Average-Case Lower Bounds and Satisfiability Algorithms for Small Threshold Circuits

Ruiwen Chen*1, Rahul Santhanam*2, and Srikanth Srinivasan3

- 1 University of Oxford, Oxford, UK ruiwen.chen@cs.ox.ac.uk
- 2 University of Oxford, Oxford, UK rahul.santhanam@cs.ox.ac.uk
- 3 Indian Institute of Technology Bombay, Mumbai, India srikanth@math.iitb.ac.in

Abstract -

We show average-case lower bounds for explicit Boolean functions against bounded-depth threshold circuits with a superlinear number of wires. We show that for each integer d>1, there is $\varepsilon_d>0$ such that Parity has correlation at most $1/n^{\Omega(1)}$ with depth-d threshold circuits which have at most $n^{1+\varepsilon_d}$ wires, and the Generalized Andreev Function has correlation at most $1/2^{n^{\Omega(1)}}$ with depth-d threshold circuits which have at most $n^{1+\varepsilon_d}$ wires. Previously, only worst-case lower bounds in this setting were known [22].

We use our ideas to make progress on several related questions. We give satisfiability algorithms beating brute force search for depth-d threshold circuits with a superlinear number of wires. These are the first such algorithms for depth greater than 2. We also show that Parity cannot be computed by polynomial-size AC^0 circuits with $n^{o(1)}$ general threshold gates. Previously no lower bound for Parity in this setting could handle more than $\log(n)$ gates. This result also implies subexponential-time learning algorithms for AC^0 with $n^{o(1)}$ threshold gates under the uniform distribution. In addition, we give almost optimal bounds for the number of gates in a depth-d threshold circuit computing Parity on average, and show average-case lower bounds for threshold formulas of any depth.

Our techniques include adaptive random restrictions, anti-concentration and the structural theory of linear threshold functions, and bounded-read Chernoff bounds.

1998 ACM Subject Classification F.1.3 Complexity Measures and Classes

Keywords and phrases threshold circuit, satisfiability algorithm, circuit lower bound

Digital Object Identifier 10.4230/LIPIcs.CCC.2016.1

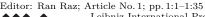
1 Introduction

One of the main goals in complexity theory is to prove circuit lower bounds for explicit functions in P or NP. We seem quite far from being able to prove that there is a problem in NP that requires superlinear Boolean circuits. We have some understanding, via formulations such as the relativization barrier [5], the "natural proofs" barrier [39] and the algebrization barrier [1], of why current techniques are inadequate for this purpose.

^{*} Work done when the author was at University of Edinburgh, Edinburgh, UK, and supported by the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013)/ ERC Grant Agreement no. 615075







However, the community has had more success proving explicit lower bounds against bounded-depth circuits of various kinds. Thanks to pioneering work of Ajtai [2], Furst-Saxe-Sipser [13], Yao [49] and Hastad [19], we know that the Parity and Majority functions require bounded-depth unbounded fan-in circuits of exponential size if only AND and OR gates are allowed. Later Razborov [38] and Smolensky [46] showed that Majority requires exponential size even when MODp gates are allowed in addition to AND and OR gates, for any prime p. The case of bounded-depth circuits with AND, OR and MODm gates, where m is a composite, has been open for nearly thirty years now, even though Majority is conjectured to be hard for such circuits. Williams [47] recently made significant progress by showing that non-deterministic exponential time does not have super-polynomial size circuits with AND, OR and MODm gates, for any m.

For all the bounded-depth circuit classes above, Majority is either known or conjectured to be hard. How about circuit classes which incorporate majority gates, or more generally, gates that are arbitrary linear threshold functions? Note that such gates generalize AND and OR, though not MODp. In the 90s, there was some work on studying the power of bounded-depth threshold circuits. Paturi and Saks [34] showed that depth-2 circuits with majority gates computing Parity require $\tilde{\Omega}(n^2)$ wires; there is also a nearly matching upper bound for Parity. Impagliazzo, Paturi and Saks [22] considered bounded-depth threshold circuits with arbitrary linear threshold gates, and showed that for each depth d, there is a constant $\epsilon_d > 0$ such that Parity requires $n^{1+\epsilon_d}$ wires to compute with depth d threshold circuits.

These lower bounds are worst case lower bounds - they show that for any sequence of small circuits, there exist inputs of every length on which the circuits fail to compute Parity. There are several reasons to be interested in average case lower bounds under the uniform distribution, or equivalently, in correlation upper bounds¹. For one, average-case lower bounds show that a randomly chosen input is likely to be hard, and thus give a way to generate hard instances efficiently. Second, average-case lower bounds are closely tied to pseudo-random generators via the work of Nisan-Wigderson [30], and are indeed a pre-requisite for obtaining pseudo-random generators with non-trivial seed length for a circuit class. Third, recent work on satisfiability algorithms [42, 20, 7] indicates that the design and analysis of non-trivial satisfiability algorithms is closely tied to proving average-case lower bounds, though there is no formal connection. Fourth, the seminal work of Linial-Mansour-Nisan [26] shows that average-case lower bounds for Parity against a circuit class are tied to non-trivially learning the circuit class under the uniform distribution.

With these different motivations in mind, we systematically study average-case lower bounds for bounded-depth threshold circuits. Our first main result shows correlation upper bounds for Parity and another explicit function known as the Generalized Andreev function with respect to threshold circuits with few wires. No correlation upper bounds for explicit functions against bounded-depth threshold circuits with superlinear wires was known before our work.

▶ Theorem 1.1. For each depth $d \ge 1$, there is a constant $\epsilon_d > 0$ such that for all large enough n, no threshold circuit of depth d with at most $n^{1+\epsilon_d}$ wires agrees with Parity on more than $1/2 + 1/n^{\epsilon_d}$ fraction of inputs of length n, and with the Generalized Andreev function on more than $1/2 + 1/2^{n^{\epsilon_d}}$ fraction of inputs of length n.

Theorem 1.1 captures the content of Theorem 4.4 and Theorem 4.7 in Section 4.

¹ By contraposition, if any circuit agreeing with a function f on $1/2 + \varepsilon$ of the inputs has size at least s, then size-s circuits have correlation at most 2ε with f, and vice versa.

We constructivize the ideas of the proof of the strong correlation upper bounds for the Generalized Andreev function to get non-trivial satisfiability algorithms for bounded-depth threshold circuits with few wires. Previously, such algorithms were only known for depth 2 circuits, due to Impagliazzo-Paturi-Schneider [21] and Tamaki (unpublished).

▶ Theorem 1.2. For each depth $d \ge 1$, there is a constant $\epsilon_d > 0$ such that the satisfiability of depth-d threshold circuits with at most $n^{1+\epsilon_d}$ wires can be solved in randomized time $2^{n-n^{\epsilon_d}}\operatorname{poly}(n)$.

Theorem 1.2 is re-stated and proved as Theorem 5.4 in Section 5.

Using our ideas, we also show correlation bounds against AC⁰ circuits with a few threshold gates, as well as learning algorithms under the uniform distribution for such circuits.

▶ **Theorem 1.3.** For each constant d, there is a constant $\gamma > 0$ such that Parity has correlation at most $1/n^{\Omega(1)}$ with AC^0 circuits of depth d and size at most $n^{\log(n)^{0.4}}$ augmented with at most n^{γ} threshold gates. Moreover, the class of AC^0 circuits of size at most $n^{\log(n)^{0.4}}$ augmented with at most n^{γ} threshold gates can be learned to constant error under the uniform distribution in time $2^{n^{1/4+o(1)}}$.

Theorem 1.3 captures the content of Corollary 7.4 and Theorem 7.6 in Section 7.

Having summarized our main results, we now describe related work and our proof techiques in more detail.

1.1 Related work

There has been a large body of work proving upper and lower bounds for constant-depth threshold circuits. Much of this work has focused on the setting of small gate complexity, which seems to be the somewhat easier case to handle. A distinction must also be drawn between work that has focused on the setting where the threshold gates are assumed to be majority gates (i.e. the linear function sign representing the gate has integer coefficients that are bounded by a polynomial in the number of variables) and work that focuses on general threshold gates, since analytic tools such as rational approximation that are available for majority gates do not work in the setting of general threshold gates.

We discuss the work on wire complexity first, followed by the results on gate complexity.

Wire complexity

Paturi and Saks [34] considered depth-2 Majority circuits and showed an $\tilde{\Omega}(n^2)$ lower bound on the wire complexity required to compute Parity; this nearly matches the upper bound of $O(n^2)$. They also showed that there exist majority circuits of size $n^{1+\Theta(\varepsilon_1^d)}$ and depth d computing Parity; here $\varepsilon_1 = 2/(1+\sqrt{5})$. Impagliazzo, Paturi, and Saks [22] showed a depth-d lower bound for general threshold circuits computing Parity: namely, that any such circuit must have wire complexity at least $n^{1+\varepsilon_2^d}$ where $\varepsilon_2 < \varepsilon_1$.

The proof of [22] proceeds by induction on the depth d. The main technical lemma shows that a circuit of depth d can be converted to a depth d-1 circuit of the same size by setting some of the input variables. The variables that are set are set in a random fashion, but not according to the uniform distribution. In fact, this distribution has statistical distance close to 1 from the uniform distribution and furthermore, depends on the circuit whose depth is being reduced. Therefore, it is unclear how to use this technique to prove a correlation bound with respect to the uniform distribution. In contrast, we are able to reduce the depth of the circuit by setting variables uniformly at random (though the variables that we restrict

are sometimes chosen in a way that depends on the circuit), which yields the correlation bounds we want.

Gate complexity

The aforementioned work of Paturi and Saks [34] also proved a near optimal $\tilde{\Omega}(n)$ lower bound on the number of gates in any depth-2 majority circuits computing Parity.

Siu, Roychowdhury, and Kailath [45] considered majority circuits of bounded depth and small gate complexity. They showed that Parity can be computed by depth-d circuits with $O(dn^{1/(d-1)})$ gates. Building on the ideas of [34], they also proved a near matching lower bound of $\tilde{\Omega}(dn^{1/(d-1)})$. Further, they also considered the problem of correlation bounds and showed that there exist depth-d majority circuits with $O(dn^{1/2(d-1)})$ gates that compute Parity almost everywhere and that majority circuits of significantly smaller size have o(1) correlation with Parity (i.e. these circuits cannot compute Parity on more than a 1/2 + o(1) fraction of inputs; recall that 1/2 is trivial since a constant function computes Parity correctly on 1/2 of its inputs). Impagliazzo, Paturi, and Saks [22] extended the worst case lower bound to general threshold gates, where they proved a slightly weaker lower bound of $\Omega(n^{1/2(d-1)})$. As discussed above, though, it is unclear how to use their technique to prove a correlation bound.

Beigel [6] extended the result of Siu et al. to the setting of AC^0 augmented with a few majority gates. He showed that any subexponential-sized depth-d AC^0 circuit with significantly less than some $k = n^{\Theta(1/d)}$ majority gates has correlation o(1) with Parity. The techniques of all the above works with the exception of [22] were based on the fact majority gates can be well-approximated by low-degree rational functions. However, this is not true for general threshold functions [44] and hence, these techniques do not carry over the case of general threshold gates.

A lower bound technique that does carry over to the setting of general threshold gates is that of showing that the circuit class has low-degree polynomial sign-representations. Aspnes, Beigel, Furst and Rudich [3] used this idea to prove that AC^0 circuits augmented with a single general threshold output gate – we refer to these circuits as TAC^0 circuits as in [15] – of subexponential-size and constant-depth have correlation o(1) with Parity. More recently, Podolskii [36] used this technique along with a trick due to Beigel [6] to prove similar bounds for subexponential-sized AC^0 circuits augmented with general threshold gates. However, this trick incurs an exponential blow-up with the number of threshold gates and hence, in the setting of the Parity function, we cannot handle $k > \log n$ threshold gates.

Another technique that has proved useful in handling general threshold gates is Communication Complexity, where the basic idea is to show that the circuit – perhaps after restricting some variables – has low communication complexity in some suitably defined communication model. We can then use results from communication complexity to infer lower bounds or correlation bounds. Nisan [29] used this technique to prove exponential correlation bounds for general threshold circuits (not necessarily even constant-depth) with $n^{1-\Omega(1)}$ threshold gates. Using Beigel's trick and multiparty communication complexity bounds of Babai, Nisan and Szegedy [4], Lovett and Srinivasan [27] (see also [40, 17]) proved exponential correlation bounds for any polynomial-sized AC⁰ circuits augmented with up to $n^{\frac{1}{2}-\Omega(1)}$ threshold gates.

We do not use this technique in our setting for many reasons. Firstly, it cannot be used to prove lower bounds or correlation bounds against functions such as Parity (which has small communication complexity in most models). In particular, these ideas do not yield the noise sensitivity bounds we get here. Even more importantly, it is unclear how to use these techniques to prove any sort of superlinear lower bound on wire complexity,

since there are functions that have threshold circuits with linearly many wires, but large communication complexity even after applying restrictions (take a generic read-once depth-2 Majority formula for example).

Perhaps most closely related to our work is that of Gopalan and Servedio [15] who use analytic techniques to prove correlation bounds for AC^0 circuits augmented with a few threshold gates. Their idea is to use Noise sensitivity bounds (as we do as well) to obtain correlation bounds for Parity with TAC^0 circuits and then extend these results in the same way as in the work of Podolskii [36] mentioned above. As a result, though, the result only yields non-trivial bounds when the number of threshold gates is bounded by $\log n$, whereas our result yields correlation bounds for up to $n^{1/2(d-1)}$ threshold gates.

1.2 Proof techniques

In recent years, there has been an explosion of work on the analytic properties (such as Noise Sensitivity) of linear threshold functions (LTFs) and their generalizations polynomial threshold functions (PTFs) (e.g., [43, 33, 9, 18, 10, 28, 23]). We show here that these techniques can be used in the context of constant-depth threshold circuits as well.

Our first result (Theorem 3.1 in Section 3) is a tight correlation bound for Parity with threshold circuits of depth d and gate complexity much smaller than $n^{1/2(d-1)}$. This generalizes both the results of Siu et al. [45], who proved such a result for majority circuits, and Impagaliazzo, Paturi, and Saks [22], who proved a worst case lower bound of the same order. The proof uses a fundamental theorem of Peres [35] on the noise sensitivity of LTFs; Peres' theorem has also been used by Klivans, O'Donnell, and Servedio [24] to obtain learning algorithms for functions of a few threshold gates. We use Peres' theorem to prove a noise sensitivity upper bound on small threshold circuits of constant depth.

The observation underlying the proof is that the noise sensitivity of a function is exactly the expected variance of the function after applying a suitable random restriction (see also [31]). Seen in this light, Peres' theorem says that, on application of a random restriction, any threshold function becomes quite biased in expectation and hence is well approximated by a constant function. Our analysis of the threshold circuit therefore proceeds by applying a random restriction to the circuit and replacing all the threshold gates at height 1 by the constants that they are well approximated by to obtain a circuit of depth d-1. A straightforward union bound tells us that the new circuit is a good approximation of the original circuit after the restriction. We continue this way with the depth-d-1 circuit until the entire circuit becomes a constant, at which point we can say that after a suitable random restriction, the original circuit is well approximated by a constant, which means its variance is small. Hence, the Noise Sensitivity of the original circuit must be small as well and we are done.

This technique is expanded upon in Section 7, where we use a powerful Noise Sensitivity upper bound for low degree PTFs due to Kane [23] along with standard switching arguments [19] to prove similar results for AC^0 circuits augmented with almost $n^{1/2(d-1)}$ threshold gates. This yields Theorem 1.3.

In Section 4, we consider the problem of extending the above correlation bounds to threshold circuits with small (slightly superlinear) wire complexity. The above proof breaks down even for depth-2 threshold circuits with a superlinear number of wires, since such circuits could have a superlinear number of gates and hence the union bound referred to above is no longer feasible.

In the case of depth-2 threshold circuits, we are nevertheless able to use Peres' theorem, along with ideas of [3] to prove correlation bounds for Parity with circuits with nearly $n^{1.5}$

wires. This result is tight, since by the work of Siu et al. [45], Parity can be well approximated by depth-2 circuits with $O(\sqrt{n})$ gates and hence $O(n^{1.5})$ wires. This argument is in Section B.

Unfortunately, however, this technique needs us to set a large number of variables, which renders it unsuitable for larger depths. The reason for this is that, if we set a large number of variables to reduce the depth from some large constant d to d-1, then we may be in a setting where the number of wires is much larger than the number of surviving variables and hence correlation bounds with Parity may no longer be possible at all.

We therefore use a different strategy to prove correlation bounds for larger constant depths. The lynchpin in the argument is a qualitative refinement of Peres' theorem (Lemma 4.1) that says that on application of a random restriction to an LTF, with good probability, the variance of the LTF becomes negligible (even exponentially small for suitable parameters). The proof of this argument is via anticoncentration results based on the Berry-Esseen theorem and the analysis of general threshold functions via a critical index argument as in many recent works [43, 33, 9, 28].

