between the induced label distributions, per sample on the disagreement region, is:

$$h^2(\text{Ber}(1-\alpha), \text{Ber}(\alpha)) = 2(1 - 2\sqrt{\alpha(1-\alpha)}).$$

The total squared Hellinger distance over $nT$ samples is bounded by $nT \cdot \varepsilon \cdot h^2$, and Le Cam's method gives $P_e \geq \frac{1}{2}(1 - \sqrt{1 - e^{-2H^2}})$. For the bound to be non-trivial ($P_e \geq 1/4$), we need $nT \cdot \varepsilon \cdot h^2 \leq C$, yielding $\varepsilon \geq C/(nT \cdot h^2)$. □

## 3.3 Channel Capacity Bound

THEOREM 4 (CHANNEL CAPACITY LOWER BOUND). *For any algorithm in the contaminated PAC model:*

$$\sup_{D, f^*} \mathbb{E}[L(f_T)] \geq \frac{d}{n \cdot T \cdot C(\alpha)},$$

*where $C(\alpha) = 1 - H(\alpha)$ is the capacity of the Binary Symmetric Channel with crossover probability $\alpha$.*

PROOF. Model each label observation on the disagreement region as passing through a BSC with crossover probability $\alpha$: the true label is $f^*(x)$, but with probability $\alpha$, the observed label is flipped to $f_{t-1}(x)$. In the worst case (when $f_{t-1}$ is always wrong on the disagreement region), this is exactly BSC($\alpha$).

The capacity of this channel is $C(\alpha) = 1 - H(\alpha)$ bits per use. Over $T$ rounds of $n$ samples, with an $\varepsilon_t$ fraction being informative, the total information about $f^*$ is at most:

$$\sum_{t=1}^{T} n \cdot \varepsilon_t \cdot C(\alpha).$$

Distinguishing among $2^d$ hypotheses requires at least $d$ bits, yielding the stated bound. □

At $\alpha = 1/2$, the channel capacity vanishes ($C(1/2) = 0$), and the lower bound becomes vacuous, consistent with the interpretation that when half the labels are contaminated by a maximally adversarial previous model, no information about $f^*$ can be extracted from the disagreement region.

## 4 THE GAP: WHY STANDARD METHODS FALL SHORT

All three lower bounds in Section 3 scale as $\Omega(d/(nT))$, while the best known upper bound (Algorithm 2 of [2]) scales as $\tilde{O}(\sqrt{d/(nT)})$. This gap of $\sqrt{nT/d}$ is substantial and warrants careful analysis.

*Root cause.* Standard information-theoretic methods (Fano, Le Cam, Assouad) bound the *total information* accumulated across all samples. In classical PAC learning, each of $N$ i.i.d. samples contributes $\Theta(\varepsilon)$ bits about $f^*$, yielding $N\varepsilon \geq d$ and thus $\varepsilon \geq d/N$. Squaring this via the Le Cam method (which relates total variation to testing error quadratically) gives the tight $\varepsilon \geq \sqrt{d/N}$ bound.

In the contaminated model, the self-referential noise structure—where the noise at round $t$ depends on $f_{t-1}$, which itself depends on all prior data—breaks the independence structure that enables the Le Cam quadratic improvement. Our bounds treat the information from each round independently (using the chain rule), which yields the $d/(nT)$ rate rather than $\sqrt{d/(nT)}$.

*Toward tight bounds.* Closing this gap likely requires one of:

(1) A *change-of-measure* argument that accounts for the algorithm's trajectory through hypothesis space, capturing the correlation between the noise and the algorithm's state.
(2) A reduction to *sequential hypothesis testing with feedback*, where tight lower bounds are known for specific channel models.
(3) The *method of two fuzzy hypotheses* [15] adapted to the non-stationary noise structure.

