

The distinguishability of product distributions by read-once branching programs

John Steinberger

Institute for Theoretical Computer Science, Tsinghua University
jpsteinb@gmail.com

Abstract. We improve the main result of Brody and Verbin [7] from FOCS 2010 on the power of constant-width branching programs to distinguish product distributions. Specifically, we show that a coin must have bias at least $\Omega(1/\log(n)^{w-2})$ to be distinguishable from a fair coin by a width w , length n read-once branching program (for each constant w), which is a tight bound. Our result introduces new techniques, in particular a novel “interwoven hybrid” technique and a “program randomization” technique, both of which play crucial roles in our proof. Using the same techniques, we also succeed in giving tight upper bounds on the maximum influence of monotone functions computable by width w read-once branching programs.

* This work was supported in part by the National Basic Research Program of China Grant 2011CBA00300, 2011CBA00301, the National Natural Science Foundation of China Grant 61033001, 61361136003.

1 Introduction

In [7] Brody and Verbin studied the question of distinguishing flips of a coin with a slight bias towards heads from those of a coin with a slight bias towards tails. More precisely, say that a coin is ϵ -biased if $\Pr[\text{Heads}] = \frac{1}{2} + \epsilon$. Given n flips of a coin which is either ϵ -biased or $(-\epsilon)$ -biased, the question is to determine which type of bias is present. Since taking a majority vote of the tosses constitutes an optimal distinguishing strategy this question is uninteresting when the distinguisher is powerful enough to count (in which case a bias of $\epsilon = \Omega(1/\sqrt{n})$ is both necessary¹ and sufficient to distinguish with constant advantage²). However, the problem seems both natural and interesting for space-bounded distinguishers, and in particular for distinguishers having only a constant amount of space.

As their main result, Brody and Verbin [7] give bounds on the ability of constant width *read-once branching programs* (ROBPs) to distinguish biased coins. A read-once branching program is a model of (non-uniform) space bounded computation in which each bit of input is accessed only once, in order. (We give a formal definition of read-once branching programs in Section 2. A glance at Figure 1, however, should suffice to understand the model.) They show, among others, that ROBPs of width $w \geq 3$ can distinguish coins of surprisingly small bias: by computing a recursive tribes function, a length n ROBP of width w can distinguish an ϵ -biased coin from a $(-\epsilon)$ -biased coin already for $\epsilon = 1/\log(n)^{w-2}$. (By “can distinguish” we mean, here and later, “can distinguish with constant (i.e. $\Omega(1)$) advantage”.) In particular, a width 3 ROBP of length n can distinguish a $(1/\log n)$ -biased coin from a $(-1/\log n)$ -biased coin. (This last observation was also made, essentially, by Braverman et al. [6].)

On the lower bound side, Brody and Verbin show that, for constant w , a length n width w ROBP cannot distinguish $\pm\epsilon$ -biased coins unless $\epsilon = \Omega(1/\log(n)^w)$. The lower bound is therefore off from the upper bound by a factor $\log(n)^2$, which seems substantial for programs of small width (e.g., width 3, for which the upper bound is $\epsilon = O(1/\log(n))$ and the lower bound is $\epsilon = \Omega(1/\log(n)^3)$). In this paper we give an improved lower bound that matches the upper bound of [7]. Namely, we show that, for constant³ w , the smallest bias ϵ that can be distinguished by a width w length n ROBP is $\Omega(1/\log(n)^{w-2})$. Our analysis is also shorter than Brody and Verbin’s. More interestingly than simply achieving a tight lower bound, however, is the fact that our result introduces new proof techniques that could be of independent interest for the study of ROBPs and, more generally, for the problem of derandomizing space-bounded computations. These techniques are described further below.

Note that a sequence of n independent tosses of a biased coin is a special case of a product distribution. (A sequence of random variables $X = (X_i)_{i=1}^n$ is a *product distribution* if and only the X_i ’s are (totally) independent.) One can consider, more generally, the power of a length n read-once branching program whose edges are labeled by elements of some finite alphabet Σ at distinguishing two product distributions $X = (X_i)_{i=1}^n, Y = (Y_i)_{i=1}^n$ where $X_i, Y_i \in \Sigma$ for all i .

Generalizing results on the distinguishability of biased coins to the distinguishability of arbitrary product distributions presupposes some kind of metric for measuring the closeness of two product distributions (i.e., requires generalizing the parameter ϵ). As explained in [7], it makes more sense, in this context, to measure closeness by probability ratios (bounding these to be near 1) rather than by probability differences (bounding these to be near 0). We say two product distributions $X = (X_i)_{i=1}^n \in \Sigma^n, Y = (Y_i)_{i=1}^n \in \Sigma^n$ are ϵ -close

¹ A possible approach for proving necessity (since the relevant statistical distance is not so obvious to upper bound from first principles) is to use Hellinger distance. See for example [4].

² The *advantage* of a distinguisher D at distinguishing distributions X and Y is $|\Pr[D(X) = 1] - \Pr[D(Y) = 1]|$. See Section 2 for more precise definitions.

³ The fact that w is constant in particular implies that $\Omega(\cdot)$ - and $O(\cdot)$ - notation refers exclusively to function growth with respect to n , and may hide constants that depend on w . In [7] the hidden constant is 1000^{-w} . Our hidden constant is 2^{-w} .

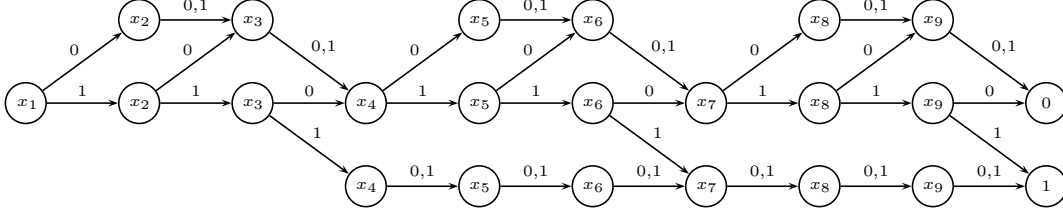


Fig. 1: A width 3 read-once branching program (computing a tribes function).

in ratio if for every $1 \leq i \leq n$ and for every $\alpha \in \Sigma$, either $\Pr[X_i = \alpha] = \Pr[Y_i = \alpha] = 0$, or else $\Pr[X_i = \alpha] \neq 0, \Pr[Y_i = \alpha] \neq 0$ and

$$\frac{\Pr[X_i = \alpha]}{\Pr[Y_i = \alpha]} \geq 1 - \epsilon, \quad \frac{\Pr[Y_i = \alpha]}{\Pr[X_i = \alpha]} \geq 1 - \epsilon.$$

It is easy to see, for example, that an ϵ -biased coin is 4ϵ -close in ratio to a $(-\epsilon)$ -biased coin.

Our results are most easily phrased and proved in the context of ϵ -close in ratio product distributions. Our main result is that a length n ROBP of constant width w cannot distinguish two ϵ -close in ratio product distributions unless $\epsilon = \Omega(1/\log(n)^{w-2})$. This directly implies our lower bound on the distinguishability of ϵ -biased coins, and also matches the upper bounds of Brody and Verbin (given by ϵ -biased coins). In fact, ϵ -close in ratio product distributions were already considered by Brody and Verbin themselves, who, by reducing to the case of ϵ -biased coins, proved that $\epsilon = \Omega(1/\log(|\Sigma|n)^{3w})$ is necessary to distinguish two ϵ -close in ratio product distributions over Σ^n , and that $\epsilon = \Omega(1/\log(n)^{2w})$ is necessary when $|\Sigma| = 2$. Our own lower bound shows there is essentially no difference between the cases $|\Sigma| = 2$ and $|\Sigma| > 2$: a larger alphabet size does not help the distinguisher.

The techniques used in our proof are roughly threefold. We use, firstly, the *collision lemma* of Brody and Verbin, which is a structural observation about ROBPs that are optimal distinguishers, and which we strengthen slightly for our purposes. A second component of the proof is a hybrid argument whose two endpoints are the product distributions $(X_i)_{i=1}^n$ and $(Y_i)_{i=1}^n$. Here the bits that change distribution from one hybrid to the next form each time an arithmetic progression (which is important for the argument). As these various arithmetic progressions are parallel and interleaved, we call our set of hybrids a set of *interwoven hybrids* (we are not aware of a similar set of such hybrids being used before). Our hybrid argument replaces a more standard random restriction argument by Brody and Verbin. Finally, the third main proof technique we use is *program randomization* which, in a nutshell, randomizes the distinguisher in order to compensate for certain helpful initial modifications made to the input distributions $(X_i)_{i=1}^n$ and $(Y_i)_{i=1}^n$ (see Section 4 for more details). Program randomization is also an original contribution of the paper. While it is crucial to the final bounds, we consider it of secondary importance compared to the collision lemma and to the interwoven hybrid technique.

In Appendix D we give a second application of the same basic set of techniques (program randomization excluded) to the upper bounding of the maximum total influence of monotone functions computable by width w ROBPs. Our main result is that a ROBP of width $w \geq 2$ and of length $n \geq 2$ that computes a monotone function has total influence at most $4\lceil 1.5 \log(n) \rceil^{w-2}$. This bound is also tight, as can be verified by considering a recursive tribes function.

2 Definitions

A branching program of width w and length n is a directed acyclic graph with n layers of w nodes each and a final layer with two nodes (accept and reject). Each non-output node is labeled by a coordinate $(k, j) \in [n] \times [w]$; output nodes are labeled by coordinates $\{n+1\} \times \{1, 2\}$, with $(n+1, 1)$ being the accept node and $(n+1, 2)$ being the reject node.

A node is *in layer* k , $1 \leq k \leq n+1$, if its label is of the form (k, \cdot) . The edges of the graph are labeled by elements of the *input alphabet* Σ (a finite set). Each node in every layer $k \leq n$ has one outgoing edge labeled α for each element $\alpha \in \Sigma$ whose endpoint is a node in layer $k+1$. The branching program has a designated *start node* in the first layer, typically the node $(1, 1)$. The computation of a branching program of length n on a string $x = x_1 \cdots x_n \in \Sigma^n$ is defined the natural way, by following the edge labeled x_i at step i , starting from the start node. We note the type of branching program just described is *read-once* since each character of x is examined at exactly one layer of the program.

