<div align="center">

## Chapter 3:
# Model-checking on sparse graphs

**Compilation date: November 10, 2019**

</div>

In this chapter we present basic results on connections between the theory of sparse graph classes and algorithmic finite model theory. This connection can be summarized in one sentence as follows: for monotone graph classes (i.e. closed under taking subgraphs), the limit of algorithmic tractability of problems definable in first-order logic is exactly marked by the notion of nowhere denseness. To understand this claim better, we need to first introduce the language.

## 1    Elements of model theory

**Structures.**    We will work with logical structures, which can be regarded as a generalization of graphs. A structure is always over some *signature* (also called *vocabulary*), which is just a set of symbols that are used to describe objects in a structure. We will mostly work with *relational structures*, where the signature consists of relation symbols. Each relation symbol has a prescribed *arity*, which is a nonnegative integer indicating how many elements are bound in a single tuple in the relation. For instance, in graphs adjacency is a binary (i.e. arity 2) relation on vertices. Note that we allow *nullary* relations, that is, relations of arity 0; these are essentially boolean flags added to the structure. Throughout this chapter we will also allow unary functions in our signatures; we note that this is not a usual setting, but it will be convenient for us.

For a signature $\Sigma$, a *structure* $\mathbb{A}$ over $\Sigma$ consists of a universe $U$ and the *interpretation* of each symbol from $\Sigma$. For a relation symbol $R$, say of arity $k$, its interpretation $R^{\mathbb{A}} \subseteq U^k$ can be any set of $k$-tuples of elements of $U$. Note that the order of entries in tuples matter: if $R$ is a binary relation, then it may be that $(x,y) \in R^{\mathbb{A}}$, but $(y,x) \notin R^{\mathbb{A}}$. For a function symbol $f$, its interpretation $f^{\mathbb{A}}$ can be any function from $U$ to $U$.

For an example, directed graphs can be modelled as relational structures over a signature consisting of one binary relation symbol indicating the existence of an edge. Here, the universe is the vertex set. Undirected graphs can be modelled in the same way by assuming that the interpretation of the symbol is always symmetric: $(x,y)$ is an edge if and only if so is $(y,x)$.

For a structure $\mathbb{A}$, the *Gaifman graph* of $\mathbb{A}$, denoted $G(\mathbb{A})$, is the undirected graph with the vertex set being the universe of $\mathbb{A}$, where two elements are considered adjacent if and only if they appear together in some tuple in some relation in $\mathbb{A}$. In case $\mathbb{A}$ contains some functions, we also add edges between every element $x$ and its image $f^{\mathbb{A}}(x)$.

When we speak about a *class of structures*, we always mean that those structures have the same signature. A class of structures has *bounded expansion* if the Gaifman graphs of structures from the class form a class of bounded expansion. Same for nowhere denseness etc.

The *size* of a structure $\mathbb{A}$, denoted $\|\mathbb{A}\|$, is the cardinality of its universe plus the sum of cardinalities of all its relations and functions, each multiplied by its arity.

<div align="center">

1

</div>

**First-order logic.** For a fixed signature $\Sigma$, say consisting only of relation symbols for now, we may consider the first order logic $\mathsf{FO}[\Sigma]$. Formulas of this logic use variables which will be later intrepreted as elements of the universe, and are built from *atomic formulas* using a few constructions. The atomic formulas are:

- Equality check: $x = y$, where $x, y$ are variables.

- Relation check: $R(x_1, \ldots, x_k)$, where $x_1, \ldots, x_k$ are variables and $k$ is the arity of relation $R$.

Formulas can be combined into larger ones using the following constructions:

- Boolean connectives: if $\varphi$ and $\psi$ are formulas, then we can consider formulas $\varphi \vee \psi$, $\varphi \wedge \psi$, and $\neg\varphi$.

- Existential and universal quantification: if $\varphi$ is a formula, then we can write a formula $\exists_x \varphi$ and $\forall_x \psi$.

Note here that a formula may use variables that are not bound by any quantifier. These are called *free variables* and can be thought as parameters fed to the formula: the formula checks whether some condition between these parameters holds. If $\bar{x}$ is the set of free variables of a formula $\varphi$, we often indicate it by writing $\varphi(\bar{x})$. A *sentence* is a formula without free variables.

A formula is *quantifier-free* if it does not involve any quantifiers. If the signature also contains function symbols, we can also apply them to variables in atomic formulas. Consequently, even in quantifier-free formulas we may use atoms like $f(g(x)) = g(f(x))$ or $R(f(x), g(y))$.

