

# Ehrenfeucht-Fraïssé Games on Random Structures

Benjamin Rossman\*

Massachusetts Institute of Technology, Cambridge MA 02139, USA

**Abstract.** Certain results in circuit complexity (e.g., the theorem that  $AC^0$  functions have low average sensitivity [5, 17]) imply the existence of winning strategies in Ehrenfeucht-Fraïssé games on pairs of random structures (e.g., ordered random graphs  $G = G(n, 1/2)$  and  $G^+ = G \cup \{\text{random edge}\}$ ). Standard probabilistic methods in circuit complexity (e.g., the Switching Lemma [11] or Razborov-Smolensky Method [19, 21]), however, give no information about how a winning strategy might look. In this paper, we attempt to identify specific winning strategies in these games (as explicitly as possible). For random structures  $G$  and  $G^+$ , we prove that the *composition of minimal strategies* in  $r$ -round Ehrenfeucht-Fraïssé games  $\mathfrak{D}_r(G, G)$  and  $\mathfrak{D}_r(G^+, G^+)$  is almost surely a winning strategy in the game  $\mathfrak{D}_r(G, G^+)$ . We also examine a result of [20] that ordered random graphs  $H = G(n, p)$  and  $H^+ = H \cup \{\text{random } k\text{-clique}\}$  with  $p(n) \ll n^{-2/(k-1)}$  (below the  $k$ -clique threshold) are almost surely indistinguishable by  $\lfloor k/4 \rfloor$ -variable first-order sentences of any fixed quantifier-rank  $r$ . We describe a winning strategy in the corresponding  $r$ -round  $\lfloor k/4 \rfloor$ -pebble game using a technique that combines strategies from several auxiliary games.

## 1 Introduction

Let  $\mathcal{A}$  be an arbitrary finite structure, let  $P$  be a uniform random subset of  $A$ , let  $q$  be a uniform random element of  $A$ , and let  $P' = P \Delta \{q\}$ . Let  $(\mathcal{A}, P)$  and  $(\mathcal{A}, P')$  denote the expansions of  $\mathcal{A}$  by a new unary relation symbol interpreted as  $P$  and  $P'$ , respectively. Results in circuit complexity [1, 8] and descriptive complexity [12, 10] imply that structures  $(\mathcal{A}, P)$  and  $(\mathcal{A}, P')$  almost surely satisfy the same first-order sentences up to a fixed quantifier-rank  $r$  (as the size of  $\mathcal{A}$  increases). Equivalently, there exists a winning strategy (for “Duplicator”) in the  $r$ -round Ehrenfeucht-Fraïssé game on these structures, which we denote by  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$ .

Standard proofs of this fact via circuit complexity merely show that a winning strategy must exist using probabilistic arguments (generally based on either the Switching Lemma [4, 11] or Razborov-Smolensky Method [19, 21]). These proofs, however, say nothing about what a winning strategy might actually look like. In this paper, we aim for an explicit characterization of a winning strategy in  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$ . Our result involves two natural notions concerning strategies:

---

\* Supported by a National Defense Science and Engineering Graduate Fellowship.

**composition** (Definition 6) given strategies  $\mathbf{S}_1$  and  $\mathbf{S}_2$  in games  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B}_1)$  and  $\mathfrak{D}_r(\mathcal{B}_2, \mathcal{C})$  (where structures  $\mathcal{B}_1$  and  $\mathcal{B}_2$  have a common universe), there is a natural composition strategy  $\mathbf{S}_1 \circ \mathbf{S}_2$  in the game  $\mathfrak{D}_r(\mathcal{A}, \mathcal{C})$

**minimal strategy** (Definition 7) for every structure  $\mathcal{A}$  with universe  $\{1, \dots, n\}$ , there is a canonical (lexicographically) minimal winning strategy in the game  $\mathfrak{D}_r(\mathcal{A}, \mathcal{A})$ .

We show (Theorem 8) that the composition  $\mathbf{S} \circ \mathbf{S}'$  is almost surely a winning strategy in the game  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$  where  $\mathbf{S}$  and  $\mathbf{S}'$  are the minimal winning strategies in games  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P))$  and  $\mathfrak{D}_r((\mathcal{A}, P'), (\mathcal{A}, P'))$ .

A second result of this paper (Proposition 13) gives a criterion for (constructively) establishing the existence of winning strategies in  $r$ -round  $k$ -pebble games among a family of  $k$ -tupled structures with a common universe. We explain how this criterion is applied in [20] to prove that  $k$ -CLIQUE is not definable in  $\lfloor k/4 \rfloor$ -variable first-order logic. This criterion involves a novel technique for combining games across multiple structures.

*Organization of the paper.* In §2, we give the relevant definitions of structures, games and strategies (including minimal strategies and composition of strategies). In §3, we identify a particular winning strategy in the game  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$ . In §4, we present a criterion for  $\equiv_r^k$ -equivalence. In §5, we explain how this criterion is used to prove that  $k$ -CLIQUE is not definable in  $\lfloor k/4 \rfloor$ -variable first-order logic. In §6, we study this criterion in the setting of games on Kripke structures. Two appendices (§A, §B) contains proofs that do not fit in the body of the paper.

## 2 Preliminaries: Structures, Games and Strategies

Throughout this paper, *structures* are finite relational structures. We assume that the universe of every structure is equipped with a linear order (not necessarily occurring as a relation in the structure); without loss of generality, every structure has universe  $\{1, \dots, n\}$  for some natural number  $n$ .

**Definition 1** (STRUCTURES). *We now make this precise.* A relational signature  $\sigma$  consists of finitely many “relation symbols”  $R_1, \dots, R_t$  with associated “arities”  $r_1, \dots, r_t \geq 1$ . A  $\sigma$ -structure  $\mathcal{A}$  consists of a set  $A$  (called the universe of  $\mathcal{A}$ ) together with interpretations for all relation symbols in  $\sigma$ , i.e., relations  $R_i^{\mathcal{A}} \subseteq A^{r_i}$  for  $i = 1, \dots, t$ . We denote the universe of  $\mathcal{A}$  by  $A$ , the universe of  $\mathcal{B}$  by  $B$ , etc. For  $S \subseteq A$ , the induced substructure of  $\mathcal{A}$  with universe  $S$  is denoted by  $\mathcal{A}|_S$ .

By default, all structures are  $\sigma$ -structure for an arbitrary fixed relational signature  $\sigma$ . If  $\mathcal{A}$  is a structure and  $R \subseteq A^r$  is an  $r$ -ary relation on  $A$ , then we denote by  $(\mathcal{A}, R)$  the  $(\sigma \cup \{\underline{R}\})$ -structure expanding  $\mathcal{A}$ , in which  $R$  interprets a new  $r$ -ary relational symbol  $\underline{R}$ .

A  $k$ -tupled structure is a pair  $(\mathcal{A}, \bar{a})$  where  $\bar{a} = (a_1, \dots, a_k) \in A^k$ . In particular, structures are 0-tupled structures.

**Definition 2** ( $\equiv_r$ - AND  $\equiv_r^k$ -EQUIVALENCE). For all  $k, r \in \mathbb{N}$ , equivalence relations  $\equiv_r$  and  $\equiv_r^k$  on the class of  $k$ -tupled structures are defined by the following induction. Let  $(\mathcal{A}, \bar{a})$  and  $(\mathcal{B}, \bar{b})$  be  $k$ -tupled structures.

- $(\mathcal{A}, \bar{a}) \equiv_0 (\mathcal{B}, \bar{b}) \iff (\mathcal{A}, \bar{a}) \equiv_0^k (\mathcal{B}, \bar{b}) \stackrel{\text{def}}{\iff} \{(a_1, b_1), \dots, (a_k, b_k)\}$  is a partial isomorphism between  $\mathcal{A}$  and  $\mathcal{B}$  (i.e., a bijective function between subsets of  $A$  and  $B$  that preserves all relations).

