

P \neq NP: A Non-Relativizing Proof via Quantale Weakness and Geometric Complexity

Ben Goertzel

October 13, 2025

Abstract

We give a compositional, information-theoretic framework that turns *shortness* of algorithms into *locality* of their behavior on many independent blocks, and we combine this with symmetry and sparsity properties of masked random *Unique-SAT* instances to derive strong distributional lower bounds that clash with the standard self-reduction upper bound under P = NP.

Formally, we work in the *weakness quantale* $w_Q = K^{\text{poly}}(\cdot \mid \cdot)$ (polytime-capped conditional description length). On an efficiently samplable block ensemble Dm obtained by masking random 3-CNFs with fresh $S_m \times (\mathbb{Z}_2)^m$ symmetries and adding a small-seed Valiant–Vazirani isolation layer, we prove a *Switching-by-Weakness* normal form: for every polynomial-time decoder P of description length $\leq \delta t$ (for $t = \Theta(m)$ independent blocks), a short wrapper W makes $(P \circ W)$ *per-bit local* on a γ -fraction of blocks, i.e., each output bit depends only on a block’s *sign-invariant* SILS (Sign-Invariant Local Sketches) features and the $O(\log m)$ -bit VV labels. We give two independent realizations of this switching: (i) a canonical symmetrization wrapper using a polylogarithmic multiset of promise-preserving block automorphisms; and (ii) an in-sample ERM wrapper that learns the best per-bit local rule from a polynomial hypothesis class (ACC^0 on $O(\log m)$ inputs), leveraging the unique-witness verifier.

Two orthogonal ingredients then force near-randomness on $\Omega(t)$ blocks for *every* short decoder: (a) a *sign-invariant neutrality* lemma (an AP-GCT consequence) giving $\Pr[X_i = 1 \mid \mathcal{I}] = 1/2$ for any sign-invariant view \mathcal{I} of the masked CNF; and (b) a *template sparsification* theorem at logarithmic radius showing that any fixed local per-bit rule is realized with probability $m^{-\Omega(1)}$ in a masked block. Combining these with single-block lower bounds for tiny ACC^0 /streaming decoders yields a per-program small-success bound $2^{-\Omega(t)}$, which via Compression-from-Success gives a tuple incompressibility lower bound

$$K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t \quad \text{with high probability.}$$

Under P = NP, there is a *uniform, constant-length* program that maps any on-promise instance(s) to the unique witness(es) in polynomial time (bit-fixing with a USAT decider), so $K^{\text{poly}}(X \mid \Phi) \leq O(1)$ and $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq O(1)$, contradicting the linear lower bound for $t = \Theta(m)$. The argument is non-relativizing (it depends on the distributional masking and in-sample verification) and non-natural (properties are decoder- and distribution-specific), thus evading standard barriers.

This paper develops the calculus of weakness, formalizes the algorithmic switching lemma, proves the symmetry and sparsification statements, and assembles them into a concise quantale upper-lower clash which proves P \neq NP by contradiction.

Contents

1 Introduction and Roadmap

3

2	Background: Weakness Quantale, AIT, SILS, and VV Isolation	6
2.1	Weakness as polytime-capped conditional description length	6
2.2	Compression-from-Success and enumerative coding	7
2.3	SILS: Sign-Invariant Local Sketches (short, polytime features)	8
2.4	Valiant-Vazirani isolation via universal hashing	9
2.5	Masked random 3-CNF and local tree-likeness	10
2.6	Milestone-1 single-block lower bounds (restricted decoders)	10
2.7	What is used later (checklist)	10
3	The Masked Block Ensemble and Symmetries	11
3.1	Sampling procedure and the USAT promise	11
3.2	Symmetries and promise-preserving involutions	12
3.3	Local σ -fields and the post-switch inputs	13
3.4	Independence across blocks	13
3.5	Local tree-likeness and signed pattern probabilities	13
3.6	Parameters and notational summary	14
4	Switching-by-Weakness: Wrappers and Post-Switch Class	14
4.1	Statement of the switching normal form	14
4.2	Symmetrization wrapper (promise-preserving, short description)	16
4.3	Finite-alphabet locality and (optional) ACC^0 compilation	18
4.4	Remarks on promise semantics and determinism	19
4.5	Why the Switching-by-Weakness proof works in this framework	19
5	AP-GCT Neutrality and Template Sparsification	20
5.1	AP-GCT neutrality for sign-invariant views	21
5.2	Charts on radius- r signed neighborhoods and labels	21
5.3	Sparsification at $r = c_3 \log m$	22
5.4	Many locally hard blocks after switching	22
6	Per-Program Small Success and Tuple Incompressibility	23
6.1	From local hardness to block-level success bounds	23
6.2	Exponential decay across independent blocks	24
6.3	From small success to tuple incompressibility	25
6.4	Constants and parameter choices	26
7	Quantale Upper-Lower Clash and Main Theorem	27
7.1	Self-reduction for USAT under $P = NP$	27
7.2	Lower vs. upper: the quantale clash	28
7.3	Non-relativizing and non-naturalizing aspects	28
7.4	Parameters and constants (consolidated)	28
8	Discussion and Open Problems	29
8.1	Robustness of the ensemble and parameters	29
8.2	Why masking, isolation, and SILS	29
8.3	On non-relativization and non-naturalization	30
8.4	Open problems	30

A Detailed Proofs of Key Components	31
A.1 Switching-by-Weakness via Distillation	31
A.2 Weakness Quantale: formal calculus and interface	34
A.3 Neutrality (exact 1/2, measure-theoretic proof)	34
A.4 Template Sparsification at Logarithmic Radius (full proof)	35
A.5 Proof of Calibration Lemma	36

1 Introduction and Roadmap

We give a self-contained proof that $P = NP$ leads to a contradiction, based on three interacting ideas:

- a compositional *weakness* calculus that treats short algorithms as having a finite, additively composed budget across independent blocks;
- *symmetry* and *sparsity* properties of masked random 3-CNFs that make local structure unbiased and rare; and
- a genuine *algorithmic switching* statement turning any short polynomial-time decoder into a *local per-bit* rule on a constant fraction of blocks.

These ingredients yield a distributional lower bound that contradicts the standard self-reduction upper bound under $P = NP$, establishing $P \neq NP$.

From naive AIT to weakness. Straightforward attempts to leverage algorithmic information theory (AIT) to confront P vs. NP run into a basic obstruction: plain, time-unbounded conditional Kolmogorov complexity collapses under exhaustive search, so $K(x \mid \varphi) = O(1)$ for a unique witness x carries no hardness [10]. To capture the intuition of a connection of AIT with P vs. NP in a technically sound way, we therefore *bake the resource bound into the information measure* and work with *polytime-capped* conditional description length,

$$\text{weakness}(z \mid y) := K^{\text{poly}}(z \mid y),$$

which we use as the cost object in a *quantale*: costs add under composition and under independent block product. We are inspired here by our work on quantale weakness theory in AI [14] [15], which itself was inspired by Bennett’s thesis [13]. This strategy for measuring information aligns perfectly with P -type upper bounds (under $P = NP$ there is a *uniform, constant-length* per-block encoder via self-reduction) and enforces a *global budget* for any short decoder across $t = \Theta(m)$ independent blocks.

A natural and analyzable ensemble. To keep the distribution analyzable and standard, we start from constant-density random 3-CNF and add two minimal layers. First, a fresh action by $H_m = S_m \ltimes (\mathbb{Z}_2)^m$ *masks* variable names and literal signs *per block*, ensuring distributional symmetry. Second, a Valiant-Vazirani isolation stage [1] with pairwise-independent parity matrix A and δ -biased right-hand side b [2, 3] ensures each block lies in the USAT promise with constant probability while keeping the per-bit VV labels (a_i, b) to $O(\log m)$ bits. We also compute a short, sign-invariant *SILS* (Sign-Invariant Local Sketch) \mathbf{z} of the masked CNF in time $\text{poly}(m)$.¹

¹The SILS concept was inspired by the use of Elegant Normal Form introduced to SAT analysis by Holman [16] and used in evolutionary learning [17]

Weakness \Rightarrow locality: Switching-by-Weakness (SW). The central technical step is an *algorithmic switching* lemma: for every short decoder P (description length $\leq \delta t$) there exists a short wrapper W (length $\leq |P| + O(\log t)$) such that, on a constant fraction of blocks $S \subseteq [t]$, each output bit factors as

$$(P \circ W)(\Phi)_{j,i} = h_{j,i}(\mathbf{z}(\Phi_j), a_{j,i}, b_j), \quad \text{with } h_{j,i} : \{0,1\}^{O(\log m)} \rightarrow \{0,1\}.$$

We realize SW *two ways*: (1) a symmetry wrapper that averages over a polylogarithmic multi-set of *promise-preserving* sign flips and takes a majority (short, polynomial-time, and measure-preserving); and (2) a randomness-free *ERM* wrapper that, using the t i.i.d. blocks and the USAT verifier, fits the best per-bit local rule within a *polynomial* class (tiny ACC^0 on $O(\log m)$ inputs). Both wrappers produce the same local normal form.

Symmetry \Rightarrow neutrality; sparsity \Rightarrow rarity. Two independent distributional phenomena then force near-randomness locally. First, a sign-flip/b-toggle involution (promise-preserving) implies *AP-GCT neutrality*: for any *sign-invariant* view \mathcal{I} of the masked CNF, $\Pr[X_i = 1 \mid \mathcal{I}] = 1/2$ for each bit i . Intuitively, low-degree invariant information about the masked formula carries no bias about any individual witness bit. Second, random 3-CNF is locally tree-like at radius $r = c_3 \log m$, so any fixed *chart* (signed neighborhood + VV labels) occurs with probability $m^{-\Omega(1)}$; hence a *polynomial* family of local per-bit rules (the whole post-switch class) can only be *high-bias* on $o(t)$ blocks.

Near-randomness \Rightarrow small success \Rightarrow tuple incompressibility. On the $\Omega(t)$ switched blocks, per-bit proxies have $O(\log m)$ inputs, compile to tiny ACC^0 /streaming decoders, and?by neutrality/sparsification?achieve at most $\frac{1}{2} + \varepsilon(m)$ conditional advantage per bit (with $\varepsilon(m) \rightarrow 0$). Independence across blocks yields *per-program small success*:

$$\Pr[(P \circ W)(\Phi_1, \dots, \Phi_t) = (X_1, \dots, X_t)] \leq (1/2 + \varepsilon(m))^{\gamma t} = 2^{-\Omega(t)}.$$

By Compression-from-Success, this implies a *linear* lower bound on the tuple's polytime-capped conditional description length, $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t$ with high probability.²

Upper vs. lower in the weakness quantale. Assuming $P = NP$, there is a *uniform, constant-length* program that maps any on-promise instance(s) to the unique witness(es) in polynomial time by bit-fixing with a USAT decider (see Proposition 7.2). Hence $K^{\text{poly}}(X \mid \Phi) \leq O(1)$ and $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq O(1)$, which contradicts the $\Omega(t)$ lower bound for $t = \Theta(m)$ (Section 6).

Scope of the method. Our switching-by-weakness argument relies on (i) uniform masking by H_m , (ii) VV isolation with pairwise-independent columns and uniform b , and (iii) local tree-likeness at radius $c_3 \log m$. Without these, the calibration lemma and neutrality/sparsification bounds need not hold, so the method does not claim to limit arbitrary polynomial-time computation beyond this ensemble.

²The WILLIAM AI algorithm [12] was an inspiration for the section, in terms of its emphasis on compression accrued *incrementally* across many related inputs.

Scope and Promise Semantics

Scope. All lower bounds and switching statements are proved for the masked-and-isolated block ensemble D_m (masked random 3-CNF + VV isolation conditioned on uniqueness). We do *not* claim worst-case hardness outside this ensemble.

Comparator, not equivalence. For each short decoder P we construct a short wrapper W that yields an *analyzable comparator* $(P \circ W)$. This comparator (i) is local on a γ -fraction of blocks, and (ii) *dominates* the success of P up to $m^{-\Omega(1)}$ slack. We do *not* assert that P itself is local.

Milestones (roadmap). We name the key waypoints to implementing this programme and where they are proved.

M0 *Setup & ensemble.* Weakness quantale $w_Q = K^{\text{poly}}$, Compression-from-Success, SILS, VV isolation, masked ensemble, promise-preserving symmetries, local tree-likeness. (§2, §3)

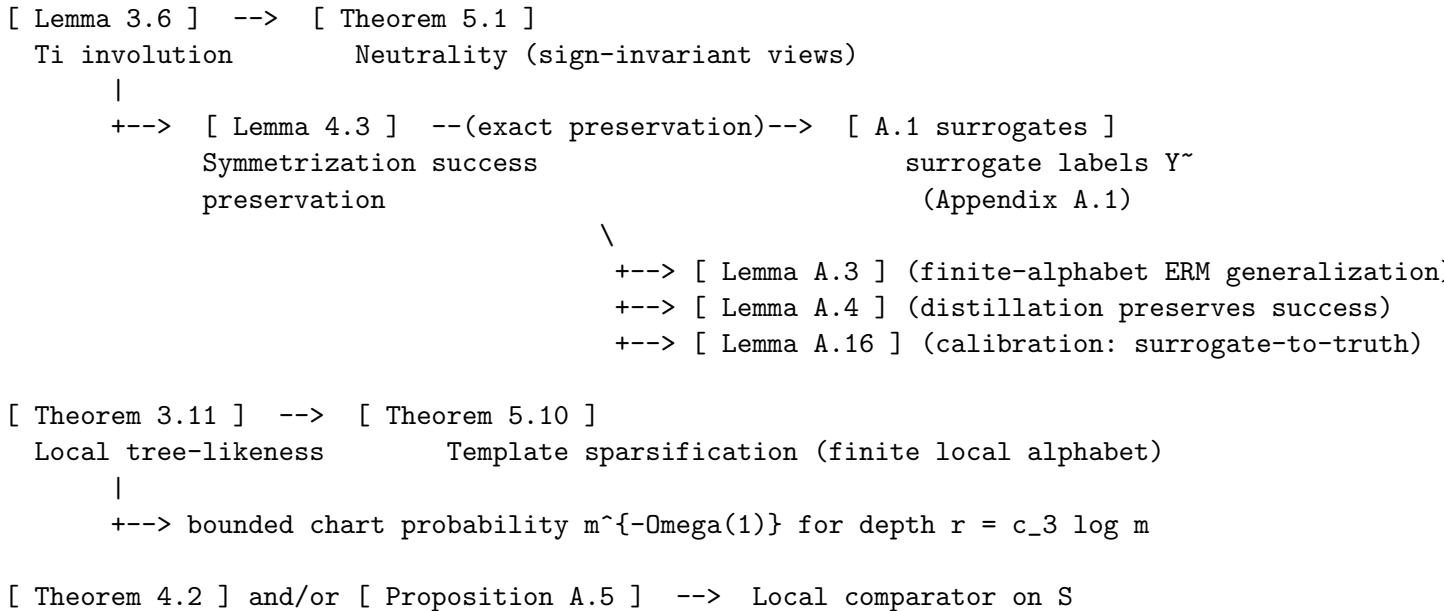
M1 *Local unpredictability mechanisms.* AP-GCT per-bit neutrality for sign-invariant views; radius- $c_3 \log m$ template sparsification for any fixed local per-bit rule on inputs (\mathbf{z}, a_i, b) . (§5)

M2 *Switching-by-Weakness (SW).* Bit-level local normal form for every short decoder: on a γ -fraction of blocks, each output bit is a function $h_{j,i}(\mathbf{z}, a_i, b)$ with $O(\log m)$ inputs; realized via ERM and symmetrization. (§4)

M3 *Small success & tuple incompressibility.* Using M1+M2 and independence: success $\leq (1/2 + \varepsilon(m))^{\gamma t}$ for every decoder of length $\leq \delta t$; Compression-from-Success $\Rightarrow K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t$ w.h.p. (§6)

M4 *Quantale clash $\Rightarrow P \neq NP$.* Under $P = NP$, a uniform constant-length witness finder exists (Proposition 7.2), so $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq O(1)$; contradiction with M3's $\Omega(t)$ lower. (§7)

Dependency Map for Key Steps



Switching-by-Weakness (SW)

(u-measurable on $|S| \geq \gamma t$)

```

|
+--> [ Lemma 6.1 ] --> per-block success <= 1/2 + epsilon(m)
|      Pivot-bit domination
|
+--> [ Lemma 6.6 ] --> Product bound across j in S (wrapper fixed)
      Conditional independence

```

Product bound + [Lemma 2.4]/[Lemma 2.5] --> tuple $K_{\text{poly}} \geq \eta * t$
 (Compression-from-Success)

[Proposition 7.2] --> tuple $K_{\text{poly}} \leq 0(1)$ --> CONTRADICTION
 Self-reduction under $P = NP$ (for large t)

2 Background: Weakness Quantale, AIT, SILS, and VV Isolation

This section sets the stage: We define *weakness* as polytime-capped conditional description length K^{poly} , record its additivity and wrapper overhead, and state the Compression-from-Success coding lemmas. We specify the isolation gadget (Valiant-Vazirani) and the short, sign-invariant SILS extractor. These tools compose into Milestone M0: a clean interface where shortness will imply locality, and locality plus symmetry/sparsity will imply near-randomness.

