

# An Existential Locality Theorem

Martin Grohe<sup>1</sup> and Stefan Wöhrle<sup>2</sup>

<sup>1</sup> Humboldt-Universität zu Berlin, Institut für Informatik, Unter den Linden 6, 10099 Berlin, Germany.

grohe@informatik.hu-berlin.de

<sup>2</sup> RWTH Aachen, Lehrstuhl für Informatik VII, 52056 Aachen, Germany.

woehrle@informatik.rwth-aachen.de

**Abstract.** We prove an existential version of Gaifman’s locality theorem and show how it can be applied algorithmically to evaluate existential first-order sentences in finite structures.

## 1 Introduction

Gaifman’s locality theorem [12] states that every first-order sentence is equivalent to a Boolean combination of sentences saying: There exist elements  $a_1, \dots, a_k$  that are far apart from one another, and each  $a_i$  satisfies some local condition described by a first-order formula whose quantifiers only range over a fixed-size neighborhood of an element of a structure. We prove that every *existential* first-order sentence is equivalent to a *positive* Boolean combination of sentences saying: There exist elements  $a_1, \dots, a_k$  that are far apart from one another, and each  $a_i$  satisfies some local condition described by an *existential* first-order formula.

The locality of first-order logic can be explored to prove that certain properties of finite structures are not expressible in first-order logic, and it seems that this was Gaifman’s main motivation. More recently, Libkin and others considered this technique of proving inexpressibility results using locality in a complexity theoretic context (see, e.g., [5, 14, 13, 15]).

A completely different application of Gaifman’s theorem has been proposed in [11]: It can be used to evaluate first-order sentences in certain finite structures quite efficiently. In general, it takes time  $n^{\Theta(l)}$  to decide whether a structure of size  $n$  satisfies a first-order sentence of size  $l$ , and under complexity theoretic assumptions, it can be proved that no real improvement is possible: The problem of deciding whether a given structure satisfies a given first-order sentence is PSPACE-complete [17, 19], and if parameterized by the size of the input sentence, it is complete for the parameterized complexity class AW[\*] [7]. The latter result implies that it is unlikely that the problem is fixed-parameter tractable (cf. [6]), i.e., that it can be solved in time  $f(l) \cdot n^c$ , for a function  $f$  and a constant  $c$ .

Gaifman’s theorem reduces the question of whether a first-order sentence holds in a structure to the question of whether the structure contains elements that are far apart from one another and satisfy some local condition expressed by a first-order formula. In certain structures, it is much easier to decide whether an element satisfies a local

first-order formula than to decide whether the whole structure satisfies a first-order sentence. An example are graphs of bounded degree: Local neighborhoods of vertices in such graphs have a size bounded by a constant only depending on the radius of the neighborhoods, so the time needed to check whether a vertex satisfies a local condition does not depend on the size of the graph. Another, less obvious example are planar graphs. To evaluate local conditions in planar graphs, we can exploit the fact that in planar graphs neighborhoods of fixed radius have bounded tree-width [16]. In general, such a locality based approach to evaluating first-order sentences in finite structures works for classes of structures that have a property called bounded local tree-width; the class of planar graphs and all classes of structures of bounded degree are examples of classes having this property. It has been proved in [11] that for each class  $C$  of structures of bounded local tree-width there is an algorithm that, given a structure  $\mathcal{A} \in C$  and a first-order sentence  $\varphi$ , decides whether  $\mathcal{A}$  satisfies  $\varphi$  in time near linear in the size of the structure  $\mathcal{A}$  (the precise statement is Theorem 7).

While a linear dependence on the size of the input structure is optimal, the dependence of these algorithms on the size of the input sentence leaves a lot to be desired: There is not even an elementary upper bound for the runtime in terms of the size of the sentence. Although the dependence of the algorithm on the structure size matters much more than the dependence on the size of the sentence, because usually we are evaluating small sentences in large structures,<sup>3</sup> it would be desirable to have a dependence on the size of the sentence that is not worse than exponential. Of course, since we are dealing with a PSPACE complete problem, we cannot expect the runtime of an algorithm to be polynomial in both the size of the input structure and the size of the input sentence.

We have observed that one of the main factors contributing to the enormous runtime of the locality based algorithms in terms of the formulas size is the number of quantifier alternations in the formula. This has motivated the present paper. We can use a variant of our existential locality theorem to improve the algorithms described above to algorithms whose runtime “only” depends doubly exponentially on the size of the input sentence.

In this paper we concentrate on the proof of our existential locality theorem, which is surprisingly complicated. This proof is presented in Section 3. The algorithmic application is outlined in Section 4.

## 2 Preliminaries

A *vocabulary* is a finite set of relation symbols. Associated with every relation symbol  $R$  is a positive integer called the *arity* of  $R$ . In the following,  $\tau$  always denotes a vocabulary.

A  $\tau$ -*structure*  $\mathcal{A}$  consists of a non-empty set  $A$ , called the *universe* of  $\mathcal{A}$ , and a relation  $R^{\mathcal{A}} \subseteq A^r$  for each  $r$ -ary relation symbol  $R \in \tau$ . For instance, we consider *graphs* as  $\{E\}$ -structures  $\mathcal{G} = (G, E^{\mathcal{G}})$ , where the binary relation  $E^{\mathcal{G}}$  is symmetric and anti-reflexive (i.e. graphs are undirected and loop-free). If  $\mathcal{A}$  is a  $\tau$ -structure and

---

<sup>3</sup> The generic example is the problem of evaluating SQL database queries against finite relational databases, which can be modeled by the problem of evaluating first-order sentences in finite structures.

$B \subseteq A$ , then  $\langle B \rangle^{\mathcal{A}}$  denotes the substructure induced by  $\mathcal{A}$  on  $B$ , that is, the  $\tau$ -structure  $B$  with universe  $B$  and  $R^B := R^{\mathcal{A}} \cap B^r$  for every  $r$ -ary  $R \in \tau$ .

The set of all free variables of a first-order formula  $\varphi$  is denoted by  $\text{free}(\varphi)$ . A *sentence* is a formula without free variables. The notation  $\varphi(x_1, \dots, x_k)$  indicates that all free variables of the formula  $\varphi$  are among  $x_1, \dots, x_k$ ; it does not necessarily mean that the variables  $x_1, \dots, x_k$  all appear in  $\varphi$ . The *weight* of a first-order formula  $\varphi$  is the number of quantifiers  $\exists x$  and  $\forall x$  occurring in  $\varphi$ .

A first-order formula is *existential* if it contains no universal quantifiers and if every existential quantifier occurs in the scope of an even number of negation symbols. A *literal* is an atom or a negated atom. A *conjunctive query with negation* is a formula of the form  $\exists \bar{x} \bigwedge_{i=1}^m \lambda_i$ , where each  $\lambda_i$  is a literal. Every existential formula  $\varphi$  of weight  $w$  and length  $l$  is equivalent to a disjunction of at most  $2^l$  conjunctive queries with negation, each of which is of weight at most  $w$  and length at most  $l$ .

We often denote tuples  $a_1 \dots a_k$  of elements of a set  $A$  by  $\bar{a}$ , and we write  $\bar{a} \in A$  instead of  $\bar{a} \in A^k$ . Similarly, we denote tuples of variables by  $\bar{x}$ .

## 2.1 Gaifman's Locality Theorem

The *Gaifman graph* of a  $\tau$ -structure  $\mathcal{A}$  is the graph  $\mathcal{G}_{\mathcal{A}}$  with vertex set  $A$  and an edge between two vertices  $a, b \in A$  if there exists an  $R \in \tau$  and a tuple  $a_1 \dots a_k \in R^{\mathcal{A}}$  such that  $a, b \in \{a_1, \dots, a_k\}$ . The *distance*  $d^{\mathcal{A}}(a, b)$  between two elements  $a, b \in A$  of a structure  $\mathcal{A}$  is the length of the shortest path in  $\mathcal{G}_{\mathcal{A}}$  connecting  $a$  and  $b$ . For  $r \geq 1$  and  $a \in A$ , we define the  *$r$ -neighborhood* of  $a$  in  $\mathcal{A}$  to be  $N_r^{\mathcal{A}}(a) := \{b \in A \mid d^{\mathcal{A}}(a, b) \leq r\}$ . For a subset  $B \subseteq A$  we let  $N_r^{\mathcal{A}}(B) := \bigcup_{b \in B} N_r^{\mathcal{A}}(b)$ .