For  $r \geq 1$ ,

- $(\mathcal{A}, \bar{a}) \equiv_r (\mathcal{B}, \bar{b}) \stackrel{\text{def}}{\iff} \begin{cases} \forall a' \in A \exists b' \in B (\mathcal{A}, \bar{a}a') \equiv_{r-1} (\mathcal{B}, \bar{b}b'), \\ \forall b' \in B \exists a' \in A (\mathcal{A}, \bar{a}a') \equiv_{r-1} (\mathcal{B}, \bar{b}b'); \end{cases}$

(note: here  $\equiv_{r-1}$  is an equivalence relation on  $(k+1)$ -tupled structures)

- $(\mathcal{A}, \bar{a}) \equiv_r^k (\mathcal{B}, \bar{b}) \stackrel{\text{def}}{\iff} (\mathcal{A}, \bar{a}) \equiv_0^k (\mathcal{B}, \bar{b})$  and for all  $i \in [k]$ ,

$$\begin{aligned} \forall a' \in A \exists b' \in B (\mathcal{A}, a_1, \dots, a_{i-1}, a', a_{i+1}, \dots, a_k) &\equiv_{r-1}^k (\mathcal{B}, b_1, \dots, b_{i-1}, b', b_{i+1}, \dots, b_k), \\ \forall b' \in B \exists a' \in A (\mathcal{A}, a_1, \dots, a_{i-1}, a', a_{i+1}, \dots, a_k) &\equiv_{r-1}^k (\mathcal{B}, b_1, \dots, b_{i-1}, b', b_{i+1}, \dots, b_k). \end{aligned}$$

This induction is clearly well-founded: the definition of  $\equiv_r$  on  $k$ -tupled structures depends on the definition of  $\equiv_{r-1}$  on  $(k+1)$ -tupled structures, etc., which eventually depends on the (base case) definition of  $\equiv_0$  on  $(k+r)$ -tupled structures. Note that  $\equiv_{r+1}$  refines  $\equiv_r$ , which refines  $\equiv_r^k$ .

**Definition 3** (EHRENFEUCHT-FRAÏSSÉ GAME). The  $r$ -round Ehrenfeucht-Fraïssé game on structures  $\mathcal{A}$  and  $\mathcal{B}$ , denoted  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B})$ , is played as follows. There are two players, Spoiler (who wishes to establish that  $\mathcal{A}$  and  $\mathcal{B}$  are non-isomorphic) and Duplicator (who attempts to prevent Spoiler from achieving his goal). The “game board” consists of structures  $\mathcal{A}$  and  $\mathcal{B}$  and the “game pieces” are  $r$  pairs of pebbles  $(p_1, q_1), \dots, (p_r, q_r)$ . The game is played in a sequence of  $r$  rounds. In round  $i$  of the game, Spoiler picks a structure (either  $\mathcal{A}$  or  $\mathcal{B}$ ) and places pebble  $p_i$  on an element of his choice in that structure. Duplicator then replies by placing pebble  $q_i$  on an element of his choice in the other structure. After  $r$  rounds, the positions of pebbles  $p_1, \dots, p_k, q_1, \dots, q_k$  describe two  $r$ -tuples  $(a_1, \dots, a_r) \in A^r$  and  $(b_1, \dots, b_r) \in B^r$ . Duplicator is declared the winner if and only if  $\{(a_1, b_1), \dots, (a_r, b_r)\}$  is a partial isomorphism from  $\mathcal{A}$  to  $\mathcal{B}$  (i.e., an isomorphism from  $\mathcal{A}|_{\{a_1, \dots, a_r\}}$  to  $\mathcal{B}|_{\{b_1, \dots, b_r\}}$ ).

The  $r$ -round  $k$ -pebble game on  $k$ -tupled structures  $(\mathcal{A}, \bar{a})$  and  $(\mathcal{B}, \bar{b})$ , denoted  $\mathfrak{D}_r^k((\mathcal{A}, \bar{a}), (\mathcal{B}, \bar{b}))$ , is similar. However, in this game there are exactly  $k$  pairs of pebbles  $(\alpha_1, \beta_1), \dots, (\alpha_k, \beta_k)$ . At the start of the game, these pairs of pebbles rest on  $(a_1, b_1), \dots, (a_k, b_k)$ . In each round of the game, Spoiler picks an index  $j \in [k]$  and a structure (either  $\mathcal{A}$  or  $\mathcal{B}$ ) and moves the  $j$ th pebble in that structure to an element of his choice. Duplicator then moves the  $j$ th pebble in the other structure to an element of his choice. Duplicator is declared the winner if and only if for every  $i \in [r]$ , the set  $\{(\alpha_1^{(i)}, \beta_1^{(i)}), \dots, (\alpha_k^{(i)}, \beta_k^{(i)})\}$  is a partial isomorphism from  $\mathcal{A}$  to  $\mathcal{B}$  where  $\alpha_j^{(i)}, \beta_j^{(i)}$  denote the positions in  $A, B$  of pebbles  $\alpha_j, \beta_j$  after round  $i$  of the game.

We are interested in strategies in games  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B})$  and  $\mathfrak{D}_r^k((\mathcal{A}, \bar{a}), (\mathcal{B}, \bar{b}))$ . By “strategy” we always mean a *deterministic strategy for Duplicator*; we are never concerned with strategies for Spoiler, who we simply assume plays optimally. For instance, the statement “there exists a winning strategy” should be read as “there exists a (deterministic) winning strategy (for Duplicator)”.

To avoid redundancy, we present definitions and lemmas concerning the game  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B})$  only. It will be obvious how to adapt the corresponding definitions and lemmas to the  $k$ -pebble game  $\mathfrak{D}_r^k((\mathcal{A}, \bar{a}), (\mathcal{B}, \bar{b}))$ .

**Definition 4** (STRATEGY). *An  $r$ -round strategy on sets  $A$  and  $B$  is a function*

$$\mathbf{S} : \bigcup_{i=1}^r (A \sqcup B)^i \longrightarrow A \sqcup B$$

*such that  $\mathbf{S}(x_1, \dots, x_i) \in A \iff x_i \in B$  for all  $(x_1, \dots, x_i) \in (A \sqcup B)^i$ . The intention is that if  $x_1, \dots, x_i$  are Spoiler’s moves in the first  $i$  rounds (i.e., which elements of which structures he plays), then  $\mathbf{S}(x_1, \dots, x_i)$  is Duplicator’s response in round  $i$  under strategy  $\mathbf{S}$ .*

*We say that  $\mathbf{S}$  is a winning strategy in the game  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B})$  if Duplicator is guaranteed to win by playing according to  $\mathbf{S}$  (no matter how Spoiler plays).*

Note that we define strategies on pairs of sets  $A$  and  $B$ , rather than structures  $\mathcal{A}$  and  $\mathcal{B}$ . For structures  $\mathcal{A}_1, \mathcal{A}_2, \mathcal{B}_1, \mathcal{B}_2$  where  $A_1 = A_2$  and  $B_1 = B_2$ , the same strategy  $\mathbf{S}$  might thus be a winning strategy in  $\mathfrak{D}_r(\mathcal{A}_1, \mathcal{B}_1)$ , but not in  $\mathfrak{D}_r(\mathcal{A}_2, \mathcal{B}_2)$ .

The simplest example of a winning strategy is the “copycat” strategy in the  $r$ -round game on two copies of a single structure  $\mathcal{A}$ . Duplicator easily defeats Spoiler by always playing the same element in the opposite structure.

The next proposition (a basic fact in model theory, see e.g. [16]) connects games,  $\equiv_r^{(k)}$ -equivalence and logic.

**Proposition 5.** *The following are equivalent:*

- i.  $\mathcal{A} \equiv_r \mathcal{B}$  (resp.  $(\mathcal{A}, \bar{a}) \equiv_r^k (\mathcal{B}, \bar{b})$ );
- ii. *there exists a winning strategy in  $\mathfrak{D}_r(\mathcal{A}, \mathcal{B})$  (resp.  $\mathfrak{D}_r^k((\mathcal{A}, \bar{a}), (\mathcal{B}, \bar{b}))$ );*
- iii.  $\mathcal{A}$  and  $\mathcal{B}$  satisfy the same first-order sentences of quantifier-rank  $r$  (resp.  $(\mathcal{A}, \bar{a})$  and  $(\mathcal{B}, \bar{b})$  satisfy the same first-order formulas  $\phi(x_1, \dots, x_k)$  of quantifier-rank  $r$  in which every subformula has at most  $k$  free variables).

The following concept of the composition of strategies is fairly intuitive.