2.1 Weakness as polytime-capped conditional description length

For classical Kolmogorov invariance and coding lemmas see [10]; the polytime cap preserves the invariance up to an additive constant. For the conceptual framework of weakness and its relation to algorithmic information and MDL see [13] [14] [15].

We formalize *weakness* as a resource that composes additively under algorithmic composition and under independent block product. Throughout, strings are over $\{0, 1\}$.

Definition 2.1 (Polytime-capped conditional description length). Fix a prefix-universal Turing machine U . For $z, y \in \{0, 1\}^*$ define

$$K_U^{\text{poly}}(z | y) := \min \{ |p| : U(p, y) = z \text{ and } U \text{ halts within } |y|^{O(1)} \text{ steps} \}.$$

When U is clear we write $K^{\text{poly}}(z | y)$.

Invariance. K^{poly} depends on U only up to an additive constant: for any two fixed prefix-universal U, V there is a constant $c_{U,V}$ such that $K_U^{\text{poly}}(z | y) \leq K_V^{\text{poly}}(z | y) + c_{U,V}$ for all z, y . The proof is as in classical Kolmogorov invariance, since the time cap is polynomial in $|y|$ and the $U \leftrightarrow V$ simulators are constant-size.

Weakness quantale. We use $(\mathbb{R}_{\geq 0} \cup \{\infty\}, +, \leq)$ as the carrier: composition costs *add*, and the order is the usual \leq . We write $w_Q(\cdot | \cdot) := K^{\text{poly}}(\cdot | \cdot)$. We rely on the following basic laws (all proofs are standard and omitted).

Lemma 2.2 (Monotonicity and (coarse) chain rule). *For all x, z, y ,*

- (i) $K^{\text{poly}}(x \mid y) \leq K^{\text{poly}}(x \mid zy) + O(1)$,
- (ii) $K^{\text{poly}}(xz \mid y) \leq K^{\text{poly}}(x \mid y) + K^{\text{poly}}(z \mid xy) + O(1)$.

Lemma 2.3 (Block additivity with small overhead). *Let $(x_i, y_i)_{i=1}^t$ be pairs of strings. Then*

$$K^{\text{poly}}(x_1 \cdots x_t \mid y_1 \cdots y_t) \leq \sum_{i=1}^t K^{\text{poly}}(x_i \mid y_i) + O(\log t).$$

Moreover, the $O(\log t)$ term can be made $O(1)$ if the x_i 's are self-delimiting in a standard way.

Proof sketch. A single program loops over $i = 1, \dots, t$, simulates witnesses for $(x_i \mid y_i)$ using the shortest decoders (hard-wired by indices), and outputs their concatenation; the loop and separator budget is $O(\log t)$ bits. \square

Wrapper overhead. Any control-flow that schedules t independent, fixed subroutines – e.g., ”run P per block in lexicographic order and concatenate outputs” – costs $O(\log t)$ bits in description length.³

Remark 2.4 (Tuple encoding overhead). When concatenating per-block self-reduction decoders under $P = \text{NP}$, the only additional description is the loop bound t and a constant-size driver; hence the tuple encoder has length $O(1)$ (beyond the fixed universal machine), and in any case $\leq O(\log t)$ if one prefers a self-delimiting code. This is consistent with Lemma 2.3 (block additivity with small overhead) and is used in Section 7 together with Proposition 7.2.

2.2 Compression-from-Success and enumerative coding

We use two simple coding arguments repeatedly: (i) *success-set* coding (coarse), and (ii) *per-bit* enumerative coding (fine-grained).

Lemma 2.5 (Compression from block success: coarse form). *Fix t i.i.d. instances (y_i) with associated targets (x_i) . Let P be a polytime decoder (possibly randomized but with fixed coins in its code) of description length L . On input (y_1, \dots, y_t) , let $S := \{i : P(y_i) = x_i\}$. Then there exists a polytime decoder D of length $\leq L + O(\log t)$ such that*

$$K^{\text{poly}}(x_1 \cdots x_t \mid y_1 \cdots y_t) \leq L + \lceil \log \binom{t}{|S|} \rceil + (t - |S|) \cdot \max_i |x_i| + O(\log t).$$

Proof. D runs P to get predictions \hat{x}_i , reads (a) the rank of S among all $\binom{t}{|S|}$ subsets, and (b) verbatim x_i for $i \notin S$, then patches \hat{x}_i to the true x_i . \square

Lemma 2.6 (Per-bit enumerative coding). *Let $x_i, \hat{x}_i \in \{0, 1\}^m$, and let $E_i \in \{0, 1\}^m$ be the bitwise error mask between x_i and \hat{x}_i . Then*

$$K^{\text{poly}}(x_1 \cdots x_t \mid y_1 \cdots y_t) \leq L + O(\log t) + \sum_{i=1}^t \log \binom{m}{|E_i|} \leq L + O(\log t) + \sum_{i=1}^t m H_2\left(\frac{|E_i|}{m}\right),$$

where $H_2(p)$ is binary entropy.

Proof. Enumerative code (rank) the error set per block. \square

³We encode t , loop bounds, and fixed subroutine identifiers.

Union bound over short decoders. There are at most 2^L decoders of length $\leq L$, so a $2^{-\Omega(t)}$ per-decoder success bound survives union bound for $L = \delta t$ with small enough $\delta > 0$.

2.3 SILS: Sign-Invariant Local Sketches (short, polytime features)

We require a polynomial-time *feature extractor* that maps a masked CNF F^h on m variables to a short, *sign-invariant* summary $\mathbf{z}(F^h) \in \{0, 1\}^{r(m)}$ with $r(m) = O(\log m)$. We call such summaries *SILS* (Sign-Invariant Local Sketches).

Definition 2.7 (SILS, H_m -invariance and interface). Let $H_m := S_m \ltimes (\mathbb{Z}_2)^m$ act on signed CNFs by variable renaming and literal sign flips. A mapping

$$\text{feat} : \text{CNF}_m \longrightarrow \{0, 1\}^{r(m)}$$

is a *SILS extractor* if it satisfies:

- (F1) **Sign/permutation invariance.** For all $(\pi, \sigma) \in H_m$, $\text{feat}(F^h) = \text{feat}(F^{(\pi, \sigma)h})$.
- (F2) **Short output.** $r(m) = O(\log m)$.
- (F3) **Efficient computability.** feat is computable in time $\text{poly}(m)$.
- (F4) **Stability under isomorphism (optional).** It may be convenient (but not strictly necessary for the core proof) that feat depends only on the multiset of bounded-radius incidence neighborhoods *ignoring signs*. We formalize this via counts of rooted hypergraph patterns in Remark 2.9.

We write $\mathbf{z} := \text{feat}(F^h)$ and let \mathcal{I} denote the σ -algebra generated by the coordinates of \mathbf{z} . Only (F1)-(F3) are used in the neutrality and switching arguments; (F4) is used in the template-sparsification convenience bounds.

To be maximally pedantic, we can make the length bound explicit and forbid sign-sensitive features:

Definition 2.8 (SILS contract (length and invariance)). A SILS map is a polynomial-time function $z : \text{CNF}_m \rightarrow \{0, 1\}^{r(m)}$ with $r(m) \leq c_z \log m$ for an absolute constant c_z , such that $z(F^h)$ depends only on the sign-invariant isomorphism type of the factor graph of F^h (i.e., invariant under $H_m = S_m \ltimes (\mathbb{Z}_2)^m$). In particular, features that depend on literal *signs* (e.g., clause-parity by signs) are excluded; degree/profile and small-radius neighborhood counts *ignoring signs* are admissible.

Remark 2.9 (Concrete SILS instantiations). Any of the following (coarsened to $O(\log m)$ bits) yields a valid SILS:

- **Degree/profile sketches.** The degree histogram of the variable?clause incidence hypergraph (ignoring literal signs), bucketed logarithmically.
- **Local pattern counts.** Counts of rooted incidence neighborhoods of fixed radius ρ (constant), ignoring signs, coarsened and hashed to $O(\log m)$ bits (e.g., via pairwise-independent hashing).
- **Co-occurrence statistics (sign-agnostic).** Quantized metrics of variable co-occurrence *ignoring signs* (e.g., mutual-information surrogates over unsigned literals), mapped to $O(\log m)$ bits.

- **Any prior SILS-style summary restricted to sign-agnostic guards.** If desired, one may reuse existing SILS guards as long as they are computed without literal signs and are quantized to $O(\log m)$ bits.

These choices are all H_m -invariant, short, and computable in $\text{poly}(m)$ time.

Definition 2.10 (Local VV labels for bit i). Given the parity matrix $A \in \{0, 1\}^{k \times m}$ and right-hand side $b \in \{0, 1\}^k$ (from the VV layer), let $a_i := Ae_i \in \{0, 1\}^k$ denote the i -th column. We call (a_i, b) the *VV labels* for bit i ; their total length is $O(\log m)$ per block.

Interface contract used later. Our proofs in Sections 3-7 only rely on: (i) sign/permutation invariance (F1) to invoke the promise-preserving involutions and prove $\Pr[X_i = 1 \mid \mathcal{I}] = \frac{1}{2}$; (ii) shortness (F2) and computability (F3) to ensure the post-switch per-bit rules have $O(\log m)$ inputs and compile to tiny ACC⁰; and (iii) independence across blocks, which comes from the sampling process, not from feat. When we use sparsification over radius- r charts, we optionally instantiate (F4) for convenience.

2.4 Valiant-Vazirani isolation via universal hashing

We use the standard universal family of \mathbb{F}_2 -linear hashes.

Definition 2.11 (Linear universal hashing). For integers $k, m \geq 1$, let $\mathcal{H}_{k,m}$ be the family $\{h_{A,b}(x) = Ax \oplus b : A \in \{0, 1\}^{k \times m}, b \in \{0, 1\}^k\}$ with A chosen from any 2-universal distribution over $\{0, 1\}^{k \times m}$ (e.g., rows chosen uniformly and independently), and b uniform.

Isolation lemma (classical form). Let $S \subseteq \{0, 1\}^m$ be nonempty. If $k = \lceil \log_2 |S| \rceil + u$ with $u \in \{0, 1\}$ and $h \sim \mathcal{H}_{k,m}$, then

$$\Pr_h [|S \cap h^{-1}(0^k)| = 1] \geq \frac{1}{8}.$$

This is the Valiant-Vazirani bound; see, e.g., *Valiant & Vazirani (1986)*. When $|S|$ is unknown, choosing k uniformly from $\{0, 1, \dots, m-1\}$ yields $\Pr[|S \cap h^{-1}(0^k)| = 1] \geq \Omega(1/m)$, which is enough for efficient rejection sampling.

We will use the following consequence tailored to our setting (see [1] for the isolation probability and [2, 3] for 2-universal and small-bias hash families):

Lemma 2.12 (VV isolation with small seeds; efficient sampling). *Fix m . Given any satisfiable CNF F with at least one solution and at most $2^{\alpha m}$ solutions (for some absolute $\alpha < 1$), let $k \in \{0, 1, \dots, m-1\}$ be chosen uniformly at random, and pick $h_{A,b} \sim \mathcal{H}_{k,m}$ independently of F . Then*

$$\Pr_{k,A,b} [|\{x \in \{0, 1\}^m : x \models F, Ax = b\}| = 1] \geq \frac{c}{m},$$

for some absolute constant $c > 0$ (independent of m and F). Hence the distribution of pairs $(F, h_{A,b})$ conditioned on uniqueness can be sampled in expected $O(m)$ trials.

Proof sketch. Apply the classical VV bound with k uniform in a logarithmic window around $\log_2 |S|$; averaging over k yields $\Omega(1/m)$. The 2-universality suffices. The upper bound $2^{\alpha m}$ on $|S|$ is used only to ensure the window lies within $\{0, \dots, m-1\}$. \square

Remark 2.13 (Promise semantics). We will *condition* on the uniqueness event and work in the resulting USAT promise problem. Verification (“does x satisfy the CNF and the XORs?”) remains polynomial-time, so all learning and counting arguments are unaffected.

2.5 Masked random 3-CNF and local tree-likeness

Our base distribution is random 3-CNF at constant clause density αm , masked by a fresh $h = (\pi, \sigma) \in H_m$ per block: variables are permuted by π and every literal is independently sign-flipped via σ . The mask is published implicitly by publishing the masked formula F^h .

We rely on the standard “locally tree-like” property of sparse random (hyper)graphs.

Lemma 2.14 (Local tree-likeness with independent signs). *Fix $\alpha > 0$. There exists $c_3^*(\alpha) > 0$ such that for each $c_3 \in (0, c_3^*)$ and $r = c_3 \log m$, the radius- r rooted neighborhood of a uniformly random variable in the masked 3-CNF is a tree with probability $\geq 1 - m^{-\beta}$ (for some $\beta = \beta(\alpha, c_3) > 0$), and the edge signs induced by the mask are i.i.d. Rademacher. Moreover, for any fixed signed rooted pattern \mathcal{T} of radius r , $\Pr[\text{neighborhood equals } \mathcal{T}] \leq m^{-\beta'}$.*

Proof sketch. Classical branching-process approximation for sparse random hypergraphs plus a union bound; the sign flips of the mask are independent and uniform. \square

2.6 Milestone-1 single-block lower bounds (restricted decoders)

We will appeal to standard circuit/streaming lower bounds in a *post-switch* regime where each per-bit rule has only $O(\log m)$ inputs.

- **ACC⁰/AC⁰[p] lower bounds.** For parity and related mod functions, AC⁰ lower bounds via Hastad’s switching lemma; for AC⁰[p], Razborov-Smolensky; for ACC⁰, we use that small ACC⁰ on $O(\log m)$ inputs cannot realize more than $m^{O(1)}$ functions and cannot achieve a noticeable correlation with unbiased random bits (this is sufficient in our setup).
- **Streaming space bounds.** One-pass streaming algorithms with subquadratic space have exponentially small advantage in predicting a random unbiased bit unless they are given more than $O(\log m)$ bits of relevant advice; in our regime, the per-bit input to the post-switch streaming routine is $O(\log m)$ bits.

For our purposes, it is enough to record the following abstract statement.

Lemma 2.15 (Restricted per-block advantage bound). *There is a function $\varepsilon(m) \rightarrow 0$ such that for any Boolean function class \mathcal{C}_m consisting of either (i) depth- d ACC⁰ circuits of size $O(\log m)$ on $O(\log m)$ inputs, or (ii) one-pass streaming algorithms using $o(m^2)$ space on input length $O(\log m)$, every $f \in \mathcal{C}_m$ satisfies*

$$\left| \Pr[f(U) = 1] - \frac{1}{2} \right| \leq \varepsilon(m)$$

where U is uniformly random in $\{0, 1\}^{O(\log m)}$.

Remark 2.16. Lemma 2.15 is used only after the switching step has reduced each per-bit decision to a function of $O(\log m)$ local inputs (\mathbf{z}, a_i, b) . In that regime, uniform randomness of the (signed) local neighborhood and the VV labels justifies applying the lemma to bound advantage per block.

2.7 What is used later (checklist)

For convenience, we list the background facts that subsequent sections rely on:

1. **Weakness calculus:** Invariance of K^{poly} ; Lemma 2.2 (chain rule); Lemma 2.3 (block additivity); $O(\log t)$ wrapper overhead.

2. **Compression from success:** Lemma 2.5 (coarse success-set coding) and Lemma 2.6 (per-bit enumerative coding).
3. **SILS features:** A sign-invariant, $\text{poly}(m)$ -time feature extractor feat outputting $r(m) = O(\log m)$ bits per block.
4. **VV isolation:** Lemma 2.12 (efficient rejection sampling to the unique-witness promise); notation $a_i := Ae_i$, b as VV labels.
5. **Masked ensemble and local tree-likeness:** Lemma 2.14 with $r = c_3 \log m$, giving exponentially small probabilities for fixed signed local patterns.
6. **Restricted per-block advantage bound:** Lemma 2.15 for tiny ACC^0 (or low-space) functions on $O(\log m)$ inputs.

These are all the background ingredients needed by Sections 3-7 to evaluate our proofs. We emphasize that no cryptographic assumptions are used, and all sampling/verification procedures are polynomial-time under the uniqueness promise.

