

Direct Sums for Parity Decision Trees

Tyler Besselman  

NYU Shanghai, China

Mika Göös 


EPFL, Lausanne, Switzerland

Siyao Guo  

NYU Shanghai, China

Gilbert Maystre 

EPFL, Lausanne, Switzerland

Weiqiang Yuan  

EPFL, Lausanne, Switzerland

Abstract

Direct sum theorems state that the cost of solving k instances of a problem is at least $\Omega(k)$ times the cost of solving a single instance. We prove the first such results in the randomised parity decision tree model. We show that a direct sum theorem holds whenever (1) the lower bound for parity decision trees is proved using the *discrepancy method*; or (2) the lower bound is proved relative to a *product distribution*.

2012 ACM Subject Classification Theory of computation \rightarrow Oracles and decision trees

Keywords and phrases direct sum, parity decision trees, query complexity

Digital Object Identifier 10.4230/LIPIcs.CCC.2025.16

Related Version *Full Version*: <https://eccc.weizmann.ac.il/report/2024/203/>

Funding *Tyler Besselman*: Supported by the National Natural Science Foundation of China Grant No.62102260, NYTP Grant No.20121201, and NYU Shanghai Boost Fund.

Mika Göös: Supported by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number MB22.00026.

Siyao Guo: Supported by the National Natural Science Foundation of China Grant No.62102260, NYTP Grant No.20121201, and NYU Shanghai Boost Fund.

Gilbert Maystre: Supported by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number MB22.00026.

Weiqiang Yuan: Supported by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number MB22.00026.

Acknowledgements We thank Farzan Byramji for useful comments on an earlier version of this paper.

1 Introduction

One of the most basic questions that can be asked for any model of computation is:

How does the cost of computing k independent instances scale with k ?

A *direct sum* theorem states that if the cost of solving a single copy is C , then solving k copies has cost at least $\Omega(k \cdot C)$, which matches the trivial algorithm that solves the k copies separately. Direct sums have been studied exhaustively for randomised query complexity R^{dt} , randomised communication complexity R^{cc} , and other concrete models of computation; see Section 1.3 for prior work. In this work, we initiate the study of direct sum problems for randomised parity decision tree complexity R^{pt} , a computational model sandwiched between the widely-studied R^{dt} and R^{cc} .



© Tyler Besselman, Mika Göös, Siyao Guo, Gilbert Maystre, and Weiqiang Yuan;

licensed under Creative Commons License CC-BY 4.0

40th Computational Complexity Conference (CCC 2025).

Editor: Srikanth Srinivasan; Article No. 16; pp. 16:1–16:38



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



Parity decision trees

Parity decision trees generalise the usual notion of decision trees by allowing *parity queries*. To compute a function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ on input $x \in \{0, 1\}^n$, a deterministic parity decision tree T performs queries of the form “what is $\langle a, x \rangle$?” where $a \in \{0, 1\}^n$ and $\langle a, x \rangle := \sum_i a_i x_i \bmod 2$. Once enough queries have been made, T outputs $f(x)$. Parity decision trees are more powerful than ordinary decision trees: We have $D^{\text{pt}}(f) \leq D^{\text{dt}}(f)$ where $D^{\text{dt}}(f)$ (resp. $D^{\text{pt}}(f)$) denotes the (parity) decision tree complexity of f , defined as the least depth of a deterministic (parity) decision tree computing f . On the other hand, the n -bit XOR function is an example where $D^{\text{dt}}(\text{XOR}) = n$ while $D^{\text{pt}}(\text{XOR}) = 1$. We define a *randomised* parity decision tree \mathcal{T} as a distribution over deterministic parity trees $T \sim \mathcal{T}$. Then $R_\varepsilon^{\text{pt}}(f)$ is defined as the worst-case depth (over both input and randomness of the tree) of the best randomised parity tree \mathcal{T} computing f with error ε , that is, $\Pr[\mathcal{T}(x) \neq f(x)] \leq \varepsilon$ for all x . As usual, we let $R^{\text{pt}} := R_{1/3}^{\text{pt}}$. To simplify notation, we drop the superscript **pt** and write $D = D^{\text{pt}}$ and $R = R^{\text{pt}}$ for short.

Our main research question is now formulated as follows. Let $f^k: (\{0, 1\}^n)^k \rightarrow \{0, 1\}^k$ denote the *direct sum* function that takes k instances $x := (x^1, \dots, x^k)$ and returns the value of f on each of them, $f^k(x) := (f(x^1), \dots, f(x^k))$. We study the following question.

► **Question 1.** *Do we have $R(f^k) \geq \Omega(k) \cdot R(f)$ for every function f ?*

We show two (incomparable) main results: We answer Question 1 affirmatively when the randomised parity decision tree lower bound is proved using the *discrepancy method* (Section 1.1), or when the lower bound is proved relative to a *product distribution* (Section 1.2).

1.1 First result: Direct sum for discrepancy

Discrepancy is one of the oldest-known methods for proving randomised communication lower bounds [56, 3]. Let us tailor its definition to the setting of randomised parity trees. Thinking of $\{0, 1\}^n$ as the vector space \mathbb{Z}_2^n , consider some affine subspace $S \subseteq \{0, 1\}^n$ and a probability distribution μ over the inputs $\{0, 1\}^n$. The discrepancy of S measures how biased f is on S . Namely, let $C_S^b := \Pr_{x \sim \mu}[f(x) = b \wedge x \in S]$. The difference $\Delta_S := |C_S^0 - C_S^1|$ is called the bias of S under μ . We define $\text{bias}(f)$ as the minimum over μ of the maximum difference Δ_S an affine subspace can attain. Finally, the *discrepancy bound* $\text{disc}(f)$ is defined as $\log(1/\text{bias}(f))$. As in communication complexity, it is not hard to see that $R(f) \geq \Omega(\text{disc}(f))$; see Section 3 for details.

► **Theorem 1.** *We have $R(f^k) \geq \Omega(k) \cdot \text{disc}(f)$ for any function f .*

In particular, if we have a function f whose randomised parity decision tree complexity is equal to its discrepancy, $R(f) = \Theta(\text{disc}(f))$, then Theorem 1 shows $R(f^k) \geq \Omega(k) \cdot R(f)$ answering Question 1 for that function. To prove Theorem 1, we first establish a particularly simple characterisation of $\text{disc}(f)$ that relies on affine spaces defined by a single constraint. We then prove a perfect direct sum (and even an XOR lemma) for discrepancy using Fourier analysis.

1.2 Second result: Direct sum for product distributions

The standard approach for proving randomised lower bounds is to use Yao’s principle [55], which states that $R(f) = \max_\mu D_{1/3}(f, \mu)$. Here $D_\varepsilon(f, \mu)$ is the *distributional* ε -error complexity of f defined as the least depth of a (deterministic) parity tree T such that

$\Pr_{\mathbf{x} \sim \mu}[T(\mathbf{x}) \neq f(\mathbf{x})] \leq \varepsilon$. We say that a distribution μ over $\{0, 1\}^n$ is *product* if it can be written as the product of n independent Bernoulli distributions. We define the best lower bound provable using a product distribution as

$$D_\varepsilon^\times(f) := \max_{\mu \text{ product}} D_\varepsilon(f, \mu) \quad \text{and} \quad D^\times := D_{1/3}^\times.$$

Our second result answers Question 1 affirmatively (modulo logarithmic factors) whenever the randomised parity decision tree lower bound is proved relative to a product distribution.

► **Theorem 2.** *We have $R(f^k) \geq \Omega(k/\log n) \cdot D^\times(f)$ for any n -bit function f .*

We show moreover that the $O(\log n)$ -factor loss in Theorem 2 can be avoided when μ is the uniform distribution (or more generally any *bounded-bias* distribution). One should compare this to the state-of-the-art in communication complexity, where the quantitatively best distributional direct sum results are also for product distributions and suffer logarithmic-factor losses [37, 4].

To prove Theorem 2, we introduce a new complexity measure tailored for product distributions, which we call *skew complexity* $S(f)$ and which we define precisely in Section 4. We prove that this new measure admits a perfect direct sum theorem, $S(f^k) = \Omega(k) \cdot S(f)$, and that it characterises the measure D^\times up to an $O(\log n)$ factor. (We also show that the logarithmic loss is necessary for our approach: there is a function f such that $S(f) = O(1)$, even though $D^\times(f) = \Theta(\log n)$.) We give a more in-depth technical overview in Section 2.

Comparison of main results

We also show that our two main results (Theorems 1 and 2) are incomparable: For some functions f , our first result gives a much stronger lower bound for f^k than the second result – and vice versa. See Section 7 for the proof.

► **Lemma 3.** *The complexity measures disc and D^\times are incomparable:*

1. *There is an n -bit function f such that $\text{disc}(f) = O(\log n)$ while $D^\times(f) = \Theta(n)$.*
2. *There is an n -bit function f such that $\text{disc}(f) = \Theta(n)$ while $D^\times(f) = O(1)$.*

1.3 Related work

Parity decision trees

Even though the direct sum problem for parity decision trees has not been studied before, the model has been studied extensively. Parity decision trees were first defined by Kushilevitz and Mansour [40] in the context of learning theory. Several prior works have studied their basic combinatorial properties [57, 46] as well as Fourier-analytic properties [28, 27], often with connections to the log-rank conjecture [54, 53, 48, 20, 32, 43]; see also the survey [39]. There are various lifting theorems involving parity decision trees: lifting from D^{pt} to D^{cc} [31], from D^{dt} to D^{pt} [19, 5, 1], and from R^{dt} to R^{pt} [52, 17]. These lifting theorems have played a central role in proving lower bounds for proof systems that can reason using parities [33, 22, 25, 11, 18, 2].

Decision trees

In the decision tree model with classical queries, a deterministic direct sum theorem, $D^{\text{dt}}(f^k) = k \cdot D^{\text{dt}}(f)$, and even the stronger *composition theorem*, $D^{\text{dt}}(g \circ f^k) = D^{\text{dt}}(g) \cdot D^{\text{dt}}(f)$, are easy to show by combining adversary strategies [50]. In the randomised case, an optimal direct

sum result, $R^{\text{dt}}(f^k) \geq \Omega(k) \cdot R^{\text{dt}}(f)$, is known [38, 36, 21]. Whether a composition theorem holds for randomised query complexity, $R^{\text{dt}}(g \circ f^k) \geq \Omega(R^{\text{dt}}(g) \cdot R^{\text{dt}}(f))$ (for total g and f), is a major open problem [10, 6, 8, 7, 49]. In the randomised setting, it is possible that the direct sum problem f^k requires strictly more than $\Theta(k) \cdot R^{\text{dt}}(f)$ queries: if one wants to succeed in computing all k copies with probability $\geq 2/3$, then a naive application of the union bound would require each copy to have error $\ll 1/k$. Results stating that one sometimes has $R^{\text{dt}}(f^k) \geq \omega(k) \cdot R^{\text{dt}}(f)$ are called “strong” direct sum theorems [12, 13] and they sometimes hold even for composed functions [9, 16, 29].

Communication complexity

The direct sum question for deterministic communication complexity was posed in [23] and it remains a notoriously difficult open problem [34]. By contrast, in the randomised setting, the direct sum problem is characterised by *information complexity* [15], which has inspired a line of works too numerous to cite here; see [35, §1.1] for an up-to-date overview. One of the key findings is that a direct sum for communication protocols is *false* in full generality in the distributional setting [26, 47]. We leave open the intriguing possibility that the information complexity approach can be adapted to parity decision trees. Historically, one of the first direct sum theorems proved for randomised communication was for the discrepancy bound [51, 41] (analogously to our Theorem 1). Here, discrepancy is known to be equivalent to the γ_2 -norm [42]. We also mention that a near-optimal direct sum theorem holds for product distributions [4] (analogously to our Theorem 2).

1.4 Open question: Deterministic direct sum

The main question left open by our work is Question 1, namely, whether $R = R^{\text{pt}}$ admits a direct sum theorem for all functions f . However, we would also like to highlight the analogous question in the deterministic case $D = D^{\text{pt}}$. As discussed above, this is a long-standing open problem in the case of deterministic communication complexity D^{cc} . The best results so far are:

1. $D^{\text{cc}}(f^k) \geq \tilde{\Omega}(k) \cdot D^{\text{cc}}(f)^{1/2}$ as proved in [23].
2. $D^{\text{cc}}(f^k) \geq \tilde{\Omega}(k) \cdot D^{\text{cc}}(f) / \log \text{rank}(f)$ as proved in [34].

We observe in Section A.1 that both approaches have analogues in the parity setting.

► **Theorem 4.** *For any function f and $k \geq 1$,*

1. $D(f^k) \geq k \cdot D(f)^{1/2}$,
2. $D(f^k) \geq k \cdot D(f) / \log \text{spar}(f)$.

We leave it as an open question whether a perfect direct sum theorem holds for deterministic parity decision trees. We think one should attack this problem before addressing the (presumably much harder) problem for deterministic communication complexity.

2 Technical overview

We focus here on our second main result in Theorem 2 stating that $R(f^k) \geq \Omega(k / \log n) \cdot D^\times(f)$ and which is technically the much more involved theorem. Our main technical result is the following direct sum result for distributional complexity. Here $\mu^k := \mu \times \dots \times \mu$ (k times).

► **Theorem 5.** *There exists a universal constant C such that the following holds. For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ over $\{0, 1\}^n$, and $k \geq 1$,*

$$D_\varepsilon(f^k, \mu^k) \geq \Omega\left(\frac{k\delta}{\log(n/\delta)}\right) \cdot (D_{\varepsilon+\delta}(f, \mu) - C \cdot \log(n/\delta)) \quad \forall \varepsilon, \delta \geq 0.$$

When $D^\times(f) \geq 6C \cdot \log n$, Theorem 2 follows by taking $\varepsilon = \delta = 1/6$. Indeed, let μ be the distribution achieving the maximum for D^\times . Using the easy direction of the minimax principle:

$$R(f^k) \geq \Omega(1) \cdot D_{1/6}(f^k, \mu^k) \geq \Omega(k/\log n) \cdot D_{1/3}(f, \mu) = \Omega(k/\log n) \cdot D^\times(f).$$

The remaining case $D^\times(f) \leq 6C \cdot \log(n)$ is handled separately using ad-hoc methods in Lemma 38. We now give an overview of the proof of Theorem 5.

Warm-up: Uniform distribution

We showcase the basic proof technique by sketching the proof in the simple case where μ is the uniform distribution. Fix an n -bit function f and let \mathcal{U} be the uniform distribution over $\{0, 1\}^n$. In the uniform (and more generally in the *bounded-bias*) case, we are actually able to avoid the $\log n$ additive/factor loss and obtain, for all $k \geq 1$,

$$D_\varepsilon(f^k, \mathcal{U}^k) \geq \Omega(k\delta) \cdot D_{\varepsilon+\delta}(f, \mathcal{U}) \quad \forall \delta \geq 0. \quad (1)$$

Fix a decision tree T of depth d computing k copies of f with error at most ε when $\mathbf{x} \sim \mathcal{U}^k$. We show how to extract a tree T^* that computes a single copy $\mathbf{y} \sim \mathcal{U}$ with error at most $\varepsilon + \delta$ and depth $\leq O(d/k\delta)$. Leaves of T correspond to affine subspaces of $(\{0, 1\}^n)^k$ of codimension $\leq d$. More generally, one can associate with any node v of T the set $C_v = \{w_1, \dots, w_{d(v)}\}$ of linear constraints that led to the node ($d(v)$ is the depth of the node v ; the root is at level 0) and the vector $b \in \{0, 1\}^k$ of desired values. The set of inputs S_v that reach node v is then given by $S_v := \{x \in (\{0, 1\}^n)^k : \langle w_j, x \rangle = b_j, \forall j \in [d(v)]\}$.

Of relevance here are the *pure constraints* one can extract from C_v . A pure constraint for copy $i \in [k]$ is some $w \in (\{0, 1\}^n)^k$ such that $w^j \neq 0^n$ if and only if $j = i$. To be more precise, the number of pure queries that can be extracted for query i at node v is defined with:

$$\text{pure}_i(C_v) := \dim(\text{span}(C_v) \cap W_i) \quad \text{where} \quad W_i := \{w \in (\{0, 1\}^n)^k : w^j = 0^n, \forall j \neq i\}.$$

We describe next two illustrative examples when there are $k = 2$ copies.

1. Node v corresponds to constraints “ $x_1^1 + x_1^2 = 0$ ” and “ $x_1^2 = 1$ ”. Then, $\text{pure}_1(C_v) = 1$ as it is possible to extract the pure parity constraint $x_1^1 = 1$ by adding the two constraints. In the same vein, $\text{pure}_2(C_v) = 1$.
2. Node v corresponds to constraints “ $x_1^1 + x_1^2 = 0$ ” and “ $x_1^2 + x_2^2 = 1$ ”. Then, $\text{pure}_1(C_v) = 0$ as it not possible to extract a pure constraint for the first copy.

► **Observation 6.** *For any node v , we have $d(v) \geq \sum_{i \in [k]} \text{pure}_i(C_v)$.*

As the second example highlights, it is possible for the inequality to be strict. This is a notable difference with classical decision trees: for any subcube $C \in (\{0, 1, *\}^n)^k$, the sum of fixed bits of each copy is the total number of fixed bits in C .

Where to plant y ?

The overarching idea of our result is that under the uniform distribution, *queries that increase the pure rank for a copy are the only ones that bring usable information*. It is thus enough to find a copy with low expected pure rank in T and plant the real instance y there. To make this precise, taking the expectation over leaves of T when $\mathbf{x} \sim \mathcal{U}$ with Observation 6 implies the existence of some copy $i \in [k]$ with low expected pure rank:

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}^k} [\text{pure}_i(C_{\ell(\mathbf{x})})] \leq O(d/k).$$

Let us fix this advantageous copy to be $i = 1$. On input $y \in \{0, 1\}^n$ we run the tree T with y planted as x^1 and delay actual querying of bits of y as much as possible. Suppose that the process has reached node v with constraint set C_v and there is a new parity query w to be answered. If $w \in \text{span}(C_v)$, the answer to that query can be found (an optimised tree would not do such a query). If $w \notin \text{span}(C_v)$, we say that w is *critical* for C_v if it would increase the pure rank for the first copy $\text{pure}_1(C_v \cup \{w\}) > \text{pure}_1(C_v)$. If w is critical, there is no way to avoid making a parity query to the real input y and our algorithm does it. If w is not critical, it is however enough to answer with a uniform bit (that is, move to a random child of v in T) without querying y at all.

To see this, further split $w = w^1 w^{-1}$, where $w^1 \in \{0, 1\}^n$ is the constraint for the first copy and $w^{-1} \in (\{0, 1\}^n)^{k-1}$ is the constraint for the rest of the copies. If w has $\text{pure}_1(C_v \cup \{w\}) = \text{pure}_1(C_v)$ and $w \notin \text{span}(C_v)$, it must be that $0^n w^{-1} \notin \text{span}(C_v)$. Since x^{-1} is drawn from the uniform distribution we thus have for any fixed y consistent with S_v :

$$\Pr_{x^{-1}} [\langle w, y x^{-1} \rangle = 0 \mid (y, x^{-1}) \in S_v] = \Pr_{x^{-1}} [\langle w^{-1}, x^{-1} \rangle = \langle w^1, y \rangle \mid (y, x^{-1}) \in S_v] = \frac{1}{2}. \quad (2)$$

Correctness and efficiency

Let us call the above randomised tree solving one copy as \mathcal{T} . Correctness can be argued by showing that the distribution of leaves attained in the process for $y \sim \mathcal{U}$ is the same as the distribution of leaves attained by $x \sim \mathcal{U}^k$ in T . On the other hand, \mathcal{T} has expected depth $O(d/k)$ as a real query to y is only ever made $\text{pure}_i(C_\ell)$ times for each leaf ℓ . In conclusion, \mathcal{T} has the following guarantees:

1. $\Pr_{y \sim \mathcal{U}, \mathcal{T} \sim \mathcal{T}} [\mathcal{T}(y) \neq f(y)] \leq \varepsilon$.
2. $\mathbb{E}_{y \sim \mathcal{U}, \mathcal{T} \sim \mathcal{T}} [\# \text{queries}(\mathcal{T}, y)] \leq d/k$.