For every  $r \geq 0$  there is an existential first-order formula  $\delta_r(x, y)$  such that for all  $\tau$ -structures  $\mathcal{A}$  and  $a, b \in A$  we have  $\mathcal{A} \models \delta_r(a, b)$  if, and only if,  $d^{\mathcal{A}}(a, b) \leq r$ . In the following, we write  $d(x, y) \leq r$  instead of  $\delta_r(x, y)$  and  $d(x, y) > r$  instead of  $\neg \delta_r(x, y)$ .

If  $\varphi(x)$  is a first-order formula, then  $\varphi^{N_r(x)}(x)$  is the formula obtained from  $\varphi(x)$  by relativizing all quantifiers to  $N_r(x)$ , that is, by replacing every subformula of the form  $\exists y \psi(x, y, \bar{z})$  by  $\exists y (d(x, y) \leq r \wedge \psi(x, y, \bar{z}))$  and every subformula of the form  $\forall y \psi(x, y, \bar{z})$  by  $\forall y (d(x, y) \leq r \rightarrow \psi(x, y, \bar{z}))$ . We usually write  $\exists y \in N_r(x) \psi$  instead of  $\exists y (d(x, y) \leq r \wedge \psi)$  and  $\forall y \in N_r(x) \psi$  instead of  $\forall y (d(x, y) \leq r \rightarrow \psi)$ .

A formula  $\psi(x)$  of the form  $\varphi^{N_r(x)}(x)$ , for some  $\varphi(x)$ , is called  *$r$ -local*. The basic property of  $r$ -local formulas  $\psi(x)$  is that it only depends on the  $r$ -neighborhood of  $x$  whether they hold at  $x$  or not, that is, for all structures  $\mathcal{A}$  and  $a \in A$  we have  $\mathcal{A} \models \psi(a)$  if, and only if,  $\langle N_r^{\mathcal{A}}(a) \rangle \models \psi(a)$ . Observe that if  $\psi(x)$  is  $r$ -local and  $s > r$ , then  $\psi(x)$  is equivalent to the  $s$ -local formula  $\psi^{N_s(x)}(x)$ . We often use this observation implicitly when considering  $r$ -local formulas as  $s$ -local for some  $s > r$ .

Sentences can never be local in the sense just defined. As a substitute, we say that a *local sentence* is a sentence of the form

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi(x_i) \right), \quad (1)$$

where  $r, k \geq 1$  and  $\psi(x)$  is  $r$ -local.

**Theorem 1 (Gaifman [12]).** *Every first-order sentence is equivalent to a Boolean combination of local sentences.*

### 3 The Existential Locality Theorems

If  $\psi(x)$  is an existential first-order formula, then for every  $r \geq 1$  the  $r$ -local formula  $\psi^{N_r(x)}(x)$  obtained from  $\psi$  is also existential. We define a local sentence as in (1) to be *existential* if the formula  $\psi$  is existential and  $r$ -local. Let us remark that, in general, an existential local sentence is *not* equivalent to an existential first-order sentence, because the formula  $d(x_i, x_j) > s$  is not existential for any  $s \geq 2$ .

**Theorem 2.** *Every existential first-order sentence is equivalent to a positive Boolean combination of existential local sentences.*

Unfortunately, neither Gaifman’s original proof of his locality theorem (based on quantifier elimination) nor Ebbinghaus and Flum’s [8] model theoretic proof can be adapted to prove this existential version of Gaifman’s theorem. Compared to these proofs, our proof is very combinatorial, which is not surprising, because there is not much “logic” left in existential sentences.

We illustrate the basic idea by a simple example:

**Example 3.** Let

$$\varphi := \exists x \exists y (\neg E(x, y) \wedge \text{RED}(x) \wedge \text{BLUE}(y))$$

(here  $E$  is a binary relation symbol and RED, BLUE are unary relation symbols). Although the syntactical form of  $\varphi$  is close to that of an existential local sentence, it is not obvious how to find a positive Boolean combination of existential local sentences equivalent to  $\varphi$ . Here is one:

$$\begin{aligned} & \exists x \exists x' \in N_2(x) \exists y \in N_2(x) (\neg E(x', y) \wedge \text{RED}(x') \wedge \text{BLUE}(y)) \\ & \vee \left( \exists x \exists y (d(x, y) > 2 \wedge (\text{RED}(x) \vee \text{BLUE}(y)) \wedge (\text{RED}(y) \vee \text{BLUE}(y))) \right) \\ & \wedge \exists x \text{RED}(x) \wedge \exists x \text{BLUE}(x) \end{aligned}$$

To understand the following proof it is worthwhile trying to extend the idea of this example to the sentence

$$\exists x \exists y \exists z (\neg E(x, y) \wedge \neg E(x, z) \wedge \neg E(y, z) \wedge \text{RED}(x) \wedge \text{BLUE}(y) \wedge \text{GREEN}(z))$$

(although it is very complicated to actually write down an equivalent positive Boolean combination of existential local sentences). Indeed, it is the main difficulty of the proof to handle sentences saying “there is an independent set of points  $x_1, \dots, x_k$  of colors  $c_1, \dots, c_k$ , respectively.” Playing with such sentences leads to the crucial observation that the basic combinatorial problem can be handled by Hall’s theorem (as it is done in Step 4 of the proof of Lemma 4).

The proof requires some preparation. We define the *rank* of a local sentence

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi(x_i) \right),$$

to be the pair  $(k + w, r)$ , where  $w$  is the weight of  $\psi$ . We partially order the ranks by saying that  $(q, r) \leq (q', r')$  if  $q \leq q'$  and  $r \leq r'$ .

**Lemma 4.** *Let  $k \geq 2$ ,  $r \geq 1$ ,  $w \geq 0$ , and let  $\mathcal{A}, \mathcal{B}$  be structures such that every existential local sentence of rank at most  $(k \cdot (w + 1), 2^{k^2} r)$  that holds in  $\mathcal{A}$  also holds in  $\mathcal{B}$ . Let*

$$\varphi := \exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2^{k^2} r \wedge \bigwedge_{i=1}^k \psi_i(x_i) \right),$$

where for  $1 \leq i \leq k$ , the formula  $\psi_i(x_i)$  is  $r$ -local, existential and of weight at most  $w$ . Suppose that  $\mathcal{A} \models \varphi$ .

Then

$$\mathcal{B} \models \exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{i=1}^k \psi_i(x_i) \right).$$

*Proof:* We prove the lemma in four steps.

*Step 1.* We show that if for some  $l, 1 \leq l \leq k$ , say  $l = k$ , there are  $b_1, \dots, b_k \in B$  such that  $d(b_i, b_j) > 4r$  for  $1 \leq i < j \leq k$ , and  $\mathcal{B} \models \psi_l(b_i)$  for  $1 \leq i \leq k$ , then it suffices to prove that

$$\mathcal{B} \models \exists x_1 \dots \exists x_{k-1} \left( \bigwedge_{1 \leq i < j \leq k-1} d(x_i, x_j) > 2r \wedge \bigwedge_{i=1}^{k-1} \psi_i(x_i) \right).$$

To see this, suppose that we have such  $b_1, \dots, b_k$  and we find  $c_1, \dots, c_{k-1}$  such that  $d(c_i, c_j) > 2r$  for all  $1 \leq i < j \leq k-1$ , and  $\mathcal{B} \models \psi_i(c_i)$  for  $1 \leq i \leq k-1$ . Then there will be at least one  $i, 1 \leq i \leq k$  such that  $b_i$  has distance greater than  $2r$  from  $c_j$  for all  $j, 1 \leq j \leq k-1$ . Thus  $c_1, \dots, c_{k-1}, b_i$  witness that

$$\mathcal{B} \models \exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{i=1}^k \psi_i(x_i) \right).$$

So without loss of generality, in the following we assume that for  $1 \leq i \leq k$ , there are at most  $(k-1)$  elements of  $B$  of pairwise distance greater than  $4r$  satisfying  $\psi_i$ .

*Step 2.* We let  $K := \{1, \dots, k\}$ , and for every set  $I \subseteq K$  we let  $\psi_I(x) := \bigvee_{i \in I} \psi_i(x)$ . Note that  $\psi_I$  is a formula of weight at most  $k \cdot w$ . Let  $C := \{c \in B \mid \mathcal{B} \models \psi_K(c)\}$ . By the assumption we made at the end of Step 1, there exist at most  $k(k-1)$  elements of  $C$  of pairwise distance greater than  $4r$ .

