

Learning a Circuit by Injecting Values

[Extended Abstract]

Dana Angluin *

James Aspnes *†

Jiang Chen *

Yinghua Wu *‡

ABSTRACT

We propose a new model for exact learning of acyclic circuits using experiments in which chosen values may be assigned to an arbitrary subset of wires internal to the circuit, but only the value of the circuit's single output wire may be observed. We give polynomial time algorithms to learn (1) arbitrary circuits with logarithmic depth and constant fan-in and (2) Boolean circuits of constant depth and unbounded fan-in over AND, OR, and NOT gates. Thus, both **AC0** and **NC1** circuits are learnable in polynomial time in this model. Negative results show that some restrictions on depth, fan-in and gate types are necessary: exponentially many experiments are required to learn AND/OR circuits of unbounded depth and fan-in; it is NP-hard to learn AND/OR circuits of unbounded depth and fan-in 2; and it is NP-hard to learn circuits of bounded depth and unbounded fan-in over AND, OR, and threshold gates, even when the target circuit is known to contain at most one threshold gate and that threshold gate has threshold 2. We also consider the effect of adding an oracle for behavioral equivalence. In this case there are polynomial-time algorithms to learn arbitrary circuits of constant fan-in and unbounded depth and to learn Boolean circuits with arbitrary fan-in and unbounded depth over AND, OR, and NOT gates. A corollary is that these two classes are PAC-learnable if experiments are available.

1. INTRODUCTION

We introduce a new model of active learning for acyclic circuits in which we may inject chosen values on an arbitrary subset of wires but can observe only the value of the circuit's output wire. Our results illuminate the relative importance of manipulation and observation in discovering the structure

*Department of Computer Science, Yale University. Email: {angluin,aspnes,criver}@cs.yale.edu, y.wu@yale.edu.

†Supported in part by NSF grants CNS-0305258 and CNS-0435201.

‡Supported by NSF grant CNS-0305258.

of networks modeled as circuits.

Gene regulatory networks are an important area in which Boolean network models have been used. In one variant of the basic model, each node in a finite network represents a gene, which has a current state of active or inactive. The states of all nodes in the network are updated synchronously; for each node there is a Boolean function giving its new state in terms of the current states of some subset of the other nodes. A key point is that the node states are *fully observable*: it is assumed that gene expression data gives the state of every node in the network at every time step. The discovery problem is to learn the updating functions of all the nodes (note that one needs to learn both the set of inputs and the functionality of each node). Of course this is difficult if the updating function may be an arbitrary Boolean function; further assumptions generally restrict the fan-in or types of the possible updating functions. One of the main difficulties in this problem is to discover the topology of the network (which nodes are inputs to which node).

Akutsu et al. [1] describe an approach to the discovery problem that models the experimental capability of multiple gene disruption and overexpression. At each time step several selected genes may be disrupted (put in the inactive state), several other selected genes may be overexpressed (put in the active state), while unaffected genes are updated as usual. In this model the states of the nodes are *fully controllable* as well as fully observable. For networks of N nodes and fan-in bounded by k , Akutsu et al. give an $O(N^{2k})$ algorithm for the discovery task. Ideker, Thorsson, and Karp [7] also consider this model and give more practical discovery methods for acyclic networks, using information theoretic criteria to select genes to disrupt or overexpress. These results show that if the class of updating functions is sufficiently restricted, the problem of learning the structure of a network in this model is tractable.

By contrast, there is ample evidence that learning Boolean circuits or formulas from their input-output behaviors may be computationally intractable. Positive learnability results include those for fairly limited classes, including propositional Horn formulas [2] general read once Boolean formulas [3], and decision trees [5], and those for specific distributions, including **AC0** circuits [12], DNF formulas [8] and **AC0** circuits with a limited number of majority gates [9]. (Note that algorithms in both papers [12] and [9] for learning **AC0** circuits and their variants run only in quasipolynomial time.) Valiant gives cryptographic evidence for the difficulty of PAC learning general Boolean circuits [15]. Kearns and Valiant [10] show that specific cryptographic

assumptions imply that **NC1** circuits and **TC0** circuits are not PAC learnable in polynomial time. These negative results have been strengthened to the setting of PAC learning with membership queries [4], even with respect to the uniform distribution [11].

For these results on learning circuits and formulas, observation and control are both restricted: values on internal wires cannot be observed or manipulated. A natural question is: *What are the relative contributions of full observation and full control to the tractability of learning Boolean networks?*

Our new model addresses this question: we postulate full control and restricted observation. Our results show that the ability to inject values into the circuit gives the learner considerable power, but not as much as would be the case with full observation. In particular, with value injection queries, **NC1** circuits and **AC0** circuits are exactly learnable in polynomial time, but our negative results show that the depth limitations are necessary. However, if behavioral equivalence queries are also available, the depth limitations can be removed, which also implies the polynomial time PAC-learnability of these classes with value injection queries.

2. THE MODEL

2.1 Circuits

We define a variant of the usual circuit model that has no distinguished inputs and permits a finite set Σ of at least two different values on wires. A *circuit* C consists of N wires, $W = \{w_1, w_2, \dots, w_N\}$, and for each wire w_i a *gate* g_i that determines the value on this wire. The *size* of the circuit is N . The wire w_N is assumed to provide the output of the circuit as a whole. A *gate* consists of a function mapping Σ^k to Σ , and a vector of k integers from $[1, N]$ specifying the *input wires* of the gate. The value k is the *fan-in* of the gate. Gates of fan-in zero compute constant functions. The maximum fan-in taken over all the gates in the circuit is the *fan-in* of the circuit. We define the *circuit graph* to have a node for each wire/gate pair and a directed edge from node i to node j if w_i is one of the input wires to gate j . Until Section 7, we assume that the graph of the circuit is acyclic. We define the *depth* of the circuit to be the number of edges in the longest path in the circuit graph.

2.2 Behavior

We focus on the behavior of a circuit in response to experiments in which we fix the values of certain wires and observe the final output of the circuit. Define an *experiment* to be a vector s in $(\Sigma \cup \{*\})^N$, where s_i specifies the value of w_i (if it is in Σ) or leaves the value of w_i unaltered (if it is $*$). If $s_i \in \Sigma$, we say w_i is *fixed* in s ; otherwise, it is *free* in s . The value of w_i given s , written $w_i(s)$, is defined as

$$w_i(s) = \begin{cases} g_i(w_{i_1}(s), w_{i_2}(s), \dots, w_{i_{k_i}}(s)) & \text{if } s_i = *, \\ s_i & \text{if } s_i \neq *. \end{cases} \quad (1)$$

where gate i has function g_i and inputs $(i_1, i_2, \dots, i_{k_i})$. We remark that gates of fan-in zero, which compute constant functions, give the base cases for the above recursive definition. The output of the circuit given an experiment s is the output of wire w_N , that is, $w_N(s)$; this is also denoted $C(s)$.

The *behavior* of a circuit is the function mapping experiments s to $C(s)$. Two circuits C and C' are *behaviorally equivalent*, if they have the same behavior, that is, if $\forall s \in (\Sigma \cup \{*\})^N, C(s) = C'(s)$. To compare our work with previous work on learning circuits, we treat the gates of fan-in zero as the *input gates* and denote the number of input gates by n . An experiment is *input-only* if it fixes every input gate and leaves every other gate free. The *input-output behavior* of a circuit C is the restriction of its behavior function to input-only experiments. Clearly behavioral equivalence implies equality of input-output behaviors but not conversely. Behavioral equivalence is more constrained by the internal structure of the circuits.