Using Markov inequality, it is possible to derandomise \mathcal{T} to get a deterministic parity tree T^* solving f with a worst-case guarantee instead of an average-case one. This step introduces a parameter δ controlling a trade-off between cost and error and yields the desired result (1).

2.1 Beyond uniform: The skew measure

Observe that (2) can fail badly for non-uniform μ . As an illustrative example suppose that two random bits a, b are generated with $a \sim \text{Ber}(1/2)$ and $b \sim \text{Ber}(1/8)$. The constraint $a \oplus b = 1$ is not pure from the point of view of a . However, since b is skewed towards being 0, the realisation of the constraint gives information about a : $\Pr[a = 0 \mid a + b = 1] = 1/8 \ll 1/2$. Thus, it seems one needs to query a to answer the query $a + b$ even though the query is not critical for a !

To circumvent this, we introduce the *skew* measure. This new measure is built around the observation that each bit of an input $x \sim \mu$ can be sampled independently in two steps. Indeed, the following process is equivalent to $\text{Ber}(1/8)$:

1. Let $\rho \in \{0, \star\}$ be “0” with probability $3/4$ and \star with probability $1/4$.
2. If $\rho = 0$, return “0”, else return a sample $\text{Ber}(1/2)$.

Note that if we are “lucky” and $\rho = \star$, we are back in the uniform case and (2) holds again. If not, we have somehow pre-emptively fixed the return bit to value 0. The skew measure explicitly splits product distributions into a *random partial fixing* ρ followed by a uniform distribution over unfixed bits of ρ . A tree computing in this model gets help from ρ because ρ reduces the complexity of the function. When those bits are unfixed, it is on the other hand easier to analyse the behaviour of the tree as it is the uniform case again.

In Sections 5 and 6, we show a perfect direct sum for the skew measure and that perhaps surprisingly, this new measure is only a $\log n$ -factor away from D^\times .

3 Direct sum for disc

The goal of this section is to prove Theorem 1, restated here for convenience.

► **Theorem 1.** *We have $R(f^k) \geq \Omega(k) \cdot \text{disc}(f)$ for any function f .*

Let us start by defining discrepancy formally. We denote by \mathcal{S}_n the set of all affine subspaces of $\{0, 1\}^n$ and $\mathcal{O}_n \subseteq \mathcal{S}_n$ the set of affine subspaces of codimension 1. Note that all spaces $S \in \mathcal{O}_n$ can be written as $S = \{x \in \{0, 1\}^n : \langle a, x \rangle = b\}$ for some $a \in \{0, 1\}^n$ and $b \in \{0, 1\}$.

► **Definition 7.** *Let $f: \{0, 1\}^n \rightarrow \{0, 1\}$ be a boolean function and μ be a distribution over $\{0, 1\}^n$. The (parity) discrepancy of f with respect to μ is defined as:*

$$\text{disc}(f, \mu) := -\log \max_{S \in \mathcal{S}_n} \text{bias}(f, \mu, S) \quad \text{where} \quad \text{bias}(f, \mu, S) := \left| \sum_{x \in S} (-1)^{f(x)} \mu(x) \right|.$$

The (parity) discrepancy of f is $\text{disc}(f) := \max_{\mu} \text{disc}(f, \mu)$ where μ ranges over all distributions.

Observe that $\text{disc}(f) \geq 1$ for all non-constant f and by standard arguments, $R(f) \geq \text{disc}(f)$ (see Lemma 42). Using the latter, the only thing left to get Theorem 1 is to prove a direct sum result for discrepancy. We do this in a very strong way by actually establishing an XOR lemma for disc. Let $f^{\oplus k}$ denote the function that takes k instance and aggregates their result under f using XOR, so that $f^{\oplus k}(x^1, \dots, x^k) := f(x^1) \oplus \dots \oplus f(x^k)$.

► **Lemma 8.** *For any function f , distribution μ and $k \geq 1$,*

$$k \cdot \text{disc}(f, \mu) \geq \text{disc}(f^{\oplus k}, \mu^k) \geq k \cdot (\text{disc}(f, \mu) - 1).$$

This result is the strongest possible. Indeed, we cannot omit the “−1” on the right because of the counterexample $f := \text{XOR}$: we have $\text{disc}(f^{\oplus k}, \mu^k) \leq 1$ for any distribution μ . In Section A.3 we revisit this XOR lemma and show that it also holds in the distribution-free setting, with $\text{disc}(f^{\oplus k}) \approx k \cdot \text{disc}(f)$. As a final comment, we note that it is easier to work with $f^{\oplus k}$ instead of f^k in the discrepancy setting, as it is somewhat tedious to define discrepancy for multi-valued functions. Before formally proving Lemma 8, we show how it is used to prove the main result Theorem 1.

Proof of Theorem 1. Any decision tree computing f^k can be converted to a decision tree computing $f^{\oplus k}$. This is achieved by replacing the label $y \in \{0, 1\}^k$ of each leaf by its parity $\langle y, 1^k \rangle$. This operation does not increase the error probability or cost and so, using the easy direction of Yao’s principle:

$$\begin{aligned} R(f^k) &\geq \max_{\mu} D(f^k, \mu^k, 1/3) && \text{(Lemma 41)} \\ &\geq \max_{\mu} D(f^{\oplus k}, \mu^k, 1/3) \\ &\geq \max_{\mu} \text{disc}(f^{\oplus k}, \mu^k) - \log_2(3) && \text{(Lemma 42)} \\ &\geq k \cdot \max_{\mu} (\text{disc}(f, \mu) - 1) - \log_2(3) && \text{(Lemma 8)} \\ &\geq k \cdot (\text{disc}(f) - 1) - \log_2(3). \end{aligned}$$

If $\text{disc}(f) \geq 10$, then the string of inequalities yields $k \cdot (\text{disc}(f) - 1) - \log_2(3) \geq k \cdot \text{disc}(f)/10$. If f is constant, the claim is vacuously true. Finally, we show that for any non-constant f , $R(f^k) \geq k - \log(3/2)$ which completes the claim. Indeed, if $\text{disc}(f) \leq 10$, then $k - \log(3/2) \geq k \cdot \text{disc}(f)/100$.

16:8 Direct Sums for Parity Decision Trees

To this end, let f be a non-constant function and μ a distribution over $\{0, 1\}^n$ which is balanced over 0-inputs and 1-inputs, i.e. $\mu(f^{-1}(0)) = \mu(f^{-1}(1)) = 1/2$. Let T be the best deterministic parity decision tree for $D_{1/3}(f, \mu)$ and suppose toward contradiction that it has strictly less than $L := 2^k \cdot (2/3)$ leaves. Let $G \subseteq \{0, 1\}^n$ be the set of solutions which appear as a label on a leaf of T . We have $|G| < L$ and since μ is balanced, any solution $y \in \{0, 1\}^k$ is equally likely so that:

$$\Pr_{x \sim \mu^k} [T(x) = f^k(x)] \leq \Pr_{x \sim \mu^k} [f^k(x) \in G] \leq |G| \cdot 2^{-k} < 2/3.$$

Thus, T errs with probability $> 1/3$: a contradiction. \blacktriangleleft

We now proceed to prove Lemma 8 in three steps.

3.1 Step 1: Characterisation of discrepancy

Much like discrepancy for communication protocols can be characterised by the γ_2 -norm of the communication matrix [51, 42], we show that the parity discrepancy of f on μ is characterised by the L_∞ -norm of the Fourier transform of a related function F_μ . This characterisation has two purposes. First, proving an XOR lemma requires exploring all the possible ways for the k copies to sum to 1. This kind of convolution operation is greatly simplified in the Fourier domain, where it simply corresponds to standard multiplication. Second, the characterisation is also quite convenient to prove lower bounds on $\text{disc}(f, \mu)$ (which we do in Sections 7 and 8): it shows that maximum bias is (almost) attained for affine spaces of codimension 1 already.

The function F_μ

We relate a real-valued boolean function $F: \{0, 1\}^n \rightarrow \mathbb{R}$ with its Fourier transform $\widehat{F}: \{0, 1\}^n \rightarrow \mathbb{R}$ using the usual basis:

$$\forall z \in \{0, 1\}^n, \quad \widehat{F}(z) := \sum_{x \in \{0, 1\}^n} F(x) \cdot (-1)^{\langle x, z \rangle} \cdot 2^{-n}; \quad [\text{Fourier transform}]$$

$$\forall x \in \{0, 1\}^n, \quad F(x) := \sum_{z \in \{0, 1\}^n} \widehat{F}(z) \cdot (-1)^{\langle z, x \rangle}. \quad [\text{Inverse Fourier transform}]$$

See also [45] for more background on Fourier analysis. We use $\|\widehat{F}\|_\infty$ to denote the maximum absolute value of a Fourier coefficient of F . To analyze $\text{disc}(f, \mu)$, we introduce an associated function $F_\mu: \{0, 1\}^n \rightarrow \mathbb{R}$ defined by $F_\mu(x) := (-1)^{f(x)} \cdot \mu(x) \cdot 2^n$ and prove the following characterisation.

► **Lemma 9.** *For every function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and distribution μ over $\{0, 1\}^n$:*

$$\max_{S \in \mathcal{O}_n} \text{bias}(f, \mu, S) \leq \max_{S \in \mathcal{S}_n} \text{bias}(f, \mu, S) \leq \|\widehat{F}_\mu\|_\infty \leq 2 \cdot \max_{S \in \mathcal{O}_n} \text{bias}(f, \mu, S).$$

Proof. The first inequality holds immediately because $\mathcal{O}_n \subseteq \mathcal{S}_n$. For the second, fix a maximizing $S \in \mathcal{S}_n$. Suppose that $\text{codim}(S) = d$ and fix its constraints $a_j \in \{0, 1\}^n$ and $b_j \in \{0, 1\}$ for $j \in [d]$ so that $S = \{x \in \{0, 1\}^n : \langle a_j, x \rangle = b_j \ \forall j \in [d]\}$. Observe that the vectors $\{a_j\}_{j \in [d]}$ are linearly independent. Let $\Phi := \sum_{x \in S} (-1)^{f(x)} \mu(x)$ so that $\text{bias}(f, \mu, S) = |\Phi|$ and observe that

$$\Phi = 2^{-n} \cdot \sum_{x \in S} F_\mu(x) = 2^{-n} \cdot \sum_{x \in S} \sum_{z \in \{0, 1\}^n} \widehat{F}_\mu(z) (-1)^{\langle z, x \rangle} = 2^{-n} \cdot \sum_{z \in \{0, 1\}^n} \widehat{F}_\mu(z) \sum_{x \in S} (-1)^{\langle z, x \rangle}.$$

We focus on analysing terms $T_z := \sum_{x \in S} (-1)^{\langle z, x \rangle}$. Let $V := \text{span}\{a_1, \dots, a_d\}$ and observe that whenever $z \in V$, $|T_z| = |S|$. Indeed, if $\beta_1, \dots, \beta_d \in \{0, 1\}$ is a linear combination of z in V :

$$T_z = \sum_{x \in S} (-1)^{\langle z, x \rangle} = \sum_{x \in S} \prod_{j \in [d]} (-1)^{\beta_j \langle a_j, x \rangle} = \sum_{x \in S} (-1)^{\sum_j \beta_j b_j} = |S| \cdot (-1)^{\sum_j \beta_j b_j}.$$

On the other hand, $T_z = 0$ for all $z \notin V$. Indeed, Letting $S^b = S \cap \{x \in \{0, 1\}^n : \langle x, z \rangle = b\}$ we have $T_z = |S^0| - |S^1|$. Because $z \notin V$, the constraint $\langle x, z \rangle = b$ splits S in half and thus $|S^0| = |S^1| = |S|/2$. Factoring in those observations, we get:

$$|\Phi| = 2^{-n} \cdot \left\| \sum_{z \in \{0, 1\}^n} \widehat{F}_\mu(z) \cdot T_z \right\| \leq 2^{-n} \cdot |S| \cdot \sum_{z \in V} \left\| \widehat{F}_\mu(z) \right\| \leq 2^{-n} \cdot |S| \cdot |V| \cdot \|\widehat{F}_\mu\|_\infty.$$

Recall that S has codimension d and as such $|S| = 2^{n-d}$ and $|V| = 2^d$, implying the desired inequality $\text{bias}(f, \mu, S) \leq \|\widehat{F}_\mu\|_\infty$. We now prove the third inequality of the lemma. Fix any maximum Fourier coefficient $y^* \in \{0, 1\}^n$ and observe:

$$|\widehat{F}_\mu(y^*)| = \left\| \sum_{x \in \{0, 1\}^n} F_\mu(x) \cdot (-1)^{\langle x, y^* \rangle} \cdot 2^{-n} \right\| \leq 2 \cdot \max_{b \in \{0, 1\}} \left\| \sum_{x: \langle x, y^* \rangle = b} (-1)^{f(x)} \mu(x) \right\|.$$

Fix the maximizing argument to b^* and define $S^* := \{x \in \{0, 1\}^n : \langle x, y^* \rangle = b^*\}$. Note that $S^* \in \mathcal{O}_n$ and as such:

$$\|\widehat{F}_\mu\|_\infty = |\widehat{F}_\mu(y^*)| \leq 2 \cdot \left\| \sum_{x \in S^*} (-1)^{f(x)} \mu(x) \right\| \leq 2 \cdot \max_{S \in \mathcal{O}_n} \text{bias}(f, \mu, S). \quad \blacktriangleleft$$

3.2 Step 2: Direct sum for the maximum Fourier coefficient

The outer-product of functions $F, G : \{0, 1\}^n \rightarrow \mathbb{R}$ is defined as the function $F \otimes G : \{0, 1\}^{2n} \rightarrow \mathbb{R}$ with $(F \otimes G)(x^1, x^2) := F(x^1) \cdot G(x^2)$. Next is a direct sum result for its max Fourier coefficient.

▷ **Claim 10.** For any $F, G : \{0, 1\}^n \rightarrow \mathbb{R}$, $\|\widehat{F \otimes G}\|_\infty = \|\widehat{F}\|_\infty \cdot \|\widehat{G}\|_\infty$.

Proof. Let $H = F \otimes G$; for any $z^1, z^2 \in \{0, 1\}^n$, the definition of Fourier transform implies

$$\begin{aligned} \widehat{H}(z^1, z^2) &= 2^{-2n} \cdot \sum_{x^1, x^2 \in \{0, 1\}^n} H(x^1, x^2) \cdot (-1)^{\langle x^1, z^1 \rangle + \langle x^2, z^2 \rangle} \\ &= 2^{-2n} \cdot \sum_{x^1, x^2 \in \{0, 1\}^n} F(x^1) \cdot G(x^2) \cdot (-1)^{\langle x^1, z^1 \rangle} \cdot (-1)^{\langle x^2, z^2 \rangle} \\ &= \widehat{F}(z^1) \cdot \widehat{G}(z^2). \end{aligned}$$

From this, the equivalence is immediate:

$$\|\widehat{H}\|_\infty = \max_{z^1, z^2} |\widehat{H}(z^1, z^2)| = \max_{z^1, z^2} |\widehat{F}(z^1)| \cdot |\widehat{G}(z^2)| = \|\widehat{F}\|_\infty \cdot \|\widehat{G}\|_\infty. \quad \blacktriangleleft$$

3.3 Step 3: Conclusion

We tie together Lemma 9 and Claim 10 and prove Lemma 8.

16:10 Direct Sums for Parity Decision Trees

Proof of Lemma 8. Let $H: (\{0, 1\}^n)^k \rightarrow \mathbb{R}$ be the function associated with $f^{\oplus k}$ and μ^k in Lemma 9. It is possible to express H as the k -fold outer-product of F_μ : $H = F_\mu \otimes \cdots \otimes F_\mu$. Indeed, for $x \in (\{0, 1\}^n)^k$, we have:

$$H(x) = 2^{-kn} \cdot (-1)^{f^{\oplus k}(x)} \mu^k(x) = \prod_{i \in [k]} 2^{-n} (-1)^{f(x^i)} \mu(x^i) = \prod_{i \in [k]} F_\mu(x^i).$$

Thus, using the characterisation of Lemma 9 and Claim 10 k times:

$$\max_{S \in \mathcal{S}_{kn}} \text{bias}(f^{\oplus k}, \mu^k, S) \leq \|\widehat{H}\|_\infty = \left(\|\widehat{F}_\mu\|_\infty \right)^k \leq 2^k \cdot \left(\max_{S \in \mathcal{S}_n} \text{bias}(f, \mu, S) \right)^k.$$

The XOR-lemma $\text{disc}(f^{\oplus k}, \mu^k) \geq k \cdot (\text{disc}(f, \mu) - 1)$ follows directly. We now show the other direction, $\text{disc}(f^{\oplus k}, \mu^k) \leq k \cdot \text{disc}(f, \mu)$. To do so, fix some $S \in \mathcal{S}_n$ maximizing $\text{bias}(f, \mu, S)$ and define $T \in \mathcal{S}_{kn}$ which is concatenation of k copies of S . Formally:

$$T = \{x \in (\{0, 1\}^n)^k : x^i \in S \quad \forall i \in [k]\}.$$

Now, it is easy to check that $\text{bias}(f^{\oplus k}, \mu^k, T) = \text{bias}(f, \mu, S)^k$ and the claim follows. \blacktriangleleft

4 Direct sum for D^\times part I: proof organisation

The goal of this section is to prepare the ground for a proof of our main technical contribution: a direct sum for parity trees in the distributional setting (restated below).

► **Theorem 5.** *There exists a universal constant C such that the following holds. For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ over $\{0, 1\}^n$, and $k \geq 1$,*

$$D_\varepsilon(f^k, \mu^k) \geq \Omega\left(\frac{k\delta}{\log(n/\delta)}\right) \cdot (D_{\varepsilon+\delta}(f, \mu) - C \cdot \log(n/\delta)) \quad \forall \varepsilon, \delta \geq 0.$$

As stated in Section 2, this is sufficient to prove Theorem 2 whenever $D^\times(f) \geq 6C \cdot \log(n)$. The remaining case $D^\times(f) \leq 6C \cdot \log(n)$ is proved in Lemma 38 in Section A.2. We thus focus on proving Theorem 5 in the next two sections (this section is devoted to introducing the necessary definitions and technical lemmas).

4.1 Two strengthenings of Theorem 5

For technical convenience, we study distributional complexity for *randomised* trees. For a deterministic parity tree T we let $q(T, x)$ be the number of queries made by T on input x . If \mathcal{T} is a randomised tree and μ is a distribution, we define $\bar{q}(\mathcal{T}, \mu)$ and $\text{err}_f(\mathcal{T}, \mu)$ in the natural way with:

$$\bar{q}(\mathcal{T}, \mu) := \mathbb{E}_{\substack{T \sim \mathcal{T} \\ x \sim \mu}} [q(T, x)] \quad \text{and} \quad \text{err}_f(\mathcal{T}, \mu) := \Pr_{\substack{T \sim \mathcal{T} \\ x \sim \mu}} [T(x) \neq f(x)].$$

Finally, we define $\bar{D}_\varepsilon(f, \mu) = \min_{\mathcal{T}} \{\bar{q}(\mathcal{T}, \mu) : \text{err}_f(\mathcal{T}, \mu) \leq \varepsilon\}$. It is clear that $\bar{D}_\varepsilon(f, \mu) \leq D_\varepsilon(f, \mu)$ but a converse result is more complicated, as the derandomisation can increase both the error and the depth simultaneously.

► **Claim 11.** For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, μ over $\{0, 1\}^n$ and $\varepsilon, \delta \geq 0$, $D_{\varepsilon+\delta}(f, \mu) \leq \bar{D}_\varepsilon(f, \mu)/\delta$.

We delay a proof of this folklore fact to Section A.4. We also refer readers to [36] which proves the analogue for ordinary decision trees. With this tool in hand, we can reduce Theorem 5 to the following theorem.

► **Theorem 12.** *There exists a universal constant C such that the following holds. For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ , and $k \geq 1$,*

$$\bar{D}_\varepsilon(f^k, \mu^k) \geq \Omega(k/\log(n/\gamma)) \cdot (\bar{D}_{\varepsilon+\gamma}(f, \mu) - C \cdot \log(n/\gamma)) \quad \forall \gamma \in (0, 1/n).$$

► **Definition 13.** *We say that a product distribution μ over $\{0, 1\}^n$ is λ -bounded for some $\lambda \in (0, 1]$ if $\Pr_{\mathbf{x} \sim \mu}[\mathbf{x}_i = 1] \in [\lambda/2, 1 - \lambda/2]$ for every $i \in [n]$.*

In the next sections, we also show the following qualitative improvement over Theorem 12 for bounded distributions.