For a set of variables $\bar{x}$ and a universe $U$, we write $U^{\bar{x}}$ for the set of all *valuations* of variables from $\bar{x}$ in $U$, which are just functions mapping variables of $\bar{x}$ to elements of $U$. We usually think of elements of $U^{\bar{x}}$ as of tuples of elements of $U$, with an implicit assignment between the entries of the tuple and the variables of $\bar{x}$.

Now for a formula $\varphi(\bar{x}) \in \mathsf{FO}[\Sigma]$, a $\Sigma$-structure $\mathbb{A}$ with universe $U$, and $\bar{a} \in U^{\bar{x}}$ we may define whether $\varphi(\bar{a})$ is true in $\mathbb{A}$ in the obvious way. We denote this fact by

$$\mathbb{A}, \bar{a} \models \varphi \qquad \text{or} \qquad \mathbb{A} \models \varphi(\bar{a}),$$

which should be read as "$\mathbb{A}$ with valuation $\bar{a}$ is a model for $\varphi$". For a structure $\mathbb{A}$ with universe $U$ and formula $\varphi(\bar{x})$, by $\varphi(\mathbb{A})$ we denote the set of all tuples $\bar{a} \in U^{\bar{x}}$ for which $\mathbb{A} \models \varphi(\bar{a})$. Thinking of structures as of databases, we sometimes say that $\varphi$ *selects* $\varphi(\mathbb{A})$ *in* $\mathbb{A}$.

The *model-checking* problem is the following: given a sentence $\varphi \in \mathsf{FO}[\Sigma]$ and a $\Sigma$-structure $\mathbb{A}$, we would like to decide whether $\mathbb{A} \models \varphi$. Whenever we have a class of structures $\mathcal{C}$, we can restrict this problem to the input structures belonging to $\mathcal{C}$.

## 2 Statement of the results

It turns out on classes of sparse structures, the model-checking problem can be solved efficiently.

**Theorem 2.1.** *For every class of structures $\mathcal{C}$ with bounded expansion, the model-checking problem on $\mathcal{C}$ can be solved in time $f(\varphi) \cdot \|A\|$, where $f$ is some computable function.*

**Theorem 2.2.** *For every nowhere dense class of structures $\mathcal{C}$ and any $\varepsilon > 0$, the model-checking problem on $\mathcal{C}$ can be solved in time $f(\varphi) \cdot \|A\|^{1+\varepsilon}$, where $f$ is some computable function.*

In fact, up to certain complexity assumptions, Theorem 2.2 constitutes the limit of algorithmic tractability of model-checking on monotone graph classes: an algorithm with a similar running time should not be expected on any monotone class that is somewhere dense. Since a precise statement of this fact would require some prior knowledge of parameterized complexity, we omit its further discussion.

In the next section we will give a full proof of Theorem 2.1. The main idea is to use low-treedepth colorings, as we did for SUBGRAPH ISOMORPHISM. The proof of Theorem 2.2 is more difficult and we will not give it.

# 3  Model-checking on classes of bounded expansion

The basic idea for the proof is to give a *quantifier elimination* procedure: an algorithm that given a formula $\varphi(\bar{x})$ and structure $\mathbb{A}$, iteratively simplifies the formula while enriching the structure with more information until the formula becomes quantifier-free and its satisfaction on a given tuple can be checked trivially. We will start with designing such a procedure for very simple structures: (colored) forests of bounded depth. Then we will lift the result to classes of bounded treedepth, and then to arbitrary classes of bounded expansion using low treedepth colorings.

## 3.1  Forests of bounded depth

A class of structures $\mathcal{C}$ is a class of *forests of bounded depth* if the signature of $\mathcal{C}$ consists only of relation symbols of arity 0 and 1, plus one function symbol which we shall denote $\mathsf{parent}(\cdot)$. Moreover, we assume that in all structures in $\mathcal{C}$, the interpretation of $\mathsf{parent}(\cdot)$ forms the parent function of some rooted forest of depth at most $d$ on the elements of the universe, where $d \in \mathbb{N}$ is a fixed constant. Here, roots of the forest can be recognized as the vertices that point to themselves in the $\mathsf{parent}(\cdot)$ function.

The following lemma provides quantifier elimination on classes of forests of bounded depth.