Let f be a (read-once) branching program of length n and width w . If α is an element of the input alphabet Σ and $k \in [n]$, the α -*transition function of f at layer k* is the function $\tau_\alpha : [w] \rightarrow [w]$ such that $\tau_\alpha(i) = z$ iff the edge labeled α leaving node (k, i) has endpoint $(k+1, z)$. We say τ_α *contains a collision* if τ_α is not a permutation, i.e. if $\tau_\alpha(i) = \tau_\alpha(j)$ for some $i \neq j$.

The k -th layer of a ROBP f *equals* the j -th layer of a ROBP g if f and g have the same width w , are defined over the same input alphabet Σ , and if the α -transition function of f at layer k is identical to the α -transition function of g at layer j for every $\alpha \in \Sigma$.

The statistical distance of two random variables X, Y of same range is written $\Delta(X, Y)$. Namely, if X and Y take values in a set S , then

$$\Delta(X, Y) = \frac{1}{2} \sum_{b \in S} |\Pr[X = b] - \Pr[Y = b]|.$$

If f is a ROBP of length n over the alphabet Σ and if $X, Y \in \Sigma^n$ are two random variables, then f 's *advantage* at distinguishing X and Y is defined as the statistical distance

$$\Delta(f(X), f(Y)).$$

(We note this is a statistical distance between two probability distributions on the output nodes of f .) This differs from the traditional definition of f 's advantage as $|\Pr[f(X) = 1] - \Pr[f(Y) = 1]|$, but it is easy to see the two definitions are equivalent.

We write $X \sim X'$ when X, X' induce identical probability distributions over their (identical) ranges.

3 Results

Our main result is an upper bound on the advantage $\Delta(f(X), f(Y))$ of a width w ROBP f at distinguishing ϵ -close product distributions $X, Y \in \Sigma^n$ for an arbitrary finite alphabet Σ . While our original interest lies with constant values of w , our main result, given by the next theorem, is slightly more general, as it also allows "small" non-constant w .

Theorem 1. *There is a function⁴ $\lambda(n) = o(1)$ such that for any positive integers n, w with $2 \leq w \leq \log n / \log \log n$, for any product distributions $X, Y \in \Sigma^n$ that are ϵ -close in ratio, and for any read-once branching program f over the alphabet Σ of width w and length n ,*

$$\Delta(f(X), f(Y)) \leq \epsilon(2 \log(n))^{w-2}(1 + \lambda(n)). \tag{1}$$

⁴ I.e., $\lim_{n \rightarrow \infty} \lambda(n) = 0$.

In particular, if w is constant, ϵ needs to be at least $\Omega(1/\log(n)^{w-2})$ in order for X and Y to be distinguishable with constant advantage, where the hidden constant⁵ (depending on w but not on n) is 2^{-w} . This lower bound on ϵ is tight up to a constant factor: as shown in [7], width w , length n ROBPs can distinguish coins of bias $\pm\epsilon$ already for $\epsilon = O(1/\log(n)^{w-2})$. (For full disclosure, the hidden constant in the latter $O(\cdot)$ is 3^w ; hence, there is still a gap between the upper and lower bounds as far as the constant factors are concerned.) Theorem 1 will be proved as a corollary of a slightly more fine-grained statement (Theorem 4) in Appendix B.

In Appendix D we also prove an upper bound on the maximum total influence of monotone ROBPs (see relevant definitions in Appendix D), which constitutes our second main result and is as follows:

Theorem 2. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a monotone boolean function computable by a ROBP of width w and length n . Then*

$$\text{Inf}(f) \leq 4[1.5 \log(n)]^{w-2}.$$

For constant w this bound is also tight up to a multiplicative factor, as can be seen using a recursive tribes function of depth $w - 1$ with the same tribe sizes as in [7]. (See also [2].)

The proof of Theorem 2 is given in Appendix D and the proof of Theorem 1 is divided between Appendices A, B and C. We proceed with an outline of the proof of Theorem 1.

4 Proof Overview

This section gives a self-contained overview of the proof of Theorem 1. For simplicity, we sketch the proof for the case $\Sigma = \{0, 1\}$ (which anyway captures the full complexity of the problem). Moreover, we first sketch the proof for the case of distinguishing $\pm\epsilon$ -biased coins and, later, discuss how to handle ϵ -close in ratio distributions (which, indeed, require an additional idea).

Let $X \in \{0, 1\}^n$ be the product distribution of an ϵ -biased coin, and let $Y \in \{0, 1\}^n$ be the product distribution of a $(-\epsilon)$ -biased coin. Let X_j be the j -th bit of X .

Let \mathcal{F}_w be the set of all (binary) ROBPs of length n and width w (the parameter n is elided for simplicity). Let

$$\delta_w = \max_{f \in \mathcal{F}_w} \Delta(f(X), f(Y))$$

be the maximum distinguishing advantage. The proof bounds δ_w by establishing the recurrence

$$\delta_w = O(\log n)\delta_{w-1} + o(1) \tag{2}$$

and by showing that $\delta_2 \leq \epsilon$. In fact the $o(1)$ term is $1/\text{poly}(n)$, so that recursively “unfolding” the inequality gives

$$\delta_w \leq O(\log n)^{w-2}\epsilon + o(1).$$

Tweaking the constants then yields Theorem 1. We now sketch how (2) is established.

Let

$$\mathcal{F}_w^{\max} = \{f \in \mathcal{F}_w : \Delta(f(X), f(Y)) = \delta_w\}$$

be the set of “best distinguishers”. Note \mathcal{F}_w^{\max} is nonempty since \mathcal{F}_w is finite. A crucial observation, due to Brody and Verbin [7], is that \mathcal{F}_w^{\max} contains an element f_0 in which every transition function is either the

⁵ In fact, for constant w , (1) can be replaced with $\Delta(f(X), f(Y)) \leq \epsilon(\delta \log(n))^{w-2} + \lambda_\delta(n)$ where $\delta > 1$ is any constant and where $\lambda_\delta(n) \rightarrow 0$ now depends on δ . Thus a sharper statement would say that the hidden constant is really “ δ^w for any $\delta > 1$ ”. We refer to Appendix B for more details.

identity from $[w]$ to $[w]$, or else is not a permutation of $[w]$ at all, but contains a collision. We call an ROBP with this property a *collision ROBP*, or cROBP for short. To upper bound δ_w it thus suffices to upper bound $\Delta(f(X), f(Y))$ for an arbitrary cROBP f of length n and width w . (A nearly identical observation is called the *collision lemma* in [7]. We maintain this terminology, even while our own collision lemma is slightly different. The difference is explained in Section 5.)

Let f , therefore, be a cROBP of length n and width w . By dropping layers of f at which both transition functions are the identity (these have no effect), one can assume that every layer of f has at least one transition function with a collision.

To upper bound $\Delta(f(X), f(Y))$ we use a hybrid argument over distributions $Z_0, \dots, Z_{c \log(n)}$ on $\{0, 1\}^n$, such that $Z_0 = X$ and $Z_{c \log(n)} = Y$. Here $c > 0$ is a constant we will set later (in fact, $c = 2$ will do). More precisely, assuming $c \log(n)$ is an integer (otherwise substitute $\lceil c \log(n) \rceil$ for $c \log(n)$ throughout), Z_i is the product distribution whose j -th coordinate $Z_{i,j}$ is given by

$$Z_{i,j} = \begin{cases} Y_j & \text{if } (j \bmod c \log(n)) < i, \\ X_j & \text{otherwise.} \end{cases}$$

For example, Z_1 is the distribution such that

$$Z_{1,j} = \begin{cases} Y_j & \text{if } j \equiv 0 \pmod{c \log(n)}, \\ X_j & \text{otherwise.} \end{cases}$$

Clearly, then, $Z_0 = X$ and $Z_{c \log(n)} = Y$.

We note that Z_i and Z_{i+1} differ on a set of bits whose indices form an arithmetic progression of step size $c \log(n)$. This is the key feature of these hybrids; in fact any sequence of $c \log(n)$ hybrids with this property, starting with X and ending with Y , would do as well (there are $(c \log n)!$ possible such sequences). Let $\mathcal{Z}_i \subseteq [n]$ be the set of bits at which (the definitions of) Z_i and Z_{i+1} differ. We call $Z_0, \dots, Z_{c \log(n)}$ a sequence of “interwoven hybrids” because $\mathcal{Z}_0, \dots, \mathcal{Z}_{c \log(n)-1}$ are interwoven arithmetic progressions of equal step size.

By a standard argument, it suffices to bound the distance $\Delta(f(Z_i), f(Z_{i+1}))$ between two neighboring hybrids. Let $\bar{Z} \in \{0, 1\}^{[n] \setminus \mathcal{Z}_i}$ be the value of Z_i, Z_{i+1} on the bits outside \mathcal{Z}_i . Fixing a value of \bar{Z} induces (in the natural way) a width w , length $|\mathcal{Z}_i|$ ROBP $f_{\bar{Z}} : \{0, 1\}^{\mathcal{Z}_i} \rightarrow \{0, 1\}$ taking as input the bits in \mathcal{Z}_i . Let $X' \in \{0, 1\}^{\mathcal{Z}_i}$ be an ϵ -biased coin, and let $Y' \in \{0, 1\}^{\mathcal{Z}_i}$ be a $(-\epsilon)$ -biased coin. Then $f(Z_i)$ is equidistributed to $f_{\bar{Z}}(X')$, and $f(Z_{i+1})$ is equidistributed to $f_{\bar{Z}}(Y')$, with randomness taken over \bar{Z}, X', Y' . By elementary properties of statistical distance, one has

$$\Delta(f(Z_i), f(Z_{i+1})) \leq \mathbb{E}_{\bar{Z}} \Delta(f_{\bar{Z}}(X'), f_{\bar{Z}}(Y')). \quad (3)$$

The crucial observation is that, in fact, $f_{\bar{Z}}$ is (equivalent to) a width $w - 1$ ROBP with high probability over \bar{Z} . This uses the fact that f is a cROBP. Consider the transition functions τ_0, τ_1 at layer k of $f_{\bar{Z}}$. By definition of $f_{\bar{Z}}$, these transition functions depend on $c \log(n) - 1$ consecutive bits of \bar{Z} . Let these $c \log(n) - 1$ bits have indices $j_1, \dots, j_{c \log(n)-1}$ in f . To picture how τ_0, τ_1 are induced by \bar{Z} , consider w (distinguishable) pebbles placed on the w nodes of f at layer j_1 . Then for a fixed value of \bar{Z} , we can assign in the natural way a path to each pebble, starting at layer j_1 and ending at layer $j_{c \log(n)-1} + 1 = j_1 + c \log(n) - 1$. Then τ_0 is the composition of the 0-transition τ'_0 at layer $j_1 - 1$ of f with the function from $[w]$ to $[w]$ given by the pebble paths, and likewise τ_1 is the composition of the 1-transition τ'_1 at layer $j_1 - 1$ of f with the same pebble paths. Moreover, note that if two pebbles collide, they cannot separate again; thus, if two pebbles collide, τ_0 and τ_1 have at most $w - 1$ nodes in the union of their ranges.