► **Theorem 14.** *For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, λ -bounded distribution μ and $k \geq 1$,*

$$\bar{D}_\varepsilon(f^k, \mu^k) \geq \Omega(k\lambda) \cdot \bar{D}_\varepsilon(f, \mu).$$

Let us highlight the difference between Theorem 12 and Theorem 14: the latter is free from both the $\log n$ factor and the extra error γ . This theorem is especially interesting when the hard distribution for the function at hand (e.g. MAJ) is close to the uniform one.

4.2 The Skew measure

For the rest of this paper, we let \mathcal{U} be the uniform distribution. Let μ be a distribution over $\{0, 1\}^n$ and $S \subseteq \{0, 1\}^n$. We use $\mu(S) := \sum_{s \in S} \mu(s)$ to denote the mass of S with respect to μ . When $\mu(S) > 0$, we let μ_S be the distribution of μ conditioned on S . Let $\rho \in \{0, \star\}^n$ be a partial assignment corresponding to the sub-cube $C_\rho = \{x \in \{0, 1\}^n : \rho_i = 0 \implies x_i = 0 \ \forall i \in [n]\}$. We use μ_ρ to denote μ_{C_ρ} .

4.2.1 Random partial fixings

Let μ be a product distribution over $\{0, 1\}^n$. We say that μ is 0-biased if $\Pr_{\mathbf{x} \sim \mu}[\mathbf{x}_i = 0] \geq 1/2$ for every $i \in [n]$. For the rest of the paper, we will assume without loss of generality that any encountered input distribution is 0-biased. Indeed, should μ not be 0-biased, we can apply the following iterative transformation. Let $f_0 := f$ and $\mu_0 := \mu$. For every $i \in [n]$, if $\Pr_{\mathbf{x} \sim \mu}[\mathbf{x}_i = 1] \leq 1/2$ – the coordinate is already biased in the right direction – we simply leave $f_i := f_{i-1}$ and $\mu_i := \mu_{i-1}$. Otherwise, let:

$$\begin{aligned} f_i(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) &:= f_{i-1}(x_1, \dots, x_{i-1}, 1 - x_i, x_{i+1}, \dots, x_n); \\ \mu_i(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) &:= \mu_{i-1}(x_1, \dots, x_{i-1}, 1 - x_i, x_{i+1}, \dots, x_n). \end{aligned}$$

Observe that μ_n is 0-biased and $\bar{D}_\varepsilon(f_n^k, \mu_n^k) = \bar{D}_\varepsilon(f^k, \mu^k)$ for every $\varepsilon \geq 0$ and $k \geq 1$. Now that we are certain that μ is 0-biased, let $\delta_i := 2\Pr_{\mathbf{x} \sim \mu}[\mathbf{x}_i = 1] \in [0, 1]$. We define next the *random partial fixing* distribution with respect to μ . The intuition comes from the observation that each bit of μ can be written as a convex combination of the fixed bit “0” and a uniform bit.

► **Definition 15** (Random Partial Fixing). *The random partial fixing with respect to μ , denoted \mathcal{R}_μ , is a distribution of partial assignments $\rho \in \{0, \star\}^n$ sampled as follows: For each $i \in [n]$, we set independently*

$$\rho_i = \begin{cases} 0 & \text{w.p. } 1 - \delta_i \\ \star & \text{w.p. } \delta_i \end{cases}.$$

Observe that the following alternative two-step process is equivalent to sampling an input directly from μ . First, sample $\rho \sim \mathcal{R}_\mu$ and then sample and return $\mathbf{x} \sim \mathcal{U}_\rho$.

4.2.2 The new measure

Given a parity decision tree T and a partial assignment ρ over the input string, let T_ρ denote the pruned T by

1. fixing all the variables in the support of ρ ,
2. removing redundant queries (those can be written as a linear combination of previous queries).

For randomised parity decision tree \mathcal{T} , we define \mathcal{T}_ρ as the distribution of \mathbf{T}_ρ , where $\mathbf{T} \sim \mathcal{T}$.

► **Definition 16.** For every randomised parity decision tree \mathcal{T} and product distribution μ , define the skew average cost $\overline{sq}(\mathcal{T}, \mu) := \mathbb{E}_{\rho \sim \mathcal{R}_\mu} [\overline{q}(\mathcal{T}_\rho, \mathcal{U}_\rho)]$. Let $f: \{0, 1\}^n \rightarrow \{0, 1\}$ be a function. For $\varepsilon \geq 0$, we define the skew measure $S_\varepsilon(f)$ with:

$$S_\varepsilon(f, \mu) := \min_{\mathcal{T}} \{ \overline{sq}(\mathcal{T}, \mu) \mid \text{err}_f(\mathcal{T}, \mu) \leq \varepsilon \}.$$

▷ **Claim 17.** For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ , and $\varepsilon \geq 0$, $\overline{D}_\varepsilon(f, \mu) \geq S_\varepsilon(f, \mu)$. Furthermore, equality holds if $\mu = \mathcal{U}$.

Proof. The claim is immediate as $\overline{sq}(\mathcal{T}, \mu) \leq \overline{q}(\mathcal{T}, \mu)$ for every randomised parity tree \mathcal{T} and product distribution μ . ◀

4.3 Proof plan

The proofs of Theorems 12 and 14 are carried out in two steps. First, we prove a perfect direct sum for the skew measure in Section 5.

► **Theorem 18.** We have $S_\varepsilon(f^k, \mu^k) \geq k \cdot S_\varepsilon(f, \mu)$ for any function f , product μ and $\varepsilon \geq 0$.

As a second step, we demonstrate in Section 6 that $\overline{D}_\varepsilon(f, \mu) \approx S_\varepsilon(f, \mu)$. We first prove a lossless conversion for product distribution which are *constant-bounded*. We then extend this to general product distributions for which we lose a $\log(n)$ -factor. Let us recall here that the $\log n$ loss for general (unbounded) product distribution is inherent to the skew measure. Indeed, we show in Section 8 the existence of some f and μ for which $\overline{D}_{1/3}(f, \mu) = \Theta(\log n)$ but $S_0(f, \mu) = \Theta(1)$.

► **Theorem 19.** For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ , $\gamma \in (0, 1/n)$, we have

$$\overline{D}_{\varepsilon+\gamma}(f, \mu) \leq O(\log(n/\gamma)) \cdot (S_\varepsilon(f, \mu) + 1) \quad \forall \varepsilon \geq 0.$$

► **Theorem 20.** For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and λ -bounded product distribution μ , we have

$$\overline{D}_\varepsilon(f, \mu) \leq O(1/\lambda) \cdot S_\varepsilon(f, \mu) \quad \forall \varepsilon \geq 0.$$

Combining the results above it is now straightforward to conclude and prove Theorems 12 and 14. For instance, the proof of the former goes as follows.

Proof of Theorem 12.

$$\begin{aligned} \overline{D}_\varepsilon(f^k, \mu^k) &\geq S_\varepsilon(f^k, \mu^k) && \text{(Claim 17)} \\ &\geq k \cdot S_\varepsilon(f, \mu) && \text{(Theorem 18)} \\ &\geq \Omega(k/\log(n/\gamma)) \cdot (\overline{D}_{\varepsilon+\gamma}(f, \mu) - C \cdot \log(n/\gamma)). && \text{(Theorem 19)} \end{aligned} \quad \blacktriangleleft$$

4.4 Some notation

Let us finish this section by defining some notations which will be useful for the rest of the paper. Let T be a parity decision tree on input $\{0, 1\}^n$. We define $\mathcal{N}(T)$ as the set of nodes of T and $\mathcal{L}(T)$ as the set of leaves of T . For each node $v \in \mathcal{N}(T)$, we define the following: (items marked with $*$ are only defined for non-leaf nodes)

- $\text{path}(v)$: the set of nodes on the root-to- v path (including the root, excluding v)
- $d(v) := |\text{path}(v)|$: the depth of v
- * $Q^v \in \{0, 1\}^n$: the query made at node v
- * $\text{child}(v, b)$ the child of v corresponding to the query outcome $\langle x, Q^v \rangle = b$, where $b \in \{0, 1\}$
- $Q^{\prec v}$: an $n \times d(v)$ boolean matrix with column vectors $\{Q^u\}_{u \in \text{path}(v)}$
- * $Q^{\preceq v} := [Q^{\prec v} \ Q^v]$ of dimension $n \times (d(v) + 1)$.
- $b^{\prec v} \in \{0, 1\}^{d(v)}$: the labels on the root-to- v path

For every boolean matrix $A \in \{0, 1\}^{n \times m}$, we use $\text{rank}(A)$ to denote the rank of A (understood as a matrix over \mathbb{F}_2) and let $\text{col}(A) \subseteq \{0, 1\}^n$ be the column space of A . For every $S \subseteq [n]$, let $A_S \in \{0, 1\}^{|S| \times m}$ stand for the sub-matrix of A consisting of row with indices in S . For every $x, y \in \{0, 1\}^n$ and $S \subseteq [n]$, we denote $\langle x_S, y_S \rangle = \sum_{i \in S} x_i y_i$ by $\langle x, y \rangle_S$.

Let μ and ν be two distributions over S . We use $d_{\text{TV}}(\mu, \nu) := \sup_{S' \subseteq S} |\mu(S') - \nu(S')|$ to denote the total variation distance between μ and ν and write $\mu \equiv \nu$ if $d_{\text{TV}}(\mu, \nu) = 0$.

5 Direct sum for D^\times part II: direct sum for S

In this section, we prove a perfect direct sum for S (restated below). A direct consequence of this fact is a perfect direct sum for distributional parity query complexity under the uniform distribution.

► **Theorem 18.** *We have $S_\varepsilon(f^k, \mu^k) \geq k \cdot S_\varepsilon(f, \mu)$ for any function f , product μ and $\varepsilon \geq 0$.*

► **Corollary 21.** *We have $\overline{D}_\varepsilon(f^k, \mathcal{U}^k) \geq k \cdot \overline{D}_\varepsilon(f, \mathcal{U})$ for any function f and $\varepsilon \geq 0$.*

Proof. Combine Claim 17 with Theorem 18. ◀

To prove Theorem 18, our overall strategy is to take a tree achieving $S_\varepsilon(f^k, \mu^k)$ and extract a tree computing a single copy of f under μ to within error ε while having cost bounded by $S_\varepsilon(f^k, \mu^k)/k$. To do so, we employ the extraction strategy hinted at in Section 2. The extractor works as long as the input distributions are uniform, which is the case after the random partial fixing step of S.

5.1 Extracting a single instance under uniform distributions

Let T be a deterministic parity tree taking inputs $x \in \mathcal{X} := \{0, 1\}^{m_1} \times \dots \times \{0, 1\}^{m_k}$ and returning labels in $\{0, 1\}^k$. We assume without loss of generality that the queries along any root-to-leaf path are linearly independent. Let $L(\ell) \in \{0, 1\}^k$ be the label associated with the leaf $\ell \in \mathcal{L}(T)$. For $i \in [k]$, we define the linear subspace $W_i \subseteq \mathcal{X}$ of query vectors that are zero everywhere except for copy i :

$$W_i := \{w \in \mathcal{X} : w^j = 0^{m_j} \iff j \neq i\}.$$

We say a node $v \in \mathcal{N}(T)$ is *critical* with respect to i if $\text{col}(Q^{\prec v}) \cap W_i \neq \text{col}(Q^{\preceq v}) \cap W_i$ and denote the set of critical indices at node v with $I_v := \{i \in [k] : v \text{ is critical w.r.t. } i\}$. Finally, we let $d_i(v) := \sum_{u \in \text{path}(v)} \mathbb{1}[i \in I_u]$ be the relative depth of v with respect to instance i and highlight that $d_i(v) = \dim(\text{col}(Q^{\prec v}) \cap W_i)$. The algorithm $\text{Ext}_i(T)$ which

Algorithm 1 $\text{Ext}_i(T)$.

Input: $y \in \{0, 1\}^{m_i}$ **Output:** $a \in \{0, 1\}$

```

1: Initialize  $v \leftarrow$  root of  $T$ 
2: while  $v$  is not a leaf do
3:   if  $i \in I^v$  then
4:     Let  $w$  be any vector in  $(\text{col}(Q^{\preceq v}) \setminus \text{col}(Q^{\prec v})) \cap W_i$ , query  $\langle y, w \rangle$ 
5:     Compute  $b^v := \langle y, Q^v \rangle$  from  $b^{\prec v}$  and  $\langle y, w \rangle$ 
6:     Move  $v \leftarrow \text{child}(v, b^v)$ 
7:   else
8:     Sample  $\xi \sim \text{Ber}(1/2)$ 
9:     Move  $v \leftarrow \text{child}(v, \xi)$ 
10:  end if
11: end while
12: return  $L_i(v)$ 

```

extracts a tree for the i -th instance out of T is described in Algorithm 1. Observe that it is indeed possible to compute the value of $\langle y, Q^v \rangle$ from $b^{\prec v}$ and $\langle y, w \rangle$ on line 5: Since $w \notin \text{col}(Q^{\prec v})$, we have $\text{rank}([Q^{\prec v} \ w]) = \text{rank}(Q^{\prec v}) + 1$. On the other hand, as $w \in \text{col}(Q^{\preceq v})$, we have $\text{rank}([Q^{\preceq v} \ w]) = \text{rank}(Q^{\preceq v}) = \text{rank}(Q^{\prec v}) + 1$. Thus $Q^v \in \text{col}([Q^{\prec v} \ w])$, which means that Q^v can be written as a linear combination of the columns of $[Q^{\prec v} \ w]$: $Q^v = Q^{u_1} + \dots + Q^{u_t} + w$ where u_1, \dots, u_t are some ancestors of v . This in turn implies that $\langle y, Q^v \rangle = \sum_{i \in [t]} \langle y, Q^{u_i} \rangle + \langle y, w \rangle$.

We stress that although T is a deterministic tree, $\text{Ext}_i(T)$ is a randomized decision tree with internal randomness inherited from the bits ξ . Our main technical claim is that for any fixed $y \in \{0, 1\}^{m_i}$, the algorithm $\text{Ext}_i(T)$ perfectly simulates a run of T when the input is on a random input $\mathbf{x} = (\mathbf{x}^1, \dots, \mathbf{x}^{i-1}, y, \mathbf{x}^{i+1}, \dots, \mathbf{x}^k)$ and $\mathbf{x}^j \sim \mathcal{U}(\{0, 1\}^{m_j})$. In a nutshell, the randomness of the other $k - 1$ instances can be substituted with the internal randomness ξ . To make this precise, we let $X_v = \{x \in \mathcal{X} : x^T Q^{\prec v} = b^{\prec v}\}$ be the set of inputs leading to the node $v \in \mathcal{N}(T)$.

▷ **Claim 22.** For any $y \in \{0, 1\}^{m_i}$, $\Pr_{\xi}[\text{Ext}_i(T) \text{ reaches node } v \text{ in its execution on } y] = \Pr_{\mathbf{x}}[\mathbf{x} \in X^v]$.

Proof. Let us fix $i := 1$ and $d := d(v)$ for simplicity. We establish an alternative description of X^v that puts pure constraints on instance 1 first. Pick $t := d_1(v)$ independent vectors $Q_1, \dots, Q_t \in \text{col}(Q^{\prec v}) \cap W_1$ and extend them arbitrarily to a basis $\{Q_j\}_{j \in [d]}$ of $Q^{\prec v}$. As each vector of this basis can be expressed as a linear combination of $\{Q^u\}_{u \in \text{path}(v)}$, it is possible to apply those linear combinations to $b^{\prec v}$ and obtain values $\{b_j\}_{j \in [d]}$ such that $X_v = \{x \in \mathcal{X} \mid \forall j \in [d] : \langle x, Q_j \rangle = b_j\}$. The set $Y^v \subseteq \{0, 1\}^{m_1}$ of inputs that can reach node v in a run of $\text{Ext}_1(T)$ thus corresponds to

$$Y^v := \{y \in \{0, 1\}^{m_1} \mid \forall j \in [t] : \langle y, Q_j^1 \rangle = b_j\}.$$

If $y \notin Y^v$, the statement follows directly as both probabilities are zero. However, if $y \in Y^v$,

$$\Pr_{\xi}[\text{Ext}_1(T) \text{ reaches node } v \text{ in its execution}] = 2^{-d+t}.$$

This is so because a node v can only be reached by having the “right” $d - t$ coin tosses of ξ (provided that $y \in Y^v$). Thus, it remains to show that $\Pr_{\mathbf{x}}[\mathbf{x} \in X^v] = 2^{-d+t}$ if $y \in Y^v$.

Let $m = \sum_{i \in [k]} m_i$ and $S = \{m_1 + 1, \dots, m\}$ be the indices of the bits of every copy but the first one. Fix the $m \times (d - t)$ boolean matrix $A = [Q_{t+1} \cdots Q_d]$ and observe that $\text{rank}(A) = d - t$ by construction. We show that $\text{rank}(A_S) = d - t$ too. If $\text{rank}(A_S) < \text{rank}(A)$, we can find a non-empty set $J \subseteq \{t + 1, \dots, d\}$ such that $\sum_{j \in J} (Q_j)_S = 0$. This implies that $Q' := \sum_{j \in J} Q_j \in W_i \cap \text{col}(Q^{\prec v})$. But Q' is linearly independent of $\{Q_1, \dots, Q_t\}$ – this contradicts $\dim(\text{col}(Q^{\prec v}) \cap W_i) = t$. Therefore, if $y \in Y^v$, we use this observation to conclude:

$$\begin{aligned}
\Pr_{\mathbf{x}}[\mathbf{x} \in X^v] &= \Pr_{\mathbf{x}}[\forall j \in [d] : \langle \mathbf{x}, Q_j \rangle = b_j] \\
&= \Pr_{\mathbf{x}}[\mathbf{x}^T A = (b_j)_{t+1 \leq j \leq d}] \\
&= \Pr_{\mathbf{z} := (\mathbf{x}^2, \dots, \mathbf{x}^k)}[\mathbf{z}^T A_S = (b_j + \langle y, Q_j^1 \rangle)_{t+1 \leq j \leq d}] \\
&= 2^{-\text{rank}(A_S)} \\
&= 2^{-d+t}.
\end{aligned}$$

◁

5.2 Proof of Theorem 18

We are now ready to show Theorem 18. Let \mathcal{T} be a randomised parity decision tree which witnesses $C := S_\varepsilon(f^k, \mu^k)$. For each $i \in [k]$, define the randomized decision tree $\mathcal{T}_i: \{0, 1\}^n \rightarrow \{0, 1\}$ with:

1. Sample $T \sim \mathcal{T}$.
2. Sample $\rho^1, \dots, \rho^{i-1}, \rho^{i+1}, \dots, \rho^k \sim \mathcal{R}_\mu$.
3. Let $\tilde{\rho} := (\rho^1, \dots, \rho^{i-1}, \star^n, \rho^{i+1}, \dots, \rho^k)$.
4. Return $\text{Ext}_i(T_{\tilde{\rho}})$.

We show in Lemma 23 that $\text{err}_f(\mathcal{T}_i, \mu) \leq \varepsilon$ simultaneously for all $i \in [k]$. On the other hand, we show in Lemma 24 that $\sum_{i \in [k]} \overline{\text{sq}}(\mathcal{T}_i, \mu) \leq C$. By an averaging argument, this shows the existence of a copy $i^* \in [k]$ with $\text{cost} \leq C/k$ and therefore $S_\varepsilon(f, \mu) \leq C/k$. The remainder of this section is devoted to proving both claims.

