# APPROXIMATING MINIMUM REPRESENTATIONS OF KEY HORN FUNCTIONS[*]

KRISTÓF BÉRCZI[†], ENDRE BOROS[‡], ONDŘEJ ČEPEK[§], PETR KUČERA[¶], AND KAZUHISA MAKINO[‖]

**Abstract.** Horn functions form an important subclass of Boolean functions and appear in many different areas of computer science and mathematics as a general tool to describe implications and dependencies. Finding minimum sized representations for such functions with respect to most commonly used measures is a computationally hard problem admitting a $2^{\log^{1-o(1)} n}$ inapproximability bound.

In this paper we consider the natural class of key Horn functions representing keys of relational databases. For this class, the minimization problems for most measures remain NP-hard. In this paper we provide logarithmic factor approximation algorithms for key Horn functions with respect to all such measures.

**Key words.** Approximation algorithms, Directed hypergraphs, Horn minimization, Implicational systems

**AMS subject classifications.** 05C65, 05C85, 68W25, 90C27

**1. Introduction.** A Boolean function of $n$ variables is a mapping from $\{0,1\}^n$ to $\{0,1\}$. Boolean functions naturally appear in many areas of mathematics and computer science and constitute a principal concept in complexity theory. In this paper we shall study an important problem connected to Boolean functions, a so called Boolean minimization problem, which aims at finding a shortest possible representation of a given Boolean function. The formal statement of the Boolean minimization problem (BM) of course depends on (i) how the input function is represented, (ii) how it is represented on the output, and (iii) the way the output size is measured.

One of the most common representations of Boolean functions are conjunctive normal forms (CNFs), the conjunctions of clauses which are elementary disjunctions of literals. There are two usual ways how to measure the size of a CNF: the number of clauses and the total number of literals (sum of clause lengths). It is easy to see that BM is NP-hard if both input and output is a CNF (for both above mentioned measures of the output size). This is an easy consequence of the fact that BM contains the CNF satisfiability problem (SAT) as its special case (an unsatisfiable formula can

[†]MTA-ELTE Momentum Matroid Optimization Research Group and MTA-ELTE Egerváry Research Group, Department of Operations Research, Eötvös University, Budapest, Hungary (kristof.berczi@ttk.elte.hu).

[‡]MSIS Department and RUTCOR, Rutgers University, New Jersey, USA (endre.boros@rutgers.edu).

[§]Charles University, Faculty of Mathematics and Physics, Department of Theoretical Computer Science and Mathematical Logic, Praha, Czech Republic (cepek@ktiml.mff.cuni.cz).

[¶]Charles University, Faculty of Mathematics and Physics, Department of Theoretical Computer Science and Mathematical Logic, Praha, Czech Republic (kucerap@ktiml.mff.cuni.cz).

[‖]Research Institute for Mathematical Sciences (RIMS), Kyoto University, Kyoto, Japan (makino@kurims.kyoto.ac.jp).

1

be trivially recognized from its shortest CNF representation). In fact, BM was shown to be probably harder than SAT: while SAT is NP-complete (i.e. $\Sigma_1^p$-complete [11]), BM is $\Sigma_2^p$-complete [29] (see also the review paper [30] for related results). It was also shown that BM is $\Sigma_2^p$-complete when considering Boolean functions represented by general formulas of constant depth as both the input and output for BM [8]. A $O(n^{1-\varepsilon})$-inapproximability result was given in [28].

Horn functions form a subclass of Boolean functions which plays a fundamental role in constructive logic and computational logic. They are important in automated theorem proving and relational databases. An important feature of Horn functions is that SAT is solvable for this class in linear time [15]. A CNF is Horn if every clause in it contains at most one positive literal, and it is pure Horn (or definite Horn in some literature) if every clause in it contains exactly one positive literal. Such a positive literal is then called the *head* of the given clause and the set of all negative literals is called the *body* of the clause (we often identify the body of a clause with the set of variables with negative occurrences especially if we view the clause as an implication in which the body implies the head). A Boolean function is (pure) Horn, if it admits a (pure) Horn CNF representation. Pure Horn functions represent a very interesting concept which was studied in many areas of computer science and mathematics under several different names. The same concept appears as directed hypergraphs in graph theory and combinatorics, as implicational systems in artificial intelligence and database theory, and as lattices and closure systems in algebra and concept lattice analysis [9].

*Example* 1.1. Consider a pure Horn CNF $\Phi = (\overline{a} \vee \overline{b} \vee \overline{c} \vee d) \wedge (\overline{d} \vee e) \wedge (\overline{d} \vee f) \wedge (\overline{d} \vee g) \wedge (\overline{e} \vee \overline{f} \vee \overline{g} \vee a) \wedge (\overline{e} \vee \overline{f} \vee \overline{g} \vee b) \wedge (\overline{e} \vee \overline{f} \vee \overline{g} \vee c)$ on variables $a, b, c, d, e, f, g$, where $\overline{a}$ stands for the negation of $a$, etc. The CNF $\Phi$ can be viewed equivalently as a directed hypergraph $\mathcal{H} = (V, \mathcal{E})$ with vertex set $V = \{a, b, c, d, e, f, g\}$ and directed hyperarcs $\mathcal{E} = \{(\{a,b,c\}, d), (\{d\}, e), (\{d\}, f), (\{d\}, g), (\{e,f,g\}, a), (\{e,f,g\}, b), (\{e,f,g\}, c)\}$. This latter can be expressed more concisely using a generalization of adjacency lists for ordinary digraphs in which all hyperarcs with the same body (also called source) are grouped together $\{a,b,c\} : d; \{d\} : e, f, g; \{e,f,g\} : a, b, c$, or can be represented as an implicational (closure) system on variables $a, b, c, d, e, f, g$ defined by rules $abc \rightarrow d, d \rightarrow efg, efg \rightarrow abc$.

Interestingly, in each of these areas the problem similar to BM, i.e. a problem of finding the shortest equivalent representation of the input data (CNF, directed hypergraph, set of rules) was studied. For example, such a representation can be used to reduce the size of knowledge bases in expert systems, thus improving the performance of the system. The above examples show that a "natural" way how to measure the size of the representation depends on the area. Six different measures and corresponding concepts of minimality were considered in [2, 12]: (B) number of bodies, (BA) body area, (TA) total area, (C) number of clauses, (BC) number of bodies and clauses, and (L) number of literals. For precise definitions, see Section 2. With a slight abuse of notation we shall use (B), (BA), (TA), (C), (BC) and (L) to denote both the measures and the corresponding minimization problems.

The only one of these six minimization problems for which a polynomial time procedure exists to derive a minimum representation is (B). The first such algorithm appeared in the database theory literature [23]. Different algorithms for the same task were then independently discovered in hypergraph theory [2], and in the theory of closure systems [18].

For the remaining five measures it is NP-hard to find the shortest representation.

There is an extensive literature on the intractability results in various contexts for these minimization problems [2, 19, 23]. It was shown that (C) and (L) stay NP-hard even when the inputs are limited to cubic (bodies of size at most two) pure Horn CNFs [6], and the same result extends to the remaining three measures. Note that if all bodies are of size one then the above problems become equivalent with the transitive reduction of directed graphs, which is tractable [1]. It should be noted that there exists many other tractable subclasses, such as acyclic and quasi-acyclic pure Horn CNFs [20], and CQ Horn CNFs [5]. There are also a few heuristic minimization algorithms for pure Horn CNFs [4].