Since f is a cROBP, there are values $b_1, \dots, b_{c \log(n)-1} \in \{0, 1\}$ such that the b_i -transition at layer j_i of f has a collision. By the above remarks, if the j_h -th bit⁶ of \bar{Z} is equal to b_h for any $1 \leq h \leq c \log(n) - 1$, then τ_0, τ_1 have joint range of size at most $w - 1$. But any coordinate of \bar{Z} is equal to a given binary value with probability at least $\frac{1}{2} - \epsilon$, since each coordinate of \bar{Z} is distributed either according to X or according to Y ; namely,

$$\Pr[\bar{Z}_{j_h} = b_h] \geq \frac{1}{2} - \epsilon \quad (4)$$

for any $1 \leq h \leq c \log(n) - 1$. Thus the probability that no collisions occur among the pebbles as they travel from layer j_1 to layer $j_{c \log(n)-1} + 1$ is, in the worst case, at most

$$\left(\frac{1}{2} + \epsilon\right)^{c \log(n)-1} \approx \frac{1}{n^c}.$$

(Where we use $\epsilon = o(1)$; we are being, here, a bit informal for the sake of the proof sketch.) By a union bound, the probability that *any* of the $n/c \log(n)$ pairs of transition functions of $f_{\bar{Z}}$ do not have joint range of size at most $w - 1$ is at most $\approx 1/n^{c-1} \log(n)$. Thus, $f_{\bar{Z}}$ can be written as a width $w - 1$ ROBP with probability at least $\approx 1 - \frac{1}{n^{c-1} \log(n)}$, with the probability taken over \bar{Z} . This allows us to upper bound (3) by

$$\mathbb{E}_{\bar{Z}} \Delta(f_{\bar{Z}}(X'), f_{\bar{Z}}(Y')) \leq O\left(\frac{1}{n^{c-1} \log(n)}\right) + \delta_{w-1}.$$

(In fact, one could even replace δ_{w-1} with the advantage of the best distinguisher of length $n/c \log(n)$ and of width $w - 1$, but such an optimization has little effect for constant-width ROBPs.) Finally, summing together the distances between the $c \log(n)$ pairs of neighboring hybrids, one thus obtains that

$$\begin{aligned} \delta_w = \Delta(f(X), f(Y)) &\leq c \log(n) O\left(\frac{1}{n^{c-1} \log(n)}\right) + c \log(n) \delta_{w-1} \\ &= O\left(\frac{1}{n^{c-1}}\right) + c \log(n) \delta_{w-1}, \end{aligned}$$

establishing (2).

Finishing the proof also requires showing that $\delta_2 \leq \epsilon$. This is not trivial and requires the collision lemma as well as a coupling argument. We refer to Sections 5 for more details.

When working with arbitrary product distributions that are ϵ -close in ratio, the above analysis breaks down in one crucial place: even when ϵ is very small, there is no guarantee that $\Pr[\bar{Z}_{j_h} = b_h]$ will be near $\frac{1}{2}$, cf. (4). Instead, $\Pr[\bar{Z}_{j_h} = b_h]$ could be arbitrarily close to 0. The probability that no collisions occur among the pebbles could thus be arbitrarily close to 1, and, therefore, $f_{\bar{Z}}$ is no longer equivalent to a width $w - 1$ program with high probability.

In view of circumventing this (apparently complete) breakdown of the argument, note first that we do not care if $\Pr[\bar{Z}_{j_h} = b_h]$ is low if *both* the 0-transitions and 1-transitions at layer j_h contain collisions; in this case, indeed, we obtain width reduction with probability 1. Assume, therefore, wlog, that the 0-transition at layer j_h contains a collision, whereas the 1-transition is the identity function. Moreover assume that $\Pr[\bar{Z}_{j_h} = 0]$ is low. To be concrete, say

$$\Pr[X_{j_h} = 0] = \frac{1}{n}, \quad \Pr[Y_{j_h} = 0] = \frac{0.9}{n}. \quad (5)$$

⁶ We index the bits of \bar{Z} by their original index in Z_i, Z_{i+1} .

Such values would be compatible with $\epsilon = 0.1$, and would imply $\Pr[\overline{Z}_{j_h} = 0] \leq \frac{1}{n}$.

The intuition is that in the case above, \overline{Z}_{j_h} is quite likely to be equal to 1, which is an identity transition function, and therefore *it is quite likely the program does nothing at all at layer j_h* . Namely, the program is, with high probability, not reacting to input bit j_h , and layer j_h is therefore “wasted with high probability” for the program.

To leverage this intuition, let f^\perp be the ROBP identical to f , but whose 0-transition function and 1-transition function at layer j_h are both the identity. Note that with high probability over the input distributions X and Y , f^\perp computes the same as f (assuming (5)). We define a random ROBP f^* to be

$$f^* = \begin{cases} f^\perp & \text{with probability } 1 - \frac{1}{\gamma}, \\ f & \text{with probability } \frac{1}{\gamma} \end{cases}$$

where $\gamma \geq 1$ is chosen as large as possible such that the distributions X^*, Y^* defined by

$$X_k^* = \begin{cases} X_k & \text{if } k \neq j_h, \\ 0 & \text{with probability } \gamma \Pr[X_{j_h} = 0] \text{ if } k = j_h, \text{ and} \\ 1 & \text{with probability } 1 - \gamma \Pr[X_{j_h} = 0] \text{ if } k = j_h \end{cases} \quad (6)$$

$$Y_k^* = \begin{cases} Y_k & \text{if } k \neq j_h, \\ 0 & \text{with probability } \gamma \Pr[Y_{j_h} = 0] \text{ if } k = j_h, \text{ and} \\ 1 & \text{with probability } 1 - \gamma \Pr[Y_{j_h} = 0] \text{ if } k = j_h \end{cases} \quad (7)$$

are ϵ -close in ratio. Note that $f^*(X^*), f^*(Y^*)$ are distributed identically to $f(X), f(Y)$, respectively, since $\Pr[f^* = f \wedge X_{j_h}^* = 0] = \Pr[X_{j_h} = 0]$ and $\Pr[f^* = f \wedge Y_{j_h}^* = 0] = \Pr[Y_{j_h} = 0]$. (Note that $X_{j_h} = 0$ exactly when the non-identity transition is used at layer j_h in the computation of f on X , and that the event $f^* = f \wedge X_{j_h}^* = 0$ occurs exactly when the non-identity transition is used at layer j_h in the computation of f^* on X^* .)

In the example above, in which $\epsilon = 0.1$, this means choosing γ as large as possible such that

$$\frac{1 - \gamma \frac{1}{n}}{1 - \gamma \frac{0.9}{n}} \geq 1 - \epsilon = 0.9.$$

A short computation shows the maximum value of γ is $\gamma = n/1.9$. Thus, in this case,

$$\Pr[X_{j_h}^* = 0] = \frac{1}{1.9}, \quad \Pr[Y_{j_h}^* = 0] = \frac{0.9}{1.9}.$$

Note the difference with (5): both probabilities have moved away from 0, and are now close to $\frac{1}{2}$.

In the proof, the above operation consisting of randomizing the program at a transition (to be the original program w.p. $\frac{1}{\gamma}$, or to be the identity w.p. $1 - \frac{1}{\gamma}$) and of simultaneously boosting by a factor γ in each distribution the probability of the input value giving a collision at that layer, is carried out for all layers of the program at once, with the value of γ individually computed for each layer. The resulting randomized program f^* is defined by choosing each layer independently to be either the identity or the original layer, with respect to the relevant probabilities. Since $f(X), f(Y)$ are distributed identically to $f^*(X^*), f^*(Y^*)$, with randomness taken also over the choice of f^* , we have that

$$\Delta(f(X), f(Y)) = \Delta(f^*(X^*), f^*(Y^*)) \leq \mathbb{E}_{f^*} \Delta(f^*(X^*), f^*(Y^*)).$$

In the righthmost expression, the statistical distance is computed solely over the randomness induced by X^* and Y^* , for a fixed value of f^* . Because of the probability boosting, one can show that if the b -transition at the k -th layer of f^* has a collision while the $(1-b)$ -transition does not (note this implies the same statement holds in f), then

$$\min(\Pr[X_k^* = b], \Pr[Y_k^* = b]) \geq \frac{1 - \epsilon}{2 - \epsilon}.$$

Since the latter probability is near $\frac{1}{2}$, the same hybrid method used for biased coins can be used to upper bound $\Delta(f^*(X^*), f^*(Y^*))$ for any fixed value of f^* .