The above refinement of Peres' theorem allows us to proceed with our argument as in the gates case. We apply a random restriction to the circuit and by the refinement, with good probability (say $1 - n^{-\Omega(1)}$) most gates end up exponentially close to constants. We can then set these "imbalanced" gates to constants and still apply a union bound to ensure that the new circuit is a good approximation to the old one. For the small number of gates that do not become imbalanced in this way, we set *all* variables feeding into them. Since the number of such gates is small, we do not set too many variables. We now have a depth d-1 circuit. Continuing in this way, we get a correlation bound of $n^{-\Omega(1)}$ with Parity. This gives part of Theorem 1.1.

We then strengthen this correlation bound to $\exp(-n^{\Omega(1)})$ for the Generalized Andreev function, which, intuitively speaking, has the following property: even after applying any restriction that leaves a certain number of variables unfixed, the function has exponentially small correlation with any LTF on the surviving variables. To prove lower bounds for larger depth threshold circuits, we follow more or less the same strategy, except that in the above argument, we need most gates to become imbalanced with very high probability $(1 - \exp(-n^{\Omega(1)}))$. To ensure this, we use a bounded read Chernoff bound due to Gavinsky, Lovett, Saks, and Srinivasan [14]. We can use this technique to reduce depth as above as long as the number of threshold gates at height 1 is "reasonably large". If the number of gates at height 1 is very small, then we simply guess the values of these few threshold gates and move them to the top of the circuit and proceed. This gives the other part of Theorem 1.1.

This latter depth-reduction lemma can be completely constructivized to design a satisfiability algorithm that runs in time $2^{n-n^{\Omega(1)}}$. The algorithm proceeds in the same way as the above argument, iteratively reducing the depth of the circuit. A subtlety arises when we replace imbalanced gates by constants, since we are changing the behaviour of the circuit on some (though very few) inputs. Thus, a circuit which was satisfiable only at one among these inputs might now end up unsatisfiable. However, we show that there is an efficient algorithm that enumerates these inputs and can hence check if there are satisfiable assignments to the circuits from among these inputs. This gives Theorem 1.2.

In Section 6, we prove correlation bounds for the Generalized Andreev function with threshold *formulas* of any arity and any depth. The proof is based on a retooling of the argument of Nečiporuk for formulas of constant arity over any basis and yields a correlation bound as long as the wire complexity is at most $n^{1.5-\Omega(1)}$.

2 Preliminaries

2.1 Basic Boolean function definitions

A Boolean function on n variables will be a function $f: \{-1,1\}^n \to \{-1,1\}$. We use the standard inner product on functions $f,g: \{-1,1\}^n \to \mathbb{R}$ defined by $\langle f,g \rangle = \mathbf{E}_{x \sim \{-1,1\}^n} [f(x)g(x)]^2$.

Given Boolean functions f, g on n variables, the *Correlation* between f and g – denoted Corr(f, g) – is defined as

$$Corr(f,g) := |\langle f,g \rangle| = \left| \underset{x \sim \{-1,1\}^n}{\mathbf{E}} [f(x)g(x)] \right| = |2 \Pr_x[f(x) = g(x)] - 1|.$$

Also, we use $\delta(f,g)$ to denote the fractional distance between f and g: i.e., $\delta(f,g) = \Pr_x[f(x) \neq g(x)]$. Then, we have $\operatorname{Corr}(f,g) = |1 - 2\delta(f,g)|$. We say that f is δ -approximated by g if $\delta(f,g) \leq \delta$.

We use Par_n to denote the parity function on n variables. I.e. $\operatorname{Par}_n(x_1,\ldots,x_n)=\prod_{i=1}^n x_i$.

- ▶ **Definition 2.1** (Restrictions). A restriction on n variables is a function $\rho:[n] \to \{-1,1,*\}$. A random restriction is a distribution over restrictions. We use \mathcal{R}_p^n to denote the distribution over restrictions on n variables obtained by setting each $\rho(x) = *$ with probability p and to 1 and p with probability p and to 1. We will often view the process of sampling a restriction as picking a pair (I,y) where $I \subseteq [n]$ is obtained by picking each element of [n] to be in I with probability p and $p \in \{-1,1\}^{n-|I|}$ uniformly at random.
- ▶ Definition 2.2 (Restriction trees and Decision trees). A restriction tree T on $\{-1,1\}^n$ of depth h is a binary tree of depth h all of whose internal nodes are labelled by one of n variables, and the outgoing edges from an internal node are labelled +1 and -1; we assume that a node and its ancestor never query the same variable. Each leaf ℓ of T defines a restriction ρ_{ℓ} that sets all the variables on the path from the root of the decision tree to ℓ and leaves the remaining variables unset. A random restriction tree T of depth h is a distribution over restriction trees of depth h.

Given a restriction tree T, the process of choosing a random edge out of each internal node generates a distribution over the leaves of the tree (note that this distribution is not uniform: the weight it puts on leaf ℓ at depth d is 2^{-d}). We use the notation $\ell \sim T$ to denote a leaf ℓ of T picked according this distribution.

A decision tree is a restriction tree all of whose leaves are labelled either by +1 or -1. We say a decision tree has size s if the tree has s leaves. We say a decision tree computes a function $f: \{-1,1\}^n \to \{-1,1\}$ if for each leaf ℓ of the tree, $f|_{\rho_{\ell}}$ is equal to the label of ℓ .

- ▶ Fact 2.3 (Facts about correlation). Let $f, g, h : \{-1, 1\}^n \to \{-1, 1\}$ be arbitrary.
- 1. $Corr(f, g) \in [0, 1]$.
- **2.** If $Corr(f,g) \le \varepsilon$ and $\delta(g,h) \le \delta$, then $Corr(f,h) \le \varepsilon + 2\delta$.
- 3. Let g_1, \ldots, g_N be Boolean functions such that no two of them are simultaneously true and let h denote their OR. Then, $\operatorname{Corr}(f,h) \leq \sum_{i=1}^N \max\{\operatorname{Corr}(f,1),\operatorname{Corr}(f,g_i)\}$, where 1 denotes the constant 1 function.
- **4.** Let \mathcal{T} be any random restriction tree. Then $\operatorname{Corr}(f,g) \leq \mathbf{E}_{T \sim \mathcal{T}, \ell \sim T}[\operatorname{Corr}(f|_{\rho_{\ell}}, g|_{\rho_{\ell}})]$.

² $x \sim \{-1, 1\}^n$ stands for that x is uniform in $\{-1, 1\}^n$.

- ▶ **Definition 2.4** (Noise sensitivity and Variance [32]). Given a Boolean function $f: \{-1,1\}^n$ $\to \{-1,1\}$ and a parameter $p \in [0,1]$, we define the *Noise sensitivity of* f *with noise parameter* p denoted $NS_p(f)$ as follows. Pick $x \in \{-1,1\}^n$ uniformly at random and $y \in \{-1,1\}^n$ by negating (i.e. flipping) each bit of x independently with probability p; we define $NS_p(f) = Pr_{(x,y)}[f(x) \neq f(y)]$. The variance of f denoted Var(f) is defined to be $2NS_{1/2}(f)$.
- ▶ **Proposition 2.5.** Let $f: \{-1,1\}^n \to \{-1,1\}$ be any Boolean function. Then,
- 1. For $p \le 1/2$, $NS_p(f) = \frac{1}{2} \mathbf{E}_{\rho \sim \mathcal{R}_{2p}^n} [Var(f|_{\rho})]$.
- 2. If $p \geq \frac{1}{n}$, then $Corr(f, Par_n) \leq O(NS_p(f))$.

The above fact is folklore, but we couldn't find explicit proofs in the literature. Therefore we present them in the appendix (see Appendix A).

▶ Fact 2.6. Let $f: \{-1,1\}^n \to \{-1,1\}$ be any Boolean function. Let $p = \min\{\Pr_x[f(x) = 1], \Pr_x[f(x) = -1]\}$ where x is chosen uniformly from $\{-1,1\}^n$. Then, $\operatorname{Var}(f) = \Theta(p)$.

2.2 Threshold functions and circuits

- ▶ Definition 2.7 (Threshold functions and gates). A Threshold gate is a gate ϕ labelled with a pair (w, θ) where $w \in \mathbb{R}^m$ for some $m \in \mathbb{N}$ and $\theta \in \mathbb{R}$. The gate computes the Boolean function $f_{\phi}: \{-1, 1\}^m \to \{-1, 1\}$ defined by $f_{\phi}(x) = \operatorname{sgn}(\langle w, x \rangle \theta)$ (we define $\operatorname{sgn}(0) = -1$ for the sake of this definition). The fan-in of the gate ϕ denoted fan-in (ϕ) is m. A Linear Threshold function (LTF) is a Boolean function that can be represented by a Threshold gate.
- ▶ Definition 2.8 (Threshold circuits). A Threshold circuit C is a Boolean circuit whose gates are all threshold gates. There are designated output gates, which compute the functions computed by the circuit. Unless explicitly mentioned, however, we assume that our threshold circuits have a unique output gate. The *gate complexity* of C is the number of (non-input) gates in the circuit, while the *wire complexity* is the sum of all the fan-ins of the various gates.

A Threshold map from n to m variables is a depth-1 threshold circuit C with n inputs and m outputs. We say that such a map is read-k if each input variable is an input to at most k of the threshold gates in C.

The proof of the following can be found for example in [41].

- ▶ Lemma 2.9 ([41]). The number of distinct linear threshold functions on n bits is at most $2^{O(n^2)}$.
- ▶ Definition 2.10 (Restrictions of threshold gates and circuits). Given a threshold gate ϕ of fan-in m labelled by the pair (w, θ) and a restriction ρ on m variables, we use ϕ_{ρ} to denote the threshold gate over the variables indexed by $\rho^{-1}(*)$ obtained in the natural way by setting variables according to ρ .

We will also need Peres' theorem, which bounds the Noise Sensitivity of threshold functions.

▶ Theorem 2.11 (Peres' theorem[35, 32]). Let $f : \{-1,1\}^n \to \{-1,1\}$ be any LTF. Then, $\underset{\rho \sim \mathcal{R}_2^n}{\mathbf{E}} [\operatorname{Var}(f|_{\rho})] = \operatorname{NS}_{\frac{p}{2}}(f) = O(\sqrt{p}).$

Using the above for p = 1/n and Proposition 2.5, we obtain

▶ Corollary 2.12. Let $f: \{-1,1\}^n \to \{-1,1\}$ be any threshold function. Then $Corr(f, Par_n) < O(n^{-1/2})$.

2.3 Description lengths and Kolmogorov Complexity

- ▶ **Definition 2.13** (Kolmogorov Complexity). The Kolmogorov complexity of an n-bit Boolean string x is the length of the shortest bit string of the form (M, w) where M is the description of a Turing Machine and w an input to M such that M(w) = x. We use K(x) to denote the Kolmogorov complexity of x.
- ▶ Fact 2.14. For any $\alpha \in (0,1)$, the fraction of n-bit strings x satisfying $K(x) \leq (1-\alpha)n$ is at most $2^{-\alpha n+1}$.
- ▶ Definition 2.15 (Descriptions of circuits). We can also talk about the description lengths of threshold circuits, which we define as follows. By Lemma 2.9, we know that the number of LTFs on n bits is $2^{O(n^2)}$, and hence we can fix some $O(n^2)$ -bit description for each such function. The description of a threshold circuit C is a description of the underlying graph theoretic structure of C followed by the descriptions of the threshold functions computed by each of its gates and the input variables labelling its input gates. We use $\sigma(C)$ to denote the length of this description of C.
- ▶ Proposition 2.16. For any threshold circuit C with wire complexity at most s on at most n variables, $\sigma(C) = O(s^2 + s \log n)$. If $s \ge n$, then the description length is at most $O(s^2)$.

Proof. Since the wire complexity is at most s, the graph underlying the circuit can be described using $O(s \log s)$ bits (for example, for each wire, we can describe the gates that it connects). Let ϕ_1, \ldots, ϕ_m be the threshold gates in the circuit. We can write down a description of the LTFs f_1, \ldots, f_m using $\sum_i O(k_i^2)$ bits where k_i is the fan-in of ϕ_i ; this is at most $O(\sum_i k_i)^2 = O(s^2)$. Finally, to describe the input variable assignments to the input gates, we need $O(s \log n)$ bits.

2.4 The Generalized Andreev function

We state here the definition of a generalization of Andreev's function, due to Komargodski and Raz, and Chen, Kabanets, Kolokolova, Shaltiel, and Zuckerman [25, 7]. This function will be used to give strong correlation bounds for constant-depth threshold circuits with slightly superlinear wire complexity.

We first need some definitions.

▶ **Definition 2.17** (Bit-fixing extractor). A function $E: \{-1,1\}^n \to \{-1,1\}^m$ is a (n,k,m,ζ) bit-fixing extractor if for every random variable X that is uniform on a subcube³ of $\{-1,1\}^n$ of dimension at least k, the function E(X) is ζ -close to uniform on $\{-1,1\}^m$.

We have the following explicit construction of a bit-fixing extractor.

▶ Theorem 2.18 ([37]). There is an absolute constant $c \ge 1$ so that the following holds. There is a polynomial-time computable function $E: \{-1,1\}^n \to \{-1,1\}^m$ that is an (n,k,m,ζ) -bit fixing extractor for any $k \ge (\log n)^c$, m = 0.9k, and $\zeta \le 2^{-k^{\Omega(1)}}$.

Also recall that a function Enc: $\{-1,1\}^a \to \{-1,1\}^b$ defines (α, L) -error-correcting code for parameters $\alpha \in [0,1]$ and $L \in \mathbb{N}$ if for any $z \in \{-1,1\}^b$, the number of elements in the image of Enc that are at relative Hamming distance at most α from z is bounded by L.

The following theorem is a folklore result, and stated explicitly in the work of Chen et al. [7].

A subcube of dimension k is a subset of $\{-1,1\}^n$ containing elements which are consistent with some restriction with k *'s.

▶ Theorem 2.19 ([7], Theorem 6.4). Let $r = n^{\beta}$ for any fixed $0 < \beta < 1$. There exists an (α, L) -error correcting code with Enc: $\{-1, 1\}^{4n} \to \{-1, 1\}^{2^r}$ where $\alpha = \frac{1}{2} - O(2^{-r/4})$ and $L = O(2^{r/2})$. Further, there is a poly(n) time algorithm, which when given as input $x \in \{-1, 1\}^n$ and $i \in [2^r]$ in binary, outputs $\operatorname{Enc}(x)_i$, the ith bit of $\operatorname{Enc}(x)$.

Now we can define the generalized Andreev function as in [7]. The function is $F: \{-1,1\}^{4n} \times \{-1,1\}^n \to \{-1,1\}$ and is defined as follows. Let $\gamma > 0$ be a constant parameter. The parameter will be fixed later according to the application at hand.

Let E be any $(n, n^{\gamma}, m = 0.9n^{\gamma}, 2^{-n^{\Omega(\gamma)}})$ extractor (we can obtain an explicit one using Theorem 2.18). We interpret the output of E as an integer from $[2^m]$ in the natural way. Let Enc: $\{-1, 1\}^{4n} \to \{-1, 1\}^{2^m}$ define a $(\frac{1}{2} - O(2^{-m/4}), 2^{m/2})$ -list decodable code as in Theorem 2.19. Then, we define $F(x_1, x_2)$ by

$$F(x_1, x_2) = \text{Enc}(x_1)_{E(x_2)}. (1)$$

Given $a \in \{-1, 1\}^{4n}$, we use $F_a(\cdot)$ to denote the resulting sub-function on n bits obtained by fixing $x_1 = a$.

The following lemma was proved as part of Theorem 6.5 in [7].

▶ Lemma 2.20 ([7], Theorem 6.5). Let C be any circuit on n^{γ} variables with binary description length $\sigma(C) \leq n$ according to some fixed encoding scheme. Let ρ be any restriction of n variables leaving n^{γ} variables unfixed. Let $f(y) := F_a|_{\rho}(y)$ for $a \in \{-1,1\}^{4n}$ satisfying $K(a) \geq 3n$. Then

$$\operatorname{Corr}(f, C) \le \exp(-n^{\Omega(\gamma)}).$$

2.5 Concentration bounds

We state a collection of concentration bounds that we will need in our proofs. The proofs of Theorems 2.21 and 2.23 may be found in the excellent book by Dubhashi and Panconesi [11].

▶ Theorem 2.21 (Chernoff bound). Let $w \in \mathbb{R}^n$ be arbitrary and x is chosen uniformly from $\{-1,1\}^n$. Then

$$\Pr_{x}[|\langle w, x \rangle| \ge t \cdot ||w||_2] \le \exp(-\Omega(t^2)).$$

▶ **Definition 2.22 (Imbalance).** We say that a threshold gate ϕ labelled by (w, θ) is t-imbalanced if $|\theta| \ge t \cdot ||w||_2$ and t-balanced otherwise.

We also need a multiplicative form of the Chernoff bound for sums of Boolean random variables.

