# Chapter 6

# Ehrenfeucht-Fraïssé Games

*We introduce combinatorial games that are useful for determining what can be expressed in various logics. Ehrenfeucht-Fraïssé games offer a semantics for first-order logic that is equivalent to, but more directly applicable than, the standard definitions.*

## 6.1 Definition of the Games

Suppose we assert that structure $\mathcal{A}$ satisfies the formula $(\forall x)\varphi(x)$. This can be understood in a game-theoretic way: an opponent may choose any element $a \in |\mathcal{A}|$. We are then obliged to show that $\mathcal{A} \models \varphi(a)$. Similarly, if we assert $(\exists y)\psi(y)$, then we move first, choosing $b \in |\mathcal{A}|$ and asserting that $\mathcal{A} \models \psi(b)$.

This operational, game-theoretic view of logical assertions is the subject of the present chapter. Ehrenfeucht-Fraïssé games offer a convenient, model-theoretic approach to logic. These games have been used extensively for proving that certain queries are not expressible in certain logics.

Using games, we present a complete methodology for proving that a boolean query is not expressible in first-order logic. We provide many examples. In later chapters we will show how to modify the games to provide complete methodologies for other languages, stronger than first-order logic.

We begin by defining the games, which were invented by Ehrenfeucht and Fraïssé.

**Definition 6.2**  The game $\mathcal{G}_m^k$ is played by two players: Samson and Delilah on a pair of structures $\mathcal{A}$ and $\mathcal{B}$ of the same vocabulary, $\tau$. $\mathcal{G}_m^k$ is played for $m$ rounds,
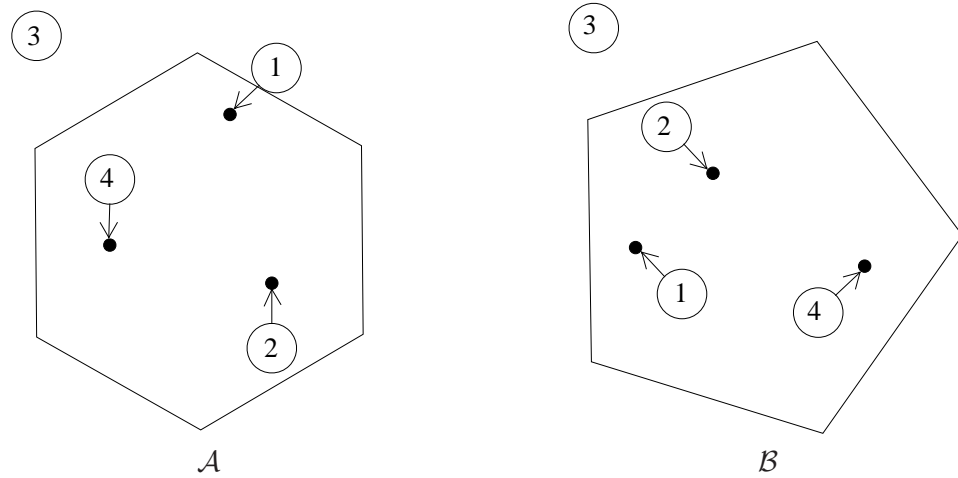
**Figure 6.1:** A four-pebble game.

using $k$ pairs of pebbles. Samson tries to point out a difference between the two structures, and Delilah tries to match his moves so that the differences between them are hidden.

At each move, Samson places one of the pebbles on an element of the universe of one of the two structures, i.e., he places pebble $i$ on an element of $|\mathcal{A}|$, or $|\mathcal{B}|$. Delilah then responds by placing the other pebble $i$ on an element of the other structure.

The position of the game immediately after move $r$ is denoted by $(\alpha_r, \beta_r)$. Such a $k$-*configuration* of $\mathcal{A}, \mathcal{B}$ is a pair of partial functions

$$
\begin{array}{rcl}
\alpha & : & (\mathrm{const}(\tau) \cup \{x_1, x_2, \ldots, x_k\}) \;\rightarrow\; |\mathcal{A}| \\
\beta & : & (\mathrm{const}(\tau) \cup \{x_1, x_2, \ldots, x_k\}) \;\rightarrow\; |\mathcal{B}|
\end{array}
\tag{6.3}
$$

where we require that the domains of the functions $\alpha$ and $\beta$ be equal, $\mathrm{dom}(\alpha) = \mathrm{dom}(\beta)$, and for all $c \in \mathrm{const}(\tau)$, $\alpha(c) = c^{\mathcal{A}}$ and $\beta(c) = c^{\mathcal{B}}$.

The meaning of $\alpha_r(x_i) = a$ and $\beta_r(x_i) = b$ is that just after move $r$, the $i^{\mathrm{th}}$ pebbles are sitting on $a \in |\mathcal{A}|$ and $b \in |\mathcal{B}|$. Some variable $x_i$ is not in the domain of $\alpha_r$ iff just after move $r$, the $i^{\mathrm{th}}$ pebbles are off the board. The valid positions of game $\mathcal{G}_m^k$ on $\mathcal{A}, \mathcal{B}$ consist of any possible $k$-configuration on $\mathcal{A}$ and $\mathcal{B}$. See Figure 6.1 in which the current configuration has $\mathrm{dom}(\alpha) = \mathrm{dom}(\beta) = \{1, 2, 4\} \cup \mathrm{const}(\tau)$, indicating that pebbles 1, 2, and 4 are currently placed on elements of $|\mathcal{A}|$ and $|\mathcal{B}|$, and both pebbles numbered 3 are off the board.

Write $\mathcal{G}_m^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$ to denote the $k$-pebble, $m$-move game on $\mathcal{A}, \mathcal{B}$, with initial configuration $(\alpha_0, \beta_0)$. $\mathcal{G}_m^k(\mathcal{A}, \mathcal{B})$ is the game in which all the pebbles start off the board, i.e., $\text{dom}(\alpha_0) = \text{dom}(\beta_0) = \text{const}(\tau)$. The reason we include the constants in the domain of every configuration is to make the statement of conditions simpler in what follows. As will be seen, in Ehrenfeucht-Fraïssé games constants behave exactly like pebbles that are fixed at the beginning of the game.

At each move $r$, $1 \leq r \leq m$, Samson picks up a pair of pebbles and places one of them on an element of one of the two structures. Delilah must then answer by placing the other pebble in the pair on an element of the other structure. Thus, for some $i \in \{1, 2, \ldots, k\}$, pair $i$ of pebbles is placed on $a \in |\mathcal{A}|$ and $b \in |\mathcal{B}|$. Define the next configuration $(\alpha_r, \beta_r) = (\alpha_{r-1}[a/x_i], \beta_{r-1}[b/x_i])$,

$$\alpha_r(x_j) = \left\{ \begin{array}{ll} \alpha_{r-1}(x_j) & \text{if } i \neq j \\ a & \text{if } i = j \end{array} \right. , \qquad \beta_r(x_j) = \left\{ \begin{array}{ll} \beta_{r-1}(x_j) & \text{if } i \neq j \\ b & \text{if } i = j \end{array} \right. .$$

Just after move $r$, the configuration $\alpha_r, \beta_r$ determines a relation $\beta_r \circ \alpha_r^{-1} \subseteq |\mathcal{A}| \times |\mathcal{B}|$. We say that *Delilah wins round $r$ of the game* iff the map $\alpha_r(x_j) \mapsto \beta_r(x_j)$, for $x_j \in \text{dom}(\alpha_r)$ determines an isomorphism of the induced substructures[1].

$$\beta_r \circ \alpha_r^{-1} : \langle \text{rng}(\alpha) \rangle^{\mathcal{A}} \cong \langle \text{rng}(\beta) \rangle^{\mathcal{B}} \tag{6.4}$$

This means in particular that $\beta_r \circ \alpha_r^{-1}$ must be a 1:1 function, so $\alpha_r(x_i) = \alpha_r(x_j)$ iff $\beta_r(x_i) = \beta_r(x_j)$. Furthermore, all constants and relations of the structures must be preserved. For example, if vocabulary $\tau$ includes the constant symbol $c$ and the binary relation symbol $E$, then $\langle c^{\mathcal{A}}, \alpha_r(x_i) \rangle \in E^{\mathcal{A}}$ iff $\langle c^{\mathcal{B}}, \beta_r(x_i) \rangle \in E^{\mathcal{B}}$. This can be more easily written as follows:

$$(\mathcal{A}, \alpha_r) \models E(c, x_i) \quad \Leftrightarrow \quad (\mathcal{B}, \beta_r) \models E(c, x_i) .$$

Delilah *wins the game* iff she wins every single round. Delilah must preserve an isomorphism at all times. If a difference between the two structures is ever exposed, then Samson wins.