**Lemma 3.1.** *Let $\mathcal{C}$ be a class of forests of bounded depth and let $\varphi(\bar{x})$ be an* FO *formula over the signature of $\mathcal{C}$. There exists a class of forests $\widehat{\mathcal{C}}$ of bounded depth (same as of $\mathcal{C}$), a quantifier-free formula $\widehat{\varphi}(\bar{x})$ over the signature of $\widehat{\mathcal{C}}$, and a linear-time algorithm that, given a forest $\mathbb{T} \in \mathcal{C}$, outputs a forest $\widehat{\mathbb{T}} \in \widehat{\mathcal{C}}$ with the same universe and parent function as $\mathbb{T}$, such that*

$$\varphi(\mathbb{T}) = \widehat{\varphi}(\widehat{\mathbb{T}}).$$

*Proof.* We proceed by induction on the structure of the input formula by proving the statement for all subformulas appearing in $\varphi(\bar{x})$ in a bottom-up manner. For atomic formulas there is nothing to prove. Also the induction step for negation and disjunction is very easy. For negation, suppose we are considering a subformula of the form $\neg\psi$. Then having applied the induction assumption to $\psi$, it suffices to just negate the obtained quantifier-free formula $\widehat{\psi}$. Similarly, for a disjunction $\psi_1 \vee \psi_2$ we apply the induction assumption to $\psi_1$ and $\psi_2$, thus obtaining formulas $\widehat{\psi}_1$ and $\widehat{\psi}_2$ and, for a given forest $\mathbb{T} \in \mathcal{C}$, forests $\widehat{\mathbb{T}}_1$ and $\widehat{\mathbb{T}}_2$ with same universe and parent function as $\mathbb{T}$. Then we may take the formula $\widehat{\psi}_1 \vee \widehat{\psi}_2$ and forest $\widehat{\mathbb{T}}$ obtained by superposing $\widehat{\mathbb{T}}_1$ and $\widehat{\mathbb{T}}_1$: the signature of $\widehat{\mathbb{T}}$ is the union of signatures of $\widehat{\mathbb{T}}_1$ and $\widehat{\mathbb{T}}_1$, and in $\widehat{\mathbb{T}}$ we just inherit all relations from both $\widehat{\mathbb{T}}_1$ and $\widehat{\mathbb{T}}_2$.

The remaining induction step is the one for quantifiers. As a universal quantifier can be replaced by an existential quantifier combined with two negations, we need to prove the statement

for subformulas of the form $\exists_z \psi(\bar{y}, z)$, assuming the statement is already proved for $\psi(\bar{y}, z)$. By applying the induction assumption to $\psi(\bar{y}, z)$, we can assume that $\psi(\bar{y}, z)$ is quantifier-free.

Observe that since $\psi(\bar{y}, z)$ is quantifier-free, it is a boolean combination of atomic checks of the following form:

- $R$ holds, for some nullary relation $R$;

- $R(\mathsf{parent}^i(s))$ holds, for some unary relation $R$, $0 \leq i \leq d$, and variable $s \in \bar{y}, z$;

- $\mathsf{parent}^i(s) = \mathsf{parent}^j(t)$, for some $0 \leq i, j \leq d$ and variables $s, t \in \bar{y}, z$.

Note that the number of different atomic checks as above is bounded by $\mathcal{O}(d^2 |\bar{y}|^2)$. Let a *basic formula* be a conjunction of atomic checks and their negations, where each possible atomic check as above is taken either in positive or in negative (i.e. negated). Then $\psi(\bar{y}, z)$ can be equivalently rewritten as a disjunction over some set of basic formulas; these are exactly basic formulas which imply $\psi(\bar{y}, z)$. Since disjunction commutes with existential quantification, and we have already considered the induction step for a disjunction of two subformulas, we may apply the reasoning to each such basic formula and then combine the results. Thus, from now on we assume that $\psi(\bar{y}, z)$ is a basic formula.

Consider any forest $\mathbb{T} \in \mathcal{C}$ with universe $U$, and take any $\bar{a}, b \in U^{\bar{y}, z}$. Observe that since $\psi(\bar{y}, z)$ is basic, it is either non-satisfiable, or it exactly checks the isomorphism type of the subforest of $\mathbb{T}$ induced by the nodes of $\bar{a}, b$ and their ancestors in $\mathbb{T}$; such subforest shall be denoted by $\mathbb{T}[\bar{a}, b]$. As non-satisfiable formulas can be just removed from consideration, we may assume the following.