**Definition 6** (COMPOSITION OF STRATEGIES). *Let  $\mathbf{S}$  be an  $r$ -round strategy on sets  $A$  and  $B$ , and let  $\mathbf{T}$  be an  $r$ -round strategy on sets  $B$  and  $C$ . The composition  $\mathbf{S} \circ \mathbf{T}$  is an  $r$ -round strategy on  $A$  and  $C$  defined as follows. Given  $i \in \{1, \dots, r\}$  and  $(x_1, \dots, x_i) \in (A \sqcup C)^i$ , define  $(y_1, \dots, y_i) \in (A \sqcup B)^i$  and  $(z_1, \dots, z_i) \in (B \sqcup C)^i$  inductively for  $j = 1, \dots, i$  by*

- if  $x_j \in A$ , then  $y_j \triangleq x_j$  and  $z_j \triangleq \mathbf{T}(z_1, \dots, z_{j-1}, \mathbf{S}(y_1, \dots, y_j))$ ,
- if  $x_j \in C$ , then  $z_j \triangleq x_j$  and  $y_j \triangleq \mathbf{S}(y_1, \dots, y_{j-1}, \mathbf{T}(z_1, \dots, z_j))$ .

*Then  $(\mathbf{S} \circ \mathbf{T})(x_1, \dots, x_i) \triangleq \begin{cases} z_i & \text{if } x_i \in A, \\ y_i & \text{if } x_i \in C. \end{cases}$*

For two structures  $\mathcal{A}$  and  $\mathcal{B}$  with given well-orderings such that  $\mathcal{A} \equiv_r \mathcal{B}$ , there is a canonical way to define a winning strategy in the  $r$ -round game on  $\mathcal{A}$  and  $\mathcal{B}$ : let Duplicator always play the minimal winning move.

**Definition 7 (MINIMAL WINNING STRATEGY).** *Let  $\mathcal{A}$  and  $\mathcal{B}$  be structures with given well-orderings such that  $\mathcal{A} \equiv_r \mathcal{B}$ . The (lexicographically) minimal winning  $r$ -round strategy  $\mathbf{M}_r^{\mathcal{A}, \mathcal{B}} : \bigcup_{i=1}^r (A \sqcup B)^i \rightarrow (A \sqcup B)$  on  $\mathcal{A}$  and  $\mathcal{B}$  is defined inductively as follows. Let  $i \in \{1, \dots, r\}$  and assume  $\mathbf{M}_r^{\mathcal{A}, \mathcal{B}}$  has been defined on all sequences of length  $i-1$ . Consider any  $(x_1, \dots, x_i) \in (A \sqcup B)^i$ . For all  $j \in \{1, \dots, i-1\}$ , define  $a_j \in A$  and  $b_j \in B$  by*

$$a_j = \begin{cases} x_j & \text{if } x_j \in A, \\ \mathbf{M}_r^{\mathcal{A}, \mathcal{B}}(x_1, \dots, x_j) & \text{if } x_j \in B, \end{cases} \quad b_j = \begin{cases} x_j & \text{if } x_j \in B, \\ \mathbf{M}_r^{\mathcal{A}, \mathcal{B}}(x_1, \dots, x_j) & \text{if } x_j \in A. \end{cases}$$

Assuming that  $(\mathcal{A}, a_1, \dots, a_{i-1}) \equiv_{r-i+1} (\mathcal{B}, b_1, \dots, b_{i-1})$  (which is guaranteed by the induction), define  $\mathbf{M}_r^{\mathcal{A}, \mathcal{B}}(x_1, \dots, x_i) \in A \sqcup B$  as follows:

- if  $x_i \in A$ , then  $\mathbf{M}_r^{\mathcal{A}, \mathcal{B}}(x_1, \dots, x_i)$  is the minimal  $b_i \in B$  such that
 
$$(\mathcal{A}, a_1, \dots, a_{i-1}, x_i) \equiv_{r-i} (\mathcal{B}, b_1, \dots, b_{i-1}, b_i),$$
- if  $x_i \in B$ , then  $\mathbf{M}_r^{\mathcal{A}, \mathcal{B}}(x_1, \dots, x_i)$  is the minimal  $a_i \in A$  such that
 
$$(\mathcal{A}, a_1, \dots, a_{i-1}, a_i) \equiv_{r-i} (\mathcal{B}, b_1, \dots, b_{i-1}, x_i).$$

We write  $\mathbf{M}_r^{\mathcal{A}}$  (instead of  $\mathbf{M}_r^{\mathcal{A}, \mathcal{A}}$ ) for the minimal winning strategy in the  $r$ -round game on two copies of a single structure  $\mathcal{A}$ . (Note that  $\mathbf{M}_r^{\mathcal{A}}$  is generally not the “copycat” strategy.)

### 3 A Winning Strategy in $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$

We now return the result mentioned in §1 concerning a winning strategy in the game  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$ . As in §1, let  $\mathcal{A}$  be an arbitrary structure with a universe of size  $n$ , let  $P$  be a uniform random subset of  $A$ , let  $q$  be a uniform random element of  $A$ , and let  $P' = P \triangle \{q\}$ . Let  $(\mathcal{A}, P)$  and  $(\mathcal{A}, P')$  denote the expansions of  $\mathcal{A}$  by a new unary relation symbol interpreted as  $P$  and  $P'$ , respectively.

**Theorem 8.** *The composition of minimal strategies  $\mathbf{M}_r^{(\mathcal{A}, P)} \circ \mathbf{M}_r^{(\mathcal{A}, P')}$  is almost surely a winning strategy in the game  $\mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))$ . In fact,*

$$\Pr [\mathbf{M}_r^{(\mathcal{A}, P)} \circ \mathbf{M}_r^{(\mathcal{A}, P')} \text{ is not winning in } \mathfrak{D}_r((\mathcal{A}, P), (\mathcal{A}, P'))] \leq O((\log n)^{O(r)}/n).$$

The key notion in proving Theorem 8 is that of an  $\mathcal{A}$ -minimal  $r$ -tuple. We will show (Lemma 12) that the set of induced substructures on  $\mathcal{A}$ -minimal  $r$ -tuples contains the essential information about the strategy  $\mathbf{M}_r^{\mathcal{A}}$ .

**Definition 9 ( $\mathcal{A}$ -MINIMAL TUPLES).** *An  $r$ -tuple  $(a_1, \dots, a_r) \in A^r$  is  $\mathcal{A}$ -minimal if for all  $i \in \{1, \dots, r\}$ , there exists  $a' \in A$  such that  $a_i = \mathbf{M}_r^{\mathcal{A}}(a_1, \dots, a_{i-1}, a')$ .*

The next lemma is essentially a folklore result in model theory. (Lemma 10 is also valid when  $\mathcal{A}$  is an infinite structure with a given well-ordering.)

**Lemma 10.** *There exists a constant  $c = c(r, \sigma)$  (depending only on  $r$  and the signature  $\sigma$ ) such that for every structure  $\mathcal{A}$ , there exist  $\leq c$  distinct  $\mathcal{A}$ -minimal  $r$ -tuples.*

*Proof (sketch).* This lemma follows from the folklore fact that there are only finitely many  $\equiv_r$ -equivalence classes of  $k$ -tupled structures for all  $r, k \in \mathbb{N}$  (see [16]). (In fact, one can prove a bound of  $c \leq \prod_{j=1}^r c_j$  where  $c_j$  is the number of  $\equiv_{r-j}$ -equivalence classes of  $(k+j)$ -tupled structures.)

**Definition 11 (STRONG  $r$ -EQUIVALENCE).** *Structures  $\mathcal{A}_1$  and  $\mathcal{A}_2$  with a common universe  $A$  are strongly  $r$ -equivalent if for every  $(a_1, \dots, a_r) \in A^r$ ,*

- $(a_1, \dots, a_r)$  is  $\mathcal{A}_1$ -minimal if and only if it is  $\mathcal{A}_2$ -minimal, and
- if  $(a_1, \dots, a_r)$  is  $\mathcal{A}_1$ -minimal, then  $\mathcal{A}_1|_{\{a_1, \dots, a_r\}} = \mathcal{A}_2|_{\{a_1, \dots, a_r\}}$  (i.e., these induced substructures are identical).

*Alternatively, can replace these two conditions with the single (equivalent) condition that  $\mathcal{A}_1|_{\{a_1, \dots, a_r\}} = \mathcal{A}_2|_{\{a_1, \dots, a_r\}}$  whenever  $(a_1, \dots, a_r)$  is  $\mathcal{A}_1$ -minimal or  $\mathcal{A}_2$ -minimal.*

The next lemma characterizes strong equivalence via minimal strategies.

