

Fourier Growth of Parity Decision Trees

Uma Girish* Avishay Tal† Kewen Wu‡

Abstract

We prove that for every parity decision tree of depth d on n variables, the sum of absolute values of Fourier coefficients at level ℓ is at most $d^{\ell/2} \cdot O(\ell \cdot \log(n))^\ell$. Our result is nearly tight for small values of ℓ and extends a previous Fourier bound for standard decision trees by Sherstov, Storozhenko, and Wu (STOC, 2021).

As an application of our Fourier bounds, using the results of Bansal and Sinha (STOC, 2021), we show that the k -fold Forrelation problem has (randomized) parity decision tree complexity $\tilde{\Omega}(n^{1-1/k})$, while having quantum query complexity $\lceil k/2 \rceil$.

Our proof follows a random-walk approach, analyzing the contribution of a random path in the decision tree to the level- ℓ Fourier expression. To carry the argument, we apply a careful cleanup procedure to the parity decision tree, ensuring that the value of the random walk is bounded with high probability. We observe that step sizes for the level- ℓ walks can be computed by the intermediate values of level $\leq \ell - 1$ walks, which calls for an inductive argument. Our approach differs from previous proofs of Tal (FOCS, 2020) and Sherstov, Storozhenko, and Wu (STOC, 2021) that relied on decompositions of the tree. In particular, for the special case of standard decision trees we view our proof as slightly simpler and more intuitive.

In addition, we prove a similar bound for noisy decision trees of cost at most d – a model that was recently introduced by Ben-David and Blais (FOCS, 2020).

*Department of Computer Science, Princeton University. Email: ugirish@cs.princeton.edu

†Department of EECS, University of California at Berkeley. Email: avishay.tal@gmail.com

‡Department of EECS, University of California at Berkeley. Email: shlw.kevin@hotmail.com

1 Introduction

A common theme in the analysis of Boolean functions is proving structural results on classes of Boolean devices (e.g., decision trees, bounded-depth circuits) and then exploiting the structure to: (i) devise pseudorandom generators fooling these devices, (ii) prove lower bounds, showing that some explicit function cannot be computed by such Boolean devices of certain size, or (iii) design learning algorithms for the class of Boolean devices in either the membership-query model or the random-samples model. Such structural results can involve properties of the Fourier spectrum of Boolean functions associated with Boolean devices, like concentration on low-degree terms or concentration on a few terms (i.e., “approximate sparsity”).

In this work, we investigate the Fourier spectrum of parity decision trees. A parity decision tree (PDT) is an extension of the standard decision tree model. A PDT is a binary tree where each internal node is marked by a linear function (modulo 2) on the input variables (x_1, \dots, x_n) , with two outgoing edges marked with 0 and 1, and each leaf is marked with either 0 or 1. A PDT naturally describes a computational model: on input $x = (x_1, \dots, x_n)$, start at the root and at each step query the linear function specified by the current node on the input x and continue on the edge marked with the value of the linear function evaluated on x . Finally, when reaching a leaf, output the value specified in the leaf. PDTs naturally generalize standard decision trees that can only query the value of a single input bit in each internal node.

PDTs were introduced in the seminal paper of Kushilevitz and Mansour [KM93]. Aligned with the aforementioned theme, Kushilevitz and Mansour proved a structural result for PDTs and used it to design learning algorithms for PDTs. They showed that every PDT of size s computing a Boolean function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ has

$$L_1(f) \triangleq \sum_{S \subseteq [n]} |\hat{f}(S)| \leq s,$$

where $\hat{f}(S)$ are the Fourier coefficients of f (see Subsection 2.1 for a precise definition). Then, they gave a learning algorithm in the membership-query model, running in time $\text{poly}(t, n)$ that can learn any function f with $L_1(f) \leq t$. Combining the two results together, they obtained a $\text{poly}(s, n)$ -time algorithm for learning PDTs of size s .

Parity decision trees were also studied in relation to communication complexity and the log-rank conjecture [MO09, ZS09, ZS10, TWXZ13, STIV17, OWZ⁺14, CS16, KQS15, HHL18, San19, MS20]. Suppose Alice gets input $x \in \{0, 1\}^n$, Bob gets input $y \in \{0, 1\}^n$ and they want to compute some function $f(x, y)$. When f is an XOR-function, namely $f(x, y) = g(x \oplus y)$ for some $g: \{0, 1\}^n \rightarrow \{0, 1\}$, then any PDT for g of depth d can be translated into a communication protocol for f at cost $2d$: Alice and Bob simply traverse the PDT together, both exchanging the parity of their part of the input to simulate each query in the PDT. With this view, parity decision trees can be thought of as special cases of communication protocols for XOR functions. A surprising result by Hatami, Hosseini, and Lovett [HHL18], shows that this is not far from the optimal strategy for XOR functions. Namely, if the communication cost for computing f is c , then the parity decision tree complexity of g is at most $\text{poly}(c)$. Due to this connection, the log-rank conjecture for XOR-functions reduces to the question of whether Boolean functions with at most s non-zero Fourier coefficients can be computed by PDTs of depth $\text{polylog}(s)$ [MO09, ZS09]. The best known upper bound is that such functions can be computed by PDTs of depth $O(\sqrt{s})$ [TWXZ13] (or even non-adaptive PDTs of depth $\tilde{O}(\sqrt{s})$ [San19]).

While having small $L_1(f)$ norm implies learning algorithms and also simple pseudorandom generators fooling f [NN93], this property can be quite restrictive. In particular, very simple functions (e.g., the Tribes function) have $L_1(f)$ exponential in n . Such examples motivated Reingold,

Steinke, and Vadhan [RSV13] to study a more refined notion measuring for a given level ℓ , the sum of absolute values of Fourier coefficients of sets S of size exactly ℓ , i.e, to study

$$L_{1,\ell}(f) \triangleq \sum_{S \subseteq [n]: |S|=\ell} \left| \widehat{f}(S) \right|.$$

In particular, for $\ell = 1$, the measure $L_{1,1}(f)$ is tightly related to the total influence of f (and equals to it if f is monotone). The idea behind this more refined notion is that Fourier coefficients of different levels behave differently under standard manipulations to the function like random restrictions or noise operators. For example, when applying a noise operator with parameter γ , level- ℓ coefficients are multiplied by γ^ℓ . This motivates to establish a bound of the form $L_{1,\ell}(f) \leq t^\ell$ for some parameter t and all $\ell = 1, \dots, n$. If f satisfies such a bound, we say that $f \in \mathcal{L}_1(t)$.¹

Reingold, Steinke, and Vadhan [RSV13] showed that for read-once permutation branching programs of width w , while $L_1(f)$ could be exponential in n (even for $w = 3$), it nevertheless holds that $L_{1,\ell}(f) \leq (2w^2)^\ell$ for all $\ell = 1, \dots, n$. Then, they constructed a pseudorandom generator that fools any class of read-once branching programs for which $f \in \mathcal{L}_1(t)$ using only $t \cdot \text{polylog}(n)$ random bits. This result was significantly generalized to a pseudorandom generator that fools any class of functions $f \in \mathcal{L}_1(t)$ using only $t^2 \cdot \text{polylog}(n)$ random bits [CHHL19]. Further results established pseudorandom generators assuming $L_{1,\ell}$ bounds only on the first few levels [CHRT18, CGL⁺20].

It turns out that read-once permutation branching programs are just one example of many well-studied Boolean devices with non-trivial $L_{1,\ell}$ bounds. The following classes of Boolean functions are other examples:

1. Width- w CNF and width- w DNF formulae are in $\mathcal{L}_1(O(w))$ [Man95].
2. AC^0 circuits of size s and depth d are in $\mathcal{L}_1(O(\log(s))^{d-1})$ [Tal17].
3. Boolean functions with max-sensitivity at most s are in $\mathcal{L}_1(O(s))$ [GSTW16]
4. Read-once branching programs of width w are in $\mathcal{L}_1(O(\log(n))^w)$ [CHRT18]
5. Deterministic and randomized decision trees of depth d are in $\mathcal{L}_1\left(O\left(\sqrt{d \log(n)}\right)\right)$ [Tal20, SSW20].
6. If $f(x, y)$ is a function computed by communication protocol exchanging at most c bits, then $h(z) = \mathbb{E}_x[f(x, x \oplus z)]$ satisfies $h \in \mathcal{L}_1(O(c))$ [GRT21, GRZ20].
7. Polynomials f over $\text{GF}(2)$ of degree d have $L_{1,\ell}(f) \leq (2^{3d} \cdot \ell)^\ell$ [CHHL19].
8. Product tests, i.e., the XOR of multiple Boolean functions operating on disjoint sets of at most m bits each, are in $\mathcal{L}_1(O(m))$ [Lee19].

We remark that Items 1, 2, 4, 5 and 8 are essentially tight, Item 3 can be potentially improved polynomially [OD12, OS07], Item 6 can be potentially improved quadratically [GRT21] and Item 7 can be potentially improved exponentially [CHLT19]. Indeed, improving Item 7 exponentially would imply that $\text{AC}^0[\oplus]$ in $\mathcal{L}_1(\text{polylog}(n))$ and would give the first poly-logarithmic pseudorandom generators for this well-studied class of Boolean circuits [CHLT19].

The most relevant result to our work is the recent tight bounds on the $L_{1,\ell}$ of decision trees of depth d . Sherstov, Storozhenko and Wu [SSW20] recently proved that for any randomized decision tree of depth d computing a function f , it holds that $L_{1,\ell}(f) \leq \sqrt{\binom{d}{\ell}} \cdot O(\log(n))^{\ell-1}$. Their

¹Note that if $f \in \mathcal{L}_1(t)$ then after applying noise operator with $\gamma = 1/(2t)$, the noisy-version of f has total L_1 -norm at most $O(1)$ which makes it is quite easy to fool using small-biased distributions [NN93].

bound is nearly tight (see [Tal20, Section 7] and [O’D14, Chapter 5.3] for tightness examples). One motivation for showing such a bound for decision trees is that it demonstrates a stark difference between quantum algorithms making few queries and randomized algorithms making a few queries. Indeed, the Fourier spectrum associated with quantum query algorithms making a few queries can be far from being approximately sparse (in the sense that its $L_{1,\ell}$ is quite large). Based on that difference, both [SSW20] and [BS20] showed that there are partial functions, either k -fold Forrelation or k -fold Rorrelation, that can be correctly computed with probability at least $1/2 + \Omega(1)$ by quantum algorithms making $\lceil k/2 \rceil$ queries, but require $\tilde{\Omega}(n^{1-1/k})$ queries for any randomized algorithm. Moreover, due to the result of Aaronson and Ambainis [AA18] this is the largest possible separation between the two models.

Indeed, as suggested in [Tal20], one can show that any function with sufficiently good bounds on its $L_{1,\ell}$, for all $\ell = 1, \dots, n$, cannot solve the k -fold Rorrelation, and such bounds were obtained by [SSW20] for randomized decision trees of depth $n^{1-1/k}/\text{polylog}(n)$. Independently, Bansal and Sinha obtained the same separation but only relying on the $L_{1,\ell}$ bounds for $\ell \in \{k, k+1, \dots, k^2\}$. With this additional flexibility, they were able to obtain their separation for the simpler and explicit function called k -fold Forrelation.

For parity decision trees, the work of Blais, Tan, and Wan [BTW15] established a tight bound of $O(\sqrt{d})$ on the first level $\ell = 1$. To the best of our knowledge, bounds on higher levels were not considered previously in the literature (in fact, even for standard decision trees, such bounds were not considered prior to [Tal20]).

1.1 Our Results

We prove level- ℓ bounds for any parity decision tree of depth d .

Theorem 1.1 (Informal). *Let \mathcal{T} be a depth- d parity decision tree on n variables. Then the sum of absolute Fourier coefficients at level ℓ is bounded by $d^{\ell/2} \cdot O(\ell \cdot \log(n))^\ell$.*

See Theorem 5.5 and Theorem 5.12 for a precise statement taking into account the probability that \mathcal{T} accepts a uniformly random input. Theorem 1.1 extends the result of [SSW20] from standard decision trees to parity decision trees at the cost of an $(\ell \cdot \log(n))^{O(\ell)}$ multiplicative factor. We remark that even for standard decision tree there is a lower bound of $L_{1,\ell}(f) \geq \sqrt{\binom{d}{\ell}} \cdot (\log(n))^{\ell-1}$ [Tal20, Section 7] for constant ℓ and $L_{1,\ell}(f) \geq \frac{1}{\text{poly}(\ell)} \cdot \sqrt{\binom{d}{\ell}}$ for all ℓ [O’D14, Chapter 5.3]. Thus, our bounds are tight up to $\text{polylog}(n)$ factors for constant ℓ , and they deteriorate as ℓ grows. Nevertheless, our main application relies on the bounds for small values of ℓ (constant or at most $\log^2 n$).

Noisy Decision Trees. We also investigate the Fourier spectrum of noisy decision trees. Noisy decision trees are a different generalization of the standard model; here in each internal node v we query a noisy version of an input bit, that equals the true bit with probability $(1 + \gamma_v)/2$. Any such query costs γ_v^2 . We say that a noisy decision tree has cost at most d if the total cost in any root-to-leaf path is at most d . Recent work studied this model and established connections to the question of how randomized decision tree complexity behaves under composition [BB20].

We prove level- ℓ bounds for any noisy decision tree of cost at most d . See Theorem 6.3 for a precise statement.

Theorem 1.2 (Informal). *Let \mathcal{T} be a noisy decision tree of cost at most d on n variables. Then the sum of absolute Fourier coefficients at level ℓ is bounded by $O(d)^{\ell/2} \cdot (\ell \cdot \log(n))^{(\ell-1)/2}$.*

Extension to Randomized Query Models. It is simple to verify that if f is a convex combination of Boolean functions f_1, \dots, f_m each with $L_{1,\ell}(f_i) \leq t_\ell$ then also f satisfy $L_{1,\ell}(f) \leq t_\ell$. Thus, if we take a distribution over PDTs of depth d (resp., noisy decision trees of cost d) we get the same bounds on their $L_{1,\ell}$ as those in [Theorem 1.1](#) (resp., [Theorem 1.2](#)). This is captured in the following corollary.

Corollary 1.3. *Let \mathcal{T} be a randomized parity decision tree of depth at most d on n variables. Then,*

$$\forall \ell \in [n] : L_{1,\ell}(\mathcal{T}) \leq d^{\ell/2} \cdot O(\ell \cdot \log(n))^\ell.$$

Let \mathcal{T}' be a randomized noisy decision tree of cost at most d on n variables. Then,

$$\forall \ell \in [n] : L_{1,\ell}(\mathcal{T}') \leq O(d)^{\ell/2} \cdot (\ell \cdot \log(n))^{(\ell-1)/2}.$$

1.2 Applications

Quantum versus Randomized Query Complexity. Let $k \leq \log(n)$. Bansal and Sinha [[BS20](#)] gave a $\lceil k/2 \rceil$ versus $\tilde{\Omega}(n^{1-1/k})$ separation between the quantum and randomized query complexity of k -fold Forrelation (defined by [[AA18](#)]). For our purposes just think of k -fold Forrelation as a partial Boolean function on n input bits. Our main application is an extension of Bansal and Sinha's lower bound for the model of randomized parity decision trees. This follows from their main technical result and [Theorem 1.1](#).

Theorem 1.4 (Restatement of [[BS20](#), Theorem 3.2]). *Let $f: \{0,1\}^n \rightarrow [0,1]$ such that f and all its restrictions satisfy $L_{1,\ell}(f) \leq t^\ell$ for $\ell = \{k, \dots, k(k-1)\}$. Let $\delta = 2^{-5k}$. Suppose f is δ -close to the value of k -fold Forrelation of x for all x on which k -fold Forrelation is defined. Then, $t \geq \Omega\left(\frac{n^{(1-1/k)/2}}{k^{15}}\right)$.*

Corollary 1.5. *If \mathcal{T} is a randomized parity decision tree of depth d computing k -fold Forrelation with success probability $\frac{1}{2} + \gamma$, then $d \geq \gamma^2 \cdot \frac{n^{1-1/k}}{\text{poly}(k) \log^2 n}$.*

Proof. We can amplify the success probability of the randomized parity decision tree from $1/2 + \gamma$ to $1 - 2^{-5k}$ by repeating the query algorithm $O(k/\gamma^2)$ times independently and taking majority. This results in a randomized parity decision tree \mathcal{T}' of depth $d' = O(d \cdot k/\gamma^2)$. Now, [Corollary 1.3](#) gives $L_{1,\ell}(\mathcal{T}') \leq (d')^{\ell/2} \cdot O(\ell \cdot \log(n))^\ell$ for all ℓ . In particular, $L_{1,\ell}(\mathcal{T}') \leq t^\ell$ for all $\ell \leq k(k-1)$ where $t = O\left(\sqrt{d'} \cdot k(k-1) \cdot \log(n)\right)$. This is also true for any restriction of \mathcal{T}' , since fixing variables to constants yields another randomized parity decision tree of depth at most d' . Combining the bounds on $L_{1,\ell}(\mathcal{T}')$ for $\ell \in \{k, \dots, k(k-1)\}$ with [Theorem 1.4](#) gives $d' \geq \frac{n^{1-1/k}}{O(k^{34}) \cdot \log^2(n)}$ and thus $d \geq \gamma^2 \cdot \frac{n^{1-1/k}}{O(k^{35}) \cdot \log^2(n)}$. \square

For constant k and $\gamma = 2^{-O(k)}$, we get a $\lceil k/2 \rceil$ versus $\tilde{\Omega}(n^{1-1/k})$ separation between the quantum query complexity and the randomized parity query complexity of k -fold Forrelation.

Similarly, we can obtain the following corollary for noisy decision trees.

Corollary 1.6. *If \mathcal{T} is a randomized noisy decision tree of cost at most d computing k -fold Forrelation with success probability $\frac{1}{2} + \gamma$, then $d \geq \gamma^2 \cdot \frac{n^{1-1/k}}{\text{poly}(k) \log(n)}$.*

Towards Communication Complexity Lower Bounds. We recall an open question from [GRT21].

Conjecture 1.7. *Let $f: \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$ computed by a deterministic communication protocol of cost at most c . Let $h: \{0, 1\}^n \rightarrow [0, 1]$ defined by $h(z) = \mathbb{E}_x[f(x, x \oplus z)]$. Then, $L_{1,2}(h) \leq c \cdot \text{polylog}(n)$.*

We view Theorem 1.1 as a first step towards this conjecture. Indeed, for communication protocols that follow a parity decision tree strategy according to some tree \mathcal{T} , it is simple to verify that $h = \mathcal{T}$ (as functions), and thus $L_{1,2}(h) = L_{1,2}(\mathcal{T}) \leq c \cdot \text{polylog}(n)$. We remark that proving Conjecture 1.7 would improve the lower bounds by [GRT21] on the randomized communication complexity of the XOR-lifted version of 2-fold Forrelation from $\tilde{\Omega}(n^{1/4})$ to $\tilde{\Omega}(n^{1/2})$.

Application to Expander Random Walk. Recently, [CPT20] showed that expander random walks fool symmetric functions and also general functions in $\mathcal{L}_1(t)$. To be more precise, assume $f \in \mathcal{L}_1(t)$. Let G be an expander, with second eigenvalue $\lambda \ll \frac{1}{t^4}$, where half of G 's vertices are labeled by 0 and the rest are labeled by 1. Then the expected value of f on bits sampled by an $(m-1)$ -step random walk on G is approximately the value it would get on a uniformly random string in $\{0, 1\}^m$. Combined with our results, this shows that if f can be computed by low-depth parity decision trees then f can be fooled by the expander random walk.