An important special case is *Boolean circuits*, for which $\Sigma = \{0, 1\}$. The input-output behavior of a Boolean circuit is just the usual concept of a circuit computing a Boolean function of n inputs. **NC1** circuits are Boolean circuits with constant fan-in and depth $O(\log n)$. **AC0** circuits are Boolean circuits of constant depth and polynomial size whose gates are unbounded fan-in AND, OR, and NOT. The *threshold function* Θ_t is the Boolean function that is 1 if and only if at least t of its inputs are 1. **TC0** circuits are Boolean circuits of constant depth and polynomial size whose gates are unbounded fan-in threshold, AND, OR and NOT.

2.3 Queries

We assume that the learner can get information about the circuit by specifying an experiment s and observing $C(s)$, the output of the circuit. Such an action is termed a *value injection query*, abbreviated VIQ. We also define a *behavioral equivalence query*, abbreviated BEQ: the learner proposes a circuit C' , and, if it is not behaviorally equivalent to the target circuit C , receives in response a *counterexample*, that is, an arbitrarily chosen experiment s such that $C'(s) \neq C(s)$. To be consistent with previous usage, we use the term *membership query*, for a VIQ restricted to an input-only experiment s , although $C(s)$ may be non-binary, and the term *equivalence query*, for a query that tests whether the proposed circuit C' has the same input/output behavior as the target circuit C and returns an arbitrary input-only experiment s witnessing $C'(s) \neq C(s)$ if not. An algorithm that learns the (full) behavior of any circuit from a given class using VIQ's and BEQ's also learns the input/output behavior of any circuit from the class using VIQ's and equivalence queries. Then a standard polynomial time transformation yields a PAC learning algorithm using VIQ's for input/output behavior of circuits in the class, which implies the following.

PROPOSITION 1. *If a class of circuits is learnable in polynomial time with VIQ's and BEQ's, then the input/output behaviors of circuits in the class are PAC-learnable in polynomial time using VIQ's.*

2.4 The Problem

The learning problems we address are: by making VIQ's (respectively, VIQ's and BEQ's) to a target circuit C , find a behaviorally equivalent circuit C' . Both C and C' use gates from a specified class \mathcal{F} .

We remark that it is not possible to discover the exact structure of the target circuit, even when all wires are relevant. The following example shows that the same behavior may be exhibited by structurally distinct circuits. The three circuits C_1 , C_2 , and C_3 shown in Figure 1 are behaviorally

Table 1: Summary of our results

Depth	Fan-in	Gates	Query types	Learnability	Reference
Unbounded	Unbounded	AND/OR	VIQ	$2^{\Omega(N)}$ queries	Theorem 2
Unbounded	2	AND/OR	VIQ	NP-hard	Theorem 3
Constant	Unbounded	AND/OR/ Θ_2	VIQ/BEQ	NP-hard	Theorem 6
Logarithmic	Constant	Arbitrary	VIQ	Poly-time	Theorem 19
Constant	Unbounded	AND/OR/NOT	VIQ	Poly-time	Theorem 21
Unbounded	Constant	Arbitrary	VIQ/BEQ	Poly-time	Theorem 22
Unbounded	Unbounded	AND/OR/NOT	VIQ/BEQ	Poly-time	Theorem 25

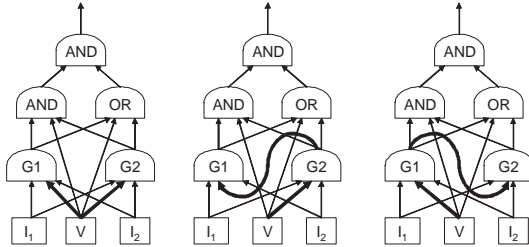


Figure 1: Three equivalent circuits

equivalent, where G_1 and G_2 are arbitrary gate functions. Only when G_2 and V both have value 1 (respectively, 0) can the value of G_1 propagate through the depth 1 AND gate (respectively, OR gate). Therefore we cannot decide which one of G_2 and V is an input of G_1 , and similarly in the other case.

3. RESULT SUMMARY

We investigate the computational tractability of these problems for classes of circuits defined by restrictions on depth, fan-in, and \mathcal{F} . Our results are summarized in Table 1.

Section 4 contains negative results for exact learning with VIQ’s and BEQ’s. In Section 5 we give an algorithm, CircuitBuilder, that takes a class of gates \mathcal{F} and a set U of experiments and constructs a circuit C' by making VIQ’s on experiments in U and their one-symbol perturbations, and then finding a gate for each wire consistent with the results. If U contains for every wire and every gate that is wrong for that wire a witness experiment that excludes the incorrect gate, then the resulting circuit is behaviorally equivalent to the target circuit. We then show how to construct appropriate sets of experiments U for the class of log-depth constant fan-in circuits and for the class of **AC0** circuits. In Section 6, we extend these methods to use BEQ’s as well as VIQ’s, and show that the limitations on circuit depth can be removed for both classes.

3.1 Learnability of the gates

What is the relationship between the learnability of circuits in our model and learnability of the class \mathcal{F} of permitted gates? A depth 1 circuit consists of n input gates and one gate g depending on some subset of the inputs; if any nontrivial circuits are to be learnable, then depth 1 circuits must be learnable in the same sense. Therefore, we assume that the class of permitted gates \mathcal{F} is learnable.

For depth 1 circuits, a VIQ reduces to a membership query. Classes \mathcal{F} of gates for which depth 1 circuits are learnable in polynomial with membership queries include

- (1) the class of gates with fan-in at most some constant k over an arbitrary finite value set Σ and
- (2) the class of all symmetric Boolean gates (which includes unbounded fan-in AND, OR, NAND, NOR, threshold and parity gates.)

Another aspect of the learnability of depth 1 circuits is the *consistency problem*, defined as follows. The input is a set E of *prohibited pairs* (s, σ) where s is an input-only experiment (recall that an input-only experiment fixes all input gates), $\sigma \in \Sigma$ is a value, and the desired output is a depth 1 circuit C' over \mathcal{F} such that $C'(s) \neq \sigma$ for every pair $(s, \sigma) \in E$. For Boolean circuits, because $C'(s) \neq 0$ implies $C'(s) = 1$ and similarly $C'(s) \neq 1$ implies $C'(s) = 0$, the consistency problem can be stated as finding C' that agrees with given values of experiments in E . We remark that the consistency problem is a computational problem and therefore each s should fix all input gates for the problem to be well defined. One of the major differences between the consistency problem and the problem of learning with membership queries is that in the consistency problem no information is provided beyond the set E and hence one may not be able to query Hamming neighbors of an experiment. For the class of arbitrary gates with fan-in at most k , the consistency problem can be solved in time $O(n^k)$. There is also a polynomial time algorithm to solve the consistency problem over the class of unbounded fan-in AND, OR, NAND, NOR, NOT and parities. However, Lemma 5 shows that the consistency problem is NP-hard over the class of unbounded fan-in AND, OR, and thresholds. CircuitBuilder makes use of algorithms for the consistency problem. Lemma 4 shows that polynomial time learnability with VIQ’s and BEQ’s implies a polynomial time algorithm for consistency in certain cases.

3.2 Relation to circuit testing

A central challenge for learning algorithms in our model is to propagate the effects of a changed value on some internal wire to the observable output of the circuit. Our methods are similar in some respects to the idea of *path sensitization* in circuit testing, used to detect whether the output of some gate is “stuck” at a constant value instead of computing the correct value of its inputs [6]. Path sensitization is not a complete method: there are examples of single stuck-at faults in acyclic circuits that cannot be detected by path sensitization, though they are detectable by other tests. In our model, the ability to inject values on internal wires gives the approach greater power. However, this power does not trivialize the problem; the negative results in Section 4 illustrate the subtle “shadowing” or “filtering” effects that limit the power of VIQ’s.