**Lemma 12.**  *$\mathcal{A}_1$  and  $\mathcal{A}_2$  are strongly  $r$ -equivalent if and only if  $\mathbf{M}_r^{\mathcal{A}_1} \circ \mathbf{M}_r^{\mathcal{A}_2}$  is a winning strategy in  $\mathcal{D}_r(\mathcal{A}_1, \mathcal{A}_2)$ .*

We omit the proof of Lemma 12, which follows easily from definitions. Note that Lemma 12 reduces Theorem 8 to the inequality

$$\Pr [(\mathcal{A}, P) \text{ and } (\mathcal{A}, P') \text{ are not strongly } r\text{-equivalent}] \leq O((\log n)^{O(r)}/n).$$

The remainder of the proof of Theorem 8 (which proves this inequality using a result from circuit complexity that  $\text{AC}^0$  functions have low average sensitivity) is given in Appendix A.

## 4 Criterion for $\equiv_r^k$ -Equivalence

We present a criterion for establishing  $\equiv_r^k$ -equivalences among a family of  $k$ -tupled structures with the same universe. The criterion speaks about a family  $\{\mathcal{G}_{\bar{v}}\}_{\bar{v} \in V^k}$  of structures indexed by  $k$ -tuples over a common universe  $V$ , plus an additional structure  $\mathcal{G}^*$  with universe  $V$ . We give two hypotheses which together imply that  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_r^k (\mathcal{G}^*, \bar{v})$  for every  $k$ -tuple  $\bar{v} \in V^k$ . In the next section, we describe how this condition is used to show that  $\ell$ -CLIQUE is not definable in  $[\ell/4]$ -variable first-order logic.

To state the criterion, we fix a set  $V$  and a symmetric binary relation  $\sim$  on  $V$  with finite diameter  $d$  (so that every two elements of  $V$  are connected by a  $\sim$ -path of length  $\leq d$ ). For  $k$ -tuples  $\bar{v}, \bar{w} \in V^k$ , let

$$\bar{v} \sim \bar{w} \stackrel{\text{def}}{\iff} \exists i \in [k] \text{ such that } v_i \sim w_i \text{ and } v_j = w_j \text{ for all } j \in [k] \setminus \{i\}.$$

(In other words, the binary relation  $\sim$  on  $k$ -tuples is the  $k$ th Cartesian power of binary relation  $\sim$  on  $V$ .)

**Proposition 13** (CRITERION FOR  $\equiv_r^k$ ). *Let  $\{\mathcal{G}_{\bar{v}}\}_{\bar{v} \in V^k}$  be a family of structures  $\mathcal{G}_{\bar{v}}$  with universe  $V$  (indexed by  $k$ -tuples  $\bar{v} \in V^k$ ) and let  $\mathcal{G}^*$  be another structure with universe  $V$ . Suppose that*

- i.  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_0 (\mathcal{G}^*, \bar{v})$  for all  $\bar{v} \in V^k$ , and
- ii.  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_{rd} (\mathcal{G}_{\bar{w}}, \bar{w})$  for all  $\bar{v}, \bar{w} \in V^k$  such that  $\bar{v} \sim \bar{w}$ .

Then  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_r^k (\mathcal{G}^*, \bar{v})$  for all  $\bar{v} \in V^k$ .

**Remark 14.** *Proposition 13 is valid with a weaker hypothesis in which  $\equiv_{rd}$  is replaced by  $\equiv_{rd}^k$ . It is even valid when  $V$  is infinite. However, Proposition 13 as stated is all that we require (for the application in §5).*

Proposition 13 can be proved by a very simple argument using the inductive definition of  $\equiv_{rd}$ . However, we prefer to understand Proposition 13 in light of the winning strategy it entails in the game  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ . An examination of the proof reveals a *deterministic strategy* for winning the game  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$  given black-box winning strategies in games  $\mathcal{D}_{rd}((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}_{\bar{w}}, \bar{w}))$  for all  $\bar{v}, \bar{w} \in V^k$  such that  $\bar{v} \sim \bar{w}$ . Under this strategy for  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ , Duplicator plays a sequence of simulated games  $\mathcal{D}_{rd}((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}_{\bar{w}}, \bar{w}))$ , in fact making  $\binom{r}{2}$  queries to these auxiliary black-box strategies.

We will explicate this strategy (and also prove Proposition 13) later in §6. First, let’s see a neat application of Proposition 13.

## 5 $\ell$ -Clique Requires $\lfloor \ell/4 \rfloor$ Variables

Fix any  $\ell \geq 2$  and let  $k = \lfloor \ell/4 \rfloor$ .

Let  $\mathcal{G} = (V, E, <)$  be an ordered Erdos-Renyi random graph where  $V = [n]$  and  $E$  is a random anti-reflexive symmetric binary relation on  $V$  which includes each pair of vertices independently with probability  $p(n) = n^{-2/(\ell-1.5)}$ . Note that  $p(n)$  is below the  $\ell$ -clique threshold, i.e.,  $\mathcal{G}$  is almost surely  $\ell$ -clique-free. Let  $A$  be a uniform random  $\ell$ -element subset of  $V$ , let  $\mathcal{G}^* = (\mathcal{G} \cup \text{clique on } A)$ , and let  $\bar{1} = (1, \dots, 1)$  denote the all 1’s  $k$ -tuple in  $V^k$ .

The following theorem is proved in [20].

**Theorem 15.** *For every  $r$ , it holds a.a.s. (asymptotically almost surely as  $n \rightarrow \infty$ ) that  $(\mathcal{G}, \bar{1}) \equiv_r^k (\mathcal{G}^*, \bar{1})$ .*

Since  $\mathcal{G}$  is almost surely  $\ell$ -clique-free, while  $\mathcal{G}^*$  contains an  $\ell$ -clique (with probability 1), it follows that  $\ell$ -clique is not definable in  $k$ -variable first-order logic.

We now explain how Theorem 15 is proved using Proposition 13. There are three steps.

*Step 1.* For a sufficiently small constant  $\varepsilon > 0$  (to be determined), we fix an arbitrary reflexive and symmetric binary relation  $\sim$  on  $V$  with degree  $\leq n^\varepsilon$  and diameter  $\leq 2/\varepsilon$  (e.g., a spanning tree).

*Step 2.* For all  $k$ -tuples  $\bar{v} \in V^k$ , define  $A_{\bar{v}} \subseteq A$  by

$$A_{\bar{v}} \triangleq \left\{ a \in A : \exists A' \subseteq A \text{ s.t. } \begin{array}{l} |A'| < 2k \text{ and} \\ (\mathcal{G} \cup \text{clique on } A') \not\equiv_{\lceil 2r/\varepsilon \rceil} (\mathcal{G} \cup \text{clique on } A' \setminus \{a\}) \end{array} \right\}$$

and let  $\mathcal{G}_{\bar{v}} \triangleq (\mathcal{G} \cup \text{clique on } A_{\bar{v}})$ .

*Step 3.* We show that for all  $\bar{v} \in V^k$ ,

$$\Pr[|A_{\bar{v}}| > 0] = o(1),$$

$$\Pr[|\bigcup_{\bar{w} \sim \bar{v}} A_{\bar{w}}| \geq 2k] = o(1/n^k).$$

Steps 1 and 2 are merely definitions. Step 3 (where all the work is done) is proved by a probabilistic argument with ingredients from circuit complexity (see [20] for the proof). However, the intuition behind Step 3 is not hard to understand. The edge probability  $p(n) = n^{-2/(\ell-1.5)}$  lies *below* the  $\ell$ -clique threshold  $\Theta(n^{-2/(\ell-1)})$ , but *above* the  $(\ell-1)$ -clique threshold  $\Theta(n^{-2/(\ell-2)})$ . In particular, for every  $\ell' < \ell$ , the random graph  $\mathcal{G}$  almost surely has *many*  $\ell'$ -cliques. In particular,  $\mathcal{G}$  almost surely has  $\omega(n^k)$  cliques of size  $2k$  ( $\approx \ell/2$ ).