It was shown that (C) and (L) are not only hard to solve exactly but even hard to approximate. More precisely, [3] shows that these problems are inapproximable within a factor $2^{\log^{1-\varepsilon}(n)}$ assuming $NP \subsetneq DTIME(n^{polylog(n)})$, where $n$ denotes the number of variables. In addition, [7] shows that they are inapproximable within a factor $2^{\log^{1-o(1)} n}$ assuming $P \subsetneq NP$ even when the input is restricted to 3-CNFs with $O(n^{1+\varepsilon})$ clauses, for some small $\varepsilon > 0$. It is not difficult to see that the same proof extends to (BC) and (TA) as well. On the positive side, (C), (BC), (BA), and (TA) admit $(n-1)$-approximations and (L) has an $\binom{n}{2}$-approximation [19]. To the best of our knowledge, no better approximations are known even for pure Horn 3-CNFs.

Given a relational database, a key is a set of attributes with the property that a value assignment to this set uniquely determines the values of all other attributes [24, 27]. The concept of a key is essential for standard database operations. A relational database uniquely defines a pure Horn function $h$ over the set of attributes, representing the so-called functional dependencies of the database. An implicate $B \to v$ of $h$ represents the fact that the knowledge of the attribute values in set $B$ uniquely defines the value for attribute $v$. If $K$ is a key of the database, then $K \to v$ is an implicate of $h$ for all attributes $v$. Motivated by this, we say that a pure Horn CNF is *key Horn* if each of its bodies implies all other variables, that is, setting all variables in any of its bodies to one forces all other variables to one. A Boolean function is called *key Horn* if it has a key Horn CNF representation. Key Horn functions are natural concepts to represent the keys of relational databases. They generalize the well studied class of *hydra functions* considered in [25]. For this special class, in which all bodies are of size two, a 2-approximation algorithm for (C) was presented in [25] while the NP-hardness for (C) was proved in [22]. The latter result implies NP-hardness for hydra functions also for (BC), (TA), and (L). It is also easy to see that (B) and (BA) are trivial in this case.

In this paper we consider the minimization problems for key Horn functions. Any irredundant representation of a key Horn function has the same set of bodies, implying that problems (B) and (BA) are in P. We show that a simple algorithm gives a $\frac{2k}{k+1}$-approximation for (TA) and a $k$-approximation for (C), (BC), and (L), where $k$ is the size of a largest body. Our paper contains two main results. The first one gives a $\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2\}$-approximation bound for key Horn functions for (C) and (BC) which is significantly better than the $(n-1)$-approximation bound known for general Horn functions. The second result improves the $\binom{n}{2}$-approximation bound for (L) to $\frac{108}{17}\lceil \log k \rceil + 2$. Table 1 summarizes the state of the art of Horn minimization and the results presented in this paper for key Horn functions.

The structure of our paper is as follows: Section 2 introduces the necessary definitions and notation, Section 3 provides lower bounds for the measures we introduced, while Section 4 contains our results on approximation algorithms. For the (L) measure, our approach in Section 4 relies on approximating a solution to a subproblem

TABLE 1
*Complexity landscape of Horn and key Horn minimization, where the bold letters represent the results obtained in this paper. Here n and k respectively denote the number of variables and the size of a largest body. All problems except those labeled by P are NP-hard. Inapproximability bounds for Horn minimization hold even when the size of the bodies are bounded by k ($\geq 2$).*

| Measure | Horn | | Key Horn | |
|---------|------|------|----------|------|
| | Inapprox. | Approx. | Inapprox. | Approx. |
| (B) | P [23] | | P [23] | |
| (BA) | $1$ [2] | $n-1$ [19] | **P** | |
| (TA) | $2^{\log^{1-o(1)} n}$ [7] | $n-1$ [19] | $1$ [22] | $\frac{2k}{k+1}$ |
| (C) | $2^{\log^{1-o(1)} n}$ [7] | $n-1$ [19] | $1$ [22] | $\mathbf{\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2, k\}}$ |
| (BC) | $2^{\log^{1-o(1)} n}$ [7] | $n-1$ [19] | $1$ [22] | $\mathbf{\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2, k\}}$ |
| (L) | $2^{\log^{1-o(1)} n}$ [7] | $\binom{n}{2}$ [19] | $1$ [22] | $\mathbf{\min\{\frac{108}{17}\lceil \log k \rceil + 2, k\}}$ |

which is shown to be NP-hard in Section 5. Finally, Section 6 discusses the relation of our approach to the problem of finding a minimum weight strongly connected subgraph.

**2. Preliminaries.** Let $V$ denote a set of variables. Members of $V$ are called *positive literals* while their negations are called *negative literals*. Throughout the paper, the number of variables is denoted by $n = |V|$. A *Boolean function* is a mapping $f : \{0,1\}^V \to \{0,1\}$. The *characteristic vector* of a set $Z$ is denoted by $\chi_Z$, that is, $\chi_Z(v) = 1$ if $v \in Z$ and 0 otherwise. We say that a set $Z \subseteq V$ is a *true set* of $f$ if $f(\chi_Z) = 1$, and a *false set* otherwise.

For a subset $\emptyset \neq B \subseteq V$ and $v \in V \setminus B$ we write $B \to v$ to denote the pure Horn clause $C = v \vee \bigvee_{u \in B} \overline{u}$. We can consider such a clause to be an implication as if all variables in $B$ are set to true in a true assignment then $v$ must be true as well. Here $B$ and $v$ are called the *body* and *head* of the clause, respectively. That is, a pure Horn CNF can be associated with a directed hypergraph where every clause $B \to v$ is considered to be a directed hyperarc oriented from $B$ to $v$. The *set of bodies* appearing in a pure Horn CNF representation $\Phi$ is denoted by $\mathcal{B}_\Phi$. We will also use the notation $B \to H$ to denote $\bigwedge_{v \in H} B \to v$. By grouping the clauses with the same body, a pure Horn CNF $\Phi = \bigwedge_{B \in \mathcal{B}_\Phi} \bigwedge_{v \in H(B)} B \to v$ can be represented as $\bigwedge_{B \in \mathcal{B}_\Phi} B \to H(B)$. The latter representation is in a one-to-one correspondence with the adjacency list representation of the corresponding directed hypergraph.

For any pure Horn function $h$ the family of its true sets is closed under taking intersection (see Lemma 4.5 in [13]) and clearly contains $V$. This implies that for any non-empty set $Z \subseteq V$ there exists a unique minimal true set containing $Z$. This set is called the *closure* of $Z$ and we denote it by $F_h(Z)$. If $\Phi$ is a pure Horn CNF representation of $h$, then $F_h(Z)$ can be computed in linear time in the size of $\Phi$ [15]. Note that the resulting closure $F_h(Z)$ depends only on the set $Z$ and the Horn function $h$, and not on the particular CNF $\Phi$ we use to represent $h$. It is important to note here that $h : \{0,1\}^V \to \{0,1\}$ is a function that exists independently of its representations.

It can be represented, in particular, by CNFs, and typically by many different ones. In our algorithmic approach to generate a better (shorter) CNF representation of a Horn function, that is represented by a given CNF on the input, we shall rely on certain invariants that in fact depend only on the function and not on its particular representation.