Fourier Bounds for Small-size Parity Decision Trees. By a simple size-to-depth reduction we obtain Fourier bounds for parity decision trees of bounded size. We defer the simple proof to Appendix A.

Corollary 1.8. *Let \mathcal{T} be a parity decision tree of size at most $s > 1$ on n variables. Then,*

$$\forall \ell \in [n] : L_{1,\ell}(f) \leq (\log(s))^{\ell/2} \cdot O(\ell \cdot \log(n))^{1.5\ell}.$$

1.3 Technical Overview

For the rest of the paper we consider Boolean functions as functions from $\{\pm 1\}^n$ to $\{0, 1\}$. This is for convenience, since most of our calculations become easier under this representation. Observe that under this view, a parity decision tree queries at each internal node the product $\prod_{i \in S} x_i$ for some $S \subseteq [n]$ and goes left/right depending on whether $\prod_{i \in S} x_i = 1$ or -1 .

Let $\ell \in \mathbb{N}_+$. For simplicity of notation, we use $\tilde{O}_\varepsilon(d^m)$ to denote $(d \cdot \text{polylog}(n^\ell/\varepsilon))^m$ for $m, n, d \in \mathbb{N}_+$ and $\varepsilon \in (0, 1/2]$. When we omit the subscript ε , it is understood that $\varepsilon = 1$. As per this notation, we show a bound of $\tilde{O}(d^{\ell/2})$ on the level- ℓ Fourier mass of parity decision trees of depth d . We first describe the proof for standard decision trees and then show how to generalize to parity decision trees.

Standard Decision Trees. Let \mathcal{T} be a decision tree and for simplicity, assume that every leaf is of depth d . Let v_0, \dots, v_d be a random root-to-leaf path in \mathcal{T} and $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(d)} \in \{-1, 0, 1\}^n$ denote the sequence of partial assignments, i.e., for $j \in [n]$ and $i \in \{0, \dots, d\}$, let

$$\mathbf{v}_j^{(i)} = \begin{cases} 1 & \text{if } x_j \text{ is fixed to } 1 \text{ before reaching } v_i, \\ -1 & \text{if } x_j \text{ is fixed to } -1 \text{ before reaching } v_i, \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

For $u \in \mathbb{R}^n$, we use u_S to denote $\prod_{j \in S} u_j$. Let $a_S = \text{sgn}(\widehat{\mathcal{T}}(S))$ for $|S| = \ell$ and 0 otherwise. Note that

$$\sum_{S:|S|=\ell} |\widehat{\mathcal{T}}(S)| = \sum_{S:|S|=\ell} a_S \cdot \widehat{\mathcal{T}}(S) = \sum_{S:|S|=\ell} a_S \cdot \mathbb{E}_{v_d} \left[\mathcal{T}(v_d) \mathbf{v}_S^{(d)} \right] = \mathbb{E}_{v_d} \left[\mathcal{T}(v_d) \cdot \left(\sum_{S:|S|=\ell} a_S \cdot \mathbf{v}_S^{(d)} \right) \right]. \quad (2)$$

Thus, to bound $\sum_{S:|S|=\ell} |\widehat{\mathcal{T}}(S)|$ it suffices to show that $\left| \sum_{S:|S|=\ell} a_S \cdot \mathbf{v}_S^{(d)} \right|$ is bounded by $\widetilde{O}(d^{\ell/2})$ in expectation. Denote by $X^{(i)} := \sum_{S:|S|=\ell} a_S \cdot \mathbf{v}_S^{(i)}$ for $i = 0, 1, \dots, d$. We write $X^{(d)}$ as a telescoping sum $X^{(d)} = \sum_{i=1}^d (X^{(i)} - X^{(i-1)})$. To analyze the difference sequence, observe that in the expression

$$X^{(i)} - X^{(i-1)} = \sum_{S:|S|=\ell} a_S \cdot \left(\mathbf{v}_S^{(i)} - \mathbf{v}_S^{(i-1)} \right),$$

if set S contributes to the sum, then S must include the bit queried at the $(i-1)$ -th step of the path. Conditioning on v_0, \dots, v_{i-1} , let x_j be the variable queried in v_{i-1} , then we have

$$X^{(i)} - X^{(i-1)} = \sum_{S:|S|=\ell, j \in S} a_S \cdot \mathbf{v}_S^{(i)} = x_j \cdot \left(\sum_{S:|S|=\ell, j \in S} a_S \cdot \mathbf{v}_{S \setminus \{j\}}^{(i-1)} \right).$$

Furthermore, we observe that the sum $\sum_{S:|S|=\ell, j \in S} a_S \cdot \mathbf{v}_{S \setminus \{j\}}^{(i-1)}$ is determined by v_{i-1} ; thus conditioning on v_0, \dots, v_{i-1} the value of $X^{(i)} - X^{(i-1)}$ is a random coin in $\{\pm 1\}$ multiplied by some fixed integer. In other words, we get that $X^{(0)}, \dots, X^{(d)}$ is a martingale with varying step sizes.

Recall that Azuma's inequality provides concentration bounds for martingales with bounded step sizes, thus now we need to bound $\left| \sum_{S:|S|=\ell, j \in S} a_S \cdot \mathbf{v}_{S \setminus \{j\}}^{(i-1)} \right|$, which is similar to our initial goal. Put differently, we wish to analyze the sum

$$\sum_{S' \subseteq [n] \setminus \{j\}: |S'|=\ell-1} a_{S' \cup \{j\}} \cdot \mathbf{v}_{S'}^{(i-1)},$$

which calls for an inductive argument on ℓ . In addition, since we eventually apply a union bound on all steps, we need to show that $\left| \sum_{S'} a_{S' \cup \{j\}} \mathbf{v}_{S'}^{(i-1)} \right|$ is bounded with high probability (and not just in expectation).

More generally, to carry an inductive argument we define for any set $T \subseteq [n], |T| \leq \ell$ and any $i \in \{0, \dots, d\}$, the random variable

$$X_T^{(i)} := \sum_{S \supseteq T: |S|=\ell} a_S \cdot \mathbf{v}_{S \setminus T}^{(i)} = \sum_{S' \subseteq \bar{T}: |S'|=\ell-|T|} a_{S' \cup T} \cdot \mathbf{v}_{S'}^{(i)}.$$

Note that our initial goal was to bound $|X_\emptyset^{(d)}| = |X^{(d)}|$, which is analyzed by (reverse) induction on $|T|$ going from larger sets to smaller sets as [Lemma 1.9](#).

Lemma 1.9. *For all $t \in \{0, \dots, \ell\}$ and $\varepsilon > 0$, the probability that there exist $i \in \{0, \dots, d\}$ and $T \subseteq [n]$ of size at least t such that $|X_T^{(i)}| \geq \widetilde{O}_\varepsilon(d^{\ell-t}/2)$ is at most $\varepsilon \cdot (\ell - t)$.*

The main observation for the proof is that $X_T^{(0)}, X_T^{(1)}, \dots, X_T^{(d)}$ is a martingale whose difference sequence consists of terms of the form $X_{T'}^{(i-1)}$ where $T \subsetneq T'$. To see this, if we are querying x_j at v_{i-1} , then

$$X_T^{(i)} - X_T^{(i-1)} = \begin{cases} 0 & j \in T, \\ x_j \cdot \left(\sum_{j \notin S \subseteq \bar{T}} a_{S \cup T \cup \{j\}} \cdot \mathbf{v}_S^{(i-1)} \right) = x_j \cdot X_{T \cup j}^{(i-1)} & j \notin T. \end{cases}$$

Note that $X_{T \cup j}^{(i-1)}$ depends only on the history until v_{i-1} , and x_j is a uniformly random bit independent of this history, thus $X_T^{(i)}$ is a martingale. The inductive hypothesis implies that with at least $1 - \varepsilon \cdot (\ell - t - 1)$ probability, $|X_{T \cup j}^{(i-1)}| \leq \tilde{O}_\varepsilon(d^{\ell-t-1}/2)$ for all T of size t and $j \in [n] \setminus T$. Whenever this happens, Azuma's inequality implies that² with probability at least $1 - \varepsilon / (d \cdot n^t)$, we have

$$|X_T^{(i)}| \leq 2\sqrt{\log(d \cdot n^t / \varepsilon)} \cdot \sqrt{\sum_{i=1}^d \tilde{O}_\varepsilon(d^{\ell-t-1})} = \tilde{O}_\varepsilon(d^{\ell-t/2}).$$

This, along with a union bound over T of size t and $i \in \{0, \dots, d\}$ completes the inductive step. The Fourier bound for noisy decision trees can be proved using a similar approach.

Parity Decision Trees. The basic approach is as before. Let \mathcal{T} be a parity decision tree. As in (1), we use v_i and $\mathbf{v}^{(i)}$ to denote the random walk and the partial assignments to the variables respectively. We say v_i is k -clean if

$$\forall S \subseteq [n], |S| \leq k, \quad \mathbf{v}_S^{(i)} = \begin{cases} 1 & \text{if } x_S \text{ is fixed to 1 before reaching } v_i, \\ -1 & \text{if } x_S \text{ is fixed to } -1 \text{ before reaching } v_i, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

For (2) to be true, we need that at least v_d is ℓ -clean. We note that this is not always true,³ but it is useful as it simplifies the study of high-level Fourier coefficients. To address this issue, we define a *cleanup* process for parity decision trees in which we make additional queries to ensure that certain key nodes are k -clean. We do this by recursively cleaning nodes in a top-down fashion so that for every node v in the original tree \mathcal{T} , any node v' in the new tree \mathcal{T}' obtained at the end of the cleanup step for v is k -clean.

The cleanup process is simple to describe: Let v_1, \dots, v_d be any root-to-leaf path in \mathcal{T} . Assume we have completed the cleanup process for v_1, \dots, v_{i-1} . We then query the parity at v_i . While there exists a (minimal) set S violating (3), we pick and query an arbitrary coordinate in S . Once (3) is satisfied, we proceed to the cleanup process for v_{i+1} . This process increases the depth by a factor of at most k . We set $k = \Theta(\ell \cdot \log(n))$ and work with the new tree \mathcal{T}' of depth $D \leq k \cdot d$.

Let v_0, \dots, v_D be a random root-to-leaf path in \mathcal{T}' and $I_i, i \in [D]$ be the set of coordinates fixed due to the query at v_{i-1} . Note that this set might be of size larger than 1.⁴ It follows from simple linear algebra that $\sum_{i=1}^D |I_i| \leq D$. Since v_D is k -clean, (2) holds. Defining $X_T^{(i)}$ exactly as before,

²Technically this is not true, since a martingale after conditioning may not still be a martingale. We handle this by truncating the martingale when a bad event happens instead of conditioning on the good event.

³For example, let $S = \{1, 2\}$ and consider the parity decision tree whose only query is $x_1 x_2$. At any leaf, the value of $x_1 x_2$ is fixed, however, the values of x_1 and x_2 are free, hence S violates (3).

⁴For example, suppose we query $x_1 x_2, x_1 x_3, x_1 x_4$ and finally x_1 . Then, the last query reveals 4 coordinates.

our goal is to prove [Lemma 1.9](#) with D instead of d . The proof is still by induction on $\ell - t$. It turns out that $X_T^{(0)}, X_T^{(1)}, \dots, X_T^{(D)}$ is no longer a martingale; instead, $X_T^{(i)} - X_T^{(i-1)} = Y_i + Z_i$ where

$$Y_i := \sum_{\substack{\emptyset \neq J \subseteq I_i \cap \bar{T} \\ |J| \text{ is even}}} x_J \cdot X_{J \cup T}^{(i-1)} \quad \text{and} \quad Z_i := \sum_{\substack{\emptyset \neq J \subseteq I_i \cap \bar{T} \\ |J| \text{ is odd}}} x_J \cdot X_{J \cup T}^{(i-1)}. \quad (4)$$

and Z_i (resp., Y_i) is an odd (resp., even) polynomial of degree at most ℓ over the newly fixed variables $\{x_j \mid j \in I_i\}$. Conditioning on v_{i-1} , every pair of random bits $(x_j, x_{j'})$ from $\{x_j \mid j \in I_i\}$ is either identical ($x_j \equiv x_{j'}$) or opposite ($x_j \equiv -x_{j'}$), which means Y_i is a constant and Z_i can be written as $z_i \cdot |Z_i|$ where $|Z_i|$ is a constant and $z_i \sim \{\pm 1\}$.

For now, let us ignore Y_i and assume that we have a martingale $X_T^{(i)}$ such that $X_T^{(i)} - X_T^{(i-1)} = z_i \cdot |Z_i|$, where $z_i \sim \{\pm 1\}$ is a uniformly random bit independent of z_0, \dots, z_{i-1} and $|Z_i|$ depends only on v_{i-1} . Combined with an adaptive version of Azuma's inequality, we only need to show the sum of squares of step sizes $\sum_{i=1}^D |Z_i|^2$ is $\tilde{O}_\varepsilon(D^{\ell-t})$ to prove $|X_T^{(i)}| = \tilde{O}_\varepsilon(D^{(\ell-t)/2})$. By the induction hypothesis, with probability at least $1 - \varepsilon \cdot (\ell - t - 1)$ the coefficients of Z_i are bounded appropriately. Since $\sum_{i=1}^D |I_i| \leq D$ and in particular $|I_i| \leq D$, we have

$$|Z_i| \leq \sum_{\text{odd } j \geq 1} \binom{|I_i|}{j} \cdot \max_{|T'|=j+t} |X_{T'}^{(i-1)}| \leq \sum_{j \geq 1}^{\ell-t} \binom{|I_i|}{j} \cdot \tilde{O}_\varepsilon(D^{(\ell-j-t)/2}) = \tilde{O}_\varepsilon(|I_i| \cdot D^{(\ell-t-1)/2})$$

and thus $\sum_{i=1}^D |Z_i|^2 \leq D^2 \cdot \tilde{O}_\varepsilon(D^{\ell-t-1})$. This is too loose for our purpose.

We instead try to bound the sum of squares of step sizes *with high probability*. Imagine for now that v_{i-1} is 2-clean.⁵ Then, the variables $\{x_j \mid j \in I_i\}$ are 2-wise independent conditioning on v_{i-1} . This gives

$$\mathbb{E} \left[|Z_i|^2 \mid v_{i-1} \right] \leq \sum_{\text{odd } j \geq 1} \binom{|I_i|}{j} \cdot \max_{|T'|=j+t} |X_{T'}^{(i-1)}|^2 \leq \sum_{j \geq 1}^{\ell-t} \binom{|I_i|}{j} \cdot \tilde{O}_\varepsilon(D^{\ell-j-t}) = \tilde{O}_\varepsilon(|I_i| \cdot D^{\ell-t-1})$$

and thus $\mathbb{E} \left[\sum_{i=1}^D |Z_i|^2 \right] \leq \tilde{O}_\varepsilon(D^{\ell-t})$. To show this bound holds with high probability, we use concentration properties of degree- ℓ polynomials under k -wise independent distributions for $k \gg \ell$.

In the actual proof, we proceed by conditioning on $C(v_{i-1})$, the nearest ancestor of v_{i-1} that is k -clean, instead of conditioning on v_{i-1} , which allows to remove the assumption that v_{i-1} is 2-clean. This is because the queries within a cleanup step are non-adaptive, thus Z_i depends only on $C(v_{i-1})$ and not on v_{i-1} .

Meanwhile, although $X_T^{(i)}$ is not quite a martingale sequence (due to Y_i) and the step sizes (i.e., $|Z_i|$) are adaptive and not always bounded, we are nonetheless able to prove an adaptive version of Azuma's inequality of the form $\Pr \left[\max_{i \in [D]} |X_T^{(i)}| \geq \mu + t \cdot \sigma \right] \leq e^{-\Omega(t^2)} + \varepsilon$ provided $\Pr \left[\left(\sum_{i=1}^D |Y_i| \leq \mu \right) \wedge \left(\sum_{i=1}^D |Z_i|^2 \leq \sigma^2 \right) \right] \geq 1 - \varepsilon$. Then it suffices to bound $\sum_{i=1}^D |Y_i|$ similarly to $\sum_{i=1}^D |Z_i|^2$ above.

1.4 Related Work

We remark that our proof for level- ℓ Fourier growth (even when specialized to the case of standard decision trees) differs from the proofs appearing in [\[Tal20\]](#) and [\[SSW20\]](#). There, the results were

⁵This assumption immediately implies that $|I_i| \leq 1$ and trivially proves our inequality, however, this type of reasoning doesn't generalize to the case when v_{i-1} is not 2-clean.

based on decompositions of decision trees. We view our martingale approach as natural and intuitive. We wonder if one can obtain the tight results from [SSW20] using this approach. It seems that the main bottleneck is a union bound on events related to all sets $T \subseteq [n]$ of size at most ℓ .

Our bounds for level-1 improve those obtained by [BTW15]. They prove that $L_{1,1}(\mathcal{T}) \leq O(\sqrt{p \cdot d})$ when $p = \Pr_x[\mathcal{T}(x) = 1]$, whereas we obtain a bound of

$$L_{1,1}(\mathcal{T}) \leq O\left(p\sqrt{d} \cdot \log(1/p)\right).$$

In particular, our bound is almost quadratically better for small values of p . It remains open whether the bound can be further improved to $O\left(p\sqrt{d \cdot \log(1/p)}\right)$, which is the optimal bound for standard decision trees.

We remark that our cleanup technique is inspired by [BTW15], which used cleanup to prove their level-1 bound. However, our proof strategies and the way we use the cleanup procedure is quite different than that of [BTW15].

Organization. We make formal definitions in Section 2. We state and prove the necessary concentration inequalities in Section 3. We present the cleanup process in Section 4. We present the Fourier bounds for parity decision trees in Section 5 and for noisy decision trees in Section 6.

2 Preliminaries

We use $\log(\cdot)$ to denote the logarithm with base 2. We use $[n]$ to denote $\{1, 2, \dots, n\}$; and $\binom{[n]}{k}$ (resp., $\binom{[n]}{\leq k}$) to denote the set of all size- k (resp., size-at-most- k) sets from $[n]$. If S is a set from universe U , then we write \bar{S} for $U \setminus S$. We use \mathcal{U}_n to denote the uniform distribution over $\{\pm 1\}^n$. We use $\text{sgn}(\text{value}) \in \{-1, 0, 1\}$ to denote the sign of value , i.e., $\text{sgn}(\text{value})$ equals -1 if $\text{value} < 0$, 1 if $\text{value} > 0$, and 0 if $\text{value} = 0$.

We use $\mathbb{F}_2 = \{0, 1\}$ to denote the binary field, $\text{Span}(\text{vectors})$ to denote the subspace spanned by vectors over \mathbb{F}_2 . For a distribution \mathcal{D} we use $x \sim \mathcal{D}$ to represent that x is a random variable sampled from \mathcal{D} . For a finite set \mathcal{X} we use $x \sim \mathcal{X}$ to denote that x is a random variable sampled uniformly from \mathcal{X} . We use the standard notion of k -wise independent distribution over $\{\pm 1\}^n$.

Definition 2.1 (k -wise independence). A distribution \mathcal{D} over $\{\pm 1\}^n$ is k -wise independent if for $x \sim \mathcal{D}$ and any k -indices $1 \leq i_1 < i_2 < \dots < i_k \leq n$, the random variables $(x_{i_1}, \dots, x_{i_k})$ are uniformly distributed over $\{\pm 1\}^k$.