▶ Theorem 2.23 (Multiplicative Chernoff bound). Let Y_1, \ldots, Y_m be independent Boolean random variables such that $\mathbf{E}[Y_i] = p_i$ for each $i \in [m]$. Let p denote the average of the p_i . Then, for any $\varepsilon > 0$

$$\Pr[|\sum_{i} Y_i - pm| \ge \varepsilon pm] \le \exp(-\Omega(\varepsilon^2 pm)).$$

Let Y_1, \ldots, Y_m be random variables defined as functions of independent random variables X_1, \ldots, X_n . For $i \in [m]$, let $S_i \subseteq [n]$ index those random variables among X_1, \ldots, X_n that influence Y_i . We say that Y_1, \ldots, Y_m are read-k random variables if any $j \in [n]$ belongs to S_i for at most k different $i \in [m]$.

The notation D(p||q) represents the KL-divergence (see, e.g., [8]) between the two probability distributions on $\{0,1\}$ where the probabilities assigned to 1 are p and q respectively.

▶ Theorem 2.24 (A read-k Chernoff bound [14]). Let Y_1, \ldots, Y_m be $\{0, 1\}$ -valued read-k random variables such that $\mathbf{E}[Y_i] = p_i$. Let p denote the average of p_1, \ldots, p_m . Then, for any $\varepsilon > 0$,

$$\Pr[\sum_{i} Y_{i} \ge pm(1+\varepsilon)] \le \exp(-D(p(1+\varepsilon)||p)m/k).$$

Using standard estimates on the KL-divergence, we get

▶ Corollary 2.25. Let $Y_1, ..., Y_m$ be as in the statement of Theorem 2.24 and assume $\mathbf{E}[\sum_i Y_i] \leq \mu$. Then,

$$\Pr[\sum_{i} Y_i \ge 2\mu] \le \exp(-\Omega(\mu/k)).$$

3 Correlation bounds for threshold circuits with small gate complexity

In this section, we show that constant-depth threshold circuits with a small number of gates cannot correlate well with the Parity function.

It should be noted that Nisan [29] already proved strong correlation bounds for the Inner Product function against any threshold circuit (not necessarily constant-depth) with a sub-linear (much smaller than $n/\log n$) number of threshold gates. The idea of the proof is to first show that each threshold gate on n variables has a δ -error randomized protocol with complexity $O(\log(n/\delta))$ [29, Theorem 1]. One can use this to show that any threshold circuit as in the theorem can be written as a decision tree of depth n/k querying threshold functions and hence has a $\exp(-\Omega(k))$ -error protocol of complexity at most n/10. Standard results in communication complexity imply that any such function can have correlation at most $\exp(-\Omega(k))$ with inner product.

However, such techniques cannot be used to obtain lower bounds or correlation bounds for the parity function, since the parity function has low communication complexity (even in the deterministic setting). An even bigger disadvantage to this technique is that it cannot be used to obtain *any* superlinear lower bound on the wire complexity, since threshold circuits with a linear number of wires can easily compute functions with high communication complexity (such as the Disjointness function).

The techniques we use here can be used to give correlation bounds for the parity function; further, these correlation bounds are nearly tight (Theorem 3.4). In fact, we prove something stronger: we upper bound the noise sensitivity of small constant-depth threshold circuits, which additionally implies the existence of non-trivial learning algorithms [24, 15]. Further, our techniques also imply noise sensitivity bounds for AC⁰ circuits augmented with a small number of threshold gates.

In this section, we illustrate our technique with the case of threshold circuits with a small number of gates. The generalizations to AC^0 circuits augmented with a small number of threshold gates are obtained in Section 7.

3.1 Correlation bounds via noise sensitivity

▶ **Theorem 3.1.** Let C be a depth d threshold circuit with at most k threshold gates. Then, for any parameters $p, q \in [0, 1]$, we have

$$NS_{p^{d-1}q}(C) \le O(k\sqrt{p} + \sqrt{q}).$$

Proof. We assume that $q \leq \frac{1}{2}$, since otherwise the statement of the theorem is trivial. We will instead prove that for $p_d := 2p^{d-1}q \in [0,1]$ and $\rho_d \sim \mathcal{R}_{p_d}^n$ (n is the number of input variables to C), we have

$$\mathbf{E}_{\rho_d}[\operatorname{Var}(C|_{\rho_d})] \le O(k\sqrt{p} + \sqrt{q}). \tag{2}$$

This will imply the theorem, since by Proposition 2.5, we have $NS_{p^{d-1}q}(C) = \frac{1}{2} \mathbf{E}_{\rho_d}[Var(C|_{\rho_d})].$

The proof of (2) is by induction on the depth d of the circuit. The base case d = 1 is just Peres' theorem (Theorem 2.11).

Now assume that C has depth d > 1. Let k_1 be the number of threshold circuits at height 1 in the circuit. We choose a random restriction $\rho \sim \mathcal{R}_p^n$ and consider the circuit $C|_{\rho}$. It is easy to check that

$$\mathbf{E}_{\rho_d}[\operatorname{Var}(C|_{\rho_d})] = \mathbf{E}_{\rho}[\mathbf{E}_{\rho_{d-1}}[\operatorname{Var}((C|_{\rho})|_{\rho_{d-1}})]],\tag{3}$$

and hence to prove (2), it suffices to bound the expectation of $Var((C|_{\rho})|_{\rho_{d-1}})$.

Let us first consider the circuit $C|_{\rho}$. Peres' theorem tells us that on application of the restriction ρ , each threshold gate at height 1 becomes quite biased on average. Formally, by Theorem 2.11 and Fact 2.6, for each threshold gate ϕ at height 1, there is a bit $b_{\phi,\rho} \in \{-1,1\}$ such that

$$\mathbf{E}[\Pr_{\rho \mid x \in \{-1,1\}^{|\rho^{-1}(*)|}} [\phi_{\rho}(x) \neq b_{\phi,\rho}]] \leq O(\sqrt{p}).$$

In particular, replacing ϕ_{ρ} by $b_{\phi,\rho}$ in the circuit $C|_{\rho}$ yields a circuit that differs from $C|_{\rho}$ on only an $O(\sqrt{p})$ fraction of inputs (in expectation). Applying this replacement to each of the k_1 threshold gates at height 1 yields a circuit C'_{ρ} with $k-k_1$ threshold gates and depth d-1 such that

$$\mathbf{E}_{\rho}[\delta(C|_{\rho}, C'_{\rho})] \le O(k_1 \sqrt{p}) \tag{4}$$

where $\delta(C|_{\rho}, C'_{\rho})$ denotes the fraction of inputs on which the two circuits differ. On the other hand, we can apply the inductive hypothesis to C'_{ρ} to obtain

$$\mathbf{E}_{\rho_{d-1}}[\text{Var}((C'_{\rho})|_{\rho_{d-1}})] \le O((k-k_1)\sqrt{p} + \sqrt{q}). \tag{5}$$

Therefore, to infer (2), we put the above together with (4) and the following elementary fact.

▶ Proposition 3.2. Say $f, g : \{-1, 1\}^m \to \{-1, 1\}$ and $\delta = \delta(f, g)$. Then, for any $r \in [0, 1]$, we have $\mathbf{E}_{\rho \sim \mathcal{R}_r^n}[\operatorname{Var}(f|_{\rho})] \leq \mathbf{E}_{\rho \sim \mathcal{R}_r^n}[\operatorname{Var}(g|_{\rho})] + 4\delta$.

Proof of Proposition 3.2. By Proposition 2.5, we know that $\mathbf{E}_{\rho \sim \mathcal{R}_r^n}[\operatorname{Var}(f|_{\rho})] = 2\operatorname{NS}_{r/2}(f)$, and similarly for g. By definition of noise sensitivity, we have $\operatorname{NS}_{r/2}(f) = \operatorname{Pr}_{(x,y)}[f(x) \neq f(y)]$ where $x \in \{-1,1\}^m$ is chosen uniformly at random and g is chosen by flipping each bit of g with probability g. Note that each of g and g is individually uniformly distributed over $\{-1,1\}^m$ and hence, both g and g and g and g both with probability at least g. This yields

$$NS_{r/2}(f) = \Pr_{(x,y)}[f(x) \neq f(y)] \le \Pr_{(x,y)}[g(x) \neq g(y)] + 2\delta = NS_{r/2}(g) + 2\delta,$$

which implies the claimed bound.

▶ Corollary 3.3. Let $d \geq 2$ and $\delta \in [0,1]$ be arbitrary parameters. Assume that C is a depth d threshold circuit over n variables with at most $\delta n^{\frac{1}{2(d-1)}}$ threshold gates. Then, $\operatorname{Corr}(C, \operatorname{Par}_n) \leq O(\delta^{(1-\frac{1}{d})})$.

Proof. Let $k \leq \delta n^{1/2(d-1)}$ be the number of gates in the threshold circuit C. We apply Theorem 3.1 with the following optimized parameters: $p = \frac{1}{n^{1/d}} \cdot \frac{1}{k^{2/d}}$ and $q \in [0, 1]$ such that $p^{d-1}q = \frac{1}{n}$. It may be verified that for this setting of parameters, Theorem 3.1 gives us

$$NS_{1/n}(C) \le O\left(\frac{k^{1-1/d}}{n^{1/(2d)}}\right) \le O(\delta^{1-\frac{1}{d}}).$$

As noted in Proposition 2.5, we have $\operatorname{Corr}(C,\operatorname{Par}_n) \leq O(\operatorname{NS}_{1/n}(C))$. This completes the proof.

▶ Remark. It is instructive to compare the above technique with the closely related work of Gopalan and Servedio [15]. The techniques of [15] applied to the setting of Theorem 3.1 show that $NS_p(C) \leq O(k2^k\sqrt{p})$, which gives a better dependence on the noise parameter p, but a much worse dependence on k. Indeed, this is not surprising since in this setting, the technique of Gopalan and Servedio does not use the fact that the circuit is of depth d. The threshold circuit is converted to a decision tree of depth k querying threshold functions and it is this tree that is analyzed.

We believe that the right answer should incorporate the best of both bounds: $NS_p(f) \le O_d(k^{d-1} \cdot \sqrt{p})$. As in Corollary 3.3, this would show that $Corr(C, Par_n) = o(1)$ if $k = o(n^{1/2(d-1)})$, but additionally, we would also get $Corr(C, Par_n) \le n^{-\frac{1}{2} + o(1)}$ as long as $k = n^{o(1)}$, which we are not able to prove currently.

It is known from the work of Siu, Roychowdhury and Kailath [45, Theorem 7] that Corollary 3.3 is tight in the sense that there do exist circuits of gate complexity roughly $n^{1/2(d-1)}$ that have significant correlation with Par_n. More formally,

▶ Theorem 3.4 (Theorem 7 in [45]). Let $\varepsilon > 0$ be an arbitrary constant. Then, there is a threshold circuit of depth d with $O(d) \cdot (n \log(1/\varepsilon))^{1/2(d-1)}$ gates that computes Par_n correctly on a $1 - \varepsilon$ fraction of inputs.

4 Correlation bounds for threshold circuits with small wire complexity

The following is a key lemma that will be used in the proofs of our correlation bounds. We state the lemma here and prove our correlation bounds. The lemma will be proved in Section 4.2.

Recall that a threshold gate ϕ with label (w, θ) is t-balanced if $|\theta| \leq t \cdot ||w||_2$.

▶ **Lemma 4.1** (Main Structural lemma for threshold gates). For any threshold gate ϕ over n variables with label (w, θ) and any $p \in [0, 1]$, we have

$$\Pr_{\rho \sim \mathcal{R}_p^n} [\phi_\rho \ is \ \frac{1}{p^{\Omega(1)}} \text{-}balanced] \leq p^{\Omega(1)}.$$

The proof of the correlation bounds proceed by iteratively reducing the depth of the circuit. In order to perform this depth-reduction for a depth d circuit, we need to analyze the threshold map defined by the threshold gates at depth d-1. The first observation, which follows from Markov's inequality, shows that we may assume (after setting a few variables) that the map reads each variable only a few times.

▶ Fact 4.2 (Small wire-complexity to small number of reads). Let C be any threshold circuit on n variables with wire complexity at most cn. Then, there is a set S of at most n/2 variables such that each variable outside S is an input variable to at most 2c many gates in C.

The second observation is that if the fan-ins of all the threshold gates are small, then depth-reduction is easy (after setting some more variables).

▶ Proposition 4.3 (Handling small fan-in gates). Let $C = (\phi_1, \ldots, \phi_m)$ be any read-k threshold map on n variables such that $\max_i \text{fan-in}(\phi_i) \leq t$. Then, there is a set S of n/kt variables such that each ϕ_i depends on at most one variable in S.

Proof. This may be done via a simple graph theoretic argument. Define an undirected graph whose vertex set is the set of n variables and two variables are adjacent iff they feed into the same threshold gate. We need to pick an S that is an independent set in this graph. Since the graph has degree at most kt, we can greedily find an independent set of size at least n/kt. Let S be such an independent set.

4.1 Proofs of correlation bounds

Let B > 2 be a constant real parameter that we will choose to satisfy various constraints in the proofs below. For $d \ge 1$, define $\varepsilon_d = B^{-(2d-1)}$ and $\delta_d = B\varepsilon_d$.

▶ Theorem 4.4 (Correlation bounds for parity). For any $d \ge 1$ and $c \le n^{\varepsilon_d}$, any depth-d threshold circuit C with at most cn wires satisfies $Corr(C, Par_n) \le O(n^{-\varepsilon_d})$ where the $O(\cdot)$ hides absolute constants (independent of d and n).

Proof. The proof is by induction on the depth d of C. The base case is d=1, which is the case when C is only a single threshold gate. In this case, Corollary 2.12 tells us that $\operatorname{Corr}(C,\operatorname{Par}_n) \leq O(n^{-1/2}) \leq n^{-\varepsilon_1}$, since B>2.

Now, we handle the inductive case when the depth d > 1. Our analysis proceeds in phases.

Phase 1

We first transform the circuit into a read-2c circuit by setting n/2 variables. This may be done by Fact 4.2. This defines a restriction tree of depth n/2. By Fact 2.3, it suffices to show that each leaf of this restriction tree, the correlation of the restricted circuit and $\operatorname{Par}_{n/2}$ remains bounded by $O(n^{-\varepsilon_d})$.

Let n_1 now denote the new number of variables and let C_1 now be the restricted circuit at some arbitrary leaf of the restriction tree. By renaming the variables, we assume that they are indexed by the set $[n_1]$.

Phase 2

Let ϕ_1, \ldots, ϕ_m be the threshold gates at depth d-1 in the circuit C_1 . We call ϕ_i large if fan-in $(\phi_i) > n^{\delta_d}$ and small otherwise. Let $L \subseteq [m]$ be defined by $L = \{i \in [m] \mid \phi_i \text{ large}\}$. Assume that $|L| = \ell$. Note that $\ell \cdot n^{\delta_d} \le n^{1+\varepsilon_d}$ and hence $\ell \le n^{1+\varepsilon_d-\delta_d} \le n$.

We restrict the circuit with a random restriction $\rho = (I, y) \sim \mathcal{R}_p^{n_1}$, where $p = n^{-\delta_d/2}$. By Lemma 4.1, we know that for each $i \in [m]$ and some $t = \frac{1}{n^{\Omega(1)}}$ and $q = p^{\Omega(1)}$,

$$\Pr_{\rho}[\phi_i|_{\rho} \ t\text{-balanced}] \le q. \tag{6}$$

Further, we also know that for each $i \in L$, the expected value of fan-in $(\phi_i|_{\rho}) = p$ -fan-in (ϕ_i) , since each variable is set to a constant with probability 1 - p. Since $i \in L$, the expected fan-in of each ϕ_i $(i \in L)$ is at least $n^{\delta_d/2}$. Hence, by a Chernoff bound (Theorem 2.23), we see that for any $i \in L$,

$$\Pr_{\rho}[\text{fan-in}(\phi_i|_{\rho}) > 2p \cdot \text{fan-in}(\phi_i)] \le \exp(-\Omega(n^{\delta_d/2})). \tag{7}$$

Finally another Chernoff bound (Theorem 2.23) tells us that

$$\Pr_{\rho=(I,y)}[|I| < \frac{n_1 p}{2}] \le \exp(-\Omega(n_1 p)) = \exp(-\Omega(n p)).$$
(8)

We call a set I generic if $|I| \ge \frac{n_1 p}{2}$ and fan-in $(\phi_i|_{\rho}) \le 2p \cdot \text{fan-in}(\phi_i)$ for each $i \in L$. Let \mathcal{G} denote the event that I is generic. By (7) and (8), we know that $\Pr_I[\neg \mathcal{G}] \le \ell \exp(-\Omega(n^{\delta_d/2})) + \exp(-\Omega(np)) \le \exp(-n^{\delta_d/4})$. In particular, conditioning on \mathcal{G} doesn't change (6) by much.

$$\Pr_{\rho=(I,y)}[\phi_i|_{\rho} \text{ t-balanced } |\mathcal{G}] \le q + \exp(-n^{\delta_d/4}) \le 2q.$$
(9)

Our aim is to further restrict the circuit by setting all the input variables to the gates ϕ_i that are t-balanced. In order to analyze this procedure, we define random variables Y_i $(i \in L)$ so that $Y_i = 0$ if $\phi_i|_{\rho}$ is t-imbalanced and fan-in $(\phi_i|_{\rho})$ otherwise. Let $Y = \sum_{i \in L} Y_i$. Note that

$$\mathop{\mathbf{E}}_{\rho}[Y_i \mid \mathcal{G}] \leq (2p \cdot \text{fan-in}(\phi_i)) \cdot \Pr_{\rho}[\phi_i|_{\rho} \text{ t-balanced } \mid \mathcal{G}] \leq 4pq \cdot \text{fan-in}(\phi_i)$$

where the first inequality follows from the fact that since we have conditioned on I being generic, we have $\operatorname{fan-in}(\phi_i|_{\rho}) \leq 2p \cdot \operatorname{fan-in}(\phi_i)$ with probability 1. Hence, we have

$$\mathbf{E}[Y \mid \mathcal{G}] \le 4pq \cdot \sum_{i} \text{fan-in}(\phi_i) \le 4pq \cdot n^{1+\varepsilon_d}. \tag{10}$$

We let $\mu := 4pq \cdot n^{1+\varepsilon_d}$. By Markov's inequality,

$$\Pr_{\rho}[Y \ge \frac{\mu}{\sqrt{q}} \mid \mathcal{G}] \le \sqrt{q}. \tag{11}$$

In particular, we can condition on a fixed generic $I \subseteq [n]$ such that for random $y \sim \{-1,1\}^{n_1-|I|}$, we have

$$\Pr_{y}[Y \ge \frac{\mu}{\sqrt{q}}] \le \sqrt{q}.$$

The above gives us a restriction tree T (that simply sets all the variables in $[n_1] \setminus I$) such that at all but $1 - 2\sqrt{q}$ fraction of leaves λ of T, the total fan-in of the large gates at depth 1 in C_1 that are t-balanced is at most $\frac{\mu}{\sqrt{q}}$; call such λ good leaves. Let n_2 denote |I|, which is the number of surviving variables.