Since $\mathcal{G}_m^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$ is a finite game of perfect information, one of the two players must have a winning strategy. We use the notation $(\mathcal{A}, \alpha_0) \sim_m^k$

---

[1]The induced substructure $\langle \text{rng}(\alpha) \rangle^{\mathcal{A}}$ has universe the closure of $\text{rng}(\alpha)$ under all the functions of $\mathcal{A}$. When $\tau$ has no function symbols, this simply means that we add all the constants to $\text{rng}(\alpha)$. For $\tau = \langle R_1^{a_1}, \ldots, R_r^{a_r}, c_1, \ldots, c_s \rangle$, $|\langle S \rangle| = S \cup \{c_1^{\mathcal{A}}, \ldots, c_s^{\mathcal{A}}\}$. The meaning of induced substructure is that the relations of $\langle \text{rng}(\alpha) \rangle^{\mathcal{A}}$ are restrictions of the relations of $\mathcal{A}$ to the universe of $\langle \text{rng}(\alpha) \rangle^{\mathcal{A}}$.
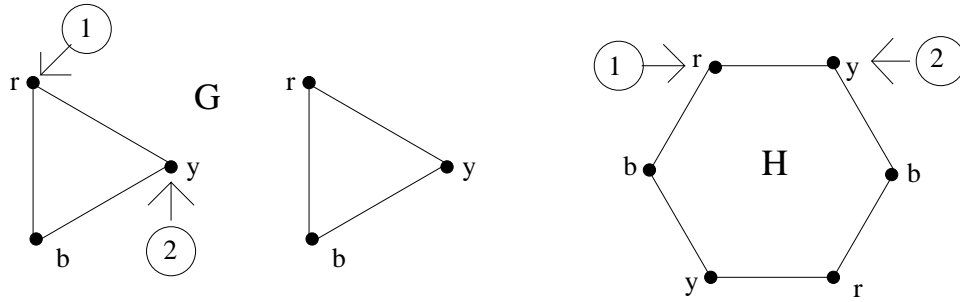
**Figure 6.6:** A two-pebble game

$(\mathcal{B}, \beta_0)$ to mean that Delilah has a winning strategy for $\mathcal{G}_m^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$. We write $(\mathcal{A}, \alpha_0) \sim^k (\mathcal{B}, \beta_0)$ to mean that for all $m$, $(\mathcal{A}, \alpha_0) \sim_m^k (\mathcal{B}, \beta_0)$. Similarly $(\mathcal{A}, \alpha_0) \sim_m (\mathcal{B}, \beta_0)$ means that for all $k$, $(\mathcal{A}, \alpha_0) \sim_m^k (\mathcal{B}, \beta_0)$.                    □

Delilah wins the game iff after every round $\beta_r \circ \alpha_r^{-1}$ is an isomorphism of the induced substructures. Samson is trying to point out a difference between the two structures, and Delilah is trying to keep them looking the same. An isomorphism preserves all the symbols of $\tau$. It is important to decide whether to include the numeric predicates $\leq$ and BIT in $\tau$. If these relations are available in the language in question and thus as part of the definition of isomorphism, then the game becomes much easier for Samson and much harder for Delilah. For this reason, *in this chapter, we assume unless otherwise noted that ordering and* BIT *are not present.*

**Exercise 6.5** Prove that $\sim_m^k$ is an equivalence relation. This is not hard to show directly from the definition of $\mathcal{G}_m^k$. It also explains why we use this notation. Please do not use Theorem 6.10!                    □

As an example, consider the two-pebble game on the colored graphs $G$ and $H$ shown in Figure 6.6. Here the vocabulary $\tau = \langle E^2, R^1, Y^1, B^1 \rangle$ consists of a binary edge relation and three unary relations which may be thought of as colorings of the vertices.

Assume that the initial configuration has all pebbles off the board, so $\alpha_0 = \beta_0 = \emptyset$. Suppose that Samson's first move is to place pebble 1 on a red vertex in $G$. Delilah may answer by putting pebble one on either of the red vertices in $H$. Now suppose Samson puts pebble 2 on an adjacent yellow vertex in $H$. Delilah has a response because in $G$, $\alpha_1(x_1)$ also has an adjacent yellow vertex. On the third move, suppose that Samson puts pebble 1 on the blue vertex in $H$ that is not

adjacent to $\beta_2(x_2)$. Delilah may answer with the blue vertex in $G$ not adjacent to $\alpha_2(x_2)$. The reader should by now be able to prove by induction that,

**Proposition 6.7** *Let $G$ and $H$ be the graphs shown in Figure 6.6. For all $m$, $G \sim_m^2 H$, i.e., $G \sim^2 H$.*

On the other hand Samson has an easy win for the game $\mathcal{G}_3^3(G, H)$. He can simply choose three points in the same triangle in $G$ on three consecutive moves. Delilah has no response because there is no triangle in $H$, and thus Samson wins.

Observe that Samson's winning strategy in this three-pebble game is to "play the sentence" $\Delta$ which says that a triangle exists. Note that $G \models \Delta$ while $H \models \neg\Delta$.

$$\Delta \quad \equiv \quad (\exists x_1)(\exists x_2)(\exists x_3)(E(x_1, x_2) \wedge E(x_2, x_3) \wedge E(x_3, x_1))$$

Sentence $\Delta$ has three variables, corresponding to the number of pebble pairs in the game. Define the *quantifier rank* ($\mathrm{qr}(\varphi)$) of a formula $\varphi$ to be the depth of nesting of quantifiers in $\varphi$. Note that sentence $\Delta$ has quantifier rank 3, corresponding to the number of moves in the game.

**Example 6.8** Another game example is on the strings $w_1 = 1101$ and $w_2 = 1011$, thought of as structures of vocabulary $\tau_s$, with the ordering relation, but not successor. Samson can win the two-move game on these two strings. He can place the $x_1$ pebble on the second 1 in $w_1$. Delilah must answer by placing $x_1$ on some 1 in $w_2$. If she answers with the first 1, then Samson can reply by placing $x_2$ on the first 1 in $w_1$, and Delilah has no reply. If Delilah instead answers with the second or third 1 in $w_2$, then Samson replies by placing $x_2$ on the 0 in $w_1$. Delilah loses because $w_2$ has no 0 to the right of $x_1$. In this case, Samson's winning strategy is to play the following formula that is true of $w_1$, but not $w_2$,

$$\varphi \quad \equiv \quad (\exists x_1)(S(x_1) \ \wedge \ (\exists x_2)(S(x_2) \wedge x_2 < x_1) \ \wedge \ (\exists x_2)(\neg S(x_2) \wedge x_1 < x_2))\square$$

**Definition 6.9** Define language $\mathcal{L}^k$ to be the restriction of language $\mathcal{L}$ in which only variables $x_1, \ldots, x_k$ occur. Define language $\mathcal{L}_m^k$ to be the restriction of $\mathcal{L}^k$ to formulas of quantifier rank at most $m$. Define $\mathcal{L}_m$ to be the set of formulas of quantifier-rank at most $m$. Let $\mathcal{A}$ and $\mathcal{B}$ be two structures of some vocabulary $\tau$.

We say that $\mathcal{A}$ and $\mathcal{B}$ are $\mathcal{L}$ *equivalent* ($\mathcal{A} \equiv_{\mathcal{L}} \mathcal{B}$) iff they agree on all formulas from $\mathcal{L}$,

$$\mathcal{A} \equiv \mathcal{B} \quad \text{iff} \quad \text{for all } \varphi \in \mathcal{L}(\tau), \ \mathcal{A} \models \varphi \ \Leftrightarrow \ \mathcal{B} \models \varphi$$

$$\mathcal{A} \equiv_m^k \mathcal{B} \quad \text{iff} \quad \text{for all } \varphi \in \mathcal{L}_m^k(\tau), \ \mathcal{A} \models \varphi \ \Leftrightarrow \ \mathcal{B} \models \varphi \ . \qquad \square$$

We can now state and prove the fundamental theorem of Ehrenfeucht-Fraïssé Games. This theorem holds for infinite as well as finite structures.

**Theorem 6.10** *Let $\mathcal{A}$ and $\mathcal{B}$ be structures of the same finite, relational vocabulary and let $\alpha_0, \beta_0$ be a $k$-configuration of $\mathcal{A}, \mathcal{B}$. Then the following are equivalent:*

*1. $(\mathcal{A}, \alpha_0) \sim_m^k (\mathcal{B}, \beta_0)$*

*2. $(\mathcal{A}, \alpha_0) \equiv_m^k (\mathcal{B}, \beta_0)$ .*

**Proof** We prove the equivalence of (1) and (2) by induction on $m$. For $m = 0$, Delilah wins the zero move game iff the relation $\beta_0 \circ \alpha_0^{-1}$ is an isomorphism of the induced substructures. This is true iff for every quantifier free formula $\gamma \in \mathcal{L}(\tau)$,

$$(\mathcal{A}, \alpha_0) \models \gamma \quad \Leftrightarrow \quad (\mathcal{B}, \beta_0) \models \gamma \ .$$

Note that $\gamma$ may have as free variables only those variables that occur in $\mathrm{dom}(\alpha_0) = \mathrm{dom}(\beta_0)$. Thus, (1) and (2) are equivalent for $m = 0$.

Assume the theorem is true for $m$, and suppose that $\mathcal{A}$ and $\mathcal{B}$ disagree on the formula $\varphi \in \mathcal{L}_{m+1}^k$. Note that if $\varphi$ is $\alpha \wedge \beta$ then $\mathcal{A}$ and $\mathcal{B}$ disagree on one of $\alpha$ and $\beta$. Similarly, if $\varphi$ is $\neg \alpha$, then they disagree on $\alpha$, so we may assume that $\varphi$ is $(\exists x_i)\psi$. Suppose that $(\mathcal{A}, \alpha_0) \models \varphi$ and $(\mathcal{B}, \beta_0) \models \neg\varphi$. Samson's first move in $\mathcal{G}_{m+1}^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$ is to place pebble $i$ on a witness for $\psi$ in $\mathcal{A}$. Wherever Delilah responds, it will not be a witness for $\psi$ because there is none in $\mathcal{B}$. Thus, after the first move, $(\mathcal{A}, \alpha_1)$ and $(\mathcal{B}, \beta_1)$ disagree on the quantifier depth $m$ formula $\psi$. By the inductive hypothesis, Samson has a winning strategy for the remaining $m$-move game. Thus we have shown that (1) implies (2).