4. WHAT CANNOT BE LEARNED EFFICIENTLY?

Figure 2 presents a *gadget*, which is an AND/OR circuit of fan-in 2 that computes $Y = \Theta_2(X, V, W)$. We can view V and W as controlling a switch: only when their values are different will the value of X be passed on to Y . Gadgets can be concatenated to get a gadget chain (see Figure 3), in which the value of X is passed on to Y only when every pair V_i and W_i have different values. The chain can be used to “hide” part of the circuit unless the values of V_i and W_i are complements of each other. In Figure 3, exactly one of each pair V_i and W_i is an input to the big AND gate. The learner has to guess which combination of them are the inputs to the AND gate, which yields the following negative result.

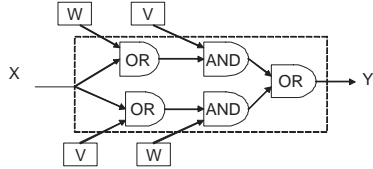


Figure 2: The gadget

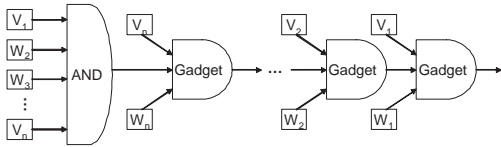


Figure 3: A hidden AND gate and the gadget chain

THEOREM 2. *Learning the class of acyclic Boolean AND/OR circuits requires $2^{\Omega(N)}$ VIQ’s.*

PROOF. Suppose an adversary reveals the gadget chain and the fact that the last gate is an AND gate with exactly one of each pair V_i and W_i as an input, but hides the exact combination. When there exists a pair V_i and W_i that are both set 0 or 1, the output of the circuit is determined by the gadget chain. In particular, it is determined by the pair with smallest index that are both set 0 and 1 (the pairs are ordered according to their distances to the output gate as in Figure 3). The adversary answers 0 if the pair are set 0 and 1 if the pair are set 1.

Only when for every pair, V_i and W_i are set differently, will the value of the big AND gate affect the output of the circuit. There are 2^n such settings of V ’s and W ’s. The adversary answers 0 until only one setting remains. The theorem then follows because $N = O(n)$. \square

Theorem 2 gives an exponential information-theoretic lower bound, using a deep circuit and a gate of large fan-in. The following construction uses the gadget chain to give a computational hardness result for deep circuits with fan-in 2.

THEOREM 3. *Learning the class of fan-in 2 AND/OR circuits using VIQ’s is NP-hard.*

PROOF. We use the implicit negation enforced by the gadget chain to construct a circuit representing a CNF formula, for which a satisfying instance must be found in order to expose a hidden part of the circuit.¹ We associate a Boolean

¹Using a NOT gate would not achieve the same effect because we can override its output in an experiment.

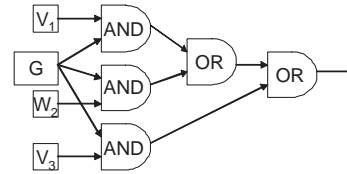


Figure 4: $G \wedge (x_1 \vee \overline{x_2} \vee x_3)$

variable x_i with each pair V_i and W_i , and let $x_i = 1$ if $V_i = 1, W_i = 0$ and 0 if $V_i = 0, W_i = 1$ (we only deal with the case that V_i and W_i are set differently thanks to the gadget chain). Then the circuit in Figure 4 computes $G \wedge (x_1 \vee \overline{x_2} \vee x_3)$. We can chain such clauses by replacing G with the output of the succeeding clause and finally connect the output of the first clause to the gadget chain shown in Figure 3 (note that the big AND gate is not part of the gadget chain). Let g be an AND gate of fan-in 2 with different inputs than V ’s and W ’s. We connect the output of g to the last clause of the clause chain (by replacing G of the last clause by the output of g). In order to learn the circuit, we have to be able to observe the output of g because we are not able to distinguish between this circuit and the circuit with g being replaced by an OR gate if we can not observe g ’s output. In this construction, the functionality of g matters if there exists an satisfying assignment to all clauses.

As in the proof of Theorem 2, when there exists a pair V_i and W_i that are both set to 0 and 1, the output of the circuit is determined by the gadget chain regardless of other parts of the circuit. Suppose every pair of V_i and W_i are set differently. In order to observe g ’s output, we must compute an assignment to V ’s and W ’s that satisfies all clauses in the clause chain, since otherwise the circuit output will simply be 0. In other words, we have to solve the 3-SAT problem. Because the gadget chain consists of AND/OR gates of fan-in 2, the theorem follows. \square

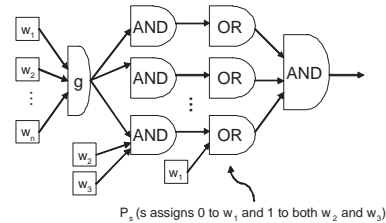


Figure 5: The “filtering” circuit and an illustration of a filtering path P_s with three input wires.

The gadget chain is a deep circuit, but in the following construction we use AND and OR gates to achieve a constant depth filter that forces a learning algorithm to solve the consistency problem for \mathcal{F} (defined in Section 3.1). Given (1) any depth 1 circuit with input wires w_1, w_2, \dots, w_n and gate g from a class \mathcal{F} of Boolean gates and (2) a set E of input-only experiments, we add the following structure to construct another circuit C as follows (see Figure 5). For each $s \in E$, which assign each w_i to 0 or 1, we add a distinct directed path P_s of length 3 consisting of g , a new AND gate, a new OR gate, and the output gate. Let each wire w_i that is set 0 in s be an input of the OR gate of P_s .

Let each wire w_i that is set 1 in s be an input of the AND gate of P_s . The construction has the property that if the assignment to w_1, w_2, \dots, w_n is not s , either the output of the AND gate or the output of the OR gate in P_s is determined, and hence the output of g can not be passed through the path P_s . Therefore, P_s “filters” out all assignments in E except s . Finally, we take the output gate of the whole circuit to be an AND gate.²

LEMMA 4. *There is a polynomial time algorithm using VIQ’s and BEQ’s to learn circuits from the above class if and only if there is a polynomial time algorithm for the consistency problem over \mathcal{F} .*

PROOF. Suppose the consistency problem is solvable in polynomial time. The structure of a target circuit C from the class, (except for the identity of gate g) can be learned using VIQ’s in a straightforward fashion, which determines the corresponding set E . For each $s \in E$ (which corresponds to a path), we construct an experiment s' that agrees with s on w_1, w_2, \dots, w_n , and sets the output gate and the OR gate and the AND gate in P_s free, sets all the other OR gates 1. It is clear that $C(s') = g(s)$. Therefore, we can learn the value $g(s)$ for each $s \in E$ and learn g using the algorithm that solves the consistency problem. Let g' be the function learned. g and g' are equivalent on *experiments* of E . Because of the filtering structure of C , the only experiments in which the value of g matters are those that assign w_1, w_2, \dots, w_n in the exactly same way as some $s \in E$ does. Therefore, we conclude that the learned circuit is equivalent to C .