One such invariant is the closure of a subset, defined above. The algorithm, computing the closure $F_h(Z)$ of a subset $Z$ using a given CNF representation $\Phi$ of $h$, is also called the *forward chaining procedure* (see e.g., [12]). Informally speaking, this algorithm starts with the set $Z$ and as long as there exists a clause in $\Phi$ with its body contained in the current set and its head outside of the current set, the head is added to the current set. More formally the procedure can be described as follows. We start with $F_\Phi^0(Z) := Z$. In a general step, if $F_\Phi^i(Z)$ is a true set then we output $F_h(Z) = F_\Phi^i(Z)$ and stop. Otherwise, let $A \subseteq V \setminus F_\Phi^i(Z)$ denote the set of all variables $v$ for which there exists a clause $B \to v$ in $\Phi$ with $B \subseteq F_\Phi^i(Z)$ and define $F_\Phi^{i+1}(Z) := F_\Phi^i(Z) \cup A$. Note that any CNF $\Phi$ uniquely defines a Horn function, and sometimes we do not have separate notation for that function. In such cases we shall also use $F_\Phi(Z)$ to denote the closure of subset $Z$ with respect to the Horn function represented by $\Phi$.

DEFINITION 2.1. *A pure Horn function $h$ is* key Horn *if it has a CNF representation of the form $\bigwedge_{B \in \mathcal{B}} B \to (V \setminus B)$ for some $\mathcal{B} \subseteq 2^V \setminus \{V\}$. In such a case we shall refer to $h$ as $h_\mathcal{B}$.*

Assume now that $\Phi$ is a pure Horn CNF of the form $\bigwedge_{i=1}^m B_i \to H_i$ where $B_i \neq B_j$ for $i \neq j$. Note that the number of clauses in the CNF is $c_\Phi = \sum_{i=1}^m |H_i|$. The size of the formula can be measured in different ways:
- **(B) number of bodies**: $|\Phi|_B := m$,
- **(BA) body area**: $|\Phi|_{BA} := \sum_{i=1}^m |B_i|$,
- **(TA) total area**: $|\Phi|_{TA} := \sum_{i=1}^m (|B_i| + |H_i|)$,
- **(C) number of clauses (i.e., hyperarcs)**: $|\Phi|_C := c_\Phi$,
- **(BC) number of bodies and clauses**: $|\Phi|_{BC} := m + c_\Phi = \sum_{i=1}^m (|H_i| + 1)$,
- **(L) number of literals**: $|\Phi|_L := \sum_{i=1}^m \big((|B_i| + 1) \cdot |H_i|\big)$.

These measures come up naturally in connection with directed hypergraphs, implicational systems, and CNF representations. For example, (L) corresponds to the size of a CNF when encoded in DIMACS format, a format that is widely accepted as the standard format for Boolean formulas in CNF. The number of clauses (C) is an important parameter for SAT solvers when the Horn formula in question encodes a constraint which is part of a larger problem. Similarly, (TA) is the space needed to store an adjacency list of the corresponding hypergraph, and might be an important parameter for an efficient implementation. The Horn minimization problem is to find a representation that is equivalent to a given Horn formula and has minimum size with respect to $|\cdot|_*$ where $*$ denotes one of the aforementioned functions.

*Example* 2.2. Consider the CNF $\Phi$ introduced in Example 1.1 written as a conjunction of implications $\Phi = (abc \to d) \wedge (d \to efg) \wedge (efg \to abc)$. Note that $\Phi$ represents the key Horn function $h_\mathcal{B}$ defined by the system of bodies $\mathcal{B} = \{\{a, b, c\}, \{d\}, \{e, f, g\}\}$. The CNF $\Phi$ has $m = 3$ different bodies, thus $|\Phi|_B = 3$. Furthermore, it has body area $|\Phi|_{BA} = 7$, total area $|\Phi|_{TA} = 14$, number of clauses $|\Phi|_C = 7$, number of bodies and clauses $|\Phi|_{BC} = 3 + 7 = 10$, and number of literals $|\Phi|_L = 22$. Since every variable occurs exactly once as a positive literal (or as a head of some clause) in $\Phi$, we can conclude that $\Phi$ has the smallest number of clauses among the representations of $h_\mathcal{B}$. However, it is not optimal with respect to the number of literals.

Consider the equivalent formula $\Phi' = (abc \to d) \wedge (efg \to d) \wedge (d \to abcefg)$ which has only $|\Phi'|_L = 20$ literals. On the other hand, $\Phi'$ consists of 8 clauses which is not optimal with respect to the number of clauses. This example demonstrates that different measures may be optimized by different CNF formulas.

**3. Lower bounds for the size of optimal solutions.** The present section provides some simple reductions of the problem and lower bounds for the size of an optimal solution. For a family $\mathcal{B} \subseteq 2^V \setminus \{V\}$, we denote by $\mathcal{B}^\perp$ the family of minimal elements of $\mathcal{B}$. Recall that $h_\mathcal{B}$ denotes the function defined by

$$(3.1) \qquad \Psi_\mathcal{B} = \bigwedge_{B \in \mathcal{B}} B \to (V \setminus B).$$

LEMMA 3.1. *For any measure (∗) and for any $\mathcal{B} \subseteq 2^V \setminus \{V\}$, there exists a $|\cdot|_*$-minimum representation of $h_\mathcal{B}$ that uses exactly the bodies in $\mathcal{B}^\perp$.*

*Proof.* Take a $|\cdot|_*$-minimum representation $\Phi$ for which $|\mathcal{B}_\Phi \setminus \mathcal{B}^\perp|$ is as small as possible. First we show $\mathcal{B}_\Phi \subseteq \mathcal{B}^\perp$. Assume that $B \in \mathcal{B}_\Phi \setminus \mathcal{B}^\perp$. As $B$ is a false set of $h_\mathcal{B}$, there must be a clause $B' \to v$ in $\Psi_\mathcal{B}$ that is falsified by $\chi_B$, implying that $B' \subseteq B$. Therefore there exists a $B'' \in \mathcal{B}^\perp$ such that $B'' \subseteq B' \subseteq B$. If we substitute every clause $B \to v$ of $\Phi$ by $B'' \to v$, then we get another representation of $h_\mathcal{B}$ since $B'' \to v$ is a clause of $\Psi_\mathcal{B}$. Meanwhile, the $|\cdot|_*$ size of the representation does not increase while $|\mathcal{B}_\Phi \setminus \mathcal{B}^\perp|$ decreases, contradicting the choice of $\Phi$.

Next we prove $\mathcal{B}_\Phi \supseteq \mathcal{B}^\perp$. If there exists a $B \in \mathcal{B}^\perp \setminus \mathcal{B}_\Phi$, then $B$ is a true set of $\Phi$ while it is a false set of $h_\mathcal{B}$, contradicting the fact that $\Phi$ is a representation of $h_\mathcal{B}$. ☐

Recall that a *Sperner family* is family of subsets of a finite set in which none of the sets contains another. Lemma 3.1 has an easy corollary.

COROLLARY 3.2. *It suffices to consider Sperner families of bodies defining key Horn functions as an input. Moreover, it is enough to consider pure Horn CNFs using bodies from the input Sperner family when searching for minimum representations.*

For non-key Horn functions, this is not the case. For example, the function defined by implications $a \to b$, $ac \to d$ has five false sets, namely $\{a\}$, $\{a, c\}$, $\{a, d\}$, $\{a, c, d\}$, $\{a, b, c\}$. Clearly, $\{a\}$ has to appear as a body in any representation of the function together with at least one of the other false sets as a body, although it is contained in the other.

From now on we assume that $\mathcal{B}$ is a Sperner family. We also assume that

$$\bigcup_{B \in \mathcal{B}} B = V \quad \text{and} \quad \bigcap_{B \in \mathcal{B}} B = \emptyset.$$

Indeed, if a variable $v \in V \setminus \bigcup_{B \in \mathcal{B}} B$ is not covered by the bodies, then there must be a clause with head $v$ and body in $\mathcal{B}$ in any minimum representation of $h_\mathcal{B}$, and actually one such clause suffices. Furthermore, if $v \in \bigcap_{B \in \mathcal{B}} B$, then we can reduce the problem by deleting it. None of these reductions affects the approximability of the problem.