Conversely, suppose that $(\mathcal{A}, \alpha_0) \equiv_{m+1}^k (\mathcal{B}, \beta_0)$. Now let Samson make his first move in the game $\mathcal{G}_{m+1}^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$. Suppose he places pebble $i$ on an element of $\mathcal{A}$, thus defining $\alpha_1$. Note that there are only finitely many inequivalent formulas in $\mathcal{L}_m^k$, cf. Exercise 6.11. Let $\Phi$ be the conjunction of all these formulas that are satisfied by $(\mathcal{A}, \alpha_1)$. Thus, we know that

$$(\mathcal{A}, \alpha_0) \models (\exists x_i)\Phi \ ,$$

so, by assumption,

$$(\mathcal{B}, \beta_0) \models (\exists x_i)\Phi \ .$$

Delilah's answer is to place the other pebble $i$ on a witness in $\mathcal{B}$ of $\Phi$. Thus, $(\mathcal{A}, \alpha_1)$ and $(\mathcal{B}, \beta_1)$ both satisfy $\Phi$, a complete description of every formula from $\mathcal{L}_m^k$ that $(\mathcal{A}, \alpha_1)$ satisfies. Thus $(\mathcal{A}, \alpha_1) \equiv_m^k (\mathcal{B}, \beta_1)$. It follows by induction that Delilah has a winning strategy for the remaining $m$ moves of the game. □

In the following exercise, you are asked to prove the lemma that there are only finitely many inequivalent formulas of a given quantifier rank. This was needed for the proof of $(2) \Rightarrow (1)$ in Theorem 6.10. As the exercise shows, this lemma holds for infinite structures as well, as long as the vocabulary is finite and has no function symbols.

**Exercise 6.11** Theorem 6.10 requires the lemma that there are only finitely many inequivalent formulas of quantifier rank $r$.

1. As usual, let $\tau$ be a finite, relational vocabulary. Prove that there are only finitely many inequivalent first-order formulas in $\mathcal{L}_r(\tau)$. [Hint: induction on $r$.]

2. Let $\tau$ be a finite vocabulary which may include function symbols. Let $\Gamma \subset \mathcal{L}(\tau)$ be a first-order theory. Suppose that for every model $\mathcal{A}$ of $\Gamma$ — including infinite models — and for any finite set $S \subseteq |\mathcal{A}|$, the induced substructure of $\mathcal{A}$ generated by $S$ $\langle S \rangle^{\mathcal{A}}$ is finite. Prove that theory $\Gamma$ admits only finitely many inequivalent formulas of quantifier rank $r$. In this case, $\varphi$ and $\psi$ are equivalent iff $\Gamma \vdash \varphi \leftrightarrow \psi$.

3. Give a counterexample to Theorem 6.10 when $\tau = \langle R_1^1, R_2^1, \ldots \rangle$ consists of infinitely many unary relation symbols.

4. Give a counterexample to Theorem 6.10 when $\tau = \langle R^1, f^1 \rangle$ consists of one unary relation symbol and one unary function symbol. □

The following exercise shows that we never have to consider a move of a game in which Samson pebbles an element that is already pebbled by another pebble or constant.

**Exercise 6.12** Prove that in any game $\mathcal{G}_m^k(\mathcal{A}, \alpha_0, \mathcal{B}, \beta_0)$, if Samson has a winning strategy, then he still has a winning strategy if he is never allowed to place a pebble on a constant or an element that already has another pebble sitting on it. □

Theorem 6.10 gives us a way to determine precisely how many variables and how much quantifier rank is needed to express a given query. Here are two examples.

**Proposition 6.13** *Let* CLIQUE$(k)$ *be the set of undirected graphs that contain a clique, i.e., a complete subgraph, of size $k$. In the language without ordering,* CLIQUE$(k)$ *is expressible with $k$ variables but not $k - 1$ variables:*

$$\text{CLIQUE}(k) \quad \in \quad \mathcal{L}^k(\tau_g)(\text{wo}\leq) \ - \ \mathcal{L}^{k-1}(\tau_g)(\text{wo}\leq) \ .$$

**Proof** It is easy to write CLIQUE$(k)$ in $\mathcal{L}^k$:

$$(\exists x_1 x_2 \ldots x_k)(\text{distinct}(x_1, \ldots, x_k) \wedge E(x_1, x_2) \wedge \ldots \wedge E(x_1, x_k) \wedge \ldots E(x_{k-1}, x_k))$$

To prove that $k$ variables are necessary, we prove that $K_k \sim^{k-1} K_{k-1}$, where $K_r$ is the complete graph on $r$ vertices. Delilah has a simple winning strategy for the game $\mathcal{G}_{k-1}(K_k, K_{k-1})$: When Samson places a pebble on an unpebbled vertex in one of the two graphs, Delilah places the corresponding pebble on any unpebbled vertex in the other graph. Since there are only $k - 1$ pebble pairs, such an unpebbled vertex is always available. Note that this is a winning strategy since edges exist between all points in each graph. Thus, any 1:1 correspondence is an isomorphism. □

As another example, we show that first-order logic without ordering is not strong enough to express any facts about counting, or even parity.

**Proposition 6.14** *In the absence of ordering, the boolean query on graphs that is true iff there are an odd number of vertices requires $n+1$ variables, for graphs with $n$ or more vertices. The same is true for the query that there are an odd number of edges.*

**Proof** Let $G_n$ be the graph on $n$ vertices that has a loop at each vertex but no other edges. We claim that $G_n \sim^n G_{n+1}$. Delilah's strategy is to match each move by Samson on a vertex not already pebbled with a vertex not already pebbled in the other graph. Since each graph has at least $n$ vertices and there are no edges between different vertices, this is a winning strategy for Delilah. It follows that $G_n \equiv^n G_{n+1}$, so the parity of the number of vertices or the number of edges is not expressible in $\mathcal{L}^n$. □

We saw in Chapter 5 that quantifier rank and number of variables are important parameters of parallel complexity. It is useful to have a game that allows us to determine how many variables and how much quantifier rank is needed to describe various queries. As an example, we now use Ehrenfeucht-Fraïssé games to determine the exact number of variables and quantifier rank needed to assert the existence of paths in a graph:

**Proposition 6.15** *Let the formula* $\mathrm{PATH}_k(x, y) \in \mathcal{L}(\tau_g)$ *mean that there is a path of length at most* $2^k$ *from* $x$ *to* $y$. *With or without ordering, quantifier rank* $k$ *is necessary and sufficient to express* $\mathrm{PATH}_k$. *Furthermore, only three variables are necessary to express* $\mathrm{PATH}_k$. *In symbols,*

$$\mathrm{PATH}_k \quad \in \quad \mathcal{L}_k^3(\tau_g)(\mathrm{wo}\leq) \; - \; \mathcal{L}_{k-1}(\tau_g) \; .$$

**Proof** For the upper bound, we express $\mathrm{PATH}_k$ inductively as follows:

$$\begin{aligned} \mathrm{PATH}_0(x, y) &\equiv& x = y \;\vee\; E(x, y) \\ \mathrm{PATH}_{k+1}(x, y) &\equiv& (\exists z)(\mathrm{PATH}_k(x, z) \wedge \mathrm{PATH}_k(z, y)) \; . \end{aligned}$$

Thus, $\mathrm{PATH}_k$ is expressible using three variables and quantifier rank $k$. Only three variables are needed because the right hand side of the inductive definition of $\mathrm{PATH}_{k+1}(x, y)$ may be written in a way that reuses variables:

$$\begin{aligned} \mathrm{PATH}_k(x, z) &\equiv& (\exists y)(\mathrm{PATH}_{k-1}(x, y) \wedge \mathrm{PATH}_{k-1}(y, z)) \\ \mathrm{PATH}_k(z, y) &\equiv& (\exists x)(\mathrm{PATH}_{k-1}(z, x) \wedge \mathrm{PATH}_{k-1}(x, y)) \end{aligned}$$

For the lower bound, we prove the following,

**Lemma 6.16** *Let* $G$ *and* $H$ *be line graphs with at least* $2^{k+1} + 1$ *vertices and constants 0 and max for the leftmost and rightmost elements. Then Delilah wins the* $k$ *move game on* $G$ *and* $H$.

We prove Lemma 6.16 by induction on $k$.

$k = 0$ : Both graphs have at least three vertices. Thus there is no edge from 0 to *max* in either graph and Delilah wins the 0 move game.

Assume true for $k - 1$ and let's play the game for $k$. Think of both $G$ and $H$ as being in three parts: the first $2^k + 1$ points, the middle part, and the last $2^k + 1$

points. The middle part starts with the rightmost vertex of the first part and ends with the leftmost vertex of the last part. (If the size of one of the graphs is exactly $2^{k+1} + 1$, then the left and right parts intersect at a point and the middle part is that intersecting point.)

**Key point:**   once the first pair of moves is made $G$ and $H$ are both split into two parts, and Delilah need only have a winning strategy on each part separately.

Delilah's first move:

1. If Samson plays pebble $p$ in the first part of $G$, say $\ell$ points from the left, then Delilah plays $\ell$ points from the left in first part of $H$. The game is now in two halves— the line graph up to and including $p$ and the line graph from $p$ on.

   For the half to the left of $p$, these sections of $G$ and $H$ are the same legnth, so Delilah wins trivially. In the half to the right of $p$, both sections have length at least $2^k + 1$, so Delilah wins by the induction hypothesis.

2. If Samson plays in last part of $G$, or in first or last part of $H$, then Delilah's strategy is similar.

3. If Samson plays the middle part of either $G$ or $H$, then Delilah plays in the middle part of the other. Both halves of the game are now inductively wins for Delilah.

This proves Lemma 6.16 and Proposition 6.15.                                          □

We remark that the three variables used to express paths in Proposition 6.15 are necessary. In Proposition 6.7, we saw a connected graph $H$ of diameter three and a disconnected graph $G$ such that $G \sim^2 H$. It follows from Theorem 6.10 that CONNECTED is not expressible using 2 variables, no matter what the quantifier rank. Suppose $\text{PATH}_k$ were expressible using only two variables for some $k \geq 2$. Then $G$ and $H$ would differ on the $\mathcal{L}^2$ formula $(\forall x_1 x_2)\text{PATH}_k(x_1, x_2)$.