For every  $k$ -tuple  $\bar{v} \in V^k$ , the set  $A_{\bar{v}}$  depends only on the  $\equiv_{\lceil 2r/\varepsilon \rceil}$ -equivalence classes of ordered graphs  $(\mathcal{G} \cup \text{clique on } A')$  for subsets  $A' \subseteq A$  of size  $|A'| < 2k$ . But because  $\mathcal{G}$  already contains huge numbers of cliques of size  $\leq 2k$ , the addition of a random  $\ell'$ -clique where  $\ell' \leq 2k$  is *unlikely* to change the  $\equiv_{\lceil 2r/\varepsilon \rceil}$ -equivalence class of  $(\mathcal{G}, \bar{v})$ : this boils down to the fact that  $\text{AC}^0$  functions have low average sensitivity (cp. Appendix 8). In fact,  $\Pr[|A_{\bar{v}}| \geq \ell']$  is roughly bounded by  $1/\mathbb{E}[\#\text{ of } \ell'\text{-cliques in } \mathcal{G}]$ . Thus, we get  $\Pr[|A_{\bar{v}}| > 0] = o(1)$  and  $\Pr[|A_{\bar{v}}| \geq 2k] < o(1/n^k)$ . Using the fact that  $\bar{v}$  has at most  $kn^\varepsilon$   $\sim$ -neighbors (and picking a sufficiently small constant  $\varepsilon$  in Step 1), we are able to show  $\Pr[|\bigcup_{\bar{w} \sim \bar{v}} A_{\bar{w}}| \geq 2k] \leq o(1/n^k)$ , proving Step 3.

Steps 1–3, together with Proposition 13, furnish a proof of Theorem 15. By Step 3, it holds almost surely that  $A_{\bar{1}} = \emptyset$  (hence  $\mathcal{G}_{\bar{1}} = \mathcal{G}$ ) and  $|\bigcup_{\bar{w} \sim \bar{v}} A_{\bar{w}}| < 2k$  for all  $\bar{v} \in V^k$  (taking a union bound over all  $k$ -tuples in  $V^k$ ). Given that these events hold, hypotheses (i) and (ii) of Proposition 13 follow directly from the definition of  $A_{\bar{v}}$  and  $\mathcal{G}_{\bar{v}}$  in Step 2. Therefore, by Proposition 13, we have  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_r^k (\mathcal{G}^*, \bar{v})$  for all  $\bar{v} \in V^k$ . In particular, we have  $(\mathcal{G}, \bar{1}) = (\mathcal{G}_{\bar{1}}, \bar{1}) \equiv_r^k (\mathcal{G}^*, \bar{1})$ , proving Theorem 15.

## 6 Winning Strategy Behind the $\equiv_r^k$ -Criterion (Prop. 13)

We now analyze the winning strategy in the game  $\mathfrak{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$  that arises from our criterion for  $\equiv_r^k$ -equivalence (Proposition 13). In order to more clearly

describe this strategy, we move over to the simpler setting of games on Kripke structures. In §6.1, we state and prove an analogous criterion for establishing  $\sim_m$ -equivalences among classes of Kripke structures with the same universe. In §6.2, we reprove this criterion in terms of games on Kripke structures to get a clear picture of the winning strategy that emerges. In Appendix B, we prove Proposition 13 from the analogous criterion on Kripke structures. In this way, we get a clear picture of the winning strategy in  $\mathfrak{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ .

For our purposes, we consider a restricted version of Kripke structures, namely colored directed graphs (digraphs) with a single distinguished vertex.

**Definition 16.** Let  $\mathcal{A} = (V, E)$  be a digraph. For  $a \in V$ , let  $N_{\mathcal{A}}(a)$ —or simply  $N(a)$  if  $\mathcal{A}$  is clear from context—denote the set of vertices  $a' \in V$  such that  $(a, a') \in E$ . A function  $f : V \rightarrow \mathbb{N}$  is called a coloring of  $\mathcal{A}$ . The pair  $(\mathcal{A}, f)$  is called a colored digraph. For  $a \in V$ , the triple  $(\mathcal{A}, f, a)$  is called a Kripke structure.

**Definition 17.** For every  $m \in \mathbb{N}$ , equivalence relation  $\approx_m$  on Kripke structures is defined inductively by

- $(\mathcal{A}, f, a) \approx_0 (\mathcal{B}, g, b) \stackrel{\text{def}}{\iff} f(a) = g(b)$ ;
- for  $m \geq 1$ ,  $(\mathcal{A}, f, a) \approx_m (\mathcal{B}, g, b) \stackrel{\text{def}}{\iff} f(a) = g(b)$  and

$$\begin{aligned} \forall a' \in N_{\mathcal{A}}(a) \exists b' \in N_{\mathcal{B}}(b) (\mathcal{A}, f, a') \approx_{m-1} (\mathcal{B}, g, b'), \\ \forall b' \in N_{\mathcal{B}}(b) \exists a' \in N_{\mathcal{A}}(a) (\mathcal{A}, f, a') \approx_{m-1} (\mathcal{B}, g, b'). \end{aligned}$$

**Remark 18.** The equivalence relation  $\approx_m$  characterizes indistinguishability of Kripke structures up to sentences of rank  $m$  in modal logic (cp. Proposition 5).

There is a simple characterization of  $\approx_m$ -equivalence in terms of appropriately defined games (cp. Definition 3).

**Definition 19.** For Kripke structures  $(\mathcal{A}, f, a)$  and  $(\mathcal{B}, g, b)$  the game  $\mathfrak{D}_m((\mathcal{A}, f, a), (\mathcal{B}, g, b))$  is defined as follows. There are two players, Spoiler and Duplicator. The “game board” consists of disjoint copies of digraphs  $\mathcal{A}$  and  $\mathcal{B}$ . There are two “game pieces”,  $\alpha$  and  $\beta$  which sit on vertices of  $\mathcal{A}$  and  $\mathcal{B}$ , respectively. Initially  $\alpha$  sits on  $a$  and  $\beta$  sits on  $b$ . The game itself takes place in a sequence of  $m$  rounds. In each round, Spoiler moves first. Spoiler’s move consists of choosing either  $\mathcal{A}$  or  $\mathcal{B}$  and moving the corresponding game piece forward along an edge. Duplicator then replies by moving the other game piece along an edge in the other graph. The “record” of a game is the sequence  $(a_0, b_0), \dots, (a_m, b_m)$  of positions before and after each of the  $m$  rounds (where  $(a_0, b_0) = (a, b)$  is the initial position of  $\alpha$  and  $\beta$ ). Duplicator is declared the winner iff  $f(a_i) = g(b_i)$  for all  $i \in \{0, \dots, m\}$ . Duplicator is said to have a winning strategy if he can play in a manner such that he wins no matter what moves Spoiler makes.

## 6.1 Criterion for $\approx_m$ -Equivalence

Analogous to Proposition 13, we have the following criterion for  $\approx_m$ -equivalence.

**Proposition 20.** *Let  $\mathcal{A} = (V, E)$  be a digraph, let  $\{f_a\}_{a \in V}$  be a family of colorings  $f_a : V \rightarrow \mathbb{N}$  indexed by vertices  $a$  of  $\mathcal{A}$ , and let  $g : V \rightarrow \mathbb{N}$  be an additional coloring of  $\mathcal{A}$ . Suppose that*

- i.  $f_a(a) = g(a)$  for all  $a \in V$ , and
- ii.  $(\mathcal{A}, f_a, b) \approx_m (\mathcal{A}, f_b, b)$  for all  $(a, b) \in E$ .

*Then  $(\mathcal{A}, f_a, a) \approx_m (\mathcal{A}, g, a)$  for all  $a \in V$ .*

*Proof.* Because we consider various Kripke structures with the same underlying graph, we simplify notation by abbreviating  $(\mathcal{A}, f_a, b)$  and  $(\mathcal{A}, g, a)$  as  $(f_a, b)$  and  $(g, a)$ , respectively.