Recall that the size of the ground set is denoted by $|V| = n$, while $|\mathcal{B}| = m$. The size of an optimal solution with respect to measure function $|\cdot|_*$ is denoted by $OPT_*(\mathcal{B})$. Using these notations Lemma 3.1 has the following easy corollary:

COROLLARY 3.3. *We have $OPT_B(\mathcal{B}) = m$ and $OPT_{BA}(\mathcal{B}) = \sum_{B \in \mathcal{B}} |B|$. Therefore the minimization problems (B) and (BA) are solvable in polynomial time.*

For the remaining measures we prove the following simple lower bound.

LEMMA 3.4. $OPT_*(\mathcal{B}) \geq m$ for all measures $*$, and $OPT_*(\mathcal{B}) \geq n$ for $* \in \{TA, C, BC, L\}$. Furthermore, we have $OPT_{TA}(\mathcal{B}) \geq m + \sum_{i=1}^{m} |B_i|$ and $OPT_L(\mathcal{B}) \geq \max\{n(\delta + 1), 2m\}$, where $\delta$ is the size of a smallest body in $\mathcal{B}$.

*Proof.* By definition, $| \cdot |_B$ is a lower bound for all the other measures, implying $OPT_*(\mathcal{B}) \geq OPT_B(\mathcal{B}) = m$.

To see the second part, observe that $| \cdot |_C$ is a lower bound for the three other measures. Therefore it suffices to prove $OPT_C(\mathcal{B}) \geq n$. By the assumption that for every $v \in V$ there exists a $B \in \mathcal{B}$ not containing $v$, we can conclude by the fact that the closure $F_{h_\mathcal{B}}(B) = V$ and by the way the forward chaining procedure works that every pure Horn CNF representation of $h_\mathcal{B}$ must contain at least one clause with $v$ as its head. This implies $OPT_C(\mathcal{B}) \geq n$.

To see the last part, note that every set $B \in \mathcal{B}$ is the body of at least one clause, verifying the lower bound for (TA). Every variable $v \in V$ is the head of at least one clause, the body of which is of at least size $\delta \geq 1$. Since all clauses are of size at least 2, the bound for (L) follows. $\square$

Let us now introduce a key concept of this paper. For a pair $S, T \subseteq V$ of sets, we denote by $price_*(S, T)$ the minimum $| \cdot |_*$-size of a pure Horn CNF $\Phi$ for which $\mathcal{B}_\Phi \subseteq \mathcal{B}$ and $T \subseteq F_\Phi(S)$, that is,

$$(3.2) \qquad\qquad price_*(S, T) = \min_\Phi \left\{ |\Phi|_* \mid \mathcal{B}_\Phi \subseteq \mathcal{B}, T \subseteq F_\Phi(S) \right\}.$$

*Example* 3.5. Let us consider the set of bodies $\mathcal{B} = \{\{a, b, c\}, \{d\}, \{e, f, g\}\}$ and let us consider $S = \{a, b, c\}$ and $T = \{e, f, g\}$. It is easy to see that $price_C(S, T) = 3$ and that it is realized by a single implication $abc \rightarrow efg$. Actually, as we will show later in Lemma 4.3, we always have that $price_C(S, T) = |T \setminus S|$ provided $S, T \in \mathcal{B}$. However, estimating $price_L(S, T)$ is a bit more tricky. Considering the above single implication $abc \rightarrow efg$ we get that $price_L(S, T) \leq 12$. We can do better by using the small body $d$. In particular, using implications $(abc \rightarrow d) \wedge (d \rightarrow efg)$ we achieve the optimum value $price_L(S, T) = 10$.

The following lemma plays a principal role in our approximability proofs.

LEMMA 3.6. *Let* $\mathcal{B} = \mathcal{B}_1 \cup \cdots \cup \mathcal{B}_q$ *be a partition of* $\mathcal{B}$ *and let* $B_i \in \mathcal{B}_i$ *for* $i = 1, \ldots, q$. *Then*

$$(3.3) \qquad\qquad OPT_*(\mathcal{B}) \geq \sum_{i=1}^{q} \min\{price_*(B_i, B) \mid B \in \mathcal{B} \setminus \mathcal{B}_i\}$$

*for all six measures* $*$.

*Proof.* Take a minimum representation $\Phi$ with respect to $| \cdot |_*$ which uses bodies only from $\mathcal{B}$. Such a representation exists by Lemma 3.1. We claim that the contribution of the clauses with bodies in $\mathcal{B}_i$ to the total size of $\Phi$ is at least $\min\{price_*(B_i, B) \mid B \in \mathcal{B} \setminus \mathcal{B}_i\}$ for each $i = 1, \ldots, q$. This would prove the lemma as the $\mathcal{B}_i$'s form a partition of $\mathcal{B}$.

To see the claim, take an index $i \in \{1, \ldots, q\}$ and let $B'$ be the first body (more precisely, one of the first bodies) not contained in $\mathcal{B}_i$ that is reached by the forward chaining procedure from $B_i$ with respect to $\Phi$. Every clause that is used to reach $B'$ from $B_i$ has its body in $\mathcal{B}_i$ and their contribution to the size of the representation is lower bounded by $price_*(B_i, B')$, thus concluding the proof. $\square$

**4. Approximability results for (TA), (C), (BC), and (L).** Given a Sperner family $\mathcal{B} \subseteq 2^V \setminus \{V\}$, we can associate with it a complete directed graph $D_{\mathcal{B}}$ by defining $V(D_{\mathcal{B}}) = \mathcal{B}$ and $E(D_{\mathcal{B}}) = \mathcal{B} \times \mathcal{B}$. We refer to $D_{\mathcal{B}}$ as the *body graph* of $\mathcal{B}$.

For any subset $E' \subseteq E(D_{\mathcal{B}})$, define

$$(4.1) \qquad \Phi_{E'} = \bigwedge_{(B,B') \in E'} B \to (B' \setminus B).$$

Note that if $E' \subseteq E(D_{\mathcal{B}})$ forms a strongly connected spanning subgraph of $D_{\mathcal{B}}$, then $\Phi_{E'}$ is a representation of $h_{\mathcal{B}}$. Let us add that not all representations arise this way, in particular, minimum representations might have significantly smaller size.

LEMMA 4.1. *If $E'$ is a Hamiltonian cycle in $D_{\mathcal{B}}$, then $\Phi_{E'}$ defined in (4.1) provides a $k$-approximation for all measures, where $k$ is an upper bound on the sizes of bodies in $\mathcal{B}$.*

*Proof.* By Lemma 3.1, there exists a minimum representation $\Phi$ of $h_{\mathcal{B}}$ such that $\mathcal{B}_\Phi = \mathcal{B}$. Since $|B' \setminus B|$ is at most $k$ for all arcs $(B, B') \in E'$, the statement follows. □

In fact, for (B) and (BA) (4.1) gives an optimal representation for any strongly connected spanning $E'$. Furthermore, if $E'$ is a Hamiltonian cycle, we get a $\frac{2k}{k+1}$-approximation for (TA) based on the fact that the total area of any representation is lower bounded by $\sum_{B \in \mathcal{B}} |B|$.