In fact, let $(\alpha_0, \beta_0)$ be the 2-configuration of graphs $G, H$ shown in Figure 6.17. Observe that $(G, \alpha_0) \sim^2 (H, \beta_0)$, but these two structures disagree on the formula $\text{PATH}_1(x_1, x_2)$. It follows that for $k \geq 1$, $\text{PATH}_k$ is not expressible in $\mathcal{L}^2$.
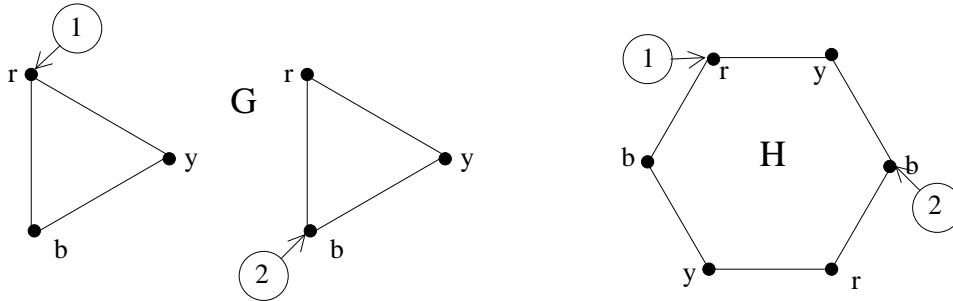
**Figure 6.17:** The game $\mathcal{G}^2(G, \alpha_0, H, \beta_0)$

## 6.2 Methodology for First-Order Expressibility

As we now show, Ehrenfeucht-Fraïssé games provide a complete methodology for
proving that a query is not first-order. We have already seen that these games are a
convenient tool for determining what can be said in first-order logic. This theorem
says that if we can show using any method that a query is not first-order expressible,
then we can show it using Ehrenfeucht-Fraïssé games.

**Theorem 6.18 (Methodology Theorem)** *Let $\mathcal{C}$ be any class of finite or infinite
structures of some finite, relational vocabulary. Let $S \subseteq \mathcal{C}$ be a boolean query on
$\mathcal{C}$. To prove that $S$ is* not *first-order describable on $\mathcal{C}$ it is necessary and sufficient
to show that for all $r \in \mathbf{N}$, there exist structures $\mathcal{A}_r, \mathcal{B}_r \in \mathcal{C}$ such that*

1. *$\mathcal{A}_r \in S$ and $\mathcal{B}_r \notin S$, and*

2. *$\mathcal{A}_r \sim_r \mathcal{B}_r$*

**Proof** It is easy to see that the methodology is sufficient. The above two conditions
imply that $\mathcal{A}_r$ and $\mathcal{B}_r$ agree on all formulas in $\mathcal{L}_r$, but disagree on $S$. Thus, $S$ is
not expressible in $\mathcal{L}_r$ for any $r$.

Conversely, suppose that $S$ is not first-order expressible over $\mathcal{C}$. Recall from
Exercise 6.11 that for any fixed $r$ there are only a bounded number of inequivalent
sentences of quantifier rank $r$.

We say that $\varphi \in \mathcal{L}_r$ is a *complete quantifier-rank $r$ sentence* iff for every
other quantifier-rank $r$ sentence $\psi$ of the same vocabulary, $\varphi \vdash \psi$ or $\varphi \vdash \neg\psi$. Let
$\varphi_1, \varphi_2, \ldots, \varphi_B$ be a list of all inequivalent, complete quantifier-rank $r$ sentences.

For every quantifier-rank $r$ sentence $\psi$, each $\varphi_i$ must assert either $\psi$ or $\neg\psi$. Observe that each structure from $\mathcal{C}$ satisfies a unique $\varphi_i$. Suppose there are structures $\mathcal{A}_r \in S$ and $\mathcal{B}_r \in \mathcal{C} - S$ that satisfy the same $\varphi_i$. Then $\mathcal{A}_r$ and $\mathcal{B}_r$ satisfy the above conditions. If there is no such pair, then the $\varphi_i$'s are partitioned by $S$. In this case, let $Y = \{i \mid (\exists \mathcal{A} \in S)(\mathcal{A} \models \varphi_i)\}$ and let

$$\varphi \quad \equiv \quad \bigvee_{i \in Y} \varphi_i \,.$$

Then $\varphi$ is a first-order formula of quantifier rank $r$ that expresses $S$.  □

The Methodology Theorem holds whether or not we include ordering in our languages. We now give a few easy applications, proving that various properties are not first-order. Most of these applications do not include ordering.

First, we introduce a general theorem that allows the use of the Methodology Theorem without constructing winning strategies for Delilah by hand.

**Definition 6.19 (Gaifman Graph, $d$-type)** Let $\mathcal{A}$ be any structure of vocabulary $\tau = \langle R_1^{a_1}, \ldots, R_r^{a_r}, c_1, \ldots, c_s \rangle$. Define the *Gaifman graph* $G_{\mathcal{A}} = (|\mathcal{A}|, E_{\mathcal{A}})$ as follows:

$$E_{\mathcal{A}} \;=\; \left\{ (a,b) \;\middle|\; a \neq b \wedge (\exists i)(\exists \langle d_1, \ldots, d_{a_i} \rangle \in R_i^{\mathcal{A}})(a, b \in \{d_1, \ldots, d_{a_i}\}) \right\} \,.$$

There is an edge between $a$ and $b$ in the Gaifman graph iff $a$ and $b$ occur in the same tuple of some relation of $\mathcal{A}$. As a simple example, if $\mathcal{A} \in \mathrm{STRUC}[\tau_g]$ is a loop-free graph, then $G_{\mathcal{A}} = \mathcal{A}$.

Let $(\mathcal{A}, \alpha_r)$ be the configuration of structure $\mathcal{A}$ after move $r$ of a game. Define the universe of the neighborhood of element $a$ at distance $d$ to be the set of elements of distance at most $d$ from $a$ in the Gaifman graph:

$$|N(a,d)| \quad = \quad \left\{ b \in |\mathcal{A}| \;\middle|\; \mathrm{DIST}(a,b) \leq d \right\}$$

$N(a,d)$ is almost an induced substructure of $(\mathcal{A}, \alpha_r)$: It inherits the relations from $\mathcal{A}$, but it contains only those constants and pebbled points that are within distance $d$ of $a$. Define the *$d$-type* of $a$ to be the isomorphism type of $N(a,d)$. Note that isomorphisms must send each constant $c_j^{\mathcal{A}}$ to $c_j^{\mathcal{B}}$ and each pebbled point $\alpha_r(x_i)$ to $\beta_r(x_i)$. (Neighborhood $N(a,d)$ and thus the $d$-type of $a$ depend on the current configuration $(\mathcal{A}, \alpha_r)$. If the configuration is not clear from the context, then we say the $d$-type of $a$ with respect to configuration $(\mathcal{A}, \alpha_r)$.)  □

The above definitions allow the following,

**Theorem 6.20 (Hanf's Theorem)** *Let $\mathcal{A}, \mathcal{B} \in \mathrm{STRUC}[\tau]$ and let $r \in \mathbf{N}$. Suppose that for each possible $2^r$-type, $t$, $\mathcal{A}$ and $\mathcal{B}$ have exactly the same number of elements of type $t$. Then $\mathcal{A} \equiv_r \mathcal{B}$.*

**Proof** We must show that Delilah wins the game $\mathcal{G}_r(\mathcal{A}, \mathcal{B})$. This is similar to the proof of Proposition 6.15. Delilah's winning strategy is to maintain the following invariant: after move $m$, $0 \le m \le r$,

$$(\mathcal{A}, \alpha_m), (\mathcal{B}, \beta_m) \text{ have same number of each } 2^{r-m}\text{-type.} \tag{6.22}$$

In $\mathcal{G}_r(\mathcal{A}, \mathcal{B})$ there is no bound on the number of pebbles. Therefore we may assume that Samson uses a new pebble at each step. Thus Delilah wins iff she wins at the last round. If Delilah preserves (6.22) then after the last move, the neighborhoods of distance one around each constant or pebbled point are isomorphic to the corresponding neighborhoods in the other structure. It follows that Delilah wins the game.

We have that (6.22) holds for $m = 0$ by assumption. Inductively, assume that it holds after move $m$. On move $m + 1$, let Samson choose some vertex $v$. Delilah should respond with any vertex $v'$ of the same $2^{r-m}$-type as $v$.

We have to show that (6.22) still holds. The inductive assumption immediately implies that,

$$(\mathcal{A}, \alpha_m), (\mathcal{B}, \beta_m) \text{ have same number of each } 2^{r-(m+1)}\text{-type.}$$

Furthermore, the neighborhood $N(a, 2^{r-(m+1)})$ of $(\mathcal{A}, \alpha_m)$ is different from the same neighborhood of $(\mathcal{A}, \alpha_{m+1})$ iff $\mathrm{DIST}(a, v) \le 2^{r-(m+1)}$. Consider the isomorphism $f : N(v, 2^{r-m}) \to N(v', 2^{r-m})$. It maps every vertex $a$ in $N(v, 2^{r-(m+1)})$ to a corresponding $a' \in N(v', 2^{r-(m+1)})$. Here is the key idea: $f$ maps $N(a, 2^{r-(m+1)})$ isomorphically onto $N(a', 2^{r-(m+1)})$ because these smaller neighborhoods lie inside $\mathrm{dom}(f) = N(v, 2^{r-m})$ (see Figure 6.21). Thus, there is a 1:1 correspondence between the isomorphism types of these neighborhoods close to $v$ and $v'$, so the 1:1 correspondence between the other neighborhoods is undisturbed.