2.1 Boolean Functions

Here we recall definitions in the analysis of Boolean functions (see [O'D14] for a detailed introduction). Let $f: \{\pm 1\}^n \rightarrow \mathbb{R}$ be any Boolean function. For any $p > 0$, the p -norm of f is defined as $\|f\|_p = (\mathbb{E}_{x \sim \mathcal{U}_n} [|f(x)|^p])^{1/p}$. For any subset $S \subseteq [n]$, x_S denotes $\prod_{i \in S} x_i$ (in particular, $x_\emptyset = 1$). It is a well-known fact that we can uniquely represent f as a linear combination of $\{x_S\}_{S \subseteq [n]}$:

$$f(x) = \sum_{S \subseteq [n]} \hat{f}(S) x_S,$$

where the coefficients $\{\hat{f}(S)\}_{S \subseteq [n]}$ are referred to as the *Fourier coefficients* of f and are given by $\hat{f}(S) = \mathbb{E}_{x \sim \mathcal{U}_n} [f(x) x_S]$. The above representation expresses f as a multilinear polynomial and is called the Fourier representation of f . We say that f is of degree at most d if its Fourier representation is a polynomial of degree at most d , i.e., if $\hat{f}(S) = 0$ for all $S \subseteq [n]$, $|S| > d$.

2.2 Parity Decision Trees

Here we formally define parity decision trees (with Boolean outputs).

Definition 2.2 (Parity decision tree). A *parity decision tree* \mathcal{T} is a representation of a Boolean function $f: \{\pm 1\}^n \rightarrow \{0, 1\}$. It consists of a rooted binary tree in which each internal node v is labeled by a non-empty set $Q_v \subseteq [n]$, the outgoing edges of each internal node are labeled by $+1$ and -1 , and the leaves are labeled by 0 and 1 .

On input $x \in \{\pm 1\}^n$, the tree \mathcal{T} constructs a *computation path* \mathcal{P} from the root to a leaf. Specifically, when \mathcal{P} reaches an internal node v we say that \mathcal{T} *queries* Q_v ; then \mathcal{P} follows the outgoing edge labeled by $\prod_{i \in Q_v} x_i$. We require that Q_v is not implied by its ancestors' queries. The output of \mathcal{T} (and hence f) on input x is the label of the leaf reached by the computation path. Conversely, we say x is *consistent with* the path \mathcal{P} if \mathcal{P} is the computation path (possibly ending before reaching a leaf) for x .

We make a few more remarks on a parity decision tree $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$.

- A node v in \mathcal{T} can be either an internal node or a leaf, and we use $\mathcal{T}(v) \in \{0, 1\}$ to denote the label on v when v is a leaf. Meanwhile, we use \mathcal{T}_v to denote the sub parity decision tree starting with node v .
- The *depth* of a node is the number of its ancestors (e.g., the root has depth 0) and the depth of \mathcal{T} is the maximum depth over all its leaves.
- We say that two parity decision trees \mathcal{T} and \mathcal{T}' are *equivalent* (denoted by $\mathcal{T} \equiv \mathcal{T}'$) if they compute the same function.

2.3 Noisy Decision Trees

Definition 2.3 (Noisy oracle). A noisy query to a bit $b \in \{\pm 1\}$ with correlation $\gamma \in [-1, 1]$ returns a bit $b' \in \{\pm 1\}$ where

$$b' = \begin{cases} b & \text{with probability } (1 + \gamma)/2, \\ -b & \text{with probability } (1 - \gamma)/2. \end{cases}$$

The cost of a noisy query with correlation γ is defined to be γ^2 .

Definition 2.4 (Noisy decision tree). A *noisy decision tree* \mathcal{T} is a rooted binary tree in which each internal node v is labeled by an index $q_v \in [n]$ and a correlation $\gamma_v \in [-1, 1]$. The outgoing edges are labeled by $+1$ and -1 and the leaves are labeled by 0 and 1 .

On input $x \in \{\pm 1\}^n$, the tree \mathcal{T} constructs a *computation path* \mathcal{P} from the root to leaf as follows. When \mathcal{P} reaches an internal node v , it makes a noisy query to x_{q_v} with correlation γ_v and follows the edge labeled by the outcome of this noisy query. The output of the tree is defined by sampling a root-to-leaf path and returning the label of the leaf. Since the computation path \mathcal{P} is probabilistic, this is an inherently randomized model of computation. We use $\mathcal{T}(x) \in \{0, 1\}$ to denote the (probabilistic) output of \mathcal{T} on input x . We also use $\mathcal{T}(v) \in \{0, 1\}$ to denote the label on v when v is a leaf. We do *not* require that the indices q_v queried along a path \mathcal{P} are distinct. The *cost* of any path is the sum of costs of the noisy queries along that path; and the cost of \mathcal{T} is the maximum cost of any root-to-leaf path.

We remark that for any noisy decision tree \mathcal{T} , its Fourier coefficient $\widehat{\mathcal{T}}(S)$ is given by $\mathbb{E}[\mathcal{T}(x)x_S]$ where the expectation is over the randomness of both $x \sim \mathcal{U}_n$ and \mathcal{T} .

3 Useful Concentration Inequalities

We describe useful concentration inequalities in this section.

3.1 Low Degree Polynomials

We use the fact that low degree polynomials satisfy strong concentration properties under k -wise independent distributions. We will find the following hypercontractive inequality useful.

Theorem 3.1 ([Bon70], see also [O'D14, (2, q)-hypercontractivity]). *Let $f: \{\pm 1\}^n \rightarrow \mathbb{R}$ be a degree- d polynomial. Then for any $q \geq 2$, we have $\|f\|_q \leq (q-1)^{d/2} \|f\|_2$.*

Lemma 3.2. *Let $f: \{\pm 1\}^n \rightarrow \mathbb{R}$ be a degree- d polynomial. Let \mathcal{D} be a $2k$ -wise independent distribution over $\{\pm 1\}^n$, where $k \geq d$. Let $\mu = \mathbb{E}_{x \sim \mathcal{D}} [f(x)]$ and $\sigma^2 = \mathbb{E}_{x \sim \mathcal{D}} [(f(x) - \mu)^2]$. Then for any $\alpha > 0$ and any integer $1 \leq \ell \leq k/d$, we have*

$$\mathbb{E}_{x \sim \mathcal{D}} [(f(x) - \mu)^{2\ell}] \leq \sigma^{2\ell} \cdot (2\ell - 1)^{d \cdot \ell}.$$

In particular we have

$$\Pr_{x \sim \mathcal{D}} [|f(x) - \mu| \geq \alpha \cdot \sigma] \leq \alpha^2 \cdot \left(\frac{2k}{d \cdot \alpha^{2/d}} \right)^k.$$

Proof. Observe that $(f(x) - \mu)^{2\ell}$ is a polynomial of degree at most $2\ell \cdot d \leq 2k$. Thus its expectation under \mathcal{D} is the same as its expectation under the uniform distribution over $\{\pm 1\}^n$. By [Theorem 3.1](#), we have

$$\|f - \mu\|_{2\ell} \leq (2\ell - 1)^{d/2} \|f - \mu\|_2 = \sigma \cdot (2\ell - 1)^{d/2}.$$

Hence by Markov's inequality, we have

$$\Pr_{x \sim \mathcal{D}} [|f(x) - \mu| \geq \alpha \cdot \sigma] \leq \frac{\mathbb{E}_{x \sim \mathcal{D}} [(f(x) - \mu)^{2\ell}]}{(\alpha \cdot \sigma)^{2\ell}} = \frac{\|f - \mu\|_{2\ell}^{2\ell}}{(\alpha \cdot \sigma)^{2\ell}} \leq \frac{(2\ell - 1)^{\ell \cdot d}}{\alpha^{2\ell}}.$$

Now we derive the second bound. We only need to focus on the case $\alpha \geq 1$ since otherwise the RHS is at least 1. Then by setting $\ell = \lfloor k/d \rfloor$, we have

$$\Pr_{x \sim \mathcal{D}} [|f(x) - \mu| \geq \alpha \cdot \sigma] \leq \frac{(2\lfloor k/d \rfloor - 1)^{\lfloor k/d \rfloor \cdot d}}{\alpha^{2\lfloor k/d \rfloor}} \leq \frac{(2k/d)^k}{\alpha^{2(k/d-1)}} = \alpha^2 \cdot \left(\frac{2k}{d \cdot \alpha^{2/d}} \right)^k. \quad \square$$

3.2 Martingales

We show an adaptive version of Azuma's inequality for martingales. The proof is similar to the inductive proof of the standard Azuma's inequality and thus deferred to [Appendix B](#).

Lemma 3.3 (Adaptive Azuma's inequality). *Let $X^{(0)}, \dots, X^{(D)}$ be a martingale and $\Delta^{(1)}, \dots, \Delta^{(D)}$ be a sequence of magnitudes such that $X^{(0)} = 0$ and $X^{(i)} = X^{(i-1)} + \Delta^{(i)} \cdot z^{(i)}$ for $i \in [D]$, where if conditioning on $z^{(1)}, \dots, z^{(i-1)}$,*

- (1) $z^{(i)}$ is a mean-zero random variable and $|z^{(i)}| \leq 1$ always holds;
- (2) $\Delta^{(i)}$ is a fixed value.

If there exists some constant $U \geq 0$ such that $\sum_{i=1}^D |\Delta^{(i)}|^2 \leq U$ always holds, then for any $\beta \geq 0$ we have

$$\Pr \left[\max_{i=0,1,\dots,D} |X^{(i)}| \geq \beta \cdot \sqrt{2U} \right] \leq 2 \cdot e^{-\beta^2/2}.$$

Next, we generalize [Lemma 3.3](#) as follows.

Lemma 3.4. *Let $m \geq 1$ be an integer. For each $t \in [m]$, let $X_t^{(0)}, \dots, X_t^{(D)}$ be a sequence of random variables and $\Delta_t^{(1)}, \dots, \Delta_t^{(D)}$ be a sequence of magnitudes such that $X_t^{(0)} = 0$ and $X_t^{(i)} = X_t^{(i-1)} + \Delta_t^{(i)} \cdot z_t^{(i)} + \mu_t^{(i)}$ for $i \in [D]$, where if conditioning on $z_t^{(1)}, \dots, z_t^{(i-1)}$,*

(1) $z_t^{(i)}$ is a mean-zero random variable and $|z_t^{(i)}| \leq 1$ always holds;

(2) $\Delta_t^{(i)}$ is a fixed value and $\mu_t^{(i)}$ is a random variable.

If there exist some constants $U, V \geq 0$ and $\eta \in [0, 1]$ such that

$$\Pr \left[\exists t \in [m], \left(\sum_{i=1}^D |\Delta_t^{(i)}|^2 > U \right) \vee \left(\sum_{i=1}^D |\mu_t^{(i)}| > V \right) \right] \leq \eta,$$

then for any $\beta \geq 0$ we have

$$\Pr \left[\exists t \in [m], \max_{i=0,1,\dots,D} |X_t^{(i)}| \geq V + \beta \cdot \sqrt{2U} \right] \leq \eta + 2m \cdot e^{-\beta^2/2}.$$

Proof. We divide the proof into the following two cases.

Case $\eta = 0$. Let $\widehat{X}_t^{(i)} = X_t^{(i)} - \sum_{j=1}^i \mu_t^{(j)}$ for each t and i . Then $|X_t^{(i)}| = |\widehat{X}_t^{(i)} + \sum_{j=1}^i \mu_t^{(j)}| \leq V + |\widehat{X}_t^{(i)}|$. By a union bound, it suffices to show for any fixed t , we have

$$\Pr \left[\max_{i=0,1,\dots,D} |\widehat{X}_t^{(i)}| \geq \beta \cdot \sqrt{2U} \right] \leq 2 \cdot e^{-\beta^2/2},$$

which follows from [Lemma 3.3](#).

Case $\eta \geq 0$. Consider $\widetilde{X}_t^{(0)}, \dots, \widetilde{X}_t^{(D)}$ defined by setting $\widetilde{X}_t^{(0)} = 0$ and $\widetilde{X}_t^{(i)} = \widetilde{X}_t^{(i-1)} + \widetilde{\Delta}_t^{(i)} \cdot z_t^{(i)} + \widetilde{\mu}_t^{(i)}$, where

$$\widetilde{\Delta}_t^{(i)} = \begin{cases} \Delta_t^{(i)} & \sum_{j=1}^i |\Delta_t^{(j)}|^2 \leq U, \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad \widetilde{\mu}_t^{(i)} = \begin{cases} \mu_t^{(i)} & \sum_{j=1}^i |\mu_t^{(j)}| \leq V, \\ 0 & \text{otherwise.} \end{cases}$$

Then Item (1) and (2) hold for $(\widetilde{X}_t^{(i)})_{t,i}$ and $(\widetilde{\Delta}_t^{(i)})_{t,i}, (\widetilde{\mu}_t^{(i)})_{t,i}$.

Note that $\Pr \left[\exists t \in [m], i \in \{0, 1, \dots, D\}, \widetilde{X}_t^{(i)} \neq X_t^{(i)} \right] \leq \eta$ and $\sum_{i=1}^D |\widetilde{\Delta}_t^{(i)}|^2 \leq U, \sum_{i=1}^D |\widetilde{\mu}_t^{(i)}| \leq V$ always. Hence from the previous case, we have

$$\begin{aligned} & \Pr \left[\exists t \in [m], \max_{i=0,1,\dots,D} |X_t^{(i)}| \geq V + \beta \cdot \sqrt{2U} \right] \\ & \leq \Pr \left[\exists t \in [m], i \in \{0, 1, \dots, D\}, \widetilde{X}_t^{(i)} \neq X_t^{(i)} \right] + \Pr \left[\exists t \in [m], \max_{i=0,1,\dots,D} |\widetilde{X}_t^{(i)}| \geq V + \beta \cdot \sqrt{2U} \right] \\ & \leq \eta + 2m \cdot e^{-\beta^2/2}. \quad \square \end{aligned}$$

4 How to Clean Up Parity Decision Trees

In this section we show how to *clean up* the given parity decision tree to make it easier to analyze.

4.1 k -cleanness

It will be useful to identify \mathbb{F}_2^n with $\{\pm 1\}^n$ by $\text{Enc}: (x_1, \dots, x_n) \mapsto ((-1)^{x_1}, \dots, (-1)^{x_n})$. For a subset $X \subseteq \mathbb{F}_2^n$ we will denote $\text{Enc}(X) = \{\text{Enc}(x) : x \in X\}$. Thus, we may think of Boolean functions also as $f: \mathbb{F}_2^n \rightarrow \{0, 1\}$. We observe that under this representation of the input, a parity decision tree $\mathcal{T}: \mathbb{F}_2^n \rightarrow \{0, 1\}$ indeed queries parity functions (i.e., linear functions over \mathbb{F}_2) of the input bits $x \in \mathbb{F}_2^n$ and decides whether to go left or right based on their outcome. Thus, the set of all possible inputs in \mathbb{F}_2^n that reach a given node in a parity decision tree is an affine subspace of \mathbb{F}_2^n .

We introduce some notation.

Notation 4.1. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a parity decision tree and let v be a node in it.

- We use $\mathcal{P}_v \subseteq \{\pm 1\}^n$ to denote the set of all points reaching node v . Note that $\mathcal{P}_v = \text{Enc}(H_v + a)$ where H_v is a linear subspace of \mathbb{F}_2^n of dimension $n - \text{depth}(v)$ and $a \in \mathbb{F}_2^n$.
- For any $S \subseteq [n]$, we define $\widehat{\mathcal{P}}_v(S) = \mathbb{E}_{x \sim \mathcal{P}_v}[x_S]$.
- We use \mathcal{S}_v to denote all fully correlated sets with \mathcal{P}_v , i.e., $\mathcal{S}_v = \left\{ S \subseteq [n] \mid \widehat{\mathcal{P}}_v(S) \in \{\pm 1\} \right\}$. We observe that if $\mathcal{P}_v = \text{Enc}(H_v + a)$, then $\mathcal{S}_v = H_v^\perp$. Additionally, if the queries on the path from root to v are $Q_{v_0}, \dots, Q_{v_{i-1}}$, then $\mathcal{S}_v = \text{Span}\langle \{Q_{v_0}, \dots, Q_{v_{i-1}}\} \rangle$.
- If v is an internal node, then define $J(v)$ as the set of newly fixed coordinates after querying Q_v , i.e., $i \in J(v)$ iff $\{i\} \notin \mathcal{S}_v$ but $\{i\} \in \text{Span}\langle \mathcal{S}_v \cup \{Q_v\} \rangle$.

The following simple fact shows that there is no “somewhat” correlated set.

Fact 4.2. For any parity decision tree \mathcal{T} and any node v in \mathcal{T} , $\widehat{\mathcal{P}}_v(S) \in \{+1, 0, -1\}$ holds for any set S .

Proof. Since $\mathcal{P}_v = \text{Enc}(H_v + a)$ where $H_v + a$ is an affine subspace, \mathcal{P}_v falls into one of the following 3 cases: (a) all points in \mathcal{P}_v satisfy $\chi_S(x) = 1$, (b) all points satisfy $\chi_S(x) = -1$, (c) exactly half of the points satisfy $\chi_S(x) = 1$. \square

Let $\mathcal{S} \subseteq \mathbb{F}_2^n$ be a subspace and $S \subseteq [n]$. For simplicity, we write $S \in \mathcal{S}$ iff the indicator vector of S is contained in \mathcal{S} . Now we describe the desired property: *k-clean*.

Definition 4.3 (*k-clean subspace and mess-witness*). Let k be a positive integer. A subspace \mathcal{S} is *k-clean* if for any set $S \in \mathcal{S}$ such that $|S| \leq k$, we have that $\{i\} \in \mathcal{S}$ holds for any $i \in S$.

Moreover, when \mathcal{S} is not *k-clean*, we say i is a *mess-witness* if there exists some $S \ni i, |S| \leq k$ such that $S \in \mathcal{S}$ but $\{i\} \notin \mathcal{S}$.

Definition 4.4 (*k-clean parity decision tree*). A parity decision tree \mathcal{T} is *k-clean* if the following holds:

- For any internal node v , either (a) \mathcal{S}_v is *k-clean*, or (b) $Q_v = \{i\}$ where i is a mess-witness for \mathcal{S}_v . Moreover, we say v is *k-clean* if (a) holds; and we say v is *cleaning* if (b) holds.
- For any leaf v , \mathcal{S}_v is *k-clean* (in such a case, we say that v is *k-clean*).

- For any k -clean internal node v , \mathcal{T}_v starts with $\ell(v)$ non-adaptive queries⁶ where $\ell(v) \geq 1$. In addition, for any $i \in \{1, \dots, \ell(v) - 1\}$, any node of depth i in \mathcal{T}_v is cleaning; and all node of depth $\ell(v)$ are k -clean.⁷

Example 4.5. If \mathcal{T} is a decision tree (i.e., $|Q_v| \equiv 1$ for any internal node v) then it is k -clean for any k , where each internal node is k -clean.

If \mathcal{T} is the depth-1 parity decision tree for $\mathcal{T}(x) = x_1x_2x_3$ (i.e., \mathcal{T} only has a root v_0 querying $Q_{v_0} = \{1, 2, 3\}$), then it is 2-clean but not 3-clean, since for either leaf v we have $\{1, 2, 3\} \in \mathcal{S}_v$ but $\{1\} \notin \mathcal{S}_v$.