The proof is a simple argument by induction: for  $\ell = 0, \dots, m$  we claim that  $(f_a, a) \approx_\ell (g, a)$  for every  $a \in V$ . The base case when  $\ell = 0$  we get by hypothesis (i). For the induction step, assume that  $\ell > 0$  and  $(f_a, a) \approx_{\ell-1} (g, a)$  for every  $a \in V$ . Consider any  $(a, b) \in E$ . We have  $(f_a, b) \approx_m (f_b, b)$  by hypothesis (ii). We have  $(f_b, b) \approx_{\ell-1} (g, b)$  by the induction hypothesis. Therefore, we have  $(f_a, b) \approx_{\ell-1} (g, b)$  by the fact that  $\approx_m$  refines  $\approx_{\ell-1}$  and by transitivity of  $\approx_{\ell-1}$ . We now check that

- $f_a(a) = g(a)$  (by hypothesis (i)),
- for all  $a' \in N(a)$  there exists  $\exists b' \in N(a)$  (namely  $b' = a'$ ) such that  $(f_a, a') \approx_{\ell-1} (g, b')$ ,
- for all  $b' \in N(a)$  there exists  $\exists a' \in N(a)$  (namely  $a' = b'$ ) such that  $(f_a, a') \approx_{\ell-1} (g, b')$ .

Thus we have  $(f_a, a) \approx_\ell (g, a)$  by definition of  $\approx_\ell$ .

## 6.2 Winning Strategy behind the $\approx_m$ -Criterion (Proposition 20)

We now unravel the induction in the proof of Proposition 20 in order to extract a winning strategy in the game  $\mathfrak{D}_m((f_a, a), (g, a))$ . Suppose that for every edge  $(a, b) \in E$ , we are given a black-box winning strategy  $\mathbf{S}(a, b)$  in the game  $\mathfrak{D}_m((f_a, b), (f_b, b))$ . Our goal is to devise a winning strategy—let's call it  $\mathbf{\Sigma}(a)$ —in the game  $\mathfrak{D}_m((f_a, a), (g, a))$ . Before describing the strategy  $\mathbf{\Sigma}(a)$  in an informal manner, we prove a key lemma that captures the basic idea behind  $\mathbf{\Sigma}(a)$ .

**Lemma 21.** *Assume the hypotheses of Proposition 20. Let  $\ell \in \{0, \dots, m\}$ , let  $x_0, \dots, x_\ell$  and  $y_0, \dots, y_\ell$  be vertices of  $\mathcal{A}$  such that  $x_\ell = y_\ell$  and*

$$(\dagger) \quad (f_{x_0}, y_0) \approx_{m-\ell} (f_{x_1}, y_1) \approx_{m-\ell+1} \cdots \approx_{m-1} (f_{x_{\ell-1}}, y_{\ell-1}) \approx_m (f_{x_\ell}, y_\ell).$$

*Then*

1.  $f_{x_0}(y_0) = g(y_\ell)$ ,
2. if  $\ell < m$ , then for every  $z_0 \in N(y_0)$  there exists  $(z_1, \dots, z_\ell) \in N(y_1) \times \cdots \times N(y_\ell)$  such that

$$(f_{x_0}, z_0) \approx_{m-\ell-1} (f_{x_1}, z_1) \approx_{m-\ell} \cdots \approx_{m-1} (f_{x_\ell}, z_\ell) \approx_m (f_{z_\ell}, z_\ell),$$

3. if  $\ell < m$ , then for every  $z_\ell \in N(y_\ell)$  there exists  $(z_0, \dots, z_{\ell-1}) \in N(y_0) \times \dots \times N(y_{\ell-1})$  such that

$$(f_{x_0}, z_0) \approx_{m-\ell-1} (f_{x_1}, z_1) \approx_{m-\ell} \dots \approx_{m-1} (f_{x_\ell}, z_\ell) \approx_m (f_{z_\ell}, z_\ell).$$

*Proof.* We first prove statement (1). Hypothesis  $(\dagger)$  implies that  $(f_{x_0}, y_0) \approx_{m-\ell} (f_{x_\ell}, y_\ell)$ . Therefore  $f_{x_0}(y_0) = f_{x_\ell}(y_\ell)$ . By hypothesis (i) of Proposition 20 we have  $f_{y_\ell}(y_\ell) = g(y_\ell)$ . Since  $x_\ell = y_\ell$  we get  $f_{x_0}(y_0) = g(y_\ell)$  as required.

For statement (2), assume  $\ell < m$  and let  $z_0 \in N(y_0)$ . The fact that  $(f_{x_0}, y_0) \approx_{m-\ell} (f_{x_1}, y_1)$  by  $(\dagger)$  implies there exists  $z_1 \in N(y_1)$  such that  $(f_{x_0}, z_0) \approx_{m-\ell-1} (f_{x_1}, z_1)$ . Next, the fact that  $(f_{x_1}, y_1) \approx_{m-\ell+1} (f_{x_2}, y_2)$  by  $(\dagger)$  implies there exists  $z_2 \in N(y_2)$  such that  $(f_{x_1}, z_1) \approx_{m-\ell} (f_{x_2}, z_2)$ . Continuing in this fashion we obtain a sequence  $(z_1, \dots, z_\ell) \in N(y_1) \times \dots \times N(y_\ell)$  such that

$$(f_{x_0}, z_0) \approx_{m-\ell-1} (f_{x_1}, z_1) \approx_{m-\ell} \dots \approx_{m-2} (f_{x_{\ell-1}}, z_{\ell-1}) \approx_{m-1} (f_{x_\ell}, z_\ell).$$

Finally, we have  $(f_{x_\ell}, z_\ell) \approx_m (f_{z_\ell}, z_\ell)$  by hypothesis (ii) of Proposition 20 since  $z_\ell \in N(y_\ell) = N(x_\ell)$ .

For statement (3), again assume  $\ell < m$  and this time let  $z_\ell \in N(y_\ell)$ . From  $(f_{x_{\ell-1}}, y_{\ell-1}) \approx_m (f_{x_\ell}, y_\ell)$ , it follows that there exists  $z_{\ell-1} \in N(y_{\ell-1})$  such that  $(f_{x_{\ell-1}}, z_{\ell-1}) \approx_{m-1} (f_{x_\ell}, z_\ell)$ . Continuing, we get a sequence  $(z_0, \dots, z_{\ell-1}) \in N(y_0) \times \dots \times N(y_{\ell-1})$  such that

$$(f_{x_0}, z_0) \approx_{m-\ell-1} (f_{x_1}, z_1) \approx_{m-\ell} \dots \approx_{m-2} (f_{x_{\ell-1}}, z_{\ell-1}) \approx_{m-1} (f_{x_\ell}, z_\ell).$$

As before, we get  $(f_{x_\ell}, z_\ell) \approx_m (f_{z_\ell}, z_\ell)$  by Proposition 20(ii) since  $z_\ell \in N(y_\ell) = N(x_\ell)$ .

Lemma 21 implicitly contains a winning strategy  $\Sigma(a)$  for Duplicator in the game  $\mathcal{D}_m((f_a, a), (g, a))$  for every  $a \in V$ , given a family  $\{\mathbf{S}(u, v)\}_{(u,v) \in E}$  of black-box winning strategies in games  $\mathcal{D}_m((f_u, v), (f_v, v))$  for all edges  $(u, v) \in E$ . Suppose  $(a_0, b_0), \dots, (a_\ell, b_\ell)$  is the sequence of positions before and after each of the first  $\ell$  rounds of  $\mathcal{D}_m((f_a, a), (g, a))$ . (In particular,  $(a_0, b_0) = (a, a)$  is the initial position.) Under the strategy  $\Sigma(a)$ , Duplicator maintains sequences  $x_0, \dots, x_\ell$  and  $y_0^{(\ell)}, \dots, y_\ell^{(\ell)}$  such that

- $x_0 = a$  and  $y_0^{(\ell)} = a_\ell$  and  $x_\ell = y_\ell^{(\ell)} = b_\ell$ ,
- $x_0, \dots, x_\ell$  is a path in  $A$ , and
- $(\dagger)$  holds, that is,

$$(f_{x_0}, y_0^{(\ell)}) \approx_{m-\ell} (f_{x_1}, y_1^{(\ell)}) \approx_{m-\ell+1} \dots \approx_{m-1} (f_{x_{\ell-1}}, y_{\ell-1}^{(\ell)}) \approx_m (f_{x_\ell}, y_\ell^{(\ell)}).$$

Assuming  $\ell < m$ , Duplicator plays under strategy  $\Sigma(a)$  in round  $\ell + 1$  as follows. First, suppose Spoiler plays  $a_{\ell+1} \in N(a_\ell)$  in the colored digraph  $(\mathcal{A}, f_a)$ . Then Duplicator sets  $z_0 = a_{\ell+1}$ , using strategies  $\mathbf{S}(x_0, x_1), \dots, \mathbf{S}(x_{\ell-1}, x_\ell)$  (note that  $(x_{i-1}, x_i) \in E$  for all  $1 \leq i \leq \ell$ ), constructs the sequence  $z_1, \dots, z_\ell$  as in Lemma 21(2). Duplicator then plays  $z_\ell \in N(b_\ell)$  in  $(\mathcal{A}, g)$  and updates his internal bookkeeping by setting  $x_{\ell+1} = z_\ell$  and  $y_0^{(\ell+1)} = z_0, \dots, y_\ell^{(\ell+1)} = z_\ell$  and  $y_{\ell+1}^{(\ell+1)} = z_\ell$ . If instead Spoiler plays  $b_{\ell+1} \in N(b_\ell)$  in  $(\mathcal{A}, g)$ , then Duplicator responds in a similar manner using Lemma 21(3).