Conversely, suppose that there is a polynomial time learning algorithm. We can construct a circuit that corresponds to the consistency problem and simulate VIQ’s by evaluating the circuit. We simulate BEQ’s only when the proposed gate for g does not solve the consistency problem since otherwise the learning algorithm solves the consistency problem and we are done. Denote the proposed gate by g' . Let C' be the proposed circuit. We show that whenever g' does not solve the consistency problem, we can find a counterexample to C' . Therefore, the learning algorithm must propose a gate that solves the consistency problem.

There are two cases. When the inputs of g' are contained in the set $\{w_1, w_2, \dots, w_n\}$, there must exist an $s \in E$ such that $g'(s) \neq g(s)$ as g' does not solve the consistency problem. In fact, we can find such an s in polynomial time by simulating a membership query algorithm on g' (recall that we assume that gates are learnable with membership queries). Among all the queries the algorithm makes, a query must be made to distinguish g and g' and the query can serve as s . Suppose w.l.o.g. that $g(s) = 1$ and $g'(s) = 0$. We construct s' as before. We give s' as a counterexample if $C(s') \neq C'(s')$. Otherwise, we construct two experiments s'_0 and s'_1 that set in s' the corresponding wire of g to 0 and 1 respectively. By our construction, we know that $C(s'_0) = 0$ and $C(s'_1) = 1$. However, we have that $C(s') = C(s'_1)$ (since $g(s) = 1$) and that $C'(s') = C'(s'_0)$ (since $g'(s) = 0$). Therefore $C'(s'_0) = C'(s') = C(s') = C(s'_1) = 1 \neq C(s'_0)$. Thus s'_0 is a counterexample to C' .

If g' depends on an AND/OR gate h that lies in a path P_s for some $s \in E$, construct s', s'_0 and s'_1 as before. Let w be the corresponding wire of g and w' be the corresponding wire

of h . Suppose that C and C' agree on all three experiments. Compare s'_0 and s'_1 . In the proposed circuit C' , w' must take the same value on s'_0 and s'_1 , because the fact that g' depends on w' and that C' is acyclic implies that in C' , w' is not reachable from w and therefore changing the value on w will not affect the value on w' (see also Proposition 9). However, by construction, in C , w' takes value 0 on s'_0 and value 1 on s'_1 , which means that, on either of s'_0 and s'_1 , w' takes different values in C and C' . Suppose w.l.o.g. it is s'_0 . Then for s'_0 , w takes value 0 in C and value 1 in C' . Set w' to 1 in s'_0 and let the resulting experiment be s'_2 . We have that $C'(s'_2) = C'(s'_0) = C(s'_0) = 0$ but $C(s'_2) = 1$. Thus, s'_2 is a counterexample to C' . \square

LEMMA 5. *The consistency problem is NP-hard for the class of unbounded fan-in AND, OR, and Θ_2 gates.*

PROOF. We reduce a 3-SAT instance ϕ over the variables x_i for $i \in [1, n]$ to the consistency problem for this class. The input wires are $\{I_1, I_2, I_3, V_1, W_1, V_2, W_2, \dots, V_n, W_n\}$. Let G denote the output wire of the unknown gate. We define a correspondence between literals of ϕ and wires: literal x_i corresponds to wire V_i literal \bar{x}_i corresponds to wire W_i . We design the set of experiments and their outputs as follows, so as to constrain the unknown gate to be a Θ_2 gate with its inputs corresponding to a satisfying assignment of ϕ .

- For each of the eight experiments assigning 0 and 1 to I_1, I_2 , and I_3 , with all $V_i = W_i = 0$, the output value is $\Theta_2(I_1, I_2, I_3)$. This guarantees that G cannot be AND or OR, and must therefore be a Θ_2 gate whose inputs include I_1, I_2 , and I_3 .
- For each i , on the experiment with $I_1 = V_i = W_i = 1$ and all other input wires assigned 0, the output value is 1. This implies at least one of V_i and W_i is an input wire of G .
- For each i , on the experiment with $V_i = W_i = 1$ and all other input wires assigned 0, the output value is 0. This implies not both V_i and W_i are inputs to G .
- For each clause of ϕ , on the experiment that sets I_1 and the three wires corresponding to the three literals in the clause to 1 and all other wires to 0, the output is 1. This ensures that at least one wire corresponding to a literal in the clause is an input to G .

It is easily verified that ϕ is satisfiable if and only if there is a gate G from the specified class of gate functions consistent with these experiment/value pairs. \square

Lemma 4 and Lemma 5 establish the following theorem.

THEOREM 6. *Learning constant depth AND/OR/ Θ_2 circuits with VIQ’s and BEQ’s is NP-hard.*

5. LEARNING WITH EXPERIMENTS

In this section, we give algorithms for arbitrary circuits with logarithmic depth and constant fan-in and Boolean circuits of constant depth and unbounded fan-in over AND, OR and NOT gates. Therefore, we show that both **AC0** and **NC1** circuits are learnable in polynomial time with VIQ’s. One of the main issues is to learn the set of inputs for each

²Note that we cannot use the same method to replace the gadget chain because it would require $|E|$ to be exponential.

gate. We remark that the gates in the circuit that the algorithms output may not have the same set of inputs as the target circuit (see Figure 1).

First we develop some tools. A *partial experiment* is a partial function from $[1, N]$ to $\Sigma \cup \{*\}$, where some wires are unspecified. Let s be an experiment and τ be a partial experiment. Define $s|_\tau$ to be the experiment obtained by replacing in s the settings of all wires that are specified in τ by the corresponding settings in τ . Let s and t be two experiments. We say that $t \preceq s$ if the set of free wires in t is a subset of the set of free wires in s . We say $t < s$ if $t \preceq s$ and there is at least one free wire in s that is fixed in t . Let s be an experiment with wire w set free. We call $s|_{w=\sigma}$, where $\sigma \in \Sigma$, the (w, σ) -perturbation of s . If $C(s) \neq C(s|_{w=\sigma})$, we say s is (w, σ) -exposing.

Consider any gate g with inputs (i_1, i_2, \dots, i_l) . We overload g to take an experiment s as an argument. That is, let $g(s) = g(w_{i_1}(s), w_{i_2}(s), \dots, w_{i_l}(s))$, where $w_i(s)$ is the value of wire w_i on s in the target circuit C . The following useful facts are easily verified.

PROPOSITION 7. *Let s and t be two experiments with the output wire set free. If s and t agree on every wire that is either free in s or an input to a wire that is free in s then $C(s) = C(t)$.*

PROPOSITION 8. $C(s) = C(s|_{w=w(s)})$. *(This is meaningful only when w is set free in s . In this case, $w(s)$ is the value the corresponding gate computes. The proposition simply says that if we fix w to the value it takes on an experiment s , the circuit output stays the same.)*

PROPOSITION 9. *Let w and u be two wires and suppose there is no path from w to u in the graph of the circuit. Then $u(s) = u(s|_{w=\sigma})$ for any experiment s and $\sigma \in \Sigma$.*

The main task of our learning algorithms is to find a “correct” gate function for each wire. Formally speaking, a gate g is *wrong* for a wire w , if there exists an experiment s that fixes all of g ’s inputs and is $(w, g(s))$ -exposing. We call such an s a *witness* experiment for g and w . Otherwise, we say that g is *correct* for w .