*Claim:* There are  $p, l, 1 \leq p \leq k(k-1) + 1, 1 \leq l \leq k(k-1)$ , and elements  $c_1, \dots, c_l \in C$  such that  $d^{\mathcal{B}}(c_i, c_j) > 2^{p+1}r$  for  $1 \leq i < j \leq l$ , and for all  $c \in C$  there exists an  $i \leq l$  such that  $d^{\mathcal{B}}(c, c_i) \leq 2^p r$ .

*Proof:* We construct  $c_1, \dots, c_l$  inductively: As the inductive basis, let  $c_1$  be an arbitrary element of  $C$ . If  $c_1, \dots, c_i$  are constructed, we choose  $c_{i+1} \in C$  such that for  $1 \leq j \leq i$  we have  $d^{\mathcal{B}}(c_{i+1}, c_j) > 2^{k(k-1)+1-(i-1)}r$ . If no such  $c_{i+1}$  exists, we let  $l := i, p := k(k-1) + 1 - (l-1)$  and stop.

Our construction guarantees that for  $1 \leq i < j \leq l$  we have

$$d^{\mathcal{B}}(c_i, c_j) > 2^{k(k-1)+1-(j-2)}r. \quad (2)$$

For  $j \leq k(k-1) + 1$ , this implies  $d^{\mathcal{B}}(c_i, c_j) > 4r$ . Since there are at most  $k(k-1)$  elements of  $C$  of pairwise distance greater than  $4r$ , this guarantees that  $l \leq k(k-1)$ . (2) also guarantees that for  $1 \leq i < j \leq l$  we have  $d^{\mathcal{B}}(c_i, c_j) > 2^{k(k-1)+1-(l-2)}r = 2^{p+1}r$ .

Since we stopped at  $l = i$ , for all  $c \in C$  there exists an  $i \leq l$  such that  $d^{\mathcal{B}}(c, c_i) \leq 2^{k(k-1)+1-(l-1)} = 2^p r$ . This proves the claim.

*Step 3.* Let  $p, l, c_1, \dots, c_l$  be as stated in the claim in Step 2. For  $I \subseteq K$ , let

$$\varphi_I := \exists x_1 \dots \exists x_k \left( \bigwedge_{\substack{i, j \in I \\ i < j}} d(x_i, x_j) > 2^{p+1}r \wedge \bigwedge_{i \in I} \psi_I(x_i) \right).$$

Since  $\mathcal{A} \models \varphi$ , we have  $\mathcal{A} \models \varphi_I$ . Thus, since  $\varphi_I$  is an existential local sentence of rank at most  $(k \cdot (w+1), 2^{k^2}r)$ , we also have  $\mathcal{B} \models \varphi_I$ .

*Step 4.* Let  $L := \{1, \dots, l\}$ . We define a relation  $R \subseteq K \times L$  as follows: For  $i \in K, j \in L$  we let  $iRj$  if there is a  $b \in B$  such that  $\mathcal{B} \models \psi_i(b)$  and  $d^{\mathcal{B}}(b, c_j) \leq 2^p r$ .

*Claim:* For every  $I \subseteq K$  the set  $R(I) := \{j \in L \mid \exists i \in I : iRj\}$  contains at least as many elements as  $I$ .

*Proof:* Recall that  $\mathcal{B} \models \varphi_I$ . For  $i \in I$ , let  $b_i \in B$ , such that for all  $i, j \in I$  with  $i < j$  we have  $d^{\mathcal{B}}(b_i, b_j) > 2^{p+1}r$  and for all  $i \in I$  we have  $\mathcal{B} \models \psi_i(b_i)$ . Then  $b_i \in C$ , and thus there exist a  $j \in L$  such that  $d^{\mathcal{B}}(b_i, c_j) \leq 2^p r$ . Since  $d^{\mathcal{B}}(b_i, b_j) > 2^{p+1}r$ , for every  $j \in L$  there can be at most one  $i \in I$  such that  $d^{\mathcal{B}}(b_i, c_j) \leq 2^p r$ . This proves the claim.

By Hall's theorem, there exists a one-to-one mapping  $f$  of  $K$  into  $L$  such that for all  $i \in K$  we have  $iRf(i)$ . In other words, there exist  $b_1, \dots, b_k$  such that for  $1 \leq i \leq k$  we have  $\mathcal{B} \models \psi_i(b_i)$  and  $d^{\mathcal{B}}(b_i, c_{f(i)}) \leq 2^p r$ . Since  $d^{\mathcal{B}}(c_{f(i)}, c_{f(j)}) > 2^{p+1}r$ , the latter implies  $d^{\mathcal{B}}(b_i, b_j) > 2r$ . Thus

$$\mathcal{B} \models \exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{i=1}^k \psi_i(x_i) \right).$$

□

**Lemma 5.** *There is a function  $f(k)$ , such that the following holds for all  $k \geq 1$ : Let  $\mathcal{A}, \mathcal{B}$  be structures such that every existential local sentence of rank  $(k(k+1), f(k))$  that holds in  $\mathcal{A}$  also holds in  $\mathcal{B}$ . Then every existential sentence of weight at most  $k$  that holds in  $\mathcal{A}$  also holds in  $\mathcal{B}$ .*

*Proof:* Since every existential sentence is equivalent to a disjunction of conjunctive queries with negation of the same weight, it suffices to prove that every conjunctive query with negation of weight  $k$  that holds in  $\mathcal{A}$  also holds in  $\mathcal{B}$ . Let

$$\varphi := \exists x_1 \dots \exists x_k \psi(x_1, \dots, x_k)$$

with

$$\psi(x_1, \dots, x_k) := \left( \bigwedge_{i=1}^p \alpha_i \wedge \bigwedge_{i=1}^q \beta_i \right),$$

where all the  $\alpha_i$  are atoms and the  $\beta_i$  are negated atoms. Suppose that  $\mathcal{A} \models \varphi$ . We shall prove that  $\mathcal{B} \models \varphi$ .

We define the *positive graph of  $\varphi$*  to be the graph  $\mathcal{G}$  with universe  $G := \text{var}(\varphi) = \{x_1, \dots, x_k\}$  and

$$E^{\mathcal{G}} := \{xy \mid \exists i, 1 \leq i \leq p : x, y \in \text{var}(\alpha_i)\}.$$

Let  $\mathcal{H}_1, \dots, \mathcal{H}_r$  be the connected components of  $\mathcal{G}$ . Without loss of generality, we may assume that for  $1 \leq i \leq r$  we have  $x_i \in H_i$  (recall that  $H_i$  denotes the universe of the structure  $\mathcal{H}_i$ ). Then we know that  $H_i \subseteq N_k^{\mathcal{G}}(x_i)$ . If  $r = 1$ , then this means that  $\text{var}(\varphi) \subseteq N_k^{\mathcal{G}}(x_1)$ , and  $\varphi$  is equivalent to the  $k$ -local sentence

$$\exists x_1 \exists x_2 \in N_k(x_1) \dots \exists x_k \in N_k(x_1) \psi$$

of rank  $(k, k)$ . If we choose  $f$  such that  $f(k) \geq k$ , then  $\mathcal{A} \models \varphi$  implies  $\mathcal{B} \models \varphi$ . In the following, we assume that  $r \geq 2$ .

Let  $c_0 := 0$  and  $c_{i+1} := 2^{k^2}(c_i + k + 1)$  for  $i \geq 0$ . We let  $R := \{\{i, j\} \mid 1 \leq i < j \leq r\}$ ,  $h := |R| = \binom{r}{2}$  and

$$f(k) = 2^{k^2}(c_h + k + 1). \quad (3)$$

For  $\bar{a} = a_1 \dots a_r \in A^r$ , the *distance pattern* of  $\bar{a}$  is the mapping  $\Delta_{\bar{a}} : R \rightarrow \{0, \dots, h\}$  defined by

$$\Delta_{\bar{a}}(\{i, j\}) := \begin{cases} 0 & \text{if } d^A(a_i, a_j) = 0 \\ t & \text{if } c_t < d^A(a_i, a_j) \leq c_{t+1} \text{ for some } t \text{ such that } 0 \leq t < h \\ h & \text{if } d^A(a_i, a_j) > c_h \end{cases}$$

By the pigeonhole principle, for every distance pattern  $\Delta$  there is an integer  $\text{gap}(\Delta)$  such that  $0 \leq \text{gap}(\Delta) \leq h$  and  $\Delta(\{i, j\}) \neq \text{gap}(\Delta)$  for all  $\{i, j\} \in R$ .