## References

- [1] Ajtai, M.:  $\Sigma_1^1$  formulae on finite structures. *Annals of Pure and Applied Logic* 24, 1–48 (1983)
- [2] Amano, K., Maruoka, A.: A superpolynomial lower bound for a circuit computing the clique function with at most  $(1/6) \log \log n$  negation gates. *SIAM J. Comput.* 35(1), 201–215 (2005)
- [3] Beame, P.: Lower bounds for recognizing small cliques on CRCW PRAM's. *Discrete Appl. Math.* 29(1), 3–20 (1990)
- [4] Beame, P.: A switching lemma primer. Technical Report UW-CSE-95-07-01, Department of Computer Science and Engineering, University of Washington (November 1994)
- [5] Boppana, R.B.: The average sensitivity of bounded-depth circuits. *Inf. Process. Lett.* 63(5), 257–261 (1997)
- [6] Dawar, A.: How many first-order variables are needed on finite ordered structures? In: *We Will Show Them: Essays in Honour of Dov Gabbay*, pp. 489–520 (2005)
- [7] Denenberg, L., Gurevich, Y., Shelah, S.: Definability by constant-depth polynomial-size circuits. *Information and Control* 70(2/3), 216–240 (1986)
- [8] Furst, M.L., Saxe, J.B., Sipser, M.: Parity, circuits, and the polynomial-time hierarchy. *Mathematical Systems Theory* 17, 13–27 (1984)
- [9] Goldmann, M., Håstad, J.: A simple lower bound for the depth of monotone circuits computing clique using a communication game. *Information Processing Letters* 41(4), 221–226 (1992)
- [10] Gurevich, Y., Lewis, H.R.: A logic for constant-depth circuits. *Information and Control* 61(1), 65–74 (1984)
- [11] Håstad, J.: Almost optimal lower bounds for small depth circuits. In: *STOC 1986: Proceedings of the Eighteenth Annual ACM Symposium on Theory of Computing*, pp. 6–20 (1986)
- [12] Immerman, N.: Upper and lower bounds for first order expressibility. *J. Comput. Syst. Sci.* 25(1), 76–98 (1982)
- [13] Immerman, N.: *Descriptive Complexity*. In: *Graduate Texts in Computer Science*. Springer, New York (1999)
- [14] Immerman, N., Buss, J., Barrington, D.M.: Number of variables is equivalent to space. *Journal of Symbolic Logic* 66 (2001)
- [15] Koucky, M., Lautemann, C., Poloczek, S., Thérien, D.: Circuit lower bounds via Ehrenfeucht-Fraïssé games. In: *CCC 2006: Proceedings of the 21st Annual IEEE Conference on Computational Complexity*, pp. 190–201 (2006)
- [16] Libkin, L.: *Elements of Finite Model Theory*. Springer, Heidelberg (2004)
- [17] Linial, N., Mansour, Y., Nisan, N.: Constant depth circuits, fourier transform, and learnability. *J. ACM* 40(3), 607–620 (1993)
- [18] Lynch, J.F.: A depth-size tradeoff for boolean circuits with unbounded fan-in. In: *Structure in Complexity Theory Conference*, pp. 234–248 (1986)
- [19] Razborov, A.A.: Lower bounds on the size of bounded depth networks over a complete basis with logical addition. *Matematicheskie Zametki* 41, 598–607 (1987); English translation in *Mathematical Notes of the Academy of Sciences of the USSR* 41, 333–338 (1987) (in Russian)
- [20] Rossman, B.: On the constant-depth complexity of  $k$ -clique. In: *Proceedings of the 40th Annual ACM Symposium on Theory of Computing*, pp. 721–730 (2008)
- [21] Smolensky, R.: Algebraic methods in the theory of lower bounds for boolean circuit complexity. In: *Proceedings of the 19th Annual ACM Symposium on Theory of Computing*, pp. 77–82 (1987)

## A Proof of Theorem 8

We begin by defining circuits, the complexity class  $AC^0(\text{depth } d)$  and the average sensitivity of functions with domain  $\{0, 1\}^n$ .

**Definition 22.** A circuit on  $n$  inputs is an acyclic directed graph  $\mathcal{C}$  in which sources are labeled by elements of  $[n] \times \{+, -\}$  (where  $(i, +)$  corresponds to the literal  $x_i$  and  $(i, -)$  corresponds to the literal  $\bar{x}_i$ ) and all other nodes are labeled by  $\wedge$  or  $\vee$ . Sinks in  $\mathcal{C}$  are called outputs. The circuit  $\mathcal{C}$  computes a Boolean function  $\{0, 1\}^n \rightarrow \{0, 1\}^m$  where  $m = |\{\text{output nodes in } \mathcal{C}\}|$ .

**Definition 23.** Let  $\bar{\mathcal{C}} = (\mathcal{C}_n)_{n \in \mathbb{N}}$  be a sequence of circuits  $\mathcal{C}_n$  on  $n$  inputs. We say that  $\bar{\mathcal{C}} \in AC^0(\text{depth } d)$  if

- $\mathcal{C}_n$  has size  $n^{O(1)}$  and
- there exists a constant  $c$  such that for all  $n$ , every directed path in  $\mathcal{C}_n$  contains at most  $d$  nodes of fan-in (i.e., in-degree)  $> c$ .

**Definition 24.** For a function  $f$  with domain  $\{0, 1\}^n$  and an element  $x \in \{0, 1\}^n$ , the sensitivity of  $f$  at  $x$  is defined by

$$\text{sens}(f, x) \triangleq |\{i \in [n] : f(x) \neq f(x_1, \dots, x_{i-1}, 1 - x_i, x_{i+1}, \dots, x_n)\}|.$$

The average sensitivity of  $f$  is defined by

$$\text{ave-sens}(f) \triangleq 2^{-n} \sum_{x \in \{0, 1\}^n} \text{sens}(f, x).$$

The next lemma is a fundamental result in circuit complexity.

**Lemma 25 ([17, 5]).** Suppose  $f : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is computed by a circuit in  $AC^0(\text{depth } d)$ . Then  $\text{ave-sens}(f) = O(m \log^{d-1} n)$ .

We now make an observation about strong  $r$ -equivalence akin to well-known results in descriptive complexity characterizing first-order logic in terms of  $AC^0$  (see [7, 10, 12, 13]).

**Lemma 26.** Let  $\mathcal{A}$  be a structure with universe  $[n]$  (really, a sequence of structures, one for each natural number  $n$ ). There exists a function (really, a sequence of functions)  $f_{\mathcal{A}} : \{0, 1\}^n \rightarrow \{0, 1\}^{O(\log n)}$  computed by circuits in  $AC^0(\text{depth } O(r))$  such that for all  $P, Q \subseteq [n]$ , structures  $(\mathcal{A}, P)$  and  $(\mathcal{A}, Q)$  are strongly  $r$ -equivalent iff  $f_{\mathcal{A}}(P) = f_{\mathcal{A}}(Q)$  (where we view  $P$  and  $Q$  as elements of  $\{0, 1\}^n$ ).

That is,  $f_{\mathcal{A}}(P)$  completely describes the set of  $(\mathcal{A}, P)$ -minimal sequences  $(a_1, \dots, a_r)$  of length  $r$  as well as substructures  $(\mathcal{A}, P)|_{\{a_1, \dots, a_r\}}$ .