► **Lemma 23.** *For every $i \in [k]$, $\text{err}_f(\mathcal{T}_i, \mu) \leq \text{err}_{f^k}(\mathcal{T}, \mu^k)$.*

Proof. It is enough to prove the statement assuming \mathcal{T} is a deterministic parity tree T and $i = 1$. Let \mathcal{R} be the distribution of $\tilde{\rho}$ in the step 3 of generating \mathcal{T}_1 . Fix some $\rho \in \text{supp}(\mathcal{R})$ and note that $\rho^1 = \star^n$. We also define $\mathcal{U}^{-1} := \mathcal{U}_{\rho^2} \times \cdots \times \mathcal{U}_{\rho^k}$. Using Claim 22 on a leaf $\ell \in \mathcal{L}(T_\rho)$ yields:

$$\begin{aligned}
&\Pr_{\mathbf{y}, \xi}[\text{Ext}_1(T_\rho) \text{ reaches } \ell \text{ on } \mathbf{y} \wedge L_1(\ell) \neq f(\mathbf{y})] \\
&= \mathbb{E}_{\mathbf{y}} \left[\Pr_{\xi}[\text{Ext}_1(T_\rho) \text{ reaches } \ell \text{ on } \mathbf{y}] \cdot \mathbb{1}[L_1(\ell) \neq f(\mathbf{y})] \right] \\
&= \mathbb{E}_{\mathbf{y}} \left[\Pr_{\mathbf{x}^{-1} \sim \mathcal{U}^{-1}}[(\mathbf{y}, \mathbf{x}^{-1}) \in X^\ell] \cdot \mathbb{1}[L_1(\ell) \neq f(\mathbf{y})] \right] \\
&= \Pr_{\mathbf{x} \sim \mu \times \mathcal{U}^{-1}}[\mathbf{x} \in X^\ell \wedge L_1(\ell) \neq f(\mathbf{x}^1)].
\end{aligned}$$

16:16 Direct Sums for Parity Decision Trees

Thus:

$$\begin{aligned}
\text{err}_f(\mathcal{T}_1, \mu) &= \mathbb{E}_{\tilde{\rho} \sim \mathcal{R}} \left[\Pr_{\mathbf{y} \sim \mu, \xi} [\text{Ext}_1(T_{\tilde{\rho}})(\mathbf{y}) \neq f(\mathbf{y})] \right] \\
&= \mathbb{E}_{\tilde{\rho}} \left[\sum_{\ell \in \mathcal{L}(T_{\tilde{\rho}})} \Pr_{\mathbf{x} \sim \mu \times \mathcal{U}^{-1}} [\mathbf{x} \in X^\ell \wedge L_1(\ell) \neq f(\mathbf{x}^1)] \right] \\
&\leq \mathbb{E}_{\tilde{\rho}} \left[\sum_{\ell \in \mathcal{L}(T_{\tilde{\rho}})} \Pr_{\mathbf{x} \sim \mu \times \mathcal{U}^{-1}} [\mathbf{x} \in X^\ell \wedge L(\ell) \neq f(\mathbf{x})] \right] \\
&= \mathbb{E}_{\tilde{\rho}} \left[\text{err}_{f^k}(T_{\tilde{\rho}}, \mu \times \mathcal{U}^{-1}) \right].
\end{aligned}$$

Observe now that for any $x \in \text{supp}(\mu \times \mathcal{U}^{-1})$, we have $T_{\tilde{\rho}}(x) = T(x)$. Using the definition of \mathcal{R}_μ thus yields:

$$\text{err}_f(\mathcal{T}_1, \mu) \leq \mathbb{E}_{\tilde{\rho} \sim \mathcal{R}} [\text{err}_{f^k}(T_{\tilde{\rho}^{-1}}, \mu \times \mathcal{U}^{-1})] = \text{err}_{f^k}(T, \mu^k). \quad \blacktriangleleft$$

► **Lemma 24.** $\sum_{i \in [k]} \overline{sq}(\mathcal{T}_i, \mu) \leq \overline{sq}(\mathcal{T}, \mu^k).$

Proof. It is sufficient to prove this for the case where \mathcal{T} is a deterministic tree T . We have:

$$\begin{aligned}
\sum_{i \in [k]} \overline{sq}(\mathcal{T}_i, \mu) &= \sum_{i \in [k]} \mathbb{E}_{\rho^i \sim \mathcal{R}_\mu} [\overline{q}((\mathcal{T}_i)_{\rho^i}, \mathcal{U}_{\rho^i})] \\
&= \sum_{i \in [k]} \mathbb{E}_{\substack{\rho^i \sim \mathcal{R}_\mu \\ \tilde{\rho} \sim \mathcal{R}}} \left[\overline{q} \left(\left(\text{Ext}_i(T_{\tilde{\rho}}) \right)_{\rho^i}, \mathcal{U}_{\rho^i} \right) \right] \\
&= \sum_{i \in [k]} \mathbb{E}_{\rho \sim \mathcal{R}_\mu^k} [\overline{q}(\text{Ext}_i(T_\rho), \mathcal{U}_{\rho^i})] \\
&= \mathbb{E}_{\rho \sim \mathcal{R}_\mu^k} \left[\sum_{i \in [k]} \overline{q}(\text{Ext}_i(T_\rho), \mathcal{U}_{\rho^i}) \right].
\end{aligned}$$

where the third equality is due to the fact that the operations of applying Ext and fixing variables are commutable. Let $\rho \in (\{0, \star\}^n)^k$ be a partial fixing and $\ell \in \mathcal{L}(T_\rho)$. The probability that node ℓ is visited during the process $\text{Ext}_i(T_\rho)$ when the input is $\mathbf{x}^i \sim \mathcal{U}_{\rho^i}$ is $2^{-d(\ell)}$. Observe that $\text{Ext}_i(T_\rho)$ only makes $d_i(\ell)$ queries to \mathbf{x}^1 to reach ℓ . As such, we have:

$$\begin{aligned}
\sum_{i \in [k]} \overline{q}(\text{Ext}_i(T_\rho), \mathcal{U}_{\rho^i}) &= \sum_{i \in [k]} \sum_{\ell \in \mathcal{L}(T')} 2^{-d(\ell)} d_i(\ell) \\
&\leq \sum_{\ell \in \mathcal{L}(T')} 2^{-d(\ell)} d(\ell) \\
&= \overline{q}(T_\rho, \mathcal{U}_\rho).
\end{aligned}$$

The inequality is due to the fact that $\sum_{i \in [k]} d_i(v) \leq d(v)$. This is because $\dim(W_i \cap W_j) = 0$ for each $i \neq j$ and so

$$\sum_{i \in [k]} d_i(v) = \sum_{i \in [k]} \dim(\text{col}(Q^{\prec v}) \cap W_i) \leq \dim(\text{col}(Q^{\prec v})) = d(v).$$

To conclude, we have

$$\sum_{i \in [k]} \overline{sq}(\mathcal{T}_i, \mu) = \mathbb{E}_{\rho \sim \mathcal{R}_\mu^k} \left[\sum_{i \in [k]} \overline{q}(\text{Ext}_i(T_\rho), \mathcal{U}_{\rho^i}) \right] \leq \mathbb{E}_{\rho \sim \mathcal{R}_\mu^k} [\overline{q}(T_\rho, \mathcal{U}_\rho)] = \overline{sq}(T, \mu^k). \quad \blacktriangleleft$$

6 Direct sum for D^\times part III: from S to D^\times

In this section, we show how to convert parity tree of the S_ε model to the more common \overline{D}_ε model and prove Theorems 19 and 20. Let us fix for this section a boolean function $f: \{0,1\}^n \rightarrow \{0,1\}$ together with some 0-biased product distribution μ over $\{0,1\}^n$. Let T be a deterministic parity tree trying to solve f against μ . We begin by establishing an alternative view of the quantity $\overline{sq}(T, \mu)$. For any fixed $x \in \{0,1\}^n$, define the product distribution R_μ^x over $\{0, \star\}^n$ with:

$$\Pr_{\rho \sim R_\mu^x}[\rho_i = \star] = \begin{cases} \delta_i/(2 - \delta_i) & \text{if } x_i = 0 \\ 1 & \text{if } x_i = 1 \end{cases} \quad \text{where } \delta_i := 2 \cdot \Pr_{x \sim \mu}[x_i = 1] \in [0, 1]. \quad (3)$$

Sampling $\rho \sim R_\mu$, $x \sim \mathcal{U}_\rho$ and completing $x_j = 0$ for all $\rho = 0$ is equivalent to first sampling $x \sim \mu$ and then some $\rho \sim R_\mu^x$. One can therefore see the process of $\overline{sq}(T, \mu)$ as follows:

1. Sample $x \sim \mu$, $\rho \sim R_\mu^x$.
2. Run T on x .
3. Every time T attempts to make a query, check if ρ simplifies the query: $\rho_i = 0 \implies x_i = 0$.

We describe this alternative view in detail in Algorithm 2. With this new interpretation, we can recast the quantity $\overline{sq}(T, \mu)$ with

$$\overline{sq}(T, \mu) = \mathbb{E}_{x \sim \mu, \rho \sim R_\mu^x}[\text{Number of times Section 6 is executed in Algorithm 2}]. \quad (4)$$

The idea to convert S_ε algorithms to \overline{D}_ε ones is to simulate the process of Algorithm 2 by maintaining an incomplete but consistent view $p \in \{0, \star, ?\}^n$ of ρ . Initially, $p = ?^n$ – i.e. nothing is known about ρ – and we gradually update p based on the queries we get. For instance, if $x_i = 1$, then (3) asserts $\rho_i = \star$. This scheme helps to relate the cost of the converted \overline{D}_ε algorithm with $\overline{sq}(T, \mu)$. The description of the converted algorithm is given in Algorithm 3.

► **Definition 25.** Let $p \in \{0, \star, ?\}^n$ be a fixing. The following are subsets of indices:

$$S_\star^p = \{j \in [n] : p_j = \star\} \quad S_0^p = \{j \in [n] : p_j = 0\} \quad S_?^p = \{j \in [n] : p_j = ?\} \quad S_{\neq 0}^p = S_\star^p \cup S_?^p$$

We also write $S(p, \star)$ to mean S_\star^p and likewise for other sets.

Let $P^v \subseteq \{0, \star, ?\}^n$ be the set of all possible p that could be at the start of an iteration of Algorithm 3 at node v . We now prove an invariant of Algorithm 3 and then its correctness.

► **Lemma 26.** For any state $v \in \mathcal{N}(T)$ and $p \in P^v$ that Algorithm 3 could be in at the start of a while iteration (Section 6), it holds that:

$$\text{rank}(Q_{S(p, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p, \star)}^{\prec v}) = |S(p, \star)|.$$

Proof. We prove the claim by induction on T . The statement is true when v is the root because both $Q^{\prec v}$ and S_\star^p are empty. Let us now assume that the statement is true for some v and $p \in P^v$ and prove that the invariant carries over to the next iteration regardless of the query outcomes and the randomness η of the process. If p' is the updated value of p at Section 6, this amounts to showing that $\text{rank}(Q_{S(p', \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p', \star)}^{\prec v}) = |S(p', \star)|$. We consider three cases.

Case $D^{v,p} = \emptyset$. Then, there is no update for p and $p' = p$. Since $\text{rank}(Q_{S(p, \star)}^{\prec v}) = \text{rank}(Q_{S(p, \star)+j}^{\prec v})$ for all $j \in S(p, \neq 0)$, we have $\text{rank}(Q_{S(p, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p, \star)}^{\prec v}) = |S(p, \star)|$, as desired.

16:18 Direct Sums for Parity Decision Trees

■ **Algorithm 2** an alternative view of $\overline{sq}(T, \mu)$.

Input: $x \in \{0, 1\}^n, \rho \in \{0, \star\}^n$

Output: $a \in \{0, 1\}$

```

1:  $v \leftarrow$  root of  $T$ 
2: while  $v$  is not a leaf do
3:   if  $\text{rank}(Q_{\bar{S}(\rho, \star)}^{\prec v}) = \text{rank}(Q_{S(\rho, \star)}^{\prec v}) + 1$  then
4:     Query  $b^v \leftarrow \langle x, Q^v \rangle$ 
5:   else
6:     Infer  $b^v \leftarrow \langle x, Q^v \rangle$  from the fact that  $(Q^{\prec v})^T x = b^{\prec v}$  and  $x_j = 0$  for all  $\rho_j = 0$ 
7:   end if
8:   Move  $v \leftarrow \text{child}(v, b^v)$ .
9: end while
10: return  $L(v)$ 

```

■ **Algorithm 3** converts an algorithm T for S_ε to \overline{D}_ε .

Input: $x \in \{0, 1\}^n$

Output: $a \in \{0, 1\}$

```

1:  $v \leftarrow$  root of  $T$ 
2:  $p \leftarrow ?^n$ 
3: while  $v$  is not a leaf do
4:    $D^{v,p} \leftarrow \{j \in [n] : p_j = ? \text{ and } \text{rank}(Q_{\bar{S}(p, \star)+j}^{\prec v}) = \text{rank}(Q_{S(p, \star)}^{\prec v}) + 1\}$ 
5:   if  $D^{v,p} = \emptyset$  then
6:     Infer  $b^v \leftarrow \langle x, Q^v \rangle$  from the fact that  $(Q^{\prec v})^T x = b^{\prec v}$  and  $x_j = 0$  for all  $p_j = 0$ 
7:   else
8:     for  $j \in D^{v,p}$  do
9:       Query  $x_j$ 
10:      Sample  $\eta \sim \text{Ber}(\delta_j / (2 - \delta_j))$ 
11:      if  $x_j = 1$  or  $\eta = 1$  then
12:         $p_j \leftarrow \star$ 
13:        break
14:      end if
15:       $p_j \leftarrow 0$ 
16:    end for
17:    Query  $b^v \leftarrow \langle x, Q^v \rangle$ 
18:  end if
19:  Move  $v \leftarrow \text{child}(v, b^v)$ 
20: end while
21: return  $L(v)$ 

```

Case $D^{v,p} \neq \emptyset$ and $p'_j = 0$ for all $j \in D^{v,p}$. Then, $S_\star^{p'} = S_\star^p$ and $S(p', \neq 0) = S(p, \neq 0) \setminus D^{v,p}$. By definition of $D^{v,p}$, we still have $\text{rank}(Q_{S(p,\neq 0)}^{\prec v}) = \text{rank}(Q_{S(p,\neq 0)+j}^{\prec v})$ for all $j \in S(p', \neq 0)$, so $\text{rank}(Q_{S(p', \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p', \neq 0)}^{\prec v}) = |S(p', \neq 0)|$.

Case $D^{v,p} \neq \emptyset$ and $p'_j = \star$ for some $j \in D^{v,p}$. Then $S_\star^{p'} = S_\star^p + j$ and it must hold that $\text{rank}(Q_{S(p', \neq 0)}^{\prec v}) = |S(p', \neq 0)|$. On the other hand,

$$\text{rank}(Q_{S(p', \neq 0)}^{\prec v}) \leq \text{rank}(Q_{S(p, \neq 0)}^{\prec v}) + 1 = |S(p, \neq 0)| + 1 = |S(p', \neq 0)|.$$

Where the inequality follows from the fact that $S(p', \neq 0) \subseteq S(p, \neq 0)$. Finally, this implies $\text{rank}(Q_{S(p', \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p, \neq 0)}^{\prec v}) = |S(p, \neq 0)|$. \blacktriangleleft

► **Lemma 27.** For any $x \in \{0, 1\}^n$, $\Pr_\eta[\text{Algorithm 3 outputs } 1] = \mathbb{1}[T(x) = 1]$.

Proof. It is not hard to see that if Algorithm 3 gets the correct value of $\langle x, Q^v \rangle$ at each iteration of the while loop, it perfectly simulates T . Thus, it suffices to show that whenever $D^{v,p} = \emptyset$, the algorithm can compute the value of $\langle x, Q^v \rangle$ from the previous query outcomes. Lemma 26 and its proof implies that if $D^{v,p} = \emptyset$, then $\text{rank}(Q_{S(p, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p, \neq 0)}^{\prec v}) = |S(p, \neq 0)|$. Thus $Q_{S(p, \neq 0)}^v$ can be written as a linear combination of column vectors of $Q_{S(p, \neq 0)}^{\prec v}$. Namely, $Q_{S(p, \neq 0)}^v = \sum_{j \in [t]} Q_{S(p, \neq 0)}^{v_j}$, where v_1, \dots, v_t are some ancestors of v . On the other hand, we know that $x_j = p_j = 0$ for all $j \in S_0^p$. Consequently, we have

$$\langle x, Q^v \rangle = \langle x, Q^v \rangle_{S(p, \neq 0)} = \sum_{j \in [t]} \langle x, Q^{v_j} \rangle_{S(p, \neq 0)} = \sum_{j \in [t]} b^{v_j}.$$

Thus, Algorithm 3 follows the same path of vertices as T , irrespective of the randomness η . Consequently, its outputs correspond to the one of T . \blacktriangleleft

We now turn our attention to the efficiency of Algorithm 3. We shall start with the special case of μ being a constant-bounded distribution. In this particular case, we obtain a lossless conversion. We then turn our attention to general product distributions, for which Algorithm 3 suffers a $\log(n)$ factor. This loss factor is inherent to reducing S_ϵ to \bar{D} as Section 8 shows.

6.1 Conversion for constant-bounded distribution

We now prove a strong efficiency result for Algorithm 3 in the special case where μ is λ -bounded (see Definition 13). A proof of our goal (Theorem 20) then follows easily.

► **Lemma 28.** We have $\bar{q}(\text{Algorithm 3 on } T, \mu) \leq (2/\lambda) \cdot \bar{sq}(T, \mu)$.

Before proving this, we need an alternative view of the randomness used in the for-loop of Algorithm 3 (line 8 to 16). At the start of the process, a random partial fixing $\rho \sim \mathcal{R}_\mu^x$ is generated. The algorithm is then deterministic: whenever some x_j is queried in the for-loop, this is replaced by a query to ρ_j . The algorithm updates p_j with ρ_j and exits the loop if $\rho_j = \star$. This process is given in detail in Algorithm 4. Note that as \mathcal{R}_μ^x is a product distribution, one can actually implement Algorithm 4 without querying all of x at the start. Indeed, it is enough to query x_j whenever one needs the value of ρ_j , similarly to Algorithm 3. This implies that both processes are equivalent.

Suppose one runs Algorithm 4 on $x \sim \mu$ and $\rho \sim \mathcal{R}_\mu^x$. Fix some state (v, p) the algorithm could be in at the start of the while loop (Section 6.1). We let $\mathcal{X}^{v,p}$ be the distribution of x conditioned on reaching state (v, p) . Furthermore, for a fixed $x \in \{0, 1\}^n$ and (v, p) reachable with x we let $\mathcal{R}^{v,p,x}$ be the marginal distribution of ρ conditioned on reaching state (v, p) and $x = x$. We now develop explicit formulations for those distributions.

■ **Algorithm 4** an alternative view of Algorithm 3 where the randomness is fixed at the start.

Input: $x \in \{0, 1\}^n$
Output: $a \in \{0, 1\}$

```

1:  $v \leftarrow \text{root of } T$ 
2:  $p \leftarrow ?^n$ 
3: Sample  $\rho \sim R_\mu^x$ 
4: while  $v$  is not a leaf do
5:    $D^{v,p} \leftarrow \{j \in [n] : p_j = ? \wedge \text{rank}(Q_{S(p,*)+j}^{\prec v}) = \text{rank}(Q_{S(p,*)}^{\prec v}) + 1\}$ 
6:   if  $D^{v,p} = \emptyset$  then
7:     Infer  $b^v \leftarrow \langle x, Q^v \rangle$  from the fact that  $(Q^{\prec v})^T x = b^{\prec v}$  and  $x_j = 0$  for all  $p_j = 0$ 
8:   else
9:     for  $j \in D^{v,p}$  do
10:       $p_j \leftarrow \rho_j$ 
11:      if  $p_j = \star$  then
12:        break
13:      end if
14:    end for
15:    Query  $b^v \leftarrow \langle x, Q^v \rangle$ 
16:  end if
17:  Move  $v \leftarrow \text{child}(v, b^v)$ .
18: end while
19: return  $L(v)$ 

```

Explicit definition of $\mathcal{X}^{v,p}$

Let $\hat{\mathcal{X}}^{v,p}$ be the distribution over $\{0, 1\}^n$ defined as follows:

1. For all $j \in S_0^p$, fix $x_j = 0$.
2. For all $j \in S_?^p$, sample $x_j \sim \text{Ber}(\delta_j/2)$.
3. Determine $\{x_j : j \in S_\star^p\}$ by solving $\{\langle x, Q^u \rangle_{S(p,*)} = \langle x, Q^u \rangle_{S(p,\neq\star)} + b^u\}_{u \in \text{path}(v)}$

Explicit definition of $\mathcal{R}^{v,p,x}$

Let $\hat{\mathcal{R}}^{p,x}$ be the product distribution over $\{0, \star\}^n$ defined as follows:

1. For all $j \in S_?^p$ such that $x_j = 0$, let $\rho_j = \star$ with probability $\delta_j/(2 - \delta_j)$ and $\rho_j = 0$ else.
2. For all $j \in S_?^p$ such that $x_j = 1$, fix $\rho_j = \star$.
3. For all $j \in S(p, \neq ?)$, fix $\rho_j = p_j$.