We note that the (central) idea of obtaining width reduction of the program via collisions originates in [7]. There, random restrictions are used to obtain collisions and width-reduction. Our paper swaps random restrictions for a hybrid argument, which has the advantage that one can control the position of the restricted bits (these being, in the hybrid argument, the bits common to two neighboring hybrids). Having long intervals of consecutive restricted bits augments the chance of obtaining at least one collision in each of these intervals, and thus improves the chance of obtaining width-reduction. (On the other hand, longer intervals means more hybrids, implying a tradeoff.)

As another point of comparison, we note that [7] eschews program randomization in favor of a (lossy) reduction from the problem of distinguishing “well-behaved” input distributions (with probabilities near $\frac{1}{2}$) to the problem of distinguishing “troublesome” input distributions (with probabilities near 0). It is partly this reduction which causes the alphabet size $|\Sigma|$ to appear in the final bound of [7] (whereas our bounds are independent of $|\Sigma|$).

5 Some further proof details: width two branching programs and the collision lemma

As explained, the proof of Theorem 1 relies on an inductive argument whose base case is an upper bound on the distinguishing power of width 2 branching programs. Intuitively, a ROBP of width 2 (say, when distinguishing a $\pm\epsilon$ -biased coin) cannot do better than to determine acceptance based on the outcome of a single coin flip, given its limited memory—e.g., by ignoring all coin flips except for the last. Note that such a ROBP has advantage $(\frac{1}{2} + \epsilon) - (\frac{1}{2} - \epsilon) = 2\epsilon$, the statistical distance in a single coin flip.

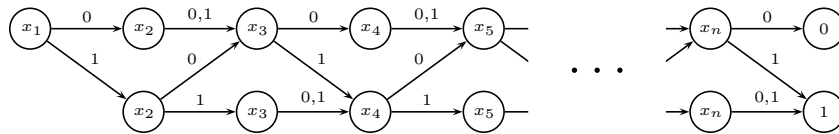


Fig. 2: Width 2 branching program obtaining better than 2ϵ advantage at distinguishing a $(\frac{1}{2} + \epsilon, \frac{1}{2} - \epsilon)$ -biased coin from an $(\frac{1}{2} - \epsilon, \frac{1}{2} + \epsilon)$ -biased coin (n odd, $n \geq 3$).

However this intuition is incorrect. Indeed, a width 2 ROBP can distinguish a $(\pm\epsilon)$ -biased coin with advantage approaching

$$2\epsilon \left(\frac{3}{4} + \epsilon^2 \right)^{-1} \quad (8)$$

as the length n of the program goes to infinity, which is close to $\frac{4}{3}$ as large as 2ϵ for small ϵ . The program whose distinguishing advantage approaches this value is shown in Fig. 2. The program of Fig. 2 is, convec-

turally, the best width 2 distinguisher of length $n = 2k + 1$ for $(\pm\epsilon)$ -biased coins, but we do not have a proof. (Also conjecturally, width two ROBPs of length $n = 2k + 2$ do no better at distinguishing $(\pm\epsilon)$ -biased coins than length $2k + 1$ ROBPs.) Our own bound shows that width 2 ROBPs cannot distinguish $(\pm\epsilon)$ -biased coins with advantage better than

$$1 - \frac{\frac{1}{2} - \epsilon}{\frac{1}{2} + \epsilon} = 2\epsilon \left(\frac{1}{2} + \epsilon \right)^{-1} \quad (9)$$

regardless of their length. For small ϵ , this is roughly 1.5 times as large as the conjectured optimal advantage (8), but this constant-factor discrepancy is unimportant for our final bound.

The general theorem which we prove on width two branching programs is the following:

Theorem 3. *Let f be a width 2 ROBP of length n and let $X = (X_i)_{i=1}^n \in \Sigma^n$, $Y = (Y_i)_{i=1}^n \in \Sigma^n$ be two ϵ -close in ratio product distributions. Then $\Delta(f(X), f(Y)) \leq \epsilon$.*

(We note that (9) is the direct application of Theorem 3.) The proof of Theorem 3 is in fact nontrivial and uses many ideas from the inductive proof described in Section 4, including (a strengthened version of) Brody and Verbin’s collision lemma, program randomization, and a coupling argument. It would be interesting to know if an “easy” proof exists.

To state the collision lemma, which plays a key role in our results, we start by giving the formal definition of cROBPs.

Definition 1. *A width w read-once branching program f is called a collision read-once branching program (cROBP) if every transition function τ_α of f is either the identity from $[w]$ to $[w]$ or else is not injective (i.e. is not a permutation).*

Collision Lemma. (After [7].) Let $X, Y \in \Sigma^n$ be product distributions and let $w \geq 1$. Then there exists a cROBP f of length n and width w whose distinguishing advantage $\Delta(f(X), f(Y))$ is at least as great as the distinguishing advantage $\Delta(g(X), g(Y))$ of any length n width w ROBP g .

The collision lemma found in [7] states that the optimal distinguishing advantage can be achieved by a program f with the following property: at every layer of f , either all the transition functions are the identity, or else at least one of the transition functions contains a collision. This is weaker than our lemma, which implies that *all* non-identity transition functions contain a collision (i.e. are not permutations). While our version may seem much stronger at first glance, we comment that its proof only requires a minor modification of the proof of [7].

The proof of Theorem 3 is in Appendix A and the proof of the collision lemma is in Appendix C.

Acknowledgments. I would like to thank Kevin Matulef, Joshua Brody and Elad Verbin for helpful conversations at all stages of this work.

References

1. Scott Aaronson, Andrew Drucker. Advice coins for classical quantum computation. *Proc. of the 38th International Colloquium on Automata, Languages and Programming*, 2011.
2. Kazuyuki Amano. Bounds on the size of small depth circuits for approximating majority. In *Proc. of the 36th International Colloquium on Automata, Languages and Programming*, 2009.

3. Anindya De. Improved pseudorandomness for regular branching programs. In *Conference on Computational Complexity*, 2011.
4. Boaz Barak, Ishay Haviv, Moritz Hardt, Anup Rao, Oded Regev, and David Steurer. Rounding parallel repetitions of unique games. In *Proc. of the 49th Annual ACM Symposium on the Foundations of Computer Science*, 2008.
5. Andrej Bogdanov, Zeev Dvir, Elad Verbin, Amir Yehudahoff. Pseudorandom generators for width two branching programs. *ECCC*, 2009.
6. Mark Braverman, Anup Rao, Ran Raz, Amir Yehudahoff. Pseudorandom generators for regular branching programs. In *Proc. of the 51st Annual ACM Symposium on the Foundations of Computer Science*, 2010.
7. Joshua Brody, Elad Verbin. The coin problem and pseudorandomness for branching programs. In *Proc. of the 51st Annual ACM Symposium on the Foundations of Computer Science*, 2010.
8. Bill Fefferman, Ronen Shaltiel, Christopher Umans, Emanuele Viola. On beating the hybrid argument. *ITCS* 2012.
9. Johan Håstad. A slight sharpening of LMN. *Journal of Computer and System Sciences*, **63**, 2001.
10. Martin Hellman, Thomas Cover. Learning with finite memory. *Ann. of Math. Stat.*, **41**, 1970.
11. Michal Koucký, Prajata Nimbhorkar, Pavel Pudlak. Pseudorandom generators for group products. In *Proc. of the 43rd Annual ACM Symposium on the Theory of Computing*, 2011.
12. Anup Rao, David Zuckerman. Pseudorandom generators for polynomial threshold functions. In *Proc. of the 51st Annual ACM Symposium on the Foundations of Computer Science*, 2010.
13. Jiří Šíma, Stanislav Žák. Almost k -wise independent sets establish hitting sets for width-3, 1-branching programs. *Computer Science Theory And Applications: 6th International Computer Science Symposium in Russia*, St. Petersburg, 2011.
14. Emanuele Viola. Randomness buys depth for approximate counting. *Electronic Colloquium on Computational Complexity (ECCC)*, 17:175, 2010.

A The base case: proof of Theorem 3

In this section we prove Theorem 3 concerning the power of width two branching programs. This theorem serves as the “base case” for the proof of Theorem 1.

Let U, V be random variables taking values in some set Σ and defined over a common probability space. We say U, V are *maximally coupled* if for every $\alpha \in \Sigma$,

$$\begin{aligned} (\Pr[U = \alpha] \leq \Pr[V = \alpha]) &\implies (U = \alpha \implies V = \alpha), \\ (\Pr[V = \alpha] \leq \Pr[U = \alpha]) &\implies (V = \alpha \implies U = \alpha). \end{aligned}$$

It is an easy fact that for any two random variables U, V there exist random variables $\tilde{U} \sim U, \tilde{V} \sim V$ such that \tilde{U}, \tilde{V} are maximally coupled. Moreover, $\Delta(U, V) = \Delta(\tilde{U}, \tilde{V}) = \Pr[\tilde{U} \neq \tilde{V}]$ when \tilde{U}, \tilde{V} are maximally coupled copies of U, V .

Proof of Theorem 3. By the Collision Lemma, we can assume without loss of generality that f is a cROBP. We can also assume that $n \geq 1$ is chosen to be the smallest value for which the Theorem does not hold.

For each $i, 1 \leq i \leq n$, let $\mathcal{C}_i \subseteq \Sigma$ be the set of values α such that the transition function τ_α at layer i of f has a collision, and let $\mathcal{I}_i = \Sigma \setminus \mathcal{C}_i$ be the set of values α such that the transition function τ_α at layer i is the identity function. Let

$$\sigma_i = \max(\Pr[X_i \in \mathcal{C}_i], \Pr[Y_i \in \mathcal{C}_i]).$$

By the minimality of n we can assume $\sigma_i \neq 0$ for all i (otherwise we could remove level i and the i -th input of the program without affecting the program’s output, and obtain a smaller counterexample). Note $\sigma_i \neq 0$ implies $\Pr[X_i \in \mathcal{C}_i] \neq 0, \Pr[Y_i \in \mathcal{C}_i] \neq 0$ by the definition of ϵ -close in ratio. We define new random variables \bar{X}_i, \bar{Y}_i such that

$$\begin{aligned} \Pr[\bar{X}_i = \alpha] &= \Pr[X_i = \alpha] / \sigma_i \\ \Pr[\bar{Y}_i = \alpha] &= \Pr[Y_i = \alpha] / \sigma_i \end{aligned}$$

for all $\alpha \in \mathcal{C}_i$, and such that if $\mathcal{I}_i \neq \emptyset$, then there exists a $\beta \in \mathcal{I}_i$ such that

$$\begin{aligned}\Pr[\overline{X}_i = \beta] &= 1 - \Pr[X_i \in \mathcal{C}_i]/\sigma_i \\ \Pr[\overline{Y}_i = \beta] &= 1 - \Pr[Y_i \in \mathcal{C}_i]/\sigma_i.\end{aligned}$$

We assume moreover that $\overline{X}_i, \overline{Y}_i$ are maximally coupled. We note that $\overline{X}_i, \overline{Y}_i$ are *not* constructed analogously to $X_{j_h}^*, Y_{j_h}^*$, as defined in (6), (7). For example, $\overline{X}_i, \overline{Y}_i$ are in general not ϵ -close in ratio.