Thus, Delilah's strategy preserves Equation (6.22) and wins the game.     $\square$

As a sample application of Theorem 6.20, we prove the following:

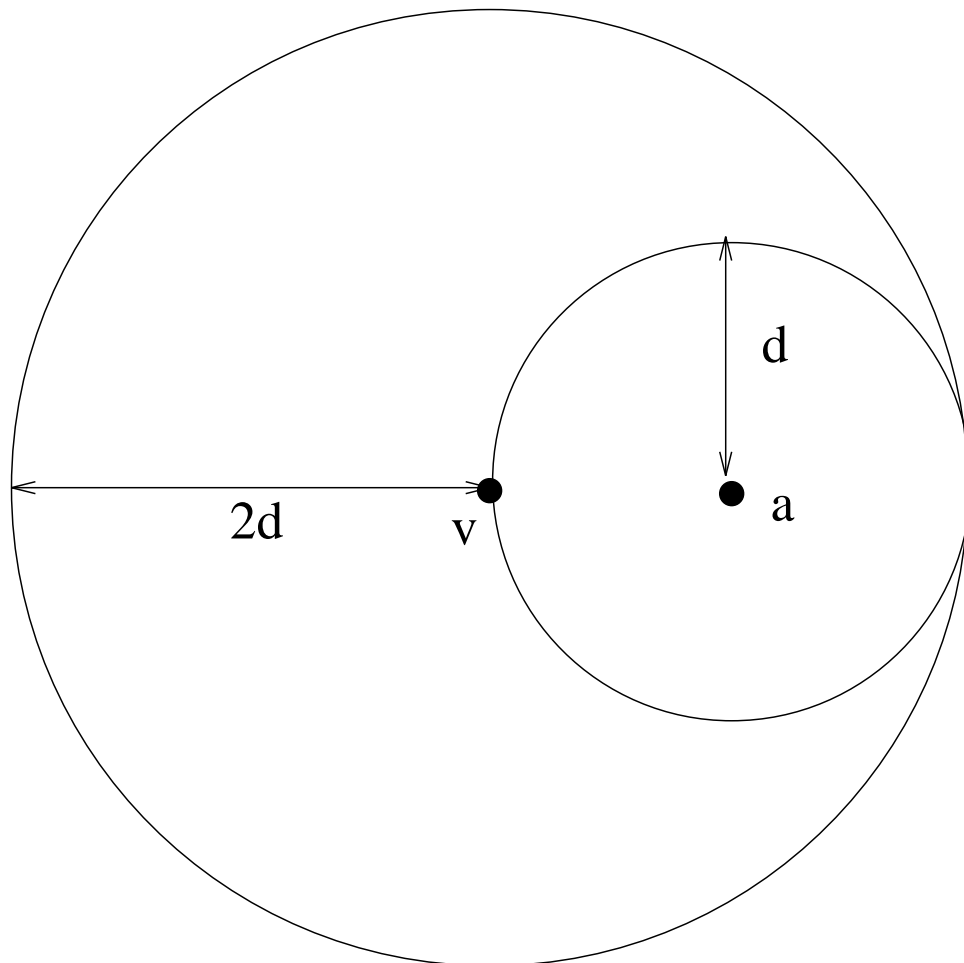**Proposition 6.23** *Acyclicity is not first-order expressible.*

**Figure 6.21:** Inductive step in proof of Hanf's Theorem: $d = 2^{r-(m+1)}$

**Figure 6.24:** Proof of Proposition 6.23.

**Proof** Let $\mathcal{A}_r$ be a line segment on $2^{r+3}$ vertices. Let $\mathcal{B}_r$ be the union of a line segment on $2^{r+2}$ vertices and a cycle on $2^{r+2}$ vertices. See Figure 6.24. Observe that $\mathcal{A}_r$ and $\mathcal{B}_r$ both have the same number of each $2^r$-type. Therefore, by Theorem 6.20, $\mathcal{A}_r \equiv_r \mathcal{B}_r$. It follows that Acyclicity is not first-order expressible. Note that the same proof works for directed or undirected graphs. $\qquad\square$

**Exercise 6.25** Using Hanf's Theorem, prove that the following boolean queries are not first-order expressible in the language without ordering.

1. Two-colorability of graphs (cf. Exercise 4.27). [Hint: use $\mathcal{A}_r$ a cycle of size $6d$, and $\mathcal{B}_r$ two cycles of size $3d$ each, with $d = 3^r$.]

2. Consider the following boolean query,

$$\begin{aligned} \text{CONNECTED} \quad &= \quad \left\{ G \mid G \text{ is a connected, undirected graph} \right\} \\ &\equiv \quad (\forall xy)(\text{PATH}(x,y) \quad \wedge \quad (E(x,y) \rightarrow E(y,x))) \end{aligned}$$

   Prove that

$$\text{CONNECTED} \in \mathcal{L}^3_{2+\lceil \log(n-1) \rceil}(\tau_g)(\text{wo}\leq) - \mathcal{L}_{\lceil \log(n-2)-2 \rceil}(\tau_g)(\text{wo}\leq)$$

   [Hint: for the upper bound use Proposition 6.15. For the lower bound, use Hanf's theorem with $G_r$ a pair of disjoint cycles of $2^{r+1}+1$ vertices each and $H_r$ a cycle of $2^{r+2}+2$ vertices.]

3. Show that REACH is not first-order. [Hint: this is very similar to (2). You just need to place the constants $s$ and $t$ appropriately. Note that $s$ and $t$ may be thought of as the first two pebbled points in the game. Thus, you need $\mathcal{A}_r = G_{r+2}$ and $\mathcal{B}_r = H_{r+2}$ with the appropriately placed $s$ and $t$.]    □

## 6.3  First-Order Properties are Local

Hanf's theorem implies that every first-order property is local in the sense that it only concerns neighborhoods of a fixed radius. We have seen that this locality is a useful tool in proving that certain queries are not first-order.

The *degree* of a graph is the maximum number of edges adjacent to any vertex. The *degree* of a structure $\mathcal{A}$ is the degree of its Gaifman graph. We now prove a strengthening of Hanf's theorem for graphs of bounded degree. In this generalization, the number of instances of a given $r$-type in the two structures need not be equal as long as both numbers are sufficiently large.

Let $\mathcal{A}$ and $\mathcal{B}$ be structures and let $n, s$ be integers. We say that $\mathcal{A}$ and $\mathcal{B}$ are $(n, s)$-*equivalent* iff for each $n$-type, $\sigma$, $\mathcal{A}$ and $\mathcal{B}$ have the same number of neighborhoods of type $\sigma$ or they both have more than $s$ neighborhoods of type $\sigma$. The following is a generalization of Hanf's theorem for structures of bounded degree.

**Theorem 6.26  (Bounded-Degree Hanf Theorem)**  *Let $r$ and $d$ be fixed. There is an integer $s$ such that for all structures $\mathcal{A}$ and $\mathcal{B}$ of degree at most $d$, if $\mathcal{A}$ and $\mathcal{B}$ are $(2^r, s)$-equivalent, then $\mathcal{A} \equiv_r \mathcal{B}$.*

**Proof**  This proof is similar to the proof of Theorem 6.20. We must show that Delilah wins the game $\mathcal{G}_r(\mathcal{A}, \mathcal{B})$. Let $s = rd^{2^r} + 1$. Delilah's winning strategy is to maintain the following invariant: after move $m$, $0 \leq m \leq r$,

$$(\mathcal{A}, \alpha_m), (\mathcal{B}, \beta_m) \text{ have the same number of each } 2^{r-m}\text{-type.}$$
$$\text{or both have over } (r - m)d^{2^r} + 1 \text{ elements of this type.} \tag{6.27}$$

We have that (6.27) holds for $m = 0$ by assumption. Inductively, assume that it holds after move $m$. On move $m + 1$, let Samson choose some vertex $v$. Delilah should respond with any vertex $v'$ of the same $2^{r-m}$-type as $v$.

We have to show that Equation (6.27) still holds. The inductive assumption immediately implies that,

$(\mathcal{A}, \alpha_m), (\mathcal{B}, \beta_m)$ have the same number of each $2^{r-(m+1)}$-type. or both have over $(r - (m + 1))d^{2^r} + 1$ elements of this type.

Just as in the proof of Theorem 6.20, the only neighborhoods that change are those within distance $2^{r-(m+1)}$ of $v$. Furthermore, the same number of neighborhoods change in the same way in $\mathcal{A}$ as in $\mathcal{B}$. The only harm that can be done to Equation (6.27) is that the number of some types can be reduced by the same amount in $\mathcal{A}$ and in $\mathcal{B}$. The number of vertices within distance $\rho = 2^{r-(m+1)}$ of $v$ is at most $d^{\rho+1}/(d - 1)$ which is less than $d^{2^r}$. Thus we have (6.27) holds for $m + 1$ as desired. $\qquad\square$

A striking application of Theorem 6.26 is the following theorem of Seese. The definition of linear time in the following is linear time on a unit-cost RAM with $O(\log n)$ bit word size,

**Theorem 6.28** *Let $\varphi \in$ FO. Then over bounded degree structures, $\varphi$ is recognizable in linear time.*

**Proof** For simplicity, assume that the structures in question are bounded degree graphs and let them be given via adjacency lists.[2] Let $r$ be the quantifier rank of $\varphi$ and let $d$ be the degree of the graphs in question.

There are a large but bounded number of possible $2^r$-types in degree $d$ graphs. The linear-time algorithm is to determine the $2^r$ type of each vertex and count — up to $s$ — how many of each type occurs. This information is what we can call the $2^r, s$ description of $G$. By Theorem 6.26, the $2^r, s$ description of $G$ determines whether $G$ satisfies $\varphi$. We could in principle build — once and for all — a table that lists for each of the finitely many possible $2^r, s$ descriptions, whether or not a graph with this description satisfies $\varphi$. From $G$'s description, we can use the table to check in constant additional time whether $G$ satisfies $\varphi$. $\qquad\square$

## 6.4 Bounded Variable Languages

A theory $\Sigma$ satisfies the *k-variable property* iff every first-order formula is equivalent with respect to $\Sigma$ to a first-order formula that has only $k$ bound variables.

---

[2]Adjacency lists are linked lists, one for each vertex, listing all the vertices adjacent to the given vertex, see [AHU74].

Gabbay has shown that the set of models of $\Sigma$ has the $k$-variable property for some $k$ iff there exists a finite basis for the set of all temporal-logic connectives over these models [Gab81]. Kozen and Immerman used Ehrenfeucht-Fraïssé games to prove that the theories of linear order and of bounded degree trees have the $k$-variable property, for appropriate $k$, (Fact 12.32).

It is interesting and useful to know when a set of structures has the property that all first-order formulas can be expressed using only a bounded number of bound variables. In this section we give one example, showing that the theory of linear order has the 3-variable property.