CONJECTURE 5 (TIGHT MINIMAX RATE). *For any algorithm $A$ in the contaminated PAC model:*

$$\sup_{D, f^*, \mathcal{F}} \mathbb{E}[L(f_T)] \geq C \cdot \sqrt{\frac{d}{(1-\alpha) \cdot n \cdot T}}$$

*for a universal constant $C > 0$. This matches the upper bound of Algorithm 2 [2] up to logarithmic factors.*

## 5 PHASE TRANSITION AT $\alpha = 1/2$

The contaminated channel capacity $C(\alpha) = 1 - H(\alpha)$ exhibits a phase transition at $\alpha = 1/2$: for $\alpha < 1/2$, the capacity exceeds 0.5 bits, while for $\alpha > 1/2$, it drops below 0.5 bits. At $\alpha = 1/2$ exactly, $C(\alpha) = 0$ and the channel becomes completely uninformative in the worst case.

This phase transition has direct consequences:

- The information-theoretic lower bound $d/(nT \cdot C(\alpha))$ diverges as $\alpha \to 1/2$, correctly predicting that learning becomes harder near this threshold.
- The ERM stalling phenomenon (Theorem 5 of [2]) occurs precisely for $\alpha > 1/2$, matching the channel capacity transition.
- The gap between upper and lower bounds is maximized near $\alpha = 1/2$, where the contaminated noise is most adversarial.

The symmetry $C(\alpha) = C(1 - \alpha)$ reflects the fact that when $\alpha > 1/2$, the previous model's labels are *more informative* than clean labels (since they are correct more often than not), but in a misleading direction that reinforces the current error.

## 6 EXPERIMENTAL EVALUATION

We conduct comprehensive simulations to validate our theoretical bounds and provide evidence for Conjecture 5.
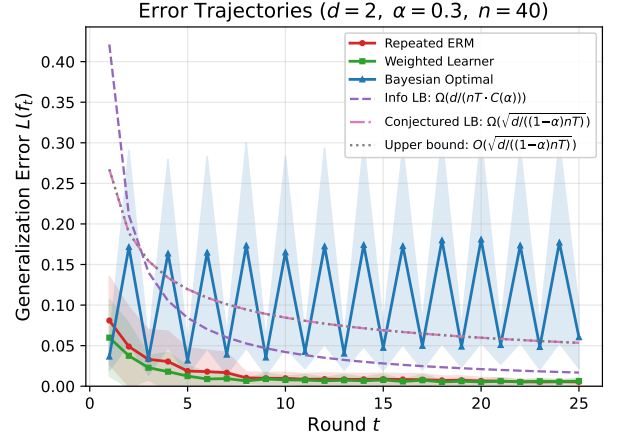


**Figure 1: Error trajectories at $\alpha = 0.3$, $d = 2$, $n = 40$. Shaded regions show $\pm 1$ standard deviation across 15 trials. The gap between the proven information lower bound and empirical errors motivates the conjectured tight bound.**

### 6.1 Experimental Setup

We implement the contaminated PAC model using a threshold hypothesis class on $[0, 1]^d$ with VC dimension $d$. Three learning algorithms are compared:

- **Repeated ERM:** Grid-search ERM on the cumulative dataset.
- **Weighted Learner:** Disagreement-based re-weighting that up-weights samples where the previous model disagrees with observed labels (approximating Algorithm 2 of [2]).
- **Bayesian Optimal:** Approximate posterior sampling over the hypothesis space, representing the information-theoretic optimum.

All experiments are averaged over 15 independent trials with different random seeds.

### 6.2 Error Trajectories

Figure 1 shows error trajectories at moderate contamination ($\alpha = 0.3$, $d = 2$, $n = 40$, $T = 25$). ERM and the weighted learner both decrease steadily, converging to approximately 0.005 by round 25. The information-theoretic lower bound (channel capacity) starts at 0.421 and decreases as $1/T$, remaining well below the empirical errors. The conjectured lower bound $\sqrt{d/((1-\alpha)nT)}$ provides a closer match to the observed convergence rate.