Let  $\bar{a} = a_1 \dots a_k \in A^k$  such that  $\mathcal{A} \models \psi(\bar{a})$ . Let  $\Delta := \Delta_{a_1 \dots a_r}$ , and  $g := \text{gap}(\Delta)$ . Then for all  $\{i, j\} \in R$  we either have  $d(a_i, a_j) \leq c_g$  or  $d(a_i, a_j) > 2^{k^2}(c_g + k + 1)$ . This implies that the relation on  $\{a_1, \dots, a_r\}$  defined by  $d^A(a_i, a_j) \leq c_g$  is an equivalence relation. Without loss of generality, we may assume that  $a_1, \dots, a_s$  form a system of representatives of the equivalence classes.

We let  $l := c_g + k$ . For  $1 \leq i \leq s$ , we let  $I_i := \{j \mid 1 \leq j \leq k, d^A(a_i, a_j) \leq l\}$ . Then  $(I_i)_{1 \leq i \leq s}$  is a partition of  $\{1, \dots, k\}$ . To see this, first recall that for  $1 \leq j \leq r$

there is an  $i, 1 \leq i \leq s$  such that  $d^A(a_i, a_j) \leq c_g$ . For  $t$  with  $r+1 \leq t \leq k$  there exist a  $j, 1 \leq j \leq r$  such that  $x_t \in H_j$ , the connected component of  $x_j$  in the positive graph of  $\varphi$ . Since  $\mathcal{A} \models \psi(\bar{a})$ , this implies that  $d^A(a_j, a_t) \leq k$ . Thus there exists an  $i, 1 \leq i \leq s$  such that  $d^A(a_i, a_t) \leq c_g + k$ .

For  $1 \leq i \leq s$ , we let

$$\psi_i(x_i) := \exists \bar{x}^i \in N_l(x_i) \bigwedge_{\text{var}(\alpha_i) \subseteq I_i} \alpha_i \wedge \bigwedge_{\text{var}(\beta_i) \subseteq I_i} \beta_i, \quad (4)$$

where  $\bar{x}^i$  consists of all variables  $x_j$  with  $j \in I_i \setminus \{i\}$ . Then for  $1 \leq i \leq s$  we have  $\mathcal{A} \models \psi_i(a_i)$ , because  $\mathcal{A} \models \psi(\bar{a})$ . Thus

$$\mathcal{A} \models \exists x_1 \dots \exists x_s \left( \bigwedge_{1 \leq i < j \leq s} d(x_i, x_j) > 2^{k^2}(l+1) \wedge \bigwedge_{1 \leq i \leq s} \psi_i(x_i) \right).$$

Since  $f(k) = 2^{k^2}(c_h + k + 1) \geq 2^{k^2}(l+1)$ , by Lemma 4, this implies

$$\mathcal{B} \models \exists x_1 \dots \exists x_s \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2(l+1) \wedge \bigwedge_{1 \leq i \leq s} \psi_i(x_i) \right).$$

Thus there exist  $b_1, \dots, b_s \in B$  such that for  $1 \leq i < j \leq s$  we have  $d^B(b_i, b_j) > 2(l+1)$  and for  $1 \leq i \leq s$  we have  $\mathcal{B} \models \psi_i(b_i)$ . Since  $I_1, \dots, I_s$  is a partition of  $\{1, \dots, k\}$ , there are  $b_{s+1}, \dots, b_k \in B$  such that:

- (i)  $d^B(b_i, b_j) \leq l$  for all  $j \in I_i$ .
- (ii)  $\mathcal{B} \models \alpha_j(\bar{b})$  for all  $j, 1 \leq j \leq p$  such that  $\text{var}(\alpha_j) \subseteq I_j$ .
- (iii)  $\mathcal{B} \models \beta_j(\bar{b})$  for all  $j, 1 \leq j \leq q$  such that  $\text{var}(\beta_j) \subseteq I_j$ .

We claim that  $\mathcal{B} \models \psi(\bar{b})$ . Since for each connected component  $H_j$  of the positive graph of  $\varphi$  there is an  $i, 1 \leq i \leq s$  such that  $t \in I_i$  whenever  $x_t \in H_j$ , (ii) implies that  $\mathcal{B} \models \alpha_j(\bar{b})$  for  $1 \leq j \leq p$ . It remains to prove that  $\mathcal{B} \models \beta_j(\bar{b})$  for  $1 \leq j \leq q$ . If  $\text{var}(\beta_j) \subseteq I_i$  for some  $i$ , then  $\mathcal{B} \models \beta_j(\bar{b})$  by (iii). Otherwise,  $\beta_j$  has variables  $x_u, x_v$  such that there exist  $i \neq i'$  with  $x_u \in I_i, x_v \in I_{i'}$ . Then by (i),  $d^B(b_i, b_u) \leq l$  and  $d^B(b_{i'}, b_v) \leq l$ . Since  $d^B(b_i, b_{i'}) > 2l+1$ , this implies  $d^B(b_u, b_v) > 1$ . Since  $\beta_j$  is a negated atom, this implies  $\mathcal{B} \models \beta_j(\bar{b})$ .

Thus  $\mathcal{B} \models \varphi$ . □

*Proof (of Theorem 2):* Let  $\varphi$  be an existential sentence of weight  $k$  and  $\mathcal{K} := \{\mathcal{A} \mid \mathcal{A} \models \varphi\}$  the class of all structures satisfying  $\varphi$ . Let  $\Psi$  be the set of all existential local sentences of rank at most  $(k(k+1), f(k))$ , where  $f$  is the function from Lemma 5. Let

$$\varphi' := \bigvee_{\mathcal{A} \in \mathcal{K}} \bigwedge_{\substack{\psi \in \Psi \\ \mathcal{A} \models \psi}} \psi.$$

We claim that  $\varphi$  is equivalent to  $\varphi'$ . The forward implication is trivial, and the backward implication follows from Lemma 5. Since up to logical equivalence, the set  $\Psi$  is finite



and therefore  $\varphi'$  contains at most  $2^{|\Psi|}$  non-equivalent disjuncts, this proves the theorem.  $\square$

Our proof of the existential version of Gaifman's theorem does not give us good bounds on the size and rank of the local formulas to which we translate a given existential formula. Therefore, for the algorithmic applications, it is preferable to work with the following weaker version of Theorem 2, which gives us better bounds.

An *asymmetric local sentence* is a sentence  $\varphi$  of the form

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi_i(x_i) \right),$$

where  $r, k \geq 1$  and  $\psi_1(x), \dots, \psi_k(x)$  are  $r$ -local.  $\varphi$  is an *existential asymmetric local sentence*, if in addition  $\psi_1(x), \dots, \psi_k(x)$  are existential.

An  *$r$ -local conjunctive query with negation*, for some  $r \geq 1$ , is a formula  $\psi(x)$  of the form  $\exists y_1 \in N_r(x) \dots \exists y_n \in N_r(x) \bigwedge_{i=1}^m \lambda_i$ , where each  $\lambda_i$  is a literal.

**Theorem 6.** *Every existential first-order sentence  $\varphi$  is equivalent to a disjunction  $\varphi'$  of existential asymmetric local sentences.*

*More precisely, if  $k$  is the weight of  $\varphi$  and  $l$  its size, then  $\varphi'$  is a disjunction of  $2^{O(l+k^4)}$  asymmetric local sentences of the form*

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi_i(x_i) \right),$$

where  $\psi_1, \dots, \psi_k$  are  $r$ -local conjunctive queries with negation. The rank of each of these local sentences is at most  $(k, 2^{k^2+1})$ , and their size is in  $O(l)$ .

Furthermore, there is a polynomial  $p$  and an algorithm translating  $\varphi$  to  $\varphi'$  in time  $O(2^{p(l)})$ .