LEMMA 10. *Let C' be a circuit with the same set of wires as C . If C' is acyclic and every gate of C' is correct for the corresponding wire in C , C' is behaviorally equivalent to C .*

PROOF. Suppose to the contrary that C' is not behaviorally equivalent to C . Let s be a minimal experiment such that $C'(s) \neq C(s)$. Let w be a free wire in experiment s and g be its corresponding gate in C' , chosen so that all g ’s inputs are fixed in s (such a wire exists because C' is acyclic and the input gates of C' are considered to have fixed inputs because they have no inputs). By Proposition 8, we have $C'(s) = C'(s|_{w=g(s)})$. By the minimality of s , we have $C'(s|_{w=g(s)}) = C(s|_{w=g(s)})$, which then implies that $C(s|_{w=g(s)}) \neq C(s)$. This contradicts the fact that g is correct for w . \square

5.1 Constructing a circuit

Let \mathcal{F} be a class of gates containing all the gates in the target circuit C . We describe an important subroutine, *CircuitBuilder* (Algorithm 1), that takes a set of experiments U and constructs an acyclic circuit C' using gates from \mathcal{F} . *CircuitBuilder* builds C' from the bottom up, starting with

gates of fan-in zero. At each iteration, *CircuitBuilder* attempts to add another wire to C' by choosing a gate in \mathcal{F} among those that depend only on wires that are already in C' . This method has the advantage of building an acyclic circuit, which is crucial because the dependence between gates is not always clear, as in Figure 1.

Define U to be a *sufficient set of tests* for C and \mathcal{F} if for every wire w_i in C and every gate $g \in \mathcal{F}$ that is wrong for w_i , U contains at least one witness for g and w_i . In the remainder of this section we prove the following.

THEOREM 11. *If U is a sufficient set of tests for C and \mathcal{F} then C' is behaviorally equivalent to C , where C' is the circuit constructed by *CircuitBuilder*. Moreover, the queries can be made non-adaptively.*

Let U_w denote the set of experiments in U with w set free. In *CircuitBuilder*, since before we replace s by $s|_{w=g(s)}$, we check whether $C(s) = C(s|_{w=g(s)})$ for s in U_w , we do not need to make queries on the replacing experiments in U . However, we may have to make queries on the perturbations of replacing experiments. This would require the algorithm to make queries adaptively. Instead, in *CircuitBuilder*, we maintain another set of experiments V which contains all possible perturbations of experiments in U at the beginning of the algorithm. Similarly, let V_w denote the set of experiments in V with w set free. In Lemma 13, we will show that after replacement, V still contains all necessary perturbations of experiments in U . In Lemma 14, we show that even in V , a replacing experiment will have the same circuit output as the original one. Therefore, we only need to make queries on $U \cup V$ at the beginning of the algorithm. The algorithm is non-adaptive.

Algorithm 1 *CircuitBuilder*

INPUT: U and \mathcal{F} .

OUTPUT: C' .

- 1: $\forall w, \forall s \in U_w, \forall \sigma \in \Sigma$, let V contain the (w, σ) -perturbation of s .
 - 2: Make a VIQ on every experiment $s \in U \cup V$.
 - 3: $C' \leftarrow \emptyset$. $Z \leftarrow W$.
 - 4: **while** Z is not empty **do**
 - 5: **for** $w \in Z$ **do**
 - 6: **if** there exists a function $g \in \mathcal{F}$ that depends only on wires in C' , such that $\forall s \in U_w, C(s) = C(s|_{w=g(s)})$ **then**
 - 7: Add w and g to C' and remove w from Z .
 - 8: $\forall s \in U_w \cup V_w$, replace s by $s|_{w=g(s)}$.
 - 9: **break**
-

At each iteration of the algorithm, experiments in $U \cup V$ may be replaced. We make the following claims about the replacements.

LEMMA 12. *At any iteration, for any $s \in U \cup V$, no wire in C' is set free in s .*

PROOF. This is because we fix the wire to a value whenever we add it to C' . \square

The following lemma says that if $s \in U$ and $t \in V$ are an experiment and perturbation pair, and w is the corresponding wire, they will continue to be such a pair until w is added to C' .

LEMMA 13. Consider any iteration, any $w \in Z$ and any $s \in U_w$, and let s_0 be the version of s at the start of the algorithm. For any $\sigma \in \Sigma$, let t_0 be the (w, σ) -perturbation of s_0 at the start of the algorithm and t be the replacement of t_0 at the iteration considered. Then t is a (w, σ) -perturbation of s .

PROOF. The statement is clearly true at the start of the algorithm. At each previous iteration, at most one setting of s and t will be changed. We only need to show that each replacement will replace the same value for s and t . The lemma then follows by observing that the replaced value is an output of a function that depends on wires that do not include w (w has not been added to C') and that s and t differ only at their settings of w . \square

Together with Lemma 13, the following lemma shows that although we need to compare the circuit outputs of replacement experiments and their perturbations, we do not need to make any further VIQ's.

LEMMA 14. Suppose U is a sufficient set of tests for C and \mathcal{F} . At any iteration consider any $s \in U \cup V$ and let s_0 be the version of s at the start of the algorithm. Then we have $C(s) = C(s_0)$.

If $s \in U$, there is nothing to prove since the algorithm checks the equality before making the replacement. The case that $s \in V$ is a little bit trickier. We prove an even more general lemma, from which the case $s \in V$ follows.

LEMMA 15. Suppose U is a sufficient set of tests for C and \mathcal{F} . Let g be the function that the algorithm chooses for gate w . Then g is correct for w . That is, $\forall s$ with g 's input wires fixed, $C(s) = C(s|_{w=g(s)})$.

PROOF. W.l.o.g., let w be the first wire added to C' for which the statement in the lemma does not hold. Suppose g is wrong for w . By the assumption that U is sufficient for C and \mathcal{F} , there exists an experiment $s_0 \in U$ at the start of the algorithm such that s_0 fixes all g 's inputs, and $C(s_0) \neq C(s_0|_{w=g(s_0)})$. Let $s \in U$ be the replacement of s_0 at the iteration w is added to C' . We have that $C(s) = C(s_0) \neq C(s_0|_{w=g(s_0)}) = C(s|_{w=g(s_0)})$, by the assumption that w is the first wire violating the condition. Moreover, $g(s) = g(s_0)$ because s_0 and s both fixes all g 's input wires and therefore should agree on them. Therefore, we have

$$C(s) \neq C(s|_{w=g(s)})$$

which contradicts the choice of g . \square

Lemma 15 together with Lemma 10 show that if U is a sufficient set, C' is behaviorally equivalent to C . Lemma 13 and 14 show that all the queries can be made at the beginning of the algorithm, which establishes Theorem 11. Lemma 12 validates the operation of picking a function g , which amounts to solving the following consistency problem (defined in Section 3.1). Let E be the projection of U_w to C' (note that all wires in C' are fixed) and let the prohibited pairs (t, σ) be those $t \in E$ and $\sigma \in \Sigma$ such that there is an experiment in U_w that agrees with t and is (w, σ) -exposing.

The next lemma shows that the algorithm terminates in N iterations.

LEMMA 16. At each iteration, the algorithm adds one wire to C' .

PROOF. First we observe that there is at least one wire w in Z whose input wires are all contained in C' , because the circuit graph of C is acyclic. The true gate of w in C will survive every *if-test* in the algorithm. \square

5.2 Test paths

One of the key ideas at the heart of our algorithms is to use test paths. A *test path* is an experiment whose free wires are a directed path from some wire w to w_N , through which w is exposed. Let a *side wire* of a test path be a fixed wire that is an input to a gate whose corresponding wire is set free. The meaning of test paths is made clear in the following lemma, which says that using test paths is sufficient.

LEMMA 17. Let s be a (w, σ) -exposing experiment, where $\sigma \in \Sigma$. Let $s^* \preceq s$ be a minimal (w, σ) -exposing experiment. Then the free wires in s^* are a directed path in the graph of C , which starts with w and ends with the output wire w_N . (s^* is a test path.)