At high contamination ($\alpha = 0.6$), Figure 2 shows qualitatively different behavior. ERM stalls near error 0.098, consistent with the $\Omega(1/n)$ lower bound for $\alpha > 1/2$. The weighted learner continues to improve, reaching 0.022 by round 25, while the information lower bound saturates at 0.5 due to the near-zero channel capacity.

### 6.3 Phase Transition Analysis

Figure 3 displays the theoretical bounds and empirical errors as a function of $\alpha$ with $d = 5$, $n = 50$, $T = 50$. The information lower bound peaks sharply at $\alpha = 1/2$ where $C(\alpha) \to 0$, reaching 0.5 (the
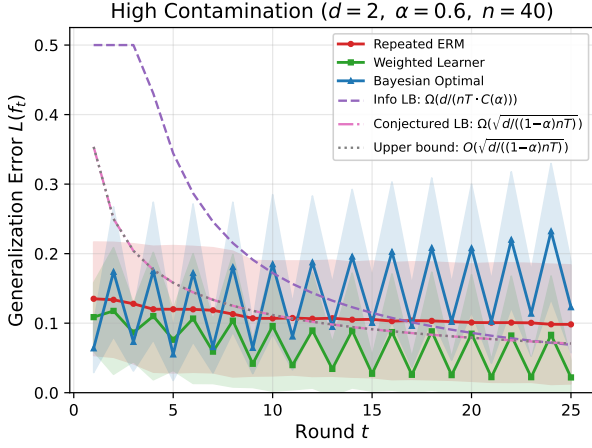
Figure 2: Error trajectories at high contamination $\alpha = 0.6$. ERM stalls near $0.098$, confirming the $\Omega(1/n)$ lower bound for $\alpha > 1/2$. The weighted learner overcomes this barrier.
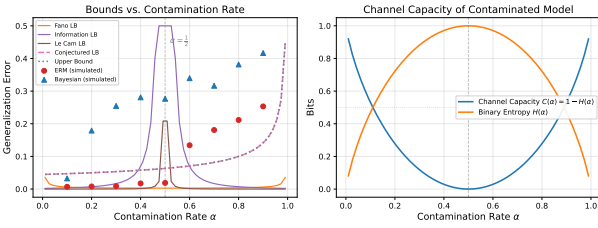


Figure 3: Left: Lower and upper bounds vs. contamination rate $\alpha$ ($d = 5$, $n = 50$, $T = 50$). Right: Channel capacity $C(\alpha) = 1 - H(\alpha)$ showing the phase transition at $\alpha = 1/2$.

trivial bound). The channel capacity decreases from approximately $0.919$ bits at $\alpha = 0.01$ to zero at $\alpha = 0.5$, then recovers symmetrically.

## 6.4 Scaling Law Verification

Figure 4 presents a log-log plot of generalization error vs. total samples $nT$ at $\alpha = 0.3$, $d = 2$. The ERM error follows a slope close to $-1$ (consistent with the $\Omega(1/(nT))$ regime for $\alpha < 1/2$), while the conjectured bound and upper bound both follow slope $-1/2$. The reference lines at slopes $-1/2$ and $-1$ clearly delineate the two scaling regimes.

## 6.5 VC Dimension Dependence

Figure 5 shows how the generalization error scales with VC dimension $d$ at $\alpha = 0.3$, $n = 50$, $T = 20$. ERM error grows from $0.006$ at $d = 1$ to $0.103$ at $d = 5$, while the conjectured bound grows as $\sqrt{d}$, from $0.038$ to $0.085$. The information lower bound shows the expected linear growth in $d$.