*Proof:* We first assume that  $\varphi$  is a conjunctive query with negation, say,

$$\varphi := \exists x_1 \dots \exists x_k \psi(x_1, \dots, x_k)$$

with

$$\psi(x_1, \dots, x_k) := \left( \bigwedge_{i=1}^p \alpha_i \wedge \bigwedge_{i=1}^q \beta_i \right),$$

where all the  $\alpha_i$  are atoms and the  $\beta_i$  are negated atoms. Without loss of generality, we may assume that  $k \geq 2$ , because for  $k = 1$  there is nothing to prove. Let  $\mathcal{H}_1, \dots, \mathcal{H}_r$  be the connected components of the positive graph of  $\mathcal{G}$  (cf. proof of Lemma 5). Again we may assume that  $r \geq 2$ , and that for  $1 \leq i \leq r$  we have  $x_i \in H_i$ . Hence we know that  $H_i \subseteq N_k^{\mathcal{G}}(x_i)$ .

Let  $c_0 := 0$  and  $c_{i+1} := 2(c_i + k + 1)$  for  $i \geq 0$ . Let  $R := \{\{i, j\} \mid 1 \leq i < j \leq r\}$  and  $h := |R| = \binom{r}{2}$ . It is not difficult to prove that  $c_h + k + 1 \leq 2^{k^2+1}$ .

Let  $\bar{a} \in A^r$ ,  $\Delta := \Delta_{\bar{a}}$  the distance pattern of  $\bar{a}$ , and  $g := \text{gap}(\Delta)$ . The relation on  $\{a_1, \dots, a_r\}$  defined by  $d^A(a_i, a_j) \leq c_g$  is an equivalence relation. Without loss

of generality, we may assume that  $a_1, \dots, a_s$  form a system of representatives of the equivalence classes.

Now suppose that we extend  $a_1 \dots a_r$  to a  $k$ -tuple  $\bar{a} = a_1 \dots a_k \in A^k$  such that  $\mathcal{A} \models \psi(\bar{a})$ . Let  $l := c_g + k$  and  $I_i, \psi_i(x_i)$  as in the proof of Lemma 5. Define

$$\psi_\Delta(x_1, \dots, x_s) := \bigwedge_{1 \leq i < j \leq s} d(x_i, x_j) > 2(l+1) \wedge \bigwedge_{1 \leq i \leq s} \psi_i(x_i).$$

Then  $\mathcal{A} \models \psi_\Delta(a_1, \dots, a_s)$ . Furthermore, for every tuple  $a'_1 \dots a'_s \in A^s$  with  $\mathcal{A} \models \psi_\Delta(a'_1, \dots, a'_s)$  there exists an extension  $\bar{a}' := a'_1 \dots a'_k$  such that  $\mathcal{A} \models \psi(\bar{a}')$ . To see this, observe that every positive literal  $\alpha_j$  occurs in an  $I_i$  and thus in  $\psi_\Delta$ , and every negative literal  $\beta_j$  either occurs in an  $I_i$  or has variables with indices in two distinct  $I_i, I_{i'}$  and is thus automatically satisfied, because the variables are forced to be far apart.

The formula  $\psi_\Delta(x_1, \dots, x_k)$  only depends on the distance pattern  $\Delta$  and not on the tuple  $\bar{a}$  realizing it. So for every distance pattern  $\Delta$  we obtain a formula  $\psi_\Delta(\bar{x}^\Delta)$ , whose free variables  $\bar{x}^\Delta$  are among  $x_1, \dots, x_r$ , with the following properties:

- $\exists \bar{x}^\Delta \psi_\Delta$  is an existential asymmetric local sentence of rank at most  $(k, 2^{k^2+1})$ .
- For every tuple  $\bar{a} \in A^k$  with  $\mathcal{A} \models \psi(\bar{a})$  and  $\Delta_{\bar{a}} = \Delta$  we have  $\mathcal{A} \models \psi_\Delta(\bar{a}^\Delta)$ , where  $\bar{a}^\Delta$  consists of the same entries of  $\bar{a}$  as  $\bar{x}^\Delta$  of  $\bar{x}$ .
- Every tuple  $\bar{a}^\Delta$  with  $\mathcal{A} \models \psi_\Delta(\bar{a}^\Delta)$  can be extended to a tuple  $\bar{a} = a_1 \dots a_k$  such that  $\mathcal{A} \models \psi(\bar{a})$ .

The last two items imply that  $\varphi$  is equivalent to the formula

$$\varphi' := \bigvee_{\Delta \text{ distance pattern}} \exists \bar{x}^\Delta \psi_\Delta.$$

It is not hard to see that the number of distance pattern is in  $2^{O(k^4)}$ , thus  $\varphi$  is a disjunction of  $2^{O(k^4)}$  existential asymmetric local sentences of rank at most  $(k, 2^{k^2+1})$  and size in  $O(l)$  (where  $l$  denotes the length of  $\varphi$ ).

If  $\varphi$  is an arbitrary existential sentence, we first transform it to a disjunction of at most  $2^l$  conjunctive queries with negation of the same weight as  $\varphi$ .

Finally, we observe that the translation from  $\varphi$  to the disjunction of asymmetric local formulas is effective within the desired time bound: Given  $\varphi$ , we first translate it to a disjunction of conjunctive queries with negation. This is possible in time  $2^{O(l)}$ . Then we treat each of the conjunctive queries with negation separately. We compute the positive graph and all possible patterns. For each of pattern  $\Delta$ , we compute the gap and then the formula  $\varphi_\Delta$ . Since  $k \leq l$ , this is clearly possible in time  $2^{p(l)}$  for a suitable polynomial  $p$ .  $\square$

## 4 An Algorithmic Application

Our underlying model of computation is the standard RAM-model with addition and subtraction as arithmetic operations (cf. [1, 18]). In our complexity analysis we use

the uniform cost measure. Structures are represented on a RAM in a straightforward way by listing all elements of the universe and then all tuples in the relations. For details we refer the reader to [10]. We define the *size* of a  $\tau$ -structure  $\mathcal{A}$  to be  $\|\mathcal{A}\| := |A| + \sum_{R \in \tau} r\text{-ary } R \cdot |R^{\mathcal{A}}|$ ; this is the length of a reasonable representation of  $\mathcal{A}$  (if we suppress details that are inessential for us). We fix some reasonable encoding for first-order formulas and denote by  $\|\varphi\|$  the size of the encoding of a formula  $\varphi$ .

The appropriate structural notion for the algorithmic applications of locality is *bounded local tree-width*. We assume that the reader is familiar with the definition of *tree-width* of graphs (see e.g. [4]). The tree-width of a structure  $\mathcal{A}$ , denoted by  $\text{tw}(\mathcal{A})$ , is the tree-width of its Gaifman graph. The *local tree-width* of a structure  $\mathcal{A}$  is the function  $\text{ltw}_{\mathcal{A}} : \mathbb{N} \rightarrow \mathbb{N}$  defined by

$$\text{ltw}_{\mathcal{A}}(r) := \max \left\{ \text{tw}(\langle N_r^{\mathcal{A}}(a) \rangle) \mid a \in A \right\}.$$

A class  $\mathcal{C}$  of structures has *bounded local tree-width* if there is a function  $\lambda : \mathbb{N} \rightarrow \mathbb{N}$  such that  $\text{ltw}_{\mathcal{A}}(r) \leq \lambda(r)$  for all  $\mathcal{A} \in \mathcal{C}$ ,  $r \in \mathbb{N}$ . Many well-known classes of structures have bounded local tree-width, among them the class of planar graphs and all classes of structures of bounded degree.