▷ **Claim 29.** For every reachable state (v, p) and $x \in \text{supp}(\mathcal{X}^{v,p})$ in Algorithm 4, we have

1. $\mathcal{R}^{v,p,x} \equiv \hat{\mathcal{R}}^{p,x}$;
2. $\mathcal{X}^{v,p} \equiv \hat{\mathcal{X}}^{v,p}$.

We delay the proof of this technical lemma to Section A.5. We can now prove the efficiency of our algorithm for λ -bounded distributions.

Proof of Lemma 28. To relate Algorithm 2 with Algorithm 4, it is helpful to insert the book-keeping of p in Algorithm 2 (lines 5 to 16, without 10) in between Sections 6 and 6.1 of Algorithm 2. This doesn't change the number of queries or guarantees of Algorithm 2 but now both processes share the same state space over (v, p) . For $x \in \{0, 1\}^n$ and $\rho \in \{0, \star\}^n$, define $A(x, \rho)$ and $B(x, \rho)$ as the number of queries each process makes:

$A(x, \rho) :=$ number of times Section 6 is executed in Algorithm 2 on input (x, ρ) ;

$B(x, \rho) :=$ number of times Section 6.1 are executed in Algorithm 4 on input (x, ρ) .

Using (4), it is thus enough to prove that $\mathbb{E}_{\mathbf{x}, \rho} [A(\mathbf{x}, \rho)] \geq \Omega(\lambda) \cdot \mathbb{E}_{\mathbf{x}, \rho} [B(\mathbf{x}, \rho)]$ when $\mathbf{x} \sim \mu$ and $\rho \sim R_\mu^x$. We have:

$$\mathbb{E}_{\mathbf{x}, \rho} [A(\mathbf{x}, \rho)] = \sum_{(v, p)} \Pr[(v, p) \text{ is reached}] \cdot \Pr_{\substack{\mathbf{x} \sim \mathcal{X}^{v, p} \\ \rho \sim \mathcal{R}^{v, p, \mathbf{x}}}} \left[\text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) + 1 \right].$$

As both algorithms follow the same path in the state space, this expectation can be computed with respect to the code of Algorithm 4. Fix some state (v, p) and observe that if there exists some $j \in D^{v, p}$ such that $\rho_j = \star$, then by Lemma 26,

$$\text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(p, \star) + j}^{\prec v}) = \text{rank}(Q_{S(p, \star)}^{\prec v}) + 1 = \text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) + 1.$$

Therefore, for $\mathbf{x} \sim \mathcal{X}^{v, p}$ and $\rho \sim \mathcal{R}^{v, p, \mathbf{x}}$, we have

$$\begin{aligned} \Pr_{\mathbf{x}, \rho} \left[\text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) = \text{rank}(Q_{S(\rho, \neq 0)}^{\prec v}) + 1 \right] &\geq \Pr_{\mathbf{x}, \rho} [\exists j \in D^{v, p} : \rho_j = \star] \\ &= 1 - \Pr_{\mathbf{x}, \rho} [\forall j \in D^{v, p} : \rho_j = \mathbf{x}_j = 0]. \end{aligned}$$

The last equality is due to the fact that for all $j \in D^{v, p}$, if $\rho_j = 0$ then $\mathbf{x}_j = 0$. Let $D := D^{v, p}$. We can now substitute $\hat{\mathcal{X}}^{v, p}$ for $\mathcal{X}^{v, p}$ and $\hat{\mathcal{R}}^{p, \mathbf{x}}$ for $\mathcal{R}^{v, p, \mathbf{x}}$ using Claim 29:

$$\begin{aligned} &\Pr_{\mathbf{x}, \rho} [\forall j \in D : \rho_j = \star \wedge \mathbf{x}_j = 0] \\ &= \Pr_{\mathbf{x}, \rho} [\forall j \in D : \mathbf{x}_j = 0] \cdot \Pr_{\mathbf{x}, \rho} [\forall j \in D : \rho_j = \star \mid \forall j \in D : \mathbf{x}_j = 0] \\ &= \prod_{j \in D} (1 - \delta_j/2) \cdot \prod_{j \in D} \frac{2 - 2\delta_j}{2 - \delta_j} \\ &= \prod_{j \in D} (1 - \delta_j) \\ &\leq (1 - \lambda)^{|D|}. \end{aligned}$$

Thus, if $\mathbf{x} \sim \mu$ and $\rho \sim R_\mu^x$, we have

$$\mathbb{E}_{\mathbf{x}, \rho} [A(\mathbf{x}, \rho)] \geq \sum_{(v, p)} \Pr_{\mathbf{x}, \rho} [\text{state } (v, p) \text{ is reached}] \cdot \left(1 - (1 - \lambda)^{|D^{v, p}|} \right).$$

We now bound the expected number of queries made by \mathcal{T} . When $D^{v, p} = \emptyset$, \mathcal{T} skips making a query at v . On the other hand, when $D^{v, p} \neq \emptyset$, the algorithm goes over $j \in D^{v, p}$ and stops making queries as soon as it hits some $\rho_j = \star$. This probability is independent for each $j \in D^{v, p}$ and can be computed explicitly using Claim 29. For $\mathbf{x} \sim \mathcal{X}^{v, p}$ and $\rho \sim \mathcal{R}^{v, p, \mathbf{x}}$:

$$\Pr_{\mathbf{x}, \rho} [\rho_j = \star] = \Pr_{\mathbf{x}} [\mathbf{x}_j = 0] \cdot \Pr_{\mathbf{x}, \rho} [\rho_j = \star \mid \mathbf{x}_j = 0] + \Pr_{\mathbf{x}} [\mathbf{x}_j = 1] \cdot \Pr_{\mathbf{x}} [\rho_j = \star \mid \mathbf{x}_j = 1] = \delta_j \geq \lambda.$$

Therefore, if $\mathbf{x} \sim \mu$ and $\rho \sim R_\mu^x$,

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}, \rho} [B(\mathbf{x}, \rho)] \\ &\leq \sum_{(v, p)} \Pr_{\mathbf{x}, \rho} [\text{state } (v, p) \text{ is reached}] \cdot \left(\mathbb{1}[D^{v, p} \neq \emptyset] + \sum_{j=0}^{|D^{v, p}|-1} (1 - \lambda)^j \right) \\ &\leq \sum_{(v, p)} \Pr_{\mathbf{x}, \rho} [\text{state } (v, p) \text{ is reached}] \cdot \left(\mathbb{1}[D^{v, p} \neq \emptyset] + \left(1 - (1 - \lambda)^{|D^{v, p}|} \right) / \lambda \right) \\ &\leq \sum_{(v, p)} \Pr_{\mathbf{x}, \rho} [\text{state } (v, p) \text{ is reached}] \cdot (2/\lambda) \cdot \left(1 - (1 - \lambda)^{|D^{v, p}|} \right). \end{aligned} \quad \blacktriangleleft$$

16:22 Direct Sums for Parity Decision Trees

With this in hand, we can now prove Theorem 20, which we restate below for convenience.

► **Theorem 20.** *For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and λ -bounded product distribution μ , we have*

$$\overline{D}_\varepsilon(f, \mu) \leq O(1/\lambda) \cdot S_\varepsilon(f, \mu) \quad \forall \varepsilon \geq 0.$$

Proof. Let \mathcal{T} be a randomised parity tree such that $\overline{sq}(\mathcal{T}, \mu) = S_\varepsilon(f, \mu)$ and $\text{err}_f(\mathcal{T}, \mu) \leq \varepsilon$. Define \mathcal{T}' to be the randomised algorithm obtained by sampling $\mathbf{T} \sim \mathcal{T}$ and returning Algorithm 3 applied to \mathbf{T} . Using Lemma 27, we immediately obtain that $\text{err}(\mathcal{T}', \mu) \leq \varepsilon$. On the other hand:

$$\overline{q}(\mathcal{T}', \mu) = \mathbb{E}_{\mathbf{T}} [\overline{q}(\text{Algorithm 3 on } \mathbf{T}, \mu)] \leq (2/\lambda) \cdot \mathbb{E}_{\mathbf{T}} [\overline{sq}(\mathbf{T}, \mu)] = (2/\lambda) \cdot S_\varepsilon(f, \mu).$$

Thus, $\overline{D}_\varepsilon(f, \mu) \leq O(1/\lambda) \cdot S_\varepsilon(f, \mu)$, as desired. ◀

6.2 Conversion for general product distribution

Algorithm 3 is not efficient for arbitrary product distribution since queries can be very biased so that $\prod_{j \in D^{v,p}} (1 - \delta_j) = 1 - o(1)$. In such cases, we cannot even afford to pay one query as the corresponding expected increment for \overline{sq} is $o(1)$.

To overcome this obstacle, we introduce the following idea. Run the algorithm as if every query x_j returned 0, i.e. assuming $x_j = \rho_j = 0$ for all $j \in S(p, ?)$ (this is likely to happen for very biased distributions). This generates a list of indices for which we assume $x_j = 0$. Upon reaching a leaf, we check efficiently whether one of those x_j is actually 1. If no such j exists, we're done – at the cost of *no real* queries! On the other hand, if a 1 is found, we backtrack to this state and restart the procedure. Since we've found $x_j = 1$, it must be that $\rho_j = \star$ and the S_ε algorithm has to pay one query there.

The process **BuildList** that “runs assuming $x_j = 0$ ” and produces a list of indices to check is described in Section 6.2. Then, the updated algorithm for converting an S_ε algorithm to a \overline{D}_ε one is formulated in Algorithm 5.

How to run Section 6.2?

This problem can be formulated as follows. Let $\text{FFO}_n : \{0, 1\}^n \rightarrow [n] \cup \perp$ be the search problem that asks for the index of the first (running from left to right) ‘1’ in x or \perp if $x = 0^n$. Even though a simple adversary argument shows that one cannot perfectly compute FFO_n by making $< n$ parity queries, a folklore result [24, 44, 30], proves that there is a randomised protocol making $O(\log n)$ queries that computes FFO_n with some small error.

► **Lemma 30.** *For any $\alpha > 0$, $R_\alpha(\text{FFO}_n) \leq O(\log n + \log(1/\alpha))$.*

Proof. This folklore fact is discussed for the parity context in Section A.4. ◀

We let \mathcal{T}'_γ be the parity tree obtained by running Algorithm 5 with error parameter $\alpha := \gamma/n$ on Section 6.2. Given two indices $i, j \in J$, we say $i \prec_J j$ if i appears strictly earlier than j in J , and $i \preceq_J j$ if $i \prec_J j$ or $i = j$. Fix any $x \in \text{supp}(\mathcal{X}^{v,p})$. Let i^* denote the first index i in J such that $x_i = 1$ and suppose that i^* is added to J when $u = u^*$. Observe that if such i^* exists, $x_j = 0$ for all $j \prec_J i^*$. As a consequence, we know that u^* must be reached. Moreover, we can immediately get the values of ρ_j by flipping biased coins for all $j \preceq_J i^*$. Therefore, given i^* , one can perfectly simulate Algorithm 3 by going over J and updating p , until finding the first index $j^* \preceq_J i^*$ such that $\rho_{j^*} = \star$. We are now ready to prove the correctness and efficiency of \mathcal{T}'_γ .

■ **Algorithm 5** converts an algorithm for S_ε to \overline{D}_ε for general product distributions.

Input: $x \in \{0, 1\}^n$
Output: $a \in \{0, 1\}$

- 1: Initialize $v \leftarrow \text{root of } T$, $p \leftarrow ?^n$
- 2: **while** v is not a leaf **do**
- 3: $(J, \ell) \leftarrow \text{BuildList}(v, p)$
- 4: Find the first element $i^* \in J$ with $x_{i^*} = 1$ or set $i^* = \perp$ if none exists
- 5: $\text{FOUND} \leftarrow 0$
- 6: **for** $j \in J$ **do**
- 7: Sample $\eta \sim \text{Ber}(\delta_j / (2 - \delta_j))$
- 8: **if** $j = i^*$ or $\eta = 1$ **then**
- 9: $p_j \leftarrow \star$
- 10: $u \leftarrow w_j$
- 11: $\text{FOUND} \leftarrow 1$
- 12: **break**
- 13: **end if**
- 14: $p_j \leftarrow 0$
- 15: **end for**
- 16: **if** $\text{FOUND} = 1$ **then**
- 17: Query $\langle x, Q^u \rangle$ and set b^u as the outcome
- 18: Move $v \leftarrow \text{child}(u, b^u)$
- 19: **else**
- 20: Update $v \leftarrow \ell$
- 21: **end if**
- 22: **end while**
- 23: **return** $L(v)$

■ **Algorithm 6** the subroutine BuildList.

Input: $v \in \mathcal{N}(T)$, $p \in \{0, \star, ?\}^n$
Output: a list of indices J and a leaf ℓ

- 1: Initialize $J \leftarrow []$, $u \leftarrow v$, $p' \leftarrow p$
- 2: **while** u is not a leaf **do**
- 3: $D^{v, p'} \leftarrow \{j \in [n] : p'_j = ? \wedge \text{rank}(Q_{S(p', \star)+j}^{\prec u}) = \text{rank}(Q_{S(p', \star)}^{\prec u}) + 1\}$
- 4: **for** $j \in D^{v, p'}$ **do** ▷ in arbitrary order
- 5: $p'_j \leftarrow 0$
- 6: $w_j \leftarrow u$
- 7: $J \leftarrow [J, j]$
- 8: **end for**
- 9: Infer $b^u \leftarrow \langle x, Q^u \rangle$ assuming $(Q^{\prec v})^T x = b^{\prec v}$ and $x_j = 0$ for all $j \in S(p', 0)$
- 10: Move $u \leftarrow \text{child}(u, b^u)$
- 11: **end while**
- 12: **return** (J, u)

► **Lemma 31.** For any fixed $x \in \{0, 1\}^n$, $\Pr[\mathcal{T}'_\gamma(x) = 1] \in \mathbb{1}[T(x) = 1] \pm \gamma$.

Proof. The randomness of \mathcal{T}'_γ stems from η and the randomness involved in running the FFO algorithm at Section 6.2. To analyse the latter, observe that Section 6.2 is called at most n times and each call fails with probability at most $\alpha = \gamma/n$, hence:

$$d_{\text{TV}}(\mathcal{T}'_0(x), \mathcal{T}'_\gamma(x)) \leq \Pr[\text{at least one oracle call at line 4 gives a wrong index}] \leq n \cdot \frac{\gamma}{n} = \gamma.$$

If no call fails the discussion above implies that Algorithm 5 behaves identically to the earlier Algorithm 3. Hence, correctness of the former (Lemma 27) implies $\Pr[\mathcal{T}'_0(x) = 1] = \mathbb{1}[T(x) = 1]$. ◀

► **Lemma 32.** We have $\bar{q}(\mathcal{T}'_\gamma, \mu) \leq O(\log(n/\gamma)) \cdot (\bar{sq}(T, \mu) + 1) + \gamma \cdot n$.

Proof. We first prove that the expected number of iterations of the outer while-loop is low assuming that the algorithm always gets the correct index i^* at line 4. Similar to what we did in Section 6.1, we view the randomness used in the for-loop (line 6 to 15) in Algorithm 5 as a pre-generated partial assignment $\rho \sim \mathcal{R}_\mu^x$. Note that the bits of ρ are independent. If i^* is the first index in J with $x_{i^*} = 1$, we know that $x_{i^*} = 1$ and $x_j = 0$ for all $j \prec_J i^*$. At the same time, ρ_j for all $j \preceq_J i^*$ are revealed to the algorithm one by one. As soon as some $\rho_j = \star$ is found, the algorithm quits the loop.

For each $x \in \{0, 1\}^n$ and $\rho \in \text{supp}(\mathcal{R}_\mu^x)$, consider running \mathcal{T}'_γ on input x using randomness ρ . Define $K(x, \rho)$ as the number of iterations of the outer while loop when \mathcal{T}'_γ always gets the correct i^* on line 4. Let p^* denote the final state of p . Since in each iteration except for the last one, we update some p_j as \star , we have $K(x, \rho) \leq |S(p^*, \star)| + 1$. By Lemma 26, we further have

$$K(x, \rho) \leq \text{rank}\left(Q_{S(p^*, \star)}^{\prec \ell(x)}\right) + 1 = \text{rank}\left(Q_{S(p^*, \neq 0)}^{\prec \ell(x)}\right) + 1,$$

where $\ell(x) \in \mathcal{L}(T)$ is the unique leaf at which T terminates given x . Since for all $p_j \neq ?$, $p_j = \rho_j$, we have $S_\star^{p^*} \subseteq S_\star^\rho \subseteq S_{\neq 0}^{p^*}$, hence $K(x, \rho) \leq \text{rank}(Q_{S(\rho, \star)}^{\prec \ell(x)}) + 1$. On the other hand, by definition we have

$$\bar{sq}(T, \mu) = \mathbb{E}_{\substack{x \sim \mu \\ \rho \sim \mathcal{R}_\mu^x}} \left[\text{rank}\left(Q_{S(\rho, \star)}^{\prec \ell(x)}\right) \right] \implies \mathbb{E}_{\substack{x \sim \mu \\ \rho \sim \mathcal{R}_\mu^x}} [K(x, \rho)] \leq \bar{sq}(T, \mu) + 1.$$

Lemma 30 asserts that line Section 6.2 can be implemented to error γ/n using $O(\log(n/\gamma))$ parity queries. Since all those calls are completed successfully with probability $\geq \gamma$, we finally have:

$$\bar{q}(\mathcal{T}'_\gamma, \mu) \leq (1-\gamma) \cdot \mathbb{E}_{\substack{x \sim \mu \\ \rho \sim \mathcal{R}_\mu^x}} [K(x, \rho)] \cdot O(\log(n/\gamma)) + \gamma \cdot n \leq O(\log(n/\gamma)) \cdot (\bar{sq}(T, \mu) + 1) + \gamma \cdot n. \blacktriangleleft$$

► **Theorem 19.** For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$, product distribution μ , $\gamma \in (0, 1/n)$, we have

$$\bar{D}_{\varepsilon+\gamma}(f, \mu) \leq O(\log(n/\gamma)) \cdot (S_\varepsilon(f, \mu) + 1) \quad \forall \varepsilon \geq 0.$$

Proof. Let \mathcal{T} be a randomised parity decision tree such that $\bar{sq}(\mathcal{T}, \mu) = S_\varepsilon(f, \mu)$ and $\text{err}_f(\mathcal{T}, \mu) \leq \varepsilon$. Define \mathcal{T}^* to be the randomised algorithm obtained by sampling $\mathbf{T} \sim \mathcal{T}$ and returning the corresponding \mathcal{T}'_γ . Using Lemma 31, we immediately obtain that $\text{err}(\mathcal{T}^*, \mu) \leq \varepsilon + \gamma$. By Lemma 32 and the range of parameters allowed for γ , we get

$$\bar{q}(\mathcal{T}^*, \mu) = \mathbb{E}_{\mathbf{T}} [\bar{q}(\mathcal{T}'_\gamma, \mu)] \leq O(\log(n/\gamma)) \cdot \mathbb{E}_{\mathbf{T}} [\bar{sq}(T, \mu) + 1] = O(\log(n/\gamma)) \cdot (\bar{sq}(\mathcal{T}, \mu) + 1). \blacktriangleleft$$

7 Separations I: disc vs. D^\times

In this section we prove Lemma 3, restated here for convenience.