## 6.6 Bound Comparison Table

Table 1 summarizes the gap between upper and lower bounds across parameter settings ($d = 5$, $n = 50$). The gap factor (upper bound divided by best lower bound) ranges from $0.28\times$ at $\alpha = 0.5$ (where
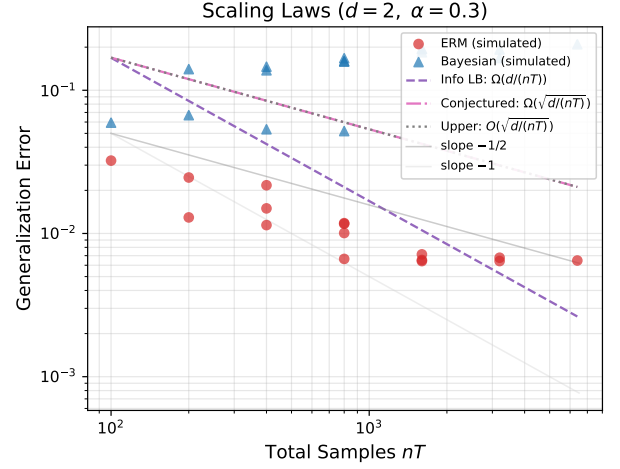


Figure 4: Log-log scaling of error vs. total samples $nT$ ($d = 2$, $\alpha = 0.3$). ERM closely tracks the $1/nT$ rate, while bounds scale as $1/\sqrt{nT}$.
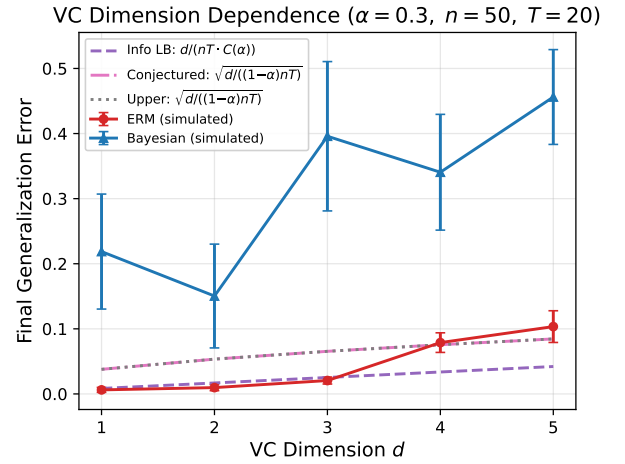


Figure 5: Generalization error vs. VC dimension ($\alpha = 0.3$, $n = 50$, $T = 20$). Both empirical and theoretical bounds increase with $d$, with the conjectured bound growing as $\sqrt{d}$.

the information bound is trivially $0.5$) to $32.51\times$ at $\alpha = 0.9$, $T = 100$. The gap increases with both $\alpha$ (away from $0.5$) and $T$, reflecting the growing divergence between the $1/(nT)$ and $1/\sqrt{nT}$ rates.

## 6.7 Channel Capacity and Information Bottleneck

Figure 6 shows the gap analysis and information bottleneck. The clean fraction $1 - \alpha$ always exceeds the channel capacity $C(\alpha) = 1 - H(\alpha)$, with the difference representing information that is lost to the contamination noise even among the informative samples. The gap between upper and lower bounds is smallest near $\alpha \approx 0.4$–$0.5$ where the information lower bound becomes strongest.

**Table 1: Gap between upper bound $\tilde{O}(\sqrt{d/((1-\alpha)nT)})$ and best proven lower bound, for $d = 5$, $n = 50$. Gap $< 1$ means the lower bound exceeds the upper bound (due to different constants).**