**Theorem 7 (Frick and Grohe [11]).** *Let  $\mathcal{C}$  be a class of structures of bounded local tree-width. Then there is a function  $f$  and, for every  $\epsilon > 0$ , an algorithm deciding in time  $O(f(\|\varphi\|)|A|^{1+\epsilon})$  whether a given structure  $\mathcal{A} \in \mathcal{C}$  satisfies a given first-order sentence  $\varphi$ .*

The algorithm proceeds as follows: Given a structure  $\mathcal{A}$  and a sentence  $\varphi$ , they first translate  $\varphi$  to a Boolean combination of local sentences. Then they evaluate each local sentence and combine the results. To evaluate a local sentence, say,

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi(x_i) \right),$$

they first compute the set  $\psi(\mathcal{A})$  of all  $a \in A$  such that  $\mathcal{A} \models \psi(a)$ . Since  $\psi$  is local and the class  $\mathcal{C}$  has bounded local tree-width this is possible quite efficiently. (In the special case of structures of bounded degree, this is easy to see, because  $\psi$  only has to be evaluated in substructures of  $\mathcal{A}$  of bounded size.) Finally, the algorithm tests whether there are  $a_1, \dots, a_k \in \psi(\mathcal{A})$  of pairwise distance greater than  $2r$ . This is possible in linear time by the following lemma:

**Lemma 8 (Frick and Grohe [11]).** *Let  $\mathcal{C}$  be a class of structures of bounded local tree-width. Then there is a function  $g$  and an algorithm that, given a structure  $\mathcal{A}$ , a subset  $P \subseteq A$ , and integers  $k, r$ , decides in time  $O(g(k, r)|A|)$  whether there are  $a_1, \dots, a_k \in P$  of pairwise distance greater than  $2r$ .*

The drawback of this algorithm is that we cannot even give an elementary upper bound for the function  $f$  in Theorem 7. The main reason for the enormous runtime of the algorithm in terms of the formula size is that to evaluate the local formulas, it translates them to tree-automata, and in the worst case the size of these automata

grows exponentially with each quantifier alternation. Therefore, it is a natural idea to bound the number of quantifier alternations in order to obtain smaller automata. But this would require that the translation of first-order sentences into local sentences preserves the quantifier structure. Unfortunately, the known proofs of Gaifman's theorem do not preserve the quantifier structure of the input formula.

These considerations motivated the present paper. Indeed, Theorem 2 shows that existential first-order sentences can be translated into Boolean combinations of existential local formulas. The price we pay for this is that these Boolean combinations of existential local formulas can get enormously large. Therefore, we use Theorem 6, because this theorem at least gives us an exponential upper bound on the size of the resulting formula. To evaluate an asymmetric local sentence, say

$$\exists x_1 \dots \exists x_k \left( \bigwedge_{1 \leq i < j \leq k} d(x_i, x_j) > 2r \wedge \bigwedge_{1 \leq i \leq k} \psi_i(x_i) \right),$$

where the  $\psi_i$  are conjunctive queries with negation, we first compute the sets  $\psi_1(\mathcal{A})$ ,  $\dots$ ,  $\psi_k(\mathcal{A})$ . This can be done as in the algorithm described above, but is actually faster since the  $\psi_i$  are conjunctive queries with negation. We use Lemma 10. Then we have to decide whether there are  $a_1 \in \psi_1(\mathcal{A}), \dots, a_k \in \psi_k(\mathcal{A})$  of pairwise distance greater than  $2r$ . Lemma 11 is an analogue of Lemma 8 for this more general situation.

If we now combine Lemma 10 and Lemma 11 together with Theorem 6 and plug them in the algorithms described in [11], we obtain the following theorem.

**Theorem 9.** *Let  $C$  be a class of structures whose local tree-width is bounded by a function  $\lambda : \mathbb{N} \rightarrow \mathbb{N}$  (i.e., for all  $\mathcal{A} \in C$  and  $r \geq 0$  we have  $\text{ltw}_{\mathcal{A}}(r) \leq \lambda(r)$ ). Then there are polynomials  $p, q$  such that for every  $\epsilon > 0$  there is an algorithm that, given a structure  $\mathcal{A}$  and an existential first-order sentence  $\varphi$ , decides if  $\mathcal{A} \models \varphi$  in time*

$$O\left(2^{2^{p(\lambda(q(\|\varphi\| + (1/\epsilon))) + \|\varphi\| + (1/\epsilon))}} \cdot |\mathcal{A}|^{1+\epsilon}\right),$$

i.e., in time doubly exponential in  $\|\varphi\|, (1/\epsilon), \lambda(q(\|\varphi\| + (1/\epsilon)))$  and near linear in  $|\mathcal{A}|$ .

For many interesting classes of structures of bounded local tree-width, such as planar graphs, the local tree-width is bounded by a linear function  $\lambda$ .

Let us complete the proof of Theorem 9 by showing the aforementioned lemmas.

**Lemma 10.** *There is a polynomial  $p$  and an algorithm that solves the following problem in time  $O(2^{p(\|\varphi\| + \text{tw}(\mathcal{A}))} \cdot |\mathcal{A}|)$ .*

*Input:* Structure  $\mathcal{A}$ , conjunctive query with negation  $\varphi$ .  
*Problem:* Decide if  $\mathcal{A} \models \varphi$ .

Lemma 10 can easily be proved using the standard dynamic programming techniques on graphs of bounded tree-width.

**Lemma 11.** *There is a polynomial  $p$  and an algorithm that solves the following problem in time  $O(2^{p(\text{ltw}_{\mathcal{A}}((k+1)r) + r + k)} \cdot |\mathcal{A}|)$ :*

*Input:* Structure  $\mathcal{A}$ , sets  $P_1, \dots, P_k \subseteq A$ , integer  $r \geq 1$ .  
*Problem:* Decide if there are  $a_1 \in P_1, \dots, a_k \in P_k$  of pairwise distance greater than  $r$ .

The proof of Lemma 11 requires some preparation. Let  $E$  be a binary relation symbol, and for each  $i \geq 1$ , let  $P_i$  be a unary relation symbol. A  $k$ -colored graph, for some  $k \geq 1$ , is an  $\{E, P_1, \dots, P_k\}$ -structure  $\mathcal{A}$  such that  $(A, E^{\mathcal{A}})$  is a graph and the sets  $P_1^{\mathcal{A}}, \dots, P_k^{\mathcal{A}}$  are disjoint. For a set  $B \subseteq A$  we let  $\text{col}^{\mathcal{A}}(B) = \{i \mid B \cap P_i^{\mathcal{A}} \neq \emptyset\}$ .

**Definition 12.** Let  $\mathcal{A}$  be a  $k$ -colored graph.

- (1) A set  $B \subseteq A$  of elements of a structure  $\mathcal{A}$  is  *$r$ -scattered* if  $|\text{col}^{\mathcal{A}}(B)| = |B|$  and for all  $b, c \in B$  with  $b \neq c$  we have  $d^{\mathcal{A}}(b, c) > r$ .
- (2) A set  $B \subseteq A$  is *maximum  $r$ -scattered* if  $B$  is  $r$ -scattered, and for all  $i \in \{1, \dots, k\} \setminus \text{col}^{\mathcal{A}}(B)$  we have  $P_i^{\mathcal{A}} \subseteq N_r^{\mathcal{A}}(B)$ .

For a structure  $\mathcal{A}$ , let

$$s_r(\mathcal{A}) := \max\{|\langle N_r^{\mathcal{A}}(a) \rangle| \mid a \in A\}.$$

Since we have  $||\mathcal{A}|| \in O(|A| \cdot \text{tw}(\mathcal{A}))$  [3], we have  $s_r(\mathcal{A}) \in O(|A| \cdot \text{ltw}_{\mathcal{A}}(r))$ .

**Lemma 13.** *There is an algorithm that solves the following problem in time  $O(|A| + k \cdot s_r(\mathcal{A}))$ .*

*Input:*  $k$ -colored graph  $\mathcal{A}$ , integer  $r \geq 1$ .  
*Problem:* Compute a maximum  $r$ -scattered subset of  $\mathcal{A}$ .

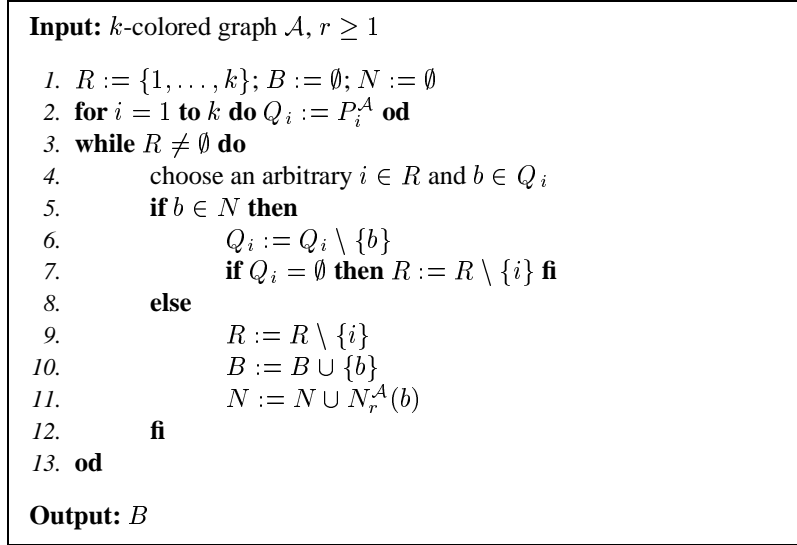
*Proof:* We solve the problem by a simple greedy algorithm, which is shown in Figure 1. It is easy to see that after each execution of the loop  $B$  is an  $r$ -scattered set,  $N = N_r^{\mathcal{A}}(B)$ , and  $R$  is the set of all  $i$  such that  $i \notin \text{col}^{\mathcal{A}}(B)$  and  $P_i^{\mathcal{A}} \not\subseteq N$ . Moreover, for all  $i \in R$  we have  $P_i^{\mathcal{A}} \setminus N \subseteq Q_i$ . The loop is executed until for all  $i \notin \text{col}^{\mathcal{A}}(B)$  we have  $Q_i = \emptyset$ , that is,  $P_i^{\mathcal{A}} \subseteq N = N_r^{\mathcal{A}}(B)$ . This shows that the algorithm is correct.