The benefit of having a k -clean parity decision tree is that it makes the expression of Fourier coefficients simpler.

Lemma 4.6. *Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a k -clean parity decision tree and let S be a set of size $\ell \leq k$. Let v_0, \dots, v_d be a random root-to-leaf path. Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(d)} \in \{-1, 0, +1\}^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each i, j . Recall that $\mathbf{v}_S^{(d)} = \prod_{j \in S} \mathbf{v}_j^{(d)}$. Then we have*

$$\widehat{\mathcal{T}}(S) = \mathbb{E}_{v_0, \dots, v_d} [\mathcal{T}(v_d) \cdot \mathbf{v}_S^{(d)}].$$

Proof. Observe that for any $j \in J(v_i) \subseteq J$, the j -th coordinate is fixed after querying Q_{v_i} . Therefore we have

$$\widehat{\mathcal{T}}(S) = \mathbb{E}_{y \sim \mathcal{U}_n} [\mathcal{T}(y) \cdot y_S] = \mathbb{E}_{v_0, \dots, v_d} \left[\mathcal{T}(v_d) \cdot \mathbb{E}_{y \sim \mathcal{P}_{v_d}} [y_S] \right] = \mathbb{E}_{v_0, \dots, v_d} [\mathcal{T}(v_d) \cdot \widehat{\mathcal{P}}_{v_d}(S)]$$

By [Fact 4.2](#), $\widehat{\mathcal{P}}_{v_d}(S) \neq 0$ iff $S \in \mathcal{S}_{v_d}$, which, due to $\ell \leq k$ and v_d being a k -clean leaf, is equivalent to all coordinates in S being fixed along this path. Hence $\widehat{\mathcal{P}}_{v_d}(S) = \prod_{j \in S} \mathbf{v}_j^{(d)}$. \square

4.2 Cleanup Process

We first analyze the cleanup process for a subspace.⁸

Lemma 4.7 (Clean subspace). *Let $k \geq 2$ be an integer and \mathcal{S} be a subspace of rank at most d . We construct a new subspace \mathcal{S}' (initialized as \mathcal{S}) as follows: while \mathcal{S}' is not k -clean, we continue to update $\mathcal{S}' \leftarrow \text{Span} \langle \mathcal{S}' \cup \{\{i\}\} \rangle$ with some mess-witness i . Then $\text{rank}(\mathcal{S}') \leq d \cdot k$ and any update choice of mess-witnesses will result in the same final subspace \mathcal{S}' .*

Proof. Assume \mathcal{S} is a subspace of \mathbb{F}_2^n . Then first note that the number of updates is finite, since we can update for at most n times.

Next we show that the number of updates and the final \mathcal{S}' does not depend on the choice of mess-witnesses. We do so by an exchange argument. Let i_1, \dots, i_r and $i'_1, \dots, i'_{r'}$ be two rounds of execution using different mess-witnesses. Then there exists some $t < \min\{r, r'\}$ such that $i_j = i'_j$ for all $j \leq t$, but $i_{t+1} \neq i'_{t+1}$. Let $\mathcal{S}_t = \text{Span} \langle \mathcal{S} \cup \{\{i_1\}, \dots, \{i_t\}\} \rangle$. Then there exist $S \ni i_{t+1}$ and $S' \ni i'_{t+1}$ (possibly $S = S'$) such that $S, S' \in \mathcal{S}_t$ but $\{i_{t+1}\}, \{i'_{t+1}\} \notin \mathcal{S}_t$. Since the final subspace is k -clean, we know there exists some $T \geq t$ such that

$$\{i_{t+1}\} \notin \text{Span} \langle \mathcal{S} \cup \{\{i'_1\}, \dots, \{i'_T\}\} \rangle \quad \text{but} \quad \{i_{t+1}\} \in \text{Span} \langle \mathcal{S} \cup \{\{i'_1\}, \dots, \{i'_{T+1}\}\} \rangle,$$

⁶This means for any $i \in \{0, 1, \dots, \ell(v) - 1\}$, all nodes of depth i in \mathcal{T}_v make the same query.

⁷This “leveled adaptive” condition is required just for convenience of proofs. In fact, one can show that the first few queries in \mathcal{T}_v can be rearranged to make sure they are non-adaptive until we reach a k -clean node. See [Lemma 4.7](#).

⁸The $k = 2$ case of [Lemma 4.7](#) is essentially [[BTW15](#), Proposition 3.5]. However there is a gap in their proof. For example, if the parity decision tree non-adaptively queries $x_1x_2x_3x_4, x_1x_5, x_2x_6$ in order, then their analysis fails.

which means $\{i'_{T+1}, i_{t+1}\} \in \text{Span} \langle \mathcal{S} \cup \{\{i'_1\}, \dots, \{i'_T\}\} \rangle$. Hence we can safely replace i'_{T+1} with i_{t+1} , and then swap i_{t+1} with i'_{t+1} . We can perform this process as long as $(i_1, \dots, i_r) \neq (i'_1, \dots, i'_{r'})$, which means $r = r'$ and the final \mathcal{S}' is always the same.

For any subspace \mathcal{H} , we define $\text{rank}_1(\mathcal{H}) = |\{i \mid \{i\} \in \mathcal{H}\}|$ and thus $\text{rank}(\mathcal{H}) - \text{rank}_1(\mathcal{H}) \geq 0$. Now we analyze the following particular way to construct \mathcal{S}' : We initialize \mathcal{S}' as \mathcal{S} . While \mathcal{S}' is not k -clean, we find a minimal $S = \{i_1, \dots, i_s\} \in \mathcal{S}'$ such that i_1 is a mess-witness; then we update $\mathcal{S}' \leftarrow \text{Span} \langle \mathcal{S}' \cup \{\{i_1\}, \dots, \{i_{s-1}\}\} \rangle$. Note that before the update, $1 < s \leq k$ and $\{i_j\} \notin \mathcal{S}'$ holds for each $j \in [s]$, since S is minimal and \mathcal{S}' is not k -clean. Thus after the update, $\text{rank}(\mathcal{S}')$ grows by $s - 1 \leq k - 1$ and $\text{rank}_1(\mathcal{S}')$ grows by s , which means $\text{rank}(\mathcal{S}') - \text{rank}_1(\mathcal{S}')$ shrinks by 1. Hence we have at most $\text{rank}(\mathcal{S}) - \text{rank}_1(\mathcal{S}) \leq d$ updates before \mathcal{S}' is k -clean; and the final \mathcal{S}' has rank at most $\text{rank}(\mathcal{S}) + (k - 1) \cdot d \leq d \cdot k$. \square

We now show how to convert an arbitrary parity decision tree into a k -clean parity decision tree which still has a small depth and fixes a small number of variables along each path. The latter quantity is in fact bounded by the depth as shown in [Fact 4.8](#).

Fact 4.8. *Let \mathcal{T} be a depth- d parity decision tree. Let $v_0, \dots, v_{d'}$ be any root-to-leaf path. Then we have $\sum_{i=0}^{d'-1} |J(v_i)| \leq d'$.*

Proof. Observe that $\sum_{i=0}^{d'-1} |J(v_i)| = \left| \left\{ i \mid \{i\} \in \text{Span} \langle Q_{v_0}, \dots, Q_{v_{d'-1}} \rangle \right\} \right| \leq d'$. \square

Corollary 4.9. *Let \mathcal{T} be a depth- D k -clean parity decision tree. Let $v_0, \dots, v_{D'}$ be any root-to-leaf path where at most d of the nodes $v_0, \dots, v_{D'-1}$ are k -clean. Then $\sum_{i: |J(v_{i-1})| > 1} |J(v_i)| \leq 2d$.*

Proof. By [Fact 4.8](#) we have $\sum_{i=0}^{D'-1} |J(v_i)| - 1 \leq 0$. Since any v_i with $J(v_i) = \emptyset$ is not cleaning and therefore must be k -clean. Thus

$$\sum_{i: |J(v_i)| > 1} |J(v_i)| - 1 \leq |\{i : J(v_i) = \emptyset\}| \leq d.$$

For $|J(v_i)| > 1$, we have $|J(v_i)| - 1 \geq |J(v_i)|/2$ and thus $\sum_{i: |J(v_i)| > 1} |J(v_i)| \leq 2d$. \square

Lemma 4.10 (Clean parity decision tree). *Let $k \geq 2$ be an integer. Let \mathcal{T} be an arbitrary depth- d parity decision tree. Then there exists a k -clean parity decision tree \mathcal{T}' of depth at most $d \cdot k$ equivalent to \mathcal{T} . Moreover, any root-to-leaf path in \mathcal{T}' has at most d nodes that are k -clean.*

Proof. We build \mathcal{T}' by the following recursive algorithm. An example of the algorithm is provided in [Figure 1](#)

We now prove the correctness of [Algorithm 1](#), which is guaranteed by the following claims.

- For any internal node $v' \in \mathcal{T}'$, $Q_{v'}$ is not implied by its ancestors' queries. By [Fact 4.2](#), this is equivalent to $Q_{v'} \notin \mathcal{S}_{v'}$, which follows from the conditions in Line 8/9/13.
- The depth of \mathcal{T}' is at most $d \cdot k$. Let $v_0, \dots, v_{d'}$ be any root-to-leaf path of \mathcal{T} and let \mathcal{P}' be its corresponding path in \mathcal{T}' . Then the construction process of \mathcal{P}' corresponds to the cleanup process for $\text{Span} \langle Q_{v_0}, \dots, Q_{v_{d'-1}} \rangle$ in [Lemma 4.7](#); hence the depth of \mathcal{T}' equals $\text{rank}(\mathcal{S}') \leq d' \cdot k \leq d \cdot k$ where \mathcal{S}' is the k -clean subspace produced by applying [Lemma 4.7](#).
- $\mathcal{T} \equiv \mathcal{T}'$ and any root-to-leaf path in \mathcal{T}' has at most d k -clean nodes. This is evident from the algorithm, as \mathcal{T}' only refines \mathcal{T} by inserting cleaning nodes.

Algorithm 1: Clean parity decision tree: build \mathcal{T}' from \mathcal{T}

Input: an arbitrary depth- d parity decision tree \mathcal{T}

Output: a parity decision tree \mathcal{T}' with desired properties

```

1  $r \leftarrow$  root of  $\mathcal{T}$ 
2 Initialize the root of  $\mathcal{T}'$  as  $r'$ 
3 Build( $r, r', 1$ )
4 Procedure Build( $v, v', \ell$ )
   /* ( $v, v'$ ) are the current nodes on  $(\mathcal{T}, \mathcal{T}')$ ;  $\ell$  is the recursion depth. */
5   if  $v$  is a leaf then Label  $v'$  with the label of  $v$ 
6   else
7      $(v_-, v_+) \leftarrow$  the left and right child of  $v$ 
8     if  $\widehat{\mathcal{P}}_{v'}(Q_v) = -1$  then Build( $v_-, v', \ell + 1$ )
9     else if  $\widehat{\mathcal{P}}_{v'}(Q_v) = +1$  then Build( $v_+, v', \ell + 1$ )
10    else /*  $\widehat{\mathcal{P}}_{v'}(Q_v) = 0$  due to Fact 4.2 */
11       $Q_{v'} \leftarrow Q_v$ 
12       $(v'_-, v'_+) \leftarrow$  the left and right child of  $v'$ 
13      Initialize  $O \leftarrow \emptyset$ 
14      while Span  $\langle \mathcal{S}_{v'} \cup \{Q_{v'}\} \cup O \rangle$  is not  $k$ -clean do
15        | Update  $O \leftarrow O \cup \{\{i\}\}$ , where  $i$  is a mess-witness
16      end
17       $\mathcal{T}'$  non-adaptively queries every set (which is a singleton) in  $O$  under  $v'$  in
        arbitrary order
18      foreach leaf  $\widehat{v}$  under  $v'_-$  do Build( $v_-, \widehat{v}, \ell + 1$ )
19      foreach leaf  $\widehat{v}$  under  $v'_+$  do Build( $v_+, \widehat{v}, \ell + 1$ )
20
21    end
22 end

```

- Whenever we call Build(\cdot, v', \cdot), v' is k -clean. We prove by induction on ℓ . The base case Line 3 is obvious. For Line 8/9, we recurse on the same v' , which is k -clean by induction. For Line 17/18, note that $\mathcal{S}_{\widehat{v}} = \text{Span} \langle \mathcal{S}_{v'} \cup \{Q_{v'}\} \cup O \rangle$; hence from the condition in Line 13, it is k -clean.
- Nodes created in Line 16 are cleaning. Let $o = |O|$ and let i_1, i_2, \dots, i_o be the query order. For any $j \in [o]$, let v'_j be any one of the nodes created for i_j , then

$$\mathcal{S}_{v'_j} = \text{Span} \langle \mathcal{S}_{v'} \cup \{Q_{v'}\} \cup \{\{i_1\}, \dots, \{i_{j-1}\}\} \rangle,$$

which is not k -clean by Line 13; hence v'_j is cleaning by the condition in Line 13. □

5 Fourier Bounds for Parity Decision Trees

Our goal in this section is to prove [Theorem 1.1](#) with detailed bounds provided.

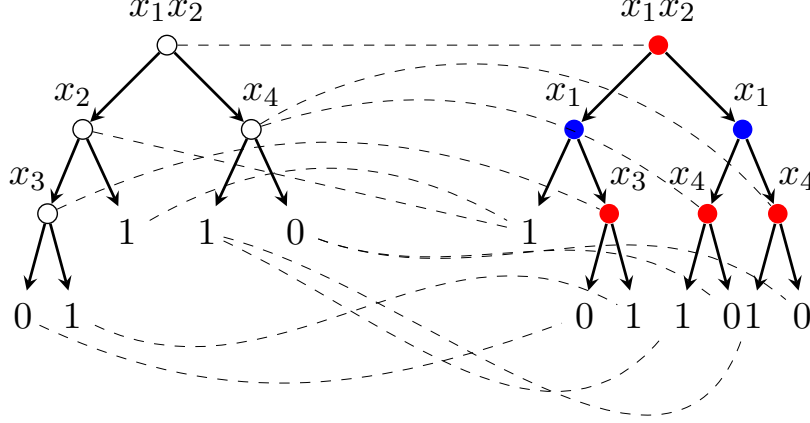


Figure 1: An example of the cleanup process with $k = 2$ where the LHS is \mathcal{T} and the RHS is \mathcal{T}' . All the left (resp., right) outgoing edges are labeled with -1 (resp., $+1$). **Red nodes** and leaves are k -clean, and **blue nodes** are cleaning (i.e., non-adaptive queries). Nodes connected with dashed curves are invoked by Build.

5.1 Level-1 Bound

We first prove the concentration result for level-1. We start with the following simple bound for general parity decision trees.

Lemma 5.1. *Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a depth- D parity decision tree. Let $v_0, \dots, v_{D'}$ be any root-to-leaf path. Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(D')} \in \{-1, 0, +1\}^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq D'$ and $j \in [n]$. Then for any $a_1, \dots, a_n \in \{-1, 0, 1\}$, we have $\left| \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(D')} \right| \leq D' \leq D$.*

Proof. Note that the set of non-zero coordinates in $\mathbf{v}^{(D')}$ is exactly $\bigcup_{i=0}^{D'-1} J(v_i)$. Hence by [Fact 4.8](#), we have

$$\left| \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(D')} \right| \leq \sum_{j=1}^n \left| \mathbf{v}_j^{(D')} \right| = \sum_{i=0}^{D'-1} |J(v_i)| \leq D' \leq D. \quad \square$$

Now we give an improved bound for k -clean parity decision trees. To do so, we need one more notation which will be crucial in our analysis.

Notation 5.2. Let \mathcal{T} be a k -clean parity decision tree. For any node v , we define $C(v)$ as the nearest ancestor of v (including itself) that is k -clean.

Lemma 5.3. *There exists a universal constant $\kappa \geq 1$ such that the following holds. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a depth- D $2k$ -clean parity decision tree where $k \geq 1$ and any root-to-leaf path has at most d nodes that are $2k$ -clean.*

Let $v_0, \dots, v_{D'}$ be a random root-to-leaf path. Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(D')} \in \{-1, 0, +1\}^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq D'$ and $j \in [n]$. Then for any $a_1, \dots, a_n \in \{-1, 0, 1\}$ and any $\varepsilon \leq 1/2$, we have $\Pr \left[\left| \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(D')} \right| \geq R(D, d, k, \varepsilon) \right] \leq \varepsilon$, where

$$R(D, d, k, \varepsilon) = \kappa \cdot \sqrt{\left(D + dk \left(\frac{1}{\varepsilon} \right)^{\frac{1}{k}} \right) \log \left(\frac{1}{\varepsilon} \right)}.$$

In the proof of [Lemma 5.3](#) we will use the following simple claim.

Fact 5.4. *Let p_1, \dots, p_n be a sub-probability distribution, i.e., $p_i \geq 0$ and $\sum_{i=1}^n p_i \leq 1$. Let $a_1, \dots, a_n \in \mathbb{R}$. Then for any $k \in \mathbb{N}$, we have $\sum_{i=1}^n p_i a_i^{2k} \geq (\sum_{i=1}^n p_i a_i^2)^k$.*

Proof. We add $p_{n+1} = 1 - (\sum_{i=1}^n p_i)$ and $a_{n+1} = 0$ so p is a probability distribution. Then the claim follows from $\mathbb{E}[X^k] \geq \mathbb{E}[X]^k$, where random variable X gets value a_i^2 with probability p_i . \square

Proof of [Lemma 5.3](#). Extend $\mathbf{v}^{(D'+1)} = \dots = \mathbf{v}^{(D)}$ to equal $\mathbf{v}^{(D')}$. For each $0 \leq i \leq D$, let $X^{(i)} = \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(i)}$. We define $\delta^{(i)} = 0$ for $D' < i \leq D$. For $1 \leq i \leq D'$, we let

$$\delta^{(i)} = X^{(i)} - X^{(i-1)} = \sum_{j=1}^n a_j \cdot (\mathbf{v}_j^{(i)} - \mathbf{v}_j^{(i-1)}) = \sum_{j \in J(v_{i-1})} a_j \cdot \mathbf{v}_j^{(i)},$$

where $J(v_{i-1})$ depends only on $C(v_{i-1})$ since $\mathcal{T}_{C(v_{i-1})}$ performs non-adaptive queries before (and possibly even after) reaching v_i . Note that for the two possible outcomes of querying Q_{v_i} , $\mathbf{v}_j^{(i)}$ is fixed to ± 1 respectively for each $j \in J(v_{i-1})$. Thus $\delta^{(i)} = \Delta^{(i)} \cdot z^{(i)}$ where $\Delta^{(i)}$ is a fixed value given $z^{(1)}, \dots, z^{(i-1)}$ and $z^{(1)}, \dots, z^{(D')}$ are independent unbiased coins in $\{\pm 1\}$.