Since either $\Pr[\overline{X}_i \in \mathcal{C}_i] = 1$ or $\Pr[\overline{Y}_i \in \mathcal{C}_i] = 1$, the coupling implies that $\Pr[\overline{X}_i = \overline{Y}_i \in \mathcal{C}_i] \geq 1 - \epsilon$. Indeed, say wlog that $\Pr[\overline{X}_i \in \mathcal{C}_i] = 1$; then

$$\begin{aligned}\Pr[\overline{X}_i = \overline{Y}_i \in \mathcal{C}_i] &= \Pr[\overline{X}_i = \overline{Y}_i] \\ &= \sum_{\alpha \in \Sigma} \Pr[\overline{X}_i = \alpha] \Pr[\overline{Y}_i = \alpha | \overline{X}_i = \alpha] \\ &\geq \sum_{\alpha \in \Sigma} \Pr[\overline{X}_i = \alpha] (1 - \epsilon) = 1 - \epsilon.\end{aligned}$$

Say that a layer of a branching program of width w is *the identity* if each of its transition functions $\tau_\alpha : [w] \rightarrow [w]$ is the identity. Let f^* be a randomized branching program of width two and length n , whose i -th layer is the identity with probability $1 - \sigma_i$, and equals the i -th layer of f with probability σ_i , with all layers determined independently. It is clear that

$$\begin{aligned}\Pr[f(X) = 1] &= \Pr[f^*(\overline{X}) = 1] \\ \Pr[f(Y) = 1] &= \Pr[f^*(\overline{Y}) = 1]\end{aligned}$$

where the probabilities are taken over f^* as well as over \overline{X} and \overline{Y} . (To see this it can be helpful to view the randomized f^* as a “filter” which modifies $\overline{X}_i, \overline{Y}_i$ and then passes on the modified values to f . In this view \overline{X}_i and \overline{Y}_i are left unchanged with probability σ_i and are reassigned some value in \mathcal{I}_i with probability $1 - \sigma_i$.) In particular, $\Delta(f(X), f(Y)) = \Delta(f^*(\overline{X}), f^*(\overline{Y}))$. We have

$$\Delta(f^*(\overline{X}), f^*(\overline{Y})) \leq \mathbb{E}_{f^*} \Delta(f^*(\overline{X}), f^*(\overline{Y})) \quad (10)$$

where the statistical distances on the right-hand side are computed for a fixed value of f^* , over the randomness of $\overline{X}, \overline{Y}$. (More generally, $\Delta(X, Y) \leq \mathbb{E}_Z \Delta((X|Z), (Y|Z)) = \sum_z \Pr[Z = z] \cdot \Delta((X|Z = z), (Y|Z = z))$ for any random variables X, Y, Z .)

Let g be some fixed value of f^* . We upper bound $\Delta(g(\overline{X}), g(\overline{Y}))$. More precisely, by (10), it is sufficient to show $\Delta(g(\overline{X}), g(\overline{Y})) \leq \epsilon$ to conclude the proof.

Say that a layer of a branching program “is the identity” if all the transition functions of that layer are the identity. If all layers of g are the identity then g is oblivious to its input and $\Delta(g(\overline{X}), g(\overline{Y})) = 0$. Otherwise let $m, 1 \leq m \leq n$, be the last non-identity layer of g . Note that layer m of g equals layer m of f and that $g(\overline{X}) = g(\overline{Y})$ if $\overline{X}_m = \overline{Y}_m \in \mathcal{C}_m$. (This uses the fact that g has width two.) Thus

$$\Pr[g(\overline{X}) = g(\overline{Y})] \geq \Pr[\overline{X}_m = \overline{Y}_m \in \mathcal{C}_m] \geq 1 - \epsilon$$

and so $\Delta(g(\overline{X}), g(\overline{Y})) \leq \Pr[g(\overline{X}) \neq g(\overline{Y})] \leq \epsilon$ (where the fact that $\Delta(g(\overline{X}), g(\overline{Y})) \leq \Pr[g(\overline{X}) \neq g(\overline{Y})]$ is a generic property of statistical distance). This concludes the proof. \square

B The main result: proof of Theorem 1

Theorem 1 is derived as a corollary of the following more general result (stated this way because it conveniently allows for induction), which is really the paper's main result. Here r corresponds to the number hybrids $c \log n$, in the notation of Section 4:

Theorem 4. *Let f be a width $w \geq 2$ ROBP of length n and let $X = (X_i)_{i=1}^n \in \Sigma^n$, $Y = (Y_i)_{i=1}^n \in \Sigma^n$ be two ϵ -close in ratio product distributions. Then for any integer $r \geq 1$,*

$$\Delta(f(X), f(Y)) \leq \epsilon r^{w-2} + \left(n(w-2) + \sum_{j=3}^{w-1} (w-j)r^{j-2} \right) \left(\frac{1}{2-\epsilon} \right)^{r-1}. \quad (11)$$

Proof. We use lexicographic induction on the pair (w, n) . That is, assuming the theorem is true for all ROBPs of width w' and length n' such that either $w' < w$ or such that $w' = w$ and $n' < n$, we prove the theorem holds for all ROBPs f of length n and width w . Note we can assume $w \geq 3$ since the case $w = 2$ follows from Theorem 3. We can also assume f is a cROBP by the Collision Lemma.

Let $\mathcal{C}_i \subseteq \Sigma$ be the set of α 's for which the α -transition function τ_α of f at layer i is not the identity, and let $\mathcal{I}_i = \Sigma \setminus \mathcal{C}_i$, $1 \leq i \leq n$. By induction on n we can assume that $\Pr[X_i \in \mathcal{C}_i], \Pr[Y_i \in \mathcal{C}_i] > 0$ for all i , since otherwise $\Pr[X_i \in \mathcal{C}_i] = \Pr[Y_i \in \mathcal{C}_i] = 0$ for some i and the i -th layer of f can simply be removed. (Note that (11) is monotonic in n , for fixed w .)

We define new product distributions $X^* = (X_i^*)_{i=1}^n$, $Y^* = (Y_i^*)_{i=1}^n$. If $\mathcal{I}_i = \emptyset$, then X_i^*, Y_i^* are distributed identically to X_i, Y_i . Otherwise let β_i be any element of $\mathcal{I}_i \neq \emptyset$, and define the distributions X^*, Y^* by

$$X_i^* = \begin{cases} \alpha & \text{with probability } \gamma_i \Pr[X_i = \alpha] \text{ if } \alpha \in \mathcal{C}_i, \text{ and} \\ \beta_i & \text{with probability } 1 - \gamma_i \Pr[X_i \in \mathcal{C}_i] \end{cases}$$

$$Y_i^* = \begin{cases} \alpha & \text{with probability } \gamma_i \Pr[Y_i = \alpha] \text{ if } \alpha \in \mathcal{C}_i, \text{ and} \\ \beta_i & \text{with probability } 1 - \gamma_i \Pr[Y_i \in \mathcal{C}_i] \end{cases}$$

where $\gamma_i \geq 1$ is the largest number such that X_i^*, Y_i^* are ϵ -close in ratio. One can easily check that, for all $1 \leq i \leq n$,

$$\min(\Pr[X_i^* \in \mathcal{C}_i], \Pr[Y_i^* \in \mathcal{C}_i]) \geq \frac{1-\epsilon}{2-\epsilon}. \quad (12)$$

(Indeed if (say) $\Pr[X_i^* \in \mathcal{C}_i] \geq \Pr[Y_i^* \in \mathcal{C}_i]$ and $\Pr[Y_i^* \in \mathcal{C}_i] < \frac{1-\epsilon}{2-\epsilon}$ then

$$\frac{\Pr[X_i^* \in \mathcal{I}_i]}{\Pr[Y_i^* \in \mathcal{I}_i]} = \frac{1 - \Pr[X_i^* \in \mathcal{C}_i]}{1 - \Pr[Y_i^* \in \mathcal{C}_i]} \geq \frac{1 - \frac{1}{1-\epsilon} \Pr[Y_i^* \in \mathcal{C}_i]}{1 - \Pr[Y_i^* \in \mathcal{C}_i]} > \frac{1 - \frac{1}{2-\epsilon}}{1 - \frac{1-\epsilon}{2-\epsilon}} = 1 - \epsilon$$

where the strict inequality uses the monotonicity of the function $\frac{1-ct}{1-t}$, $c < 1$, on $0 < t < 1$. This contradicts the maximality of γ_i .)

We let f^* be a randomized width w length n ROBP whose i -th layer equals the i -th layer of f with probability $1/\gamma_i$ (putting $\gamma_i = 1$ if $\mathcal{I}_i = \emptyset$) and equals the identity layer with probability $1 - 1/\gamma_i$, with every layer determined independently. One can easily check that $f(X) \sim f^*(X^*)$, $f(Y) \sim f^*(Y^*)$ (with randomness taken over f^*, X, Y), so that

$$\Delta(f(X), f(Y)) = \Delta(f^*(X^*), f^*(Y^*)) \leq \mathbb{E}_{f^*} \Delta(f^*(X^*), f^*(Y^*)) \quad (13)$$

where the statistical distance in the rightmost expression is measured over the randomness induced by X^* , Y^* , for fixed value f^* .