**Claim 1.** *There is a forest $\mathbb{S}$ of depth at most $d$ over the signature of $\mathcal{C}$, and a mapping $\tau$ from $\bar{y}, z$ to the nodes of $\mathbb{S}$ such that every leaf of $\mathbb{S}$ is in the image of $\tau$, such that for every $\mathbb{T} \in \mathcal{C}$ and $\bar{a}, b \in U^{\bar{y}, z}$, we have $\mathbb{T} \models \psi(\bar{a}, b)$ if and only if $\mathbb{T}[\bar{a}, b]$ is isomorphic to $\mathbb{S}$, where the isomorphism sends nodes mapped to variables $\bar{y}, z$ by the valuation $\bar{a}, b$ to respective nodes assigned by $\tau$ to variables $\bar{y}, z$.*

Now let us analyze the forest $\mathbb{S}$. By somewhat abusing the notation, for a variable $s \in \bar{y}, z$, we identify $s$ with the node $\tau(s)$ of $\mathbb{S}$.

Consider first the case that in $\mathbb{S}$, $z$ is an ancestor of some variable $s \in \bar{y}$. Then we can finish the proof very easily: $\psi(\bar{y}, z)$ forces the equality $z = \mathsf{parent}^i(s)$, for some $i \leq d$. Hence, formula $\exists_z \psi(\bar{y}, z)$ is equivalent to $\widehat{\psi}(\bar{y})$ obtained from $\psi(\bar{y}, z)$ by replacing every occurrence of $z$ with $\mathsf{parent}^i(s)$. Obviously, we then have $\widehat{\mathcal{C}} = \mathcal{C}$ and $\widehat{\mathbb{T}} = \mathbb{T}$.

Hence, we proceed with the other case: $z$ is not an ancestor of any of the variables of $s \in \bar{y}$. From now on we shall suppose that in $\mathbb{S}$, the node $z$ has a common ancestor with some variable $s \in \bar{y}$; in the end we will argue that the remaining case can be treated in a similar way. Let then $i \in \{1, \ldots, d\}$ be minimal such that $\mathsf{parent}^i(z) = \mathsf{parent}^j(s)$ for some $j \in \{0, 1, \ldots, d\}$ and $s \in \bar{y}$. Note that then $z$ is a leaf in $\mathbb{S}$ and $z = \mathsf{parent}^0(z), \mathsf{parent}^1(z), \ldots, \mathsf{parent}^{i-1}(z), \mathsf{parent}^i(z)$ form the path from $z$ to the lowest ancestor of $z$ that has at least two children in $\mathbb{S}$ (or belongs to $\bar{y}$ itself). Letting $k \leq d$ be the depth of $z$ in $\mathbb{S}$, note that $\mathsf{parent}^i(z)$ is at depth $\ell := k - i$.

For a forest $\mathbb{T} \in \mathcal{C}$, we shall call a node $q$ at depth $\ell + 1$ *good* if there is a descendant $r$ of $q$ at depth $k$ such that the path from $r$ to $q$ in $\mathbb{T}$ is isomorphic to the path from $z$ to $\mathsf{parent}^{i-1}(z)$ in $\mathbb{S}$ (this isomorphism should preserve satisfaction of unary relations on the paths). Further, for a node $p$ at depth $\ell$, let $\zeta(p) \in \{0, 1, \ldots, |\bar{y}|, \infty\}$ be the number of good children of $p$, or $\infty$ if their number is larger than $|\bar{y}|$.

We are ready to define the forest $\widehat{\mathbb{T}}$. It is obtained from $\mathbb{T}$ by adding $|\bar{y}| + 3$ unary relations:

4

- There is one new unary relation that selects all good nodes at depth $\ell + 1$.

- For each $m \in \{0, 1, \ldots, |\bar{y}|, \infty\}$ there is a new unary relation that selects all nodes $p$ at depth $\ell$ for which $\zeta(p) = m$.

Note that all nullary and unary relations present in $\mathbb{T}$ are preserved intact in $\widehat{\mathbb{T}}$. It is easy to compute $\widehat{\mathbb{T}}$ from $\mathbb{T}$ in linear time using a bottom-up dynamic programming to distinguish all good nodes. We take $\widehat{\mathcal{C}} = \{\widehat{\mathbb{T}} \colon \mathbb{T} \in \mathcal{C}\}$, so the signature of $\widehat{\mathcal{C}}$ is obtained from the signature of $\mathcal{C}$ by adding all the new unary relations described above.