THEOREM 4.2. *If $E'$ is a Hamiltonian cycle in $D_{\mathcal{B}}$, then $\Phi_{E'}$ defined in (4.1) provides a $\frac{2k}{k+1}$-approximation for (TA), where $k$ is an upper bound on the sizes of bodies in $\mathcal{B}$.*

*Proof.* By Lemma 3.4, $OPT_{TA}(\mathcal{B}) \geq m + \sum_{i=1}^{m} |B_i|$. Recall that $|B_i| \leq k$ for $i = 1, \ldots, m$. The total area of $\Phi_{E'}$ is $|\Phi_{E'}|_{TA} = \sum_{i=1}^{m}(|B_i| + |B_{i+1} \setminus B_i|) \leq m + \sum_{i=1}^{m} |B_i| + \sum_{i=1}^{m}(|B_i| - 1) \leq OPT_{TA}(\mathcal{B}) + \frac{k-1}{k+1} OPT_{TA}(\mathcal{B}) = \frac{2k}{k+1} OPT_{TA}(\mathcal{B})$, concluding the proof. □

The observation that a strongly connected subgraph of the body graph corresponds to a representation of $h_{\mathcal{B}}$, as in (4.1), suggests the reduction of our problem to the problem of finding a minimum weight strongly connected spanning subgraph in a directed graph with arc-weight $price_*(B, B')$ for $(B, B') \in E(D_{\mathcal{B}})$. The optimum solution to this problem (MWSCS) is an upper bound for the minimum $| \cdot |_*$-size of a representation of $h_{\mathcal{B}}$. As there are efficient constant-factor approximations for MWSCS [17], this approach may look promising. There are, however, two difficulties. First, in Section 5 we show that computing $price_L$ is NP-complete. Second, even when $price_*$ is efficiently computable (for measures (C) and (BC)), the upper bound obtained in this way may be off by a factor of $\Omega(n)$ from the optimum, see Section 6 for a construction.

In what follows, we overcome these difficulties. An *in-arborescence* is a directed, rooted tree in which all edges point towards the root. An in-arborescence is called *spanning* if the underlying tree is spanning. A *branching* is a directed forest in which every connected component forms an in-arborescence. For (C), instead of a strongly connected spanning subgraph, we compute a minimum weight spanning in-arborescence and extend that to a representation of $h_{\mathcal{B}}$. The same approach works for (BC) as well. For (L), the situation is more complicated. First, we develop an efficient approximation algorithm for $price_L$. Next, we compute a minimum weight spanning in-arborescence where its root is pre-specified. Finally, we extend the corresponding pure Horn CNF to a representation of $h_{\mathcal{B}}$. We show that the cost of the arborescences

built is at most a multiple of the optimum by a logarithmic factor, which in turn ensures the improved approximation factor.

**4.1. Clause and body-clause minimum representations.** In this section we consider (C) and (BC) and show that the simple algorithm described in Procedure 1 provides the stated approximation factor. We note that a minimum weight spanning in-arborescence of a directed graph can be found in polynomial time, see [10, 16].

---

**Procedure 1:** Approximation of (C) and (BC)

> **1** Determine a minimum $price_C$-weight spanning in-arborescence $\overline{T}$ of $D_{\mathcal{B}}$.
>    /* Denote by $B_0$ the body corresponding to the root of $\overline{T}$. */
> **2** Output $\Phi = \Phi_{\overline{T}} \wedge B_0 \to (V \setminus B_0)$.
>    /* Here $\Phi_{\overline{T}}$ is defined as in (4.1). */

---

Observe that $price_C$ is easy to compute.

LEMMA 4.3. $price_C(B, B') = |B' \setminus B|$ for $B, B' \in \mathcal{B}$.

*Proof.* Take a pure Horn CNF $\Phi$ attaining the minimum in (3.2). As every variable in $B' \setminus B$ is reached by the forward chaining procedure from $B$ with respect to $\Phi$, each such variable must be a head of at least one clause in $\Phi$. That is, $\Phi$ contains at least $|B' \setminus B|$ clauses. On the other hand, $B \to (B' \setminus B)$ uses exactly $|B' \setminus B|$ clauses, hence $price_C(B, B') = |B' \setminus B|$ as stated. □

LEMMA 4.4. *Let $\overline{T}$ denote a minimum $price_C$-weight spanning in-arborescence in $D_{\mathcal{B}}$. Then*

$$|\Phi_{\overline{T}}|_C \leq \lceil \log k \rceil OPT_C(\mathcal{B}) + \max\{0, m - k\},$$

*where $k$ is an upper bound on the sizes of bodies in $\mathcal{B}$.*

*Proof.* We construct a subgraph $T$ of $D_{\mathcal{B}}$ such that (i) it is a spanning in-arborescence, and (ii) $|\Phi_T|_C \leq \lceil \log k \rceil OPT_C(\mathcal{B}) + \max\{0, m - k\}$. This proves the lemma as the weight of $T$ upper bounds the weight of $\overline{T}$.

We start with the digraph $T_1$ on node set $\mathcal{B}$ that has no arcs. In a general step of the algorithm, $T_i$ will denote the graph constructed so far. We maintain the property that $T_i$ is a branching, that is, a collection of node-disjoint in-arborescences spanning all nodes. In an iteration, for each such in-arborescence we choose an arc of minimum weight with respect to $price_C$ that goes from the root of the in-arborescence to some other component. We add these arcs to $T_i$, and for each directed cycle created, we delete one of its arcs. This results in a graph $T_{i+1}$ with at most half the number of weakly connected components that $T_i$ has, all being in-arborescences. We repeat this until the number of components becomes at most $\max\{1, m/k\}$. To reach this, we need at most $\lceil \log k \rceil$ iterations. Finally, we choose one of the roots of the components and add an arc from all the other roots to this one, obtaining a spanning in-arborescence $T$.

It remains to show that $T$ also satisfies (ii). In the final stage, we add at most $\max\{1, m/k\} - 1$ arcs to $T$, which corresponds to at most $k(\max\{1, m/k\} - 1) \leq \max\{0, m-k\}$ clauses in $\Phi_T$. Now we bound the rest of $\Phi_T$. In iteration $i$, components of $T_i$ define a partition $\mathcal{B} = \mathcal{B}_1 \cup \cdots \cup \mathcal{B}_q$. Let us denote by $B_j$ the body corresponding to the root of the arborescence with node-set $\mathcal{B}_j$. Let us consider the arcs $\{(B_j, B'_j) \mid$

$j = 1, \ldots, q\}$ chosen to be added in the $i$th iteration. Now we obtain

$$|\Phi_{T_{i+1} \setminus T_i}|_C \leq \sum_{j=1}^{q} price_C(B_j, B'_j) = \sum_{j=1}^{q} \min_{B \in \mathcal{B} \setminus \mathcal{B}_j} price_C(B_j, B) \leq OPT_C(\mathcal{B}).$$

The first inequality follows from the construction of $T$. The equality follows from the criterion to choose the arcs to be added. The last inequality follows from Lemma 3.6. Since we have at most $\lceil \log k \rceil$ iterations, the lemma follows.                                                      □

THEOREM 4.5. *For key Horn functions, there exists a polynomial time* $\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2, k\}$-*approximation algorithm for (C) and (BC), where $k$ is an upper bound on the sizes of bodies in $\mathcal{B}$.*

*Proof.* We first show that $\Phi$ provided by Procedure 1 is a $\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2\}$-approximation for (C) and (BC). Note that $\Phi$ is a subformula of $\Psi_{\mathcal{B}}$ defined by (3.1) since all bodies in $\Phi$ are from $\mathcal{B}$. Furthermore, by our construction, $F_\Phi(B) = V$ for all $B \in \mathcal{B}$. This implies that the output $\Phi$ represents $h_{\mathcal{B}}$. Using Lemma 4.4 and the fact that we added $|V \setminus B_0| \leq n$ clauses to $\Phi_T$ in Step 2, we obtain

$$|\Phi|_C \leq \lceil \log k \rceil OPT_C(\mathcal{B}) + \max\{0, m - k\} + n.$$

By Lemma 3.4, this gives a $(\lceil \log k \rceil + 2)$-approximation, while setting $k = n$ gives a $(\lceil \log n \rceil + 1)$-approximation. By Lemma 3.1, $OPT_{BC}(\mathcal{B}) = |\mathcal{B}| + OPT_C(\mathcal{B})$. Since $|\Phi|_{BC} = |\mathcal{B}| + |\Phi|_C$, the same approximation ratios as above follow for (BC) as well.