Let g be some fixed value of f^* . By (13), it suffices to show

$$\Delta(g(X^*), g(Y^*)) \leq \epsilon r^{w-2} + \left(n(w-2) + \sum_{j=3}^{w-1} (w-j)r^{j-2} \right) \left(\frac{1}{2-\epsilon} \right)^{r-1} \quad (14)$$

to conclude the proof.

Say g has an identity layer. If g 's length is 1 (i.e. $n = 1$) then all layers of g are the identity and $\Delta(g(X^*), g(Y^*)) = 0$, so (14) holds. Otherwise, we can cut out this identity layer from g (removing at the same time the relevant coordinate in X^* , Y^*), and (14) follows by induction on the length n , since X^* , Y^* are ϵ -close in ratio product distributions. Thus, we can assume none of the layers of g are the identity. This implies $g = f^* = f$, so all that remains is to upper bound $\Delta(f(X^*), f(Y^*))$.

Let $Z_0, \dots, Z_r \in \Sigma^n$ be hybrid distributions defined by

$$Z_{i,j} = \begin{cases} Y_j^* & \text{if } (j \bmod r) < i, \\ X_j^* & \text{otherwise,} \end{cases}$$

where $Z_{i,j}$ is the j -th bit of Z_i . We have $Z_0 = X^*$ and $Z_r = Y^*$, so $\Delta(f(X^*), f(Y^*)) \leq \sum_{i=0}^{r-1} \Delta(f(Z_i), f(Z_{i+1}))$. Fixing i , $0 \leq i \leq r-1$, we now focus on upper bounding $\Delta(f(Z_i), f(Z_{i+1}))$.

Let $\mathcal{Z}_i = \{j \in [n] : j \bmod r \equiv i\}$ be the set of bit positions at which (the definitions of) Z_i and Z_{i+1} differ. Let \bar{Z} and $f_{\bar{Z}}$ be defined as in Section 4. Let $X' \in \Sigma^{\mathcal{Z}_i}$ be the restriction of X^* to the bits in \mathcal{Z}_i , and likewise let $Y' \in \Sigma^{\mathcal{Z}_i}$ be the restriction of Y^* to the bits in \mathcal{Z}_i . Then $f(Z_i) = f_{\bar{Z}}(X')$ and $f(Z_{i+1}) = f_{\bar{Z}}(Y')$. We therefore have

$$\Delta(f(Z_i), f(Z_{i+1})) = \Delta(f_{\bar{Z}}(X'), f_{\bar{Z}}(Y')) \leq \mathbb{E}_{\bar{Z}} \Delta(f_{\bar{Z}}(X'), f_{\bar{Z}}(Y')). \quad (15)$$

Let $\ell = |\mathcal{Z}_i|$ be the length of $f_{\bar{Z}}$. Note that $\ell \leq \lceil \frac{n}{r} \rceil$. Fix $h \in [\ell-1]$ and consider the probability, taken over \bar{Z} , that the union $\{\tau_\alpha(j) : j \in [w], \alpha \in \Sigma\} \subseteq [w]$ of the ranges of the transition functions $\{\tau_\alpha : \alpha \in \Sigma\}$ of $f_{\bar{Z}}$ at layer h has size w . By (12), one can argue via a pebble argument as in Section 4 that this probability is at most $(1 - \frac{1-\epsilon}{2-\epsilon})^{r-1} = (\frac{1}{2-\epsilon})^{r-1}$. Union bounding over $h \in [\ell-1]$ (the transition functions at layer ℓ having range $\{1, 2\}$), the probability that $f_{\bar{Z}}$ is not equivalent to a width $w-1$ branching program of length ℓ is at most $\frac{n}{r} (\frac{1-\epsilon}{2-\epsilon})^{r-1}$, using $\ell-1 \leq \frac{n}{r}$. Thus, using the inductive hypothesis for programs of width $w-1$ and length⁷ $\ell \leq \frac{n}{r} + 1$, $\mathbb{E}_{\bar{Z}} \Delta(f_{\bar{Z}}(X'), f_{\bar{Z}}(Y'))$ can be upper bounded by

$$\epsilon r^{w-3} + \left(\left(\frac{n}{r} + 1 \right) (w-3) + \sum_{j=3}^{w-2} (w-1-j)r^{j-2} \right) \left(\frac{1}{2-\epsilon} \right)^{r-1} + \frac{n}{r} \left(\frac{1}{2-\epsilon} \right)^{r-1} \quad (16)$$

which simplifies to

$$\epsilon r^{w-3} + \left(\frac{n}{r} (w-2) + \sum_{j=3}^{w-1} (w-j)r^{j-3} \right) \left(\frac{1-\epsilon}{2-\epsilon} \right)^{r-1}.$$

⁷ Technically, note that $\ell \leq \lceil \frac{n}{r} \rceil$ could be $\geq n$ if $n = 1$. However, we are using lexicographic induction over (w, n) , so we can use the inductive hypothesis on any ROBP of width $w-1$ (even of length greater than n).

This is, by (15), an upper bound for $\Delta(f(Z_i), f(Z_{i+1}))$. Since this upper bound holds for any $0 \leq i \leq r-1$, we have

$$\begin{aligned} \Delta(f(X^*), f(Y^*)) &\leq \sum_{i=0}^{r-1} \Delta(f(Z_i), f(Z_{i+1})) \\ &\leq \epsilon r^{w-2} + \left(n(w-2) + \sum_{j=3}^{w-1} (w-j)r^{j-2} \right) \left(\frac{1}{2-\epsilon} \right)^{r-1} \end{aligned}$$

which concludes the proof. \square

Corollary 1. *Let f be a width $w \geq 2$ ROBP of length n and let $X, Y \in \Sigma^n$ be two ϵ -close in ratio product distributions. Then for any integer $r \geq 1$,*

$$\Delta(f(X), f(Y)) \leq \epsilon r^{w-2} + \left(n + r^{w-2} \right) (w-2) \left(\frac{1}{2-\epsilon} \right)^{r-1}.$$

Proof of Theorem 1. Since the case $w = 2$ follows already from Theorem 3, assume $3 \leq w \leq \log n / \log \log n$. Set $r = \lceil 2 \log n \rceil$. We have

$$r^{w-2} \leq (2 \log n + 1)^{w-2} = \left(2 \log n \left(1 + \frac{1}{2 \log n} \right) \right)^{w-2} \leq (2 \log n)^{w-2} e^{\frac{w-2}{2 \log n}}.$$

Since $w \leq \log n / \log \log n$, we have

$$e^{\frac{w-2}{2 \log n}} \leq 1 + \lambda_1(n)$$

where $\lambda_1(n) = o(1)$, and

$$(2 \log n)^{w-2} \leq n^{1+\lambda_2(n)}$$

where $\lambda_2(n) = o(1)$. Thus, by Corollary 1,

$$\Delta(f(X), f(Y)) \leq \epsilon (2 \log n)^{w-2} (1 + \lambda_1(n)) + (n + n^{1+\lambda_2(n)}) (\log n) \left(\frac{1}{2-\epsilon} \right)^{r-1}. \quad (17)$$

We define

$$\lambda(n) = \lambda_1(n) + (n + n^{1+\lambda_2(n)}) (\log n) \left(\frac{1}{2-0.5} \right)^{r-1} \frac{1}{\epsilon (2 \log n)^{w-2}}.$$

Since $r = \lceil 2 \log n \rceil$, it is easy to check that $\lambda(n) = o(1)$. Next

$$\Delta(f(X), f(Y)) \leq \epsilon (2 \log n)^{w-2} (1 + \lambda(n)) \quad (18)$$

follows from (17), since (18) holds trivially if $\epsilon (2 \log n)^{w-2} \geq 1$ if $\epsilon \geq 0.5$. \square

C Proof of the Collision Lemma

In this appendix we give a proof of the Collision Lemma stated in Section 5. The proof follows [7] rather closely.

Let $X, Y \in \Sigma^n$ be two product distributions. We claim that for every ROBP g of length n and width w over the alphabet Σ there exists a cROBP f of length n and width w , also over Σ , such that $\Delta(f(X), f(Y)) \geq \Delta(g(X), g(Y))$ (cf. Definition 1).

It suffices to consider a ROBP g (of width w) maximizing $\Delta(g(X), g(Y))$; such a ROBP exists since there are only finitely many ROBPs possible for a given width, length, and input alphabet. Say that a node of a ROBP is *unreachable* if there is no path of nonzero probability from the start node to this node. (“Nonzero probability” is measured w.r.t. the distributions X and Y ; note that since $\Pr[X_i = \alpha] = 0 \iff \Pr[Y_i = \alpha] = 0$ for all $i \in [n]$, $\alpha \in \Sigma$, unreachability is well-defined.) We can moreover assume that among all ROBPs achieving the optimal distinguishing advantage, g has as many unreachable nodes as possible. We can finally assume that, among such ROBPs (with said distinguishing advantage and said number of reachable nodes), g maximizes $\Pr[g(X) = 1]$.

For every vertex v of g , let $v(X)$ be the output node reached when g is “started” at vertex v with input X . More precisely, if v is a vertex at layer $k \in [n + 1]$, $v(X)$ denotes the output reached by reading the $n - k + 1$ last characters of X starting from node v , namely by following the edges labeled X_k, \dots, X_n starting from node v . (If v is an output node then $v(X) = v$ for all X .) Define $v(Y)$ likewise.

For every node v of g , let $\mathbb{E}[v(X)]$ denote the probability, taken over X , that $v(X)$ is the accept node of g . As a piece of special-purpose notation, we also define $\mathbb{E}[\bar{v}(X)] = 1 - \mathbb{E}[v(X)]$ to be the probability that $v(X)$ is the reject node of g . Also let $\mathbb{E}[g(X)]$ and $\mathbb{E}[\bar{g}(X)]$ be the probabilities that g accepts and rejects X , respectively (this agrees with our previous notation if we identify g with its start node). We can assume without loss of generality that $\mathbb{E}[g(X)] \geq \mathbb{E}[g(Y)]$. Note, then, that

$$\Delta(g(X), g(Y)) = \mathbb{E}[g(X)] - \mathbb{E}[g(Y)] = \mathbb{E}[g(X)] + \mathbb{E}[\bar{g}(Y)] - 1.$$

Thus g maximizes $\mathbb{E}[g(X)] + \mathbb{E}[\bar{g}(Y)]$ (among all ROBPs of length n and width w).