PROOF. When $w = w_N$, the directed path is just w_N itself. Suppose the claim is true for any free wire whose corresponding gate has w as an input. First we claim that only those wires that w can reach (in the underlying digraph) can be free in s^* . This is because these wires take the same values in s^* and the perturbations $s^*|_{w=\sigma}$ (see Proposition 9). Thus, we can set them to the corresponding values and the resulting experiment is still (w, σ) -exposing, which contradicts the minimality of s^* .

Let u be a free wire in s^* whose only free input wire is w . u must exist, because the underlying digraph is acyclic. Let $\sigma_0 = w(s^*)$ and $\beta_0 = u(s^*)$ and $\beta = u(s^*|_{w=\sigma})$. We claim that $s^*|_{w=\sigma_0}$ is a minimal (u, β) -exposing experiment. Let us view the circuit as a function of values of w and u . That is, let $F(x, y) = C(s^*|_{w=x, u=y})$. By the assumption, we have $F(\sigma_0, \beta_0) \neq F(\sigma, \beta)$. By the minimality of s^* , we have $F(\sigma_0, \beta) = F(\sigma, \beta)$. Thus, we have

$$F(\sigma_0, \beta) \neq F(\sigma_0, \beta_0)$$

which implies that $s^*|_{w=\sigma_0}$ is (u, β) -exposing. Suppose on the contrary that there exists $s' \prec s^*$ is (u, β) -exposing. We set both w and u free in s' and let the resulting experiment be s'' . Let $F''(x, y) = C(s''|_{w=x, u=y})$. Again, by the assumption and the minimality of s^* , we have

$$F''(\sigma_0, \beta_0) \neq F''(\sigma_0, \beta) = F''(\sigma, \beta)$$

Therefore, we conclude s'' is (w, σ) -exposing by observing that $w(s'') = \sigma_0$, $u(s'') = \beta_0$ and $u(s''|_{w=\sigma}) = \beta$, which leads to a contradiction.

Therefore, we conclude that $s^*|_{w=\sigma_0}$ is a minimal (u, β) -exposing experiment. By induction, its free wires consist of a directed path starting with u and ending with w_N . We append w to this path to obtain the directed path in s^* . \square

5.3 Learning log depth circuits with constant fan-in

We use CircuitBuilder to give an algorithm that learns an arbitrary log depth bounded fan-in circuit. The algorithm does not perform any additional queries and hence is non-adaptive. The main idea is based on the observation that in an acyclic circuit of depth d and fan-in k , there are at most d free wires and at most kd side wires in a test path. There are at most $|\Sigma|^{O(kd)}$ settings of these wires. If we randomly assign one of $\Sigma \cup \{*\}$ to each wire with equal probabilities, the

probability that we generate one of the settings is $1/|\Sigma|^{O(kd)}$. We can generate all settings using $|\Sigma|^{O(kd)} \log \frac{1}{\delta}$ random experiments, which succeeds with probability at least $1 - \delta$. We can also generate them deterministically using a universal set construction. The following definition of universal set is adapted from Seroussi and Bshouty [14]. An experiment set U is called (N, l) -universal if for every set of indices $R = \{r_1, r_2, \dots, r_l\} \subseteq [N]$, the projection of U to R contains all $(|\Sigma| + 1)^l$ l -tuples. It is shown in [13] that a (N, l) -universal set of size $2^{O(l \log |\Sigma|)} \log N$ can be constructed in polynomial time.³ Let U be a $(N, (d+1)(k+1))$ -universal set and let \mathcal{F} be the class of all gates of fan-in at most k . We show that U is sufficient for C and \mathcal{F} .

LEMMA 18. U is a sufficient set for C and \mathcal{F} .

PROOF. Let s be a witness experiment that g is wrong for w ; all g 's inputs are fixed in s . Let $s^* \preceq s$ be a minimal $(w, g(s))$ -exposing experiment. According to Lemma 17, there are at most d free wires and dk side wires in s . Since U is a universal set, there exists an experiment $s_0 \in U$ at the beginning of the algorithm, such that s_0 agrees with s^* in all s^* 's free wires and side wires and also all g 's inputs (there are at most $(d+1)(k+1)$ wires). Proposition 7 shows that s_0 is a witness experiment that g is wrong for w . \square

Whenever $kd = O(\log N)$, the size of U is polynomial in N and so is that of V . It takes time $N^{O(k)}$ to solve the consistency problem for a function of k variables, by checking every possible combination of k inputs. When k is a constant, this is polynomial.

THEOREM 19. *Log depth bounded fan-in circuits can be learned non-adaptively in polynomial time using VIQ's.*

5.4 Learning AC0 circuits

Theorem 6 precludes polynomial time algorithms for learning constant depth unbounded fan-in circuits with fairly simple gates. In this section, we show that if we allow only AND and OR gates (it is easy to extend it to NAND, NOR and NOT), constant depth unbounded fan-in Boolean circuits are learnable. Thus we show that **AC0** circuits are learnable with VIQ's.

We are not able to use a universal set, since k can be as large as $\Omega(N)$. Instead, we use Algorithm 2 to gather the necessary test paths adaptively. Algorithm 2 begins with learning the output gate, which can be easily done for AND and OR gates. It then sets one of its input wires free and fixes the other input wires so that the free input wire is still relevant. In particular, it sets the other input wires to 1 if the output gate is an AND gate, or 0 if the output gate is an OR gate. This partial experiment is then used to find (some of) the inputs of the corresponding gate. The algorithm goes on exploring the whole circuit. We remark that Algorithm 2 alone is not sufficient, because some input wires may be fixed as side wires and therefore hidden to the learner.

Let U contain all tests that are made in Algorithm 2 and \mathcal{F} be all AND and OR gates. The following lemma shows the correctness of the learning algorithm.

LEMMA 20. U is sufficient for C and \mathcal{F} .

³The paper [13] only deals with binary vectors. But it can be easily extended to the non-binary case by viewing each non-binary literal in $\Sigma \cup \{*\}$ as a binary vector of size $\log(|\Sigma|+1)$.

Algorithm 2 Gathering test paths for a constant depth AND/OR circuit

- 1: Let Γ contain the partial experiment that sets the output wire free, $\{w_N = *\}$.
 - 2: Let $\mathbf{1}(\mathbf{0})$ be an experiment that sets all wires to 1 (0),
 - 3: **while** Γ is not empty **do**
 - 4: Pick $\tau \in \Gamma$ and remove it from Γ .
 - 5: **if** $C(\mathbf{1}|\tau) \neq C(\mathbf{0}|\tau)$ **then**
 - 6: Let $Z = \{w \mid w \text{ unspecified in } \tau, \text{ and } C(\mathbf{1}|_{\tau, w=0}) \neq C(\mathbf{1}|\tau) \text{ or } C(\mathbf{0}|_{\tau, w=1}) \neq C(\mathbf{0}|\tau)\}$.
 - 7: **for** $w \in Z$ **do**
 - 8: **if** $C(\mathbf{1}|_{\tau, w=0}) \neq C(\mathbf{1}|\tau)$ **then**
 - 9: Add $\tau|_{w=*, \forall w' \in Z \setminus \{w\}, w'=1}$ to Γ . $\{\text{AND gate}\}$
 - 10: **else if** $C(\mathbf{0}|_{\tau, w=1}) \neq C(\mathbf{0}|\tau)$ **then**
 - 11: Add $\tau|_{w=*, \forall w' \in Z \setminus \{w\}, w'=0}$ to Γ . $\{\text{OR gate}\}$
-

PROOF. Suppose s is a witness experiment that g is wrong for w and $s^* \preceq s$ is a minimal $(w, g(s))$ -exposing experiment. Let u be the successor of w in the directed path from w to w_N (Lemma 17). We define two partial experiments τ_u and τ_w . τ_u sets all free wires in the directed path before u and their side wires as in s^* and sets u free. τ_w is similarly defined. We claim that τ_w is added to Γ in Algorithm 2. We assume inductively τ_u has been added to Γ .