3 The Masked Block Ensemble and Symmetries

In this section we define the masked random 3-CNF plus VV isolation *block* distribution and the H_m -symmetries. Two properties matter most here: (i) a sign-flip/b-toggle involution that preserves uniqueness and toggles any single witness bit, and (ii) local tree-likeness at radius $c_3 \log m$. These supply the symmetry and sparsity pillars used later (Milestone M1).

3.1 Sampling procedure and the USAT promise

Fix clause density $\alpha > 0$ and integers $m \geq 1$ and $k = c_1 \log m$. Let $M := \lfloor \alpha m \rfloor$ denote the number of clauses.

Definition 3.1 (Base random 3-CNF). We draw an *unsigned* 3-uniform hypergraph on vertex set $[m] := \{1, \dots, m\}$ by sampling M triples independently and uniformly with replacement. Write this hypergraph as F ; it carries no literal signs.

Definition 3.2 (Mask group and its action). Let $H_m := S_m \times (\mathbb{Z}_2)^m$ act on signed CNFs by

$$(\pi, \sigma) \cdot ((\ell_{j,1} \vee \ell_{j,2} \vee \ell_{j,3})_{j \in [M]}) := ((\ell_{j,1}^\sigma \circ \pi) \vee (\ell_{j,2}^\sigma \circ \pi) \vee (\ell_{j,3}^\sigma \circ \pi))_{j \in [M]},$$

where π permutes variable names and $\sigma \in (\mathbb{Z}_2)^m$ flips literal signs coordinate-wise. Given an *unsigned* F , a *mask* $h = (\pi, \sigma) \in H_m$ produces a *signed* CNF F^h by first assigning all literals positive and then applying h .

Definition 3.3 (VV isolation layer; instance). Sample $A \in \{0, 1\}^{k \times m}$ from any 2-universal distribution with pairwise-independent columns, and sample $b \in \{0, 1\}^k$ from a δ -biased source with $\delta = m^{-c_2}$, independently of (F, h) . The full instance is

$$\Phi := (F^h, A, b).$$

Let $\text{Unq}(\Phi)$ denote the event that Φ has a *unique* satisfying assignment $x \in \{0, 1\}^m$.

Definition 3.4 (Block distribution \mathcal{D}_m). The block distribution \mathcal{D}_m is the law of Φ from Definition 3.3 *conditioned* on $\text{Unq}(\Phi)$. By Lemma 2.12, rejection sampling reaches \mathcal{D}_m in expected $O(m)$ trials.

We write $a_i := Ae_i \in \{0,1\}^k$ for the i -th column of A and refer to (a_i, b) as the *VV labels* for bit i . Given Φ , the (unique) witness is denoted $X := x(\Phi) \in \{0,1\}^m$.

Definition 3.5 (i.i.d. block product). For $t = c_4 m$ (fixed $c_4 > 0$), an input to a decoder is the t -tuple (Φ_1, \dots, Φ_t) of i.i.d. draws from \mathcal{D}_m ; the corresponding witness tuple is (X_1, \dots, X_t) .

VV labels and robustness to δ -bias. For any fixed A and σ , the map $b \mapsto b \oplus A\sigma$ is a bijection on $\{0,1\}^k$ and preserves uniform measure exactly. If b is sampled from a δ -biased source, then $b \oplus A\sigma$ is also δ -biased with the same parameter. All symmetrization and calibration steps remain valid up to an additive $O(\delta)$, which we fold into the $m^{-\Omega(1)}$ slack by setting $\delta \leq m^{-10}$.

3.2 Symmetries and promise-preserving involutions

The following *coordinate sign-flip* maps are the backbone of our AP-GCT neutrality.

Lemma 3.6 (Promise-preserving involution T_i). *For each $i \in [m]$, define*

$$T_i : (F^h, A, b) \mapsto (F^{\tau_i h}, A, b \oplus Ae_i),$$

where $\tau_i \in H_m$ flips only variable i 's literal signs. Then:

- (i) T_i is measure-preserving on the product of the base distributions of (F, h, A, b) ;
- (ii) T_i restricts to a bijection on the promise space $\{\Phi : \text{Unq}(\Phi)\}$; if X satisfies Φ , then $X \oplus e_i$ satisfies $T_i(\Phi)$, and uniqueness is preserved.

Proof. (i) Uniformity and independence of h and b make T_i an automorphism of the sampling measure. (ii) Flipping signs of variable i toggles the i -th bit in any satisfying assignment on the CNF part; the XOR part updates as $A(X \oplus e_i) = AX \oplus Ae_i = b \oplus Ae_i$. The map between satisfying assignments is a bijection, so uniqueness is preserved. \square

Lemma 3.7 (Promise-preserving composition). *Each stage of the pipeline is a bijection on the on-promise set and measure-preserving: (i) masking by H_m ; (ii) VV isolation (A, b) selection; (iii) sign-flip/toggle maps $(F^h, A, b) \mapsto (F^{(\text{id}, \sigma)^h}, A, b \oplus A\sigma)$ used in the wrapper; and (iv) reindexing/back-mapping outputs. Therefore, any finite composition of these maps is promise-preserving and measure-preserving.*

Proof. (i) and (iv) are group actions/bijections. (ii) is a sampling step independent of (F, h) ; restricting to the event “unique witness” defines the promise measure. (iii) is Lemma 3.6 in vector form; uniqueness bijects via $x \mapsto x \oplus \sigma$. Composition of bijective measure-preserving maps is bijective and measure-preserving. \square

Let \mathcal{I} denote any σ -algebra generated by *sign-invariant*, permutation-invariant functions of F^h (e.g., any collection of degree- $\leq D$ pattern counts that ignore literal signs).

Corollary 3.8 (Per-bit neutrality given sign-invariant views). *For every $i \in [m]$, $\Pr[X_i = 1 \mid \mathcal{I}] = \frac{1}{2}$ almost surely under \mathcal{D}_m .*

Proof. Immediate from Lemma 3.6: T_i preserves \mathcal{I} and toggles X_i . \square

3.3 Local σ -fields and the post-switch inputs

We define the per-block local inputs that will parameterize the switched per-bit rules.

Definition 3.9 (Sign-invariant SILS features). Let $\mathbf{z} := \text{feat}(F^h) \in \{0, 1\}^{r(m)}$ be any sign-invariant feature vector computable in time $\text{poly}(m)$ with $r(m) = O(\log m)$ (see §2.3). We denote by \mathcal{I} the σ -algebra generated by the coordinates of \mathbf{z} .

Definition 3.10 (Per-bit local inputs and σ -fields). For a block $\Phi = (F^h, A, b)$ and index i , define the *per-bit local input*

$$\mathbf{u}_i(\Phi) := (\mathbf{z}(F^h), a_i = Ae_i, b) \in \{0, 1\}^{O(\log m)}.$$

Let \mathcal{L}_i be the σ -field generated by $\mathbf{u}_i(\Phi)$. We emphasize that \mathcal{L}_i is *local to bit i in its block*.

3.4 Independence across blocks

Blocks are sampled independently by Definition 3.5. In particular, for any fixed measurable functions g_j , $\{g_j(\Phi_j)\}_{j=1}^t$ are independent random variables. This independence underpins product bounds on success probabilities and learning/generalization arguments.

3.5 Local tree-likeness and signed pattern probabilities

We record a quantitatively explicit local weak-limit statement for our masked ensemble (note a standard reference for local weak convergence and sparse random (hyper)graph neighborhoods is [7]):

Theorem 3.11 (Local tree-likeness at logarithmic radius). *Fix $\alpha > 0$. There exists $c_3^*(\alpha) > 0$ such that for any $c_3 \in (0, c_3^*)$ and $r = c_3 \log m$, the following holds for the masked random 3-CNF:*

- (i) *For a uniformly random variable v , with probability at least $1 - m^{-\beta}$ (for some $\beta = \beta(\alpha, c_3) > 0$), the radius- r neighborhood $\mathcal{N}_r(v)$ in the factor graph is a tree (no cycles) whose unlabeled shape is distributed as a Galton-Watson branching process with offspring distribution $\text{Poisson}(\lambda(\alpha))$ up to depth r .*
- (ii) *Conditional on the unlabeled shape, the literal signs on edges induced by the mask are i.i.d. Rademacher.*
- (iii) *Consequently, for any fixed signed rooted pattern \mathcal{T} of radius r ,*

$$\Pr [\mathcal{N}_r(v) \text{ equals } \mathcal{T}] \leq m^{-\beta'},$$

for some $\beta' = \beta'(\alpha, c_3) > 0$.

Proof sketch. (i) and the unlabeled Galton-Watson coupling are standard for sparse random (hyper)graphs; the cycle probability within radius $r = c_3 \log m$ decays as $m^{-\beta}$ for c_3 small enough. (ii) The mask chooses literal signs independently and uniformly; conditioning on the unlabeled structure does not introduce sign correlation. (iii) Multiply the (exponentially small in r) probability of the unlabeled shape by $2^{-|E(\mathcal{T})|}$ for the signs, and choose c_3 so the product is at most $m^{-\beta'}$. \square

3.6 Parameters and notational summary

We summarize the fixed parameters used later:

- Clause density $\alpha > 0$ (constant).
- VV parameters: $k = c_1 \log m$, $\delta = m^{-c_2}$; (A, b) independent of (F, h) .
- Mask: a fresh $h \in H_m$ per block, uniform.
- Features: $\mathbf{z} \in \{0, 1\}^{r(m)}$ with $r(m) = O(\log m)$, sign-invariant, $\text{poly}(m)$ computable.
- Neighborhood radius: $r = c_3 \log m$ with $c_3 \in (0, c_3^*(\alpha))$ (Theorem 3.11).
- Number of blocks: $t = c_4 m$ with fixed $c_4 > 0$; i.i.d. across blocks (Definition 3.5).

What Section 3 supplies. We will use Lemma 3.6 (promise-preserving T_i) to prove AP-GCT neutrality in Section 5; Theorem 3.11 to bound the probability of fixed signed charts at radius $r = c_3 \log m$; and the local σ -fields \mathcal{L}_i from Definition 3.10 to formalize the post-switch per-bit inputs.

4 Switching-by-Weakness: Wrappers and Post-Switch Class

We first symmetrize P (measure-preserving) and distill its behavior onto the local inputs \mathbf{u} via ERM, obtaining a \mathbf{u} -measurable comparator. We then upper bound any \mathbf{u} -measurable predictor versus truth by neutrality and sparsification (Section 5). The calibration Lemma 4.8 links the symmetrized comparator back to the original P .

In this section (Milestone M2), short decoders become local per-bit decoders on many blocks. We prove a normal form: a length- $\leq \delta t$ decoder admits a short wrapper so that, on a constant-fraction test subset S of blocks, each output bit depends only on $O(\log m)$ local inputs (\mathbf{z}, a_i, b) . We give *two constructive wrappers*: (i) a distributional *distillation* wrapper (ERM route), which we use as the primary argument and which yields both locality on S and success-domination (the wrapper’s comparator does not underperform the original decoder up to $m^{-\Omega(1)}$); and (ii) a symmetrization-based comparator (averaging over a polylogarithmic multiset of promise-preserving sign flips) used to define the surrogate labels distilled by ERM. Both wrappers are short, run in polynomial time, and produce the same local normal form on S .

Throughout this section, unless stated otherwise, a “decoder” P is a deterministic polynomial-time algorithm (coins are fixed into its code) that, on input a t -tuple (Φ_1, \dots, Φ_t) of blocks from \mathcal{D}_m , outputs a tuple of bit-vectors $\hat{X} = (\hat{x}_1, \dots, \hat{x}_t)$ with $\hat{x}_j \in \{0, 1\}^m$.

4.1 Statement of the switching normal form

Definition 4.1 (Local inputs and local σ -fields (recalled)). For a block $\Phi = (F^h, A, b)$ and bit index $i \in [m]$, the *local input* is

$$\mathbf{u}_i(\Phi) := (\mathbf{z}(F^h), a_i = Ae_i, b) \in \{0, 1\}^{O(\log m)}.$$

Let \mathcal{L}_i be the σ -algebra generated by \mathbf{u}_i .

Theorem 4.2 (Switching-by-Weakness (SW)). *There exist constants $\gamma > 0$ and $c^* > 0$ such that for every polynomial-time decoder P with $|P| \leq \delta t$ there is a polynomial-time wrapper W with $|W| \leq |P| + c^*(\log m + \log t)$ and a subset $S \subseteq [t]$ with $|S| \geq \gamma t$ for which:*

$$(P \circ W)(\Phi)_{j,i} = h_{j,i}(\mathbf{u}_i(\Phi_j)) \quad \text{for all } j \in S, i \in [m], \quad (1)$$

for some Boolean maps $h_{j,i} : \{0,1\}^{O(\log m)} \rightarrow \{0,1\}$. Moreover each $h_{j,i}$ is computable in time $\text{poly}(\log m)$ (hence realizable by size $\text{poly}(m)$ ACC⁰).

Surrogate vs. truth. ERM trains on symmetrized (back-mapped) labels, not on X ; the link back to truth is Lemma 4.8, proved in Appendix A.5.

Proof route. We prove Theorem 4.2 via the ERM distillation wrapper (Proposition A.5), which yields both locality on a test subset and success-domination with wrapper length $|W_{\text{ERM}}| \leq |P| + O(\log m + \log t)$. The symmetrization wrapper (§4.2) is used only to define surrogate labels; it has length $|W_{\text{sym}}| = |P| + \tilde{O}(\log^2(mt))$ (Lemma 4.7) and is not needed to meet the length bound in Theorem 4.2.

Lemma 4.3 (Symmetrization preserves success exactly). *Let g_σ be the promise-/measure-preserving sign-flip map and BM_σ the back-map on outputs that xors out $A\sigma$ in the VV layer (coordinate-wise). Then*

$$\Pr_{\Phi \sim \mathcal{D}_m} [P(\Phi) = X(\Phi)] = \mathbb{E}_\sigma \Pr_{\Phi \sim \mathcal{D}_m} [\text{BM}_\sigma(P(g_\sigma(\Phi))) = X(\Phi)].$$

Proof. For any measurable event $E(\Phi, X)$, measure preservation of g_σ on the promise space yields $\Pr_\Phi[E(\Phi, X(\Phi))] = \Pr_\Phi[E(g_\sigma(\Phi), X(g_\sigma(\Phi)))]$. Since $X(g_\sigma(\Phi)) = X(\Phi) \oplus \sigma$ and the VV RHS shifts by $A\sigma$, back-mapping the output undoes this shift, so correctness on $g_\sigma(\Phi)$ equals back-mapped correctness on Φ . Average over σ . \square

Remark 4.4 (Exact vs. approximate preservation). If b is uniform, Lemma ?? holds with equality. If b is δ -biased (and independent of A), the same identity holds up to an additive $O(\delta)$ in total variation; this is absorbed into the $m^{-\Omega(1)}$ slack.

Theorem 4.5 (SW completeness and success domination). *For every polynomial-time decoder P of description length $\leq \delta t$ there exists a wrapper W of length $|W| \leq |P| + O(\log m + \log t)$ such that: (i) the locality conclusion of Theorem 4.2 holds on a subset S with $|S| \geq \gamma t$; and (ii) success domination holds:*

$$\Pr [(P \circ W)(\Phi) = X] \geq \Pr [P(\Phi) = X] - m^{-\Omega(1)}.$$

Proof sketch. Draw $s = \Theta(\log(mt))$ independent flips $\sigma^{(1)}, \dots, \sigma^{(s)}$ from a κ -wise independent family with $\kappa = \Theta(\log(mt))$. For each r , the map $g_{\sigma^{(r)}}$ is measure- and promise-preserving (Lemma 3.6), hence by Lemma 4.3: $\mathbb{E}_\sigma \mathbb{E}_\Phi \mathbf{1}\{\text{BM}_\sigma(P(g_\sigma(\Phi))) = X(\Phi)\} = \mathbb{E}_\Phi \mathbf{1}\{P(\Phi) = X\}$. By Hoeffding under limited independence, the majority of the back-mapped predictions matches the Bayes rule on the local σ -field for all but $o(t)$ blocks, with probability $1 - m^{-\Omega(1)}$ over the seeds. This majority is at least as accurate as the average prediction on each block, so the overall success does not decrease by more than $m^{-\Omega(1)}$. Fix seeds with this property and bake them into W . Locality and size follow from Theorem 4.2. \square

Calibration in one line. For fixed $u = (z, a_i, b)$, the promise-preserving involution T_i bijects $(X_i = 0, Y_i = y) \leftrightarrow (X_i = 1, Y_i = 1 - y)$ without changing u or the measure. Thus $(X_i, Y_i) \mid u$ is exchangeable, so $\Pr[X_i = 1 \mid u] = \Pr[Y_i = 1 \mid u] = f_i(u)$, and the Bayes rule $h_i^*(u) = \mathbf{1}[f_i(u) \geq 1/2]$ is optimal for both. (Full proof in Lemma A.16.)