► **Lemma 3.** *The complexity measures disc and D^\times are incomparable:*

1. *There is an n -bit function f such that $\text{disc}(f) = O(\log n)$ while $D^\times(f) = \Theta(n)$.*
2. *There is an n -bit function f such that $\text{disc}(f) = \Theta(n)$ while $D^\times(f) = O(1)$.*

Proof. For the first item, we can consider the n -bit majority function $f := \text{MAJ}_n$. It follows from [14, Theorem 1.2] that $D^\times(\text{MAJ}_n) \geq \Omega(n)$ where the hard distribution is uniform. By contrast, it is not hard to see that $\text{disc}(\text{MAJ}_n) \leq O(\log n)$ (if we query x_i for a random $i \in [n]$, it will have bias $\geq \Omega(1/\sqrt{n})$ toward predicting $\text{MAJ}_n(x)$). We prove the second item by a probabilistic argument. Consider a random function \mathbf{f} , which is set with $\mathbf{f}(x) \sim \text{Ber}(2^{-0.9n})$ independently for each $x \in \{0, 1\}^n$. In Claim 33, we show that $\text{disc}(\mathbf{f}) = \Theta(n)$ and in Claim 34 that $D^\times(\mathbf{f}) = O(1)$ with high probability. ◀

► **Claim 33.** With probability $1 - 2^{-2^{\Omega(n)}}$, $\text{disc}(\mathbf{f}) \geq 0.01n$.

Proof. For each non-constant function $f : \{0, 1\}^n \rightarrow \{0, 1\}$, we define the “hard” distribution μ_f as

$$\mu_f(x) := \begin{cases} 1/(2|f^{-1}(0)|) & \text{if } f(x) = 0 \\ 1/(2|f^{-1}(1)|) & \text{if } f(x) = 1 \end{cases}.$$

To prove the claim, it suffices to show $\Pr_{\mathbf{f}}[\text{disc}(\mathbf{f}, \mu_{\mathbf{f}}) \geq 0.01n] \geq 1 - 2^{-2^{\Omega(n)}}$. Using Lemma 9, this can be further simplified to prove:

$$\Pr_{\mathbf{f}}[\max_{S \in \mathcal{O}_n} \text{bias}(\mathbf{f}, \mu_{\mathbf{f}}, S) \leq 2^{-0.01n-1}] \geq 1 - 2^{-2^{\Omega(n)}}.$$

To that end, fix any $S \in \mathcal{O}_n$, note that $|S| = |\{0, 1\} \setminus S| = 2^{n-1}$ and observe that by a Chernoff bound,

$$\begin{aligned} \Pr_{\mathbf{f}}[|\mu(\mathbf{f}^{-1}(1) \cap S) - 1/4| \geq 2^{-0.02n}] &\leq \Pr_{\mathbf{f}}[|\mathbf{f}^{-1}(1)| < 2^{0.1n-1}] \\ &\quad + \Pr_{\mathbf{f}}[||\mathbf{f}^{-1}(1) \cap S| - 2^{0.1n-1}| > 2^{0.07n}] \\ &\quad + \Pr_{\mathbf{f}}[||\mathbf{f}^{-1}(1) \setminus S| - 2^{0.1n-1}| > 2^{0.07n}] \\ &\leq 3e^{-2^{0.03n}}. \end{aligned}$$

Using a similar argument, we can also show $\Pr_{\mathbf{f}}[|\mu(\mathbf{f}^{-1}(0) \cap S) - 1/4| \geq 2^{-0.02n}] \leq 3e^{-2^{0.03n}}$. By definition, $\text{bias}(\mathbf{f}, \mu_{\mathbf{f}}, S) = |\mu(\mathbf{f}^{-1}(0) \cap S) - \mu(\mathbf{f}^{-1}(1) \cap S)|$, we thus have $\Pr[\text{bias}(\mathbf{f}, \mu_{\mathbf{f}}, S) \geq 2^{-0.01n-1}] \leq 6e^{-2^{0.03n}}$. Finally, observe that $|\mathcal{O}_n| \leq 2^n$ and so using a union bound,

$$\begin{aligned} \Pr_{\mathbf{f}}[\text{disc}(\mathbf{f}) \geq 0.01n] &\geq \Pr_{\mathbf{f}}[\max_{S \in \mathcal{O}_n} \text{bias}(\mathbf{f}, \mu_{\mathbf{f}}, S) \leq 2^{-0.01n-1}] \\ &\geq 1 - 2^n \Pr[\text{bias}(\mathbf{f}, \mu_{\mathbf{f}}, S) \geq 2^{-0.01n-1}] \\ &\geq 1 - 2^{-2^{\Omega(n)}}. \end{aligned}$$

◀

► **Claim 34.** With probability $1 - 2^{-\Omega(n)}$, $D^\times(f) \leq 20000$.

16:26 Direct Sums for Parity Decision Trees

Proof. Let $\mathcal{D}^\times := \{\text{Ber}(p_1, \dots, p_n) \mid p_1, \dots, p_n \in [0, 1/2]\}$ denote the set of 0-biased product distributions, where $\text{Ber}(p_1, \dots, p_n) := \text{Ber}(p_1) \times \dots \times \text{Ber}(p_n)$. As observed in Section 4, it suffices to show $\Pr_{\mathbf{f}}[\max_{\mu \in \mathcal{D}^\times} D_{1/3}(\mathbf{f}, \mu) \leq 20000] \geq 1 - 2^{-\Omega(n)}$.

As a first attempt, one might want to prove that $D_{1/3}(\mathbf{f}, \mu) = O(1)$ with sufficiently high probability for any fixed μ and then apply union bound over all $\mu \in \mathcal{D}^\times$. However, this cannot be done directly since \mathcal{D}^\times is infinite. Luckily, we can circumvent this barrier by discretizing \mathcal{D}^\times . Let us define $\mathcal{D}_{\mathbb{Z}}^\times := \{\text{Ber}(a_1/10n, \dots, a_n/10n) \mid a_1, \dots, a_n \in \{0, \dots, 5n\}\}$. For every $\mu = \text{Ber}(p_1, \dots, p_n) \in \mathcal{D}_{\mathbb{Z}}^\times$ and $f : \{0, 1\}^n \rightarrow \{0, 1\}$, consider the following two cases:

- If $\sum_i p_i \geq 10$, then $M := \max_{x \in \{0, 1\}^n} \mu(x) \leq e^{-\sum_i p_i} \leq 1000 \sum_i p_i$. Observe that

$$\Pr_{\mathbf{f}} \left[\sum_{x \in \{0, 1\}^n} f(x) \mu(x) \geq 1/5 \right] \leq 2^M \cdot (2^{-0.9n})^{M/5} \leq 2^{-150 \sum_i p_i n},$$

thus $\Pr_{\mathbf{f}}[D_{1/4}(\mathbf{f}, \mu) = 0] \geq 1 - 2^{-150 \sum_i p_i n}$.

- Otherwise, we devise the following protocol: Sort $\mu(x_1) \geq \dots \geq \mu(x_{2^n})$. Pick the top 1000 inputs $X = \{x_1, \dots, x_{1000}\}$, then we check if our input x is in X . If yes, we output $f(x)$, otherwise we output 0. Formally, we define the function $g : \{0, 1\}^n \rightarrow \{0, 1\}$ where

$$g(x) := \begin{cases} f(x) & \text{if } x \in X \\ 0 & \text{if } x \notin X \end{cases}.$$

Since testing whether $x = x_i$ can be done with m queries with success probability $1 - 2^{-m}$, by choosing $m = 20$ and running the testing for every $i \in [1000]$, one can show $\mathcal{R}_{0.01}(g) \leq 20000$. It remains to prove that $\Pr_{\mathbf{f}}[\mathbf{f}(x) = g(x)] \geq 4/5$ with high probability. Observe that for each $x \notin X$, $\mu(x) \leq 1/1000$. Therefore:

$$\Pr_{\mathbf{f}} \left[\sum_{x \notin X} [\mu(x) \mathbf{f}(x)] \leq 1/5 \right] \geq 1 - 2^{1000} \cdot (2^{-0.9n})^{200} \geq 1 - 2^{-150n}.$$

For those \mathbf{f} , we have $\Pr_{\mathbf{f}}[\mathbf{f}(x) = g(x)] \geq 4/5$, which implies that $D_{0.22}(\mathbf{f}, \mu) \leq 20000$. By union bound over $\mu \in \mathcal{D}_{\mathbb{Z}}^\times$, we can deduce that

$$\begin{aligned} & \Pr_{\mathbf{f}} \left[\max_{\mu \in \mathcal{D}_{\mathbb{Z}}^\times} D_{0.22}(\mathbf{f}, \mu) > 20000 \right] \\ & \leq \sum_{\mu \in \mathcal{D}_{\mathbb{Z}}^\times} \Pr_{\mathbf{f}}[D_{0.22}(\mathbf{f}, \mu) > 20000] \\ & \leq \sum_{a_1=0}^{5n} \dots \sum_{a_n=0}^{5n} \mathbb{1} \left[\sum_i a_i \geq 100n \right] \cdot e^{-150 \sum_i a_i} \\ & \quad + \sum_{a_1=0}^{5n} \dots \sum_{a_n=0}^{5n} \mathbb{1} \left[\sum_i a_i < 100n \right] \cdot 2^{-150n} \\ & \leq \sum_{a_1=0}^{5n} \dots \sum_{a_n=0}^{5n} e^{-100(a_1+5)} \dots e^{-100(a_n+5)} + 2^{101n} \cdot 2^{-150n} \\ & \leq \left(\sum_{a_1=0}^{5n} e^{-100(a_1+5)} \right)^n + 2^{-49n} \\ & \leq 2^{-\Omega(n)}. \end{aligned}$$

Consider now any product distribution $\mu = \text{Ber}(p_1, \dots, p_n) \in \mathcal{D}^\times$, define its rounded version $[\mu]$:

$$[\mu] := \left(\text{Ber} \left(\frac{\lceil 10n \cdot p_1 \rceil}{10n} \right), \dots, \text{Ber} \left(\frac{\lceil 10n \cdot p_n \rceil}{10n} \right) \right) \in \mathcal{D}_{\mathbb{Z}}^\times.$$

Observe that $d_{TV}(\mu, \lceil \mu \rceil) \leq 1 - (1 - 1/10n)^n \leq 1 - 1/e^{-1/10} < 0.1$, thus we have $\text{err}_f(T, \lceil \mu \rceil) \leq \text{err}_f(T, \mu) + 0.1$ for any parity tree T and $f: \{0, 1\}^n \rightarrow \{0, 1\}$. Together with the string of inequalities developed above, we conclude that with probability at least $1 - 2^{-\Omega(n)}$,

$$\max_{\mu \in \mathcal{D}^\times} D_{1/3}(\mathbf{f}, \mu) \leq \max_{\mu \in \mathcal{D}_{\mathbb{Z}}^\times} D_{1/3-0.1}(\mathbf{f}, \mu) \leq \max_{\mu \in \mathcal{D}_{\mathbb{Z}}^\times} D_{0.22}(\mathbf{f}, \mu) \leq 20000. \quad \triangleleft$$

8 Separations II: \mathbf{S} vs. \mathbf{D}^\times

The goal of this section is to provide the following example of a function.

► **Theorem 35.** *There exists a function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and a product distribution μ such that $\mathbf{D}^\times(f) = \Theta(\text{disc}(f, \mu)) = \Theta(\log n)$ and $\mathbf{S}_0(f, \mu) = \Theta(1)$.*

Recall that by Theorem 19, this is the largest possible gap between \mathbf{S} and \mathbf{D}^\times . To prove the separation, we use the function $\text{FPE}: \{0, 1\}^{2n} \rightarrow \{0, 1\}$ which takes two inputs $x, y \in \{0, 1\}^n$ and returns the value y_i associated with the location i of the first “1” in x . More precisely, we let $\text{FO}(x) \in [n]$ be the location (from left to right) of the first “1” in x and $\text{FO}(x) = 1$ if $x = 0^n$ and let $\text{FPE}(x, y) = y_{\text{FO}(x)}$. We choose as hard distribution the product distribution $\mu := \mathcal{X} \times \mathcal{Y}$ where for each $i \in [n]$:

$$\mathcal{X}_i \sim \text{Ber}(1/\sqrt{n}) \quad \text{and} \quad \mathcal{Y}_i \sim \text{Ber}(1/2).$$

Let us note that the choice of $1/\sqrt{n}$ in the distribution \mathcal{X} is arbitrary: any $p = n^a$ for constant $a \in (-1, 0)$ is enough to guarantee that $x \neq 0^n$ with high probability and get the $\Omega(\log n)$ lower bound.

Proof of Theorem 35. We first prove that $\mathbf{S}_0(\text{FPE}, \mu) = \Theta(1)$. Consider the following simple brute-force query algorithm T that computes f : Query the bits of x one by one from left to right, until finding the first index i such that $x_i = 1$. Then query y_i and return y_i if such i exists. Otherwise ($x = 0^n$), simply return 1.

Observe that $\text{err}_{\text{FPE}}(T, \mu) = 0$. Thus we only need to show $\overline{\text{sq}}(f, \mu) = \Theta(1)$. Let $X_i := \{x \mid x_i = 1, x_j = 0, \forall j < i\}$ denote the set of $x \in \{0, 1\}^n$ for which $\text{FO}(x) = i$. Note that $\{0, 1\}^n = X_1 \sqcup \dots \sqcup X_n \sqcup \{0^n\}$ forms a partition of $\{0, 1\}^n$. By the definition of μ , we have $\mu(X_i) = (1 - 1/\sqrt{n})^{i-1}/\sqrt{n}$. For all $x \in X_i$, T queries the same set of variables $\{x_1, \dots, x_{i-1}, x_i, y_i\}$ on x . Moreover, sample $\rho \sim \mathcal{R}_\mu^x$ and for each $1 \leq j < i$, since $x_j = 0$, we have that $\Pr[\rho_j = \star] = 1/(\sqrt{n} - 1)$. Therefore,

$$h(x) := \mathbb{E}_{\rho \sim \mathcal{R}_\mu^x}[q(T_\rho, x)] \leq \frac{i-1}{\sqrt{n}-1} + 2.$$

We conclude that

$$\begin{aligned} \overline{\text{sq}}(T, \mu) &= \mathbb{E}_{\mathbf{x} \sim \mu}[h(\mathbf{x})] \\ &\leq \sum_{i=1}^n \mu(X_i) \cdot \mathbb{E}_{\mathbf{x} \sim \mu_{X_i}}[h(\mathbf{x})] + (n+1) \cdot \mu(0^n) \\ &\leq \frac{1}{n-\sqrt{n}} \cdot \sum_{i=1}^n (i-1)(1-1/\sqrt{n})^{i-1} + n \cdot (1-1/\sqrt{n})^n + 2 \\ &< \frac{2}{n} \cdot \sum_{i=0}^{\infty} i(1-1/\sqrt{n})^i + 3 \\ &= \Theta(1). \end{aligned}$$

16:28 Direct Sums for Parity Decision Trees

Let us now turn our attention to $\text{disc}(\text{FPE}, \mu)$. The lower-bound $\text{disc}(\text{FPE}, \mu) \geq \Omega(\log n)$ is covered in Claim 36. The upper bound $\text{disc}(\text{FPE}, \mu) \leq O(\log n)$ is a direct consequence of $\text{bias}(\text{FPE}, \mu, S) \geq n^{-1/2}$ for $S = \{(x, y) \in \{0, 1\}^n : y_1 = 1\}$. More interestingly, one can actually show the stronger statement $D_{1/3}(f, \mu) \leq O(\log n)$. Indeed, $\mathbf{x} \sim \mathcal{X}$ has exactly one “1” in the first $\lceil \sqrt{n} \rceil$ bits with probability $\geq e^{-1.01} \geq 1/3$ for n large enough. In that case, a simple binary search amongst the first $\lceil \sqrt{n} \rceil$ bits of x using parity queries is enough to find that location and return the corresponding bit of y . \blacktriangleleft

\triangleright **Claim 36.** $\text{disc}(\text{FPE}, \mu) \geq \Omega(\log n)$

Proof. Using the characterisation of the bias with codimension-1 subspace Lemma 9, it is enough to show:

$$\max_{S \in \mathcal{O}^n} \text{bias}(\text{FPE}, \mu, S) \leq n^{-1/3}.$$

Fix an affine space S^* of codimension 1 that maximize the above expression, i.e. some $\alpha, \beta \in \{0, 1\}^n$ and $\gamma \in \{0, 1\}$ such that $S^* = \{(x, y) \in \{0, 1\}^{2n} : \alpha \cdot x + \beta \cdot y = \gamma\}$. To simplify notation, we assume in what follows that $\gamma = 0$ but the proof is similar for the case $\gamma = 1$. Let us partition S in two sets:

$$\begin{aligned} S^0 &:= \{(x, y) \in \{0, 1\}^{2n} : \alpha \cdot x = 0 \text{ and } \beta \cdot y = 0\}; \\ S^1 &:= \{(x, y) \in \{0, 1\}^{2n} : \alpha \cdot x = 1 \text{ and } \beta \cdot y = 1\}. \end{aligned}$$

We have:

$$\max_{S \in \mathcal{O}^n} \text{bias}(\text{FPE}, \mu, S) = \text{bias}(\text{FPE}, \mu, S^*) \leq \text{bias}(\text{FPE}, \mu, S^0) + \text{bias}(\text{FPE}, \mu, S^1).$$

Let us suppose without loss of generality that $\text{bias}(\text{FPE}, \mu, S^0) \geq \text{bias}(\text{FPE}, \mu, S^1)$ so that it is enough to show $\text{bias}(\text{FPE}, \mu, S^0) \leq 2n^{-1/2}$. Note that if $\Pr_{\mathbf{x}, \mathbf{y} \sim \mu}[(\mathbf{x}, \mathbf{y}) \in S^0] = 0$, we’re done. If not, we can re-express the bias in the language of probability:

$$\begin{aligned} \text{bias}(\text{FPE}, \mu, S^0) &= \left| \sum_{(x, y) \in S^0} (-1)^{\text{FPE}(x)} \mu(x) \right| \\ &= \left| \sum_{b \in \{0, 1\}} (-1)^b \cdot \Pr_{\mathbf{x}, \mathbf{y}} [\text{FPE}(x) = b \wedge (\mathbf{x}, \mathbf{y}) \in S^0] \right| \\ &= \Pr_{\mathbf{x}, \mathbf{y}} [(\mathbf{x}, \mathbf{y}) \in S^0] \cdot \left| \sum_{b \in \{0, 1\}} (-1)^b \cdot \Pr_{\mathbf{x}, \mathbf{y}} [\text{FPE}(x) = b \mid (\mathbf{x}, \mathbf{y}) \in S^0] \right|. \end{aligned}$$

Let us denote the quantity within the absolute value by Φ and the event $(\mathbf{x}, \mathbf{y}) \in S^0$ by E . Observe that S^0 can be conveniently decomposed as $S^0 = S^X \times S^Y$ where $S^X := \{x \in \{0, 1\}^n : \alpha \cdot x = 0\}$ and $S^Y := \{y \in \{0, 1\}^n : \beta \cdot y = 0\}$. With this, we have:

$$\begin{aligned} \Phi &= \sum_{b \in \{0, 1\}} (-1)^b \cdot \Pr_{\mathbf{x}, \mathbf{y}} [\text{FPE}(\mathbf{x}, \mathbf{y}) = b \mid E] \\ &= \sum_{i \in [n]} \sum_{b \in \{0, 1\}} (-1)^b \cdot \Pr_{\mathbf{x}, \mathbf{y}} [\text{FO}(\mathbf{x}) = i \mid E] \Pr_{\mathbf{x}, \mathbf{y}} [\text{FPE}(\mathbf{x}, \mathbf{y}) = b \mid E \wedge \text{FO}(\mathbf{x}) = i] \\ &= \sum_{i \in [n]} \Pr_{\mathbf{x} \sim \mathcal{X}} [\text{FO}(\mathbf{x}) = i \mid \mathbf{x} \in S^X] \cdot \sum_{b \in \{0, 1\}} (-1)^b \cdot \underbrace{\Pr_{\mathbf{y} \sim \mathcal{Y}} [\mathbf{y}_i = b \mid \mathbf{y} \in S^Y]}_{:= p_i^b}. \end{aligned}$$

Recall that S^Y is a codimension-1 space and \mathcal{Y} is the uniform distribution over $\{0, 1\}^n$. Thus, if $|\alpha|$ (the number of non-zero entries in α) is zero or ≥ 2 , it must be that $p_i^b = 1/2$ for all $i \in [n]$ and $b \in \{0, 1\}$. In that case, the claim is proven because $\Phi = 0$ and so