Since $C(v_{i-1})$ is $2k$ -clean, the collection of random variables $\{\mathbf{v}_j^{(i)} \mid j \in J(v_{i-1})\}$ is $2k$ -wise independent conditioning on $C(v_{i-1})$. Note that δ_i is a linear function and

$$\mathbb{E}[\delta^{(i)} \mid C(v_{i-1})] = 0 \quad \text{and} \quad \mathbb{E}\left[\left(\delta^{(i)}\right)^2 \mid C(v_{i-1})\right] = \sum_{j \in J(v_{i-1})} a_j^2 \leq |J(v_{i-1})|.$$

By the first bound in [Lemma 3.2](#), we have

$$\mathbb{E}\left[\left(\delta^{(i)}\right)^{2k} \mid C(v_{i-1})\right] \leq (2k-1)^k \cdot |J(v_{i-1})|^k. \quad (5)$$

Meanwhile, $|\delta^{(i)}| \leq |J(v_{i-1})|$ always. Our first goal is to bound $\Pr\left[\sum_{i=1}^D (\delta^{(i)})^2 > D + 2\alpha^2 d\right]$. Observe that whenever the event $\sum_{i=1}^D (\delta^{(i)})^2 > D + 2\alpha^2 d$ happens, it must be the case that $\sum_{i:|J(v_{i-1})|>1} (\delta^{(i)})^2 > 2\alpha^2 d$. Thus,

$$\begin{aligned} \Pr\left[\sum_{i=1}^D (\delta^{(i)})^2 > D + 2\alpha^2 d\right] &\leq \Pr\left[\sum_{i:|J(v_{i-1})|>1} (\delta^{(i)})^2 > 2\alpha^2 d\right] \\ &= \Pr\left[\sum_{i:|J(v_{i-1})|>1} \frac{|J(v_{i-1})|}{2d} \cdot \frac{(\delta^{(i)})^2}{|J(v_{i-1})|} > \alpha^2\right] \\ &\leq \Pr\left[\sum_{i:|J(v_{i-1})|>1} \frac{|J(v_{i-1})|}{2d} \cdot \frac{(\delta^{(i)})^{2k}}{|J(v_{i-1})|^k} > \alpha^{2k}\right] \\ &\hspace{15em} \text{(by [Fact 5.4](#) and [Corollary 4.9](#))} \\ &= \Pr\left[\sum_{i:|J(v_{i-1})|>1} \frac{(\delta^{(i)})^{2k}}{|J(v_{i-1})|^{k-1}} > 2d \cdot \alpha^{2k}\right] \\ &\leq \mathbb{E}\left[\sum_{i:|J(v_{i-1})|>1} \frac{(\delta^{(i)})^{2k}}{|J(v_{i-1})|^{k-1}}\right] \cdot \frac{1}{2d \cdot \alpha^{2k}} \quad \text{(by Markov's inequality)} \end{aligned}$$

On the other hand,

$$\begin{aligned}
\mathbb{E} \left[\sum_{i:|J(v_{i-1})|>1} \frac{(\delta^{(i)})^{2k}}{|J(v_{i-1})|^{k-1}} \right] &= \sum_{i=1}^D \mathbb{E}_{C(v_{i-1})} \left[\frac{\mathbf{1}_{|J(v_{i-1})|>1}}{|J(v_{i-1})|^{k-1}} \cdot \mathbb{E} \left[(\delta^{(i)})^{2k} \mid C(v_{i-1}) \right] \right] \\
&\leq \sum_{i=1}^D \mathbb{E}_{C(v_{i-1})} \left[\mathbf{1}_{|J(v_{i-1})|>1} \cdot (2k-1)^k \cdot |J(v_{i-1})| \right] \quad (\text{by (5)}) \\
&= (2k-1)^k \cdot \mathbb{E} \left[\sum_{i:|J(v_{i-1})|>1} |J(v_{i-1})| \right] \\
&\leq (2k-1)^k \cdot 2d. \quad (\text{by Corollary 4.9})
\end{aligned}$$

Overall, we have

$$\Pr \left[\sum_{i=1}^D (\delta^{(i)})^2 > D + 2\alpha^2 d \right] \leq \frac{(2k-1)^k}{\alpha^{2k}}.$$

Then by Lemma 3.4 with $m = 1$, we have

$$\Pr \left[\left| X^{(D)} \right| = \left| \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(D)} \right| \geq \beta \sqrt{2 \cdot (D + 2\alpha^2 d)} \right] \leq 2 \cdot e^{-\beta^2/2} + \frac{(2k-1)^k}{\alpha^{2k}}.$$

The desired bound follows from setting

$$\alpha = \left(\frac{2}{\varepsilon} \right)^{\frac{1}{2k}} \sqrt{2k-1}, \quad \text{and} \quad \beta = \Theta \left(\sqrt{\log \left(\frac{1}{\varepsilon} \right)} \right). \quad \square$$

Now we prove the complete level-1 bound for parity decision trees.

Theorem 5.5. *Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a depth- d parity decision tree. Let $p = \Pr[\mathcal{T}(x) = 1] \in [2^{-d}, 1/2]$.⁹ Then we have*

$$\sum_{j=1}^n \left| \widehat{\mathcal{T}}(j) \right| \leq p \cdot \min \left\{ d, O \left(\sqrt{d} \cdot \log \left(\frac{1}{p} \right) \right) \right\} = O(\sqrt{d}).$$

Proof. For any $i \in [n]$, let $a_i = \text{sgn}(\widehat{\mathcal{T}}(i))$. Now we prove the two bounds separately.

First Bound. Let $v_0, \dots, v_{d'}$ be a random root-to-leaf path in \mathcal{T} . Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(d')} \in \{-1, 0, +1\}^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq d'$ and $j \in [n]$. Since \mathcal{T} is 1-clean in itself, by Lemma 4.6 we have

$$\sum_{j=1}^n \left| \widehat{\mathcal{T}}(j) \right| = \sum_{j=1}^n a_j \cdot \widehat{\mathcal{T}}(j) = \mathbb{E}_{v_0, \dots, v_{d'}} \left[\mathcal{T}(v_{d'}) \cdot \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(d')} \right] \leq \mathbb{E}_{v_0, \dots, v_{d'}} [\mathcal{T}(v_{d'}) \cdot |V|], \quad (6)$$

where $V = \sum_{j=1}^n a_j \cdot \mathbf{v}_j^{(d')}$. Hence by Lemma 5.1, we have (6) $\leq d \cdot \mathbb{E}[\mathcal{T}(v_{d'})] = p \cdot d$.

⁹If $p < 2^{-d}$, then $p = 0$ and $\mathcal{T} \equiv 0$. If $p > 1/2$, we can consider $\widetilde{\mathcal{T}} = 1 - \mathcal{T}$ by symmetry.

Second Bound. By Lemma 4.10, we construct a $2k$ -clean parity decision tree \mathcal{T}' of depth $D \leq 2d \cdot k$ equivalent to \mathcal{T} , where $k = \Theta(\log(1/p))$. Let $U = \sum_{j=1}^n a_j \cdot \mathbf{u}_j^{(D')}$. Then we have

$$\sum_{j=1}^n |\widehat{\mathcal{T}}(j)| = \sum_{j=1}^n |\widehat{\mathcal{T}'}(j)| = \mathbb{E}_{u_0, \dots, u_{D'}} \left[\mathcal{T}'(u_{D'}) \cdot \sum_{j=1}^n a_j \cdot \mathbf{u}_j^{(D')} \right] \leq \mathbb{E}_{u_0, \dots, u_{D'}} [\mathcal{T}'(u_{D'}) \cdot |U|]. \quad (7)$$

Lemma 5.3 implies that for all $\varepsilon > 0$, $\Pr[|U| \geq R(\varepsilon)] \leq \varepsilon$ where

$$R(\varepsilon) = R(D, d, k, \varepsilon) = O \left(\sqrt{dk \cdot \left(\frac{1}{\varepsilon}\right)^{\frac{1}{k}} \cdot \log \left(\frac{1}{\varepsilon}\right)} \right).$$

For integer $i \geq 1$, let $I_i = [R(p/2^i), R(p/2^{i+1})]$ and $I_0 = [0, R(p/2)]$ be intervals. Then for each $i \geq 1$, $\Pr[|U| \in I_i] \leq p/2^i$. We also know that $\mathbb{E}_{u_0, \dots, u_{D'}}[\mathcal{T}'(u_{D'})] \leq p$. Thus,

$$\begin{aligned} (7) &= \mathbb{E}_{u_0, \dots, u_{D'}} \left[\mathcal{T}'(u_{D'}) \cdot |U| \cdot \sum_{i=0}^{+\infty} \mathbf{1}_{|U| \in I_i} \right] \\ &\leq R\left(\frac{p}{2}\right) \cdot \mathbb{E}_{u_0, \dots, u_{D'}}[\mathcal{T}'(u_{D'})] + \sum_{i=1}^{+\infty} R\left(\frac{p}{2^{i+1}}\right) \cdot \mathbb{E}_{u_0, \dots, u_{D'}}[\mathbf{1}_{|U| \in I_i}] \\ &\leq \sum_{i=0}^{+\infty} R\left(\frac{p}{2^{i+1}}\right) \cdot \frac{p}{2^i} \\ &= \sum_{i=0}^{+\infty} O \left(p \cdot \sqrt{dk \cdot \left(\frac{2^{i+1}}{p}\right)^{\frac{1}{k}} \cdot \left(\log\left(\frac{1}{p}\right) + i + 1\right)} \right) \cdot \frac{1}{2^i} \\ &= O \left(p \cdot \sqrt{dk \cdot \log\left(\frac{1}{p}\right)} \right) = O \left(p \cdot \sqrt{d} \cdot \log\left(\frac{1}{p}\right) \right). \quad \square \end{aligned}$$

5.2 Level- ℓ Bound

Now we turn to the general levels.

Lemma 5.6. *There exists a universal constant $\tau \geq 1$ such that the following holds. Let $\ell \geq 1$ be an integer. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a depth- D $2k$ -clean parity decision tree where $k \geq 4 \cdot \ell$ and $n \geq \max\{\tau, k, D\}$ and any root-to-leaf path has at most d nodes that are $2k$ -clean.*

Let $v_0, \dots, v_{D'}$ be a random root-to-leaf path. Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(D')} \in \{-1, 0, +1\}^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq D'$ and $j \in [n]$. Extend $\mathbf{v}^{(D'+1)} = \dots = \mathbf{v}^{(D)}$ to equal $\mathbf{v}^{(D')}$. Then for any sequence $a_S \in \{-1, 0, 1\}$, $S \in \binom{[n]}{\ell}$, any $\varepsilon \leq 1/2$ and $t \in \{0, \dots, \ell\}$, we have

$$\Pr \left[\exists t' \in \{0, \dots, t\}, \exists T \in \binom{[n]}{\ell - t'}, \exists i \in [D], \left| \sum_{S \subseteq T, |S|=t'} a_{S \cup T} \cdot \mathbf{v}_S^{(i)} \right| \geq M(D, d, k, \ell, t', \varepsilon) \right] \leq \varepsilon \cdot t,$$

where we recall that $\mathbf{v}_S^{(i)} = \prod_{j \in S} \mathbf{v}_j^{(i)}$ and where

$$M(D, d, k, \ell, t', \varepsilon) = \left(\tau \cdot (D + dk) \cdot \left(\frac{n^\ell}{\varepsilon}\right)^{\frac{6}{k}} \log\left(\frac{n^\ell}{\varepsilon}\right) \right)^{t'/2}.$$

Proof. We prove the bound by induction on $t = 0, 1, \dots, \ell$ and show $\tau = 10^4$ suffices. The base case $t = 0$ is trivial, since for any fixed T and i , we always have $\left| a_T \cdot \mathbf{v}_\emptyset^{(i)} \right| \leq 1 = M(D, d, k, \ell, 0, \varepsilon)$.

Now we focus on the case where $1 \leq t \leq \ell$. For each $0 \leq i \leq D$ and $T \in \binom{[n]}{\ell-t}$, let

$$X_T^{(i)} = \sum_{S \subseteq \bar{T}, |S|=t} a_{S \cup T} \cdot \mathbf{v}_S^{(i)}.$$

For $1 \leq i \leq D'$, we have

$$\begin{aligned} X_T^{(i)} - X_T^{(i-1)} &= \sum_{S \subseteq \bar{T}, |S|=t, S \cap J(v_{i-1}) \neq \emptyset} a_{S \cup T} \cdot \mathbf{v}_S^{(i)} \\ &= \sum_{r=1}^t \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{V \subseteq \bar{T} \cup J(v_{i-1}), \\ |U|+|V|=t}} a_{T \cup U \cup V} \cdot \mathbf{v}_V^{(i)} \\ &= \sum_{r=1}^t \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{V \subseteq \bar{T} \cup J(v_{i-1}), \\ |U|+|V|=t}} a_{T \cup U \cup V} \cdot \mathbf{v}_V^{(i-1)} \\ &\quad \text{(since } \mathbf{v}_j^{(i)} = \mathbf{v}_j^{(i-1)} \text{ for all } j \notin J(v_{i-1})) \\ &= \underbrace{\sum_{r=1}^t \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{V \subseteq \bar{T} \cup \bar{U}, \\ |U|+|V|=t}} a_{T \cup U \cup V} \cdot \mathbf{v}_V^{(i-1)}}_{A(T, r, i)} \\ &\quad \text{(since } \mathbf{v}_j^{(i-1)} = 0 \text{ for all } j \in J(v_{i-1})) \end{aligned}$$

Observe that conditioning on v_{i-1} ,

- if r is an even number, then $A(T, r, i)$ is a fixed value independent of $\mathbf{v}^{(i)}$;
- if r is an odd number, then $A(T, r, i)$ is an unbiased coin with magnitude independent of $\mathbf{v}^{(i)}$.

Therefore, trying to apply [Lemma 3.4](#), we write $X_T^{(i)} - X_T^{(i-1)} = \mu_T^{(i)} + \Delta_T^{(i)} \cdot z_T^{(i)}$, where $z_T^{(1)}, \dots, z_T^{(D)}$ are independent unbiased coins in $\{\pm 1\}$ and $\mu_T^{(i)} = \Delta_T^{(i)} = 0$ for $D' < i \leq D$ and

$$\mu_T^{(i)} = \sum_{\substack{r=2, \\ \text{even}}}^t A(T, r, i) \quad \text{and} \quad \Delta_T^{(i)} = \left| \sum_{\substack{r=1, \\ \text{odd}}}^t A(T, r, i) \right| \quad \text{for } 1 \leq i \leq D'. \quad (8)$$

First Bound on $A(T, r, i)$. Let \mathcal{E}_1 be the following event:

$$\mathcal{E}_1 = \text{“ } \exists \hat{t} \in \{0, \dots, t-1\}, \exists T' \in \binom{[n]}{\ell-\hat{t}}, \exists i' \in [D], \left| X_{T'}^{(i')} \right| \geq M(D, k, \ell, \hat{t}, \varepsilon) \text{”}.$$

By the induction hypothesis, we have

$$\Pr[\mathcal{E}_1] \leq (t-1) \cdot \varepsilon. \quad (9)$$

We first derive a simple bound, that will be effective for small values of $|J(v_{i-1})|$.

Claim 5.7. When \mathcal{E}_1 does not happen, $|A(T, r, i)| \leq |J(v_{i-1})|^r \cdot M(D, d, k, \ell, t - r, \varepsilon)$ holds for all $r \in [t], i \in [D], T \in \binom{[n]}{\ell-t}$.

Proof. Since \mathcal{E}_1 does not happen, by union bound we have

$$\begin{aligned} |A(T, r, i)| &= \left| \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{V \subseteq \bar{T} \cup \bar{U}, \\ |U|+|V|=t}} a_{T \cup U \cup V} \cdot \mathbf{v}_V^{(i-1)} \right| \leq |J(v_{i-1})|^r \max_{U \subseteq \bar{T}, |U|=r} |X_{T \cup U}^{(i-1)}| \\ &\leq |J(v_{i-1})|^r \cdot M(D, d, k, \ell, t - r, \varepsilon). \quad \square \end{aligned}$$

Second Bound on $A(T, r, i)$. The second bound requires a more refined decomposition on $A(T, r, i)$.

Assume that $c(i-1)$ is the index of $C(v_{i-1})$ in $v_0, \dots, v_{D'}$, i.e., $v_{c(i-1)} = C(v_{i-1})$. This means that $v_{c(i-1)}$ is the closest ancestor to v_{i-1} that is $2k$ -clean. Then define

$$L(v_{i-1}) = \bigcup_{c(i-1) \leq i' < i-1} J(v_{i'}).$$

The elements of $L(v_{i-1})$ are precisely the coordinates fixed by the queries from $Q_{v_{c(i-1)}}$ to $Q_{v_{i-1}}$, excluding the latter. Since $\mathcal{T}_{C(v_{i-1})}$ makes non-adaptive queries before (and possibly even after) reaching v_i , $L(v_{i-1})$ and $J(v_{i-1})$ depend only on $C(v_{i-1})$ and i . We now expand $A(T, r, i)$ by also grouping terms based on the number of coordinates in $L(v_{i-1})$ as follows:

$$\begin{aligned} A(T, r, i) &= \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{V \subseteq \bar{T} \cup \bar{U}, \\ |U|+|V|=t}} a_{T \cup U \cup V} \cdot \mathbf{v}_V^{(i-1)} \\ &= \sum_{r'=0}^{t-r} \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{W \subseteq L(v_{i-1}) \cap \bar{T}, \\ |W|=r'}} \mathbf{v}_W^{(i-1)} \sum_{\substack{W' \subseteq \bar{T} \cup \bar{U} \cup L(v_{i-1}), \\ |W'|=t-r-r'}} a_{T \cup U \cup W \cup W'} \cdot \mathbf{v}_{W'}^{(i-1)} \\ &= \sum_{r'=0}^{t-r} \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{W \subseteq L(v_{i-1}) \cap \bar{T}, \\ |W|=r'}} \mathbf{v}_W^{(i-1)} \sum_{\substack{W' \subseteq \bar{T} \cup \bar{U} \cup L(v_{i-1}), \\ |W'|=t-r-r'}} a_{T \cup U \cup W \cup W'} \cdot \mathbf{v}_{W'}^{c(i-1)} \\ &\quad (\text{since } \mathbf{v}_j^{(i-1)} = \mathbf{v}_j^{c(i-1)} \text{ for all } j \notin L(v_{i-1})) \\ &= \sum_{r'=0}^{t-r} \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{W \subseteq L(v_{i-1}) \cap \bar{T}, \\ |W|=r'}} \mathbf{v}_W^{(i-1)} \sum_{\substack{W' \subseteq \bar{T} \cup \bar{U} \cup W, \\ |W'|=t-r-r'}} a_{T \cup U \cup W \cup W'} \cdot \mathbf{v}_{W'}^{c(i-1)} \\ &\quad (\text{since } \mathbf{v}_j^{c(i-1)} = 0 \text{ for all } j \in L(v_{i-1})) \\ &= \underbrace{\sum_{r'=0}^{t-r} \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \mathbf{v}_U^{(i)} \sum_{\substack{W \subseteq L(v_{i-1}) \cap \bar{T}, \\ |W|=r'}} \mathbf{v}_W^{(i-1)}}_{\Gamma_T^{(i)}(r, r')} \cdot X_{T \cup U \cup W}^{c(i-1)}. \end{aligned}$$

Since $C(v_{i-1})$ is $2k$ -clean, by [Fact 4.2](#), the collection of random variables

$$\left\{ \mathbf{v}_j^{(i)} \mid j \in J(v_{i-1}) \right\} \cup \left\{ \mathbf{v}_j^{(i-1)} \mid j \in L(v_{i-1}) \right\}$$

is $2k$ -wise independent conditioning on $C(v_{i-1})$. Note that $\Gamma_T^{(i)}(r, r')$ is a polynomial of degree at most $r + r' \leq \ell < k$, that $\mathbb{E} \left[\Gamma_T^{(i)}(r, r') \mid C(v_{i-1}) \right] = 0$, and

$$\begin{aligned} \sigma_T^2(r, r', C(v_{i-1}), i) &:= \mathbb{E} \left[\left(\Gamma_T^{(i)}(r, r') \right)^2 \mid C(v_{i-1}) \right] = \sum_{\substack{U \subseteq J(v_{i-1}) \cap \bar{T}, \\ |U|=r}} \sum_{\substack{W \subseteq L(v_{i-1}) \cap \bar{T}, \\ |W|=r'}} \left(X_{TUUVW}^{c(i-1)} \right)^2 \\ &\leq (|J(v_{i-1})|)^r (|L(v_{i-1})|)^{r'} \left(\max_{|T'|=r+r'+\ell-t, i' \in [D]} \left| X_{T'}^{(i')} \right| \right)^2 \\ &\leq (|J(v_{i-1})|)^r D^{r'} \left(\max_{|T'|=r+r'+\ell-t, i' \in [D]} \left| X_{T'}^{(i')} \right| \right)^2. \end{aligned}$$

(since $|L(v_{i-1})| \leq D$ by [Fact 4.8](#))

We also have the following claim, the proof of which follows from [Lemma 3.2](#) applied to the low degree polynomial $\Gamma_T^{(i)}$. The proof is deferred to [Appendix C](#).