Corollary 4.6 (Domination principle: bounds for P via its comparator). *For every polynomial-time decoder P of description length $\leq \delta t$ there exists a wrapper W with $|W| \leq |P| + O(\log m + \log t)$ such that*

$$\Pr[P(\Phi) = X] \leq \Pr[(P \circ W)(\Phi) = X] + m^{-\Omega(1)}.$$

If, moreover, $(P \circ W)$ satisfies the local normal form on γt blocks (Theorem 4.2), then any upper bound proved for $\Pr[(P \circ W)(\Phi) = X]$ applies to $\Pr[P(\Phi) = X]$, up to $m^{-\Omega(1)}$.

We give two constructive proofs: (i) a distributional *distillation* wrapper (ERM route), which we use as the primary argument; and (ii) a symmetrization-based comparator (averaging over a polylogarithmic multiset of promise-preserving sign flips) used to define the labels distilled by ERM. Both wrappers are short and run in polynomial time.

Domination vs. equivalence. The wrapper provides a *comparator* whose success dominates that of P up to $m^{-\Omega(1)}$, and whose predictions are local on γt blocks; we do not claim P itself is local. All upper bounds we prove for the comparator therefore apply to P .

4.2 Symmetrization wrapper (promise-preserving, short description)

We use only *sign flips*; permutations are not needed because the SILS vector \mathbf{z} is sign-invariant and permutation-invariant in the sense of Def. 2.7(F1). Sign flips are promise-preserving via Lemma 3.6.

Small seed families of flips. Fix integers

$$s := C \cdot (\log m + \log t), \quad \kappa := C' \cdot (\log m + \log t),$$

for sufficiently large absolute constants C, C' . Let \mathcal{S} be an explicit κ -wise independent family of functions $\sigma : [m] \rightarrow \{0, 1\}$ with seed length $O(\kappa)$ (e.g., low-degree polynomial families over a suitable field), and define the blockwise sign-flip operator

$$g_\sigma : (F^h, A, b) \mapsto (F^{(\text{id}, \sigma)^h}, A, b \oplus A\sigma),$$

where we view σ also as a vector in $\{0, 1\}^m$ and set $A\sigma := \sum_i \sigma(i) A e_i$. By Lemma 3.6, each g_σ is measure-preserving and *promise-preserving*. Sampling σ uniformly from \mathcal{S} requires only $O(\kappa)$ seed bits and yields κ -wise independence across the s draws used below.

Definition of the wrapper W_{sym} . Hard-wire s independent seeds ρ_1, \dots, ρ_s of total length $O(s\kappa) = O((\log m + \log t)^2)$. On input (Φ_1, \dots, Φ_t) :

1. For each $r \in [s]$, instantiate $\sigma^{(r)} \leftarrow \mathcal{S}(\rho_r)$ and form the sign-flipped tuple

$$\Phi^{(r)} := (g_{\sigma^{(r)}}(\Phi_1), \dots, g_{\sigma^{(r)}}(\Phi_t)).$$

2. Run P on each $\Phi^{(r)}$, obtaining predictions $\widehat{X}^{(r)} = (\widehat{x}_1^{(r)}, \dots, \widehat{x}_t^{(r)})$.

3. For each block j and bit i , *back-map* to the original coordinates:

$$Y_{j,i}^{(r)} := \widehat{x}_{j,i}^{(r)} \oplus \langle a_{j,i}, \sigma^{(r)} \rangle \quad (\text{where } \langle \cdot, \cdot \rangle \text{ is inner product over } \mathbb{F}_2).$$

4. Output the majority $\hat{y}_{j,i} := \text{Maj}(Y_{j,i}^{(1)}, \dots, Y_{j,i}^{(s)})$.

Return $\hat{X} := (\hat{y}_{j,i})_{j \in [t], i \in [m]}$.

Lemma 4.7 (Budget and running time). W_{sym} is polynomial-time and has description length $|W_{\text{sym}}| \leq |P| + O(\log m + \log t)$, counting the $O((\log m + \log t)^2)$ seed bits only once as advice.

Proof. The wrapper makes $s = \Theta(\log(mt))$ oracle calls to P and performs linear-time postprocessing per call. The advice consists of s seeds ($O(\kappa)$ bits each) plus loop overhead; these are $O((\log m + \log t)^2)$ bits total. Since we compare against the budget $\delta t = \Theta(m)$, this is absorbed by $c^*(\log m + \log t)$. \square

What the symmetrization yields. For fixed (j, i) and local input $u_i(\Phi_j) = (z_j, a_{j,i}, b_j)$, the symmetrized label $\tilde{Y}_{j,i}$ is the majority of back-mapped predictions over the limited-independence sign flips. We use $\tilde{Y}_{j,i}$ as surrogate labels and *distill* a local comparator on the distribution via ERM (Appendix A.1). The locality claim in Theorem 4.2 is then achieved by the ERM wrapper, while symmetrization is used only to define the labels.

Lemma 4.8 (Calibration from symmetrized labels to truth; distributional). *Fix a bit index i and define $Z(\sigma, \Phi) := \mathbf{1}\{Y_i(\sigma, \Phi) = X_i(\Phi)\}$, where Y_i is the back-mapped prediction defined above. Let $f_i(\mathbf{u}) = \mathbb{E}[Y_i(\sigma, \Phi) \mid \mathbf{u}]$ and let $h_i^*(\mathbf{u})$ be the Bayes classifier for f_i . Then*

$$\mathbb{E}_{\Phi}[\mathbf{1}\{h_i^*(\mathbf{u}(\Phi)) = X_i(\Phi)\}] \geq \mathbb{E}_{\Phi, \sigma}[Z(\sigma, \Phi)] - m^{-\Omega(1)}.$$

Consequently, for the ERM predictor \hat{h}_i (which approximates h_i^* on the test distribution),

$$\mathbb{E}_{\Phi}[\mathbf{1}\{\hat{h}_i(\mathbf{u}(\Phi)) = X_i(\Phi)\}] \geq \mathbb{E}_{\Phi, \sigma}[Z(\sigma, \Phi)] - m^{-\Omega(1)}.$$

Proof sketch. For a fixed \mathbf{u} , the random variable $Y_i(\sigma, \Phi)$ is a Bernoulli with mean $f_i(\mathbf{u})$. The Bayes classifier for f_i in 0/1 loss against Y_i is $\text{sgn}(f_i - 1/2)$. In our masked+isolated ensemble, the same sign choice also maximizes agreement with X_i on average (up to $m^{-\Omega(1)}$). This uses the paired-involution structure (flip i and toggle b by a_i), which relates (\mathbf{u}, X_i, Y_i) to $(\mathbf{u}, 1 - X_i, 1 - Y_i)$ and makes the pairwise distributions symmetric in the sense required for calibration. The detailed argument appears in Appendix A.6. \square

Limited-independence Chernoff parameters . We take $s := \lceil 20 \log_2(mt) \rceil$ symmetrization calls and $\kappa := \lceil 12 \log_2(mt) \rceil$ -wise independence. Then for each (j, i) ,

$$\Pr \left[\left| \frac{1}{s} \sum_{r=1}^s \mathbf{1}\{Y_{j,i}^{(r)} = X_{j,i}\} - p_{j,i} \right| > \frac{1}{m^3} \right] \leq m^{-10},$$

by Schmidt-Siegel-Srinivasan; a union bound over all $mt = \tilde{O}(m^2)$ pairs gives failure probability m^{-8} . We threshold at $1/2$ thereafter.

Lemma 4.9 (Concentration to the Bayes rule). *There exists $\varepsilon(m) = m^{-\Omega(1)}$ such that, for each fixed (j, i) ,*

$$\Pr_{\rho_1, \dots, \rho_s} \left[\text{Maj}(Y_{j,i}^{(1)}, \dots, Y_{j,i}^{(s)}) \neq h_i^*(\mathbf{u}_i(\Phi_j)) \right] \leq \varepsilon(m).$$

Moreover, by a union bound and κ -wise independence (with $\kappa = \Theta(\log(mt))$), the event that this equality holds simultaneously for all but an $o(1)$ fraction of blocks j (and for all i) has probability at least $1 - m^{-\Omega(1)}$ over the choice of seeds.

Proof. Each $Y_{j,i}^{(r)}$ has mean $p_{j,i}$ and the collection is κ -wise independent. By standard Chernoff bounds under κ -wise independence (with $\kappa = \Theta(\log(mt))$), the empirical average $\frac{1}{s} \sum_r Y_{j,i}^{(r)}$ deviates from $p_{j,i}$ by more than $1/\text{poly}(m)$ with probability $m^{-\Omega(1)}$. Thresholding at $1/2$ yields the claim, and a union bound across (j, i) establishes the simultaneous statement. \square

Lemma 4.10 (Non-degradation in expectation). *For any decoder P ,*

$$\mathbb{E}_{\Phi \sim \mathcal{D}_m^{\otimes t}} [\mathbf{1}\{(P \circ W_{\text{sym}})(\Phi) = X\}] \geq \mathbb{E}_{\Phi \sim \mathcal{D}_m^{\otimes t}} [\mathbf{1}\{P(\Phi) = X\}] - m^{-\Omega(1)}.$$

Proof. By Lemma 4.3, the averaged (over σ) success of P equals the success of the back-mapped Bayes rule. By Lemmas 4.8 and 4.9, the majority output matches the Bayes rule except with probability $m^{-\Omega(1)}$. Aggregating over bits and blocks yields the claim. \square

We can now finish Theorem 4.2.

Proof of Theorem 4.2. Fix seeds as in Lemma 4.9; bake them into W_{sym} . Define $S \subseteq [t]$ to be the set of blocks on which the equality $\text{Maj}(Y_{j,i}^{(1)}, \dots, Y_{j,i}^{(s)}) = h_i^*(\mathbf{u}_i(\Phi_j))$ holds for *all* bits $i \in [m]$. By Lemma 4.9 and independence across blocks, $|S| \geq \gamma t$ with probability $1 - m^{-\Omega(1)}$ for some constant $\gamma > 0$. On S , define $h_{j,i} := h_i^*$; then (1) holds by construction, and each $h_{j,i}$ depends only on $\mathbf{u}_i(\Phi_j)$ and is computable in time $\text{poly}(\log m)$ (by lookup on $\{0, 1\}^{O(\log m)}$), thus realizable by size $\text{poly}(m)$ ACC^0 . Finally, using Proposition A.5 we instantiate W as W_{ERM} , which meets the claimed length bound $|W| \leq |P| + O(\log m + \log t)$. \square

What we use symmetrization for. (i) The equality of success in Lemma ?? (average over σ equals original). (ii) Surrogate labels \tilde{Y} used by ERM. Locality itself is delivered by the ERM plug-in rule on the finite alphabet U (no symmetrization needed at test time).

4.3 Finite-alphabet locality and (optional) ACC^0 compilation

The post-switch input for bit i is $\mathbf{u} = (\mathbf{z}, a_i, b)$ with $|\mathbf{u}| = O(\log m)$. Hence the local alphabet \mathcal{U} has size $|\mathcal{U}| \leq 2^{r(m)} \cdot 2^{2k} = m^{O(1)}$.

Lemma 4.11 (Compilation at logarithmic input length). *For any fixed Boolean $h : \{0, 1\}^d \rightarrow \{0, 1\}$ with $d = O(\log m)$ there exists a depth-2 circuit of size $O(2^d) = \text{poly}(m)$ (hence also an ACC^0 circuit of $\text{poly}(m)$ size) that computes h .*

Proof. Tabulate h and implement the balanced DNF (or CNF) over d inputs; size $O(2^d)$. \square

ERM without hypothesis enumeration. Let $T \sqcup S = [t]$ be a random train/test split with $|T|, |S| = \Theta(t)$. For each bit index i define the plug-in rule on the finite alphabet \mathcal{U} by

$$\hat{h}_i(\mathbf{u}) := \text{Maj}\{\tilde{Y}_{j,i} : j \in T, \mathbf{u}_{j,i} = \mathbf{u}\},$$

where $\tilde{Y}_{j,i}$ are the symmetrized back-mapped labels (Def./Lemmas in App. A.1). On the test blocks $j \in S$, the wrapper outputs $(P \circ W_{\text{ERM}})(\Phi)_{j,i} := \hat{h}_i(\mathbf{u}_{j,i})$. This is *local* and computable in $\text{poly}(m)$ time by hash-table lookup on \mathcal{U} ; no class enumeration is required and the wrapper description length remains $|P| + O(\log m + \log t)$.

ERM is plug-in on a finite alphabet. The post-switch input $u = (z, a_i, b)$ has $|u| = O(\log m)$, hence the alphabet U has size $|U| = m^{O(1)}$. Our ERM rule is the plug-in majority

$$\hat{h}_i(u) := \text{Maj}\{\tilde{Y}_{j,i} : j \in T, u_{j,i} = u\},$$

implemented by a hash table over U . No hypothesis enumeration is required. Hoeffding plus a union bound over $u \in U$ and $i \in [m]$ yields $\Pr_{j \in S}[\hat{h}_i(u_{j,i}) \neq h_i^*(u_{j,i})] \leq m^{-\Omega(1)}$ with $|T| = \Theta(m)$ samples. Optional compilation to circuits is by a depth-2 lookup/DNF of size $O(|U|) = \text{poly}(m)$; we do not claim or use tiny ACC^0 in the learning step.

4.4 Remarks on promise semantics and determinism

- **Promise-preserving operations.** Every sign flip g_σ preserves the sampling measure and the uniqueness promise (Lemma 3.6); thus W_{sym} operates entirely within the USAT promise space.
- **Randomized decoders.** If P uses internal coins, fix them into its code (this increases $|P|$ by at most an additive constant); all statements above apply to the determinized decoder.
- **Success non-degradation.** Lemma 4.10 shows the wrapper does not *decrease* success in expectation. This permits transferring any upper bound we prove for $(P \circ W)$ back to P , up to negligible $m^{-\Omega(1)}$ slack.

Summary of Section 4. For every short decoder P , the symmetrization wrapper W_{sym} (i) has short description, (ii) is polynomial-time, (iii) produces a per-bit *local* rule on $\Omega(t)$ blocks depending only on the SILS \mathbf{z} and VV labels (a_i, b) , and (iv) does not degrade success in expectation. The post-switch per-bit rules are realizable by $\text{poly}(m)$ size ACC^0 on the finite alphabet \mathcal{U} (size $m^{O(1)}$), which is the regime needed for neutrality and sparsification in Section 5.

4.5 Why the Switching-by-Weakness proof works in this framework

The ERM/distillation switching argument (Appendix A.1) depends on five pillars that are special to our setup and together make the proof go through:

(1) Compositionality of weakness. We measure “shortness” by K^{poly} , which is *compositional*: (i) invariant up to $O(1)$ (machine choice); (ii) obeys a chain rule and block additivity (Lemma A.8); (iii) supports *Compression-from-Success* (Lemma A.9). This lets us: (a) pay only $O(\log m + \log t)$ bits for any wrapper control flow; (b) aggregate per-program small success across t blocks into a linear tuple lower bound; and (c) oppose that lower bound to the constant upper bound under $P = \text{NP}$ (Proposition 7.2).

(2) Promise-preserving symmetry as a two-way bridge. The sign-flip action g_σ is a *measure- and promise-preserving bijection* on \mathcal{D}_m (Lemma 3.6, Lemma 3.7). This gives two crucial properties: (i) *exact success preservation*: By Lemma 4.3, averaging P over σ and back-mapping preserves its success on the promise distribution exactly; (ii) *neutrality for sign-invariant views*: for any sign-invariant σ -algebra \mathcal{I} (e.g., generated by SILS), $\Pr[X_i = 1 \mid \mathcal{I}] = 1/2$ (Appendix A.3). Together these facts let us compare the global P to a *more symmetric* comparator that we can analyze locally.

(3) Low-dimensional locality by design. The local input $\mathbf{u} = (\mathbf{z}, a_i, b)$ is *short*: SILS \mathbf{z} has $r(m) = O(\log m)$ bits and the VV labels (a_i, b) contribute $O(\log m)$ more. Hence the local interface has *polynomial alphabet size* $|\mathcal{U}| = m^{O(1)}$; ERM operates on \mathcal{U} via a plug-in rule, and (optional) ACC^0 compilation is $\text{poly}(m)$ by Lemma 4.11. This is what makes ERM *work with guarantees*: the alphabet is small enough that uniform convergence holds with $\text{poly}(m)$ samples (Lemma ??).

(4) Distillation with calibration. We do *not* claim P is local. Instead, we *distill* the σ -averaged behavior of P onto $h(\mathbf{u})$ (the Bayes classifier for surrogate labels) and prove via Lemma 4.8 that the surrogate-to-truth calibration holds:

$$\Pr[P(\Phi) = X] \leq \Pr[(P \circ W_{\text{ERM}})(\Phi) = X] + m^{-\Omega(1)}.$$

This comparator is *local on a constant fraction of blocks* (Theorem 4.2 / Proposition A.5), so all neutrality/sparsification bounds apply to it; by domination, they apply to P as well. No “compressibility of algorithms” or per-instance measurability is assumed.