Phase 3

We will show that for any good leaf λ , we have

$$\operatorname{Corr}(C_{\lambda}, \operatorname{Par}_{n_2}) \le n^{-\varepsilon_d}$$
 (12)

where C_{λ} denotes $C_{1}|_{\rho_{\lambda}}$. This will prove the theorem, since we have by Fact 2.3,

$$\begin{aligned} \operatorname{Corr}(C_1, \operatorname{Par}_{n_1}) &\leq \underset{\lambda \sim T}{\mathbf{E}} [\operatorname{Corr}(C_\lambda, \operatorname{Par}_{n_2})] \\ &\leq \Pr[\lambda \text{ not good}] + \max_{\lambda \text{ good}} \operatorname{Corr}(C_\lambda, \operatorname{Par}_{n_2}) \\ &\leq 2\sqrt{q} + n^{-\varepsilon_d} \leq 2n^{-\varepsilon_d} \end{aligned}$$

where we have used the fact that $\operatorname{Par}_{n_1}|_{\rho_{\lambda}} = \pm \operatorname{Par}_{n_2}$ for each leaf λ , and also that $2\sqrt{q} \leq n^{-\varepsilon_d}$ for a large enough choice of the constant B.

It remains to prove (12). We do this in two steps.

In the first step, we set all large t-imbalanced gates to their most probable constant values. Formally, for a t-imbalanced threshold gate ϕ labelled by (w,θ) , we have $|\theta| \geq t \cdot ||w||_2$. We replace ϕ by a constant b_{ϕ} which is 1 if $\theta \geq t \cdot ||w||_2$ and by -1 if $-\theta \geq t \cdot ||w||_2$. This turns the circuit C_{λ} into a circuit C'_{λ} of at most the wire complexity of C_{λ} . Further, note that for any $x \in \{-1,1\}^{n_1}$, $C_{\lambda}(x) = C'_{\lambda}(x)$ unless there is a t-imbalanced threshold gate ϕ such that $\phi(x) \neq b_{\phi}(x)$. By the Chernoff bound (Theorem 2.21) the probability that this happens for any fixed imbalanced threshold gate is at most $\exp(-\Omega(t^2)) \leq \exp(-n^{\Omega(\delta_d)})$. By a union bound over the $\ell \leq n$ large threshold gates, we see that $\Pr_x[C_{\lambda}(x) \neq C'_{\lambda}(x)] \leq n \exp(-n^{\Omega(\delta_d)})$. In particular, we get by Fact 2.3

$$\operatorname{Corr}(C_{\lambda}, \operatorname{Par}_{n_2}) \leq \operatorname{Corr}(C'_{\lambda}, \operatorname{Par}_{n_2}) + n \exp(-\Omega(n^{\delta_d})) \leq \operatorname{Corr}(C'_{\lambda}, \operatorname{Par}_{n_2}) + \exp(-n^{\varepsilon_d}).$$
 (13)

In the second step, we further define a restriction tree T_{λ} such that C'_{λ} becomes a depth-(d-1) circuit with at most cn wires at all the leaves of T_{λ} . We first restrict by setting all variables that feed into any of the t-balanced gates. The number of variables set in this way is at most

$$\frac{\mu}{\sqrt{q}} \le 4p\sqrt{q} \cdot n^{1+\varepsilon_d} \le (pn) \cdot (4\sqrt{q}n^{\varepsilon_d}) \le \frac{pn}{8} \le \frac{n_2}{2}$$

for a large enough choice of the constant B. This leaves $n_3 \geq \frac{n_2}{2}$ variables still alive. Further, all the large t-balanced gates are set to constants with probability 1. Finally, by Proposition 4.3, we may set all but a set S of $n_4 = n_3/2cn^{\delta_d}$ variables to ensure that with probability 1, all the small gates depend on at most one input variable each. At this point, the circuit C'_{λ} may be transformed to a depth-(d-1) circuit C''_{λ} with at most as many wires as C'_{λ} , which is at most cn.

Note that the number of unset variables is $n_4 \geq pn/8cn^{\delta_d} \geq n^{1-2\delta_d}$, for large enough B. Hence, the number of wires is at most $cn \leq n_4^{\frac{1+\varepsilon_d}{1-2\delta_d}} \leq n_4^{(1+\varepsilon_d)(1+3\delta_d)} \leq n_4^{1+\varepsilon_{d-1}}$ for suitably large B. Thus, by the inductive hypothesis, we have

$$\operatorname{Corr}(C_{\lambda}^{"}, \operatorname{Par}_{n_4}) \leq O(n_4^{-\varepsilon_{d-1}}) \leq n^{-\varepsilon_d}/2$$

with probability 1 over the choice of the variables restricted in the second step. Along with (13) and Fact 2.3, this implies (12) and hence the theorem.

4.1.1 Strong correlation bounds for the generalized Andreev function

We now prove an exponentially strong correlation bound for the generalized Andreev function defined in Section 2.4 with any $\gamma < 1/6$. As in the case of Theorem 4.4, the proof proceeds by an iterative depth reduction. We prove the depth-reduction in a separate lemma.

▶ **Definition 4.5** (Simplicity). We call a threshold circuit C (t, d, w)-simple if there is a set R of $r \leq t$ threshold functions g_1, \ldots, g_r such that for every setting of these threshold functions to bits b_1, \ldots, b_r , the circuit C can be represented on the corresponding inputs (i.e., inputs x satisfying $g_i(x) = b_i$ for each $i \in [r]$) by a depth-d threshold gate of wire complexity at most w.

In particular, note that a (t, d, w)-simple circuit C may be expressed as

$$C(x) = \bigvee_{b_1, \dots, b_r \in \{-1, 1\}} \left(C_{b_1, \dots, b_r} \wedge \bigwedge_{i: b_i = -1} g_i \wedge \bigwedge_{i: b_i = 1} (\neg g_i) \right)$$

$$(14)$$

where each $C_{b_1,...,b_r}$ is a depth d circuit of wire complexity at most w. Further, note that the OR appearing in the above expression is disjoint (i.e. no two terms in the OR can be simultaneously true).

▶ Lemma 4.6. Let $d \ge 1$ be any constant and assume that ε_d , δ_d are defined as above. Say we are given any depth d threshold circuit C on n variables with at most $n^{1+\varepsilon_d}$ wires.

There is a restriction tree T of depth $n-n^{1-2\delta_d}$ with the following property: for a random leaf $\lambda \sim T$, let $\mathcal{E}(\lambda)$ denote the event that the circuit $C|_{\rho_\lambda}$ is $\exp(-n^{\varepsilon_d})$ -approximated by an $(n^{\delta_d}, d-1, n^{1+\varepsilon_d})$ -simple circuit. Then, $\Pr_{\lambda}[\neg \mathcal{E}(\lambda)] \leq \exp(-n^{\varepsilon_d})$.

Proof. Let ϕ_1, \ldots, ϕ_m be the threshold gates appearing at height 1 in the circuit C. We say that ϕ_i is large if $\operatorname{fan-in}(\phi_i) \geq n^{\delta_d}$ and small otherwise. Let $L = \{i \mid \phi_i \text{ large}\}$ and $S = [m] \setminus L$. Let $\ell = |L|$. Note that $\ell \leq n^{1+\varepsilon_d-\delta_d} \leq n$. Let $\ell = n^{\varepsilon_d}$.

As in the inductive case of Theorem 4.4, our construction proceeds in phases.

Phase 1

This is identical to Phase 1 in Theorem 4.4. We thus get a restriction tree of depth n/2 such that at *all* leaves of this tree, the resulting circuit is a read-2c circuit with at most cn wires. Let C_1 denote the circuit obtained at some arbitrary leaf of the restriction tree and let n_1 denote the number of variables.

Phase 2

This basic idea here is similar to Phase 2 from Theorem 4.4. However, there are technical differences from Theorem 4.4 since we apply a concentration bound to ensure that the circuit simplifies with high probability.

We restrict the circuit with a random restriction $\rho = (I, y) \sim \mathcal{R}_p^{n_1}$, where $p = n^{-\delta_d/2}$. As in Theorem 4.4, we have for some $t = \frac{1}{n^{\Omega(1)}}$, $q = p^{\Omega(1)}$, and for each $i \in [m]$,

$$\Pr_{\rho}[\phi_i|_{\rho} \ t\text{-balanced}] \le q \tag{15}$$

$$\Pr_{\rho}[\text{fan-in}(\phi_i|_{\rho}) > 2p \cdot \text{fan-in}(\phi_i)] \le \exp(-\Omega(n^{\delta_d/2}))$$
(16)

$$\Pr_{\rho=(I,y)}[|I| < \frac{n_1 p}{2}] \le \exp(-\Omega(np)) \tag{17}$$

Now, we partition L as $L = L_1 \cup \cdots \cup L_a$, where $a \leq \frac{1}{\varepsilon_d}$, as follows. The set L_j indexes all threshold gates of fan-in at least $n^{\delta_d + (j-1)\varepsilon_d}$ and less than $n^{\delta_d + j\varepsilon_d}$. We let ℓ_j denote $|L_j|$. For each $i \in L$, let Y_i be a random variable that is 1 if $\phi_i|_{\rho}$ is t-balanced and 0 otherwise. Note that this defines a collection of read-2c Boolean random variables (the underlying independent random variables are $\rho(k)$ for each $k \in [n_1]$).

Let $Z_j = \sum_{i \in L_j} Y_i$, the number of t-balanced gates in L_j . We have $\mathbf{E}[Z_j] = \sum_{i \in L_j} \mathbf{E}[Y_i] \le q\ell_j$ by (15). Thus, by an application of the read-2c Chernoff bound in Theorem 2.24, we have

$$\Pr[Z_i \ge 2q\ell_i] \le \exp\{-\Omega(q\ell_i/c)\}.$$

Assuming that $\ell_j \geq n^{3\delta_d/4}$ and $B = \delta_d/\varepsilon_d$ is a large enough constant, the right hand side of the above inequality is upper bounded bounded by $\exp\{-2n^{\varepsilon_d}\}$. On the other hand if $\ell_j < n^{3\delta_d/4}$, then $Z_j < n^{3\delta_d/4}$ with probability 1. Hence, we have $\Pr_{\rho=(I,y)}[Z_j \geq \max\{2q\ell_j, n^{3\delta_d/4}\}] \leq \exp\{-2n^{\varepsilon_d}\}$ and by a union bound

$$\Pr_{\rho=(I,y)}[\exists j \in [a], Z_j \ge \max\{2q\ell_j, n^{3\delta_d/4}\}] \le a \exp\{-2n^{\varepsilon_d}\}.$$
(18)

We call a set I generic if $|I| \geq \frac{n_1 p}{2}$ and fan-in $(\phi_i|_{\rho}) \leq 2p \cdot \text{fan-in}(\phi_i)$ for each $i \in L$. Let \mathcal{G} denote the event that I is generic. By (16) and (17), we know that $\Pr_I[\neg \mathcal{G}] \leq \ell \exp(-\Omega(n^{\delta_d/2})) + \exp(-\Omega(np)) \leq \exp(-n^{\delta_d/4})$. In particular, similar to Theorem 4.4, we get,

$$\Pr_{\rho=(I,y)}[\exists j \in [a], Z_j \ge \max\{2q\ell_j, n^{3\delta_d/4}\} \mid \mathcal{G}] \le a \exp\{-2n^{\varepsilon_d}\} + \exp(-n^{\delta_d/4}) \le 2a \exp\{-2n^{\varepsilon_d}\}.$$
(19)

We fix any generic I such that

$$\Pr_{u}[\exists j \in [a], Z_{j} \ge \max\{2q\ell_{j}, n^{3\delta_{d}/4}\}] \le 2a \exp\{-2n^{\varepsilon_{d}}\}.$$
(20)

Consider the restriction tree T that sets all the variables not in I. The tree leaves $n_2 \ge pn_1/2 = pn/4$ variables unfixed. We call a leaf λ of the tree good if for each $j \in [a]$ we have $Z_j < \max\{2q\ell_j, n^{3\delta_d/4}\}$ and bad otherwise. We have

$$\Pr_{\lambda \sim T}[\lambda \text{ a bad leaf}] \le 2a \exp(-2n^{\varepsilon_d}) \tag{21}$$

For good leaves λ , we show how to approximate $C_{\lambda} := C_1|_{\rho_{\lambda}}$ as claimed in the lemma statement

For the remainder of the argument, fix any good leaf λ . We partition $[a] = J_1 \cup J_2$ where $J_1 = \{j \in [a] \mid Z_j < 2q\ell_j\}$. Note that for any $j \in J_1$, we have

$$\begin{split} \sum_{i \in L_j} Y_i \cdot \text{fan-in}(\phi_i|_{\rho_\lambda}) &\leq \sum_{i \in L_j} Y_i \cdot 2p \cdot \text{fan-in}(\phi_i) \\ &\leq 2p \cdot n^{\delta_d + j \cdot \varepsilon_d} \cdot Z_j \leq n^{\delta_d + j \cdot \varepsilon_d} \cdot 4pq\ell_j \\ &= 4pqn^{\varepsilon_d} \cdot \ell_j \cdot n^{\delta_d + (j-1) \cdot \varepsilon_d} \leq 4pqn^{\varepsilon_d} n^{1+\varepsilon_d} \\ &= 4nqn^{1+2\varepsilon_d} \end{split}$$

where for the last inequality, we have used the fact that since we have ℓ_j gates of fan-in at least $n^{\delta_d+(j-1)\varepsilon_d}$ each, we must have $\ell_j \cdot n^{\delta_d+(j-1)\varepsilon_d} \leq n^{1+\varepsilon_d}$, the total wire complexity of the circuit.

In particular, we can bound the total fan-in of all the t-balanced gates indexed by $\bigcup_{j \in J_1} L_j$ by

$$\sum_{j \in J_1} \sum_{i \in L_j} Y_i \cdot \text{fan-in}(\phi_i|_{\rho_\lambda}) \le \frac{4pqn^{1+2\varepsilon_d}}{\varepsilon_d}.$$
 (22)

Phase 3

We proceed in two steps as in Theorem 4.4. Since the steps are very similar, we just sketch the arguments. In the first step, we replace all large t-imbalanced gates by their most probable values. This yields a circuit C'_{λ} of at most the wire complexity of C_{λ} and such that

$$\Pr_{x}[C_{\lambda}(x) \neq C_{\lambda}'(x)] \leq \ell \exp(-n^{\Omega(\delta_d)}) \leq n \exp(-n^{\Omega(\delta_d)}) \leq \exp(-n^{\varepsilon_d}).$$
 (23)

In the second step, we construct another restriction tree rooted at λ that simplifies the circuit to the required form. We first restrict by setting all variables that feed into the t-balanced gates that are indexed by $\bigcup_{j \in J_1} L_j$. By (22), the number of variables set is bounded by

$$\frac{4pqn^{1+2\varepsilon_d}}{\varepsilon_d} \leq \frac{4pn \cdot n^{\varepsilon_d - \Omega(\delta_d)}}{\varepsilon_d} \leq \frac{pn}{8} \leq \frac{n_2}{2}$$

for a large enough choice of the constant B. This sets all the t-balanced gates indexed by $\bigcup_{j \in J_1} L_j$ to constants while leaving $n_3 \geq \frac{n_2}{2}$ variables still alive. Finally, by Proposition 4.3, we may set all but a set of $n_4 = n_3/2cn^{\delta_d}$ variables to ensure that with probability 1, all the small gates depend on at most one input variable each. We may replace the small gates by the unique variable they depend on or a constant (if they do not depend on any variable) without increasing the wire complexity of the circuit. Call the circuit thus obtained $C_{\lambda}^{"}$.