Finally, let $p_X(v)$ be the probability of passing through node v when g is run on input X (starting from the start node), and define $p_Y(v)$ likewise. Since X, Y are ϵ -close in ratio, we have that

$$(v \text{ is unreachable}) \iff (p_X(v) = 0 \vee p_Y(v) = 0) \iff (p_X(v) = p_Y(v) = 0).$$

We also note that if v_1, \dots, v_w are the nodes at a given layer i of f , then

$$\mathbb{E}[g(X)] = \sum_{j=1}^w p_X(v_j) \mathbb{E}[v_j(X)], \quad \mathbb{E}[\bar{g}(Y)] = \sum_{j=1}^w p_Y(v_j) \mathbb{E}[\bar{v}_j(Y)]$$

so that g ’s advantage is

$$(-1) + \sum_{j=1}^w \left(p_X(v_j) \mathbb{E}[v_j(X)] + p_Y(v_j) \mathbb{E}[\bar{v}_j(Y)] \right). \quad (19)$$

Say the k -th layer of g has at least two reachable nodes, v_1 and v_2 . We show that either

$$\mathbb{E}[v_1(X)] > \mathbb{E}[v_2(X)] \quad \text{and} \quad \mathbb{E}[\bar{v}_1(Y)] < \mathbb{E}[\bar{v}_2(Y)],$$

or

$$\mathbb{E}[v_1(X)] < \mathbb{E}[v_2(X)] \quad \text{and} \quad \mathbb{E}[\bar{v}_1(Y)] > \mathbb{E}[\bar{v}_2(Y)].$$

Indeed, assume by contradiction, say, that $\mathbb{E}[v_1(X)] \geq \mathbb{E}[v_2(X)]$, $\mathbb{E}[\bar{v}_1(Y)] \geq \mathbb{E}[\bar{v}_2(Y)]$ (the only other possible bad case being $\mathbb{E}[v_1(X)] \leq \mathbb{E}[v_2(X)]$, $\mathbb{E}[\bar{v}_1(Y)] \leq \mathbb{E}[\bar{v}_2(Y)]$). Then by re-routing all of v_2 ’s incoming edges to v_1 , we decrease $p_X(v_2)$ and $p_Y(v_2)$ to 0 and increase $p_X(v_1)$, $p_Y(v_1)$ by (the old values of) $p_X(v_2)$, $p_Y(v_2)$. This cannot decrease (19), and the new ROBP thus obtained has at least one more

unreachable node, a contradiction. The case $\mathbb{E}[v_1(X)] \leq \mathbb{E}[v_2(X)]$, $\mathbb{E}[\overline{v_1}(Y)] \leq \mathbb{E}[\overline{v_2}(Y)]$ is similarly treated by re-routing all of v_1 's incoming edges to v_2 , which proves the claim.

Note the above observation allows us to re-order the reachable vertices v_1, \dots, v_j at any layer such that $\mathbb{E}[v_1(X)] > \dots > \mathbb{E}[v_j(X)]$ and $\mathbb{E}[\overline{v_1}(Y)] < \dots < \mathbb{E}[\overline{v_j}(Y)]$. While such an ordering is not necessary for the proof, one can assume it for conceptual simplicity.

Let v be a reachable node of g at layer $i \leq n$ of g , and let $\alpha, \beta \in \Sigma$ be such that $\Pr[X_i = \alpha], \Pr[X_i = \beta]$ are nonzero (hence, ditto for Y_i), and such that

$$\frac{\Pr[X_i = \alpha]}{\Pr[Y_i = \alpha]} \geq \frac{\Pr[X_i = \beta]}{\Pr[Y_i = \beta]}. \quad (20)$$

We setup the shorthands $\Pr[X_i = \alpha] = \alpha_X$, $\Pr[X_i = \beta] = \beta_X$, and so for α_Y, β_Y . Let v_α be the node reached from v by following the α edge, and let v_β be the node reached by following the β edge. We claim that either $v_\alpha = v_\beta$, or else that $\mathbb{E}[v_\alpha(X)] > \mathbb{E}[v_\beta(X)]$ (and hence, also, that $\mathbb{E}[\overline{v_\alpha}(Y)] < \mathbb{E}[\overline{v_\beta}(Y)]$). Indeed, assume by contradiction that $\mathbb{E}[v_\alpha(X)] < \mathbb{E}[v_\beta(X)]$. Then $\mathbb{E}[\overline{v_\alpha}(Y)] > \mathbb{E}[\overline{v_\beta}(Y)]$. If we re-route the α wire leaving v from v_α to v_β , then (19) (applied to a layer containing v) changes by

$$p_X(v)\alpha_X(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)]) + p_Y(v)\alpha_Y(\mathbb{E}[\overline{v_\beta}(Y)] - \mathbb{E}[\overline{v_\alpha}(Y)]) \quad (21)$$

whereas if we re-route instead the β wire leaving v from v_β to v_α , (19) changes by

$$p_X(v)\beta_X(\mathbb{E}[v_\alpha(X)] - \mathbb{E}[v_\beta(X)]) + p_Y(v)\beta_Y(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)]). \quad (22)$$

Since g is optimal, both (21), (22) must be ≤ 0 . Therefore,

$$\begin{aligned} p_X(v)\alpha_X(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)]) &\leq p_Y(v)\alpha_Y(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)]) \\ p_Y(v)\beta_Y(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)]) &\leq p_X(v)\beta_X(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)]) \end{aligned}$$

which we can rewrite as

$$\begin{aligned} \frac{\alpha_X}{\alpha_Y} &\leq \frac{p_Y(v)(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)])}{p_X(v)(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)])} \\ \frac{\beta_X}{\beta_Y} &\geq \frac{p_Y(v)(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)])}{p_X(v)(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)])}. \end{aligned}$$

This implies $\frac{\alpha_X}{\alpha_Y} \leq \frac{\beta_X}{\beta_Y}$, but we have by assumption (cf. (20)) that $\frac{\alpha_X}{\alpha_Y} \geq \frac{\beta_X}{\beta_Y}$, so $\frac{\alpha_X}{\alpha_Y} = \frac{\beta_X}{\beta_Y}$. Thus the inequalities

$$\begin{aligned} p_X(v)\alpha_X(\mathbb{E}[v_\beta(X)] - \mathbb{E}[v_\alpha(X)]) + p_Y(v)\alpha_Y(\mathbb{E}[\overline{v_\beta}(Y)] - \mathbb{E}[\overline{v_\alpha}(Y)]) &\leq 0 \\ p_X(v)\beta_X(\mathbb{E}[v_\alpha(X)] - \mathbb{E}[v_\beta(X)]) + p_Y(v)\beta_Y(\mathbb{E}[\overline{v_\alpha}(Y)] - \mathbb{E}[\overline{v_\beta}(Y)]) &\leq 0 \end{aligned}$$

(cf. (21), (22)) are negative scalar multiples of one another, and can only be simultaneously satisfied if both left-hand sides are 0, i.e. if both (21), (22) are zero. This implies we can re-route the α wire leaving v to v_β without affecting g 's advantage (and only possibly decreasing the number of reachable nodes); but since $\mathbb{E}[v_\beta(X)] > \mathbb{E}[v_\alpha(X)]$ by assumption, this change will increase $\mathbb{E}[v(X)]$ and $\mathbb{E}[g(X)]$, contradicting that our initial assumption that g maximizes $\Pr[g(X) = 1]$ among all ROBPs having its advantage and number of reachable nodes. This concludes the proof that either $v_\alpha = v_\beta$ or $\mathbb{E}[v_\alpha(X)] > \mathbb{E}[v_\beta(X)]$.

Say that an edge has *zero probability* if there is zero probability of traversing that edge (either because its source is unreachable, or because it is labeled by an element α that has probability 0 at that layer). We

now show that, after possibly re-routing the sinks of zero probability edges, and after a possible reordering of the nodes in each layer, g is a cROBP. Note that re-routing the destination of zero probability edges has no effect on any of the measures considered so far by the proof (these being the advantage of g , the number of unreachable nodes, $\Pr[g(x) = 1] = \mathbb{E}[g(X)]$, and more generally the values $\mathbb{E}[v(X)]$, $\mathbb{E}[v(Y)]$ for all vertices v of g).

The transformation we carry out on zero probability edges is more specifically as follows: for every $\alpha \in \Sigma$ and every layer of g containing a zero probability edge of label α , we re-route that edge such that the transition function τ_α at that layer has a collision. This can obviously be done as long as $w \geq 2$. Note, now, that if an unreachable node occurs at a layer, then all transition functions at that layer contain collisions after this modification. (In particular, the presence of a permutation transition function at a layer indicates that all nodes of that layer are reachable.)

We now claim that g is a cROBP up to a reordering of the nodes in each layer. For this, it suffices to show that at each layer of g , the transition functions that are permutations are all identical. Namely, we only need to show that if τ_α, τ_β are permutation transition functions occurring at the same layer i , then $\tau_\alpha = \tau_\beta$.