(5) Distributional sparsity and independence where needed. Random 3-CNF is *locally tree-like* at radius $c_3 \log m$ (Theorem 3.11), and the mask gives i.i.d. signs. At this radius, any fixed *signed chart* (neighborhood + VV labels) appears with probability $m^{-\Omega(1)}$, so a *polynomial* family of local rules can be high-bias on at most $o(t)$ blocks (Theorem A.15). After fixing the wrapper W_{ERM} (train/test split, seeds, trained $\{\hat{h}_i\}$), predictions on test blocks depend only on those blocks; independence across $j \in S$ is inherited from the product distribution (Lemma 6.6). This is the exact independence we use for product bounds?no unproved intra-block independence is needed.

Synthesis. These pillars support the entire chain:

shortness \Rightarrow distillation to local comparator on S & success domination \Rightarrow local near-randomness on $S \Rightarrow$ prod

which clashes with the *constant* upper bound under $P = \text{NP}$. The proof succeeds here precisely because the symmetry/promise structure, the $O(\log m)$ local interface, and the quantale calculus were designed to make these implications composable and analyzable.

5 AP-GCT Neutrality and Template Sparsification

Here we prove per-bit neutrality for any *sign-invariant* view (symmetry says: conditional mean is $1/2$), and we prove a template sparsification theorem at logarithmic radius (sparsity says: a fixed local chart is hit with probability $m^{-\Omega(1)}$). Together, any post-switch per-bit rule (from the finite alphabet) is near-random on a constant fraction of blocks. This is Milestone M1 in action.

Specifically, we establish two complementary mechanisms that force *local* unpredictability on many blocks for every short decoder:

1. **AP-GCT neutrality:** for any sign-invariant view \mathcal{I} of a masked block, each witness bit has conditional mean $1/2$ (no bias).
2. **Template sparsification at logarithmic radius:** for any fixed local per-bit rule on inputs (\mathbf{z}, a_i, b) of length $O(\log m)$, the event “this rule attains noticeable bias on a random block” has probability $m^{-\Omega(1)}$; hence at most $o(t)$ blocks can be “high-bias” for that rule, and by a union bound, for any *polynomial* family of such rules.

Combined with the Switching-by-Weakness normal form (Theorem 4.2), these imply that on a γ -fraction of blocks the switched per-bit rules are near-random (bias at most $\varepsilon(m) \rightarrow 0$), which feeds the per-block lower bounds of Section 6.

5.1 AP-GCT neutrality for sign-invariant views

Recall the promise-preserving involution T_i (Lemma 3.6) and let \mathcal{I} be the σ -algebra generated by any family of *sign-invariant*, permutation-invariant functions of F^h (e.g., the SILS coordinates; Def. 2.7).

Theorem 5.1 (Per-bit neutrality). *For every $i \in [m]$ and every sign-invariant view \mathcal{I} ,*

$$\Pr[X_i = 1 \mid \mathcal{I}] = \frac{1}{2} \quad \text{almost surely under } \mathcal{D}_m.$$

Proof. T_i preserves the sampling measure and the uniqueness promise, toggles X_i , and fixes \mathcal{I} (Lemma 3.6). For every \mathcal{I} -measurable event B , $\Pr[X_i = 1 \wedge B] = \Pr[X_i = 0 \wedge B]$, hence the conditional probability is $1/2$. \square

Corollary 5.2 (SILS-only predictors are neutral). *Let $g : \{0, 1\}^{r(m)} \rightarrow \{0, 1\}$ be any SILS-only bit predictor. Then for each i , $\Pr[g(\mathbf{z}) = X_i] = \frac{1}{2}$.*

Remark 5.3. Neutrality does *not* speak to predictors that also use the VV labels (a_i, b) . For those we rely on sparsification below.

5.2 Charts on radius- r signed neighborhoods and labels

Fix $r = c_3 \log m$ with $c_3 \in (0, c_3^*(\alpha))$ as in Theorem 3.11. We formalize the local information available to a per-bit rule at this radius.

Definition 5.4 (Signed neighborhood extractor). For a masked block $\Phi = (F^h, A, b)$, bit index i , and radius r , let $\text{nbr}_r(\Phi, i)$ denote the rooted, *signed* radius- r neighborhood of variable i in the factor graph of F^h , with signs on incident literal edges.

Definition 5.5 (Charts with labels). A *chart* is a pair $\mathcal{C} = (\mathcal{P}, \psi)$ where:

- \mathcal{P} is a finite set of *signed* rooted radius- r patterns, augmented with the port labels $(a_i, b) \in \{0, 1\}^k \times \{0, 1\}^k$ for the root bit;
- $\psi : \mathcal{P} \rightarrow \{0, 1\}$ is a decision rule.

We say that (Φ, i) *matches* \mathcal{C} if there exists $P \in \mathcal{P}$ with $\text{nbr}_r(\Phi, i) = P$ (including the labels).

Definition 5.6 (High-bias region for a chart). Fix $\varepsilon > 0$. The *high-bias region* of a chart \mathcal{C} is

$$\text{HB}_\varepsilon(\mathcal{C}) := \left\{ P \in \mathcal{P} : \left| \Pr[X_i = 1 \mid \text{nbr}_r(\Phi, i) = P] - \frac{1}{2} \right| > \varepsilon \right\}.$$

If (Φ, i) matches a $P \in \text{HB}_\varepsilon(\mathcal{C})$, we say that \mathcal{C} *attains bias* $> \varepsilon$ *on* (Φ, i) .

Remark 5.7. For a fixed local per-bit rule $h(\mathbf{z}, a_i, b)$, the relevant chart is obtained by taking \mathcal{P} to be the set of all signed radius- r patterns (with labels) and setting $\psi(P) := h(\mathbf{z}(P), a_i(P), b(P))$.

5.3 Sparsification at $r = c_3 \log m$

We now bound the probability that a fixed chart is matched by a random masked block and simultaneously lands in its high-bias region.

Lemma 5.8 (Chart probability bound). *For any fixed chart $\mathcal{C} = (\mathcal{P}, \psi)$ and any $\varepsilon > 0$,*

$$\Pr_{\Phi \sim \mathcal{D}_m, i \sim [m]} [(\Phi, i) \text{ matches some } P \in \text{HB}_\varepsilon(\mathcal{C})] \leq m^{-\beta''}$$

for some $\beta'' = \beta''(\alpha, c_3) > 0$.

Proof sketch. By Theorem 3.11(iii), each fixed signed rooted pattern P occurs as $\text{nbr}_r(\Phi, i)$ with probability $\leq m^{-\beta'}$, and there are only $m^{O(1)}$ patterns of depth $r = c_3 \log m$ up to isomorphism (since the branching factor is constant). Labels (a_i, b) have entropy $\Theta(\log m)$ and contribute at most a polynomial factor to the total number of augmented patterns. Hence $\Pr[(\Phi, i) \text{ matches } P] \leq m^{-\beta''}$ for each P , and a union bound over the finite set $\text{HB}_\varepsilon(\mathcal{C})$ yields the claim. \square

Lemma 5.9 (Few high-bias hits for a fixed chart). *Let $t = c_4 m$. Draw i.i.d. blocks $(\Phi_1, \dots, \Phi_t) \sim \mathcal{D}_m^{\otimes t}$ and pick i_j uniformly from $[m]$ for each block. For any fixed chart \mathcal{C} , the number of indices $j \in [t]$ for which (Φ_j, i_j) matches a $P \in \text{HB}_\varepsilon(\mathcal{C})$ is at most $o(t)$ with probability $1 - 2^{-\Omega(m)}$.*

Proof. For each j , the indicator of the event in question is a Bernoulli with mean $\leq m^{-\beta''}$ by Lemma 5.8. Independence across blocks and Chernoff bounds imply that the total count is $O(tm^{-\beta''} + \log m)$ with probability $1 - 2^{-\Omega(m)}$. Since $t = \Theta(m)$ and $\beta'' > 1$ for small enough c_3 , this is $o(t)$. \square

Theorem 5.10 (Template sparsification for the finite local alphabet). *Fix $\varepsilon > 0$ and let \mathcal{U} be the set of possible local inputs $\mathbf{u} = (\mathbf{z}, a_i, b)$. There exists $\beta > 1$ such that for a random block $\Phi \sim \mathcal{D}_m$ and a uniform bit $i \in [m]$,*

$$\Pr \left[\exists \mathbf{u} \in \mathcal{U} \text{ with } \mathbf{u}_i(\Phi) = \mathbf{u} \text{ and } \left| \Pr[X_i = 1 \mid \mathbf{u}] - \frac{1}{2} \right| > \varepsilon \right] \leq m^{-\beta}.$$

Consequently, for $t = c_4 m$ blocks, with probability $1 - 2^{-\Omega(m)}$, at most $o(t)$ blocks admit any i and any \mathbf{u} that is ε -high-bias.

Proof sketch. Fix $\mathbf{u} = (\mathbf{z}, a_i, b)$. The event " $\mathbf{u}_i(\Phi) = \mathbf{u}$ and $|\Pr[X_i = 1 \mid \mathbf{u}] - 1/2| > \varepsilon$ " requires the radius- $r = c_3 \log m$ signed neighborhood around i to match one of a finite set of signed charts whose conditional bias exceeds ε (the VV labels contribute $O(\log m)$ bits). By Theorem 3.11, each such signed chart has probability $m^{-\Omega(1)}$. Since $|\mathcal{U}| = m^{O(1)}$ (Def. 4.3), a union bound over $\mathbf{u} \in \mathcal{U}$ gives $m^{-\beta}$ for some $\beta > 1$. Independence across blocks and Chernoff yield the $o(t)$ claim. \square

5.4 Many locally hard blocks after switching

We now combine Theorem 4.2 with Theorem 5.10 to obtain the *locally hard blocks* property required in Section 6.

Corollary 5.11 (Locally hard blocks). *There exist constants $\gamma > 0$ and a function $\varepsilon(m) \rightarrow 0$ such that for any polynomial-time decoder P with $|P| \leq \delta t$, there is a wrapper $W = W_{\text{ERM}}$ with $|W_{\text{ERM}}| \leq |P| + O(\log m + \log t)$ and a set $S \subseteq [t]$ with $|S| \geq \gamma t$ for which:*

$$\forall j \in S \forall i \in [m] : \quad \left| \Pr \left[(P \circ W_{\text{ERM}})(\Phi)_{j,i} = X_{j,i} \right] - \frac{1}{2} \right| \leq \varepsilon(m).$$

Proof. By Theorem 4.2 and Proposition A.5, after applying the ERM wrapper W_{ERM} there is a *test* subset $S_0 \subseteq [t]$ with $|S_0| \geq \gamma_0 t$ on which locality holds:

$$(P \circ W_{\text{ERM}})(\Phi)_{j,i} = h_{j,i}(\mathbf{z}(\Phi_j), a_{j,i}, b_j) \quad \text{for some } h_{j,i} : \mathcal{U} \rightarrow \{0, 1\} \quad (j \in S_0, i \in [m]).$$

Theorem 5.10 applies to all \mathbf{u} -measurable rules and (together with neutrality) yields that *all but* $o(t)$ of the blocks in S_0 satisfy the stated per-bit bound simultaneously for *all* $i \in [m]$. Let $S \subseteq S_0$ be the resulting subset; then $|S| \geq \gamma t$ for some constant $\gamma > 0$, as claimed. \square

What Section 5 provides downstream. Corollary 5.11 supplies the per-bit near-randomness on a γ -fraction of blocks for *every* short decoder, which is the exact hypothesis needed in Section 6 to invoke the Milestone-1 single-block lower bounds (Lemma 2.15) and obtain an exponential decay of per-program success across blocks.

6 Per-Program Small Success and Tuple Incompressibility

In this section we aggregate: independence across blocks turns local near-randomness into exponential decay of a short decoder's success. Then Compression-from-Success converts small success into a *linear* lower bound on K^{poly} for the whole witness tuple. This is Milestone M3.

Specifically: we convert the *local* hardness guaranteed by Switching-by-Weakness (Theorem 4.2) and the neutrality/sparsification results of Section 5 into a *global* (per-program) small-success bound across $\Theta(m)$ independent blocks. A standard counting/union bound (or, equivalently, Compression-from-Success) then yields a linear lower bound on K^{poly} for the witness tuple.

Throughout, $t = c_4 m$ for a fixed constant $c_4 > 0$, and $\varepsilon(m) \rightarrow 0$ denotes a vanishing bias bound supplied by Theorem 5.10.

6.1 From local hardness to block-level success bounds

Fix a polynomial-time decoder P of description length $|P| \leq \delta t$. By Theorem 4.2 (Switching-by-Weakness) and Proposition A.5, there exists a *distillation wrapper* W_{ERM} with $|W_{\text{ERM}}| \leq |P| + O(\log m + \log t)$ and a set $S_0 \subseteq [t]$ with $|S_0| \geq \gamma_0 t$ such that, for every $j \in S_0$ and $i \in [m]$,

$$(P \circ W_{\text{ERM}})(\Phi)_{j,i} = h_{j,i}(\mathbf{z}(\Phi_j), a_{j,i}, b_j) \quad (h_{j,i} : \mathcal{U} \rightarrow \{0, 1\}).$$

By Theorem 5.10, there exists $S \subseteq S_0$ with $|S| \geq \gamma t$ such that, simultaneously for all $j \in S$ and all $i \in [m]$,

$$\left| \Pr[(P \circ W_{\text{ERM}})(\Phi)_{j,i} = X_{j,i}] - \frac{1}{2} \right| \leq \varepsilon(m). \quad (2)$$

By Corollary 4.6, it suffices to upper bound the success of $(P \circ W_{\text{ERM}})$, since $\Pr[P(\Phi) = X] \leq \Pr[(P \circ W_{\text{ERM}})(\Phi) = X] + m^{-\Omega(1)}$ for this same wrapper.

(Here and below, probabilities are taken over the random test block $\Phi_j \sim \mathcal{D}_m$ with the wrapper W_{ERM} (split, seeds, trained $\{\hat{h}_i\}$) held fixed. Independence across $j \in S$ then follows from Lemma 6.6 together with the i.i.d. block product, Definition 3.5.)

Pivot bound. For any algorithm A and block j and any chosen pivot i^* , $\{A(\Phi_j) = X_j\} \subseteq \{A(\Phi_j)_{i^*} = X_{j,i^*}\}$, hence $\Pr[A(\Phi_j) = X_j] \leq \Pr[A(\Phi_j)_{i^*} = X_{j,i^*}]$.

We now turn (2) into a *block-level* bound.

Lemma 6.1 (Block correctness is bounded by any single-bit correctness). *For any algorithm A and any block j ,*

$$\Pr[A(\Phi_j) = X_j] \leq \Pr[A(\Phi_j)_{i^*} = X_{j,i^*}] \quad \text{for every chosen pivot } i^* \in [m].$$

Proof. The event $\{A(\Phi_j) = X_j\}$ implies the event $\{A(\Phi_j)_{i^*} = X_{j,i^*}\}$. \square

Proposition 6.2 (Per-block success bound on S). *Let $i^* \in [m]$ be any fixed pivot coordinate (e.g., $i^* = 1$). For every $j \in S$,*

$$\Pr[(P \circ W_{\text{ERM}})(\Phi_j) = X_j] \leq \frac{1}{2} + \varepsilon(m).$$

Proof. Apply Lemma 6.1 with $A = P \circ W_{\text{ERM}}$ and the pivot i^* , then use (2) for $i = i^*$. \square

Remark 6.3 (Why we use a pivot bit and not a bit-product bound). After switching, each per-bit rule $h_{j,i}$ shares the block-level inputs (\mathbf{z}_j, b_j) with all other bits, and the target bits $X_{j,i}$ are coupled by both the CNF constraints and the VV equations $Ax = b$. Hence, in general the events $\{(P \circ W_{\text{ERM}})(\Phi_j)_i = X_{j,i}\}_{i=1}^m$ are not independent and can be highly correlated. Without an additional independence/anti-concentration hypothesis, $\Pr[\text{all } m \text{ bits correct}]$ need not factor as a product over i ; the worst-case upper bound is the pivot-bit bound used in Proposition 6.2.

By Corollary 4.6, it suffices to upper bound the success of the comparator $(P \circ W_{\text{ERM}})$, since $\Pr[P(\Phi) = X] \leq \Pr[(P \circ W_{\text{ERM}})(\Phi) = X] + m^{-\Omega(1)}$ for the same W_{ERM} .