We begin with a lemma that allows us to give a game-theoretic proof that a theory has the $k$-variable property. This lemma uses the Compactness Theorem (Theorem 1.35). For this reason, *in this section we consider all structures, not just finite structures.*

**Lemma 6.29** *Let $\Sigma \subset \mathcal{L}$ be a first-order theory. Let $\mathcal{L}'$ and $\mathcal{L}''$ be subsets of $\mathcal{L}$ such that $\mathcal{L}'$ is closed under the boolean connectives. Let $k \in \mathbf{N}$. The following conditions are equivalent:*

   *1. For all models $\mathcal{A}$ and $\mathcal{B}$ of $\Sigma$ and all $k$-configurations $\alpha, \beta$ of $\mathcal{A}, \mathcal{B}$,*

$$(\mathcal{A}, \alpha) \equiv_{\mathcal{L}'} (\mathcal{B}, \beta) \quad \Rightarrow \quad (\mathcal{A}, \alpha) \equiv_{\mathcal{L}''} (\mathcal{B}, \beta)$$

   *2. For all $\varphi \in \mathcal{L}''$ with free variables among $x_1, \ldots, x_k$, there exists $\psi \in \mathcal{L}'$ such that $\Sigma \models \varphi \leftrightarrow \psi$.*

**Proof** $(2 \rightarrow 1)$: If every formula in $\mathcal{L}''$ is equivalent to a formula in $\mathcal{L}'$ and $(\mathcal{A}, \alpha)$ and $(\mathcal{B}, \beta)$ are $\mathcal{L}'$-equivalent, then they are $\mathcal{L}''$-equivalent.

$(1 \rightarrow 2)$: If $\Sigma \cup \{\varphi\}$ is inconsistent, then we may take $\psi \equiv \mathbf{false}$. Otherwise, let $T$ be the set of all complete $\mathcal{L}'$-types over the variables $x_1, \ldots, x_k$ that is consistent with $\Sigma \cup \{\varphi\}$. Let $\Gamma \in T$ be such a type. Observe that $\Sigma \cup \Gamma \models \varphi$. Otherwise, we could construct models $(\mathcal{A}, \alpha)$ and $(\mathcal{B}, \beta)$ of $\Sigma \cup \Gamma$ that disagree on $\varphi$. This is impossible by (1). It follows by the compactness theorem that there is a formula $\psi_\Gamma \in \Gamma$ such that $\Sigma \models \psi_\Gamma \rightarrow \varphi$.

Define the following set of formulas,

$$D \quad = \quad \left\{ \neg \psi_{\Gamma_i} \mid \Gamma_i \in T \right\}.$$

Then $\Sigma \cup D \cup \{\varphi\}$ is inconsistent. By compactness, there must be some finite $F \subseteq T$ such that

$$\Sigma \models \bigwedge_{\Gamma_i \in F} \neg\psi_{\Gamma_i} \ \rightarrow \ \neg\varphi$$

We can take $\psi = \bigvee_{\Gamma_i \in F} \psi_{\Gamma_i}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Let $\Sigma \subset \mathcal{L}$ be a theory, let $\mathcal{L}' = \mathcal{L}^k$, and let $\mathcal{L}'' = \mathcal{L}$. In this case, Lemma 6.29 implies that condition 1 — which may be proved using Ehrenfeucht–Fraïssé games — is sufficient to show that every formula in $\mathcal{L}$ that has at most $k$ free variables is equivalent to a formula in $\mathcal{L}^k$. To prove the $k$-variable property, we must also show that any formula with more than $k$ free variables is equivalent to a formula with at most $k$ bound variables. The following exercise explains how to do this.

**Exercise 6.30** Let $\mathcal{L}$ be a first-order relational language with no relation symbols of arity greater than $k$. Suppose that $\Sigma \subset \mathcal{L}$ is a theory and that $R_1, R_2, \ldots$ are an infinite set of monadic relation symbols from $\mathcal{L}$ that do not occur in $\Sigma$. Even though we have infinitely many $R_i$'s, we consider only structures in which only finitely many relations are non-empty. Suppose that for every pair of such structures $\mathcal{A}, \mathcal{B}$ satisfying $\Sigma$ and every pair of $k$-configurations $\alpha, \beta$, we have

$$(\mathcal{A}, \alpha) \equiv^k (\mathcal{B}, \beta) \quad \Rightarrow \quad (\mathcal{A}, \alpha) \equiv (\mathcal{B}, \beta) \ .$$

Prove that $\Sigma$ has the $k$-variable property.

[Hint: this follows essentially from Lemma 6.29. The part you must fill in is how to replace the extra free variables by new monadic relation symbols.] $\qquad\square$

The following theorem shows that for structures consisting of a linear order plus a finite number of unary relation symbols, three variables suffice to express all first-order properties. The proof of this theorem produces a winning strategy for Delilah that combines her strategies from several simpler games.

**Theorem 6.31** *The set of linear ordered structures satisfies the 3-variable property. These structures may also include any number of monadic relation symbols.*

**Proof** By Exercise 6.30 it suffices to show that for any pair of linear orders $\mathcal{A}, \mathcal{B}$ and any pair of 3-configurations $\alpha, \beta$,

$$(\mathcal{A}, \alpha) \equiv^3 (\mathcal{B}, \beta) \quad \Rightarrow \quad (\mathcal{A}, \alpha) \equiv (\mathcal{B}, \beta) \ .$$

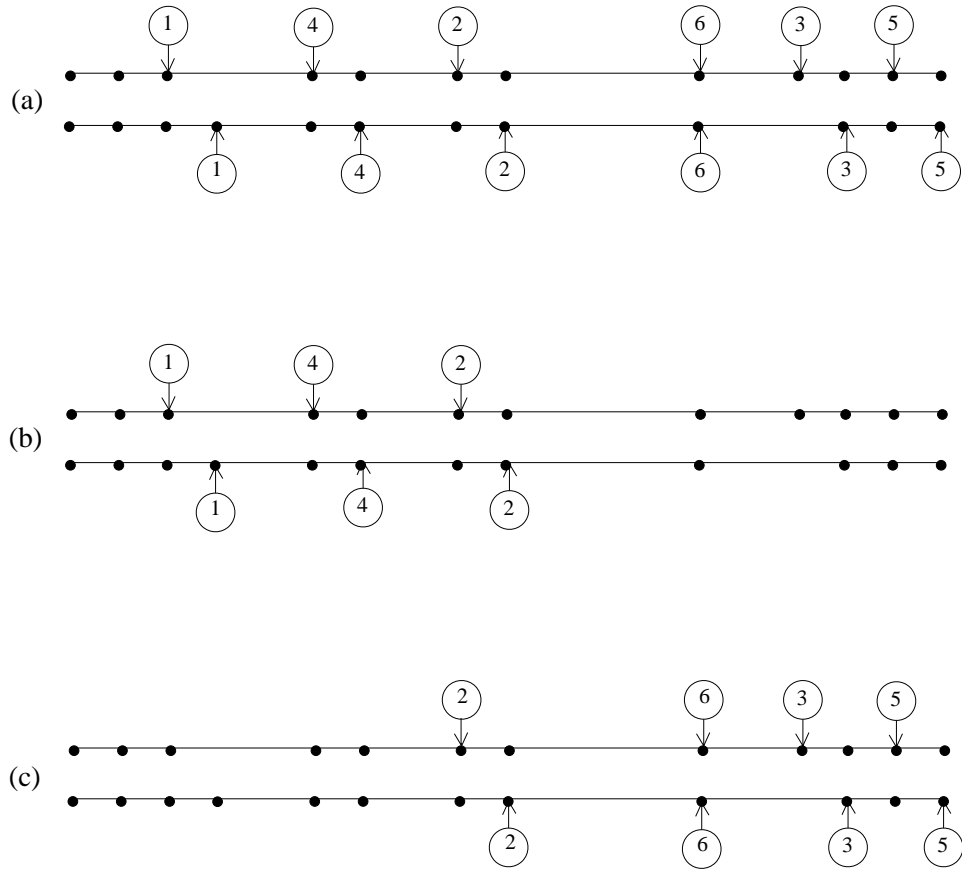We prove the slightly stronger result that for all $m$,

**Figure 6.32:**  Delilah's winning strategy in $\mathcal{G}_{m+1}(\mathcal{A}, \alpha, \mathcal{B}, \beta)$ (a) is built from her winning strategies in $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_\ell, \mathcal{B}, \beta_\ell)$ (b) and $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_r, \mathcal{B}, \beta_r)$ (c).

$$(\mathcal{A}, \alpha) \sim_m^3 (\mathcal{B}, \beta) \quad \Rightarrow \quad (\mathcal{A}, \alpha) \sim_m (\mathcal{B}, \beta) \ . \tag{6.33}$$

We prove Equation (6.33) by induction on $m$. The base case, $m = 0$, is clear because extra pebbles cannot help Samson in the zero move game.

Assume that (6.33) holds for $m$ and suppose that,

$$(\mathcal{A}, \alpha) \sim_{m+1}^3 (\mathcal{B}, \beta) \ . \tag{6.34}$$

We now describe a winning strategy for Delilah in the game $\mathcal{G}_{m+1}(\mathcal{A}, \alpha, \mathcal{B}, \beta)$. Suppose that in the initial configuration, $|\mathrm{dom}(\alpha)| < 3$, that is, fewer than three pebbles are on the board. In this case, wherever Samson plays, Delilah can answer using her winning strategy for the game $\mathcal{G}_{m+1}^3(\mathcal{A}, \alpha, \mathcal{B}, \beta)$. Let $\alpha_1, \beta_1$ be the resulting configuration. We know that

$$(\mathcal{A}, \alpha_1) \sim_m^3 (\mathcal{B}, \beta_1) \ .$$

It follows by the inductive assumption that

$$(\mathcal{A}, \alpha_1) \sim_m (\mathcal{B}, \beta_1) \ ,$$

so Delilah wins the remaining $m$ moves of the game.