Claim 5.8. **Pr** $[\mathcal{E}_2] \leq \varepsilon/3$, where \mathcal{E}_2 is the following event:

$$\text{“ } \exists T \in \binom{[n]}{\ell-t}, i, r, r', \left| \Gamma_T^{(i)}(r, r') \right| \geq \left(100 \min \left\{ k, \log \left(\frac{n^\ell}{\varepsilon} \right) \right\} \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \right)^{\frac{r+r'}{2}} \cdot \sigma_T(r, r', C(v_{i-1}), i) \text{”}.$$

On the other hand, when $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen, the following calculation holds for all $T \in \binom{[n]}{\ell-t}, i \in [D'], r \in [t], 0 \leq r' \leq t-r$:

$$\begin{aligned} \left| \Gamma_T^{(i)}(r, r') \right| &\leq M(D, k, \ell, t-r-r', \varepsilon) \cdot \sqrt{\left(100 \min \left\{ k, \log \left(\frac{n^\ell}{\varepsilon} \right) \right\} \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \right)^{r+r'} (|J(v_{i-1})|)^r \cdot D^{r'}} \\ &\leq M(D, k, \ell, t-r-r', \varepsilon) \cdot \sqrt{\left(100 \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \right)^{r+r'} (|J(v_{i-1})| \cdot k)^r \cdot \left(D \cdot \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{r'}} \\ &= \sqrt{\left(\tau(D+dk) \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{t-r-r'} \left(100 \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \right)^{r+r'} (|J(v_{i-1})| \cdot k)^r \left(D \cdot \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{r'}} \\ &\leq \sqrt{\left(\tau(D+dk) \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^t \left(\frac{100}{\tau} \right)^{r+r'} \left(\frac{|J(v_{i-1})|}{d \cdot \log(n^\ell/\varepsilon)} \right)^r} \\ &\leq \sqrt{\left(\tau(D+dk) \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^t \left(\frac{200}{\tau} \right)^{r+r'} \left(\frac{|J(v_{i-1})|}{2d} \right)^r \frac{1}{\log(n^\ell/\varepsilon)}} \\ &= M(D, d, k, \ell, t, \varepsilon) \cdot \sqrt{\left(\frac{200}{\tau} \right)^{r+r'} \left(\frac{|J(v_{i-1})|}{2d} \right)^r \frac{1}{\log(n^\ell/\varepsilon)}}. \end{aligned}$$

Hence we have a second bound on $A(T, r, i)$.

Claim 5.9. When $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen, the following holds for all $r \in [t], i \in [D], T \in \binom{[n]}{\ell-t}$:

$$\left| A(T, r, i) \right| \leq \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\left(\frac{800}{\tau} \right)^r \left(\frac{|J(v_{i-1})|}{2d} \right)^r}.$$

Proof. Since $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen, by union bound and noticing $\tau \geq 800$ we have

$$\begin{aligned} |A(T, r, i)| &\leq \sum_{r'=0}^{t-r} \left| \Gamma_T^{(i)}(r, r') \right| \leq \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\left(\frac{200}{\tau}\right)^r \left(\frac{|J(v_{i-1})|}{2d}\right)^r} \cdot \sum_{r'=0}^{+\infty} \left(\frac{200}{\tau}\right)^{r'/2} \\ &\leq \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\left(\frac{800}{\tau}\right)^r \left(\frac{|J(v_{i-1})|}{2d}\right)^r}. \quad \square \end{aligned}$$

Final Bound on $\mu_T^{(i)}$ and $\delta_T^{(i)}$. Combining [Claim 5.7](#) and [Claim 5.9](#), if $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen we have

$$|A(T, r, i)| \leq M(D, d, k, \ell, t-r, \varepsilon) + \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\left(\frac{800}{\tau}\right)^r \left(\frac{|J(v_{i-1})|}{2d}\right)^r} \cdot \mathbf{1}_{|J(v_{i-1})| > 1} \quad (10)$$

To see this, if $|J(v_{i-1})| \leq 1$, we use the bound from [Claim 5.7](#) as the first term in (10). Otherwise $|J(v_{i-1})| > 1$, in which case we use the bound from [Claim 5.9](#) as the second term in (10).

By [Corollary 4.9](#), we can now bound $\sum_{i=1}^D \left| \mu_T^{(i)} \right|$ and $\sum_{i=1}^D \left| \Delta_T^{(i)} \right|^2$ as [Claim 5.10](#). Its proof is deferred in [Appendix D](#).

Claim 5.10. When $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen, $\sum_{i=1}^D \left| \mu_T^{(i)} \right| \leq R$ and $\sum_{i=1}^D \left| \Delta_T^{(i)} \right|^2 \leq R^2$ hold for all $T \in \binom{[n]}{\ell-t}$, where

$$R = \frac{M(D, d, k, \ell, t, \varepsilon)}{5 \cdot \sqrt{\log(n^\ell/\varepsilon)}}. \quad (11)$$

Complete Induction. Let $\beta = \sqrt{2 \cdot \log(n^\ell/\varepsilon)} \geq 1$ and observe that

$$\begin{aligned} R + \beta \cdot \sqrt{2} \cdot R &\leq \beta \cdot 2\sqrt{2} \cdot R && \text{(due to } \beta \geq 1) \\ &= \frac{2\sqrt{2} \cdot \sqrt{2 \cdot \log(n^\ell/\varepsilon)}}{5 \cdot \sqrt{\log(n^\ell/\varepsilon)}} \cdot M(D, d, k, \ell, t, \varepsilon) && \text{(due to (11))} \\ &\leq M(D, d, k, \ell, t, \varepsilon). \end{aligned}$$

Then we have

$$\begin{aligned} &\Pr \left[\exists t' \in \{0, \dots, t\}, \exists T' \in \binom{[n]}{\ell-t'}, \exists i \in [D], \left| X_{T'}^{(i)} \right| \geq M(D, d, k, \ell, t', \varepsilon) \right] \\ &= \Pr \left[\mathcal{E}_1 \vee \left(\exists T \in \binom{[n]}{\ell-t}, \exists i \in [D], \left| X_T^{(i)} \right| \geq M(D, d, k, \ell, t, \varepsilon) \right) \right] \\ &\leq \Pr \left[(\mathcal{E}_1 \vee \mathcal{E}_2) \vee \left(\exists T \in \binom{[n]}{\ell-t}, \exists i \in [D], \left| X_T^{(i)} \right| \geq R + \beta \cdot \sqrt{2} \cdot R \right) \right] \\ &\leq (t-1) \cdot \varepsilon + \frac{\varepsilon}{3} + 2n^{\ell-t} \cdot e^{-\beta^2/2} && \text{(due to (9), Claim 5.8, Lemma 3.4, and Claim 5.10)} \\ &\leq (t-1) \cdot \varepsilon + \frac{\varepsilon}{3} + \frac{1}{3} \cdot n^\ell \cdot e^{-\beta^2/2} \\ &\leq t \cdot \varepsilon. \quad \square \end{aligned}$$

Before we prove the complete level- ℓ bound for parity decision trees, we first prove a simple bound for the number of vectors with a given weight in a subspace.

Lemma 5.11. *Let $\ell \geq 1$ be an integer and \mathcal{S} be a subspace of rank at most d . Let $U = \{S \mid |S| = \ell, S \in \mathcal{S}\}$, then $|U| \leq \min \left\{ \binom{d \cdot \ell}{\ell}, 2^d - 1 \right\}$.*

Proof. Let $\{S_1, \dots, S_{d'}\}$ be a maximal set of independent vectors in U . Then $d' \leq d$ and $|S_i| = \ell$ holds for all $i \in [d']$. Since $U \subseteq \text{Span} \langle S_1, \dots, S_{d'} \rangle$ and $\emptyset \notin U$, we have

$$|U| \leq |\text{Span} \langle S_1, \dots, S_{d'} \rangle| - 1 = 2^{d'} - 1 \leq 2^d - 1.$$

On the other hand, observe that $U \subseteq \binom{S_1 \cup \dots \cup S_{d'}}{\ell}$, hence we also have

$$|U| \leq \left| \binom{S_1 \cup \dots \cup S_{d'}}{\ell} \right| \leq \binom{d' \cdot \ell}{\ell} \leq \binom{d \cdot \ell}{\ell}. \quad \square$$

We remark that in [Lemma 5.11](#), it is conjectured the bound should be $\binom{d+1}{\ell}$ when $d \geq 2 \cdot \ell$ [[Kra10, BP18](#)].

Theorem 5.12. *Let $\ell \geq 1$ be an integer. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a depth- d parity decision tree where $n \geq \max \{d, \ell\}$. Let $p = \Pr[\mathcal{T}(x) = 1] \geq 2^{-d}$.¹⁰ Then we have*

$$\sum_{S \subseteq [n]: |S|=\ell} |\widehat{\mathcal{T}}(S)| \leq p \cdot \min \left\{ \binom{d \cdot \ell}{\ell}, 2^d - 1, O\left(\sqrt{d} \cdot \log\left(\frac{n^\ell}{p}\right)\right)^\ell \right\} = O\left(\sqrt{d} \cdot \ell \cdot \log(n)\right)^\ell.$$

Proof. For any $S \in \binom{[n]}{\ell}$, let $a_S = \text{sgn}(\widehat{\mathcal{T}}(S))$. Now we prove the bounds separately.

First Two Bounds. Let $v_0, \dots, v_{d'}$ be a random root-to-leaf path. Then by the definition of $\widehat{\mathcal{P}}_v$ and \mathcal{S}_v and [Fact 4.2](#), we have

$$\begin{aligned} \sum_S |\widehat{\mathcal{T}}(S)| &= \sum_S a_S \cdot \widehat{\mathcal{T}}(S) = \mathbb{E}_{v_0, \dots, v_{d'}} \left[\mathcal{T}(v_{d'}) \cdot \sum_S a_S \cdot \widehat{\mathcal{P}}_{v_{d'}}(S) \right] \\ &\leq \mathbb{E}_{v_0, \dots, v_{d'}} \left[\mathcal{T}(v_{d'}) \cdot \sum_S |\widehat{\mathcal{P}}_{v_{d'}}(S)| \right] = \mathbb{E}_{v_0, \dots, v_{d'}} [\mathcal{T}(v_{d'}) \cdot |V|], \end{aligned} \quad (12)$$

where $a_S = \text{sgn}(\widehat{\mathcal{T}}(S))$ and $V = \left\{ S \in \binom{[n]}{\ell} \mid S \in \mathcal{S}_{v_{d'}} \right\}$. Note that

$$\text{rank}(\mathcal{S}_{v_{d'}}) = \text{rank}(\text{Span} \langle Q_{v_0}, \dots, Q_{v_{d'-1}} \rangle) \leq d' \leq d.$$

Hence by [Lemma 5.11](#), we have (12) $\leq \min \left\{ \binom{d \cdot \ell}{\ell}, 2^d - 1 \right\} \cdot \mathbb{E}[\mathcal{T}(v_{d'})] = p \cdot \min \left\{ \binom{d \cdot \ell}{\ell}, 2^d - 1 \right\}$.

Third Bound. By [Lemma 4.10](#), we construct a $2k$ -clean parity decision tree \mathcal{T}' of depth $D \leq 2d \cdot k$ equivalent to \mathcal{T} , where $k = \Theta(\log(n^\ell/p)) \geq 4 \cdot \ell$. We also add dummy variables to make sure $n' = \max \{\tau, k, 6D, n\}$, where \mathcal{T}' has n' inputs and τ is the universal constant in [Lemma 5.6](#).

Let $u_0, \dots, u_{D'}$ be a random root-to-leaf path in \mathcal{T}' . Define $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(D')} \in \{-1, 0, +1\}^{n'}$ by setting $\mathbf{u}_j^{(i)} = \widehat{\mathcal{P}}_{u_i}(j)$ for each $0 \leq i \leq D'$ and $j \in [n]$. Then extend $\mathbf{u}^{(D'+1)} = \mathbf{u}^{(D'+2)} = \dots = \mathbf{u}^{(D)}$ to equal $\mathbf{u}^{(D')}$. By [Lemma 4.6](#), we have

$$\sum_S |\widehat{\mathcal{T}}(S)| = \sum_S |\widehat{\mathcal{T}'}(S)| = \mathbb{E}_{u_0, \dots, u_{D'}} \left[\mathcal{T}'(u_{D'}) \cdot \sum_S a_S \cdot \mathbf{u}_S^{(D')} \right] \leq \mathbb{E}_{u_0, \dots, u_{D'}} [\mathcal{T}'(u_{D'}) \cdot |U|], \quad (13)$$

¹⁰If $p < 2^{-d}$, then $p = 0$ and $\mathcal{T} \equiv 0$.

where $U = \sum_S a_S \cdot \mathbf{u}_S^{(D)}$.

Now we apply [Lemma 5.6](#) with $t = \ell, \varepsilon = \Theta(p/d^{\ell/2}) \leq 1/2$ to obtain the following bound¹¹

$$M = M(D, d, k, \ell, \varepsilon) = \left(O\left(\sqrt{d} \cdot \log\left(\frac{n^\ell}{p} \right) \right) \right)^\ell$$

such that $\Pr[|U| \geq M] \leq \ell \cdot \varepsilon$. Then, combining the first bound, we have

$$\begin{aligned} (13) &= \mathbb{E}[\mathcal{T}(u_{D'}) \cdot |U| \cdot (\mathbf{1}_{|U| < M} + \mathbf{1}_{|U| \geq M})] \leq M \cdot \mathbb{E}[\mathcal{T}(u_{D'})] + \ell \cdot \varepsilon \cdot \binom{d \cdot \ell}{\ell} \\ &= p \cdot \left(O\left(\sqrt{d} \cdot \log\left(\frac{n^\ell}{p} \right) \right) \right)^\ell, \end{aligned}$$

which is maximized at $p = 1$, hence (13) = $O\left(\sqrt{d} \cdot \ell \cdot \log(n)\right)^\ell$ as desired. \square

6 Fourier Bounds for Noisy Decision Trees

Let \mathcal{T} be a noisy decision tree. By adding queries with zero correlation, we assume without loss of generality each root-to-leaf path in the noisy decision tree is of the same length. Let v be any node of \mathcal{T} . We use \mathcal{P}_v to denote the uniform distribution over $\{\pm 1\}^n$ conditioning on reaching v . Note that \mathcal{P}_v is always a *product distribution*. As before, for any $S \subseteq [n]$ we define $\widehat{\mathcal{P}}_v(S) = \mathbb{E}_{x \sim \mathcal{P}_v}[x_S]$.

Claim 6.1. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a cost- d noisy decision tree. Let v_0, \dots, v_D be any root-to-leaf path in \mathcal{T} . Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(D)} \in [-1, 1]^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq D$ and $j \in [n]$. Then for any $i \in \{0, \dots, D-1\}$, $\mathbf{v}_{q_{v_i}}^{(i+1)} - \mathbf{v}_{q_{v_i}}^{(i)}$ is a mean-zero random variable with magnitude bounded by $2 \cdot |\gamma_{v_i}|$.

Proof. Fix $i \in \{0, \dots, D-1\}$. For convenience, let $j = q_{v_i}$, $\gamma = \gamma_{v_i}$, and $\alpha = \mathbf{v}_j^{(i)}$. Suppose $|\gamma| = 1$ then $\left| \mathbf{v}_j^{(i+1)} - \mathbf{v}_j^{(i)} \right| \leq 2 = 2 \cdot |\gamma_{v_i}|$ as desired. Now we turn to the case $|\gamma| < 1$.

Note that for the distribution \mathcal{P}_{v_i} , the measure of $x_j = 1$ (resp., $x_j = -1$) inputs is $(1 + \alpha)/2$ (resp., $(1 - \alpha)/2$). The measure of $x_j = 1$ (resp., $x_j = -1$) inputs that follow the edge labeled 1 is $a := (1 + \alpha)(1 + \gamma)/4$ (resp., $b := (1 - \alpha)(1 - \gamma)/4$). The total measure of inputs that take the edge labeled 1 is $a + b$ and the resulting node v_{i+1} satisfies $\mathbf{v}_j^{(i+1)} = (a - b)/(a + b)$. This implies that

$$\mathbf{v}_j^{(i+1)} = \begin{cases} \frac{\alpha + \gamma}{1 + \gamma \cdot \alpha} & \text{with probability } \frac{1 + \gamma \cdot \alpha}{2}, \\ \frac{\alpha - \gamma}{1 - \gamma \cdot \alpha} & \text{with probability } \frac{1 - \gamma \cdot \alpha}{2}. \end{cases}$$

The above calculation implies

$$\mathbf{v}_j^{(i+1)} - \mathbf{v}_j^{(i)} = \begin{cases} \gamma \cdot \frac{1 - \alpha^2}{1 + \gamma \cdot \alpha} & \text{with probability } \frac{1 + \gamma \cdot \alpha}{2}, \\ -\gamma \cdot \frac{1 - \alpha^2}{1 - \gamma \cdot \alpha} & \text{with probability } \frac{1 - \gamma \cdot \alpha}{2}, \end{cases}$$

and thus $\mathbf{v}_j^{(i+1)} - \mathbf{v}_j^{(i)}$ is a mean-zero random variable. Since $\alpha \in [-1, 1]$ and $\gamma \in (-1, 1)$, we have

$$\max \left\{ \frac{1 - \alpha^2}{1 - \gamma \cdot \alpha}, \frac{1 - \alpha^2}{1 + \gamma \cdot \alpha} \right\} \leq \frac{1 - \alpha^2}{1 - |\alpha|} = 1 + |\alpha| \leq 2,$$

which implies $\left| \mathbf{v}_j^{(i+1)} - \mathbf{v}_j^{(i)} \right| \leq 2 \cdot |\gamma|$. \square

¹¹Since $n \geq \max\{\ell, d\}$, we know $k = \Theta(\log(n^\ell/p)) = O(n^2)$ and $D \leq 2d \cdot k = O(n^3)$. Hence $n' = \max\{\tau, k, 6D, n\} = O(n^3)$. Also $n^\ell/\varepsilon \leq n^{O(\ell)}/p$ and by our choice of $k = \Theta(\log(n^\ell/p))$ we have $(n^\ell/\varepsilon)^{6/k} = O(1)$.

We now prove the general Fourier bounds. As before, for any $S \subseteq [n]$, let $\mathbf{v}_S^{(i)}$ be $\prod_{j \in S} \mathbf{v}_j^{(i)}$.