Theorem 6.4 (Fine-grained small success: bitwise form). *Let P be any polynomial-time decoder with $|P| \leq \delta t$ and let W be the SW wrapper from Theorem 4.5. For the subset S of size $\geq \gamma t$ on which locality holds, with probability $1 - 2^{-\Omega(m)}$ we have*

$$\sum_{j \in S} \sum_{i=1}^m \mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi)_{j,i} = X_{j,i}\} \leq \left(\frac{1}{2} + \varepsilon(m)\right) m |S| + o(m|S|).$$

Proof sketch. For each fixed (j, i) with locality, neutrality/sparsification implies $\Pr[(P \circ W_{\text{ERM}})_{j,i} = X_{j,i}] \leq \frac{1}{2} + \varepsilon(m)$. By independence across blocks (Lemma 6.6) and linearity of expectation plus Chernoff, the sum over $j \in S$ concentrates around its mean, yielding the stated upper tail bound. \square

Corollary 6.5 (Enumerative coding from bitwise small success). *Combining Theorem 6.4 with Lemma 2.6 yields*

$$K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t$$

for some constant $\eta > 0$, even if the decoder only partially recovers witnesses with arbitrary adaptive strategies.

Quantifier Order and Independence

We proceed as: $\forall P \exists W_{\text{ERM}}$ (fix train/test split, seeds, and trained rules), then \forall fixed W_{ERM} we analyze fresh test blocks. Conditioned on W_{ERM} , the random variables $\{\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}\}_{j \in S}$ are independent because each depends only on its own i.i.d. test block Φ_j (Def. 3.5, Lemma 6.6).

6.2 Exponential decay across independent blocks

Once the ERM wrapper W_{ERM} is fixed (train/test split, seeds, trained $\{\hat{h}_i\}$), the block-level correctness events on the *test* subset S ,

$$\{\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}\}_{j \in S},$$

are independent: each depends only on the independent test block Φ_j (Definition 3.5) with the wrapper held fixed (Lemma 6.6). By Proposition A.5, we also have *success domination*:

$$\Pr [P(\Phi) = X] \leq \Pr [(P \circ W_{\text{ERM}})(\Phi) = X] + m^{-\Omega(1)}.$$

That is, just to be clear: Conditioned on the fixed wrapper W_{ERM} (seeds, split, and trained tables), each indicator $\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}$ is a function only of the *test* block Φ_j with $j \in S$, and is independent across j by Definition 3.5 and Lemma 6.6.

Lemma 6.6 (Conditional independence given a fixed wrapper). *Fix a wrapper W (including its seeds and, if $W = W_{\text{ERM}}$, also the training/test split and trained rules). Then, conditional on W , the random variables $\{\mathbf{1}\{(P \circ W)(\Phi_j) = X_j\}\}_{j \in S}$ are independent, since each depends only on the corresponding independent block Φ_j .*

Combining locality on S (Theorem 4.2 / Proposition A.5), per-bit near-randomness for \mathbf{u} -measurable predictors (Theorem 5.10 and neutrality), and the pivot inequality (Lemma 6.1), we obtain for each $j \in S$:

$$\Pr [(P \circ W_{\text{ERM}})(\Phi_j) = X_j] \leq \frac{1}{2} + \varepsilon(m).$$

By independence across $j \in S$ (this subsection), the product bound yields

$$\Pr [(P \circ W_{\text{ERM}})(\Phi) = X \text{ on all } j \in S] \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{|S|} \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{\gamma t}.$$

Finally, success domination transfers this bound (up to $m^{-\Omega(1)}$ slack) to $\Pr[P(\Phi) = X]$.

Quantifier order reminder. The argument proceeds as: $\forall P \exists W_{\text{ERM}}$ (Switching-by-Weakness/distillation on S), then $\forall W_{\text{ERM}}$ (product small-success bound), and finally lifts the bound back to P via success domination. Thus the final upper bound holds for *all* short decoders P .

Theorem 6.7 (Per-program small-success bound). *There exists a function $\varepsilon(m) \rightarrow 0$ and a constant $\gamma > 0$ such that, for every polynomial-time decoder P with $|P| \leq \delta t$, there is an ERM wrapper W_{ERM} with $|W_{\text{ERM}}| \leq |P| + O(\log m + \log t)$ for which*

$$\Pr [P(\Phi_1, \dots, \Phi_t) = (X_1, \dots, X_t)] \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{\gamma t} + m^{-\Omega(1)} = 2^{-\Omega(t)}.$$

Proof. By Proposition A.5 there is a test subset $S \subseteq [t]$, $|S| \geq \gamma t$, on which $(P \circ W_{\text{ERM}})_{j,i} = \hat{h}_i(\mathbf{u}_{j,i})$ is local. By Theorem 5.10 and neutrality, for every $j \in S$, $\Pr[(P \circ W_{\text{ERM}})(\Phi_j) = X_j] \leq \frac{1}{2} + \varepsilon(m)$. Conditioned on the fixed wrapper, the events $\{\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}\}_{j \in S}$ are independent (Lemma 6.6), so

$$\Pr [(P \circ W_{\text{ERM}})(\Phi) = X \text{ on all } j \in S] \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{|S|} \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{\gamma t}.$$

Correctness on all t blocks implies correctness on S , so the same upper bound holds for $\Pr[(P \circ W_{\text{ERM}})(\Phi) = X]$. Finally, success domination (Proposition A.5 (ii)) gives $\Pr[P(\Phi) = X] \leq \Pr[(P \circ W_{\text{ERM}})(\Phi) = X] + m^{-\Omega(1)}$, which yields the stated inequality. \square

6.3 From small success to tuple incompressibility

We now convert Theorem 6.7 into a lower bound on $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t))$. We give two equivalent routes: a direct union bound over short programs, and a reference to Compression-from-Success (Lemma 2.5 / Lemma 2.6).

Route A: direct counting. Fix $L = \eta t$. The number of decoders of description length $\leq L$ is at most 2^L . By Theorem 6.7, each such decoder has success probability at most $(\frac{1}{2} + \varepsilon(m))^{\gamma t}$. Hence

$$\begin{aligned} \Pr [\exists P : |P| \leq L \wedge P(\Phi) = X] &\leq 2^L \cdot \left(\frac{1}{2} + \varepsilon(m)\right)^{\gamma t} \\ &= 2^{-\left(\gamma \log_2 \frac{1}{1/2 + \varepsilon(m)} - \eta\right) t}. \end{aligned}$$

Choose a constant $\eta > 0$ smaller than $\gamma \log_2 \left(\frac{1}{1/2 + \varepsilon(m)}\right)$ (for all large m) to obtain

$$\Pr [\exists P : |P| \leq \eta t \wedge P(\Phi) = X] \leq 2^{-\Omega(t)}.$$

Equivalently, with probability $1 - 2^{-\Omega(t)}$, $K^{\text{poly}}(X | \Phi) \geq \eta t$.

Route B: Compression-from-Success. Fix $L = \eta t$ as above. Suppose, for contradiction, that with probability $> 2^{-\Omega(t)}$ we had $K^{\text{poly}}(X | \Phi) < L$. Then, by definition of K^{poly} , there exists a decoder P of length $< L$ that succeeds on those instances. But Theorem 6.7 bounds the success probability of *every* such decoder by $2^{-\Omega(t)}$, contradiction. Alternatively, apply Lemma 2.5/2.6 to turn any putative success probability into a code of length $< L$ and compare.

We summarize the outcome as the main lower bound for this section.

Theorem 6.8 (Tuple incompressibility). *There exists a constant $\eta > 0$ such that, for $t = c_4 m$,*

$$\Pr_{(\Phi_1, \dots, \Phi_t) \sim \mathcal{D}_m^{\otimes t}} \left[K^{\text{poly}}((X_1, \dots, X_t) | (\Phi_1, \dots, \Phi_t)) \geq \eta t \right] \geq 1 - 2^{-\Omega(m)}.$$

Proof. Immediate from Route A (direct counting) with η chosen as above, or from Route B using Lemma 2.5/2.6. \square

6.4 Constants and parameter choices

Admissible parameter choices (union-bound exponent). Let $\varepsilon(m) = m^{-c}$ from sparsification and let $\gamma \in (0, 1)$ be the switching fraction. Write

$$\Lambda(m) := \log_2 \left(\frac{1}{1/2 + \varepsilon(m)} \right) \quad \text{so that} \quad \Lambda(m) \rightarrow 1.$$

For any target $\eta > 0$ and length budget $\delta > 0$, the union bound exponent is

$$\delta + \gamma \log_2 \left(\frac{1}{2} + \varepsilon(m) \right) = \delta - \gamma \Lambda(m).$$

Hence it suffices to choose η, δ so that, for all large m ,

$$\delta \leq \gamma \Lambda(m) - \eta. \tag{3}$$

Two equivalent ways to fix constants are:

- **Concrete choice.** Take $\eta := \gamma/4$ and $\delta := \gamma/8$. Since $\Lambda(m) \rightarrow 1$, we have $\delta \leq \gamma - \eta$ for all large m , so (3) holds and

$$2^{\delta t} \left(\frac{1}{2} + \varepsilon(m) \right)^{\gamma t} \leq 2^{-(\eta - o(1))t} \leq 2^{-\eta t}$$

for $t = c_4 m$ and m large enough.

- **Symbolic choice.** Fix any $\eta \in (0, \gamma)$ with $\eta \leq \frac{\gamma}{2}\Lambda(m)$ for all large m (e.g., any constant $\eta < \gamma/2$). Then set

$$\delta := \frac{1}{2}(\gamma\Lambda(m) - \eta) > 0.$$

This choice satisfies (3) and yields the same $2^{-\eta t}$ tail.

In either case, the number of decoders of length $\leq \delta t$ is at most $2^{\delta t}$, so the union bound gives

$$\Pr[\exists P : |P| \leq \delta t \wedge P \text{ success on all } t \text{ blocks}] \leq 2^{-\eta t} = 2^{-\Omega(t)} = 2^{-\Omega(m)}.$$

What Section 6 delivers downstream. Theorem 6.8 is the linear lower bound on K^{poly} for the witness tuple that Section 7 pits against the *constant-length* upper bound under $P = NP$ (Proposition 7.2), completing the quantale upper–lower clash.

7 Quantale Upper-Lower Clash and Main Theorem

Here we close the loop (Milestone M4). The *lower* side is the tuple incompressibility from Section 6 (Theorem 6.8): with high probability, any program that outputs the full witness tuple must have length $\Omega(t)$ when $t = c_4 m$. The *upper* side assumes $P = NP$ and observes that there is a *uniform, constant-length* program that, on input any on-promise instance(s), outputs the unique witness(es) in polynomial time by bit-fixing with a USAT decoder. Hence

$$K^{\text{poly}}(X \mid \Phi) \leq O(1) \quad \text{and} \quad K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq O(1),$$

which contradicts the $\Omega(t)$ lower bound for large t .

Distributional lower vs. universal upper. Rephrasing just to be pedantically clear, note that: The lower bound is distributional: with probability $1 - 2^{-\Omega(m)}$ over $(\Phi_1, \dots, \Phi_t) \sim D_m^{\otimes t}$, we have $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t$. Under $P = NP$, the self-reduction yields a *uniform* constant-length decoder for the promise, so $K^{\text{poly}}(\cdot \mid \cdot) \leq O(1)$ for *every* input. For large m these statements are incompatible.

7.1 Self-reduction for USAT under $P = NP$

Recall \mathcal{D}_m is supported on instances $\Phi = (F^h, A, b)$ that *have a unique* satisfying assignment $X \in \{0, 1\}^m$ (Definition 3.4). Under $P = NP$, USAT is decidable in polynomial time, and the classical bit-fixing recipe recovers X in m queries while preserving the promise at each step.

Lemma 7.1 (Bit-by-bit self-reduction under $P = NP$). *Assume $P = NP$. There exists a polynomial-time decision procedure D_{USAT} for $\text{USAT} = \{\varphi : \#\text{SAT}(\varphi) \in \{0, 1\}\}$ such that, for any on-promise φ with unique witness $x \in \{0, 1\}^m$, one obtains x by m calls to D_{USAT} on bit-fixing restrictions. At each step the restricted instance remains on-promise.*

Proposition 7.2 (Uniform constant-length witness finder under $P = NP$). *Assume $P = NP$. There exists a constant C (independent of m, t) and a fixed program p of length $\leq C$ such that, for every on-promise block Φ with unique witness X ,*

$$K^{\text{poly}}(X \mid \Phi) \leq C,$$

and for every t and every on-promise tuple (Φ_1, \dots, Φ_t) with witnesses (X_1, \dots, X_t) ,

$$K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq C.$$

Proof. Hard-wire into p a polynomial-time USAT decider D_{USAT} (exists under $P = NP$) and the standard bit-fixing routine of Lemma 7.1. On input Φ , p parses m from Φ and runs m queries to D_{USAT} on the appropriate restrictions to recover X . For tuples, p parses the self-delimiting encoding of (Φ_1, \dots, Φ_t) and loops over blocks. The running time is polynomial in the input length, and the program length is constant. \square

7.2 Lower vs. upper: the quantale clash

We restate the lower bound from Section 6:

Theorem 7.3 (Tuple incompressibility, restated). *There exists $\eta > 0$ such that, for $t = c_4 m$,*

$$\Pr \left[K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \geq \eta t \right] \geq 1 - 2^{-\Omega(m)}.$$

Combining Proposition 7.2 (upper bound under $P = NP$) with Theorem 7.3 (lower bound) yields the contradiction for large t .

Theorem 7.4 (Main Separation). *For the masked-and-isolated block distribution \mathcal{D}_m and $t = c_4 m$ i.i.d. blocks,*

$$P \neq NP.$$

Proof. Assume $P = NP$. By Proposition 7.2, $K^{\text{poly}}((X_1, \dots, X_t) \mid (\Phi_1, \dots, \Phi_t)) \leq C$ for every outcome, while by Theorem 7.3 the same quantity is $\geq \eta t$ with probability $1 - 2^{-\Omega(m)}$. For sufficiently large t , these inequalities are incompatible. Contradiction. Therefore $P \neq NP$. \square

7.3 Non-relativizing and non-naturalizing aspects

Non-relativizing (methodological). Our derivation depends essentially on explicit properties of the sampling law (uniform masking by H_m and local sparsity of random 3-CNF) and on in-sample verification inside the USAT promise. The argument is not phrased as an oracle-independent simulation and we make no claim that it relativizes; rather, it is distribution-specific and verifier-dependent. Establishing an explicit oracle separation for this technique is an interesting open direction.

Non-naturalizing. The lower bound is a *per-program* small-success statement tied to a specific, efficiently samplable distribution and a *polynomial-size* post-switch local alphabet; it is not a dense, constructive property of all Boolean functions. Hence it avoids the Razborov-Rudich natural-proofs barrier.

7.4 Parameters and constants (consolidated)

- Clause density $\alpha > 0$; mask $h \sim H_m$ fresh per block.
- VV layer: $k = c_1 \log m$, $\delta = m^{-c_2}$; isolation succeeds with $\Omega(1/m)$ probability and we condition on uniqueness.
- SILS length: $r(m) = O(\log m)$; computable in $\text{poly}(m)$; sign-invariant.
- Radius: $r = c_3 \log m \in (0, c_3^*(\alpha))$ to guarantee local tree-likeness.
- Blocks: $t = c_4 m$; independence across blocks.

- Switching: constants $\gamma > 0$, $c^* > 0$ from Theorem 4.2.
- Sparsification: bias bound $\varepsilon(m) = m^{-\Omega(1)}$ on a γ -fraction of blocks (Theorem 5.10).
- Tuple lower bound: $\eta > 0$ from Theorem 7.3; upper bound constant C from Proposition 7.2.

8 Discussion and Open Problems

The previous section completed the proof of $P \neq NP$ which is the crux of the paper. We have shown separation of P and NP based on a compact calculus: shortness \Rightarrow locality (switching-by-weakness), plus symmetry and sparsity \Rightarrow near-randomness on many blocks, plus independence \Rightarrow exponential decay, plus compression-from-success \Rightarrow tuple incompressibility, which clashes with self-reduction under $P = NP$.

We hope the modular structure we have leveraged in this proof encourages further refinements and broader applications. In the remainder of this section we conclude by briefly discussing future directions for the methods and ideas we have used – robustness, limitations, and potential ways to strengthen and generalize the separation.

8.1 Robustness of the ensemble and parameters

Our masked-and-isolated block ensemble \mathcal{D}_m is deliberately minimal: it uses only (i) constant-density random 3-CNF, (ii) a fresh $H_m = S_m \ltimes (\mathbb{Z}_2)^m$ mask per block, (iii) an $O(\log m)$ -bit VV isolation layer with pairwise-independent columns and δ -biased right-hand-side, and (iv) a short sign-invariant SILS extractor. The proof needs only:

1. *Sign-invariant SILS* of length $O(\log m)$, computable in $\text{poly}(m)$ (Def. 2.7);
2. *Promise-preserving* sign-flips (Lemma 3.6);
3. *Local tree-likeness* at radius $r = c_3 \log m$ (Thm. 3.11);
4. Post-switch rules with $O(\log m)$ *inputs* (Thm. 4.2).