At this point, the only threshold gates at height 1 in the circuit C''_{λ} are the gates indexed by the t-balanced gates in $\bigcup_{j\in J_2} L_j$. But by the definition of J_2 , there can be at most $\frac{1}{\varepsilon_d} \cdot n^{3\delta_d/4} \leq n^{\delta_d}$ of them. For every setting of these threshold gates to constants, the circuit becomes a depth-(d-1) circuit of size at most $n^{1+\varepsilon_d}$. Hence, we have a $(n^{\delta_d}, d-1, n^{1+\varepsilon_d})$ -simple circuit, as claimed.

Note that the number of variables still surviving is given by $n_4 \ge pn/8cn^{\delta_d} \ge n^{1-2\delta_d}$, for a large enough choice of the parameter B. Hence, the restriction tree constructed satisfies the required depth constraints.

For a random leaf $\nu \sim T$, the probability $\mathcal{E}(\nu)$ does not occur is at most the probability that in Phase 2, the leaf sampled is bad. By (21), this is bounded by $2a \exp(-2n^{\varepsilon_d}) \leq \exp(-n^{\varepsilon_d})$ as claimed.

We now prove the correlation bound for threshold circuits with the generalized Andreev function. For the sake of induction, it helps to prove a statement that is stronger in two ways: firstly, we consider any function $F_a = F(a, \cdot)$ where $a \in \{-1, 1\}^{4n}$ has high Kolmogorov complexity and the input to F_a is further restricted by an arbitrary restriction ρ that leaves a certain number of variables alive; secondly, we prove a correlation bound against circuits which are the AND of a small threshold circuit with a small number of threshold gates.

Formally, say that $f: \{-1,1\}^n \to \{-1,1\}$ is (N,d,t,α) -intractable if for any restriction ρ on n variables that leaves $m \geq N$ variables unset, any depth-d threshold circuit C on m variables of wire complexity at most $m^{1+\varepsilon_d}$, and any set S of at most t threshold functions, we have

$$\operatorname{Corr}(f, C \wedge \bigwedge_{g \in S} g) \le \alpha.$$

The stronger correlation bound is the following.

▶ Theorem 4.7 (Generalized version of strong correlation). Fix any constant $d \ge 1$. Let $a \in \{-1,1\}^{4n}$ be any string with $K(a) \ge 3n$. Then, the function F_a is $(n^{1-\varepsilon_d},d,n^{\varepsilon_d},\exp(-n^{\varepsilon_d}/2))$ -intractable.

The proof is by induction on d. The properties of F_a are only used to prove the base case of the theorem, which can then be used to prove the induction case using Lemma 4.6. We prove the base case separately below (we assume that the constant B > 0 is large enough so that this implies the base case of the theorem stated above).

▶ **Lemma 4.8** (Base case of induction). Let $a \in \{-1,1\}^{4n}$ be any string with $K(a) \geq 3n$. Then, the function F_a is $(\sqrt{n}, 1, \sqrt{n}, \exp(-n^{\Omega(1)}))$ -intractable.

Proof. Let $\gamma < 1/6$ in the definition of the Generalized Andreev function in Section 2.4. Let τ be any restriction of n variables leaving $m \ge \sqrt{n}$ variables unfixed. Define $f := F_a|_{\tau}$. Let C be a conjunction of $\sqrt{n} + 1$ threshold gates each on m variables. We wish to prove that

$$\operatorname{Corr}(f, C) \le \exp(-n^{\Omega(\gamma)}).$$

We build a restriction tree T for C of depth $m-n^{\gamma}$, by restricting all but n^{γ} arbitrarily chosen variables. For any leaf ℓ of T, the restricted circuit $C_{\ell} := C|_{\rho_{\ell}}$ is a conjunction of $\sqrt{n}+1$ threshold gates each on n^{γ} variables. By Lemma 2.9, each threshold function can be described using $n^{2\gamma}$ bits. Hence, the entire circuit can be described in a standard way using $(\sqrt{n}+1)\cdot O(n^{2\gamma}) < n$ bits. Then, by Lemma 2.20, we have $\operatorname{Corr}(f|_{\rho_{\ell}}, C_{\ell}) \leq \exp(-n^{\Omega(\gamma)})$.

Proof of Theorem 4.7. We only need to prove the inductive case. Assume that $d \geq 2$ is given. Fix any restriction ρ that sets all but $m \geq n^{1-\varepsilon_d}$ variables and let $f = F_a|_{\rho}$. Let C be a depth-d threshold circuit on the surviving variables of wire complexity at most $m^{1+\varepsilon_d}$. Let S be any set of at most n^{ε_d} threshold functions on the m variables. We need to show that $\operatorname{Corr}(f, C \wedge \bigwedge_{g \in S} g) \leq \exp(-n^{\varepsilon_d/2})$.

Apply Lemma 4.6 to circuit C to find a restriction tree T as guaranteed by the statement of the lemma. By Fact 2.3, we have

$$\operatorname{Corr}(f, C \wedge \bigwedge_{g \in S} g) \leq \underset{\ell \sim T}{\mathbf{E}} \left[\operatorname{Corr}(f_{\ell}, C_{\ell} \wedge \bigwedge_{g \in S} g_{\ell}) \right]$$

$$\leq \Pr_{\ell} \left[\neg \mathcal{E}(\ell) \right] + \max_{\ell : \mathcal{E}(\ell) \text{ holds}} \operatorname{Corr}(f_{\ell}, C_{\ell} \wedge \bigwedge_{g \in S} g_{\ell})$$
(24)

where f_{ℓ} denotes $f|_{\rho_{\ell}}$ and similarly for C_{ℓ} and g_{ℓ} , and $\mathcal{E}(\ell)$ is the event defined in the statement of Lemma 4.6.

Fix any leaf ℓ so that $\mathcal{E}(\ell)$ holds. We want to bound $\operatorname{Corr}(f_{\ell}, C_{\ell} \wedge \bigwedge_{g \in S} g_{\ell})$. By definition of $\mathcal{E}(\ell)$, we know that C_{ℓ} is $\exp(-m^{\varepsilon_d})$ -approximated by a $(m^{\delta_d}, d-1, m^{1+\varepsilon_d})$ -simple circuit C'_{ℓ} . This implies that $C_{\ell} \wedge \bigwedge_{g \in S} g_{\ell}$ is $\exp(-m^{\varepsilon_d})$ -approximated by $C'_{\ell} \wedge \bigwedge_{g \in S} g_{\ell}$. Hence, we have

$$\operatorname{Corr}(f_{\ell}, C_{\ell} \wedge \bigwedge_{g \in S} g_{\ell}) \leq \operatorname{Corr}(f_{\ell}, C'_{\ell} \wedge \bigwedge_{g \in S} g_{\ell}) + \exp(-m^{\varepsilon_d}). \tag{25}$$

Further, by the definition of simplicity and its consequence (14), we know that there exist $r \leq m^{\delta_d}$ threshold functions $h_1^{\ell}, \ldots, h_r^{\ell}$ such that

$$C'_{\ell} = \bigvee_{b_1, \dots, b_r \in \{-1, 1\}} C_{b_1, \dots, b_r} \wedge \bigwedge_{i: b_i = -1} h_i^{\ell} \wedge \bigwedge_{i: b_i = 1} \neg h_i^{\ell}$$

where each $C_{b_1,...,b_r}$ is a depth d-1-threshold circuit of size at most $m^{1+\varepsilon_d}$ and the OR above is disjoint. This further implies that

$$C'_{\ell} \wedge \bigwedge_{g \in S} g_{\ell} = \bigvee_{b_1, \dots, b_r \in \{-1, 1\}} \left(C_{b_1, \dots, b_r} \wedge \bigwedge_{i: b_i = -1} h_i^{\ell} \wedge \bigwedge_{i: b_i = 1} \neg h_i^{\ell} \wedge \bigwedge_{g \in S} g_{\ell} \right)$$

$$(26)$$

and the OR remains disjoint.

Note that we may apply the induction hypothesis to obtain a bound on the correlation with each term in the OR at this point, since the number of surviving variables is at least $m_1 = m^{1-2\delta_d} \geq n^{1-\varepsilon_d-2\delta_d} \geq n^{1-\varepsilon_{d-1}}$ (throughout, we assume that B is a large enough constant for many of the inequalities to hold); and the wire complexity of each depth-(d-1) circuit C_{b_1,\dots,b_r} is at most $m^{1+\varepsilon_d} \leq m_1^{(1+\varepsilon_d)/(1-2\delta_d)} \leq m_1^{1+\varepsilon_d+3\delta_d} \leq m_1^{1+\varepsilon_{d-1}}$; further, the number of threshold functions in each term is at most $n^{\varepsilon_d} + n^{\delta_d} < m^{\varepsilon_{d-1}}$. Thus, by the inductive hypothesis, we obtain for any b_1,\dots,b_r ,

$$\operatorname{Corr}(f,C_{b_1,\dots,b_r} \wedge \bigwedge_{i:b_i=-1}^{} h_i^{\ell} \wedge \bigwedge_{i:b_i=1}^{} \neg h_i^{\ell} \wedge \bigwedge_{g \in S}^{} g_{\ell}) \leq \exp(-n^{\varepsilon_{d-1}/2}).$$

Using the fact that the OR in (26) is disjoint, from Fact 2.3, we obtain

$$\operatorname{Corr}(f, C'_{\ell} \wedge \bigwedge_{q \in S} g_{\ell}) \leq 2^r \cdot \exp(-n^{\varepsilon_{d-1}/2}) \leq 2^{n^{\delta_d}} \cdot \exp(-n^{\varepsilon_{d-1}/2}) \leq \exp(-n^{\varepsilon_d}).$$

Putting the above together with (24) and (25), we obtain

$$\operatorname{Corr}(f, C \wedge \bigwedge_{g \in S} g) \leq \exp(-m^{\varepsilon_d}) + \exp(-n^{\varepsilon_d}) \leq \exp(-n^{\varepsilon_d/2}).$$

which proves the induction case and hence the theorem.

▶ Corollary 4.9 (Correlation bounds for Andreev's function). For any $d \ge 1$, any depth-d threshold circuit C of wire complexity at most $n^{1+\varepsilon_d}$ satisfies $Corr(C, F) \le 2 \exp(-n^{\varepsilon_d/2})$.

Proof. For a random $a \in \{-1, 1\}^{4n}$, we know by Fact 2.14 that $K(a) \geq 3n$ with probability $1 - \exp(-\Omega(n))$. For each such a, by Theorem 4.7, we have $\operatorname{Corr}(C_a, F_a) \leq \exp(-n^{\varepsilon_d/2})$, where C_a is the circuit obtained by substituting $x_1 = a$ in C. Hence, we have

$$\operatorname{Corr}(C,F) \leq \mathop{\mathbf{E}}_{a}[\operatorname{Corr}(C_{a},F_{a})] \leq \exp(-\Omega(n)) + \exp(-n^{\varepsilon_{d}/2}) \leq 2\exp(-n^{\varepsilon_{d}/2})$$

as claimed.

4.2 Proof of Main Structural Lemma

We need the following definitions and facts that have appeared many times before in the literature on threshold functions (see, e.g., [9]).

Let $\varepsilon \in [0,1]$ be a real parameter. We say that $w \in \mathbb{R}^n$ is ε -regular if for each $i \in [n]$, $|w_i| \leq \varepsilon \cdot ||w||_2$.

Assume for simplicity that the co-ordinates of the vector w are sorted so that $|w_1| \ge |w_2| \ge \cdots \ge |w_n|$. Let $w_{>i} \in \mathbb{R}^{n-i}$ denote the vector obtained by removing the first i co-ordinates of w. We define the ε -critical index of w be the least $K = K(\varepsilon)$ so that the vector $w_{>K}$ is ε -regular. Note that K = 0 if w is already ε -regular and we define K = n if the ε -critical index is not defined.

We say that an *n*-variable threshold gate ϕ labelled by (w, θ) is ε -regular if w is. Similarly, the ε -critical index of ϕ is defined to be the ε -critical index of w.

the ε -critical index of ϕ is defined to be the ε -critical index of w. Also, we define $L=L(\varepsilon)=\frac{100\log^2(1/\varepsilon)}{\varepsilon^3}$ for a large constant A that will be made precise later.

The Berry Esseen theorem (see, e.g., [12]) yields the following standard anticoncentration lemma for linear functions. (See [9, Corollary 2.2] for this particular statement.)

▶ Lemma 4.10 (Anticoncentration for regular linear functions). Let $w \in \mathbb{R}^n$ be ε -regular and let $J \subseteq \mathbb{R}$ be any interval. Then,

$$\left| \Pr_{x \in \{-1,1\}^n} [\langle w, x \rangle \in J] - \Phi(J) \right| \le O(\varepsilon)$$

where $\Phi(\cdot)$ denotes the cdf of the standard Gaussian with mean 0 and variance $||w||_2^2$. In particular, if |J| denotes the length of J, then

$$\Pr_x[\langle w, x \rangle \in J] \le \frac{|J|}{\|w\|_2} + O(\varepsilon).$$

We now proceed with the proof of Lemma 4.1. We start with an easier case of the lemma for regular threshold gates. Throughout, we work with random restrictions sampled from \mathcal{R}_p^n where $p \in [0,1]$ is the probability from the statement of Lemma 4.1: equivalently, we pick a pair (I,y) where $I \subseteq [n]$ and $y \in \{-1,1\}^{n-|I|}$. Let $\varepsilon = p^{\frac{1}{8}}$. Let $t = p^{-\frac{1}{16}}$.

Let the threshold gate ϕ be labelled by pair (w, θ) , where $w \in \mathbb{R}^n$. We may assume that the variables of the threshold gate have been sorted so that $|w_1| \ge |w_2| \ge \cdots |w_n|$. Note that after applying a restriction ρ , the threshold gate ϕ_{ρ} is labelled by pair (w', θ') , where w' is the restriction of w to the coordinates in I and

$$\theta' = \theta'(\rho) = \theta - \langle w'', y \rangle. \tag{27}$$

Above, we use w'' to denote the vector w restricted to the indices in $[n] \setminus I$.

For a random restriction $\rho \sim \mathcal{R}_p^n$, define the following "bad" events:

- 1. $\mathcal{B}(\rho)$: ϕ_{ρ} is t-balanced: i.e., $\theta' \leq t \cdot ||w'||_2$. This is the event whose probability we want to upper bound.
- 2. $\mathcal{B}_1(\rho)$: $\sum_{i \in I} w_i^2 \ge \sqrt{p} \cdot ||w||_2^2$.
- 3. $\mathcal{B}_2^k(\rho)$ (k a parameter): One of the first k variables x_1, \ldots, x_k is set to * by ρ .

We have the following simple upper bounds on the probabilities of some of the above bad events:

- Since each variable is set to * with probability p, we have $\mathbf{E}_{\rho}[\sum_{i \in I} w_i^2] = p \cdot ||w||_2^2$. By Markov's inequality, we have $\Pr_{\rho}[\mathcal{B}_1(\rho)] \leq \sqrt{p}$.
- By a union bound, for any k, we have $\Pr_{\rho}[\mathcal{B}_{2}^{k}(\rho)] \leq pk$.

We start with a simpler subcase of the lemma that follows almost directly from Lemma 4.10. We assume throughout that p is a small enough constant, since otherwise the statement of Lemma 4.1 is trivial.

▶ **Lemma 4.11** (The regular case). Say that w is ε -regular. Then $\Pr_{\rho}[\mathcal{B}(\rho)] \leq p^{\Omega(1)}$.

Proof. We bound $\Pr_{\rho}[\mathcal{B}(\rho)]$ as follows.

$$\Pr_{\rho}[\mathcal{B}(\rho)] \leq \Pr_{\rho}[\mathcal{B}_{1}(\rho)] + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg(\mathcal{B}_{1}(\rho))]
\leq \sqrt{p} + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg(\mathcal{B}_{1}(\rho))].$$
(28)

Now, note that the event $\neg \mathcal{B}_1(\rho)$ only depends on the choice of $\rho^{-1}(*) = I$. Hence we can condition on an I so that this event occurs; choosing ρ is now equivalent to choosing a random assignment y to the variables in $[n] \setminus I$.

We have $\theta' = \theta - \langle w'', y \rangle$. Using the fact that $\mathcal{B}_1(\rho)$ doesn't occur, we have

- $||w''||_2 \ge ||w||_2 \sqrt{1 \sqrt{p}} \ge ||w||_2 / 2$. Using the ε -regularity of w, for each $i \notin I$, we have $|w_i| \le (\varepsilon) ||w||_2 \le 2\varepsilon ||w''||_2$. Thus, w'' is 2ε -regular.
- $||w'||_2 \le p^{1/4} ||w||_2 \le 2p^{1/4} ||w''||_2,$

Using the above, we can see that the probability that

$$\begin{split} \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{1}(\rho)] &\leq \Pr_{y}[|\theta'| \leq t \cdot \|w'\|_{2}] \leq \Pr_{y}[|\theta'| \leq 2tp^{1/4} \cdot \|w''\|_{2}] \\ &\leq \Pr_{y}[\langle w'', y \rangle \in [\theta - 2tp^{1/4} \cdot \|w''\|_{2}, \theta + 2tp^{1/4} \cdot \|w''\|_{2}]] \\ &\leq 4tp^{1/4} + O(\varepsilon) = O(\varepsilon) = p^{\Omega(1)} \end{split}$$

where the final inequality uses the anti-concentration bound in Lemma 4.10. Putting the above together with (28), we are done.