*Proof (sketch).* We use the fact that the number of minimal sequences of length  $r$  in any structure is bounded by an absolute constant (by Lemma 10). We can thus represent the set of all  $(\mathcal{A}, P)$ -minimal sequences  $(a_1, \dots, a_r)$  of length  $r$  as well as substructures  $(\mathcal{A}, P)|_{\{a_1, \dots, a_r\}}$  using only  $O(\log n)$  bits. This can be achieved by a polynomial-size circuit of depth  $O(r)$  using standard arguments (cp. [7, 10, 12, 13]).

Using Lemmas 25 and 26, we easily prove Theorem 4. As remarked at the end of §3, it suffices to prove

$$\Pr [(\mathcal{A}, P) \text{ and } (\mathcal{A}, P') \text{ are not strongly } r\text{-equivalent}] \leq O((\log n)^{O(r)}/n)$$

where  $\mathcal{A}$  is an arbitrary structure with universe  $[n]$ ,  $P$  is a uniform random subset of  $[n]$ , and  $P' = P \Delta \{q\}$  where  $q$  is a uniform random element of  $[n]$ . But the statement that  $(\mathcal{A}, P)$  and  $(\mathcal{A}, P')$  are not strongly  $r$ -equivalent is equivalent to  $f_{\mathcal{A}}(P) \neq f_{\mathcal{A}}(P')$ . Completing the proof, we have

$$\Pr[f_{\mathcal{A}}(P) \neq f_{\mathcal{A}}(P')] = E[\text{sens}(f, P)/n] = \text{ave-sens}(f)/n = O((\log n)^{O(r)}/n).$$

### B Proof of Proposition 13

In order to prove Proposition 13 (the criterion for  $\equiv_r^k$ -equivalence), we first state a suitable generalization of Proposition 20 (the criterion for  $\approx_r$ -equivalence of Kripke structures). Rather than digraphs, we consider  $k$ -digraphs  $\mathcal{A} = (V, E_1, \dots, E_k)$  where each  $E_i$  is a binary relation on  $V$ . For  $a \in V$  and  $i \in [k]$ , let  $N_{\mathcal{A},i}(a)$  for the neighbor set of  $a$  under edge relation  $E_i$  and let  $N_{\mathcal{A}}(a) = N_{\mathcal{A},1}(a) \cup \dots \cup N_{\mathcal{A},k}(a)$ .

A  $k$ -Kripke structure is a triple  $(\mathcal{A}, f, a)$  where  $\mathcal{A}$  is a  $k$ -digraph and  $f : V \rightarrow \mathbb{N}$  is a coloring of  $A$  and  $a \in V$  is a distinguished vertex.  $\approx_r$ -equivalence of  $k$ -Kripke structures is defined as follows (cf. Definition 2 of  $\equiv_r^k$ -equivalence):

- $(\mathcal{A}, f, a) \approx_0 (\mathcal{B}, g, b) \stackrel{\text{def}}{\iff} f(a) = g(b)$ ;
- for  $r \geq 1$ ,  $(\mathcal{A}, f, a) \approx_r (\mathcal{B}, g, b) \stackrel{\text{def}}{\iff} f(a) = g(b)$  and for all  $i \in [k]$ ,
 
$$\forall a' \in N_{\mathcal{A},i}(a) \exists b' \in N_{\mathcal{B},i}(b) (\mathcal{A}, f, a') \approx_{r-1} (\mathcal{B}, g, b'),$$

$$\forall b' \in N_{\mathcal{B},i}(b) \exists a' \in N_{\mathcal{A},i}(a) (\mathcal{A}, f, a') \approx_{r-1} (\mathcal{B}, g, b').$$

We have the following generalization of Proposition 20 to  $k$ -Kripke structures.

**Proposition 27.** *Let  $\mathcal{A} = (V, E_1, \dots, E_k)$  be a  $k$ -digraph, let  $\{f_a\}_{a \in V}$  be a family of colorings  $f_a : V \rightarrow \mathbb{N}$  indexed by vertices  $a \in V$ , and let  $g : V \rightarrow \mathbb{N}$  be another coloring of  $A$ . Suppose that*

- i'.  $f_a(a) = g(a)$  for all  $a \in V$ , and
- ii'.  $(\mathcal{A}, f_a, b) \approx_r (\mathcal{A}, f_b, b)$  for all  $(a, b) \in E_1 \cup \dots \cup E_k$ .

*Then  $(\mathcal{A}, f_a, a) \approx_r (\mathcal{A}, g, a)$  for all  $a \in V$ .*

The proof of Proposition 27 is virtually identical to the proof of Proposition 20 in §6.1. The game version of this proof given in §6.2 likewise generalizes to the  $k$ -Kripke setting.

We are ready to derive Proposition 13 from Proposition 27. Let  $V, \sim, \{\mathcal{G}_{\bar{v}}\}_{\bar{v} \in V^k}$  and  $\mathcal{G}^*$  be as in Proposition 13. Assume that hypotheses (i) and (ii) hold:

- i.  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_0 (\mathcal{G}^*, \bar{v})$  for all  $\bar{v} \in V^k$ , and
- ii.  $(\mathcal{G}_{\bar{v}}, \bar{v}) \equiv_{rd} (\mathcal{G}_{\bar{w}}, \bar{v})$  for all  $\bar{v}, \bar{w} \in V^k$  such that  $\bar{v} \sim \bar{w}$ .

Consider the  $k$ -digraph  $\mathcal{A} = (V^k, E_1, \dots, E_k)$  with vertex set  $V^k$  and edge relations

$$E_i = \{(\bar{v}, \bar{w}) \in V^k \times V^k : v_i \sim w_i \text{ and } v_j = w_j \text{ for all } j \in [k] \setminus \{i\}\}.$$

Fix an arbitrary enumeration of the (finitely many)  $\equiv_0$ -equivalence classes of  $k$ -tupled structures. For all  $\bar{v} \in V^k$ , define  $f_{\bar{v}} : V^k \rightarrow \mathbb{N}$ , and also define  $g : V^k \rightarrow \mathbb{N}$ , by

$$f_{\bar{v}}(\bar{w}) = \text{index of the } \equiv_0\text{-equivalence class of } (\mathcal{G}_{\bar{v}}, \bar{w}),$$

$$g(\bar{v}) = \text{index of the } \equiv_0\text{-equivalence class of } (\mathcal{G}^*, \bar{v}).$$

Check that

- $f_{\bar{v}}(\bar{v}) = g(\bar{v})$  for all  $\bar{v} \in V^k$  (by hypothesis (i)),
- $f_{\bar{v}}(\bar{w}) \approx_{rd} f_{\bar{w}}(\bar{w})$  for all  $(\bar{v}, \bar{w}) \in E_1 \cup \dots \cup E_k$  (using hypothesis (ii)).

Therefore, by Proposition 27 we have  $(\mathcal{A}, f_{\bar{v}}, \bar{v}) \approx_{rd} (\mathcal{A}, g, \bar{v})$  for all  $\bar{v} \in V^k$ .

But what exactly does  $(\mathcal{A}, f_{\bar{v}}, \bar{v}) \approx_{rd} (\mathcal{A}, g, \bar{v})$  mean? It is easy to see that this is equivalent to the existence of a winning strategy in the following  $\sim$ -constrained  $rd$ -round  $k$ -pebble game, which we denote  $\tilde{\mathcal{D}}_{rd}^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ : the  $\sim$ -constrained game is just like the usual  $rd$ -round  $k$ -pebble game, except that when either player moves a pebble from its present location on an element  $v$ , he is required to place that pebble on an element  $w$  such that  $v \sim w$ . To complete the proof of Proposition 13, we claim that a winning strategy in  $\tilde{\mathcal{D}}_{rd}^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$  implies a winning strategy in the usual  $r$ -round  $k$ -pebble game  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ . Here we exploit the fact that  $\sim$  has diameter  $\leq d$ , which lets us map each move in  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$  to a sequence of  $\leq d$  moves in  $\tilde{\mathcal{D}}_{rd}^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ . We thus get a winning strategy by playing  $d$  moves in a simulation of  $\tilde{\mathcal{D}}_{rd}^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$  for each move in the game  $\mathcal{D}_r^k((\mathcal{G}_{\bar{v}}, \bar{v}), (\mathcal{G}^*, \bar{v}))$ .