Lemma 6.2. *There exists a universal constant τ such that the following holds. Let $\ell \geq 1$ be an integer. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a cost- d noisy decision tree.*

Let v_0, \dots, v_D be a random root-to-leaf path in \mathcal{T} . Define $\mathbf{v}^{(0)}, \dots, \mathbf{v}^{(D)} \in [-1, 1]^n$ by setting $\mathbf{v}_j^{(i)} = \widehat{\mathcal{P}}_{v_i}(j)$ for each $0 \leq i \leq D$ and $j \in [n]$. Then for any sequence $a_S \in \{-1, 0, 1\}$, $S \in \binom{[n]}{\ell}$, any $\varepsilon \leq 1/2$ and $t \in \{0, \dots, \ell\}$, we have

$$\Pr \left[\exists T \in \binom{[n]}{\ell-t}, \exists i \in [D], \left| \sum_{S \subseteq \bar{T}, |S|=t} a_{S \cup T} \cdot \mathbf{v}_S^{(i)} \right| \geq S(d, \ell, t, \varepsilon) \right] \leq \varepsilon \cdot t,$$

where $S(d, \ell, 0, \varepsilon) = 1$ and

$$S(d, \ell, t, \varepsilon) = \sqrt{(\tau \cdot d)^t \cdot \log\left(\frac{n^{\ell-t}}{\varepsilon}\right) \cdots \log\left(\frac{n^{\ell-1}}{\varepsilon}\right)} \quad \text{for } t \in [\ell].$$

Proof. We prove the bound by induction on t and show $\tau = 32$ suffices. The base case $t = 0$ is trivial, since for any T of size ℓ and any i , we have $|a_T \cdot \mathbf{v}_\emptyset^{(i)}| \leq 1 = S(d, \ell, 0, \varepsilon)$.

Now we focus on the case $1 \leq t \leq \ell$. For any $T \in \binom{[n]}{\leq \ell}$, define $X_T^{(0)}, \dots, X_T^{(D)}$ by $X_T^{(i)} = \sum_{S \subseteq \bar{T}, |S|+|T|=\ell} a_{S \cup T} \cdot \mathbf{v}_S^{(i)}$. Define $\delta_T^{(i)}$ for $i \in [D]$ as follows:

$$\begin{aligned} \delta_T^{(i)} &= X_T^{(i)} - X_T^{(i-1)} = \sum_{S \subseteq \bar{T}, |S|=t, S \ni q_{v_{i-1}}} a_{S \cup T} \cdot \left(\mathbf{v}_S^{(i)} - \mathbf{v}_S^{(i-1)} \right) \\ &= \left(\mathbf{v}_{q_{v_{i-1}}}^{(i)} - \mathbf{v}_{q_{v_{i-1}}}^{(i-1)} \right) \cdot \sum_{S' \subseteq \overline{T \cup \{q_{v_{i-1}\}}, |S'|=t-1} a_{S' \cup \{q_{v_{i-1}\}} \cup T} \cdot \mathbf{v}_{S'}^{(i-1)} \\ &= \left(\mathbf{v}_{q_{v_{i-1}}}^{(i)} - \mathbf{v}_{q_{v_{i-1}}}^{(i-1)} \right) \cdot X_{T \cup \{q_{v_{i-1}\}}^{(i-1)}}. \end{aligned}$$

Note that by [Claim 6.1](#) and conditioning on v_{i-1} , $\delta_T^{(i)}$ is a mean-zero random variable.

The induction hypothesis implies that with all but $\varepsilon \cdot (t-1)$ probability, for all $i \in [D]$ and $T' \in \binom{[n]}{\ell-t+1}$, we have $|X_{T'}^{(i)}| \leq S(d, \ell, t-1, \varepsilon)$. By [Claim 6.1](#), we have

$$\left| \delta_T^{(i)} \right| = \left| \mathbf{v}_{q_{v_{i-1}}}^{(i)} - \mathbf{v}_{q_{v_{i-1}}}^{(i-1)} \right| \cdot \left| X_{T \cup \{q_{v_{i-1}\}}^{(i-1)}} \right| \leq 2 \cdot |\gamma_{v_{i-1}}| \cdot S(d, \ell, t-1, \varepsilon).$$

Denote by $\Delta_T^{(i)} = 2 \cdot |\gamma_{v_{i-1}}| \cdot S(d, \ell, t-1, \varepsilon)$. We can thus express $X_T^{(i)} = X_T^{(i-1)} + \Delta_T^{(i)} \cdot z_T^{(i)}$ where $|z_T^{(i)}| \leq 1$. Then we apply [Lemma 3.4](#) to the family of martingales $X_T^{(0)}, \dots, X_T^{(D)}$, $|T| \in \binom{[n]}{\ell-t}$ with difference sequence $\delta_T^{(i)} = \Delta_T^{(i)} \cdot z_T^{(i)}$ satisfying

$$\sum_{i=1}^D \left(\Delta_T^{(i)} \right)^2 = 4 \cdot (S(d, \ell, t-1, \varepsilon))^2 \cdot \sum_{i=1}^D |\gamma_{v_{i-1}}|^2 \leq 4d \cdot (S(d, \ell, t-1, \varepsilon))^2.$$

Hence for any $\beta \geq 0$, we have

$$\Pr \left[\exists T \in \binom{[n]}{\ell-t}, \exists i \in [D], \left| X_T^{(i)} \right| \geq 2\beta \cdot \sqrt{2d} \cdot S(d, \ell, t-1, \varepsilon) \right] \leq \varepsilon \cdot (t-1) + 2 \cdot n^{\ell-t} \cdot e^{-\beta^2/2}.$$

Since $\varepsilon \leq 1/2$, we can set $\beta = 2 \cdot \sqrt{\log(n^{\ell-t}/\varepsilon)}$ so that $2 \cdot n^{\ell-t} \cdot e^{-\beta^2/2} \leq \varepsilon$, which completes the induction by noticing

$$2\beta \cdot \sqrt{2d} \cdot S(d, \ell, t-1, \varepsilon) = \sqrt{32 \cdot d \cdot \log\left(\frac{n^{\ell-t}}{\varepsilon}\right)} \cdot S(d, \ell, t-1, \varepsilon) \leq S(d, \ell, t, \varepsilon). \quad \square$$

Theorem 6.3. *Let $\ell \geq 1$ and $n \geq \max\{\ell, 2\}$ be integers. Let $\mathcal{T}: \{\pm 1\}^n \rightarrow \{0, 1\}$ be a cost- d noisy decision tree. Let $p = \Pr[\mathcal{T}(x) = 1] \in (0, 1/2]$.¹² Then we have*

$$\sum_{S \subseteq [n], |S|=\ell} |\widehat{\mathcal{T}}(S)| \leq p \cdot O(d)^{\ell/2} \cdot \sqrt{\log\left(\frac{1}{p}\right) \left(\log\left(\frac{n^\ell}{p}\right)\right)^{\ell-1}} = O(d)^{\ell/2} \cdot \sqrt{1 + (\ell \log(n))^{\ell-1}}.$$

Proof. For any $S \in \binom{[n]}{\ell}$, let $a_S = \text{sgn}(\widehat{\mathcal{T}}(S))$. Let v_0, \dots, v_D be a random root-to-leaf path in \mathcal{T} . Note that

$$\sum_S |\widehat{\mathcal{T}}(S)| = \sum_S a_S \cdot \widehat{\mathcal{T}}(S) = \mathbb{E} \left[\mathcal{T}(v_D) \cdot \sum_S a_S \cdot \mathbf{v}_S^{(D)} \right] \leq \mathbb{E}[\mathcal{T}(v_D) \cdot |V|], \quad (14)$$

where $V = \sum_S a_S \cdot \mathbf{v}_S^{(D)}$. By Lemma 6.2, we know $\Pr[|V| \geq S(\varepsilon)] \leq \varepsilon \cdot \ell$, where

$$S(\varepsilon) = S(d, \ell, \ell, \varepsilon) = \sqrt{O(d)^\ell \cdot \log\left(\frac{n^{\ell-1}}{\varepsilon}\right) \cdots \log\left(\frac{n^0}{\varepsilon}\right)} \leq \sqrt{O(d)^\ell \cdot \left(\log\left(\frac{n^{\ell-1}}{\varepsilon}\right)\right)^{\ell-1} \cdot \log\left(\frac{1}{\varepsilon}\right)}.$$

For integer $i \geq 1$, let $I_i = [S(p/(\ell 2^i)), S(p/(\ell 2^{i+1}))]$ and $I_0 = [0, S(p/\ell)]$ be intervals. Then for each $i \geq 1$, $\Pr[|V| \in I_i] \leq p/2^i$. We also know that $\mathbb{E}_{v_0, \dots, v_D}[\mathcal{T}(v_D)] \leq p$. Thus,

$$\begin{aligned} (14) &\leq \mathbb{E}_{v_0, \dots, v_D} \left[\mathcal{T}(v_D) \cdot |V| \cdot \sum_{i=0}^{+\infty} \mathbf{1}_{|V| \in I_i} \right] \\ &\leq S\left(\frac{p}{\ell}\right) \cdot \mathbb{E}[\mathcal{T}(v_D)] + \sum_{i=1}^{+\infty} S\left(\frac{p}{\ell \cdot 2^{i+1}}\right) \cdot \mathbb{E}[\mathbf{1}_{|V| \in I_i}] \\ &\leq \sum_{i=0}^{+\infty} S\left(\frac{p}{\ell \cdot 2^{i+1}}\right) \cdot \frac{p}{2^i} \\ &= \sum_{i=0}^{+\infty} p \cdot \sqrt{O(d)^\ell \cdot \left(\log\left(\frac{n^{\ell-1} \cdot \ell}{p}\right) + i + 1\right)^{\ell-1} \cdot \left(\log\left(\frac{1}{p}\right) + \log(\ell) + i + 1\right)} \cdot \frac{1}{2^i} \\ &\leq \sum_{i=0}^{+\infty} p \cdot \sqrt{O(d)^\ell \cdot \left(\left(\log\left(\frac{n^\ell}{p}\right)\right)^{\ell-1} + (i+1)^{\ell-1}\right) \cdot \left(\log\left(\frac{1}{p}\right) + i + 1\right)} \cdot \frac{1}{2^i} \\ &\quad (\text{since } n \geq \ell, \text{ and } (x+y)^b \leq 2^b \cdot (x^b + y^b) \text{ and } \sqrt{x+y} \leq \sqrt{x} + \sqrt{y} \text{ for } x, y, b \geq 0) \\ &\leq p \cdot \sqrt{O(d)^\ell \cdot \log\left(\frac{1}{p}\right) \left(\log\left(\frac{n^\ell}{p}\right)\right)^{\ell-1}}, \end{aligned}$$

where the last inequality follows from $p \leq 1/2$, $n \geq 2$ and

$$\sum_{i=0}^{+\infty} (i+1)^{\ell/2} \cdot 2^{-i} = O(\ell)^{\ell/2} \leq O(1)^\ell \cdot \ell^{(\ell-1)/2} \leq O(1)^\ell \cdot \left(\log\left(\frac{n^\ell}{p}\right)\right)^{(\ell-1)/2}.$$

¹²If $p > 1/2$, then we can consider $\widetilde{\mathcal{T}} = 1 - \mathcal{T}$ by symmetry.

Note that $p \cdot (\log(1/p))^k \leq O(k)^k$ for $p \in (0, 1)$ and $k \geq 0$, thus

$$\begin{aligned} p \cdot \sqrt{\log\left(\frac{1}{p}\right) \left(\log\left(\frac{n^\ell}{p}\right)\right)^{\ell-1}} &= p \cdot \sqrt{\log\left(\frac{1}{p}\right) \left(\ell \log(n) + \log\left(\frac{1}{p}\right)\right)^{\ell-1}} \\ &\leq O(1)^\ell \cdot \left(\sqrt{(\ell \log(n))^{\ell-1} + \ell^{\ell/2}}\right) \\ &= O(1)^\ell \cdot \sqrt{1 + (\ell \log(n))^{\ell-1}}. \quad \square \end{aligned}$$

References

- [AA18] Scott Aaronson and Andris Ambainis. Forrelation: A problem that optimally separates quantum from classical computing. *SIAM J. Comput.*, 47(3):982–1038, 2018. 4, 5
- [BB20] Shalev Ben-David and Eric Blais. A tight composition theorem for the randomized query complexity of partial functions: Extended abstract. In *FOCS*, pages 240–246. IEEE, 2020. 4
- [Bon70] Aline Bonami. Étude des coefficients de fourier des fonctions de $l^p(g)$. *Annales de l’institut Fourier*, 20(2):335–402, 1970. URL: <http://eudml.org/doc/74019>. 12
- [BP18] Joseph Briggs and Wesley Pegden. Extremal collections of k -uniform vectors. *arXiv preprint arXiv:1801.09609*, 2018. 26
- [BS20] Nikhil Bansal and Makrand Sinha. $\$k\$$ -forrelation optimally separates quantum and classical query complexity. *Electron. Colloquium Comput. Complex.*, 27:127, 2020. 4, 5
- [BTW15] Eric Blais, Li-Yang Tan, and Andrew Wan. An inequality for the fourier spectrum of parity decision trees. *CoRR*, abs/1506.01055, 2015. 4, 10, 15
- [CGL⁺20] Eshan Chattopadhyay, Jason Gaitonde, Chin Ho Lee, Shachar Lovett, and Abhishek Shetty. Fractional pseudorandom generators from any fourier level. *CoRR*, abs/2008.01316, 2020. 3
- [CHHL19] Eshan Chattopadhyay, Pooya Hatami, Kaave Hosseini, and Shachar Lovett. Pseudorandom generators from polarizing random walks. *Theory Comput.*, 15:1–26, 2019. 3
- [CHLT19] Eshan Chattopadhyay, Pooya Hatami, Shachar Lovett, and Avishay Tal. Pseudorandom generators from the second fourier level and applications to AC0 with parity gates. In *ITCS*, volume 124 of *LIPICs*, pages 22:1–22:15. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019. 3
- [CHRT18] Eshan Chattopadhyay, Pooya Hatami, Omer Reingold, and Avishay Tal. Improved pseudorandomness for unordered branching programs through local monotonicity. In *STOC*, pages 363–375. ACM, 2018. 3
- [CPT20] Gil Cohen, Noam Peri, and Amnon Ta-Shma. Expander random walks: A fourier-analytic approach. *Electron. Colloquium Comput. Complex.*, 27:163, 2020. 6
- [CS16] Gil Cohen and Igor Shinkar. The complexity of DNF of parities. In *ITCS*, pages 47–58. ACM, 2016. 2

- [GRT21] Uma Girish, Ran Raz, and Avishay Tal. Quantum versus randomized communication complexity, with efficient players. In *ITCS*, volume 185 of *LIPICs*, pages 54:1–54:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021. [3](#), [6](#)
- [GRZ20] Uma Girish, Ran Raz, and Wei Zhan. Lower bounds for XOR of correlations. *Electron. Colloquium Comput. Complex.*, 27:101, 2020. [3](#)
- [GSTW16] Parikshit Gopalan, Rocco A. Servedio, Avishay Tal, and Avi Wigderson. Degree and sensitivity: tails of two distributions. *Electron. Colloquium Comput. Complex.*, 23:69, 2016. [3](#)
- [HHL18] Hamed Hatami, Kaave Hosseini, and Shachar Lovett. Structure of protocols for XOR functions. *SIAM J. Comput.*, 47(1):208–217, 2018. [2](#)
- [KM93] Eyal Kushilevitz and Yishay Mansour. Learning decision trees using the fourier spectrum. *SIAM J. Comput.*, 22(6):1331–1348, 1993. [2](#)
- [KQS15] Raghav Kulkarni, Youming Qiao, and Xiaoming Sun. On the power of parity queries in boolean decision trees. In *TAMC*, volume 9076 of *Lecture Notes in Computer Science*, pages 99–109. Springer, 2015. [2](#)
- [Kra10] Joshua Brown Kramer. On the most weight w vectors in a dimension k binary code. *Electron. J. Comb.*, 17(1), 2010. [26](#)
- [Lee19] Chin Ho Lee. Fourier bounds and pseudorandom generators for product tests. In *Computational Complexity Conference*, volume 137 of *LIPICs*, pages 7:1–7:25. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019. [3](#)
- [Man95] Yishay Mansour. An $o(n^{\log \log n})$ learning algorithm for DNF under the uniform distribution. *J. Comput. Syst. Sci.*, 50(3):543–550, 1995. [3](#)
- [MO09] Ashley Montanaro and Tobias Osborne. On the communication complexity of XOR functions. *CoRR*, abs/0909.3392, 2009. [2](#)
- [MS20] Nikhil S. Mande and Swagato Sanyal. On parity decision trees for fourier-sparse boolean functions. In *FSTTCS*, volume 182 of *LIPICs*, pages 29:1–29:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020.
- [NN93] Joseph Naor and Moni Naor. Small-bias probability spaces: Efficient constructions and applications. *SIAM J. Comput.*, 22(4):838–856, 1993. [2](#), [3](#)
- [O’D12] Ryan O’Donnell. Open problems in analysis of boolean functions. *CoRR*, abs/1204.6447, 2012. [3](#)
- [O’D14] Ryan O’Donnell. *Analysis of Boolean Functions*. Cambridge University Press, 2014. [4](#), [10](#), [12](#)
- [OS07] Ryan O’Donnell and Rocco A. Servedio. Learning monotone decision trees in polynomial time. *SIAM J. Comput.*, 37(3):827–844, 2007. [3](#)
- [OWZ⁺14] Ryan O’Donnell, John Wright, Yu Zhao, Xiaorui Sun, and Li-Yang Tan. A composition theorem for parity kill number. In *Computational Complexity Conference*, pages 144–154. IEEE Computer Society, 2014. [2](#)

- [RSV13] Omer Reingold, Thomas Steinke, and Salil P. Vadhan. Pseudorandomness for regular branching programs via fourier analysis. In *APPROX-RANDOM*, volume 8096 of *Lecture Notes in Computer Science*, pages 655–670. Springer, 2013. 3
- [San19] Swagato Sanyal. Fourier sparsity and dimension. *Theory Comput.*, 15:1–13, 2019. 2
- [SSW20] Alexander A. Sherstov, Andrey A. Storozhenko, and Pei Wu. An optimal separation of randomized and quantum query complexity. *Electron. Colloquium Comput. Complex.*, 27:128, 2020. 3, 4, 9, 10
- [STIV17] Amir Shpilka, Avishay Tal, and Ben lee Volk. On the structure of boolean functions with small spectral norm. *Comput. Complex.*, 26(1):229–273, 2017. 2
- [Tal17] Avishay Tal. Tight bounds on the fourier spectrum of AC0. In *Computational Complexity Conference*, volume 79 of *LIPICs*, pages 15:1–15:31. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2017. 3
- [Tal20] Avishay Tal. Towards optimal separations between quantum and randomized query complexities. In *FOCS*, pages 228–239. IEEE, 2020. 3, 4, 9
- [TWXZ13] Hing Yin Tsang, Chung Hoi Wong, Ning Xie, and Shengyu Zhang. Fourier sparsity, spectral norm, and the log-rank conjecture. In *FOCS*, pages 658–667. IEEE Computer Society, 2013. 2
- [ZS09] Zhiqiang Zhang and Yaoyun Shi. Communication complexities of symmetric XOR functions. *Quantum Inf. Comput.*, 9(3&4):255–263, 2009. 2
- [ZS10] Zhiqiang Zhang and Yaoyun Shi. On the parity complexity measures of boolean functions. *Theor. Comput. Sci.*, 411(26-28):2612–2618, 2010. 2

A Proof of Corollary 1.8

Corollary (Corollary 1.8 restated). *Let \mathcal{T} be a parity decision tree of size at most $s > 1$ on n variables. Then,*

$$\forall \ell \in [n] : L_{1,\ell}(f) \leq (\log(s))^{\ell/2} \cdot O(\ell \cdot \log(n))^{1.5\ell}.$$

Proof. We approximate \mathcal{T} with error $\varepsilon = 1/n^\ell$ by another parity decision tree \mathcal{T}' of depth $d = \lceil \log(s \cdot n^\ell) \rceil$, where we simply replace all nodes of depth d in \mathcal{T} with leaves that return 0. Since there are at most s nodes in \mathcal{T} , the probability that a random input would reach one of the nodes of depth d is at most $2^{-d} \cdot s \leq 1/n^\ell$. Hence $\Pr_x[\mathcal{T}(x) \neq \mathcal{T}'(x)] \leq \varepsilon$. This implies that $|\widehat{\mathcal{T}}(S) - \widehat{\mathcal{T}'}(S)| \leq \varepsilon$ for any subset $S \subseteq [n]$. Thus,

$$L_{1,\ell}(\mathcal{T}) = \sum_{S:|S|=\ell} |\widehat{\mathcal{T}}(S)| \leq \sum_{S:|S|=\ell} (|\widehat{\mathcal{T}'}(S)| + \varepsilon) \leq L_{1,\ell}(\mathcal{T}') + 1.$$

Since \mathcal{T}' is of depth at most $d = \lceil \log(s) + \ell \cdot \log(n) \rceil = O(\log(s) \cdot \ell \cdot \log(n))$, we obtain our bound. \square

B Proof of Lemma 3.3

We will use the definition of sub-Gaussian random variables.