Proof of Lemma 4.1. The proof of the lemma is a standard case analysis based on the ε -critical index of the threshold gate ϕ (see [43, 33, 9, 28]).

The first case is when the critical index $K \leq L$. In this case, we bound the probability of $\mathcal{B}(\rho)$ by

$$\Pr_{\rho}[\mathcal{B}(\rho)] \leq \Pr_{\rho}[\mathcal{B}_{2}^{K}(\rho)] + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{K}(\rho)]
\leq pK + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{K}(\rho)] \leq \sqrt{p} + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{K}(\rho)]$$
(29)

where the final inequality follows from the fact that $pK \leq pL \leq \sqrt{p}$ by our choice of parameters. The event $\neg \mathcal{B}_2(\rho)$ only depends on the choice of the sub-restriction $\rho|_{[K]}$ and we can condition on $\rho|_{[K]}$ so that this event occurs. From now on, the random choice will be a restriction $\rho' \sim \mathcal{R}_p^{n-K}$ on the remaining variables.

Since the restricted linear function is now ε -regular by the definition of the ε -critical index, we can apply Lemma 4.11 to conclude that $\Pr_{\rho'}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_2^K(\rho)] \leq p^{\Omega(1)}$. Along with (27), this implies the lemma in the case that $K \leq L$.

The second case is when K > L. As in previous cases, we first condition on some bad event not occurring. We have

$$\Pr_{\rho}[\mathcal{B}(\rho)] \leq \Pr_{\rho}[\mathcal{B}_{2}^{L}(\rho)] + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{L}(\rho)]
\leq pL + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{L}(\rho)] \leq \sqrt{p} + \Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{L}(\rho)].$$
(30)

As in Lemma 4.11, we can condition on a fixed I so that $\neg \mathcal{B}_2^L(\rho)$ occurs (i.e. none of the first L variables belong to I). We then use the following claim that is implicit in [9].

▶ Proposition 4.12. Let $L' = \frac{10r\log(1/\varepsilon)}{\varepsilon^2}$ and assume that K > L'. Let y be a random assignment to any set of variables including the first $L' = \frac{10r\log(1/\varepsilon)}{\varepsilon}$ variables. Then, the probability over y that the restricted threshold gate is not $(\frac{1}{\varepsilon})$ -imbalanced is at most 2^{-r} .

Applying the above proposition with L' = L and $r = 10 \log(1/\varepsilon)$, we have $\Pr_{\rho}[\mathcal{B}(\rho) \mid \neg \mathcal{B}_{2}^{L}(\rho)] \le \varepsilon^{10}$. Putting this together with (30), we have the claimed upper bound on $\Pr_{\rho}[\mathcal{B}(\rho)]$ in the case that K > L.

We give a proof sketch of Proposition 4.12 in Section C.

5 Satisfiability algorithms beating brute-force search

In this section, we give satisfiability algorithms beating brute force search for bounded-depth threshold circuits with few wires. Until now, such algorithms were only known for threshold circuits of depth 2. We will assume that each threshold gate on m input bits is given as a pair (w, θ) , where $w \in \mathbb{Z}^m$ and $\theta \in \mathbb{Z}$, and θ as well as each component of w has bit complexity poly(n). Note that this assumption is without loss of generality for a threshold function, and that some assumption on representability of threshold functions is necessary in an algorithmic context.

The satisfiability algorithm relies on an algorithmic version of Lemma 4.6, along with a couple of additional ideas. Essentially, we use the algorithmized version of the lemma to reduce the satisfiability of bounded-depth circuits to satisfiability of ANDs of threshold functions, which we can then solve using a recent result of Williams, stated below.

▶ Theorem 5.1 ([48]). There is a deterministic algorithm, which given a bounded-depth circuit C on n variables of size $2^{n^{\circ(1)}}$ with ANDs, ORs and threshold gates, and with the threshold gates appearing only at the bottom layer, decides if C is satisfiable in time $2^{n-n^{\varepsilon'}}$ poly(n), where $\varepsilon' > 0$ is a constant that depends only on the depth of the circuit.

We also need the following fact about threshold gates on n input bits: the set of inputs evaluating to 1 (and dually, the set of inputs evaluating to -1) of a linear threshold gate can be enumerated in time proportional to the number of such inputs, modulo a poly(n) factor.

▶ Proposition 5.2. Let (w, θ) represent a threshold function ϕ on m input bits, where $w \in \mathbb{Z}^m$ and $\theta \in \mathbb{Z}$ are integers of bit complexity poly(m). Let S be the set of inputs on which ϕ evaluates to 1. Then S can be enumerated in time |S|poly(n).

Proof. We will show how to construct a decision tree for ϕ in time |S|poly(n), where S is the set of inputs on which ϕ takes value 1. Given a decision tree of size at most |S|poly(n), it is easy to enumerate the set of inputs on which ϕ takes value 1 in time |S|poly(n) by scanning through leaves labelled 1 and outputting all assignments corresponding to any such leaf.

The decision tree is constructed recursively as follows. Check if ϕ restricted according to the current partial assignment is satisfiable (in the sense that there is a total assignment consistent with the partial assignment for which ϕ evaluates to 1). Note that satisfiability of a linear threshold gate with polynomial bit complexity of the weights can be done trivially in polynomial time. If the satisfiability check fails, make the current node a leaf and label it with -1. If it succeeds, check if the current partial assignment is falsifiable. If this check fails, make the current node a leaf and label it with 1. Otherwise, branch on an arbitrary unassigned variable and recurse.

Clearly, this decision tree can be constructed with polynomial work at each node, and hence in time $N \operatorname{poly}(n)$, where N is the number of leaves of the tree. We show that $N \leq |S|n$. Indeed, we prove inductively that for any internal node v of the tree of height $h \geq 1$, the number of -1 leaves of the tree rooted at v is at most h times the number of 1 leaves, from which the claim follows as the height of the tree $\leq n$.

For the inductive claim, the base case h=1 is clear as any node at height 1 must have one leaf labelled 1 and the other labelled -1. Assume the claim for height h. Consider a node v at height h+1. Either one of its children is a leaf, or not. If one of the children is a leaf, then the other one v' is not and by the induction hypothesis, since it is of height h, has at most h times as many -1 leaves as 1 leaves. The number of -1 leaves of v is at most one plus the number of -1 leaves of v', and hence at most h+1 times the number of 1 leaves. In case

both children of v are internal nodes, then they are both of height at most h, and by the induction hypothesis, both have at most h times as many -1 leaves as 1 leaves, which implies that the same holds for v.

- \triangleright **Definition 5.3.** We use THR to refer to the class of linear threshold functions. We use AND \circ THR to refer to the class of polynomial-size circuits with an AND gate at the top and threshold gates at the bottom layer.
- ▶ Theorem 5.4. For each integer d > 0, there is a constant $\epsilon_d > 0$ such that satisfiability of a depth d threshold circuit with at most $n^{1+\epsilon_d}$ wires on n variables can be solved by a randomized algorithm in time $2^{n-\Omega(n^{\epsilon_d})}\operatorname{poly}(n)$.

Proof. As the proof follows the proof of Lemma 4.6 closely, we just give a sketch. Call a circuit depth-d AND \circ THR-skew if the top gate is an AND and all but one child of the top gate is a bottom-level threshold gate, with the possibly exceptional child being a depth-d-1 threshold circuit with few wires. We follow the depth reduction argument in the lemma to give a recursive algorithm which reduces satisfiability of polynomial-size depth-d AND \circ THR-skew circuits to the satisfiability of polynomial-size depth-d-1 AND \circ THR-skew circuits by appropriately restricting variables.

For the base case d=1, we simply appeal to the algorithm given by Theorem 5.1, which solves satisfiability of AND \circ THR circuits of polynomial size in time $2^{n-n^{\varepsilon'}}$ poly(n) for some constant $\varepsilon'>0$.

For the inductive case, we simulate the proof of Lemma 4.6, which performs and analyzes a certain kind of adaptive random restriction. Various bad events might happen at Phases 2 and 3 of this random restriction process, however each step of the restriction process as well as the check that a bad event happens can be implemented in polynomial time. Moreover, the probability that a bad event happens is at most $2^{-n^{\epsilon_d}}$. Whenever a bad event happens, we simply do brute force search on the remaining variables of the circuit, but thanks to the exponentially small probability that a bad event happens, with high probability, we only spend time $2^{n-n^{\epsilon_d}}$ on such brute force searches.

In Phase 3 of the restriction process, we replace imbalanced gates by their most probable values. This changes the functionality of the circuit and might lose us satisfying assignments or give us new invalid satisfying assignments. To get around this, for each such imbalanced gate, we use Proposition 5.2 to efficiently enumerate the inputs evaluating to the minority value for each imbalanced gate, and for each such input check whether it satisfies the original circuit. If it does, we just output 'yes'. We also append to the top gate of the skew circuit a child representing the assignment of the imbalanced gate to its majority value – this needs to be done so that we don't end up with "false positives" in the base case of the recursive algorithm. Although each such false positive can be tested, there might be too many of them, and this could destroy all the savings we accrue through the course of the algorithm. The total time spent in enumerating minority values of imbalanced gates is again at most $2^{n-n^{\epsilon_d}}$ poly(n), with high probability, using the efficient enumeration and the imbalance property.

Finally, there are a few balanced gates – with high probability at most $O(n^{\delta_d})$ of them – for which we need to try all possible values. This could be expensive, but is compensated for by an increased savings for depth d-1, just by setting the constant B large enough in the proof of Lemma 4.6. We also need to set B large enough so that the savings given by the application of Williams' algorithm in the base case overwhelms the loss due to branching on balanced threshold gates at depth d=2.

Thus the total running time, once B is chosen appropriately, is $2^{n-\Omega(n^{\epsilon_d})} \operatorname{poly}(n)$, using the fact that $\epsilon_d < \epsilon_{d-1} < \ldots < \epsilon_2$.

6 Threshold formulas

A threshold formula is a threshold circuit such that the fan-out of each gate is at most 1. A formula can be viewed as a tree. Note that a depth-2 threshold circuit can always be converted to a threshold formula without increasing either the wire complexity or the gate complexity (recall that the gate complexity only measures the number of *non-input* gates).

Let $F: \{-1,1\}^{4n} \times \{-1,1\}^n \to \{-1,1\}$ be the generalized Andreev function defined in Section 2.4. Recall that F is constructed with $(n,n^{\gamma},m=0.9n^{\gamma},2^{-n^{\Omega(\gamma)}})$ bit fixing extractor $E: \{-1,1\}^n \to \{-1,1\}^m$, and $(1/2-O(2^{-m/4}),2^{m/2})$ list decodable code Enc: $\{-1,1\}^{4n} \to \{-1,1\}^{2^m}$.

▶ **Theorem 6.1.** Any threshold formula on n variables with at most $n^{1.5-\gamma}$ wires for has correlation at most $\exp(-n^{\Omega(\gamma)})$ with the generalized Andreev function.

Proof. Let C be a threshold formula with n inputs and $s=n^{1.5-\gamma}$ wires. Let L be the number of leaves in the formula tree; then $L \leq s \leq 2L$. We build a restriction tree T for C up to depth n-pn for $p=n^{\gamma/2}/n$, by greedily restricting the most frequent variables appearing in the formula. Since the most frequent variable appears at least L/n times in C, after restricting one variable, the formula tree has at most L(1-1/n) leaves left. Continue until pn variables left unrestricted; then the number of remaining leaves is at most $L \cdot \frac{n-1}{n} \cdot \frac{n-2}{n-1} \cdot \cdots \cdot \frac{pn}{pn+1} = pL$. Thus, for any leaf l of T, the restricted formula $C|_{\rho_l}$ (on $pn=n^{\gamma/2}$ variables) has $s(C|_{\rho_l}) \leq 2pL \leq 2ps$ wires, and by Proposition 2.16, the description length is at most $O(p^2s^2) \leq O(n^{1-\gamma}) < n$. Let $a \in \{-1,1\}^{4n}$ be a string with Kolmogorov complexity $K(a) \geq 3n$, and let $F_a(x) := F(a,x)$. Then, by Lemma 2.20, $Corr(F_a,C) \leq \exp(-n^{\Omega(\gamma)})$.

Therefore, for any formula D with 5n inputs and $n^{1.5-\gamma}$ wires, $\Pr_x[F(a,x)=D(a,x)] \leq 1/2 + \exp(-n^{\Omega(\gamma)})$. Since a random $a \in \{-1,1\}^{4n}$ has $K(a) \geq 3n$ with probability $1-2^{-\Omega(n)}$, the correlation of D and F is at most $2^{-\Omega(n)} + \exp(-n^{\Omega(\gamma)}) = \exp(-n^{\Omega(\gamma)})$.

7 Correlation bounds for AC^0 with a few threshold gates

Following Gopalan and Servedio [15], we define $TAC^{0}[k]$ to be the class of constant-depth circuits made up of AND and OR gates and at most k arbitrary threshold gates.

We prove upper bounds on the noise sensitivity of small depth-d TAC⁰[k] circuits for k much smaller than $n^{1/2(d-1)}$. The basic idea is the same as in Theorem 3.1, but we also need to use the following powerful result of Kane [23, Corollary 3].

- ▶ **Definition 7.1** (Polynomial Threshold functions). A Boolean function $f: \{-1,1\}^n \to \{-1,1\}$ is a degree-D Polynomial Threshold function if there is a degree-D polynomial p(x) such that $f(x) = \operatorname{sgn}(p(x))$ for all $x \in \{-1,1\}^n$.
- ▶ **Lemma 7.2** (Kane [23]). Let f be a degree-D PTF. Then, for any p > 0,

$$NS_p(f) \le \sqrt{p}(\log(1/p))^{O(D\log D)} 2^{O(D^2\log D)}$$
.

The main theorem of the section is the following.

▶ **Theorem 7.3.** Fix any constant $d \ge 1$. Let C be a depth-d $TAC^0[k]$ circuit with at most M gates overall. Then, for any $p, q \in [0, 1]$ and any $D \ge 1$, we have

$$NS_{n^{d-1}q}(C) \le O(k\alpha(p, D) + \alpha(q, D) + M(10pD)^D)$$

where $\alpha(p,D) := \sqrt{p}(\log(1/p))^{O(D\log D)}2^{O(D^2\log D)}$ and $O(\cdot)$ hides an absolute constant (independent of d).

Proof. This is a standard switching argument (see, e.g., [19]) augmented with the ideas of Theorem 3.1. We assume throughout that $q \leq \frac{1}{2}$ w.l.o.g. since otherwise $\alpha(q, D) \geq q \geq \frac{1}{2}$ and the claim is trivial.

We say that a threshold gate is a true threshold gate if it is not an AND or OR gate.

For any parameters $k_1, d_1, t_1, s_1 \in \mathbb{N}$ with $d_1 \geq 2$, we define $TAC^0[k_1, d_1, t_1, s_1]$ to be the class of constant-depth circuits made up of AND, OR and threshold gates such that:

- \blacksquare The overall depth is at most d_1 ,
- The total number of gates at depth at most $d_1 2$ in the circuit is at most s_1 ,
- All the true threshold gates are at depth at most $d_1 2$ and there are at most k_1 of them, and
- The bottom fan-in of the circuit (i.e. the maximum fan-in of a gate at depth $d_1 1$) is at most t_1 .

Note that the circuit C in the statement of the theorem is in the class $TAC^0[k, d+1, 1, M]$, since we may replace the input literals with (say) AND gates of fan-in 1 at the expense of increasing the depth by 1 but in the process satisfying all the criteria of the above definition. We prove the following stronger statement: for any p, q, D as in the statement of the theorem, and any C from the class $TAC^0[k, d, D, M]$ with $d \ge 2$, we have

$$NS_{p^{d-2}q}(C) = \mathbf{E}_{\rho_d}[Var(C|_{\rho_d})] \le O(k\alpha(p, D) + \alpha(q, D) + M(10pD)^D)$$
(31)

where $\rho_d \sim \mathcal{R}_{p_d}^n$ and $p_d := 2p^{d-2}q \in [0,1]$. Proving (31) will clearly prove the theorem.

The proof is by induction on d. The base case is d=2. In this case, since there are no true threshold gates at depth d-1 by assumption, a true threshold gate can only occur as the output gate of the circuit C. Since AND and OR gates are also threshold gates, we can assume that the output gate is a threshold gate. The bottom fan-in being at most D implies that each gate at depth 1 can be represented exactly as a polynomial of degree at most D and therefore that the function computed by C is a degree-D PTF. Hence, Lemma 7.2 trivially implies the result.

Now assume d > 2. Let ψ_1, \ldots, ψ_s denote the AND and OR gates at depth exactly d - 2 in the circuit and let ϕ_1, \ldots, ϕ_m denote the true threshold gates. By assumption $m \leq k$ and $s \leq M$. We sample a random restriction $\rho \sim \mathcal{R}_p^n$ and consider the restricted circuit $C|_{\rho}$.