We now construct the formula $\widehat{\psi}(\bar{y})$ that will be equivalent to $\exists_z \psi(\bar{y}, z)$, in the sense that $\widehat{\psi}(\bar{y})$ selects same tuples in $\widehat{\mathbb{T}}$ as $\exists_z \psi(\bar{y}, z)$ selects in $\mathbb{T}$. For $\mathbb{T} \in \mathcal{C}$ with universe $U$ and $\bar{a} \in U^{\bar{y}}$, in $\widehat{\psi}(\bar{a})$ we perform the following checks:

- First, check whether the satisfaction of nullary relations in $\mathbb{T}$ is as presecribed by $\psi$, and whether $\mathbb{T}[\bar{a}]$ is appropriately isomorphic to $\mathbb{S}[\bar{a}]$.

- Second, check whether in $\mathbb{T}$ the number of good children of $\mathsf{parent}^i(\bar{a}(s))$ is larger than the number of good children of $\mathsf{parent}^i(\bar{a}(s))$ that are contained in $\mathbb{T}[\bar{a}]$.

It is easy to see that the above checks can be performed using a quantifier-free formula $\widehat{\psi}(\bar{y})$, applied to $\bar{a}$ in $\widehat{\mathbb{T}}$. For the second check, observe that the number of good children of $\mathsf{parent}^i(\bar{a}(s))$ in $\mathbb{T}$ is directly encoded using new unary relations at $\mathsf{parent}^i(\bar{a}(s))$, while all the children of $\mathsf{parent}^i(\bar{a}(s))$ that are contained in $\mathbb{T}[\bar{a}]$ can be accessed by applying the parent function appropriately many times to the nodes of $\bar{a}$ — so we can count how many of them are good. Also, note that there are at most $|\bar{y}|$ such children, so it is fine not to distinguish between the numbers above $|\bar{y}|$ when counting good children of $\mathsf{parent}^i(\bar{a}(s))$.

To see that $\widehat{\psi}(\widehat{\mathbb{T}}) = \exists_z \psi(\mathbb{T})$ observe that the second check presented above is equivalent to checking whether there is a possibility to valuate variable $z$ to some node $b$ so that $\mathbb{T}[\bar{a}, b]$ is appropriately isomorphic to $\mathbb{S}$.

We are left with arguing what happens when in $\mathbb{S}$, node $z$ has no common ancestor with any other variable from $\bar{y}$. Then $z$ and its ancestors form a separate connected component in the forest $\mathbb{S}$. We perform an analogous argumentation, with the following modification. In a given forest $\mathbb{T}$, we mark every *root* $q$ as *good* if there is a descendant $r$ of $q$ such that the ancestors of $r$ in $\mathbb{T}$ form a path that is isomorphic to the path of the ancestors of $z$ in $\mathbb{S}$. Now instead of memorizing for every node the number of good children (up to threshold $|\bar{y}|$), which would make little sense now, we use $|\bar{y}| + 2$ new nullary relations to memorize whether the number of good roots is $0, 1, \ldots, |\bar{y}|$, or larger than $|\bar{y}|$. The rest of the reasoning follows in the same way. $\square$

## 3.2 Structures of boundeed treedepth

We now lift the quantifier elimination procedure to classes of bounded treedepth. This lift will be almost automatic, but we first need to encode structures of bounded treedepth in rooted forests.

**Lemma 3.2.** *Let $\mathcal{C}$ be a class of structures of bounded treedepth and let $\varphi(\bar{x})$ be an* FO *formula over the signature of $\mathcal{C}$. There exists a class of forests of bounded depth $\widehat{\mathcal{C}}$, a formula $\widehat{\varphi}(\bar{x})$ over the signature of $\widehat{\mathcal{C}}$, and a linear-time algorithm that given $\mathbb{A} \in \mathcal{C}$ outputs a forest $\mathbb{T} \in \widehat{\mathcal{C}}$ with the following properties: $\mathbb{T}$ has the same universe as $\mathbb{A}$, $E(G(\mathbb{T})) \subseteq E(G(\mathbb{A}))$, and*

$$\varphi(\mathbb{A}) = \widehat{\varphi}(\mathbb{T}).$$

*Proof.* We may assume that the signature of $\mathcal{C}$ consists only of relation symbols. This is because every function $f$ can be replaced with a new binary relation $R^f$ that binds each argument $a$ with $f(a)$, that is, $R^f(a, b)$ is true if and only if $b = f(a)$. Note that this requires also an easy syntactic modification of the formula $\varphi(\bar{x})$: whenever $f(x)$ is used for some variable $x$, instead we use a new variable $x'$ that is introduced using existential quantification and bound to $x$ by requiring that $R^f(x, x')$ holds.