Finally, Lemma 4.1 provides a different pure Horn CNF that is a $k$-approximation for (C) and (BC).                                                      □

**4.2. Literal minimum representations.** In this section we consider (L). The first difficulty that we have to overcome is that, unlike in the case of (C) and (BC), computing $price_L$ is NP-hard as we show in Section 5. To circumvent this, we give an $O(1)$-approximation algorithm for $price_L(S, S')$ for any pair of sets $S, S' \subseteq V$. Note that if $S$ does not contain a body $B \in \mathcal{B}$ then $price_L(S, S') = \infty$, hence we assume that this is not the case.

We first analyze the structure of a pure Horn CNF $\Phi$ attaining the minimum in (3.2) for (L). Starting the forward chaining procedure from $S$ with respect to $\Phi$, let $W_i$ denote the set of variables reached within the first $i$ steps. That is, $S = W_0 \subsetneq W_1 \subsetneq \cdots \subsetneq W_t \supseteq S'$. We choose $\Phi$ in such a way that $t$ is as small as possible (among those pure Horn CNFs that already minimize (3.2) for (L)). Let $B_i \in \mathcal{B}$ be a smallest body contained in $W_i$ for $i = 0, \ldots, t-1$ and set $B_t := S'$.

PROPOSITION 4.6. $B_i \not\subseteq W_{i-1}$ *for* $i = 1, \ldots, t$.

*Proof.* Suppose to the contrary that $B_i \subseteq W_{i-1}$ for some $1 \leq i \leq t-1$. By the definition of forward chaining, every variable $v \in W_{i+1} \setminus W_i$ is reached through a clause $B \rightarrow v$ where $B \cap (W_i \setminus W_{i-1}) \neq \emptyset$. Now substitute each such clause by $B_i \rightarrow v$. As $|B_i| \leq |B|$, the $|\cdot|_L$ size of the CNF does not increase. However, the number of steps in the forward chaining procedure decreases by at least one, contradicting the choice of $\Phi$. Finally, $S' = B_t \subseteq W_{t-1}$ would contradict the minimality of $t$.                                                      □

Proposition 4.6 immediately implies that $|B_0| > |B_1| > \ldots > |B_{t-1}|$.

PROPOSITION 4.7. $W_{i+1} \setminus W_i \subseteq B_{i+1}$ *for* $i = 0, \ldots, t-1$.

*Proof.* Let $i$ be the smallest index that violates the condition. Take an arbitrary variable $v \in W_{i+1} \setminus W_i$ for which $v \notin B_{i+1}$. Then $v$ is reached in the $(i+1)$th step

of the forward chaining procedure from a body of size at least $|B_i|$. If we substitute this clause by $B_{i+1} \to v$, the resulting pure Horn CNF still satisfies $F_\Phi(B_0) \supseteq S'$ but has smaller $|\cdot|_L$ size by $|B_{i+1}| < |B_i|$, contradicting the minimality of $\Phi$.    $\square$

By Proposition 4.7, $W_{i+1} \setminus W_i = B_{i+1} \setminus (S \cup \bigcup_{j=1}^i B_j)$. Define

$$\Phi^{(1)} := \bigwedge_{i=0}^{t-1} B_i \to (B_{i+1} \setminus (S \cup \bigcup_{j=1}^i B_j)).$$

Observe that $\Phi^{(1)}$ has a simple structure which is based on a linear order of bodies $B_0, \ldots, B_t$.

PROPOSITION 4.8. $|\Phi^{(1)}|_L = |\Phi|_L$.

*Proof.* Take an arbitrary variable $v \in B_{i+1} \setminus (S \cup \bigcup_{j=1}^i B_j)$ for some $i = 0, \ldots, t-1$. By the observation above, $v \in W_{i+1} \setminus W_i$. This means that $\Phi$ has at least one clause entering $v$, say $B \to v$, for which $B \subseteq W_i$ and so $|B| \geq |B_i|$. However, $\Phi^{(1)}$ has exactly one clause entering $v$, namely $B_i \to v$. This implies that $|\Phi^{(1)}|_L \leq |\Phi|_L$, and equality holds by the minimality of $\Phi$.    $\square$

The proposition implies that $\Phi^{(1)}$ also realizes $price_L(S, S')$. As we show later in Theorem 5.8, computing $price_L(S, S')$ is NP-hard and thus we do not know any efficient algorithm to compute $\Phi^{(1)}$. Using the next two propositions, we define a pure Horn CNF that approximates $\Phi^{(1)}$ well and can be computed efficiently. We then use it to show in Theorem 4.13 that there is a polynomial time $\Theta(\log k)$ approximation algorithm for (L).

Let $i_0 = 0$ and for $j > 0$ let $i_j$ denote the smallest index for which $|B_{i_j}| \leq |B_{i_{j-1}}|/2$. Let $r-1$ be the largest value for which $B_{i_{r-1}}$ exists and set $B_{i_r} := S'$. Now define

$$\Phi^{(2)} := \bigwedge_{j=0}^{r-1} B_{i_j} \to (B_{i_{j+1}} \setminus (S \cup \bigcup_{\ell=1}^j B_{i_\ell})).$$

It is easy to see that $F_{\Phi^{(2)}}(S) \supseteq S'$.

PROPOSITION 4.9. $|\Phi^{(2)}|_L \leq 2|\Phi^{(1)}|_L$.

*Proof.* Take an arbitrary variable $v \in B_{i_{j+1}} \setminus (S \cup \bigcup_{\ell=1}^j B_{i_\ell})$ for some $j = 0, \ldots, r-1$. Then both $\Phi^{(1)}$ and $\Phi^{(2)}$ contain a single clause entering $v$. Namely, $v$ is reached from $B_{i_{j+1}-1}$ in $\Phi^{(1)}$ and from $B_{i_j}$ in $\Phi^{(2)}$. By the definition of the sequence $i_0, i_1, \ldots, i_{r-1}$, we get $|B_{i_j}| \leq 2|B_{i_{j+1}-1}|$, concluding the proof.    $\square$

Although $\Phi^{(2)}$ gives a 2-approximation for $|\Phi|_L$, it is not clear how we could find such a representation, because bodies $B_{i_j}, j = 0, \ldots, r-1$ depend on $\Phi$ which is hard to compute. Define

$$\Phi^{(3)} := \bigwedge_{j=0}^{r-1} B_{i_j} \to (B_{i_{j+1}} \setminus (S \cup B_{i_j})).$$

The only difference between $\Phi^{(2)}$ and $\Phi^{(3)}$ is that we add unnecessary clauses to the representation. The distinguishing feature of $\Phi^{(3)}$ is that each of its implications depends only on two bodies $B_{i_j}$ and $B_{i_{j+1}}$, and thus $\Phi^{(3)}$ represents a path from a body contained in $S$ to $S'$ in the body graph extended with a new node $S'$. This will allow us to obtain a CNF which is not longer than $\Phi^{(3)}$ and allows to derive $S'$ from

$S$ by forward chaining (see Lemma 4.11). The next claim shows that the size of the formula cannot increase too much.