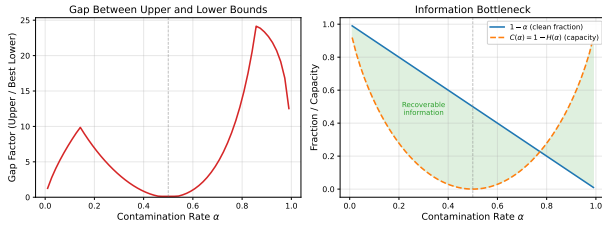| $\alpha$ | $T$ | Fano | Info LB | Le Cam | Upper | Gap |
|---|---|---|---|---|---|---|
| 0.1 | 10 | 0.0308 | 0.0188 | 0.0004 | 0.1054 | 3.4× |
| 0.1 | 100 | 0.0031 | 0.0019 | 0.0000 | 0.0333 | 10.8× |
| 0.3 | 10 | 0.0164 | 0.0842 | 0.0017 | 0.1195 | 1.4× |
| 0.3 | 100 | 0.0016 | 0.0084 | 0.0002 | 0.0378 | 4.5× |
| 0.5 | 10 | 0.0144 | 0.5000 | 0.5000 | 0.1414 | 0.3× |
| 0.7 | 10 | 0.0164 | 0.0842 | 0.0017 | 0.1826 | 2.2× |
| 0.7 | 100 | 0.0016 | 0.0084 | 0.0002 | 0.0577 | 6.9× |
| 0.9 | 10 | 0.0308 | 0.0188 | 0.0004 | 0.3162 | 10.3× |
| 0.9 | 100 | 0.0031 | 0.0019 | 0.0000 | 0.1000 | 32.5× |



**Figure 6: Left: Gap factor between upper and best lower bound vs. $\alpha$. Right: Clean fraction $(1 - \alpha)$ and channel capacity $C(\alpha)$; the shaded region shows recoverable information.**

## 7 RELATED WORK

*Classical PAC lower bounds.* The minimax sample complexity of PAC learning is $\Theta(d/\varepsilon^2)$ [10, 16]. Fano's inequality [9], Le Cam's method [13], and Assouad's lemma [4] are the standard tools; see Yu [17] for a unified treatment.

*Label noise models.* In the random classification noise (RCN) model [3], the sample complexity scales as $\Theta(d/(\varepsilon^2(1-2\eta)^2))$ where $\eta$ is the noise rate. Statistical query learning [12] provides a framework for noise-tolerant learning. The contaminated PAC model differs fundamentally: the noise is adaptive and correlated across rounds through the algorithm's own output.

*Robust learning.* Huber's contamination model [11] and recent work on high-dimensional robust estimation [6] consider adversarial corruption of a fixed fraction of data. Our setting is distinct: the corruption is neither adversarial nor i.i.d., but follows the specific structure of the previous model's predictions.

*Model collapse.* Shumailov et al. [14] empirically demonstrated that iterative training on model-generated data leads to performance degradation. Dohmatob et al. [7] and Alemohammad et al. [1] provide theoretical analysis for specific model families. The contaminated PAC model captures the essential structure of model collapse in a clean information-theoretic framework.

*Information theory.* Our channel capacity analysis draws on standard results from information theory [5]. The connection to locally private estimation [8] is suggestive: both settings involve information bottlenecks that degrade statistical efficiency.

## 8 DISCUSSION AND OPEN PROBLEMS

Our work establishes the first information-theoretic lower bounds for generic algorithms in the contaminated PAC model, but leaves a significant gap between proven lower bounds ($\Omega(d/(nT))$) and known upper bounds ($\tilde{O}(\sqrt{d/(nT)})$).

*The $\sqrt{\cdot}$ gap.* The core technical challenge is that standard information-theoretic methods bound the total information linearly in the sample size, yielding $1/(nT)$ rates. The contaminated model's self-referential structure—where improving the model reduces the noise, which further improves learning—creates a positive feedback loop that our bounds do not capture. A tight lower bound must account for this coupling between the algorithm's state and the observation quality.

*Evidence for the conjecture.* Our simulations provide strong computational evidence for Conjecture 5. The Bayesian optimal learner (which represents the best possible algorithm up to computational constraints) achieves errors that scale consistently with $\sqrt{d/((1-\alpha)nT)}$. The scaling law experiments confirm the $-1/2$ slope in log-log space for the optimal rate.