Definition B.1 (Sub-Gaussian random variables). We say a random variable x is Δ -sub-Gaussian if $\mathbb{E}[e^{t \cdot x}] \leq e^{t^2 \Delta^2}$ holds for all $t \in \mathbb{R}$.

Now we prove the following sub-Gaussian adaptive Azuma's inequality.

Lemma B.2 (Sub-Gaussian adaptive Azuma's inequality). Let $X^{(0)}, \dots, X^{(D)}$ be a martingale with respect to a filtration $(\mathcal{F}^{(i)})_{i=0}^D$ ¹³ and $\Delta^{(1)}, \dots, \Delta^{(D)}$ be a sequence of magnitudes such that $X^{(0)} = 0$ and $X^{(i)} = X^{(i-1)} + \delta^{(i)}$ for $i \in [D]$, where if conditioning on $\mathcal{F}^{(i-1)}$, $\delta^{(i)}$ is a $\Delta^{(i)}$ -sub-Gaussian random variable and $\Delta^{(i)}$ is a fixed value.

If there exists some constant $U \geq 0$ such that $\sum_{i=1}^D |\Delta^{(i)}|^2 \leq U$ always holds, then for any $\beta \geq 0$ we have

$$\Pr \left[\max_{i=0,1,\dots,D} |X^{(i)}| \geq \beta \cdot \sqrt{2U} \right] \leq 2 \cdot e^{-\beta^2/2}.$$

Proof. The bound holds trivially when $\beta = 0$, hence we assume $\beta > 0$ from now on. We construct another martingale $\widehat{X}^{(0)}, \dots, \widehat{X}^{(D)}$ as follows:

$$\widehat{X}^{(i)} = \begin{cases} X^{(i)} & 0 \leq i \leq d, \\ X^{(d)} & i > d, \end{cases} \quad \text{where } d = \min \{D\} \cup \left\{ i \in \{0, 1, \dots, D\} \mid |X^{(i)}| \geq \beta \cdot \sqrt{2U} \right\}.$$

We write $\widehat{\delta}^{(i)} = \widehat{X}^{(i)} - \widehat{X}^{(i-1)}$, then $\widehat{\delta}^{(i)} = \delta^{(i)}$ for all $i \leq d$; and $\widehat{\delta}^{(i)} \equiv 0$ for all $i > d$. Let $\widehat{\Delta}^{(i)} = \Delta^{(i)}$ for all $i \leq d$; and $\widehat{\Delta}^{(i)} \equiv 0$ for all $i > d$. Thus $\widehat{\delta}^{(i)}$ is $\widehat{\Delta}^{(i)}$ -sub-Gaussian given $\mathcal{F}^{(i-1)}$; and

$$\sum_{i=1}^D |\widehat{\Delta}^{(i)}|^2 = \sum_{i=1}^d |\Delta^{(i)}|^2 \leq U.$$

Moreover, we have

$$\Pr \left[\max_{i=0,1,\dots,D} |X^{(i)}| \geq \beta \cdot \sqrt{2U} \right] = \Pr \left[|\widehat{X}^{(D)}| \geq \beta \cdot \sqrt{2U} \right].$$

Let $t > 0$ be a parameter and we bound $\mathbb{E}[e^{t \cdot \widehat{X}^{(D)}}]$ as follows

$$\mathbb{E} \left[e^{t \cdot \widehat{X}^{(D)}} \right] = \mathbb{E}_{\mathcal{F}^{(D-1)}} \left[e^{t \cdot \widehat{X}^{(D-1)}} \cdot \mathbb{E}_{\mathcal{F}^{(D)}} \left[e^{t \cdot (\widehat{X}^{(D)} - \widehat{X}^{(D-1)})} \mid \mathcal{F}^{(D-1)} \right] \right] \quad (15)$$

$$= \mathbb{E}_{\mathcal{F}^{(D-1)}} \left[e^{t \cdot \widehat{X}^{(D-1)}} \cdot \mathbb{E}_{\mathcal{F}^{(D)}} \left[e^{t \cdot \widehat{\delta}^{(D)}} \mid \mathcal{F}^{(D-1)} \right] \right] \quad (16)$$

$$\leq \mathbb{E}_{\mathcal{F}^{(D-1)}} \left[e^{t \cdot \widehat{X}^{(D-1)}} \cdot e^{t^2 (\widehat{\Delta}^{(D)})^2} \right] \quad (\text{since } \widehat{\delta}^{(D)} \text{ is } \widehat{\Delta}^{(D)}\text{-sub-Gaussian})$$

$$\leq \mathbb{E}_{\mathcal{F}^{(D-1)}} \left[e^{t \cdot \widehat{X}^{(D-1)}} \cdot e^{t^2 (U - (\widehat{\Delta}^{(1)})^2 - \dots - (\widehat{\Delta}^{(D-1)})^2)} \right]$$

¹³ $\mathcal{F}^{(0)} \subseteq \mathcal{F}^{(1)} \subseteq \dots \subseteq \mathcal{F}^{(D)}$ is an increasing sequence of σ -algebra where each $\mathcal{F}^{(i)}$ makes $X^{(0)}, \dots, X^{(i+1)}$ measurable and $\mathbb{E}[X^{(i)} \mid \mathcal{F}^{(i-1)}] = X^{(i-1)}$. Intuitively, the filtration is the history of the martingale.

$$\begin{aligned}
&\leq \mathbb{E}_{\mathcal{F}^{(D-2)}} \left[e^{t \cdot \widehat{X}^{(D-2)}} \cdot e^{t^2 (U - (\widehat{\Delta}^{(1)})^2 - \dots - (\widehat{\Delta}^{(D-1)})^2)} e^{t^2 (\widehat{\Delta}^{(D-1)})^2} \right] \\
&\hspace{20em} \text{(similar to (15) and (16))} \\
&= \mathbb{E}_{\mathcal{F}^{(D-2)}} \left[e^{t \cdot \widehat{X}^{(D-2)}} \cdot e^{t^2 (U - (\widehat{\Delta}^{(1)})^2 - \dots - (\widehat{\Delta}^{(D-2)})^2)} \right] \\
&\leq \dots \leq \mathbb{E}_{\mathcal{F}^{(D-k)}} \left[e^{t \cdot \widehat{X}^{(D-k)}} \cdot e^{t^2 (U - (\widehat{\Delta}^{(1)})^2 - \dots - (\widehat{\Delta}^{(D-k)})^2)} \right] \leq \dots \\
&\leq e^{t^2 U}. \tag{17}
\end{aligned}$$

Setting $t = \beta/\sqrt{2U}$ implies that

$$\Pr \left[\widehat{X}^{(D)} \geq \beta \cdot \sqrt{2U} \right] \leq \frac{\mathbb{E} \left[e^{t \cdot \widehat{X}^{(D)}} \right]}{e^{t \cdot \beta \cdot \sqrt{2U}}} \leq \frac{e^{t^2 U}}{e^{\beta^2}} = e^{-\beta^2/2}.$$

Similarly we can show $\Pr \left[\widehat{X}^{(D)} \leq -\beta \cdot \sqrt{2U} \right] \leq e^{-\beta^2/2}$, which completes the proof by a union bound. \square

For our applications, we need the following fact.

Fact B.3. *Let x be a mean-zero random variable and assume $|x| \leq \Delta$ always holds. Then x is Δ -sub-Gaussian.*

Proof. Note that $e^{t \cdot x}$ is convex for all $t \in \mathbb{R}$. By Jensen's inequality, we have

$$\mathbb{E} \left[e^{t \cdot x} \right] \leq \frac{1}{2} (e^{-t\Delta} + e^{t\Delta}) = \sum_{i=0}^{+\infty} \frac{(t\Delta)^{2i}}{(2i)!} \leq \sum_{i=0}^{+\infty} \frac{(t\Delta)^{2i}}{i!} = e^{t^2 \Delta^2}. \quad \square$$

As a corollary of [Lemma B.2](#) and [Fact B.3](#), we obtain [Lemma 3.3](#).

Corollary ([Lemma 3.3](#) restated). *Let $X^{(0)}, \dots, X^{(D)}$ be a martingale and $\Delta^{(1)}, \dots, \Delta^{(D)}$ be a sequence of magnitudes such that $X^{(0)} = 0$ and $X^{(i)} = X^{(i-1)} + \Delta^{(i)} \cdot z^{(i)}$ for $i \in [D]$, where if conditioning on $z^{(1)}, \dots, z^{(i-1)}$,*

- (1) $z^{(i)}$ is a mean-zero random variable and $|z^{(i)}| \leq 1$ always holds;
- (2) $\Delta^{(i)}$ is a fixed value.

If there exists some constant $U \geq 0$ such that $\sum_{i=1}^D |\Delta^{(i)}|^2 \leq U$ always holds, then for any $\beta \geq 0$ we have

$$\Pr \left[\max_{i=0,1,\dots,D} |X^{(i)}| \geq \beta \cdot \sqrt{2U} \right] \leq 2 \cdot e^{-\beta^2/2}.$$

C Proof of Claim 5.8

Claim ([Claim 5.8](#) restated). $\Pr [\mathcal{E}_2] \leq \varepsilon/3$, where \mathcal{E}_2 is the following event:

$$\text{“ } \exists T \in \binom{[n]}{\ell-t}, i, r, r', \left| \Gamma_T^{(i)}(r, r') \right| \geq \left(100 \min \left\{ k, \log \left(\frac{n^\ell}{\varepsilon} \right) \right\} \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}} \right)^{\frac{r+r'}{2}} \cdot \sigma_T(r, r', C(v_{i-1}), i) \text{”}.$$

Proof. Let $k' = \min \{k, \lceil 6 \log(n^\ell/\varepsilon) \rceil\} \leq 12 \min \{k, \log(n^\ell/\varepsilon)\}$. Then \mathcal{T} is also a depth- D $2k'$ -clean parity decision tree. Observe that

$$\begin{aligned}
& \Pr \left[\left| \Gamma_T^{(i)}(r, r') \right| \geq \left(\frac{4k'}{\eta^{2/k'}} \right)^{(r+r')/2} \cdot \sigma_T(r, r', C(v_{i-1}), i) \right] \\
& \leq \max_{C(v_{i-1})} \Pr \left[\left| \Gamma_T^{(i)}(r, r') \right| \geq \left(\frac{4k'}{\eta^{2/k'}} \right)^{(r+r')/2} \cdot \sigma_T(r, r', C(v_{i-1}), i) \mid C(v_{i-1}) \right] \\
& \leq \underbrace{\frac{(4 \cdot k')^{r+r'}}{(2 \cdot (r+r'))^{k'}}}_{\leq 1} \cdot \underbrace{\eta^{2 - \frac{2(r+r')}{k'}}}_{\leq \eta} \quad (\text{due to the second bound in Lemma 3.2 and } k \geq 4 \cdot \ell \geq 4 \cdot (r+r')) \\
& \leq \eta.
\end{aligned}$$

Thus by union bound over all $T \in \binom{[n]}{\ell-t}, i \in [D'], r \in [t], 0 \leq r' \leq t-r$, we have

$$\Pr \left[\exists T, i, r, r', \left| \Gamma_T^{(i)}(r, r') \right| \geq \left(\frac{4k}{\eta^{2/k}} \right)^{(r+r')/2} \cdot \sigma_T(r, r', C(v_{i-1}), i) \right] \leq Dt^2 n^{\ell-t} \cdot \eta \leq \frac{n^{\ell+2} \eta}{3} \leq \frac{n^{3 \cdot \ell} \eta}{3},$$

where we use the fact $n \geq \max \{D, 3 \cdot t\}$ and $t \geq 1$. By setting $\eta = \varepsilon/n^{3 \cdot \ell}$, we have

$$\frac{4k'}{\eta^{2/k'}} = 4k' \left(\frac{n^{3 \cdot \ell}}{\varepsilon} \right)^{\frac{2}{k'}} \leq 4k' \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k'}} \leq 4 \cdot 12 \min \left\{ k, \log \left(\frac{n^\ell}{\varepsilon} \right) \right\} \cdot 2 \left(\frac{n^\ell}{\varepsilon} \right)^{\frac{6}{k}},$$

as desired. \square

D Proof of Claim 5.10

We first need the following simple bound on M .

Lemma D.1. *For any integer $s \geq 1$, we have*

$$\sum_{r=s}^t M(D, d, k, \ell, t-r, \varepsilon) \leq \frac{2 \cdot M(D, d, k, \ell, t, \varepsilon)}{(\tau D \cdot \log(n^\ell/\varepsilon))^{s/2}}.$$

Proof. We simply expand the formula of M as follows:

$$\begin{aligned}
\frac{\sum_{r=s}^t M(D, d, k, \ell, t-r, \varepsilon)}{M(D, d, k, \ell, t, \varepsilon)} &= \sum_{r=s}^t \left(\tau \cdot (D + dk) \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{6/k} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{-r/2} \\
&\leq \sum_{r=s}^{+\infty} \left(\tau \cdot (D + dk) \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{6/k} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{-r/2} \\
&\leq 2 \cdot \left(\tau \cdot (D + dk) \cdot \left(\frac{n^\ell}{\varepsilon} \right)^{6/k} \log \left(\frac{n^\ell}{\varepsilon} \right) \right)^{-s/2} \quad (\text{due to } \tau \geq 4 \text{ and } s \geq 1) \\
&\leq 2 \cdot \left(\tau D \cdot \log \left(n^\ell/\varepsilon \right) \right)^{-s/2}. \quad \square
\end{aligned}$$

Now we prove Claim 5.10.

Claim (Claim 5.10 restated). When $\mathcal{E}_1 \vee \mathcal{E}_2$ does not happen, $\sum_{i=1}^D |\mu_T^{(i)}| \leq R$ and $\sum_{i=1}^D |\delta_T^{(i)}|^2 \leq R^2$ hold for all $T \in \binom{[n]}{\ell-t}$, where

$$R = \frac{M(D, d, k, \ell, t, \varepsilon)}{5 \cdot \sqrt{\log(n^\ell/\varepsilon)}}.$$

Proof. We verify for each $T \in \binom{[n]}{\ell-t}$ as follows:

$$\begin{aligned} \sum_{i=1}^D |\mu_T^{(i)}| &= \sum_{i=1}^{D'} |\mu_T^{(i)}| \leq \sum_{i=1}^{D'} \sum_{\substack{r=2, \\ \text{even}}}^t |A(T, r, i)| && \text{(due to (8))} \\ &\leq \sum_{i=1}^{D'} \sum_{\substack{r=2, \\ \text{even}}}^t \left(M(D, d, k, \ell, t-r, \varepsilon) + \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\left(\frac{800}{\tau}\right)^r \left(\frac{|J(v_{i-1})|}{2d}\right)^r} \cdot \mathbf{1}_{|J(v_{i-1})|>1} \right) \\ &&& \text{(due to (10))} \\ &\leq \sum_{i=1}^{D'} \sum_{\substack{r=2, \\ \text{even}}}^t \left(M(D, d, k, \ell, t-r, \varepsilon) + \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \left(\frac{|J(v_{i-1})|}{2d}\right) \left(\frac{800}{\tau}\right)^{r/2} \cdot \mathbf{1}_{|J(v_{i-1})|>1} \right) \\ &&& \text{(Since } |J(v_{i-1})| \leq 2d \text{ from Corollary 4.9)} \\ &\leq \frac{2 \cdot M(D, d, k, \ell, t, \varepsilon)}{\tau \cdot \log(n^\ell/\varepsilon)} + \frac{1.1 \cdot 800 \cdot M(D, d, k, \ell, t, \varepsilon)}{\tau \cdot \sqrt{\log(n^\ell/\varepsilon)}} \\ &&& \text{(due to Lemma D.1 and Corollary 4.9 and } \tau = 10^4\text{)} \\ &\leq \frac{M(D, d, k, \ell, t, \varepsilon)}{5 \cdot \sqrt{\log(n^\ell/\varepsilon)}} = R \end{aligned}$$

and with similar calculation, we have

$$\begin{aligned} \sum_{i=1}^D |\delta_T^{(i)}|^2 &\leq \sum_{i=1}^{D'} \left(\sum_{\substack{r=1, \\ \text{odd}}}^t \left(M(D, d, k, \ell, t-r, \varepsilon) + \frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\frac{|J(v_{i-1})|}{2d}} \left(\frac{800}{\tau}\right)^{r/2} \cdot \mathbf{1}_{|J(v_{i-1})|>1} \right) \right)^2 \\ &\leq \sum_{i=1}^{D'} \left(\frac{2 \cdot M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\tau D \cdot \log(n^\ell/\varepsilon)}} + \frac{1.1 \cdot \sqrt{800} \cdot M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\tau} \cdot \sqrt{\log(n^\ell/\varepsilon)}} \cdot \sqrt{\frac{|J(v_{i-1})|}{2d}} \cdot \mathbf{1}_{|J(v_{i-1})|>1} \right)^2 \\ &&& \text{(due to } \tau = 10^4\text{)} \\ &\leq \left(\frac{M(D, d, k, \ell, t, \varepsilon)}{\sqrt{\log(n^\ell/\varepsilon)}} \right)^2 \sum_{i=1}^{D'} 2 \cdot \left(\frac{4}{\tau D} + \frac{968}{\tau} \cdot \frac{|J(v_{i-1})|}{2d} \cdot \mathbf{1}_{|J(v_{i-1})|>1} \right) \\ &&& \text{(due to } (a+b)^2 \leq 2(a^2 + b^2)\text{)} \\ &\leq \left(\frac{2000 \cdot M(D, d, k, \ell, t, \varepsilon)}{\tau \cdot \sqrt{\log(n^\ell/\varepsilon)}} \right)^2 = R^2. \quad \square \end{aligned}$$