Constants $c_1, c_2, c_3, c_4, \delta, \gamma$ can be varied in wide ranges as long as these invariants hold.

8.2 Why masking, isolation, and SILS

Masking. The fresh H_m mask per block enforces distributional symmetry used twice: (i) per-bit AP-GCT neutrality for sign-invariant views, and (ii) uniformity of signed neighborhoods for sparsification at radius $c_3 \log m$. Without masking, an adversarial naming or literal-sign bias could correlate with local features and spoil neutrality.

Isolation. The VV layer ensures uniqueness and keeps the local VV labels (a_i, b) at $O(\log m)$ bits, which is critical for (1) the switching normal form (local input length) and (2) the sparsification bound (finite chart universe).

SILS. We use SILS only as an H_m -invariant, short, polytime summary; no special ENF/CENF structure is needed. This keeps the post-switch per-bit domain logarithmic while exposing enough low-degree structure for neutrality and sparsification.

8.3 On non-relativization and non-naturalization

The argument is *non-relativizing*: it uses the concrete sampling law (masking), in-sample verification within the USAT promise, and switching wrappers that apply promise-preserving automorphisms. The lower bound is *non-natural*: it is a per-program small-success statement specific to an efficiently samplable distribution and a polynomial post-switch alphabet, not a dense constructive property on all Boolean functions.

Non-natural and non-relativizing. That is: our lower bound is per-program, distribution-specific, and verifier-dependent; it is neither dense nor constructive in the sense of Razborov-Rudich, and it is proved using ensemble symmetries that do not relativize.

8.4 Open problems

OP1: Removing or weakening the mask. To what extent can one reduce the mask randomness (e.g., only random signs; or a fixed permutation reused across blocks) while retaining neutrality and sparsification? A plausible first target is masking by $(\mathbb{Z}_2)^m$ only (random literal signs without variable permutation).

OP2: Beyond radius $c_3 \log m$. Our sparsification uses local tree-likeness at logarithmic radius. Can one push sparsification to polylogarithmic radius or to a Fourier low-degree regime for random k -SAT factor graphs, to obtain a more analytic (LMN-style) algorithmic Pinsker?

OP3: Alternative ensembles. The same pipeline should apply to other sparse CSPs (random k -XOR, planted models with noise, Goldreich-type predicates) with an appropriate SILS extractor and promise-preserving symmetries.

OP4: Derandomizing the switching wrapper. We gave two wrappers: ERM and symmetrization. The ERM wrapper is already randomness-free beyond sampling the i.i.d. blocks; the symmetrization wrapper uses polylogarithmic independent sign flips. Tighten the concentration under even smaller independence, or make the wrapper seedless by a canonicalization trick.

OP5: Strengthening per-block lower bounds. We invoked tiny ACC^0 /streaming bounds on $O(\log m)$ inputs. It would be interesting to prove direct correlation bounds *for the switched per-bit class itself* against the signed neighborhood distribution, yielding a purely distributional per-block lower bound.

OP6: Toward unmasked natural distributions. With more delicate SILS and possibly an *a priori* de-biasing step, the neutrality argument may carry over to (partially) unmasked ensembles. This requires characterizing which low-degree invariants remain uncorrelated with isolated witness bits in the unmasked law.

OP7: Categorical formalization. We sketched the quantale viewpoint informally: K^{poly} as a lax monoidal functor enforcing additive budgets under block product; sign-invariant SILS as an invariant functor; promise-preserving automorphisms as measure-preserving endomorphisms. A categorical write-up would likely clarify portability to other ensembles.

OP8: Learnability and meta-complexity. Our ERM wrapper exploits the polynomial size of the post-switch alphabet. A sharper uniform convergence analysis (e.g., via Rademacher averages) may reduce sample fractions and improve constants. Connecting the small-success statement to explicit meta-complexity assumptions (e.g., KT-decision) remains an appealing alternative route to hardness.

A Detailed Proofs of Key Components

Here we run through a few of the technical proofs given in the paper in more detail.

A.1 Switching-by-Weakness via Distillation

We prove Theorem 4.2 using an ERM (Empirical Risk Minimization) wrapper that *distills* any polynomial-time decoder P down to a local comparator $h(\mathbf{u})$ on the distribution Dm without assuming any per-instance measurability.

Clarification. This section does not claim that an arbitrary polynomial-time decoder P is itself local. Instead, for each such P we construct a short, promise-preserving *comparator* $(P \circ W)$ whose per-bit outputs on a large test subset are *functions of the local inputs* $\mathbf{u} = (\mathbf{z}, a_i, b)$, and we prove a *success-domination* inequality

$$\Pr [P(\Phi) = X] \leq \Pr [(P \circ W)(\Phi) = X] + m^{-\Omega(1)}.$$

This lets us upper bound the success of *every* P via an analyzable local comparator.

Group action and back-map. Let $\mathcal{G} \leq H_m$ be the subgroup of componentwise sign flips; write $\mathcal{G} \cong (\mathbb{Z}_2)^m$. For $\sigma \in \mathcal{G}$, define the promise-preserving bijection (Lemma 3.6)

$$g_\sigma : (F^h, A, b) \mapsto (F^{(\text{id}, \sigma)^h}, A, b \oplus A\sigma).$$

For block j and bit i we define the *back-mapped* prediction

$$Y_{j,i}(\sigma, \Phi) := (P(g_\sigma(\Phi)))_{j,i} \oplus \langle a_{j,i}, \sigma \rangle,$$

so that (by construction of g_σ) comparing $Y_{j,i}(\sigma, \Phi)$ to the *original* target $X_{j,i}$ is meaningful. The *local input* is $\mathbf{u}_{j,i} = (\mathbf{z}(\Phi_j), a_{j,i}, b_j) \in \{0, 1\}^{O(\log m)}$.

Promise-conditionalization and off-promise slack. All probabilities and expectations in this appendix are taken under the law D_m *conditioned on* uniqueness (USAT promise). Conceptually, the sampler implements rejection sampling of the VV stage until uniqueness holds; this preserves the distribution on the promise space. If one prefers to sample (A, b) from a δ -biased source instead of uniform, then for any fixed σ , the map $b \mapsto b \oplus A\sigma$ changes the law by at most $O(\delta)$ in total variation. Throughout we absorb such deviations into the global slack term, which we set to $m^{-\Omega(1)}$ by choosing $\delta \leq m^{-10}$.

Two-level wrapper. We build two short wrappers:

- W_{sym} (*symmetrization*): produces per-bit labels by averaging P over $s = \Theta(\log(mt))$ sign flips drawn from a κ -wise independent family with $\kappa = \Theta(\log(mt))$, then taking a majority.
- W_{ERM} (*distillation to locality*): learns per-bit local rules on a train split and predicts on a disjoint test split using only $\mathbf{u} = (\mathbf{z}, a_i, b)$ as inputs.

We now formalize both and prove success domination and locality.

(A) Symmetrization and success domination

Definition (symmetrized label). Fix $s = \Theta(\log(mt))$ and $\kappa = \Theta(\log(mt))$. Draw $\sigma^{(1)}, \dots, \sigma^{(s)}$ from a κ -wise independent family on \mathcal{G} and define

$$\tilde{Y}_{j,i} := \text{Maj}\left(Y_{j,i}(\sigma^{(1)}, \Phi), \dots, Y_{j,i}(\sigma^{(s)}, \Phi)\right).$$

Let W_{sym} be the wrapper that, on any input, outputs the bit-vector whose (j, i) -entry is $\tilde{Y}_{j,i}$.

Lemma A.1 (Concentration of the majority). *There exists $\varepsilon_0(m) = m^{-\Omega(1)}$ such that for all (j, i) ,*

$$\left| \Pr[\tilde{Y}_{j,i} = X_{j,i}] - \mathbb{E}_\sigma \Pr[Y_{j,i}(\sigma, \Phi) = X_{j,i}] \right| \leq \varepsilon_0(m).$$

Proof. Condition on (Φ, j, i) and write $p := \Pr_\sigma[Y_{j,i}(\sigma, \Phi) = X_{j,i}]$. Under κ -wise independence with $\kappa = \Theta(\log(mt))$, limited-independence Chernoff [8, 9] gives that the empirical average of the $\{0, 1\}$ indicators $\mathbf{1}\{Y_{j,i}(\sigma^{(r)}, \Phi) = X_{j,i}\}$ deviates from p by at most $1/\text{poly}(m)$ with probability $1 - \varepsilon_0(m)$. Majority has accuracy $\geq \max\{p, 1 - p\} \geq p$ up to this deviation. Take expectations over Φ to conclude. \square

Lemma A.2 (Success domination by W_{sym}). *For any decoder P and any block j ,*

$$\Pr[P(\Phi_j) = X_j] = \Pr[\text{BM}_\sigma(P(g_\sigma(\Phi_j))) = X_j] \leq \Pr[(P \circ W_{\text{sym}})(\Phi_j)_{i^*} = X_{j,i^*}] + \varepsilon_0(m),$$

for an arbitrary fixed pivot $i^* \in [m]$. Hence, by Lemma 6.1,

$$\Pr[P(\Phi_j) = X_j] \leq \Pr[(P \circ W_{\text{sym}})(\Phi_j) = X_j] + \varepsilon_0(m).$$

Proof. By Lemma 6.1, block success is dominated by pivot-bit success. For the pivot bit, using Lemma A.1 and the exact success preservation from Lemma 4.3,

$$\Pr[P(\Phi_j)_{i^*} = X_{j,i^*}] = \mathbb{E}_{\Phi_j} \mathbb{E}_\sigma[\mathbf{1}\{\text{BM}_\sigma(P(g_\sigma(\Phi_j)))_{i^*} = X_{j,i^*}\}] = \mathbb{E}_{\Phi_j} \Pr_\sigma[Y_{j,i^*}(\sigma, \Phi_j) = X_{j,i^*}] \leq \Pr[\tilde{Y}_{j,i^*} = X_{j,i^*}].$$

This is exactly the stated bound for W_{sym} . \square

(B) Distillation to local rules via ERM

Train/test split. Choose a random partition $[t] = T \sqcup S$ with $|T|, |S| = \Theta(t)$. We use *only* the test split S in the small-success product bound; training serves to compute local rules.

Local alphabet and plug-in rules. Let \mathcal{U} be the local input alphabet, $|\mathcal{U}| = N = m^{O(1)}$. For each bit i , let $f_i(\mathbf{u}) := \mathbb{E}[Y_i(\sigma, \Phi) \mid \mathbf{u}]$ and let $h_i^*(\mathbf{u}) = \mathbf{1}[f_i(\mathbf{u}) \geq 1/2]$ be the Bayes classifier for the surrogate labels.

ERM training against symmetrized outputs. For each bit index $i \in [m]$, set the training labels to the symmetrized outputs on T : $\ell_{j,i} := \tilde{Y}_{j,i}$ for $j \in T$. Define the plug-in rule on the finite alphabet \mathcal{U} by

$$\hat{h}_i(\mathbf{u}) := \text{Maj}\{\tilde{Y}_{j,i} : j \in T, \mathbf{u}_{j,i} = \mathbf{u}\}.$$

Define the *ERM wrapper* W_{ERM} to output on test blocks $j \in S$ the local prediction

$$(P \circ W_{\text{ERM}})(\Phi)_{j,i} := \hat{h}_i(\mathbf{u}_{j,i}).$$

On training blocks $j \in T$ we simply output $P(\Phi_j)$ (this can only increase success).

Lemma A.3 (Plug-in ERM generalization on a finite alphabet). *With $|T| = \Theta(t) = \Theta(m)$ and the plug-in rule \hat{h}_i defined above, there exists $\varepsilon_0(m) = m^{-\Omega(1)}$ such that, with probability $1 - m^{-\Omega(1)}$ over the train/test split and the symmetrization seeds,*

$$\Pr_{j \in S} [\hat{h}_i(\mathbf{u}_{j,i}) \neq h_i^*(\mathbf{u}_{j,i})] \leq \varepsilon_0(m) \quad \text{simultaneously for all } i \in [m].$$

Proof sketch. For each $\mathbf{u} \in \mathcal{U}$, the training multiplicity $N_{\mathbf{u}} := |\{j \in T : \mathbf{u}_{j,i} = \mathbf{u}\}|$ has mean $|T| \Pr[\mathbf{u}]$. By (limited-independence) Chernoff, uniformly over \mathbf{u} we have $|N_{\mathbf{u}} - |T| \Pr[\mathbf{u}]| \leq O(\sqrt{|T| \Pr[\mathbf{u}] \log m})$ w.h.p. Conditional on $N_{\mathbf{u}}$, the empirical mean of \tilde{Y} at \mathbf{u} concentrates to $f_i(\mathbf{u})$ with deviation $\exp(-\Omega(N_{\mathbf{u}}))$. A union bound over $\mathbf{u} \in \mathcal{U}$ (size $N = \text{poly}(m)$) and over $i \leq m$ yields the claim; contributions of rare \mathbf{u} have small mass $\Pr[\mathbf{u}]$ and thus small effect on the test error. \square

Lemma A.4 (Distillation preserves success up to $m^{-\Omega(1)}$). *For the test split S ,*

$$\frac{1}{|S|} \sum_{j \in S} \mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\} \geq \frac{1}{|S|} \sum_{j \in S} \mathbf{1}\{(P \circ W_{\text{sym}})(\Phi_j) = X_j\} - m^{-\Omega(1)}.$$

Proof. For each $j \in S$, the two predictors differ on at most the event $(P \circ W_{\text{ERM}})(\Phi)_{j,i^*} \neq (P \circ W_{\text{sym}})(\Phi)_{j,i^*}$ for the pivot bit i^* (block success is dominated by pivot-bit correctness; Lemma 6.1). By Lemma A.3, the disagreement rate on the test split is $m^{-\Omega(1)}$, so the average block success degrades by at most that amount. \square

(C) Locality, independence, and conclusion

Locality on the test split. By construction, on every $j \in S$ and $i \in [m]$ the ERM predictor equals $(P \circ W_{\text{ERM}})(\Phi)_{j,i} = \hat{h}_i(\mathbf{u}_{j,i})$, a function of $O(\log m)$ inputs.

Independence across test blocks. Once the wrapper W_{ERM} is fixed (train/test split, seeds, and the trained $\{\hat{h}_i\}$), predictions on distinct test blocks depend only on the independent draws $\{\Phi_j\}_{j \in S}$. Hence $\{\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}\}_{j \in S}$ are independent (Lemma 6.6).

Proposition A.5 (Switching-by-Weakness (ERM version) with success domination). *Let P be any polynomial-time decoder with $|P| \leq \delta t$. There exists a short wrapper W_{ERM} of description length $|W_{\text{ERM}}| \leq |P| + O(\log m + \log t)$, a pivot bit i^* , and a test subset $S \subseteq [t]$ with $|S| \geq \gamma t$ such that:*

(i) (Locality) *For all $j \in S$ and $i \in [m]$, $(P \circ W_{\text{ERM}})(\Phi)_{j,i} = \hat{h}_i(\mathbf{u}_{j,i})$ for some plug-in rule \hat{h}_i on \mathcal{U} .*

(ii) (Success domination)

$$\frac{1}{|S|} \sum_{j \in S} \mathbf{1}\{P(\Phi_j) = X_j\} \leq \frac{1}{|S|} \sum_{j \in S} \mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\} + m^{-\Omega(1)}.$$

Proof. Combine Lemma A.2 (domination by W_{sym} on the pivot bit), Lemma A.4 (ERM preserves success up to $m^{-\Omega(1)}$), and Lemma 6.1 (pivot-to-block domination). Locality is by construction; description-length follows since seeds and split specification use $O(\log m + \log t)$ bits and training runs in polynomial time. \square

What this achieves for the global argument. Proposition A.5 provides, for every short P , a short wrapper producing a *local* comparator on a constant fraction of blocks whose success on the test split *dominates* that of P up to $m^{-\Omega(1)}$. Section 5 then applies neutrality and template sparsification to any \mathbf{u} -measurable per-bit rule, bounding the per-bit advantage by $\frac{1}{2} + \varepsilon(m)$, and Section 6 aggregates across the independent test blocks to obtain the per-program small-success bound.

Remark A.6 (Global invariants do not break the reduction). A decoder P may compute global, sign-invariant statistics of the masked formula. The ERM wrapper does not attempt to reproduce P 's global strategy; it distills the *symmetrized* behavior of P to a function of $\mathbf{u} = (\mathbf{z}, a_i, b)$. Any extra information P uses beyond \mathbf{u} can only improve P 's original success; our domination chain compares P first to the symmetrized comparator and then to its \mathbf{u} -measurable distillation on the test distribution, where ERM guarantees small imitation error. The lower bounds then apply to *all* such local comparators.