If $|\alpha| = |\beta| = 3$, then renumber the variables if necessary so that $(\mathcal{A}, \alpha)$ and $(\mathcal{B}, \beta)$ both satisfy $x_1 < x_2 < x_3$. Let $\alpha_\ell, \beta_\ell$ and $\alpha_r, \beta_r$ be the restrictions of $\alpha, \beta$ to the domains $\{x_1, x_2\}$ and $\{x_2, x_3\}$ respectively.

By Equation (6.34), Delilah wins the three-variable, $m + 1$-move games on these reduced configurations, that is,

$$(\mathcal{A}, \alpha_\ell) \sim_{m+1}^3 (\mathcal{B}, \beta_\ell) \quad \text{and} \quad (\mathcal{A}, \alpha_r) \sim_{m+1}^3 (\mathcal{B}, \beta_r) \ .$$

Since the domains of these configurations have size less than three, we know by the previous case that,

$$(\mathcal{A}, \alpha_\ell) \sim_{m+1} (\mathcal{B}, \beta_\ell) \quad \text{and} \quad (\mathcal{A}, \alpha_r) \sim_{m+1} (\mathcal{B}, \beta_r) \ .$$

We now combine Delilah's winning strategies for the games $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_\ell, \mathcal{B}, \beta_\ell)$ and $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_r, \mathcal{B}, \beta_r)$ to give her a winning strategy for the game $\mathcal{G}_{m+1}(\mathcal{A}, \alpha, \mathcal{B}, \beta)$. Notice that we are playing a game with an unlimited number of pebbles, so Samson need never reuse a pebble. See Figure 6.32.

Delilah's strategy is as follows: If Samson places a pebble to the left of pebble two, then Delilah answers according to her winning strategy in $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_\ell, \mathcal{B}, \beta_\ell)$.

If he places a pebble to the right of pebble two, then she answers according to her winning strategy in $\mathcal{G}_{m+1}(\mathcal{A}, \alpha_r, \mathcal{B}, \beta_r)$. After the $m + 1$ moves, Delilah has won both of the subgames. Thus, the map from the chosen points of $\mathcal{A}$ to the chosen points of $\mathcal{B}$ in the left subgame is an isomorphism, and similarly for the right subgame. Furthermore, all the chosen points in the left subgame are less than $x_2$ and all the chosen points in the right subgame are greater than $x_2$. Thus, the map from *all* the pebbled points in $\mathcal{A}$ to the pebbled points in $\mathcal{B}$ is an isomorphism and Delilah has won game $\mathcal{G}_{m+1}(\mathcal{A}, \alpha, \mathcal{B}, \beta)$. $\qquad\square$

**Exercise 6.35** Show that linear order does not have the two variable property. $\quad\square$

## 6.5   Zero-One Laws

In this section, we see that with very high probability, structures chosen at random are very simple from a first-order point of view. We begin by writing some sentences called "extension axioms". For any finite vocabulary, with no function or constant symbols, the extension axioms form a complete theory that is true in almost all structures.

The extension axioms can be written for any finite relational vocabulary. We first write them for undirected graphs. Consider the following sentence $\gamma_k$, whose meaning is that there are least $k - 1$ distinct vertices and any $k - 1$ tuple of distinct vertices may be extended to a $k$ tuple in any conceivable way:

$$
\begin{aligned}
\gamma_k \;\equiv\;& \\
& (\exists x_1 \ldots x_{k-1} \,.\, \mathrm{distinct}(x_1, \ldots, x_{k-1})) \quad \wedge \\
& (\forall x_1 \ldots x_{k-1} \,.\, \mathrm{distinct}(x_1, \ldots, x_{k-1})) \\
& \Big( (\exists x_k \,.\, \mathrm{distinct}(x_1, \ldots, x_k))(E(x_1, x_k) \wedge E(x_2, x_k) \wedge \cdots \wedge E(x_{k-1}, x_k)) \\
& \wedge\, (\exists x_k \,.\, \mathrm{distinct}(x_1, \ldots, x_k))(E(x_1, x_k) \wedge E(x_2, x_k) \wedge \cdots \wedge \neg E(x_{k-1}, x_k)) \\[4pt]
& \wedge \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad\qquad (6.36) \\[4pt]
& \wedge\, (\exists x_k \,.\, \mathrm{distinct}(x_1, \ldots, x_k))(E(x_1, x_k) \wedge E(x_2.x_k) \wedge \cdots \\
& \qquad\qquad\qquad \wedge\, E(x_{i-1}, x_k) \wedge \neg E(x_i, x_k) \wedge \cdots \wedge \neg E(x_{k-1}, x_k)) \\[4pt]
& \wedge \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\[4pt]
& (\exists x_k \,.\, \mathrm{distinct}(x_1, \ldots, x_k))(\neg E(x_1, x_k) \wedge \neg E(x_2, x_k) \wedge \cdots \wedge \neg E(x_{k-1}, x_k)) \Big)
\end{aligned}
$$

A simple counting argument shows that almost all graphs satisfy $\gamma_k$. Define $\mu_n(\varphi)$ to be the fraction of (ordered) structures of size $n$ that satisfy $\varphi$,

$$\mu_n(\varphi) = \frac{\left|\{G \mid n = \|G\|;\ G \models \varphi\}\right|}{\left|\{G \mid n = \|G\|\}\right|}$$

**Lemma 6.37**  *For any fixed $k > 0$,*

$$\lim_{n \to \infty} \mu_n(\gamma_k) = 1 \quad .$$

**Proof** Let $v_1, \ldots, v_{k-1}$ be a $k - 1$ tuple of distinct vertices from a random graph $G$ on $n$ vertices, let $x$ be any of the $n + 1 - k$ remaining vertices, and let $c$ be any of the $k$ conjuncts of the sentence $\gamma_k$. Conjunct $c$ asserts $k - 1$ independent conditions on the existence of edges from $x$, each of which has probability $1/2$. For this reason, the probability that $x$ does not meet condition $c$ for $v_1, \ldots, v_{k-1}$ is $\alpha = (1 - 2^{-(k-1)})$. Thus, the probability that none of the $(n + 1 - k)$ $x$'s satisfies condition $c$ is $\alpha^{n+1-k}$. It follows that the probability that $G$ does not satisfy $\gamma_k$ is less than

$$k \cdot n^{k-1} \alpha^{n+1-k} \ .$$

This expression goes quickly to $0$ as $n$ goes to infinity.  □

The sentence $\gamma_k$ says that any next move in the game $\mathcal{G}^k$ can be matched by Delilah. Thus we have the following,

**Lemma 6.38**  *Let $G$ and $H$ be undirected graphs satisfying $\gamma_k$. Then $G \sim^k H$.*

**Proof** We show by induction on $m$ that $G \sim^k_m H$. In the base case when $m = 0$, there are no chosen points so $G \sim^k_0 H$ holds vacuously.

Suppose that $G \sim^k_m H$ and let Delilah play the $m + 1$ move game as follows: For the first $m$ moves Delilah follows her winning strategy for the $m$ move game. Thus she has not lost yet. On the last move, suppose that Samson picks up pair $k$ of pebbles and places one of them on some vertex $v$ of $G$. We may assume that the previously pebbled points are all distinct (Exercise 6.12). Since $H \models \gamma_k$, there exists a vertex $v'$ of $H$ such that for all $j < k$, there is an edge from $v'$ to $\beta_m(x_j)$ in $H$ iff there is an edge from $v$ to $\alpha_m(x_j)$ in $G$. Thus, Delilah answers by putting her pebble $k$ on $v'$ and wins the game.  □

We can generalize $\gamma_k$ as follows: Let $\tau = \langle R_1^{a_1}, \ldots, R_r^{a_r} \rangle$ be a vocabulary with no constant symbols. Let $A$ be the set of all atomic formulas of the form $R_i(y_1, \ldots, y_{a_i})$ such that

$$x_k \in \{y_1, \ldots, y_{a_i}\} \subseteq \{x_1, \ldots, x_k\}$$

Define $\gamma_k(\tau)$ to be the following conjunction, which says that every $(k-1)$-tuple may be extended to a $k$-tuple in any conceivable way.

$$\gamma_k(\tau) \quad \equiv \quad (\forall x_1 \ldots x_{k-1} . \operatorname{distinct}(x_1, \ldots, x_{k-1}))$$
$$\bigwedge_{S \subseteq A} \Big( (\exists x_k . \operatorname{distinct}(x_1, \ldots, x_k)) (\bigwedge_{\alpha \in S} \alpha \ \wedge \bigwedge_{\alpha \in A - S} \neg \alpha) \Big)$$

It is easy to see that Lemmas 6.37 and 6.38 go through for any such $\gamma_k(\tau)$. The following theorem tells us that any property expressible by a set of sentences from $\mathcal{L}^k(\tau)$ is true in almost all structures, or false in almost all structures. This is sometimes known as the zero-one law for $\mathcal{L}_{\infty\omega}^\omega$.

**Theorem 6.39 (Zero-One Law)** *Let $S \subseteq \mathcal{L}^k$ be any set of $k$ variable sentences over a finite vocabulary $\tau$ with no constant or function symbols. Then the following limit exists and is equal to zero or one.*

$$\lim_{n \to \infty} \mu_n(S) \ .$$

**Proof** By Lemma 6.38, for every sentence $\varphi \in S$, either $\gamma_k \vdash \varphi$, or $\gamma_k \vdash \neg\varphi$. Thus Lemma 6.37 tells us that the above limit exists and (a) is equal to one if $\gamma_k$ implies every sentence in $S$ and (b) is equal to 0 if $\gamma_k$ implies the negation of some sentence in $S$.                                                                                                                       $\square$

**Corollary 6.40** *Assume that no constant symbols occur. Then a zero-one law holds for the language* FO(wo$\leq$)*. Furthermore, a zero-one law holds for all of the following languages:* FO(wo$\leq$)(TC)*,* FO(wo$\leq$)(LFP)*, and* FO(wo$\leq$)(PFP)*. (The operators* TC *and* PFP *are defined in Chapters 9 and 10, respectively.)*

**Proof** We see in later chapters that any sentence in one of these languages is equivalent to an infinite disjunction of sentences from $\mathcal{L}^k(\tau)$ for some $k$ and $\tau$. Since $\gamma_k$

determines the truth of any sentence in $\mathcal{L}^k(\tau)$, it also determines the truth of any infinite disjunction of such sentences. □

Suppose first-order logic has a zero-one law for the class $\mathcal{C}$ of structures. We see in Exercise 12.52 that this means that for each $k$, $\mathcal{L}^k$ has *bounded expressive power on average* in the following sense: There is a fixed bound $b$ such that almost all elements of $\mathcal{C}$ fall in one of $b$ $\mathcal{L}^k$-equivalence classes. When talking about typical structures, $\mathcal{L}^k$ can express only a bounded number of facts.