Compare τ_u and τ_w . Those wires unspecified by τ_u but specified by τ_w are side wires that are inputs only to u . They are set to 1 if u is an AND gate and 0 if u is an OR gate so as to keep w relevant. Furthermore, we observe that

1. If u is an AND gate, $C(\mathbf{1}|\tau_u) \neq C(\mathbf{0}|\tau_u)$ and $C(\mathbf{1}|\tau_u, w=0) \neq C(\mathbf{1}|\tau_u)$;
2. If u is an OR gate, $C(\mathbf{1}|\tau_u) \neq C(\mathbf{0}|\tau_u)$ and $C(\mathbf{0}|\tau_u, w=1) \neq C(\mathbf{0}|\tau_u)$.

Therefore, τ_w must be added. Thus U must contain the following sets of experiments $\{\mathbf{0}|\tau_w, \mathbf{1}|\tau_w\}$, $\{\mathbf{1}|\tau_w, w'=0 \mid w' \text{ unspecified in } \tau_w\}$, and $\{\mathbf{0}|\tau_w, w'=1 \mid w' \text{ unspecified in } \tau_w\}$. Let g^* be the gate w in the target circuit C . The two projected functions $g|\tau_w$ and $g^*|\tau_w$ (fixing some inputs of the functions) must be different, because otherwise it contradicts the fact that s^* is a witness experiment. By case analysis, we can show that there exists at least one of the above-mentioned experiments s_0 , such that $g(s_0) \neq g^*(s_0)$. This s_0 serves our purpose. \square

It is clear that each partial experiment collected by Algorithm 2 corresponds to a directed path in the circuit C . Thus the number of partial experiments is bounded by $O(N^d) = \text{poly}(N)$ when the depth d is a constant. The size of U and the number of tests are hence polynomially bounded. The theorem then follows from the fact that the consistency problem for AND/OR gates can be solved in polynomial time.

THEOREM 21. *AC0 circuits are learnable in polynomial time using VIQ's.*

6. LEARNING WITH EXPERIMENTS AND COUNTEREXAMPLES

BEQ's overcome the obstacles of Theorem 2 and Theorem 3, because the counterexample has to give away the

combination when an appropriate hypothesis circuit is presented. However, the result in Theorem 6 still applies. Assuming both VIQ's and BEQ's are available, we give polynomial time algorithms to learn both arbitrary constant fan-in circuits and AND/OR circuits with unbounded depth.

Both algorithms repeatedly make a BEQ on a candidate circuit C' until C' is behaviorally equivalent to the target circuit. Each counterexample s is processed to give a minimal counterexample $s^* \preceq s$ such that $C'(s^*) \neq C(s^*)$. This process, *Minimize*, can easily be done with $O(N)$ VIQ's. The minimal counterexample is then used in rebuilding the candidate circuit C' . As in the proof of Lemma 10, a minimal counterexample is a witness experiment that a candidate gate g is wrong for a wire w . Therefore, each counterexample eliminates at least one candidate gate for at least one wire. For constant fan-in circuits, this immediately leads to a learning algorithm, in which we use a refinement of CircuitBuilder, which instead of checking each gate with respect to U , just picks a gate that is not eliminated for the corresponding wire. The algorithm is polynomial because there are at most $|\mathcal{F}| = N^{O(k)}$ gates to eliminate.

THEOREM 22. *Bounded fan-in circuits are learnable in polynomial time using VIQ's and BEQ's.*

However, the same method does not work for AND/OR circuits, as $|\mathcal{F}|$ is exponential. But each minimal counterexample still proves a gate wrong for a wire w . Let us denote each counterexample by a pair indicating the outputs of the true gate and the proposed gate. A (1,0) counterexample eliminates the constant 0 gate and similarly a (0,1) counterexample eliminates the constant 1 gate. Counterexamples for AND/OR gates can be divided into two types, namely, *input removing* and *input demanding* counterexamples. For instance, if the proposed gate and the true gate are both AND gates, a (1,0) counterexample says that the 0-inputs (inputs that are set 0) of the proposed gate are not inputs of the true gate. Such a counterexample causes the removal of the 0-inputs from the potential input set and is called input removing. Let R_w^\wedge contain all inputs removed by input removing counterexamples for wire w .

A (0,1) counterexample implies that inputs of the proposed gate does not include all inputs of the true gate. Such a counterexample is called input demanding. Let T_w^\wedge be a collection of sets of wires. We add the set of inputs of the proposed gate to T_w^\wedge when an input demanding counterexample is received. Any AND gate whose inputs are completely contained in any set of T_w^\wedge can not be the true gate and is therefore eliminated. We can view T_w^\wedge as constraints on candidate gates for wire w . An analogous argument applies when both gates are OR gates. Let R_w^\vee be analogous to R_w^\wedge and T_w^\vee be analogous to T_w^\wedge .

When the true gate and the proposed gate are of different types, we will process each counterexample as if they were the same type. It can be verified that even when the two gates are of different types, each counterexample will either add some wires to R_w^\wedge or R_w^\vee or add a set to T_w^\wedge or T_w^\vee and therefore make some progress. For example, if the true gate is an OR gate and the proposed gate is an AND gate, we treat the counterexample as if the true gate was an AND gate. Thus, we will add 0-inputs of the proposed gate to R_w^\wedge upon receiving a (1,0) counterexample, and add the whole set of inputs of the proposed gate to T_w^\wedge upon receiving a (0,1) counterexample. An important fact is that *the true*

gate will never be eliminated. This is because only when the proposed gate is of the same type as the true gate will there be any restriction on the true gate, and the true gate will not be eliminated in these cases as we have already discussed.

Algorithm 3 Building the proposed circuit

INPUT: $\forall w \in W, R_w^\wedge, T_w^\wedge, R_w^\vee, T_w^\vee$ and the set of constant functions \mathcal{C}_w for wire w that are not eliminated.

OUTPUT: C' .

- 1: $C' \leftarrow \emptyset, Z \leftarrow (w_1, w_2, \dots, w_N)$.
 - 2: **while** Z is not empty **do**
 - 3: Pop the first wire w in Z .
 - 4: **if** $\mathcal{C}_w \neq \emptyset$ **then**
 - 5: Add w to C' with one of the functions in \mathcal{C}_w .
 - 6: **else if** $\forall t \in T_w^\wedge, C' \setminus R_w^\wedge \not\subseteq t$ **then**
 - 7: Add w to C' with AND of wires in $C' \setminus R_w^\wedge$.
 - 8: **else if** $\forall t \in T_w^\vee, C' \setminus R_w^\vee \not\subseteq t$ **then**
 - 9: Add w to C' with OR of wires in $C' \setminus R_w^\vee$.
 - 10: **else**
 - 11: Put w at the end of Z .
-

We use Algorithm 3 to build the proposed circuit C' . At each iteration, we try to add a wire in Z to C' . We put Z in a queue and let the initial order be w_1, w_2, \dots, w_N . At each iteration, the first wire in the queue will be considered. If it is not added to C' , it will be put at the end of the queue. We add the wire to C' only if one of the following is true. One of the two constant functions is not eliminated; $C' \setminus R_w^\wedge$ is not contained in any set of T_w^\wedge ; $C' \setminus R_w^\vee$ is not contained in any set of T_w^\vee . We add the wire with a constant function, AND of wires in $C' \setminus R_w^\wedge$, or OR of wires in $C' \setminus R_w^\vee$, respectively.