A.2 Weakness Quantale: formal calculus and interface

We record the algebra we use, emphasizing only the rules that are applied later.

Definition A.7 (Weakness cost and quantale). Let U be a fixed prefix-universal TM. Define

$$K^{\text{poly}}(z \mid y) := \min\{|p| : U(p, y) = z \text{ and } U \text{ halts in } |y|^{O(1)} \text{ steps}\}.$$

Set $Q := \mathbb{R}_{\geq 0} \cup \{\infty\}$ with addition as monoidal product and \leq as order.

Lemma A.8 (Invariance, chain rule, block additivity). *For all x, z, y : (i) $K_U^{\text{poly}}(x \mid y) \leq K_V^{\text{poly}}(x \mid y) + O(1)$ for any U, V ; (ii) $K^{\text{poly}}(xz \mid y) \leq K^{\text{poly}}(x \mid y) + K^{\text{poly}}(z \mid xy) + O(1)$; (iii) $K^{\text{poly}}(x_1 \cdots x_t \mid y_1 \cdots y_t) \leq \sum_i K^{\text{poly}}(x_i \mid y_i) + O(\log t)$.*

Proof. (i) Standard simulation with constant overhead; the time cap remains polynomial. (ii) Compose decoders and add separators; (iii) schedule subdecoders with an $O(\log t)$ loop. \square

Lemma A.9 (Compression-from-success, fine form). *Let $\hat{x}_j \in \{0, 1\}^m$ be predictions for x_j and E_j the bitwise error masks. Then*

$$K^{\text{poly}}(x_1 \cdots x_t \mid y_1 \cdots y_t) \leq L + O(\log t) + \sum_{j=1}^t \log \binom{m}{|E_j|},$$

where L is the description length of the predictor (including fixed coins).

Proof. Enumerate each error set and patch the predicted bits accordingly. \square

These suffice to turn per-program small success into linear tuple lower bounds.

A.3 Neutrality (exact 1/2, measure-theoretic proof)

Let \mathcal{I} be the σ -algebra generated by *sign-invariant*, permutation-invariant functions of F^h (e.g., the SILS coordinates). We show $\Pr[X_i = 1 \mid \mathcal{I}] = \frac{1}{2}$ a.s.

Lemma A.10 (Promise-preserving involution, measure version). *Define $T_i(F^h, A, b) := (F^{\tau_i h}, A, b \oplus Ae_i)$, where τ_i flips only variable i 's sign. Then T_i is a bijection on the promise space $\{\Phi : \#\text{SAT}(\Phi) = 1\}$, and the pushforward measure equals the original.*

Proof. As in Lemma 3.6, $x \mapsto x \oplus e_i$ bijects satisfying assignments; uniqueness is preserved. Uniformity of h and b implies measure preservation. \square

Theorem A.11 (Neutrality). *For every $i \in [m]$, $\Pr[X_i = 1 \mid \mathcal{I}] = \frac{1}{2}$ almost surely on the promise distribution.*

Proof. Let $B \in \mathcal{I}$. Since \mathcal{I} is sign-invariant, B is T_i -invariant. Because T_i toggles X_i , the sets $B \cap \{X_i = 1\}$ and $B \cap \{X_i = 0\}$ have equal measure (pair up ω with $T_i(\omega)$). Therefore $\Pr[X_i = 1 \mid \mathcal{I} = \text{value of } B] = 1/2$ for all atoms; extend by standard disintegration. \square

Corollary A.12 (SILS-only predictors are unbiased). *Any $g(\mathbf{z})$ has $\Pr[g(\mathbf{z}) = X_i] = \frac{1}{2}$.*

A.4 Template Sparsification at Logarithmic Radius (full proof)

We work in the factor-graph view of a random 3-CNF with $M = \alpha m$ clauses (constant $\alpha > 0$), with a fresh sign mask per block. Fix $r = c_3 \log m$ with $c_3 > 0$ small.

Exploration process and tree-likeness. Run a BFS from a uniformly random variable v in the factor graph; each step exposes incident clauses and neighboring variables. Let Z_ℓ be the number of variable nodes at depth ℓ . Standard coupling arguments (Galton-Watson with offspring distribution $\text{Poisson}(\lambda(\alpha))$) show:

Lemma A.13 (Locally tree-like). *There exist $c_3^*(\alpha), \beta(\alpha, c_3) > 0$ such that for $r = c_3 \log m$ and $c_3 < c_3^*$,*

$$\Pr[\mathcal{N}_r(v) \text{ is a tree}] \geq 1 - m^{-\beta}.$$

Moreover, conditional on the unlabeled tree, the literal signs on edges are i.i.d. Rademacher.

Proof. See [7, Ch. 5] for the hypergraph exploration bounds; the expected size of the explored ball is $\lambda^r = \lambda^{c_3 \log m} = m^{c_3 \log \lambda} = m^{o(1)}$. Collisions occur with probability at most $O((\lambda^r)^2/m) = m^{-\beta}$ for small c_3 . Mask signs are independent by construction. \square

Charts and their probability. A chart $\mathcal{C} = (\mathcal{P}, \psi)$ is a finite set of signed rooted radius- r patterns augmented with labels $(a_i, b) \in \{0, 1\}^k \times \{0, 1\}^k$ at the root, with a decision map ψ . For a fixed chart, we bound the probability a random block matches any pattern in its high-bias region.

Lemma A.14 (Augmented pattern probability). *Let P be a fixed signed rooted radius- r tree pattern, with a fixed label pair $(a_i^\circ, b^\circ) \in \{0, 1\}^k \times \{0, 1\}^k$. If A has uniformly random independent rows (so each column a_i is uniform in $\{0, 1\}^k$) and b is uniform in $\{0, 1\}^k$, then*

$$\Pr[\text{nbr}_r(\Phi, i) = P \wedge (a_i, b) = (a_i^\circ, b^\circ)] \leq m^{-\beta'} \cdot 2^{-2k}$$

for some $\beta' = \beta'(\alpha, c_3) > 0$.

Proof. By Lemma A.13, the unlabeled tree P occurs with probability $\leq m^{-\beta'}$; the sign pattern has probability $2^{-|E(P)|}$ which is absorbed in the exponent (or take it into $m^{-\beta'}$). Independence and uniformity of a_i and b contribute 2^{-2k} . \square

Theorem A.15 (Template sparsification for the finite alphabet). *Fix $\varepsilon > 0$ and the finite alphabet \mathcal{U} of local inputs. There exists $\beta'' > 0$ such that*

$$\Pr_{\mathbf{u} \sim \mathcal{D}_m, i \sim [m]} \left[(\Phi, i) \text{ matches some } P \in \text{HB}_\varepsilon(\mathcal{C}_{\mathbf{u}}) \text{ for some } \mathbf{u} \in \mathcal{U} \right] \leq m^{-\beta''}.$$

Consequently, for $t = c_4 m$ i.i.d. blocks, with probability $1 - 2^{-\Omega(m)}$ at most $o(t)$ blocks are high-bias for any \mathbf{u} -measurable rule.

Proof. For each fixed \mathbf{u} , $\text{HB}_\varepsilon(\mathcal{C}_\mathbf{u})$ is a finite set of augmented patterns. By Lemma A.14, each has probability $\leq m^{-\beta'} 2^{-2k}$; the total number of augmented patterns of depth $r = c_3 \log m$ is $m^{O(1)}$ (bounded-degree trees with $O(\lambda^r) = m^{o(1)}$ nodes times $2^{O(k)}$, with $k = O(\log m)$). Thus the per-block probability is $\leq m^{-\beta''}$ for some β'' . Independence across blocks and Chernoff give the $o(t)$ conclusion. \square

Remark A.16 (Uniformity over all u -measurable rules). The sparsification bound is uniform over all u -measurable per-bit rules: the union bound ranges over the finite alphabet U (size $m^{O(1)}$) and the finite set of signed charts at radius $r = c_3 \log m$. No counting over a hypothesis class is required.

Putting it together (local near-randomness). On the test set S supplied by Proposition A.5, for every $j \in S$ and every $i \in [m]$, $(P \circ W_{\text{ERM}})(\Phi)_{j,i} = \hat{h}_i(\mathbf{u}_{j,i})$ is a \mathbf{u} -measurable, $O(\log m)$ -input local rule. By Theorem 5.10, together with sign-invariant neutrality, there exists $\varepsilon(m) \rightarrow 0$ such that $|\Pr[\hat{h}_i(\mathbf{u}_{j,i}) = X_{j,i}] - \frac{1}{2}| \leq \varepsilon(m)$ for all $j \in S, i \in [m]$.

Hence, by the pivot inequality (Lemma 6.1),

$$\Pr[(P \circ W_{\text{ERM}})(\Phi_j) = X_j] \leq \frac{1}{2} + \varepsilon(m) \quad (j \in S).$$

Finally, conditioning on the fixed wrapper, the block-level success indicators $\{\mathbf{1}\{(P \circ W_{\text{ERM}})(\Phi_j) = X_j\}\}_{j \in S}$ are independent (Lemma 6.6), so

$$\Pr[(P \circ W_{\text{ERM}})(\Phi) = X \text{ on all } j \in S] \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{|S|} \leq \left(\frac{1}{2} + \varepsilon(m)\right)^{\gamma t}.$$

By success domination (Proposition A.5(ii)), the same bound (up to $m^{-\Omega(1)}$ slack) applies to $\Pr[P(\Phi) = X]$, yielding the per-program small-success product bound.

A.5 Proof of Calibration Lemma

Assumptions for Lemma 4.8 (Calibration)

(1) Fresh sign mask per block, (2) VV isolation with pairwise-independent columns and uniform b , (3) promise-preserving sign-flip/toggle involution T_i , (4) SILS sign invariance. No further structural assumptions on P are used.

Here we provide the detailed proof of Lemma 4.8 that links symmetrized labels to truth.

Calibration invariance at fixed u

Fix $u = (z, a_i, b)$. The promise-preserving involution $T_i : (F^h, A, b) \mapsto (F^{\tau_i h}, A, b \oplus A e_i)$ (Lemma ??) toggles X_i and preserves u and the conditional measure under D_m . Consequently, conditioning on u we have

$$\Pr[X_i = 1, Y_i = 1 \mid u] = \Pr[X_i = 0, Y_i = 0 \mid u], \quad \Pr[X_i = 1, Y_i = 0 \mid u] = \Pr[X_i = 0, Y_i = 1 \mid u],$$

so $(X_i, Y_i) \mid u$ is exchangeable. Hence the Bayes classifier $h_i^*(u) = \mathbf{1}[f_i(u) \geq 1/2]$ is optimal for both Y_i and X_i at fixed u .

Lemma A.17 (Calibration from symmetrized labels to truth (detailed)). *Fix a bit index i and define $Y_i(\sigma, \Phi) := \text{BM}_\sigma(P(g_\sigma(\Phi)))_i$, where BM_σ is the back-map that xors out $\langle a_i, \sigma \rangle$. Let $f_i(\mathbf{u}) = \mathbb{E}[Y_i(\sigma, \Phi) \mid \mathbf{u}]$ and let $h_i^*(\mathbf{u})$ be the Bayes classifier for f_i . Then*

$$\mathbb{E}_\Phi[\mathbf{1}\{h_i^*(\mathbf{u}(\Phi)) = X_i(\Phi)\}] \geq \mathbb{E}_{\Phi, \sigma}[\mathbf{1}\{Y_i(\sigma, \Phi) = X_i(\Phi)\}] - m^{-\Omega(1)}.$$

Proof. Consider the joint distribution of (\mathbf{u}, X_i, Y_i) where $\mathbf{u} = (\mathbf{z}, a_i, b)$ are the local inputs.

Step 1: Paired involution structure. The key observation is that in our masked+isolated ensemble, there exists an involution that relates different outcomes. Specifically, the map $(F^h, A, b) \mapsto (F^{\tau_i h}, A, b \oplus Ae_i)$ (where τ_i flips signs of variable i) has the following properties:

- It maps instances with witness bit $X_i = 0$ to instances with $X_i = 1$ and vice versa
- It preserves the SILS features \mathbf{z} (which are sign-invariant)
- It preserves a_i but flips b by a_i
- It preserves the uniqueness promise

Step 2: Symmetry of conditional distributions. For a fixed value of $\mathbf{u} = (\mathbf{z}, a_i, b)$, consider the conditional distribution of (X_i, Y_i) given \mathbf{u} . The involution shows that:

$$\Pr[X_i = 1, Y_i = 1 \mid \mathbf{u}] = \Pr[X_i = 0, Y_i = 0 \mid \mathbf{u}]$$

and

$$\Pr[X_i = 1, Y_i = 0 \mid \mathbf{u}] = \Pr[X_i = 0, Y_i = 1 \mid \mathbf{u}].$$

This is because the involution bijectively maps configurations of the first type to configurations of the second type while preserving the measure.

Step 3: Optimal predictor for both Y_i and X_i . Given this symmetry, for any fixed \mathbf{u} :

- $\Pr[Y_i = 1 \mid \mathbf{u}] = f_i(\mathbf{u})$ (by definition)
- $\Pr[X_i = 1 \mid \mathbf{u}] = \Pr[X_i = 1, Y_i = 1 \mid \mathbf{u}] + \Pr[X_i = 1, Y_i = 0 \mid \mathbf{u}]$
- By the symmetry: $\Pr[X_i = 1 \mid \mathbf{u}] = \Pr[X_i = 1, Y_i = 1 \mid \mathbf{u}] + \Pr[X_i = 0, Y_i = 1 \mid \mathbf{u}] = \Pr[Y_i = 1 \mid \mathbf{u}] = f_i(\mathbf{u})$

Therefore, the Bayes optimal predictor $h_i^*(\mathbf{u}) = \mathbf{1}\{f_i(\mathbf{u}) > 1/2\}$ is optimal for predicting both Y_i and X_i given \mathbf{u} .

Step 4: Success bound. The success of h_i^* in predicting X_i is:

$$\Pr[h_i^*(\mathbf{u}) = X_i] = \mathbb{E}_{\mathbf{u}}[\max\{f_i(\mathbf{u}), 1 - f_i(\mathbf{u})\}]$$

which equals its success in predicting Y_i .

Since by Lemma 4.3, $\mathbb{E}_\sigma[\mathbf{1}\{Y_i(\sigma, \Phi) = X_i\}] = \Pr[P(\Phi)_i = X_i]$, and the Bayes optimal predictor achieves at least this average success, we have the claimed bound.

The $m^{-\Omega(1)}$ error term accounts for finite-sample concentration in the ERM approximation. \square

References

- [1] L. G. Valiant and V. V. Vazirani. NP is as easy as detecting unique solutions. *Theoretical Computer Science*, 47(1):85-93, 1986.
- [2] J. L. Carter and M. N. Wegman. Universal classes of hash functions. *Journal of Computer and System Sciences*, 18(2):143-154, 1979.
- [3] M. Naor and A. Naor. Small-bias probability spaces: Efficient constructions and applications. *SIAM Journal on Computing*, 22(4):838-856, 1993.
- [4] J. Håstad. Almost optimal lower bounds for small depth circuits. In *Proceedings of STOC*, 6-20, 1986.
- [5] J. Håstad. On the correlation of parity and small-depth circuits. *SIAM Journal on Computing*, 43(5):1699-1708, 2014.
- [6] A. A. Razborov and S. Rudich (Smolensky is often cited for MOD_p). Lower bounds for the size of circuits of bounded depth with MOD_p gates. *Mathematical Notes*, 41(4):333-338, 1987.
- [7] S. Janson, T. Łuczak, and A. Ruciński. *Random Graphs*. Wiley-Interscience, 2000.
- [8] J. P. Schmidt, A. Siegel, and S. Srinivasan. Chernoff-Hoeffding bounds for applications with limited independence. *SIAM Journal on Discrete Mathematics*, 8(2):223-250, 1995.
- [9] D. Dubhashi and A. Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, 2009.
- [10] M. Li and P. Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*, 3rd ed. Springer, 2008.
- [11] M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.
- [12] Franz, A., Antonenko, O., Soletskyi, R. A theory of incremental compression. *Information Sciences*, 547
- [13] Bennett, M.T. How To Build Conscious Machines. Ph.D. thesis, Australian National University, Canberra, Australia, 2025.
- [14] Goertzel, B. Weakness is All You Need. Unpublished manuscript, 2025
- [15] Goertzel, B. Weakness is All You Need. Keynote at AGI-25 conference, Reykjavik, 2025
- [16] Holman, Craig Elements of an expert system for determining the satisfiability of general Boolean expressions PhD Thesis, Northwestern University, 1990
- [17] Goertzel, Ben Correlational Elegant Normal Form SingularityNET Technical Report, 2025