$\text{bias}(\text{FPE}, \mu, S^0) = 0$. We can thus assume that $|\alpha| = 1$ and fix $i^* \in [n]$ to be the unique coordinate such that $\alpha_{i^*} = 1$. Now, observe that $p_i^b = 1/2$ for all $i \neq i^*$ and $b \in \{0, 1\}^n$, $p_{i^*}^0 = 1$ and $p_{i^*}^1 = 0$ so that:

$$\Phi = \sum_{i \in [n]} \Pr_{\mathbf{x} \sim \mathcal{X}} [\text{FO}(\mathbf{x}) = i \mid \mathbf{x} \in S^X] \cdot (p_i^0 - p_i^1) = \Pr_{\mathbf{x} \sim \mathcal{X}} [\text{FO}(\mathbf{x}) = i^* \mid \mathbf{x} \in S^X].$$

Finally, we use the fact that the event $\text{FO}(\mathbf{x}) = i^*$ with $\mathbf{x} \sim \mathcal{X}$ is unlikely to happen if S^X has large mass under \mathcal{X} . In any case, the probability is maximized for $i^* = 1$ and hence:

$$\begin{aligned} \text{bias}(\text{FPE}, \mu, S^0) &= \Pr_{\mathbf{x} \sim \mathcal{X}} [\mathbf{x} \in S^X] \cdot \Pr_{\mathbf{y} \sim \mathcal{Y}} [\mathbf{y} \in S^Y] \cdot |\Phi| \\ &\leq \Pr_{\mathbf{x}} [\text{FO}(\mathbf{x}) = i^* \wedge \mathbf{x} \in S^X] \\ &\leq \Pr_{\mathbf{x}} [\text{FO}(\mathbf{x}) = 1]. \end{aligned}$$

The event $\text{FO}(\mathbf{x}) = 1$ can happen because $\mathbf{x}_1 = 1$ or $\mathbf{x} = 0^n$, thus we bound the bias with

$$\Pr_{\mathbf{x} \sim \mathcal{X}} [\text{FO}(\mathbf{x}) = 1] \leq \Pr_{\mathbf{x}} [\mathbf{x}_1 = 1] + \Pr_{\mathbf{x}} [\mathbf{x} = 0^n] \leq n^{-1/2} + e^{-\sqrt{n}} \leq 2n^{-1/2}. \quad \triangleleft$$

References

- 1 Yaroslav Alekseev, Yuval Filmus, and Alexander Smal. Lifting Dichotomies. In *39th Computational Complexity Conference (CCC 2024)*, volume 300 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 9:1–9:18. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.CCC.2024.9.
- 2 Yaroslav Alekseev and Dmitry Itsykson. Lifting to bounded-depth and regular resolutions over parities via games. Technical Report TR24-128, ECCC, 2024. URL: <https://eccc.weizmann.ac.il/report/2024/128/>.
- 3 Laszlo Babai, Peter Frankl, and Janos Simon. Complexity classes in communication complexity theory. In *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*, pages 337–347, 1986. doi:10.1109/SFCS.1986.15.
- 4 Boaz Barak, Mark Braverman, Xi Chen, and Anup Rao. How to compress interactive communication. *SIAM Journal on Computing*, 42(3):1327–1363, 2013. doi:10.1137/100811969.
- 5 Paul Beame and Sajin Koroth. On Disperser/Lifting Properties of the Index and Inner-Product Functions. In *14th Innovations in Theoretical Computer Science Conference (ITCS 2023)*, volume 251 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 14:1–14:17. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPIcs.ITCS.2023.14.
- 6 Shalev Ben-David and Eric Blais. A tight composition theorem for the randomized query complexity of partial functions: Extended abstract. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 240–246, 2020. doi:10.1109/FOCS46700.2020.00031.
- 7 Shalev Ben-David and Eric Blais. A new minimax theorem for randomized algorithms. *J. ACM*, 70(6), 2023. doi:10.1145/3626514.
- 8 Shalev Ben-David, Eric Blais, Mika Göös, and Gilbert Maystre. Randomised Composition and Small-Bias Minimax. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 624–635. IEEE Computer Society, 2022. doi:10.1109/FOCS54457.2022.00065.
- 9 Shalev Ben-David, Mika Göös, Robin Kothari, and Thomas Watson. When Is Amplification Necessary for Composition in Randomized Query Complexity? In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2020)*, volume 176 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 28:1–28:16. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2020. doi:10.4230/LIPIcs.APPROX/RANDOM.2020.28.

- 10 Shalev Ben-David and Robin Kothari. Randomized query complexity of sabotaged and composed functions. *Theory of Computing*, 14(5):1–27, 2018. doi:10.4086/toc.2018.v014a005.
- 11 Sreejata Kishor Bhattacharya, Arkadev Chattopadhyay, and Pavel Dvořák. Exponential Separation Between Powers of Regular and General Resolution over Parities. In *39th Computational Complexity Conference (CCC 2024)*, volume 300 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 23:1–23:32. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.CCC.2024.23.
- 12 Eric Blais and Joshua Brody. Optimal separation and strong direct sum for randomized query complexity. In *Proceedings of the 34th Computational Complexity Conference, CCC '19*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2019. doi:10.4230/LIPIcs.CCC.2019.29.
- 13 Guy Blanc, Caleb Koch, Carmen Strassle, and Li-Yang Tan. A Strong Direct Sum Theorem for Distributional Query Complexity. In *39th Computational Complexity Conference (CCC 2024)*, volume 300 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 16:1–16:30. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.CCC.2024.16.
- 14 Mark Braverman, Ankit Garg, Denis Pankratov, and Omri Weinstein. Information lower bounds via self-reducibility. *Theory of Computing Systems*, 59(2):377–396, 2015. doi:10.1007/s00224-015-9655-z.
- 15 Mark Braverman and Anup Rao. Information equals amortized communication. *IEEE Transactions on Information Theory*, 60(10):6058–6069, 2014. doi:10.1109/TIT.2014.2347282.
- 16 Joshua Brody, Jae Tak Kim, Peem Lerdputtipongporn, and Hariharan Srinivasulu. A strong XOR lemma for randomized query complexity. *Theory of Computing*, 19(11):1–14, 2023. doi:10.4086/toc.2023.v019a011.
- 17 Farzan Byramji and Russell Impagliazzo. Lifting to randomized parity decision trees. Technical Report TR24-202, ECCC, 2024. URL: <https://eccc.weizmann.ac.il/report/2024/202/>.
- 18 Arkadev Chattopadhyay and Pavel Dvorak. Super-critical trade-offs in resolution over parities via lifting. Technical Report TR24-132, ECCC, 2024. URL: <https://eccc.weizmann.ac.il/report/2024/132/>.
- 19 Arkadev Chattopadhyay, Nikhil Mande, Swagato Sanyal, and Suhail Sherif. Lifting to Parity Decision Trees via Stifling. In *14th Innovations in Theoretical Computer Science Conference (ITCS 2023)*, volume 251 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 33:1–33:20. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPIcs.ITCS.2023.33.
- 20 Tsun Ming Cheung, Hamed Hatami, Rosie Zhao, and Itai Zilberstein. Boolean functions with small approximate spectral norm. *Discrete Analysis*, 2024. doi:10.19086/da.122971.
- 21 Andrew Drucker. Improved direct product theorems for randomized query complexity. *Comput. Complex.*, 21(2):197–244, 2012. doi:10.1007/s00037-012-0043-7.
- 22 Klim Efremenko, Michal Garlík, and Dmitry Itsykson. Lower bounds for regular resolution over parities. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, volume 41 of *STOC '24*, pages 640–651. ACM, 2024. doi:10.1145/3618260.3649652.
- 23 Tomás Feder, Eyal Kushilevitz, Moni Naor, and Noam Nisan. Amortized communication complexity. *SIAM Journal on Computing*, 24(4):736–750, 1995. doi:10.1137/S0097539792235864.
- 24 U. Feige, D. Peleg, P. Raghavan, and E. Upfal. Computing with unreliable information. In *Proceedings of the Twenty-Second Annual ACM Symposium on Theory of Computing*, STOC '90, pages 128–137. Association for Computing Machinery, 1990. doi:10.1145/100216.100230.
- 25 Yuval Filmus, Edward Hirsch, Artur Riazanov, Alexander Smal, and Marc Vinyals. Proving Unsatisfiability with Hitting Formulas. In *15th Innovations in Theoretical Computer Science Conference (ITCS 2024)*, volume 287 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 48:1–48:20. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.ITCS.2024.48.
- 26 Anat Ganor, Gillat Kol, and Ran Raz. Exponential separation of information and communication for boolean functions. *J. ACM*, 63(5), 2016. doi:10.1145/2907939.

- 27 Uma Girish, Makrand Sinha, Avishay Tal, and Kewen Wu. Fourier growth of communication protocols for XOR functions. In *2023 IEEE 64th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 721–732, 2023. doi:10.1109/FOCS57990.2023.00047.
- 28 Uma Girish, Avishay Tal, and Kewen Wu. Fourier growth of parity decision trees. In *Proceedings of the 36th Computational Complexity Conference, CCC '21*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2021. doi:10.4230/LIPIcs.CCC.2021.39.
- 29 Mika Göös and Gilbert Maystre. A majority lemma for randomised query complexity. In *Proceedings of the 36th Computational Complexity Conference, CCC '21*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2021. doi:10.4230/LIPIcs.CCC.2021.18.
- 30 Nathaniel Harms and Artur Riazanov. Better Boosting of Communication Oracles, or Not. In Siddharth Barman and Sławomir Lasota, editors, *44th IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS 2024)*, volume 323 of *Leibniz International Proceedings in Informatics (LIPIcs)*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.FSTTCS.2024.25.
- 31 Hamed Hatami, Kaave Hosseini, and Shachar Lovett. Structure of protocols for XOR functions. *SIAM Journal on Computing*, 47(1):208–217, 2018. doi:10.1137/17M1136869.
- 32 Hamed Hatami, Kaave Hosseini, Shachar Lovett, and Anthony Ostuni. Refuting Approaches to the Log-Rank Conjecture for XOR Functions. In *51st International Colloquium on Automata, Languages, and Programming (ICALP 2024)*, volume 297 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 82:1–82:11. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPIcs.ICALP.2024.82.
- 33 Dmitry Itsykson and Dmitry Sokolov. Resolution over linear equations modulo two. *Annals of Pure and Applied Logic*, 171(1):102722, 2020. doi:10.1016/j.apal.2019.102722.
- 34 Siddharth Iyer and Anup Rao. An XOR lemma for deterministic communication complexity, 2024. doi:10.48550/arXiv.2407.01802.
- 35 Siddharth Iyer and Anup Rao. XOR lemmas for communication via marginal information. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing, STOC 2024*, pages 652–658. Association for Computing Machinery, 2024. doi:10.1145/3618260.3649726.
- 36 Rahul Jain, Hartmut Klauck, and Miklos Santha. Optimal direct sum results for deterministic and randomized decision tree complexity. *Information Processing Letters*, 110(20):893–897, 2010. doi:10.1016/j.ipl.2010.07.020.
- 37 Rahul Jain, Jaikumar Radhakrishnan, and Pranab Sen. A direct sum theorem in communication complexity via message compression. In Jos C. M. Baeten, Jan Karel Lenstra, Joachim Parrow, and Gerhard J. Woeginger, editors, *Automata, Languages and Programming*, pages 300–315. Springer Berlin Heidelberg, 2003. doi:10.1007/3-540-45061-0_26.
- 38 Hartmut Klauck, Robert Špalek, and Ronald de Wolf. Quantum and classical strong direct product theorems and optimal time-space tradeoffs. *SIAM Journal on Computing*, 36(5):1472–1493, 2007. doi:10.1137/05063235X.
- 39 A. Knop, S. Lovett, S. McGuire, and W. Yuan. Guest column: Models of computation between decision trees and communication. *SIGACT News*, 52(2):46–70, 2021. doi:10.1145/3471469.3471479.
- 40 Eyal Kushilevitz and Yishay Mansour. Learning decision trees using the fourier spectrum. *SIAM Journal on Computing*, 22(6):1331–1348, 1993. doi:10.1137/0222080.
- 41 Troy Lee, Adi Shraibman, and Robert Špalek. A direct product theorem for discrepancy. In *2008 23rd Annual IEEE Conference on Computational Complexity*, pages 71–80, 2008. doi:10.1109/CCC.2008.25.
- 42 Nati Linial and Adi Shraibman. Learning complexity vs. communication complexity. In *2008 23rd Annual IEEE Conference on Computational Complexity*, pages 53–63, 2008. doi:10.1109/CCC.2008.28.
- 43 Nikhil Mande and Swagato Sanyal. On parity decision trees for fourier-sparse boolean functions. *ACM Trans. Comput. Theory*, 16(2), 2024. doi:10.1145/3647629.

- 44 Noam Nisan. The communication complexity of threshold gates. *Proc. of Combinatorics, Paul Erdős is Eighty*, 1993.
- 45 Ryan O'Donnell. *Analysis of Boolean Functions*. Cambridge University Press, 2014. doi:10.1017/CB09781139814782.
- 46 Ryan O'Donnell, John Wright, Yu Zhao, Xiaorui Sun, and Li-Yang Tan. A composition theorem for parity kill number. In *2014 IEEE Conference on Computational Complexity (CCC)*, pages 144–154. IEEE Computer Society, 2014. doi:10.1109/CCC.2014.22.
- 47 Anup Rao and Makrand Sinha. Simplified separation of information and communication. *Theory of Computing*, 14(20):1–29, 2018. doi:10.4086/toc.2018.v014a020.
- 48 Swagato Sanyal. Fourier sparsity and dimension. *Theory of Computing*, 15(11):1–13, 2019. doi:10.4086/toc.2019.v015a011.
- 49 Swagato Sanyal. Randomized query composition and product distributions. In *41st International Symposium on Theoretical Aspects of Computer Science (STACS)*, volume 289 of *LIPIcs*, pages 56:1–56:19. Schloss Dagstuhl, 2024. doi:10.4230/LIPIcs.STACS.2024.56.
- 50 Petr Savický. On determinism versus unambiguous nondeterminism for decision trees. Technical Report TR02-009, Electronic Colloquium on Computational Complexity (ECCC), 2002. URL: <http://eccc.hpi-web.de/report/2002/009/>.
- 51 Ronen Shaltiel. Towards proving strong direct product theorems. *computational complexity*, 12(1):1–22, 2003. doi:10.1007/s00037-003-0175-x.
- 52 Alexander Shekhovtsov and Vladimir Podolskii. Randomized lifting to semi-structured communication complexity via linear diversity. In *16th Innovations in Theoretical Computer Science Conference (ITCS)*, *LIPIcs*, pages 78:1–78:21. Schloss Dagstuhl, 2025. doi:10.4230/LIPIcs.ITCS.2025.78.
- 53 Amir Shpilka, Avishay Tal, and Ben Volk. On the structure of boolean functions with small spectral norm. *computational complexity*, 26(1):229–273, 2017. doi:10.1007/s00037-015-0110-y.
- 54 Hing Yin Tsang, Chung Hoi Wong, Ning Xie, and Shengyu Zhang. Fourier sparsity, spectral norm, and the log-rank conjecture. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 658–667, 2013. doi:10.1109/FOCS.2013.76.
- 55 Andrew Yao. Probabilistic computations: Toward a unified measure of complexity. In *Proceedings of the 18th Annual Symposium on Foundations of Computer Science, SFCS '77*, pages 222–227. IEEE Computer Society, 1977. doi:10.1109/SFCS.1977.24.
- 56 Andrew Yao. Lower bounds by probabilistic arguments. In *Proceedings of the 24th Annual Symposium on Foundations of Computer Science, SFCS '83*, pages 420–428. IEEE Computer Society, 1983. doi:10.1109/SFCS.1983.30.
- 57 Zhiqiang Zhang and Yaoyun Shi. On the parity complexity measures of boolean functions. *Theoretical Computer Science*, 411(26-28):2612–2618, 2010. doi:10.1016/j.tcs.2010.03.027.

A Appendix

A.1 Direct sums for D

In this appendix, we prove that the best-known direct sum results in the context of deterministic communication complexity can be obtained in the parity decision tree setting. We restate our theorem for convenience below.

► **Theorem 4.** *For any function f and $k \geq 1$,*

1. $D(f^k) \geq k \cdot D(f)^{1/2}$,
2. $D(f^k) \geq k \cdot D(f) / \log \text{spar}(f)$.

Let us first introduce a couple of definitions. Fix a function $f: \{0, 1\}^n \rightarrow \{0, 1\}$. A parity certificate for $f(x)$ is an affine space $S \subseteq \{0, 1\}^n$ such that $x \in S$ and for any $x' \in S$, $f(x) = f(x')$. Similarly to the classical case, the parity certificate complexity $C(f)$ is the smallest codimension of a space that certifies the value $f(x)$ – where the hardest possible

$x \in \{0, 1\}^n$ is taken. We also define $\text{spar}(f) := \|\hat{f}\|_0 = |\{z \mid \hat{f}(z) \neq 0\}|$ for the number of non-zero Fourier coefficients of f . To prove Theorem 4, it is enough to prove a direct sum for parity certificate complexity and employ the following two results:

1. $C(f) \geq D(f)^{1/2}$ [57]
2. $C(f) \geq D(f)/\log \text{spar}(f)$ [54]

► **Lemma 37.** *For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and $k \geq 1$, $C(f^k) \geq k \cdot C(f)$.*

Proof of Lemma 37. Fix an input $x \in \{0, 1\}^n$ attaining $d := C(f)$ and suppose towards contradiction that $C(f^k) < dk$. This implies in particular that there exists an affine space $S \subseteq (\{0, 1\}^n)^k$ described by $m < dk$ equations $Q^T x = b$ (where $Q \in \{0, 1\}^{n \times m}$, $b \in \{0, 1\}^m$) that certifies the value of the input $y \in (\{0, 1\}^n)^k$ which is composed of k copies of x . Define d_i for $i \in [k]$ with:

$$d_i := \dim(\text{col}(Q) \cap W_i) \quad W_i := \{w \in (\{0, 1\}^n)^k : w^j = 0^n \iff j \neq i\}.$$

Observe that $\sum_{i \in [k]} d_i \leq m < dk$ and as such there must be some i^* with $d_{i^*} < k$. Fix for simplicity $i^* = 1$. Using Gaussian elimination, one can re-express $S = S_1 \cap S_2$ where

1. the constraints in S_1 are exclusively on bits of the first copy and
2. any constraint in S_2 has at least one bit of a copy other than the first.

Since S_1 is about the first copy only, it can be identified with a single-copy affine space $S^* \subseteq \{0, 1\}^n$ where $\text{codim}(S^*) = d_1 < k$ in a natural way. Observe that $x \in S^*$ as $y \in S$. Because the codimension of S^* is strictly less than k , there must be some $x' \in S^*$ with $f(x) \neq f(x')$. Note that fixing $x^1 := x'$ leaves the system of linear constraints S_2 feasible and as such there exists $x^2, \dots, x^k \in \{0, 1\}^n$ such that $y' := (x', x^2, \dots, x^k) \in S$: a contradiction since $f(y) \neq f(y')$. ◀

A.2 Omitted case of Theorem 2

► **Lemma 38.** *If $D^\times(f) \leq 6C \cdot \log(n)$, we have $R(f^k) \geq \Omega(k/\log n) \cdot D^\times(f)$.*

Proof. Fix a hard product distribution μ for $D^\times(f)$. If $D_{1/3}(f, \mu) = 0$, the claim follows trivially. Else, we have $D_{1/3}(f, \mu) > 0$ and using Claim 39 with $\varepsilon := 1/6$, it must be that $S_{1/6}(f) \geq 1/6$. Using Claim 17 and Theorem 18, we thus have:

$$R(f^k) \geq D_{1/6}(f^k, \mu^k) \geq S_{1/6}(f^k, \mu^k) \geq k \cdot S_{1/6}(f, \mu) \geq k/6 \geq \Omega(k/\log n) \cdot D^\times(f) \quad \blacktriangleleft$$

▷ **Claim 39.** For any f , product distribution μ and $\varepsilon \geq 0$, we have $D_{\varepsilon+S_\varepsilon(f, \mu)}(f, \mu) = 0$.