PROPOSITION 4.10. $|\Phi^{(3)}|_L \leq \frac{27}{17}|\Phi^{(2)}|_L$.

*Proof.* Take an arbitrary variable $v$ that appears as the head of a clause in the representation $\Phi^{(3)}$. Let $j$ be the smallest index for which $v \in B_{i_{j+1}} \setminus (S \cup \bigcup_{\ell=1}^{j} B_{i_\ell})$. Then $\Phi^{(2)}$ contains a single clause entering $v$, namely $B_{i_j} \to v$. On the other hand, the set $\{B_{i_j} \to v\} \cup \{B_{i_\ell} \to v \mid \ell = j+2, \ldots, r-1\}$ contains all the clauses of $\Phi^{(3)}$ that enter $v$. By the definition of the sequence $i_0, i_1, \ldots, i_{r-1}$, we get $\sum_{\ell=j+2}^{r-1}(|B_{i_\ell}| + 1) = (r - j - 2) + \sum_{\ell=j+2}^{r-1} |B_{i_\ell}| \leq \lfloor \log |B_{i_{j+1}}| \rfloor + |B_{i_j}|/2 - 1 \leq \lfloor \log |B_{i_j}| \rfloor + |B_{i_j}|/2 - 2$. We get at most this many extra literals in $\Phi^{(3)}$ on top of the $|B_{i_j}| + 1$ literals in $\Phi^{(2)}$. As $\lfloor \log x \rfloor/(x+1) + x/(2(x+1)) - 2/(x+1) \leq 10/17$ for $x \in \mathbb{Z}_+$, the statement follows. $\square$

By Propositions 4.8, 4.9 and 4.10,

$$(4.2) \qquad |\Phi^{(3)}|_L \leq \frac{27}{17}|\Phi^{(2)}|_L \leq \frac{54}{17}|\Phi^{(1)}|_L = \frac{54}{17}|\Phi|_L.$$

LEMMA 4.11. *There exists an efficient algorithm to construct a pure Horn CNF* $\Lambda(S, S')$ *such that* $|\Lambda(S, S')|_L \leq \frac{54}{17} \, price_L(S, S')$, $\mathcal{B}_{\Lambda(S,S')} \subseteq \mathcal{B}$, *and* $F_{\Lambda(S,S')}(S) \supseteq S'$.

*Proof.* We consider an extension of the body graph by adding $S'$ to $V(D_\mathcal{B})$. We also define arc-weights by setting $w(B, B') := |B' \setminus (S \cup B)|(|B|+1)$ for $B, B' \in \mathcal{B} \cup \{S'\}$. Let $B_0$ be a smallest body contained in $S$ (as defined before Proposition 4.6). Compute a shortest path $P$ from $B_0$ to $S'$ and define

$$(4.3) \qquad \Lambda(S, S') = \bigwedge_{(B,B')\in P} B \to (B' \setminus (S \cup B)).$$

Note that, by definition, $|\Lambda(S, S')|_L$ is the weight of the shortest path $P$, while $|\Phi^{(3)}|_L$ is the length of one of the paths from $S$ to $S'$. By (4.2), $|\Lambda(S, S')|_L \leq |\Phi^{(3)}|_L \leq \frac{54}{17}|\Phi|_L$. That is, $\Lambda(S, S')$ provides a $\frac{54}{17}$-approximation for $price_L(S, S')$ as required, finishing the proof of the lemma. $\square$

We prove that the algorithm described in Procedure 2 provides the stated approximated factor for (L). We note that a minimum weight spanning in-arborescence of a directed graph rooted at a fixed node can be found in polynomial time, see [10,16]. Let $B_{\min}$ be a smallest body in $\mathcal{B}$, let $\delta := |B_{\min}|$ and denote $\mathcal{B}' = \mathcal{B} \setminus \{B_{\min}\}$. We define the weight of an arc $(B, B') \in E(D_\mathcal{B})$ in the body graph to be $w(B, B') = |\Lambda(B, B')|_L$.

---

**Procedure 2:** Approximation of (L)

1. Let $B_{\min}$ be a smallest body in $\mathcal{B}$.
2. Set $w(B, B') = |\Lambda(B, B')|_L$ for $(B, B') \in E(D_\mathcal{B})$.
3. Determine a minimum $w$-weight spanning in-arborescence $\overline{T}$ of $D_\mathcal{B}$ such that $\overline{T}$ is rooted at $B_{\min}$.
4. Output $\Phi = \bigwedge_{(B,B')\in\overline{T}} \Lambda(B, B') \wedge (B_{\min} \to (V \setminus B_{\min}))$.
   /* Here $\Lambda(B, B')$ is defined as in (4.3). */

---

The proof of the following lemma is very similar to the proof of Lemma 4.4. There are a few differences: The first one is that we use a different cost function on the edges (the approximation value $|\Lambda(B, B')|_L$ given by Lemma 4.11 instead of

$price_C(B, B')$). We also have a slightly different terminating condition ($m/k^2$ instead of $m/k$). Finally, in the last step of the construction we do not use an arbitrary root, but we make sure that $B_{\min}$ is the root of the constructed in-arborescence.

LEMMA 4.12. *Let $\overline{T}$ denote a minimum $w$-weight spanning in-arborescence in $D_{\mathcal{B}}$ such that $\overline{T}$ is rooted at $B_{\min}$. Then*

$$\left| \bigwedge_{(B,B')\in\overline{T}} \Lambda(B, B') \right|_L \leq \left( \frac{108}{17}\lceil \log k\rceil + 1 \right) OPT_L(\mathcal{B}),$$

*where $k$ is the size of a largest body in $\mathcal{B}$.*

*Proof.* We construct a subgraph $T$ of $D_{\mathcal{B}}$ such that (i) it is a spanning in-arborescence, and (ii) $|\bigwedge_{(B,B')\in T} \Lambda(B, B')|_L \leq (\frac{108}{17}\lceil \log k\rceil + 1)OPT_L(\mathcal{B})$. This clearly proves the lemma as the weight of $T$ upper bounds the weight of $\overline{T}$.

We start with the directed graph $T_1$ on node set $\mathcal{B}$ that has no arcs. In a general step of the algorithm, $T_i$ will denote the graph constructed so far. We maintain the property that $T_i$ is a branching, that is, a collection of node-disjoint in-arborescences spanning all nodes. In an iteration, for each such in-arborescence we choose an arc of minimum weight with respect to $w$ that goes from the root of the in-arborescence to some other component. We add these arcs to $T_i$, and for each directed cycle created, we delete one of its arcs. This results in a graph $T_{i+1}$ with at most half the number of weakly connected components that $T_i$ has, all being in-arborescences. We repeat this until the number of components becomes at most $\max\{1, m/k^2\}$. To reach this, we need at most $\lceil \log k^2\rceil \leq 2\lceil \log k\rceil$ iterations. Finally, we add an arc from all the other roots to $B_{\min}$ and delete all the arcs leaving $B_{\min}$, obtaining a spanning in-arborescence $T$ rooted at $B_{\min}$.