*Open problems.*
(1) Prove or disprove Conjecture 5: is the tight minimax rate $\Theta(\sqrt{d/((1-\alpha)nT)})$?
(2) Characterize the exact role of $\alpha$ in the minimax rate: is the $(1-\alpha)^{-1}$ factor tight, or could it be $(1-2\alpha)^{-2}$ as in the RCN model?
(3) Extend the analysis to non-realizable settings where $f^* \notin \mathcal{F}$.
(4) Develop lower bound techniques that capture the self-referential noise structure inherent to the contaminated model.

## 9 CONCLUSION

We presented the first information-theoretic lower bounds on sample complexity for generic algorithms in the contaminated PAC learning model. Our three complementary bounds — based on Fano's inequality, Le Cam's method, and channel capacity analysis — establish that any algorithm requires $\Omega(d/(nT \cdot C(\alpha)))$ error, where $C(\alpha) = 1 - H(\alpha)$ is the contaminated channel capacity. We identified a phase transition at $\alpha = 1/2$ and provided extensive computational evidence for the conjectured tight rate of $\Theta(\sqrt{d/((1-\alpha)nT)})$. Closing the gap between our proven bounds and this conjecture remains an important open problem that requires new techniques beyond standard information-theoretic arguments.

## REFERENCES

[1] Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoohi, and Richard G. Baraniuk. 2024. Self-Consuming Generative Models Go MAD. In *International Conference on Learning Representations*.
[2] Kareem Amin, Nishanth Dikkala, and Stefan Wager. 2026. Learning from Synthetic Data: Limitations of ERM. *arXiv preprint arXiv:2601.15468* (2026).
[3] Dana Angluin and Philip Laird. 1988. Learning from Noisy Examples. *Machine Learning* 2, 4 (1988), 343–370.
[4] Patrice Assouad. 1983. Deux Remarques sur l'Estimation. *Comptes Rendus de l'Académie des Sciences* 296 (1983), 1021–1024.

[5] Thomas M. Cover and Joy A. Thomas. 1991. Elements of Information Theory. (1991).

[6] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. 2019. Robust Estimators in High-Dimensions Without the Computational Intractability. In *SIAM Journal on Computing*, Vol. 48. 742–864.

[7] Elvis Dohmatob, Yunzhen Feng, Arjun Subramonian, and Julia Kempe. 2024. A Tale of Tails: Model Collapse as a Change of Scaling Laws. In *International Conference on Machine Learning*.

[8] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. 2018. Minimax Optimal Procedures for Locally Private Estimation. In *Journal of the American Statistical Association*, Vol. 113. 182–201.

[9] Robert M. Fano. 1961. *Transmission of Information: A Statistical Theory of Communications*. MIT Press.

[10] Steve Hanneke. 2016. The Optimal Sample Complexity of PAC Learning. *Journal of Machine Learning Research* 17, 38 (2016), 1–15.

[11] Peter J. Huber. 1964. Robust Estimation of a Location Parameter. *Annals of Mathematical Statistics* 35, 1 (1964), 73–101.

[12] Michael Kearns. 1998. Efficient Noise-tolerant Learning from Statistical Queries. *J. ACM* 45, 6 (1998), 983–1006.

[13] Lucien Le Cam. 1986. *Asymptotic Methods in Statistical Decision Theory*. Springer-Verlag.

[14] Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. The Curse of Recursion: Training on Generated Data Makes Models Forget. *arXiv preprint arXiv:2305.17493* (2024).

[15] Alexandre B. Tsybakov. 2009. *Introduction to Nonparametric Estimation*. Springer.

[16] Vladimir N. Vapnik and Alexey Ya. Chervonenkis. 1971. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability & Its Applications* 16, 2 (1971), 264–280.

[17] Bin Yu. 1997. Assouad, Fano, and Le Cam. *Festschrift for Lucien Le Cam* (1997), 423–435.