Håstad's switching lemma [19] tells us that for each $i \in [s]$, we have

$$\Pr_{\rho}[\text{DT-depth}(\psi_i|_{\rho}) \ge D] \le (10pD)^D, \tag{32}$$

and hence by a union bound,

$$\Pr_{o}[\exists i \in [s] : \mathrm{DT\text{-}depth}(\psi_i|_{\rho}) \ge D] \le s(10pD)^D. \tag{33}$$

Also, as in the base case, we see that each ϕ_j computes a degree-D PTF. Hence, Lemma 7.2 gives us

$$\mathbf{E}\left[\sum_{\rho} \operatorname{Var}(\phi_j|_{\rho})\right] \le m\alpha(p, D). \tag{34}$$

Consider the circuit C'_{ρ} obtained from $C|_{\rho}$ as follows: if there is an $i \in [s]$ such that $\mathrm{DT\text{-}depth}(\psi_i|_{\rho}) \geq D$, then C'_{ρ} is defined to be a trivial circuit that always outputs 1; otherwise, C'_{ρ} is the depth-d-1 circuit obtained from $C|_{\rho}$ as follows:

■ We replace each $\phi_j|_{\rho}$ by a bit $b_{j,\rho} \in \{-1,1\}$ so that by Fact 2.6, we have

$$\Pr_{x \in \{-1,1\}^{|\rho^{-1}(*)|}} [\phi_j(x) \neq b_{j,\rho}] \le O(\text{Var}(\phi_j)),$$

- Since each $\psi_i|_{\mathcal{O}}$ is a depth-D decision tree, we can write it as a D-DNF or D-CNF or as a disjoint sum of terms of size at most D each. For each gate χ at depth at most d-3that takes ψ_i as an input, we do the following:
 - If χ is an OR gate, then we take the D-DNF representing $\psi_i|_{\rho}$ and feed the terms of the DNF directly into χ , eliminating the output OR gate of the D-DNF.
 - If χ is an AND gate, we do the same as above, except that we use the D-CNF representation of $\psi_i|_{\rho}$ and eliminate the output AND gate.
 - If χ is a threshold gate, then we write $\psi_i|_{\rho}$ as a disjoint sum of terms of size at most D each and feed each of the terms directly to χ . The gate χ now has many inputs in the place of $\psi_i|_{\rho}$, and the weight given to each of these inputs is the same as the weight given to $\psi_i|_{\rho}$.

Note that the above operations do not increase the number of gates at depth at most d-3 in the circuit.

Note that C'_{ρ} has depth d-1 and bottom fan-in at most D. Further, the number of gates at depth at most d-3 in C'_{ρ} is at most M-s. Hence, C'_{ρ} is a circuit from the class $TAC^{0}[k-m,d-1,D,M]$. We can thus apply the induction hypothesis and obtain

$$\underset{\rho_{d-1}}{\mathbf{E}} \left[\text{Var}(C_{\rho}'|_{\rho_{d-1}}) \right] \le O((k-m)\alpha(p,D) + \alpha(q,D) + (M-s)(10pD)^D). \tag{35}$$

To obtain (31), we use

$$\mathbf{E}[\operatorname{Var}(C)] = \mathbf{E}_{\rho_{d-1}}[\mathbf{E}[\operatorname{Var}(C|\rho)|_{\rho_{d-1}}]] \leq \mathbf{E}_{\rho_{d-1}}[\mathbf{E}[\operatorname{Var}(C'_{\rho})|_{\rho_{d-1}}]] + O\left(\mathbf{E}[\delta(C|\rho, C'_{\rho})]\right) \\
= \mathbf{E}[\mathbf{E}_{\rho}[\operatorname{Var}(C'_{\rho})|_{\rho_{d-1}}]] + O\left(\mathbf{E}[\delta(C|\rho, C'_{\rho})]\right) \quad (36)$$

where the inequality follows from Proposition 3.2. Inequality (35) allows us to bound the first term on the right hand size.

It remains to analyze the last term on the right hand side of (36). Define a Boolean random variable $Z = Z(\rho)$ which is 1 iff there is an $i \in [s]$ such that ϕ_i is not a depth-D decision tree. Let $\Delta = \Delta(\rho)$ be the random variable defined by $\Delta := Z + \sum_{j \in [m]} \operatorname{Var}(\phi_j|_{\rho})$.

It easily follows from the definition of C'_{ρ} that for any choice of ρ , either Z=1 – in which case we can trivially bound $\delta(C'_{\rho}, C|_{\rho})$ by 1 – or $\delta(C'_{\rho}, C|_{\rho}) \leq \sum_{i} \delta(\phi_{i}|_{\rho}, b_{j,\rho}) =$ $\sum_{i} \Pr_{x}[\phi_{j}|_{\rho}(x) \neq b_{j,\rho}]$. Hence, for any choice of ρ , we get

$$\delta(C_{\rho}',C|_{\rho}) \leq Z + \sum_{j \in [m]} \Pr_{x \in \{-1,1\}^{|\rho^{-1}(*)|}} [\phi_{j}|_{\rho}(x) \neq b_{j,\rho}] \leq O(\Delta).$$

Further, by (33) and (34), we have $\mathbf{E}_{\rho}[\Delta] \leq O(m\alpha(p,D) + s(10pD)^D)$. Putting this together with (35) and $(36)^4$ gives the claimed bound. This completes the induction.

Of course, we need to be judicious in our choice of constants in the $O(\cdot)$. We leave this matter to the interested reader.

This yields the following correlation bound as in Corollary 3.3.

- ▶ Corollary 7.4. The following is true for any constant $d \geq 2$. Say C is a depth-d $TAC^0[k]$ circuit with at most M gates where $k \leq \delta \cdot n^{1/2(d-1)}$ and $M = n^{o(\sqrt{\log n/\log\log n})}$. Then $Corr(C, Par_n) \leq n^{o(1)} \cdot \delta^{1-\frac{1}{d}}$. In particular, if $\delta = n^{-\Omega(1)}$, then $Corr(C, Par_n) = n^{-\Omega(1)}$.
- **Proof.** We choose a $D = o(\sqrt{\log n/\log\log n})$ so that $M \le n^{o(D)}$ and p,q as in Corollary 3.3. We can then use Theorem 7.3 to obtain $\mathrm{NS}_{1/n}(C) \le n^{o(1)} \cdot \delta^{1-\frac{1}{d}} + M \cdot (10pD)^D$. Since the latter term is $\frac{1}{n^{\omega(1)}}$, we get $\mathrm{NS}_{1/n}(C) \le n^{o(1)}\delta^{1-\frac{1}{d}}$. By Proposition 2.5, we have $\mathrm{Corr}(C, \mathrm{Par}_n) \le O(\mathrm{NS}_{1/n}(C))$, which proves the claim.
- ▶ Remark. The above corollary can be strengthened considerably if a widely believed strengthening of Lemma 7.2 named the Gotsman-Linial conjecture [16] is known to hold. The Gotsman-Linial conjecture is a conjecture about the average sensitivity of low-degree PTFs. We do not recall the exact statement of the conjecture here, and refer the reader to the work of Gopalan and Servedio [15] instead. As noted by [15, Corollary1 13], the Gotsman-Linial conjecture implies that for any p and any degree p PTF, we have p NSp (p) p Plugging in this bound in place of Lemma 45, it is not hard to see that we can obtain p Correction p PTF, we have p where p PTF, we have p PTF, we have p PTF, we can obtain p Pugging in this bound in place of Lemma 45, it is not hard to see that we can obtain p PTF, we have p

7.1 Learning algorithms for $TAC^0[k]$ circuits

Theorem 7.3 also allows us to obtain an algorithm to learn small $TAC^0[k]$ circuits under the uniform distribution via an observation of Klivans, O'Donnell, and Servedio [24]. We have the following lemma that can be obtained by putting together Fact 9 and Corollary 15 in [24].

▶ Lemma 7.5. Let \mathcal{F} be a class of Boolean functions defined on $\{-1,1\}^n$. Assume that we know that for some $\varepsilon > 0$ and $f \in \mathcal{F}$, there is a $\gamma > 0$ such that $\mathrm{NS}_{\gamma}(f) \leq \varepsilon/3$. Then, there is an algorithm that learns \mathcal{F} with error ε in time $n^{O(1/\gamma)}$.

Using the above lemma and Theorem 7.3, we get subexponential-time (i.e. $2^{o(n)}$ -time) learning algorithms for TAC⁰[k] circuits of small size.

- ▶ Theorem 7.6. Let d be any fixed constant. The class of $TAC^0[k]$ circuits of depth d and size M where $M = n^{o(\sqrt{\log n/\log\log n})}$ and $k = \delta n^{1/2(d-1)}$ for some $\delta > 0$ can be learned to within error $\varepsilon > 0$ in time $n^{O(m)}$ where $m = \max\{n^{1+o(1)}\delta^{2(d-1)}/\varepsilon^{2d}, n^{1/4+o(1)}/\varepsilon^2\}$. In particular, if $\delta = n^{-\Omega(1)}$ and $\varepsilon = \Omega(1)$, then the running time of the algorithm is $2^{o(n)}$.
- **Proof.** We can assume that $\varepsilon \geq 1/n^{1/2d}$ since otherwise, we can just run a brute force algorithm that takes time $2^{O(n)}$. We choose a $D = o(\sqrt{\log n/\log\log n})$ so that $M = n^{o(D)}$. Theorem 7.3 tells us that for any $p, q \geq \frac{1}{n}$ any C from the class of circuits described in the theorem statement, we have

$$\operatorname{NS}_{p^{d-1}q}(C) \leq Ak\sqrt{p} + B\sqrt{q} + O(M(10pD)^D)$$

where A and B are $n^{o(1)}$.

We choose p,q so that the first two terms above are each bounded by $\varepsilon/10$. This requires $p \leq \varepsilon^2/O(k^2A)$ and $q \leq \varepsilon^2/O(B)$. Further, to ensure that the last term is at most $\varepsilon/10$, it suffices to choose $p \leq n^{-\Omega(1)}$ (in fact, this ensures that the third term is $n^{-\omega(1)}$ whereas $\varepsilon \geq n^{-1/2d}$ by assumption). Thus, we fix $p = \min\{\varepsilon^2/O(k^2A), n^{-1/4d}\}$ and $q = \varepsilon^2/O(B)$ so that all the above conditions are satisfied. This gives

$$NS_{\gamma}(C) \leq \varepsilon/3$$

where $\gamma = p^{d-1}q$. Hence, by Lemma 7.5, we obtain the statement of the theorem.

▶ Remark. Assuming the Gotsman Linial conjecture, the above technique yields subexponential time constant-error learning algorithms as long as $M \leq 2^{n^{o(1)}}$ and $\delta = n^{-\Omega(1)}$. To contrast again with the work of Gopalan and Servedio [15], the results of [15] – even assuming the Gotsman Linial conjecture – only yield subexponential time learning algorithms in the setting where $k < \log n$. However, the dependence on the error parameter in [15] is better than the dependence we obtain here (the running time there has a ε^3 in place of the ε^{2d} that we obtain here).

Acknowledgements. Work of the first and second authors supported by the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013) / ERC Grant Agreement No. 615075.

References -

- Scott Aaronson and Avi Wigderson. Algebrization: A new barrier in complexity theory. In Proceedings of the 40th Annual ACM Symposium on Theory of Computing (STOC'08), 2008.
- Miklos Ajtai. Σ_1^1 -formulae on finite structures. Annals of Pure and Applied Logic, 24:1–48, 1983.
- 3 James Aspnes, Richard Beigel, Merrick L. Furst, and Steven Rudich. The expressive power of voting polynomials. *Combinatorica*, 14(2):135–148, 1994. doi:10.1007/BF01215346.
- 4 László Babai, Noam Nisan, and Mario Szegedy. Multiparty protocols, pseudorandom generators for logspace, and time-space trade-offs. *J. Comput. Syst. Sci.*, 45(2):204–232, 1992. doi:10.1016/0022-0000(92)90047-M.
- 5 Theodore Baker, John Gill, and Robert Solovay. Relativizations of the P =? NP question. SIAM Journal on Computing, 4(4):431–442, 1975.
- 6 Richard Beigel. When do extra majority gates help? polylog(N) majority gates are equivalent to one. Computational Complexity, 4:314–324, 1994. doi:10.1007/BF01263420.
- 7 Ruiwen Chen, Valentine Kabanets, Antonina Kolokolova, Ronen Shaltiel, and David Zuckerman. Mining circuit lower bound proofs for meta-algorithms. *Computational Complexity*, 24(2):333–392, 2015. doi:10.1007/s00037-015-0100-0.
- 8 Thomas M. Cover and Joy A. Thomas. Elements of Information Theory. Wiley-Interscience, New York, USA, 2000.
- 9 Ilias Diakonikolas, Parikshit Gopalan, Ragesh Jaiswal, Rocco A. Servedio, and Emanuele Viola. Bounded independence fools halfspaces. SIAM J. Comput., 39(8):3441–3462, 2010. doi:10.1137/100783030.
- 10 Ilias Diakonikolas, Prasad Raghavendra, Rocco A. Servedio, and Li-Yang Tan. Average sensitivity and noise sensitivity of polynomial threshold functions. *SIAM J. Comput.*, 43(1):231–253, 2014. doi:10.1137/110855223.

- Devdatt P. Dubhashi and Alessandro Panconesi. Concentration of Measure for the Analysis of Randomized Algorithms. Cambridge University Press, 2009. URL: http://www.cambridge.org/gb/knowledge/isbn/item2327542/.
- W. Feller. An Introduction to Probability Theory and its Applications. John Wiley & Sons, New York, 3 edition, 1968.
- Merrick Furst, James Saxe, and Michael Sipser. Parity, circuits, and the polynomial-time hierarchy. *Mathematical Systems Theory*, 17(1):13–27, April 1984.
- Dmitry Gavinsky, Shachar Lovett, Michael E. Saks, and Srikanth Srinivasan. A tail bound for read-k families of functions. *Random Struct. Algorithms*, 47(1):99–108, 2015. doi: 10.1002/rsa.20532.
- Parikshit Gopalan and Rocco A. Servedio. Learning and lower bounds for ac⁰ with threshold gates. *Electronic Colloquium on Computational Complexity (ECCC)*, 17:74, 2010. URL: http://eccc.hpi-web.de/report/2010/074.
- Craig Gotsman and Nathan Linial. Spectral properties of threshold functions. *Combinatorica*, 14(1):35–50, 1994. doi:10.1007/BF01305949.
- 17 Kristoffer Arnsfelt Hansen and Peter Bro Miltersen. Some meet-in-the-middle circuit lower bounds. In *Mathematical Foundations of Computer Science 2004, 29th International Symposium, MFCS 2004, Prague, Czech Republic, August 22-27, 2004, Proceedings*, pages 334–345, 2004. doi:10.1007/978-3-540-28629-5_24.
- 18 Prahladh Harsha, Adam Klivans, and Raghu Meka. Bounding the sensitivity of polynomial threshold functions. *Theory of Computing*, 10:1-26, 2014. URL: http://theoryofcomputing.org/articles/v010a001/.
- Johan Håstad. Almost optimal lower bounds for small depth circuits. In Proceedings of the 18th Annual ACM Symposium on Theory of Computing, May 28-30, 1986, Berkeley, California, USA, pages 6-20, 1986. doi:10.1145/12130.12132.
- 20 Russell Impagliazzo, William Matthews, and Ramamohan Paturi. A satisfiability algorithm for ACO. In *Proceedings of 23rd Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 961–972, 2012.
- 21 Russell Impagliazzo, Mohan Paturi, and Stefan Schneider. A satisfiability algorithm for sparse depth two threshold circuits. In *Proceedings of 54th Annual IEEE Symposium on Foundations of Computer Science*, pages 479–488, 2013.
- 22 Russell Impagliazzo, Ramamohan Paturi, and Michael E. Saks. Size-depth trade-offs for threshold circuits. SIAM J. Comput., 26(3):693-707, 1997. doi:10.1137/S0097539792282965.
- Daniel M. Kane. The correct exponent for the gotsman-linial conjecture. *Computational Complexity*, 23(2):151–175, 2014. doi:10.1007/s00037-014-0086-z.
- Adam R. Klivans, Ryan O'Donnell, and Rocco A. Servedio. Learning intersections and thresholds of halfspaces. *J. Comput. Syst. Sci.*, 68(4):808-840, 2004. doi:10.1016/j.jcss.2003.11.002.
- 25 Ilan Komargodski and Ran Raz. Average-case lower bounds for formula size. In *Symposium on Theory of Computing Conference, STOC'13, Palo Alto, CA, USA, June 1-4, 2013*, pages 171–180, 2013. doi:10.1145/2488608.2488630.
- Nathan Linial, Yishay Mansour, and Noam Nisan. Constant depth circuits, Fourier transform and learnability. *Journal of the ACM*, 40(3):607–620, 1993.
- Shachar Lovett and Srikanth Srinivasan. Correlation bounds for poly-size ${\rm AC^0}$ circuits with $n^{1-o(1)}$ symmetric gates. In Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques 14th International Workshop, APPROX 2011, and 15th International Workshop, RANDOM 2011, Princeton, NJ, USA, August 17-19, 2011. Proceedings, pages 640–651, 2011. doi:10.1007/978-3-642-22935-0_54.