Having this assumption, we present an algorithm that for a given structure $\mathbb{A} \in \mathcal{C}$ computes a suitable forest $\mathbb{T}$. The formula $\widehat{\varphi}(\bar{x})$ will be defined on the way, while the class $\widehat{\mathcal{C}}$ will comprise all forests $\mathbb{T}$ produced for all input $\mathbb{A} \in \mathcal{C}$.

For a given $\mathbb{A} \in \mathcal{C}$, we first compute the Gaifman graph $G := G(\mathbb{A})$. Then we compute any depth-first search forest $T$ of $G$; we treat $T$ as a rooted forest on the same vertex set as $G$, which is the universe of $\mathbb{A}$. Observe that if $d$ is the bound on the treedepth of the Gaifman graphs of structures in $\mathcal{C}$, then these Gaifman graphs will not contain paths longer than $d' := 2^d$, hence the depth of $T$ will be bounded by $d'$. We will now define a forest $\mathbb{T}$ by taking $T$, treating it as a logical structure with parent function, and adding some unary relations that will encode the relations from $\mathbb{A}$. As clearly $E(T) \subseteq E(G)$, we will have $E(G(\mathbb{T})) = E(T) \subseteq E(G) = E(G(\mathbb{A}))$.

Consider any relation $R$ present in $\mathbb{A}$, say of arity $k$. Observe that for every tuple $\bar{a} = (a_1, \ldots, a_k) \in R$, the elements of $\bar{a}$ form a clique in $G$, hence they are pairwise bound by the ancestor-descendant relation in $T$. Therefore, there exists an element of $\bar{a}$, say $a_i$, such that all $a_j$ for $j \in \{1, \ldots, k\}$ are ancestors of $a_j$. Let $\pi \colon \{1, \ldots, k\} \to \{0, \ldots, d\}$ be such that $a_j = \mathsf{parent}^{\pi(j)}(a_i)$, for $j \in \{1, \ldots, k\}$. Now, for every $\pi$ as above we introduce a new unary relation $R^\pi$ in $\mathbb{T}$ which selects all elements $e$ such that

$$(\mathsf{parent}^{\pi(1)}(e), \mathsf{parent}^{\pi(2)}(e), \ldots, \mathsf{parent}^{\pi(k)}(e)) \in R.$$

Now, every atomic check of the form $R(x_1, \ldots, x_k)$ in $\mathbb{A}$ can be replaced by the following equivalent formula in $\mathbb{T}$. We make a disjunction over all $i \in \{1, \ldots, k\}$ and all $\pi \colon \{1, \ldots, k\} \to \{0, \ldots, d\}$ satisfying $\pi(i) = 0$ of the following assertions:

$$R^\pi(x_i) \wedge \bigwedge_{j=1}^{k} \left( x_j = \mathsf{parent}^j(x_i) \right).$$

By performing this construction for every relation $R$ in $\mathbb{A}$ we obtain a forest $\mathbb{T}$ and a formula $\widehat{\varphi}(\bar{x})$ over the signature of this forest such that $\varphi(\mathbb{A}) = \widehat{\varphi}(\mathbb{T})$. $\qquad\square$

By composing the transformations provided by Lemmas 3.2 and 3.1 in order we obtain the following corollary, which is just a strengthening of Lemma 3.2 by requiring that the output formula is quantifier-free.

**Lemma 3.3.** *Let $\mathcal{C}$ be a class of structures of bounded treedepth and let $\varphi(\bar{x})$ be an* FO *formula over the signature of $\mathcal{C}$. There exists a class of forests of bounded depth $\widehat{\mathcal{C}}$, a quantifier-free formula $\widehat{\varphi}(\bar{x})$ over the signature of $\widehat{\mathcal{C}}$, and a linear-time algorithm that given $\mathbb{A} \in \mathcal{C}$ outputs a forest $\mathbb{T} \in \widehat{\mathcal{C}}$ with the following properties: $\mathbb{T}$ has the same universe as $\mathbb{A}$, $E(G(\mathbb{T})) \subseteq E(G(\mathbb{A}))$, and*

$$\varphi(\mathbb{A}) = \widehat{\varphi}(\mathbb{T}).$$

## 3.3 Classes of bounded expansion

We can now state and prove the result about quantifier elimination on classes of bounded expansion.