Proof. Fix a deterministic decision tree T and consider the zero-query decision tree T' that comes out of applying Algorithm 7 to T . To relate T and T' , we go through Algorithm 5. Again, let \mathcal{T}_0 be the tree obtained by applying Algorithm 5 to T with error zero on Section 6.2. We stress that \mathcal{T}_0 is a randomised decision tree depending on η . On the other hand, T' can be seen as \mathcal{T}_0 with fewer instructions executed. Using Lemma 31, we have:

$$\begin{aligned} \Pr_{x \sim \mu}[T(x) \neq T'(x)] &= \Pr_{x, \eta}[\mathcal{T}_0(x) \neq T'(x)] \\ &\leq \Pr_{x, \eta}[\text{Section 6.2 is executed while running } \mathcal{T}_0(x)] \\ &= \Pr_{x, \rho \sim \mathcal{R}_\mu^x}[T_\rho(x) \text{ makes a query}] \\ &\leq \mathbb{E}_{x, \rho}[q(T_\rho, x)] \\ &= \overline{sq}(T, \mu) \end{aligned}$$

Now, let \mathcal{T} be a randomised parity tree such that $\overline{sq}(\mathcal{T}, \mu) = S_\varepsilon(f, \mu)$ and $\text{err}_f(\mathcal{T}, \mu) \leq \varepsilon$. Let \mathcal{T}' be the randomised parity tree obtained as follows:

1. Sample $T \sim \mathcal{T}$
2. Return T' obtained by applying T to Algorithm 7.

With the analysis above, we obtain $\Pr_{\mathbf{x}, T}[T'(\mathbf{x}) \neq T(\mathbf{x})] \leq \mathbb{E}_{T \sim \mathcal{T}}[\overline{sq}(T, \mu)] = \overline{sq}(\mathcal{T}, \mu)$. We remark that T' makes no queries and has the following error probability:

$$\text{err}_f(T', \mu) \leq \text{err}_f(T, \mu) + \Pr_{\mathbf{x}, T}[T'(\mathbf{x}) \neq T(\mathbf{x})] \leq \varepsilon + \overline{sq}(\mathcal{T}, \mu) = \varepsilon + S_\varepsilon(f, \mu). \quad \triangleleft$$

■ **Algorithm 7** converts an algorithm for $S_\varepsilon < 1$ to a zero-query algorithm.

Input: $x \in \{0, 1\}^n$

Output: $a \in \{0, 1\}$

- 1: Initialize $v \leftarrow \text{root of } T$, $p \leftarrow ?^n$
- 2: $(J, \ell) \leftarrow \text{BuildList}(v, p)$
- 3: **return** $L(\ell)$

A.3 Direct sum for distribution-free discrepancy

► **Theorem 40.** For every function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and $k \geq 1$,

$$k \cdot \text{disc}(f) + 1 \geq \text{disc}(f^{\oplus k}) \geq k \cdot (\text{disc}(f) - 1).$$

Proof. The lower bound is a simple consequence of Lemma 8 by fixing μ to be a distribution such that $\text{disc}(f) = \text{disc}(f, \mu)$ and observing that $\text{disc}(f^{\oplus k}) \geq \text{disc}(f^{\oplus k}, \mu^k)$. The other direction is more interesting as it says that the hardest distribution for $f^{\oplus k}$ is basically k products of the hardest distribution for a single copy f . Let $\|\widehat{f}\|_\infty^* := \min_\mu \|\widehat{F}_\mu\|_\infty$ where μ ranges over all distributions. Using Lemma 9, we obtain the following relation between $\text{disc}(f)$ and $\|\widehat{f}\|_\infty^*$:

$$-\log \|\widehat{f}\|_\infty^* + 1 \geq \text{disc}(f) \geq -\log \|\widehat{f}\|_\infty^*.$$

Therefore, to prove the upper bound, it is enough to show a perfect direct product for $\|\widehat{f}\|_\infty^*$ and apply it k time. To this end, fix some other function $g: \{0, 1\}^n \rightarrow \{0, 1\}$ and let us show that

$$\|\widehat{f \oplus g}\|_\infty^* \geq \|\widehat{f}\|_\infty^* \cdot \|\widehat{g}\|_\infty^*.$$

Where we recall that $f \oplus g: \{0, 1\}^{2n} \rightarrow \{0, 1\}$. We can write $\|\widehat{f}\|_\infty^*$ as the value of the following linear program where the variables describe a distribution μ :

$$\begin{aligned} \min. \quad & c \\ \text{s.t.} \quad & \left| \sum_{x \in \{0, 1\}^n} (-1)^{f(x)} \cdot \mu_x \cdot (-1)^{\langle x, z \rangle} \right| \leq c \quad \forall z \in \{0, 1\}^n \\ & \sum_{x \in \{0, 1\}^n} \mu_x = 1 \\ & \mu_x \geq 0 \quad \forall x \in \{0, 1\}^n \end{aligned} \tag{5}$$

The dual of (5) is:

$$\begin{aligned} \max. \quad & d \\ \text{s.t.} \quad & \sum_{z \in \{0, 1\}^n} (-1)^{f(x)} \cdot \beta_z \cdot (-1)^{\langle x, z \rangle} \geq d \quad \forall x \in \{0, 1\}^n \\ & \sum_{z \in \{0, 1\}^n} |\beta_z| = 1 \end{aligned} \tag{6}$$

Let (β^f, d^f) and (β^g, d^g) be the optimal feasible solutions to (6) with respect to f and g . By the strong duality theorem, it holds that $\|\widehat{f}\|_\infty^* = d^f$ and $\|\widehat{g}\|_\infty^* = d^g$. We now extract a feasible solution for (6) with respect to the function $f \oplus g$. Let $\beta \in \{0, 1\}^{2n}$ be defined with $\beta_{(z_1, z_2)} = \beta_{z_1}^f \cdot \beta_{z_2}^g$ and observe that $(\beta, d^f \cdot d^g)$ is a feasible solution for the dual of $\|\widehat{f \oplus g}\|_\infty^*$. By applying the strong duality theorem again, we have $\|\widehat{f \oplus g}\|_\infty^* \geq d^f \cdot d^g = \|\widehat{f}\|_\infty^* \cdot \|\widehat{g}\|_\infty^*$, as desired. ◀

A.4 Some facts about parity decision trees

Yao's minimax principle is a powerful technique to analyse randomised algorithms – we adapt here the statement to parity trees, but the proof is exactly the same as the original one [55].

► **Lemma 41.** *For any $f: \{0, 1\}^n \rightarrow \{0, 1\}$ and distribution μ over $\{0, 1\}^n$, $R_\varepsilon(f) \geq D_\varepsilon(f, \mu)$.*

The following is a folklore fact relating randomised parity tree complexity and discrepancy [56, 3] which we re-prove in the parity context.

► **Lemma 42.** $D_\varepsilon(f, \mu) \geq \text{disc}(f, \mu) + \log(1 - 2\varepsilon)$ for any $\varepsilon \in [0, 1/2)$.

Proof. Fix a parity decision tree T of depth $d := D_\varepsilon(f, \mu)$ which makes error $\text{err}_f(T, \mu) \leq \varepsilon$, note that

$$\begin{aligned} 1 - 2\varepsilon &\leq \Pr_{x \sim \mu} [T(x) = f(x)] - \Pr_{x \sim \mu} [T(x) \neq f(x)] \\ &= \sum_{S \in \mathcal{L}} \Pr_{x \sim \mu} [T(x) = f(x) \wedge x \in S] - \Pr_{x \sim \mu} [T(x) \neq f(x) \wedge x \in S]. \end{aligned}$$

As $|\mathcal{L}(T)| \leq 2^d$, there exists some $S \in \mathcal{L}(T)$ – an affine subspace – with large correlation:

$$\text{bias}(f, \mu, S) = \left| \Pr_{x \sim \mu} [T(x) = f(x) \wedge x \in S] - \Pr_{x \sim \mu} [T(x) \neq f(x) \wedge x \in S] \right| \geq \frac{1 - 2\varepsilon}{2^d}. \quad \blacktriangleleft$$

► **Lemma 30.** *For any $\alpha > 0$, $R_\alpha(\text{FFO}_n) \leq O(\log n + \log(1/\alpha))$.*

Proof. Let $\text{NOR}_n: \{0, 1\}^n \rightarrow \{0, 1\}$ be the function that checks whether the input is 0^n and rejects otherwise. Observe that one iteration of the sumcheck protocol can be performed in one parity query. More precisely for any $x \in \{0, 1\}^n$, if $s \sim U(\{0, 1\}^n)$ then:

$$\Pr_s[\langle x, s \rangle = 1] = \begin{cases} 1/2 & \text{if } x \neq 0^n \\ 0 & \text{if } x = 0^n \end{cases}.$$

Performing two random checks independently shows that $R(\text{NOR}_n, 1/4) \leq O(1)$. It is a folklore result that a (classical) randomised decision tree can solve FFO_n with probability ε using $O(\log n + \log(1/\varepsilon))$ oracle NOR -queries even if the oracle fails with probability $1/3$ [24, 44]. We highlight that this is an improvement over the naive method that boosts the noisy NOR queries and yields complexity $O(\log(n)^2 \log(1/\varepsilon))$. Recent work [30, §3] revisits this trick in depth for communication complexity and their result can be re-interpreted in the context of parity decision trees as follows:

$$\forall f: R(f, \varepsilon) \leq O(D^{\text{NOR}}(f) + \log(1/\varepsilon)).$$

Thus, plugging in $f = \text{FFO}$ and noting that $D^{\text{NOR}}(\text{FFO}_n) \leq \log n$ (with binary search), we get the desired result. ◀

▷ **Claim 11.** For any $f: \{0,1\}^n \rightarrow \{0,1\}$, μ over $\{0,1\}^n$ and $\varepsilon, \delta \geq 0$, $D_{\varepsilon+\delta}(f, \mu) \leq \overline{D}_\varepsilon(f, \mu)/\delta$.

Proof. Let \mathcal{T} be a randomised PDT satisfying that $d := \bar{q}(\mathcal{T}, \mu) = \overline{D}_\varepsilon(f, \mu)$ and $\text{err}_f(\mathcal{T}, \mu) \leq \varepsilon$. To prove the lemma, it suffices to construct a deterministic parity tree T of depth $T \leq d/\gamma$ with $\text{err}_f(T, \mu) \leq \varepsilon + \gamma$. Sample $\mathbf{T} \sim \mathcal{T}$. We construct a new tree \mathbf{T}' by pruning \mathbf{T} as follows: We remove all the nodes of \mathbf{T} of depth greater than d/δ . If any node of depth d/δ becomes a leaf, we label it with an arbitrary bit. Note that \mathbf{T}' has depth $\leq d/\delta$. Finally, let \mathcal{T}' denote the distribution over \mathbf{T}' inherited from \mathcal{T} .

We observe that for each $x \in \{0,1\}^n$, both $\mathbf{T}(x) = f(x)$ and $\mathbf{T}'(x) \neq f(x)$ happen only if $q(\mathbf{T}, x) > d/\gamma$. Moreover, by Markov's inequality,

$$\Pr_{\substack{\mathbf{T} \sim \mathcal{T} \\ \mathbf{x} \sim \mu}}[q(\mathbf{T}, \mathbf{x}) > d/\gamma] \leq \frac{\bar{q}(\mathbf{T}, \mu)}{d/\gamma} = \gamma.$$

Therefore, $\text{err}_f(\mathcal{T}', \mu) \leq \text{err}_f(\mathcal{T}, \mu) + \gamma \leq \varepsilon + \gamma$. By an averaging argument, there exists some $T \in \text{supp}(\mathcal{T}')$ of depth $\leq d/\delta$ that computes f with error $\text{err}_f(T, \mu) \leq \varepsilon + \gamma$, as desired. ◁

A.5 Omitted proofs of Section 6

In this appendix, we prove Claim 29, an alternative description for the distributions of Section 6.1. Let $p^1, p^2 \in \{0, \star, ?\}^n$. We write $p^1 \sim p^2$ if p^1 and p^2 are consistent over their non-? entries. That is, $p^1 \sim p^2$ if for all $j \in [n]$, if $p_j^1 \neq ?$ and $p_j^2 \neq ?$, then $p_j^1 = p_j^2$. Claim 29 follows from Claims 43 and 44.

▷ **Claim 43.** For every reachable state (v, p) , consistent $x \in \{0,1\}^n$ and $\rho \in \{0, \star\}^n$, $\mathcal{R}^{v,p,x} \equiv \widehat{\mathcal{R}}^{p,x}$.

Proof. Upon inspection of $\widehat{\mathcal{R}}^{v,p}$, it is enough to prove that for all $x \in \{0,1\}^n$:

$$\Pr_{\rho \sim \mathcal{R}^{v,p,x}}[\rho = \rho] = \prod_{j \in S_{\neq ?}^p} \mathbb{1}[\rho_j = p_j] \times \prod_{\substack{j \in S_{?}^p \\ x_j = 1}} \mathbb{1}[\rho_j = \star] \times \prod_{\substack{j \in S_{?}^p \\ x_j = 0}} \begin{cases} \delta_j/(2 - \delta_j) & \text{if } \rho_j = \star \\ 1 - \delta_j/(2 - \delta_j) & \text{if } \rho_j = 0 \end{cases}.$$

Fix $x \in \{0,1\}^n$. We prove this by induction on the state space (v, p) consistent with x . The entry-point of the state space is $(\text{root}(T), ?^n)$. In this case, the statement holds by definition. Suppose now that the statement is true for state (v, p) . Depending on the value of ρ , there are several next state (v', p') possible. Observe however that the next vertex of T to be visited does not depend on ρ , as it is fixed to be $v' := \text{child}(v, \langle x, Q^v \rangle)$. For any fixed $\rho \in \{0, \star\}^n$, we have:

$$\begin{aligned} \Pr_{\rho \sim \mathcal{R}^{v',p',x}}[\rho = \rho] &= \Pr_{\mathbf{x} \sim \mu, \rho \sim R_\mu^{\mathbf{x}}}[\rho = \rho \mid (v', p') \text{ is reached and } \mathbf{x} = x] \\ &= \frac{\Pr_{\mathbf{x}, \rho}[\rho = \rho \text{ and } (v', p') \text{ is reached and } \mathbf{x} = x]}{\Pr_{\mathbf{x}, \rho}[(v', p') \text{ is reached and } \mathbf{x} = x]}. \end{aligned}$$

Note that there can be only one state from which (v', p') can be reached, namely (v, p) . Indeed, suppose that there is another state (v, p^*) from which (v', p') can be reached. Then (v, p) and (v, p^*) have a common ancestor (u, q) . Since the paths diverged after (u, q) , it must be that $p \approx p^*$ and thus $p^* \approx p'$: a contradiction. Thus, we have the following equivalence:

$$(v', p') \text{ is reached} \iff (v, p) \text{ is reached and } \rho \sim p'.$$

Therefore, we have:

$$\Pr_{\rho \sim \mathcal{R}^{v',p',x}}[\rho = \rho] = \frac{\Pr_{\rho \sim \mathcal{R}^{v,p,x}}[\rho = \rho] \cdot \mathbb{1}[\rho \sim p']}{\Pr_{\rho \sim \mathcal{R}^{v,p,x}}[\rho \sim p']}. \quad (7)$$

We can now use the inductive hypothesis on (v, p) . Since $\rho \sim p'$ implies $\rho \sim p$, the numerator of (7) simplifies to:

$$\prod_{j \in S_{\neq}^{p'}} \mathbb{1}[\rho_j = p_j] \times \prod_{\substack{j \in S_{\neq}^p \\ x_j=1}} \mathbb{1}[\rho_j = \star] \times \prod_{\substack{j \in S_{\neq}^p \\ x_j=0}} \begin{cases} \delta_j/(2 - \delta_j) & \text{if } \rho_j = \star \\ 1 - \delta_j/(2 - \delta_j) & \text{if } \rho_j = 0 \end{cases}.$$

Let $\Delta = S_{\neq}^p \setminus S_{\neq}^{p'}$ and observe that the denominator of (7) is equal to:

$$\prod_{\substack{j \in \Delta \\ x_j=1}} \mathbb{1}[\rho_j = \star] \times \prod_{\substack{j \in \Delta \\ x_j=0}} \begin{cases} \delta_j/(2 - \delta_j) & \text{if } \rho_j = \star \\ 1 - \delta_j/(2 - \delta_j) & \text{if } \rho_j = 0 \end{cases}. \quad \triangleleft$$

▷ **Claim 44.** For every reachable state (v, p) and $x \in \{0, 1\}^n$, $\mathcal{X}^{v,p} \equiv \widehat{\mathcal{X}}^{v,p}$.

Proof. Fix some (v, p) and $x \in \{0, 1\}^n$. Upon inspection of $\widehat{\mathcal{X}}^{v,p}$, it is enough to prove that

$$\Pr_{\mathbf{x} \sim \mathcal{X}^{v,p}}[\mathbf{x} = x] = M(x, v, p) \cdot \prod_{j \in S_{\neq}^p} 1 - \delta_j/2 - x_j \cdot (1 - \delta_j),$$

where $M(x, v, p)$ is an indicator set to 1 if and only if for all $j \in [n]$, $p_j = 0$ implies $x_j = 0$ and $\langle x, Q^u \rangle = b^u$ for all $u \in \text{path}(v)$. By Baye's rule we have:

$$\Pr_{\substack{\mathbf{x} \sim \mathcal{X}^{v,p} \\ \rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}}[\mathbf{x} = x] = \Pr_{\substack{\mathbf{x} \sim \mu \\ \rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}}[\mathbf{x} = x \mid (v, p) \text{ is reached on } (\mathbf{x}, \rho)] = \frac{p(x)}{\sum_{x' \in \{0,1\}^n} p(x')},$$

where

$$p(x) := \Pr_{\mathbf{x}, \rho}[\mathbf{x} = x] \cdot \Pr_{\mathbf{x}, \rho}[(v, p) \text{ is reached on } (\mathbf{x}, \rho) \mid \mathbf{x} = x].$$

To analyse $p(x)$, we have:

$$\Pr_{\substack{\mathbf{x} \sim \mu \\ \rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}}[\mathbf{x} = x] = \Pr_{\mathbf{x} \sim \mu}[\mathbf{x} = x] = \prod_{j \in [n]} \Pr_{\mathbf{x} \sim \mu}[\mathbf{x}_j = x_j] = \prod_{j \in [n]} 1 - (\delta_j/2) - x_j \cdot (1 - \delta_j).$$

On the other hand, the second component of $p(x)$ is clearly zero if $M(x, v, p) = 0$. For instance, v cannot be reached if x does not satisfy all equations on the path to v . Thus, we have:

$$\begin{aligned} & \Pr_{\substack{\mathbf{x} \sim \mu \\ \rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}}[(v, p) \text{ is reached on } (\mathbf{x}, \rho) \mid \mathbf{x} = x] \\ &= \Pr_{\rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}[(v, p) \text{ is reached on } (\mathbf{x}, \rho)] \\ &= M(x, v, p) \cdot \Pr_{\rho \sim \mathcal{R}_{\mu}^{\mathbf{x}}}[\rho \sim p] \\ &= M(x, v, p) \cdot \prod_{j \in S_0^p} \frac{2 - 2\delta_j}{2 - \delta_j} \cdot \prod_{j \in S_{\neq}^p} \left(\frac{\delta_j}{2 - \delta_j} \right)^{1-x_j}. \end{aligned}$$

Combining those two observations, we get:

$$p(x) = M(v, p, x) \cdot \prod_{j \in S_{\neq}^p} (1 - \delta_j/2 - x_j \cdot (1 - \delta_j)) \cdot \prod_{j \in S_0^p} (1 - \delta_j) \cdot \prod_{j \in S_{\neq}^p} \delta_j/2.$$

16:38 Direct Sums for Parity Decision Trees

Observe that the last two products do not involve x at all and can thus be cancelled in the initial expression:

$$\Pr_{\mathbf{x} \sim \mathcal{X}^{v,p}}[\mathbf{x} = x] = \frac{p'(x)}{\sum_{x'} p'(x)} \text{ where } p'(x) = M(x, v, p) \cdot \prod_{j \in S_i^p} (1 - \delta_j/2 - x_j \cdot (1 - \delta_j)).$$

Finally, observe that $M(x, v, p)$ fixes the value of all the bits of x except for S_i^p . Thus, the summation in the denominator equals 1 and the claim follows. \triangleleft