Let $\alpha_X = \Pr[X_i = \alpha]$, $\beta_X = \Pr[X_i = \beta]$ and so for α_Y, β_Y . Note $\alpha_X, \beta_X, \alpha_Y, \beta_Y > 0$ or else τ_α, τ_β would not be permutations. Assume wlog that $\alpha_X/\beta_X \geq \alpha_Y/\beta_Y$. Let v_1, \dots, v_w be the nodes at layer i and let $v_{j,\alpha}, v_{j,\beta}$ be the vertices reached by following the α and β edges leaving v_j , respectively. Then v_1, \dots, v_w are all reachable. From our observation above, since $\alpha_X/\beta_X \geq \alpha_Y/\beta_Y$, we have that $\mathbb{E}[v_{j,\alpha}(X)] \geq \mathbb{E}[v_{j,\beta}(X)]$ for all j . Moreover, $\mathbb{E}[v_{j,\alpha}(X)] > \mathbb{E}[v_{j,\beta}(X)]$ if $v_{j,\alpha} \neq v_{j,\beta}$. Thus if $v_{j,\alpha} \neq v_{j,\beta}$ for some j , then

$$\sum_{j=1}^w \mathbb{E}[v_{j,\alpha}(X)] > \sum_{j=1}^w \mathbb{E}[v_{j,\beta}(X)]$$

which is impossible since τ_α, τ_β are permutations. Therefore, we have $v_{j,\alpha} = v_{j,\beta}$ for all j , i.e. $\tau_\alpha = \tau_\beta$, which finishes the claim.

D The maximum total influence of monotone programs

In this appendix we use similar techniques to give tight upper bounds on the maximum influence of a monotone read-once branching program of width w . A ROBP $f : \{0, 1\}^n \rightarrow \{0, 1\}$ is *monotone* if

$$x \leq y \implies f(x) \leq f(y)$$

for every $x, y \in \{0, 1\}^n$, where we write $x \leq y$ if $x_i \leq y_i$ for $i = 1 \dots n$, x_i being the i -th bit of x .

The *influence* $\text{Inf}_i(f)$ of the i -th bit of f 's input is the probability

$$\Pr_x[f(x) \neq f(x + e_i)]$$

where $x + e_i$ denotes the bitwise xor of x with the i -th unit vector. This probability is taken over x uniformly chosen in $\{0, 1\}^n$. The *total influence* $\text{Inf}(f)$ is the sum of the individual influences of the bits:

$$\text{Inf}(f) = \sum_{i=1}^n \text{Inf}_i(f).$$

The parity (or xor) function exhibits the maximum possible total influence, being n . However, the xor functions is non-monotone. The total influence of any monotone function is upper bounded by \sqrt{n} (or by $(1 + o(1))\sqrt{2n/\pi}$ more exactly, which is matched by the majority function). We will see the maximum

total influence of a width w ROBP is $O(\log(n)^{w-2})$. In fact, we will prove that every monotone ROBP is equivalent to some cROBP of equal width, and that every width w cROBP has total influence $O(\log(n)^{w-2})$.

We first establish some notation and terminology. Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a boolean ROBP. A node v of f is *unreachable* if there is no path from the start node of f to v . Each node v of f induces a function on $n - k + 1$ bits in the natural way; we identify v with this function, and write $v : \{0, 1\}^{n-k+1} \rightarrow \{0, 1\}$. Two nodes u, v in the same layer of f are *equivalent* if they are equal as functions.

The following is the equivalent of the Collision Lemma for monotone RBPs (and is potentially of independent interest):

Lemma 1. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a monotone ROBP of length n and width w . Then there exists a cROBP g of length n and width w such that $g(x) = f(x)$ for all $x \in \{0, 1\}^n$.*

Proof. By re-routing edges we can assume no two reachable nodes of f are equivalent. We now order the nodes in each layer working inductively by layers, starting at the output layer and working backwards. We order the output layer such that the reject node (output 0) has index 1 and the accept node (output 1) has index 2.

Assuming layer $k + 1$ has been ordered we now describe how to order the reachable nodes in layer k , and later extend this ordering to include the unreachable nodes in layer k . For a node u in layer k , let $z(u)$ be the node reached by following the outgoing 0 edge from u , and let $o(u)$ be the node reached by following the 1 edge from u . We order the reachable nodes in layer k lexicographically according to the pair $(z(\cdot), o(\cdot))$, where coordinates are compared by the ordering of layer $k + 1$; that is, u comes before v if either $z(u) < z(v)$ or else if $z(u) = z(v)$ and $o(u) < o(v)$, where the relations $z(u) < z(v)$, $o(u) < o(v)$ refer to the ordering in layer $k + 1$. Since we have eliminated equivalent reachable nodes, it is easy to see this establishes a total ordering on the reachable nodes at layer k . Finally, we arbitrarily extend this ordering to include the unreachable nodes as well (say, by putting the unreachable nodes of each layer last, in some arbitrary order). This completes the description of how nodes are ordered in each layer.

For two reachable nodes u, v in the same layer, we write $u < v$ if u comes before v in the ordering just described. We now claim that if $u < v$, with u, v in layer k , then there is some $x \in \{0, 1\}^{n-k+1}$ such that $u(x) < v(x)$. We prove this by reverse induction on k . The base case is $k = n + 1$, which is obvious since then u, v are the two output nodes. Otherwise, for the induction step, just note that if $z(u) < z(v)$ then the claim follows by applying the induction hypothesis to the pair $z(u), z(v)$. So we can assume $z(u) = z(v)$. But then $o(u) < o(v)$, since $u < v$, so the claim follows by applying the induction hypothesis to the pair $o(u), o(v)$. This completes the proof of the claim.

It directly follows from this observation that for every non-output reachable node u , $z(u) \leq o(u)$ (or else f would not be monotone). We additionally modify the outgoing edges of unreachable nodes in f such that $z(v) \leq o(v)$ also for all unreachable nodes v . This does not affect f 's computation, and since we are not reordering nodes or changing the outgoing edges of reachable nodes, we still have $z(u) \leq o(u)$ for all reachable nodes u after this step.

It now suffices to show that if the transition functions τ_0, τ_1 at some layer k of f are both permutations, then $\tau_0 = \tau_1$. (The result then follows by a final permutation of each layer, working either from back-to-front or from front-to-back through the program.) That is, we need to show that $z(u) = o(u)$ for all nodes u in the k -th layer, assuming the two transition functions at the k -th layer are permutations. However, this is obvious from the fact that $z(u) \leq o(u)$ and that τ_0, τ_1 are permutations. (In a little more detail, if we associate $z(u)$ to the *index* of node $z(u)$ in layer $k + 1$ —so $z(u) \in [w]$ —and do the same for $o(u)$ —then if u_1, \dots, u_w are

the nodes in layer k and if $z(u_j) < o(u_j)$ for some j , we have

$$\sum_{i=1}^w o(u_i) - z(u_i) > 0$$

which contradicts the fact that τ_0, τ_1 are bijective.) \square

As for the distinguishability of product distributions, our results on influence require a preliminary result for width 2 branching programs.

Lemma 2. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a width 2 cROBP. Then $\text{Inf}(f) \leq 2$.*

Proof. We can assume without loss of generality that f has no identity layers. Let n be the length of f . For each i , $1 \leq i \leq n$, choose a value $b_i \in \{0, 1\}$ such that τ_{b_i} is not a permutation; such a b_i always exists by the fact that f is a cROBP. Then, trivially, $\text{Inf}_i(f) \leq 2^{i-n}$, since $f(x) = f(x + e_i)$ if $x_j = b_j$ for some $j > i$. Thus $\text{Inf}(f) \leq \sum_{i=1}^n 2^{i-n} \leq 2$. \square

We note that Lemma 2, like Theorem 3, is apparently not tight. The width 2 influence champion is again, conjecturally, the program of Fig. 2, whose total influence approaches $4/3$ as the length grows.

Theorem 5. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a cROBP of width $w \geq 2$ and length n . Then for every integer $r \geq 2$,*

$$\text{Inf}(f) \leq 2r^{w-2} + \frac{n^{1.5} r^{w-2} - 1}{2^{r-1} r - 1}.$$

Proof. The case of width $w = 2$ follows from Lemma 2.

Let $\mathcal{Z}_0, \dots, \mathcal{Z}_{r-1} \subseteq [n]$ be like in the proof of Theorem 4; namely,

$$\mathcal{Z}_i = \{j \in [n] : j \equiv i \pmod{r}\}.$$

Given a string $z_i \in \{0, 1\}^{[n] \setminus \mathcal{Z}_i}$, let f_{z_i} be the width w length $n - |\mathcal{Z}_i|$ ROBP induced by fixing the bits in $[n] \setminus \mathcal{Z}_i$ to z_i . Clearly, we have

$$\sum_{j \in \mathcal{Z}_i} \text{Inf}_j(f) = \mathbb{E}_{z_i}[\text{Inf}(f_{z_i})].$$

Since f is a cROBP, the probability f_{z_i} is not equivalent to a width $w - 1$ ROBP of same length is at most $\frac{n}{r} \cdot \frac{1}{2^{r-1}}$, following the same reasoning as in Section 4 and as in the proof of Theorem 4. It is also easy to see that f_{z_i} is monotone. Thus, since a monotone ROBP of length n can have total influence at most \sqrt{n} (like any monotone function on n bits), and since f_{z_i} has length at most $\lceil \frac{n}{r} \rceil \neq n$, we obtain by induction on the width that

$$\mathbb{E}_{z_i}[\text{Inf}(f_{z_i})] \leq 2r^{w-3} + \frac{n^{1.5} r^{w-3} - 1}{2^{r-1} r - 1} + \frac{n}{r 2^{r-1}} \sqrt{n}.$$

Therefore,

$$\begin{aligned} \text{Inf}(f) &= \sum_{i=0}^{r-1} \sum_{j \in \mathcal{Z}_i} \text{Inf}_j(f) = \sum_{i=0}^{r-1} \mathbb{E}_{z_i}[\text{Inf}(f_{z_i})] \\ &\leq r \left(2r^{w-3} + \frac{n^{1.5} r^{w-3} - 1}{2^{r-1} r - 1} + \frac{n^{1.5}}{r 2^{r-1}} \right) \end{aligned}$$

which is the desired bound. \square

Setting $r = \lceil 1.5 \log(n) \rceil$ in Theorem 5 (for which we need to assume $n \geq 2$) immediately implies:

Corollary 2. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a cROBP of width $w \geq 2$ and length $n \geq 2$. Then,*

$$\text{Inf}(f) \leq 4 \lceil 1.5 \log(n) \rceil^{w-2}.$$

We note that Corollary 2 is not an asymptotic statement about “large n and constant w ”; it holds for all combinations of n and w with $n, w \geq 2$.