**Theorem 3.4.** *Let $\mathcal{C}$ be a class of structures of bounded expansion and let $\varphi(\bar{x})$ be an* FO *formula over the signature of $\mathcal{C}$. There exists a class of bounded expansion $\widehat{\mathcal{C}}$, a quantifier-free formula $\widehat{\varphi}(\bar{x})$ over the signature of $\widehat{\mathcal{C}}$, and a linear-time algorithm that given $\mathbb{A} \in \mathcal{C}$ outputs a structure $\widehat{\mathbb{A}} \in \widehat{\mathcal{C}}$ with the following properties: $\widehat{\mathbb{A}}$ has the same universe as $\mathbb{A}$, $E(G(\widehat{\mathbb{A}})) \subseteq E(G(\mathbb{A}))$, and*

$$\varphi(\mathbb{A}) = \widehat{\varphi}(\widehat{\mathbb{A}}).$$

*Proof.* We again proceed by induction on the structure of the formula $\varphi(\bar{x})$, proving the statement for all subformulas in a bottom-up manner. Similarly as in the proof of Lemma 3.1, there is nothing to prove for atomic subformulas, and the induction step for boolean connectives is easy. We are left with giving the induction step for subformulas of the form $\exists_z \psi(\bar{y}, z)$. Again, by applying the induction assumption to $\psi(\bar{y}, z)$, we may assume that this formula is quantifier-free.

Similarly as in the proof of Lemma 3.2, we may remove functions from the structure $\mathbb{A}$ by replacing each function $f$ with a new binary relation $R^f$ binding arguments with values. For this, we also repeatedly replace every usage of a term of the form $f(x)$ in $\psi(\bar{y}, z)$ with a new variable $x'$, quantified existentially, and add a constraint $R^f(x, x')$ to $\psi$. Note that this operation does not change the Gaifman graph of $\mathbb{A}$. Thus, from now on we may assume that the input structure $\mathbb{A}$ does not contain any functions, but this comes at the cost of introducing new existentially quantified variables. Hence, we are working with a formula of the form $\exists_{\bar{z}} \psi(\bar{y}, \bar{z})$, where now $\bar{z}$ is a tuple of variables, rather than one variable.

Let $k = |\bar{y}| + |\bar{z}|$ be the total number of variables present. Since $k$ is considered a constant and $\mathcal{C}$ has bounded expansion, in linear time we can compute a treedepth-$k$ coloring $\lambda$ of $G := G(\mathbb{A})$, that is, a coloring of the universe $U$ of $\mathbb{A}$ with $M$ colors, where $M$ is a constant depending only on $\mathcal{C}$ and $k$, such that every $k$ color classes induce in $G$ a subgraph of treedepth at most $k$. Let $\Xi$ be the set of all $k$-sized subsets of colors under $\lambda$; then $|\Xi| = \binom{M}{k}$. For every set $X$ of $k$ colors under $\lambda$, let $\mathbb{A}[X]$ be the substructure of $\mathbb{A}$ induced by the vertices with colors in $X$. That is, the universe of $\mathbb{A}[X]$ is $\lambda^{-1}(X)$ and in $\mathbb{A}[X]$ we preserve all tuples from all relations from $\mathbb{A}$ whose all entries are in $\lambda^{-1}(X)$. Letting $\mathcal{D}$ be the class of all structures over the signature of $\mathcal{C}$ whose Gaifman graphs have treedepth at most $k$, we see that $\mathbb{A}[X] \in \mathcal{D}$ for all $X \in \Xi$.

Now comes the key claim, which follows immediately from the fact that $\exists_{\bar{z}} \psi(\bar{y}, \bar{z})$ is an existential formula that involves only $k$ variables, while $\Xi$ contains all $k$-tuples of colors in $\lambda$.