To estimate the running time, we observe that the loop is executed at most  $|A|$  times. Each line of the pseudocode except Line 11 can be executed in constant time. Line 11 is executed at most  $k$  times and requires time  $O(s_r(\mathcal{A}))$ . Together, this yields the desired time bound.  $\square$

Let SCATTERED SET be the following problem.

*Input:*  $k$ -colored graph  $\mathcal{A}$ , integer  $r \geq 1$ .  
*Problem:* Decide if there is an  $r$ -scattered set of size  $k$  in  $\mathcal{A}$ .

**Lemma 14.** *There is a polynomial  $p$  and an algorithm solving SCATTERED SET in time  $O(r^{p(\text{tw}(\mathcal{A})+k)} \cdot |A|)$ .*



**Fig. 1.**

*Proof:* In this proof, we assume that the reader is familiar with tree-decompositions of graphs and the typical dynamic programming algorithms on graphs of bounded tree-width.

Our algorithm proceeds as follows: It first computes a tree-decomposition of the input graph  $\mathcal{A}$  of width  $w := \text{tw}(\mathcal{A})$ . Bodlaender [3] proved that this is possible in time  $O(2^{p_1(w)}|A|)$ , where  $p_1$  is a suitable polynomial.

Let  $(\mathcal{T}, (B_t)_{t \in T})$  be this tree-decomposition; we can assume that  $\mathcal{T}$  is a binary rooted tree and that each block  $B_t$  contains exactly  $w + 1$  vertices, say,  $b_1^t, \dots, b_{w+1}^t$ . For every  $t \in T$ , we let

$$C_t := \bigcup_{\substack{u \in T \\ u=t \text{ or} \\ u \text{ descendant of } t}} B_t,$$

and  $\mathcal{C}_t := \langle C_t \rangle^{\mathcal{A}}$ .

Starting from the leaves of the tree, for every node  $t \in T$  our algorithm computes tables  $D_t, E_t$  that store the following information:

- $D_t$  stores  $d^{\mathcal{C}_t}(b_i, b_j)$  for  $1 \leq i < j \leq w + 1$ ,
- For each subset  $L \subseteq \{1, \dots, k\}$ , say of size  $l$ , and for all matrices  $(s_{ij})_{\substack{i \in L \\ 1 \leq j \leq w+1}}$  with  $0 \leq s_{ij} \leq r + 1$ , the table  $E_t$  stores whether there are vertices  $c_i \in P^{\mathcal{C}_t}$ , for  $i \in L$ , such that the set  $\{c_i \mid i \in L\}$  is  $r$ -scattered in  $\mathcal{C}_t$ , and for  $i \in L, 1 \leq j \leq w + 1$  we have

$$d^{\mathcal{C}_t}(c_i, b_j^t) \begin{cases} = s_{ij} & \text{if } 0 \leq s_{ij} \leq r, \\ > r & \text{if } s_{ij} = r + 1. \end{cases}$$

The size of these tables is in  $r^{O(k^2 w)}$ , and it is easy to see that for a node  $t \in T$  with children  $u, v$ , we can compute  $D_t$  and  $E_t$  from  $D_u, E_u, D_v, E_v$  in time  $r^{p_2(k+w)}$ , for a suitable polynomial  $p_2$ .  $\square$

**Lemma 15.** *There is a polynomial  $p$  and an algorithm solving SCATTERED SET in time  $O(2^{p(\text{tw}_{\mathcal{A}}((k+1)r+r+k)} \cdot |A|))$ .*

*Proof:* The algorithm is shown in Figure 2. To prove that it is correct, it suffices to verify

**Input:**  $k$ -colored graph  $\mathcal{A}$ , integer  $r \geq 1$

1.  $D := \{1, \dots, k\}; C := \emptyset; t := 1$
2. **for**  $1 \leq i \leq k$  **do**  $Q_i := P_i$  **od**
3. **while**  $D \neq \emptyset$  **do**
4. compute a maximum  $r$ -scattered set  $B^t$   
in the colored graph  $(A, E^{\mathcal{A}}, (Q_i)_{i \in D})$
5. **if**  $|B^t| = |D|$  **then accept** **fi**
6. **for**  $1 \leq i \leq k$  **do**  $Q_i := P_i \setminus N_{tr}(B^1 \cup \dots \cup B^t)$  **od**
7.  $C := C \cup \{i \mid Q_i = \emptyset\}, D := \{1, \dots, k\} \setminus C$
8. **if** there exists an  $r$ -scattered  $B$  with  $\text{col}^{\mathcal{A}}(B) = C$   
in  $\langle N_{(t+1)r}(B^1 \cup \dots \cup B^t) \rangle$  **then**  
 $t := t + 1$
9. **for**  $1 \leq i \leq k$  **do**  $Q_i := P_i \setminus N_{tr}(B^1 \cup \dots \cup B^t)$  **od**
10. **else**
11. **reject**
12. **fi**
13. **fi**
14. **od**
15. **accept**

**Fig. 2.**

the following loop conditions:

- (1) Before executing the while-loop the  $t$ -th time we have:
  - (a)  $C = \{i \mid P_i \subseteq N_{(t-1)r}(B^1 \cup \dots \cup B^{t-1})\}, D = \{1, \dots, k\} \setminus C$ , and  $|C| \geq t-1$ .
  - (b) If  $k_1, \dots, k_n$  is an enumeration of  $C$  then

$$\mathcal{A} \models \exists x_{k_1} \dots \exists x_{k_n} \in N_{(t-1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^{t-1}) \left( \bigwedge_{1 \leq i < j \leq n} d(x_{k_i}, x_{k_j}) > r \wedge \bigwedge_{1 \leq i \leq n} P_{k_i} x_{k_i} \right).$$

- (c) For all  $i$ , for all  $a \in Q_i$ , and for all  $b \in N_{(t-1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^{t-1})$  we have  $d^{\mathcal{A}}(a, b) > r$ .

- (2) If conditions (a) – (c) are met, then the algorithm accepts if, and only if, the desired  $r$ -scattered set exists.

Before entering the while loop for the first time we have  $C = \emptyset, D = \{1, \dots, k\}$  and conditions (b) and (c) are trivially satisfied. In Line 5 we start the algorithm of Lemma 13 with input  $\mathcal{A}, r$ . It returns some maximum  $r$ -scattered set  $B^1$ . If  $|B^1| = k = |D|$  we are done. Otherwise, there exists an  $i \notin B^1$  with  $P_i \subseteq N_r^{\mathcal{A}}(B^1)$ . After executing Lines 6 and 7 we have  $Q_i = \emptyset$  and  $i \in C$ . Without loss of generality, we may assume that  $C = \{1, \dots, s\}$ . For all  $b_1, b_2 \in N_r^{\mathcal{A}}(B^1)$  we have

$$d^{\mathcal{A}}(b_1, b_2) \leq r \iff d^{\langle N_{2r}^{\mathcal{A}}(B^1) \rangle}(b_1, b_2) \leq r.$$

Therefore,  $P_1, \dots, P_s \subseteq N_r^{\mathcal{A}}(B^1)$  implies that there exists an  $r$ -scattered set  $B$  with  $\text{col}^{\mathcal{A}}(B) = C$  in  $\langle N_{2r}^{\mathcal{A}}(B^1) \rangle$  if, and only if, there exists such a set  $B$  in  $\mathcal{A}$ . Thus if in Line 8 the algorithm finds out that there is no such set  $B$ , it correctly rejects in Line 12, because if there is not even an  $r$ -scattered set  $B$  with  $\text{col}^{\mathcal{A}}(B) = C$ , then there is certainly no  $r$ -scattered set  $B'$  with  $\text{col}^{\mathcal{A}}(B') = \{1, \dots, k\}$ . On the other hand, if there is such a set  $B$ , then there is still hope, and the algorithm continues. Lines 9 and 10 ensure that (c) is satisfied before returning to Line 3.