The zero-one laws do not hold for ordered structures or for structures with constants. This can be seen from the following equation,

$$\mu_n(E(0,0)) \quad = \quad \frac{1}{2} \, . \tag{6.41}$$

For ordered structures $\mathcal{A}$ and $\mathcal{B}$, if $\mathcal{A} \equiv^2 \mathcal{B}$, then $\|\mathcal{A}\| = \|\mathcal{B}\|$. Thus, for $k \geq 2$, $\mathcal{L}^k$ is not bounded.

## 6.6 Ehrenfeucht-Fraïssé Games with Ordering

When the ordering relation is present, Ehrenfeucht-Fraïssé games become much more difficult for Delilah. We present a clear explanation for this. Then we include a few game lower bounds for languages including ordering. It should not be surprising that the lower bounds are more difficult with ordering because here we are really talking about computation. For the same reason, such lower bounds are quite deserving of the effort required. Later, we show some more sophisticated lower bounds with ordering.

The following provides an upper bound on the complexity lower bounds we can prove using games for ordered structures.

**Proposition 6.42** *Let $G$ and $H$ be ordered graphs and let $n = \max(\|G\|, \|H\|)$. If $G \sim^3_{\lceil \log(n-1) \rceil + 1} H$, then $G = H$.*

**Proof** Suppose for the sake of contradiction that $G \sim^3_{\lceil \log(n-1) \rceil + 1} H$ but $G \neq H$. Let $n = \|G\|$ and $m = \|H\|$. Suppose $n > m$. Let PATH $<_d (x, y)$ mean that there is a path of length at least $d$ from $x$ to $y$, where each step in the path is given by the less than relation. Thus $G \models$ PATH $<_{n-1} (0, max)$, but $H \models \neg$PATH $<_{n-1}$ $(0, max)$. From Proposition 6.15 we know that PATH $<_n \in \mathcal{L}^3_{\lceil \log(n-1) \rceil}$. Thus, $n = \|G\| = \|H\|$.

Since $G \neq H$, there must be a pair of vertices $i, j$ such that there is an edge from vertex $i$ to vertex $j$ in one of the graphs but not the other. In the game $\mathcal{G}^3_{\lceil \log(n-1) \rceil + 1}(G, H)$, Samson should play the vertices $i$ and $j$ in $G$ in his first two moves. If Delilah answers with vertices $i$ and $j$ from $H$, then she loses immediately. If she does not answer with these vertices then $(G, \alpha_2)$ and $(H, \beta_2)$ disagree on a formula of the form $\text{PATH} <_d (x_k, c)$, for $k \in \{1, 2\}$, $c \in \{max, 0\}$, and $d \leq (n-1)/2$. It follows that Samson wins the remaining $\lceil \log(n-1) \rceil - 1$ move game. $\qquad \square$

**Exercise 6.43** Show that life is even worse when BIT is present: If $G$ and $H$ are ordered graphs that are $\mathcal{L}^2_{\lceil \log_*(n-1) \rceil + 2}$ equivalent in the presence of BIT, then $G = H$. (Recall that $\log_* n$ is the smallest $k$ such that $\log$ applied $k$ times to $n$ is at most 1.)

[Hint: in the presence of BIT, if Delilah matches vertex number $i$ from one graph, with vertex number $i' \neq i$ from the other graph, how long can she avoid a contradiction?] $\qquad \square$

We saw in Proposition 6.14 that without ordering, we needed $n + 1$ variables to say that a structure's universe is size exactly $n$ or that it has even cardinality. Thus, without ordering, even languages as strong as FO(LFP) cannot express very simple queries.

We now show that with ordering but still without BIT, quantifier rank $\log n$ is necessary to count even mod 2 (cf. Exercise 4.18 where BIT and thus counting mod 2, etc., are shown to be expressible in IND$[\log n]$).

**Proposition 6.44** *The sentence EVEN, meaning that the cardinality of the universe is even, is not expressible in quantifier rank* $\lceil \log(n-1) \rceil - 1$ *with ordering, but without* BIT.

**Proof** This proposition follows immediately from Lemma 6.16. Let $G$ be the line graph on $n = 2^{k+1} + 1$ and let $H$ be the line graph on $n+1$ vertices. Then $G$ and $H$ are $k$-equivalent but disagree on property EVEN. Note that $k = \lceil \log(n-1) \rceil - 1$. $\square$

**Corollary 6.45** *Boolean query* REACH *is not expressible in quantifier rank* $\lceil \log(n-1) \rceil - 1$ *with ordering, but without* BIT.

**Proof** Define $G_n$ and $G_{n+1}$ to be graphs that have the same universe and ordering relation as $L_n$ and $L_{n+1}$, respectively. Let $s = 0$ and $t = max$. Replace the edge predicate by the following relation, meaning that the points are two steps apart in the ordering,

$$E(x, y) \quad \equiv \quad (\exists z)(\text{SUC}(x, z) \ \wedge \ \text{SUC}(z, y)) .$$

Thus REACH holds for one of $G_n, G_{n+1}$ and not the other. However, $G_n$ and $G_{n+1}$ are still $\lceil \log(n-1) \rceil - 2$ equivalent because any win by Samson in $\mathcal{G}_r(G_n, G_{n+1})$ can be converted in one more move to a win by Samson in $\mathcal{G}_{r+1}(L_n, L_{n+1})$. □

Another corollary of Proposition 6.44 is that the language $(aa)^\star$, i.e., the set of even-length strings over a single letter alphabet, is not first-order without BIT. By Theorem 1.37, this is equivalent to the fact that $(aa)^\star$ is not a star-free regular language.

**Exercise 6.46** Prove that in the genealogical database, the ancestor relation is not a first-order query, cf. Examples 1.2, 1.27. This theorem holds even in the presence of the ordering relation. □

**Open Problem 6.47** Extend Corollary 6.45 to show that REACH is not expressible in quantifier rank $o[\log n]$ even in the presence of BIT. Class $\text{NC}_1$ is contained in $\text{IND}[\log n / \log \log n]$ (Theorem 5.35). Thus, a solution of this problem would prove the very interesting result $\text{NC}_1 \neq \text{NL}$.

## Historical Notes and Suggestions for Further Reading

Games $\mathcal{G}_m(\mathcal{A}, \mathcal{B})$ are called "Ehrenfeucht-Fraïssé games" in honor of their inventors [Ehr61, Fra54]. Barwise invented games that measure the number of variables used as well as the quantifier rank [Bar77]. Immerman reinvented this game, using pebbles [I80]. As seen in later chapters, variations of Ehrenfeucht-Fraïssé games for most logics have been developed. See [KV95] for the games that characterize languages with arbitrary generalized quantifiers. Originally, Samson was referred to as "Player I" and Delilah as "Player II". Joel Spencer coined the names "spoiler" and "duplicator" [ASE92]. From there it was a short step to the more memorable Samson and Delilah. Theorem 6.10 is due to the inventors of the various versions of these games.

We use Ehrenfeucht-Fraïssé games in later chapters to determine that certain queries are not expressible in first-order logic and other, more powerful, languages.

The "methodology" for proving lower bounds (Theorem 6.18) was adapted from a lecture by Phokion Kolaitis in [IKL95].

Thanks to Bill Gasarch for pointing out a simpler and more elegant proof of the lower bound in Propostion 6.15 than the one that appeared in the first printing of this book. Gasarch attributes the statement and proof of Lemma 6.16 to [Ros82].

Lemma 6.29 and Theorem 6.31 are due to Kozen and Immerman [IK87]. In addition, exact bounds are proved there on the numbers $k$ such that theories of bounded trees have the $k$-variable property (Fact 12.32).

Theorems 6.20 and 6.26 are from [FSV95]. Gaifman graphs (Definition 6.19) are named for Haim Gaifman. See [Gai81] for Gaifman's theorem on the locality of first-order properties.

Recently a significant strengthening of Gaifman's locality theorem for first-order logic was proved by Schwentick and Barthelmann [SB98]. They showed that every first-order formula is equivalent to a formula of the form $(\exists x_1 \cdots x_k)(\forall y)\varphi$ where $\varphi$ is $r$-local around $y$. That is, all quantification in $\varphi$ is restricted to elements of distance at most $r$ from $y$ in the Gaifman graph. Here $r$ depends only on $\varphi$.

Theorem 6.28 is from [See95]. See [AHU74] for a discussion of the unit cost RAM.

The Zero-One Law for First-Order Logic is due to Fagin [Fag73, Fag76]. Fagin's original proof used extension axioms (6.36). A good reference for Theorem 6.39 and Corollary 6.40 is [BGK85] by Blass, Gurevich and Kozen. See also [I80].

In our opinion, zero-one laws are inimical to computation. In order to compute we seem to need an ordering on the universe. As we have seen, having an ordering, or even having a single constant, eliminates the possibility of a zero-one law (6.41). Furthermore, the interested reader may look ahead to Fact 12.53, which shows that a language must be very weak indeed — on almost all structures — in order to support a zero-one law. Zero-one laws have very little attention in this book, but there is a large literature on the subject. See [Co88, KV92a, Spe93] for surveys.