Now we bound the number of counterexamples each wire can receive. There are at most 2 counterexamples that eliminate constant functions. There are at most $O(N)$ input removing counterexamples. The most subtle case is input demanding counterexamples. We identify the phase number of the learning algorithm with the number of counterexamples it has received (recall that the algorithm rebuilds C' each time a counterexample is received). In Algorithm 3, let *the round number of an iteration* be the number of times the wire being considered has been added to the queue. Let $I_w(t)$ be the round number of the iteration that w is added to C' at phase t . We will show that $I_w(t)$ will never decrease in the following. The intuition is that we add more constraints on w as we receive more counterexamples.

Let $C'_{(w,i)}(t)$ be the set C' when w is considered at round i in phase t . If we order pairs in $W \times [1, N]$ first by their round number and then by the order of wires in W , we have that $C'_{(w,i)}(t) = \{w' | (w', I_{w'}(t)) \leq (w, i)\}$. Together with R_w 's and T_w 's, $C'_{(w,i)}(t)$ decides whether w can be added to the circuit at round i in phase t . In the following, we show that $C'_{(w,i)}(t)$ never gets bigger. In other words, the set of wires C' that are available when w is considered to be added at round i is a subset of the corresponding set of wires in the previous phase. Therefore, if w is not added at round i in the previous phase, it is likely that it will not be added at round i at this phase. The key observation in the reasoning is that if a set C' cannot pass the test with R_w^\wedge and T_w^\wedge , and R_w^\vee and T_w^\vee ,

1. no subset of C' can pass the test;
2. C' cannot pass the test if we add more wires to R_w^\wedge and R_w^\vee or more sets to T_w^\wedge and T_w^\vee .

LEMMA 23. For any wire w , if $C'_{(w,i)}(t)$ exists and the algorithm does not end at phase t , $C'_{(w,i)}(t+1)$ exists and $C'_{(w,i)}(t+1) \subseteq C'_{(w,i)}(t)$.

PROOF. We do induction on the pairs (w, i) . The lemma clearly holds for $(w_1, 1)$ as $C'_{(w_1,1)}(t)$ is always empty. Suppose it holds for all pairs that precede (w, i) . Suppose there exists a wire w' in $C'_{(w,i)}(t+1) \setminus C'_{(w,i)}(t)$. We have that $(w', j = I_{w'}(t+1)) \leq (w, i)$. Therefore, w' is added at the j^{th} round at phase $t+1$ but after the j^{th} round at phase t . This implies (using the observations 1 and 2) that $C'_{(w',j)}(t+1) \not\subseteq C'_{(w',j)}(t)$, because R_w^\wedge and R_w^\vee , and T_w^\wedge and T_w^\vee only grow when more counterexamples are received. This leads to a contradiction. \square

It follows (also using the observations 1 and 2) that

COROLLARY 24. $I_w(t)$ is non-decreasing.

Now let us bound $I_w(t)$. Recall that the true gate will never be eliminated. Therefore, whenever C' contains all inputs of the true gate, w will be added. Thus, if the true gate of w is a constant gate, $I_w(t) = 1$. If the true gate of w depends only on constant gates, $I_w(t) \leq 2$. In general, we have that $I_w(t) \leq \max(I_{w_{i_1}}(t), I_{w_{i_2}}(t), \dots, I_{w_{i_k}}(t)) + 1$, where w_{i_j} ($j \in [1, k]$) is an input of a true gate of w . Therefore, $I_w(t) \leq d$.

Each input demanding counterexample for wire w at phase $t+1$ will eliminate the gate that Algorithm 3 picked for w at round $I_w(t)$ in phase t , by adding all inputs of the gate to T_w^\wedge or T_w^\vee . By Lemma 23 and the observations 1 and 2, this means that at phase $t+1$, we cannot pick gates of the same type again for w at or before round $I_w(t)$. Thus $I_w(t)$ will increase when w receives at most 2 input demanding counterexamples, one for AND gates and the other for OR gates. Therefore, we can bound the number of input demanding counterexamples by $O(d)$ per wire.

THEOREM 25. AND/OR/NOT circuits with unbounded fan-in and unbounded depth are learnable in polynomial time using VIQ's and BEQ's.

7. OTHER RESULTS

We also consider an extension of our model, namely, the *synchronous model*, where time is quantized and we can inject values as well as observe the output of the circuit at each time step. The circuits may be cyclic in this new model. A generalization of the methods of this paper shows that the classes of bounded fan-in circuits and AND/OR/NOT circuits are also learnable in polynomial time in the synchronous model.

8. REFERENCES

[1] T. Akutsu, S. Kuhara, O. Maruyama, and S. Miyano. Identification of gene regulatory networks by strategic gene disruptions and gene overexpressions. In *SODA '98: Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 695–702, Philadelphia, PA, USA, 1998. Society for Industrial and Applied Mathematics.

[2] D. Angluin, M. Frazier, and L. Pitt. Learning conjunctions of Horn clauses. *Machine Learning*, 9:147–164, 1992.

[3] D. Angluin, L. Hellerstein, and M. Karpinski. Learning read-once formulas with queries. *J. ACM*, 40:185–210, 1993.

[4] D. Angluin and M. Kharitonov. When won't membership queries help? *J. Comput. Syst. Sci.*, 50(2):336–355, 1995.

[5] N. H. Bshouty. Exact learning boolean functions via the monotone theory. *Inf. Comput.*, 123(1):146–153, 1995.

[6] H. Fujiwara. *Logic Testing and Design for Testability*. MIT Press, 1986.

[7] T. Ideker, V. Thorsson, and R. Karp. Discovery of regulatory interactions through perturbation: Inference and experimental design. In *Pacific Symposium on Biocomputing 5*, pages 302–313, 2000.

[8] J. C. Jackson. An efficient membership-query algorithm for learning DNF with respect to the uniform distribution. *J. Comput. Syst. Sci.*, 55(3):414–440, 1997.

[9] J. C. Jackson, A. R. Klivans, and R. A. Servedio. Learnability beyond AC0. In *STOC '02: Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 776–784, New York, NY, USA, 2002. ACM Press.

[10] M. Kearns and L. Valiant. Cryptographic limitations on learning boolean formulae and finite automata. *J. ACM*, 41(1):67–95, 1994.

[11] M. Kharitonov. Cryptographic hardness of distribution-specific learning. In *STOC '93: Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*, pages 372–381, New York, NY, USA, 1993. ACM Press.

[12] N. Linial, Y. Mansour, and N. Nisan. Constant depth circuits, Fourier transform, and learnability. *Journal of the ACM*, 40(3):607–620, 1993.

[13] J. Naor and M. Naor. Small-bias probability spaces: Efficient constructions and applications. In *ACM Symposium on Theory of Computing*, pages 213–223, 1990.

[14] G. Seroussi and N. H. Bshouty. Vector sets for exhaustive testings of logic circuits. In *IEEE Transactions on Information Theory*, volume 34, pages 513–522, May 1988.

[15] L. G. Valiant. A theory of the learnable. *Commun. ACM*, 27:1134–1142, 1984.