**Claim 2.** *For any $\bar{a} \in U^{\bar{y}}$, we have that $\mathbb{A} \models \exists_{\bar{z}} \psi(\bar{a}, \bar{z})$ if and only if there exists $X \in \Xi$ such that $\bar{a} \subseteq \lambda^{-1}(X)$ and $\mathbb{A}[X] \models \exists_{\bar{z}} \psi(\bar{a}, \bar{z})$.*

We now apply Lemma 3.3 to the class $\mathcal{D}$, formula $\exists_{\bar{z}} \psi(\bar{y}, \bar{z})$, and structures $\mathbb{A}[X]$ for $X \in \Xi$. This provides us with a class of bounded-depth forests $\widehat{\mathcal{D}}$, formula $\widehat{\psi}(\bar{y})$, and structures $\widehat{\mathbb{A}}[X]$ for $X \in \Xi$ (computable in linear time) such that

$$\widehat{\psi}(\widehat{\mathbb{A}}[X]) = \exists_{\bar{z}} \psi(\mathbb{A}[X]) \qquad \text{for all } X \in \Xi.$$

Now we modify the forests $\widehat{\mathbb{A}}[X]$ for $X \in \Xi$ by renaming symbols so that their signatures become disjoint; in particular, each forest $\widehat{\mathbb{A}}[X]$ uses its own parent function $\mathsf{parent}^X(\cdot)$. This in particular

means that now $\widehat{\mathbb{A}}[X] \in \widehat{\mathcal{D}}^X$, where $\widehat{\mathcal{D}}^X$ is obtained from $\widehat{\mathcal{D}}$ by the same renaming of the signature, and we may similarly modify $\widehat{\psi}(\bar{y})$ to a formula $\widehat{\psi}^X(\bar{y})$ so that

$$\widehat{\psi}^X(\widehat{\mathbb{A}}[X]) = \exists_{\bar{z}}\psi(\mathbb{A}[X]) \qquad \text{for all } X \in \Xi. \tag{1}$$

We now define the final structure $\widehat{\mathbb{A}}$ as the superposition of structures $\widehat{\mathbb{A}}[X]$ for $X \in \Xi$. That is, the universe of $\widehat{\mathbb{A}}$ is $U$ and for every $X \in \Xi$, we add all relations from $\widehat{\mathbb{A}}[X]$ to $\widehat{\mathbb{A}}$ (note that these relations are restricted to elements of $\lambda^{-1}(X)$). Observe that since the signatures of structures $\widehat{\mathbb{A}}[X]$ are disjoint, the signature of $\widehat{\mathbb{A}}$ is the (disjoint) union of those signatures and relations originating in different structures $\widehat{\mathbb{A}}[X]$ do not mix with each other. In addition, we add $M$ unary predicates, each selecting vertices of a different color in $\lambda$.

Note that thus, we have

$$E(G(\widehat{\mathbb{A}})) = \bigcup_{X \in \Xi} E(G(\widehat{\mathbb{A}}[X])),$$

however by Lemma 3.3 we know that

$$E(G(\widehat{\mathbb{A}}[X])) \subseteq E(G(\mathbb{A}[X])) \subseteq E(G) \qquad \text{for all } X \in \Xi.$$

Hence we have

$$E(G(\widehat{\mathbb{A}})) \subseteq E(G(\mathbb{A})),$$

as required.

Finally, we define the formula $\widehat{\psi}(\bar{y})$ as follows:

$$\widehat{\psi}(\bar{y}) = \bigvee_{X \in \Xi} \left[ \bigwedge_{s \in \bar{y}} (\lambda(s) \in X) \wedge \widehat{\psi}^X(\bar{y}) \right].$$

Here, $\lambda(s) \in X$ is a shorthand for a disjunction of size $|X| = k$ that checks whether $s$ is of any of the colors in $X$. It follows directly from Claim 2 and (1) that

$$\widehat{\psi}(\widehat{\mathbb{A}}) = \exists_{\bar{z}}\psi(\mathbb{A}).$$

Since $\widehat{\psi}(\bar{y})$ is quantifier-free, we are done. $\qquad\qquad\square$

We observe that Theorem 2.1 trivially follows from Theorem 3.4. Indeed, by applying Theorem 3.4 to the input sentence $\varphi$ and structure $\mathbb{A}$, we obtain a quantifier-free sentence $\widehat{\varphi}$ and a structure $\widehat{\mathbb{A}}$, computable in linear time from $\mathbb{A}$, such that $\varphi$ is true in $\mathbb{A}$ if and only if $\widehat{\varphi}$ is true in $\widehat{\mathbb{A}}$. But $\widehat{\varphi}$ is a quantifier-free sentence, which means that it is just a boolean combination of nullary relation checks. Therefore, its satisfaction in $\widehat{\mathbb{A}}$ can be checked in constant time and we are done.