- 28 Raghu Meka and David Zuckerman. Pseudorandom generators for polynomial threshold functions. SIAM J. Comput., 42(3):1275–1301, 2013. doi:10.1137/100811623.
- 29 Noam Nisan. The communication compelxity of threshold gates. Combinatorics: Paul Erdös is Eighty, Bolyai Society Mathematical Studies, pages 301–315, 1993.
- 30 Noam Nisan and Avi Wigderson. Hardness vs randomness. *Journal of Computer and System Sciences*, 49(2):149–167, 1994.
- 31 Ryan O'Donnell. Hardness amplification within np. *J. Comput. Syst. Sci.*, 69(1):68-94, 2004. doi:10.1016/j.jcss.2004.01.001.
- 32 Ryan O'Donnell. Analysis of Boolean Functions. Cambridge University Press, 2014. URL: http://www.cambridge.org/de/academic/subjects/computer-science/algorithmics-complexity-computer-algebra-and-computational-g/analysis-boolean-functions.
- 33 Ryan O'Donnell and Rocco A. Servedio. The chow parameters problem. *SIAM J. Comput.*, 40(1):165–199, 2011. doi:10.1137/090756466.
- Ramamohan Paturi and Michael E. Saks. Approximating threshold circuits by rational functions. *Inf. Comput.*, 112(2):257–272, 1994. doi:10.1006/inco.1994.1059.
- Yuval Peres. Noise stability of weighted majority, December 19 2004. Comment: six pages. URL: http://arxiv.org/abs/math/0412377.
- Vladimir V. Podolskii. Exponential lower bound for bounded depth circuits with few threshold gates. Inf. Process. Lett., 112(7):267-271, 2012. doi:10.1016/j.ipl.2011.12. 011.
- 37 Anup Rao. Extractors for low-weight affine sources. In *Proceedings of the 24th Annual IEEE Conference on Computational Complexity, CCC 2009, Paris, France, 15-18 July 2009*, pages 95–101, 2009. doi:10.1109/CCC.2009.36.
- 38 Alexander Razborov. Lower bounds on the size of bounded-depth networks over the complete basis with logical addition. *Mathematical Notes of the Academy of Sciences of the USSR*, 41(4):333–338, 1987.
- 39 Alexander Razborov and Steven Rudich. Natural proofs. *Journal of Computer and System Sciences*, 55(1):24–35, 1997.
- 40 Alexander A. Razborov and Avi Wigderson. n^omega(log n) lower bounds on the size of depth-3 threshold circuits with AND gates at the bottom. *Inf. Process. Lett.*, 45(6):303–307, 1993. doi:10.1016/0020-0190(93)90041-7.
- 41 Michael E. Saks. Slicing the hypercube. Surveys in Combinatorics, 1993, pages 211-255, 1993. URL: http://dl.acm.org/citation.cfm?id=164558.164579.
- 42 Rahul Santhanam. Fighting perebor: New and improved algorithms for formula and QBF satisfiability. In *Proceedings of 51st Annual IEEE Symposium on Foundations of Computer Science*, pages 183–192, 2010.
- Rocco A. Servedio. Every linear threshold function has a low-weight approximator. *Computational Complexity*, 16(2):180–209, 2007. doi:10.1007/s00037-007-0228-7.
- 44 Alexander A. Sherstov. Optimal bounds for sign-representing the intersection of two halfspaces by polynomials. *Combinatorica*, 33(1):73–96, 2013. doi:10.1007/s00493-013-2759-7.
- Kai-Yeung Siu, Vwani P. Roychowdhury, and Thomas Kailath. Rational approximation techniques for analysis of neural networks. *IEEE Transactions on Information Theory*, 40(2):455–466, 1994. doi:10.1109/18.312168.
- 46 Roman Smolensky. Algebraic methods in the theory of lower bounds for boolean circuit complexity. In *Proceedings of the 19th Annual Symposium on Theory of Computing*, pages 77–82, 1987.
- 47 Ryan Williams. Non-uniform ACC circuit lower bounds. In *Proceedings of 26th Annual IEEE Conference on Computational Complexity*, pages 115–125, 2011.

- Ryan Williams. New algorithms and lower bounds for circuits with linear threshold gates. In Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 June 03, 2014, pages 194–202, 2014. doi:10.1145/2591796.2591858.
- 49 Andrew Yao. Separating the polynomial-time hierarchy by oracles. In *Proceedings of 26th Annual IEEE Symposium on Foundations of Computer Science*, pages 1–10, 1985.

A Proof of Proposition 2.5

Proof. For point 1., we know that $NS_p(f) = Pr_{(x,y)}[f(x) \neq f(y)]$ where x and y are sampled as in Definition 2.4. Alternately, we may also think of sampling (x,y) in the following way: choose $\rho = (I,z) \sim \mathcal{R}_{2p}^n$ and for the locations indexed by I we choose $x', y' \in \{-1,1\}^{|I|}$ independently and uniformly at random to define strings x and y respectively. Hence, we have

$$NS_{p}(f) = \Pr_{x,y}[f(x) \neq f(y)] = \mathbf{E}[\Pr_{\rho}[f(x') \neq f|_{\rho}(x') \neq f|_{\rho}(y')]] = \mathbf{E}[\frac{1}{2}Var(f|_{\rho})].$$

We now proceed with point 2.. As $NS_p(f)$ is a decreasing function of p [32], we may assume that $p = \frac{1}{n} \leq \frac{1}{2}$ and hence we have $NS_{1/n}(f) = \frac{1}{2} \mathbf{E}_{\rho \sim \mathcal{R}_{2/n}^n}[Var(f|_{\rho})]$. Note that for $\rho = (I, y)$ chosen as above, the probability that $I \neq \emptyset$ is $\Omega(1)$. Hence we have

$$\operatorname{NS}_{1/n}'(f) := \frac{1}{2} \mathop{\mathbf{E}}_{\rho \sim \mathcal{R}_{2/n}^n} [\operatorname{Var}(f|_{\rho}) \mid I \neq \emptyset] \le \frac{\operatorname{NS}_{1/n}(f)}{\operatorname{Pr}_I[I \neq \emptyset]} = O(\operatorname{NS}_{1/n}(f)).$$

Further, note that for any $m \geq 1$ and any Boolean function $g : \{-1,1\}^m \to \{-1,1\}$, its distance from either the constant function 1 or the constant function -1 is at most Var(g)/2. Since Par_m has correlation 0 with any constant function, using Fact 2.3, we have $Corr(Par_m, g) \leq Var(g)/2$.

Using Fact 2.3 again, we get

$$\operatorname{Corr}(\operatorname{Par}_{n}, f) \leq \underset{\rho \sim \mathcal{R}_{2/n}^{n}}{\mathbf{E}} \left[\operatorname{Corr}(\operatorname{Par}_{n}|_{\rho}, f|_{\rho}) \mid I \neq \emptyset \right] = \underset{\rho \sim \mathcal{R}_{2/n}^{n}}{\mathbf{E}} \left[\operatorname{Corr}(\operatorname{Par}_{|I|}, f|_{\rho}) \mid I \neq \emptyset \right]$$

$$\leq \underset{\rho \sim \mathcal{R}_{2/n}^{n}}{\mathbf{E}} \left[\frac{1}{2} \operatorname{Var}(g) \mid I \neq \emptyset \right] = \operatorname{NS}'_{1/n}(f) = O(\operatorname{NS}_{1/n}(f)).$$

B Correlation bounds for depth-2 threshold circuits

In this section, we prove near optimal correlation bounds for depth-2 threshold circuits computing Parity.

▶ Theorem B.1 (Main). Fix any constant $\varepsilon < \frac{1}{2}$. Let $\gamma = \frac{1}{2} - \varepsilon$. Any depth-2 threshold circuit on n variables with at most $n^{1+\varepsilon}$ wires has correlation at most $n^{-\Omega(\gamma)}$ with the parity function on n variables.

Note that the above theorem is tight, since by Theorem 3.4, there is a depth-2 circuit with $O(\sqrt{n})$ gates (and hence $O(n^{3/2})$ wires) that computes Parity correctly with high probability. The proof is based on the following two subclaims:

▶ **Theorem B.2** (Aspnes, Beigel, Furst, and Rudich [3]). Any degree-t polynomial threshold function (PTF) has correlation at most $O(t/\sqrt{m})$ with the parity function on m variables.

We say that a circuit C is δ -approximated by a circuit C' if $\Pr_x[C(x) \neq C'(x)] \leq \delta$.

▶ Claim B.3. Let ε, γ be as in the statement of Theorem B.1 and let α denote $\gamma/3$. Say C denotes a depth-2 threshold circuit of wire complexity $n^{1+\varepsilon}$ and let f_1, \ldots, f_t be the LTFs computed by C at depth-1. Under a random restriction ρ with *-probability $p = \frac{1}{n^{1-\alpha}}$, with probability at least $1 - n^{-\Omega(\gamma)}$, the circuit $C|_{\rho}$ is $n^{-\Omega(\gamma)}$ -approximated by a circuit \tilde{C}_{ρ} which is obtained from C by replacing each of the $f_i|_{\rho}s$ by an $O(n^{\alpha/2-\Omega(\gamma)})$ -junta g_i .

Assuming the above two claims, we can finish the proof of Theorem B.1 easily as follows. Let C be a circuit of wire complexity $n^{1+\varepsilon}$. We apply a random restriction ρ with *-probability $p=\frac{1}{n^{1-\alpha}}$ as in Claim B.3. Call the restriction good if there is a circuit \tilde{C}_{ρ} as in the Claim that $n^{-\Omega(\gamma)}$ -approximates $C|_{\rho}$ and bad otherwise. The probability that we have a bad restriction is at most $n^{-\Omega(\gamma)}$.

Say ρ is a good restriction. The circuit \tilde{C}_{ρ} can be represented by an $O(n^{\alpha/2-\Omega(\gamma)})$ -degree PTF and hence by Theorem B.2 has correlation at most $n^{-\Omega(\gamma)}$ with parity (on the remaining n^{α} variables). Moreover, then $C|_{\rho}$ is well-approximated by \tilde{C}_{ρ} and hence has correlation at most $n^{-\Omega(\gamma)} + n^{-\Omega(\gamma)}$ with parity.

Upper bounding the correlation by 1 for bad restrictions, we see that the overall correlation is at most $n^{-\Omega(\gamma)}$.

We now prove Claim B.3.

Proof of Claim B.3. Let f_1, \ldots, f_t be the LTFs appearing at depth 1 in the circuit. We will divide the analysis based on the fan-ins of the f_i s (i.e. the number of variables they depend on).

We denote by β the quantity $\frac{3}{4} + \frac{\varepsilon}{2}$. It can be checked that we have both

$$\beta = \frac{1}{2} + \varepsilon + \frac{\alpha}{2} + \Omega(\gamma) \text{ and } 1 - \beta = \frac{\alpha}{2} + \Omega(\gamma).$$
 (37)

Consider any f_i of fan-in at most n^{β} . When hit with a random restriction with *-probability $n^{-(1-\alpha)}$, we see that the expected number of variables of f_i that survive is at most $n^{\beta-(1-\alpha)} = n^{\alpha-(1-\beta)} = n^{\alpha/2-\Omega(\gamma)}$ by (37) above. By a Chernoff bound, the probability that this number exceeds twice its expectation is exponentially small. Union bounding over all the gates of small fan-in, we see that with probability $1 - \exp(-n^{\Omega(1)})$, all the low fan-in gates depend on at most $2n^{\alpha/2-\Omega(\gamma)}$ many variables after the restriction. We call this high probability event \mathcal{E}_1 .

Now, we consider the gates of fan-in at least n^{β} . W.l.o.g., let f_1, \ldots, f_r be these LTFs. Since the total number of wires is at most $n^{1+\varepsilon}$, we have $r \leq n^{1+\varepsilon-\beta} = n^{\frac{1}{2}-\frac{\alpha}{2}-\Omega(\gamma)}$ by (37).

By Theorem 2.11, we know that for any
$$f_i$$
,

$$\mathop{\mathbf{E}}_{\rho}[\operatorname{Var}(f_i|_{\rho})] \le O(\sqrt{p}) = O(\frac{1}{n^{(1-\alpha)/2}}).$$

By linearity of expectation, we have

$$E := \mathop{\mathbf{E}}_{\rho}[\sum_{i=1}^{r} \operatorname{Var}(f_{i}|_{\rho})] \le O(r \cdot \frac{1}{n^{(1-\alpha)/2}}) = O(n^{(1-\alpha)/2 - \Omega(\gamma)} \cdot \frac{1}{n^{(1-\alpha)/2}}) = O(n^{-\Omega(\gamma)}).$$

By Markov's inequality, we see that the probability that $\sum_{i=1}^r \operatorname{Var}(f_i|_{\rho}) > \sqrt{E}$ is at most $\sqrt{E} = n^{-\Omega(\gamma)}$. We let \mathcal{E}_2 denote the event that $\sum_{i=1}^r \operatorname{Var}(f_i|_{\rho}) \leq \sqrt{E}$.

Consider the event $\mathcal{E} = \mathcal{E}_1 \wedge \mathcal{E}_2$. A union bound tells us that the probability of \mathcal{E} is at least $1 - n^{-\Omega(\gamma)}$. When this event occurs, we construct the circuit \tilde{C}_{ρ} from the statement of the claim as follows.

When the event \mathcal{E} occurs, the LTFs of low arity are already $n^{\alpha/2-\Omega(\gamma)}$ -juntas, so there is nothing to be done for them.

Now, consider the LTFs of high fan-in, which are f_1, \ldots, f_r . For each $f_i|_{\rho}$ $(i \in [r])$, replace f_i by a bit $b_i \in \{-1,1\}$ such that $\Pr_x[f_i|_{\rho}(x) = b_i] \geq \frac{1}{2}$. In the circuit \tilde{C}_{ρ} , these gates thus become constants, which are 0-juntas. The circuit \tilde{C}_{ρ} now has the required form. We now analyze the error introduced by this operation.

We know that $\Pr_x[f_i|_{\rho}(x) \neq b_i] \leq 2\text{Var}(f_i|_{\rho})$ and thus the overall error introduced is at most $2\sum_{i\in[r]} \text{Var}(f_i|_{\rho}) \leq O(\sqrt{E}) = n^{-\Omega(\gamma)}$ (since \mathcal{E}_2 is assumed to occur). Thus, the circuit \tilde{C}_{ρ} is an $n^{-\Omega(\gamma)}$ -approximation to C.

C Proof of Proposition 4.12

Proof of Proposition 4.12. Let J be the set of variables being set and let $y \in \{-1,1\}^{|J|}$ denote the random assignment chosen. Let $L_0 = \frac{1}{\varepsilon^2} \cdot 3\log(1/\varepsilon)$. It can be checked that for any $i < L' - L_0$, we have

$$||w_{>(i+L_0)}||_2^2 \le \frac{\varepsilon^2}{9} \cdot ||w_{>i}||^2 \le \frac{w_i^2}{9}.$$

Hence, we can choose indices $i_1 = 1, i_2 = 1 + L_0, \dots, i_{r+2} = 1 + (r+1)L_0 \le L'$ such that $|w_{i_{j+1}}| \le \frac{|w_{i_j}|}{3}$ and $||w_{i_{j+1}}||_2^2 \le \frac{\varepsilon^2}{9} \cdot ||w_{i_j}||_2^2$. Further, we have

$$\sum_{i \not\in J} w_i^2 \le \|w_{>L'}\|^2 \le \|w_{>i_{r+2}}\|^2 \le \frac{\varepsilon^2}{9} \cdot \|w_{>i_{r+1}}\|_2^2 \le \frac{\varepsilon^2}{81} \cdot w_{i_r}^2.$$

We condition on any setting of variables other than y_{i_1}, \ldots, y_{i_r} . This means that the constant term of the restricted threshold gate θ' is given by

$$\theta' = \theta'' - \sum_{j \in [r]} w_{i_j} y_{i_j}$$

for some $\theta'' \in \mathbb{R}$. The probability that the threshold gate is not $\frac{1}{\epsilon}$ -imbalanced is at most

$$\begin{split} \Pr_{y_{i_1},\dots,y_{i_r}}[|\theta'| &\leq \frac{1}{\varepsilon^2} \cdot \sqrt{\sum_{i \notin J} w_i^2}] \\ &\leq \Pr_{y_{i_1},\dots,y_{i_r}}[|\theta'| \leq \frac{1}{9} \cdot |w_{i_r}|] \\ &= \Pr_{y_{i_1},\dots,y_{i_r}}[\sum_j w_{i_j} y_{i_j} \in [\theta'' - \frac{1}{9} \cdot |w_{i_r}|, \theta'' + \frac{1}{9} \cdot |w_{i_r}|]] \end{split}$$

Now, as a result of the exponentially decreasing nature of the $|w_{i_j}|$, it follows that for any interval of length at most $|w_{i_r}|/2$, there can be at most one choice of y_{i_1}, \ldots, y_{i_r} such that the $\sum_j w_{i_j} y_{i_j}$ lies in that interval. Thus, we have the given bound.