Now let  $t > 1$ , and suppose (a) – (c) are satisfied. Without loss of generality, we may assume that  $C := \{1, \dots, n\}$ . Then by (b), there exist  $b_1, \dots, b_n \in N_{(t-1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^{t-1})$  of pairwise distance greater than  $r$  such that for  $1 \leq i \leq n$  we have  $b_i \in P_i^{\mathcal{A}}$ .

Let  $B^t$  be the maximum  $r$ -scattered set of vertices in colors  $Q_i$  computed in Line 4. By (c), we have  $d^{\mathcal{A}}(a, b) > r$  for all  $a \in B^t$  and all  $b \in N_{(t-1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^{t-1})$ . If  $|B^t| = |D|$  the algorithm correctly accepts in Line 5, because  $B^t \cup \{b_1, \dots, b_n\}$  is an  $r$ -scattered set of size  $k$ .

If  $|B^t| < |D|$  there exists an  $i \in D$  with  $Q_i \subseteq N_r^{\mathcal{A}}(B^t)$ . After executing Lines 6 and 7, for the new set  $C$  and  $D$  we have  $|C| \geq t + 1, D = \{1, \dots, k\} \setminus C$ , and  $P_i \subseteq N_{tr}^{\mathcal{A}}(B^1 \cup \dots \cup B^t)$  for all  $i \in C$ . Hence condition (a) holds for  $t + 1$ . We argue similarly to the case  $t = 1$ : For all  $b_1, b_2 \in N_{tr}^{\mathcal{A}}(B^1 \cup \dots \cup B^t)$  we have

$$d^{\mathcal{A}}(b_1, b_2) \leq r \iff d^{\langle N_{(t+1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^t) \rangle}(b_1, b_2) \leq r.$$

Therefore  $P_i \subseteq N_{tr}^{\mathcal{A}}(B^1 \cup \dots \cup B^t)$  for all  $i \in C$  implies that there exists an  $r$ -scattered set  $B$  with  $\text{col}^{\mathcal{A}}(B) = C$  in  $\langle N_{(t+1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^t) \rangle$  if, and only if, there exists such a set  $B$  in  $\mathcal{A}$ . Thus if in Line 8 the algorithm finds out that there is no such set  $B$ , it correctly rejects in Line 12. If there is such a set  $B$ , then the algorithm continues. Again, Lines 9 and 10 ensure that (c) is satisfied before returning to Line 3.

Condition (a) implies that the while loop is iterated at most  $k + 1$  times. If the algorithm terminates within the loop, we have already seen that the answer is correct. If the algorithm reaches Line 15, then  $C = \{1, \dots, k\}$  is computed in Line 7 of the last, say  $t$ -th, iteration. Since the algorithm did not reject in Line 12, there must be an  $r$ -scattered set  $B$  with  $\text{col}^{\mathcal{A}}(B) = C = \{1, \dots, k\}$  in  $\langle N_{(t+1)r}^{\mathcal{A}}(B^1 \cup \dots \cup B^t) \rangle$  and thus in  $\mathcal{A}$ . Thus the algorithm accepts correctly.

This shows that the algorithm is correct.

To analyze the running time of the algorithm we notice that the auxiliary sets in Lines 1 and 2 can be computed in time  $O(k \cdot |\mathcal{A}|)$ . The while loop in lines 3–14 is



iterated at most  $k + 1$  times. By Lemma 13 the maximum  $r$ -scattered set in Line 4 can be computed in time  $O(k \cdot |A|)$ . Set manipulation in Lines 6, 7 and 11 again takes time  $O(k \cdot |A|)$ . For the test in Line 8, we use the algorithm of Lemma 14, we obtain a runtime of  $O(r^{p(\text{tw}_{\mathcal{A}}((k+1)r)+k)} \cdot |A|)$ . Putting everything together, we obtain the desired runtime.  $\square$

*Proof (of Lemma 11):* Essentially we can just apply Lemma 15 to the Gaifman graph of the input structure. The only problem is that in Lemma 11, the sets  $P_1, \dots, P_k$  are not necessarily disjoint. To solve this problem, we replace each vertex  $a$  of the Gaifman graph that is contained in  $l$  of the sets  $P_1, \dots, P_k$  by an  $l$ -clique; outside the clique, each vertex of this clique has the same neighbors as  $a$  had in the original graph, and each of these vertices belongs to exactly one  $P_i$  that  $a$  belonged to. We obtain a colored graph  $\mathcal{G}$ . It is easy to see that, for any  $r \geq 1$ , there is an  $r$ -scattered set of size  $k$  in  $\mathcal{G}$  if, and only if, there exist vertices  $a_1 \in P_1, \dots, a_k \in P_k$  of pairwise distance greater than  $r$  in  $\mathcal{A}$ .  $\square$

## References

1. A.V. Aho, J.E. Hopcroft, and J.D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974.
2. H.L. Bodlaender. NC-algorithms for graphs with small treewidth. In J. van Leeuwen, editor, *Proceedings of the 14th International Workshop on Graph Theoretic Concepts in Computer Science (WG'88)*, volume 344 of *Lecture Notes in Computer Science*, pages 1–10. Springer-Verlag, 1988.
3. H.L. Bodlaender. A linear-time algorithm for finding tree-decompositions of small treewidth. *SIAM Journal on Computing*, 25:1305–1317, 1996.
4. R. Diestel. *Graph Theory*. Springer-Verlag, second edition, 2000.
5. G. Dong, L. Libkin, and L. Wong. Local properties of query languages. *Theoretical Computer Science*, 239:277–308, 2000.
6. R.G. Downey and M.R. Fellows. *Parameterized Complexity*. Springer-Verlag, 1999.
7. R.G. Downey, M.R. Fellows, and U. Taylor. The parameterized complexity of relational database queries and an improved characterization of  $W[1]$ . In D.S. Bridges, C. Calude, P. Gibbons, S. Reeves, and I.H. Witten, editors, *Combinatorics, Complexity, and Logic – Proceedings of DMTCS '96*, pages 194–213. Springer-Verlag, 1996.
8. H.-D. Ebbinghaus and J. Flum. *Finite Model Theory*. Springer-Verlag, 2nd edition, 1999.
9. D. Eppstein. Diameter and treewidth in minor-closed graph families. *Algorithmica*, 27:275–291, 2000.
10. J. Flum, M. Frick, and M. Grohe. Query evaluation via tree-decompositions. *Journal of the ACM*, 49:716–752, 2002.
11. M. Frick and M. Grohe. Deciding first-order properties of locally tree-decomposable structures. *Journal of the ACM*, 48:1184 – 1206, 2001.
12. H. Gaifman. On local and non-local properties. In *Proceedings of the Herbrand Symposium, Logic Colloquium '81*. North Holland, 1982.
13. L. Hella, L. Libkin, J. Nurmonen, and L. Wong. Logics with aggregate operators. *Journal of the ACM*, 48:880–907, 2001.
14. L. Hella, L. Libkin, and Y. Nurmonen. Notions of locality and their logical characterizations over finite models. *Journal of Symbolic Logic*, 64:1751–1773, 1999.

15. L. Libkin. Logics with counting and local properties. *ACM Transactions on Computational Logic*, 1:33–59, 2000.
16. N. Robertson and P.D. Seymour. Graph minors III. Planar tree-width. *Journal of Combinatorial Theory, Series B*, 36:49–64, 1984.
17. L.J. Stockmeyer. *The Complexity of Decision Problems in Automata Theory*. PhD thesis, Department of Electrical Engineering, MIT, 1974.
18. P. van Emde Boas. Machine models and simulations. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, volume 1, pages 1–66. Elsevier Science Publishers, 1990.
19. M.Y. Vardi. The complexity of relational query languages. In *Proceedings of the 14th ACM Symposium on Theory of Computing*, pages 137–146, 1982.