It remains to show that $T$ also satisfies (ii). In the final stage, we add at most $\max\{1, m/k^2\}$ arcs to $T$ whose total weight is upper bounded by $(k+1)\delta \max\{1, m/k^2\}$. ∎ Since $k+1 \leq n$, we have that $(k+1)\delta \leq n\delta$. We have that $\frac{(k+1)\delta m}{k^2} = \frac{k+1}{k} \cdot \frac{\delta}{k} \cdot m \leq 2m$ where the inequality holds, because $(k+1)/k \leq 2$ for $k \geq 1$ and $\delta \leq k$. Together we get that the total weight of arcs added in the last step is upper bounded by $(k+1)\delta \max\{1, m/k^2\} \leq \max\{n\delta, 2m\} \leq OPT_L(\mathcal{B})$ where the last inequality follows by Lemma 3.4. Now we bound the rest of $\bigwedge_{(B,B')\in T} \Lambda(B, B')$. In iteration $i$, components of $T_i$ define a partition $\mathcal{B} = \mathcal{B}_1 \cup \cdots \cup \mathcal{B}_q$. Let us denote by $B_j$ the body corresponding to the root of the arborescence with node-set $\mathcal{B}_j$. Let us consider the arcs $\{(B_j, B'_j) \mid j = 1, \ldots, q\}$ chosen to be added in the $i$th iteration. Now we obtain

$$\left| \bigwedge_{(B,B')\in T_{i+1}\setminus T_i} \Lambda(B, B') \right|_L = \sum_{j=1}^q w(B_j, B'_j) = \sum_{j=1}^q \min_{B\in\mathcal{B}\setminus\mathcal{B}_j} w(B_j, B)$$

$$\leq \frac{54}{17} \sum_{j=1}^q \min_{B\in\mathcal{B}\setminus\mathcal{B}_j} price_L(B_j, B) \leq \frac{54}{17}OPT_L(\mathcal{B}),$$

where the first and second inequalities follow by Lemmas 4.11 and 3.6, respectively. Since we have at most $2\lceil \log k\rceil$ iterations, the lemma follows. □

THEOREM 4.13. *For key Horn functions, there exists a polynomial time $\min\{\frac{108}{17}\lceil \log k\rceil + 2, k\}$-approximation algorithm for (L), where $k$ is the size of a largest body in $\mathcal{B}$.*

*Proof.* We first show that $\Phi$ provided by Procedure 2 is a $(\frac{108}{17}\lceil \log k \rceil + 2)$-approximation for (L). Note that $\Phi$ is a subformula of $\Psi_{\mathcal{B}}$ defined by (3.1) since all bodies in $\Phi$ are from $\mathcal{B}$. Furthermore, by our construction, $F_\Phi(B) = V$ for all $B \in \mathcal{B}$. This implies that the output $\Phi$ represents $h_{\mathcal{B}}$. By Lemma 3.4, we add at most $n(\delta + 1) \leq OPT_L(\mathcal{B})$ literals to $\bigwedge_{(B,B') \in T} \Lambda(B, B')$ in Step 4. This, together with Lemma 4.12, implies the theorem. $\square$

**5. Hardness of computing** $price_L$. In this section we prove that computing $price_L$ is NP-hard. Given a sequence $\mathcal{S} = (S_0, S_1, ..., S_s)$ of sets we associate to it a pure Horn CNF

$$(5.1) \qquad \Phi_{\mathcal{S}} = \bigwedge_{i=0}^{s-1} \left( S_i \rightarrow \left( S_{i+1} \setminus \bigcup_{j \leq i} S_j \right) \right).$$

We denote by $cost_L(\mathcal{S}) = cost_L(S_0, ..., S_s)$ the $L$-measure (number of literals) of $\Phi_{\mathcal{S}}$, i.e.,

$$cost_L(\mathcal{S}) = cost_L(S_0, ..., S_s) = \sum_{i=0}^{s-1} (|S_i| + 1) \cdot \left| S_{i+1} \setminus \left( \bigcup_{j \leq i} S_j \right) \right|.$$

Let us note that we view $\mathcal{S}$ as a sequence of subsets. This is because in this section we are concerned with sequences between given sets $S_0$ and $S_s$ that minimize $cost_L(\mathcal{S})$ over all possible sequences $\mathcal{S}$ that start at $S_0$ and end at $S_s$.

By Proposition 4.6 we can assume for such sequences that $|S_0| > |S_1| > \cdots > |S_{s-1}|$. Note also that $cost_L(\mathcal{S}) = cost_L(\mathcal{S}, \emptyset)$. In other words, concatenating/deleting empty sets from the end of the sequence does not change the $cost_L$ value.

We will show NP-hardness for computing $price_L$ by a reduction from 3-SAT. Consider a 3-CNF (exactly 3 literals in each clause) $\Phi = \bigwedge_{k=1}^{m} C_k$ in which every variable $x_i$, $i = 1, ..., n$ appears at most 4 times. SAT is NP-complete for this family of CNFs [26]. For a clause $C \in \Phi$, let us denote by $\mathcal{C}(C)$ the set of eight possible clauses consisting of the three variables in $C$. For example, if $C = (\overline{x}_1 \vee x_2 \vee x_4)$, then $\mathcal{C}(C) = \{(x_1 \vee x_2 \vee x_4), (\overline{x}_1 \vee x_2 \vee x_4), (x_1 \vee \overline{x}_2 \vee x_4), (x_1 \vee x_2 \vee \overline{x}_4), (\overline{x}_1 \vee \overline{x}_2 \vee x_4), (\overline{x}_1 \vee x_2 \vee \overline{x}_4), (x_1 \vee \overline{x}_2 \vee \overline{x}_4), (\overline{x}_1 \vee \overline{x}_2 \vee \overline{x}_4)\}$. Furthermore, let

$$M = \bigcup_{C \in \Phi} \mathcal{C}(C).$$

We regard $M$ as a multiset, that is, if two clauses $C$ and $C'$ share the same three variables then $\mathcal{C}(C)$ and $\mathcal{C}(C')$ are considered to be disjoint, and so the corresponding eight clauses are added for both of them. Accordingly, $\Phi \cap \mathcal{C}(C)$ is defined to be $C$. Let us denote by $\delta_i$ the number of clauses in $M$ containing positive literal $x_i$. Note that for all $i = 1, ..., n$, the negative literal $\overline{x}_i$ also appears in $\delta_i$ clauses of $M$, and $\delta_i \leq 16$.

Let us introduce

$$M(x_i) = \{C \in M \mid x_i \in C\} \quad \text{and} \quad M(\bar{x}_i) = \{C \in M \mid \bar{x}_i \in C\}.$$

Let us next define sets $T$, $B_j$ $(j = 0, ..., n)$ and $A_j$ $(j = 1, ..., n+1)$ to be pairwise disjoint and disjoint from $M$, such that for some integer parameters $\alpha$, $\beta$ and $\tau$ we have $|T| = \tau$, $|A_j| = \alpha$ $(j = 1, ..., n+1)$, and $|B_j| = \beta$ $(j = 0, ..., n)$.

Let us further introduce

$$X_i = \left( \bigcup_{j=i}^{n} B_j \right) \cup \left( \bigcup_{j=1}^{i} A_j \right) \cup M(x_i), \qquad \text{and}$$

$$Y_i = \left( \bigcup_{j=i}^{n} B_j \right) \cup \left( \bigcup_{j=1}^{i} A_j \right) \cup M(\bar{x}_i),$$

for $i = 0, ..., n+1$. Note that since $x_0$ and $x_{n+1}$ are not variables of $\Phi$, we have $X_0 = Y_0 = B_0 \cup \cdots \cup B_n$ and $X_{n+1} = Y_{n+1} = A_1 \cup \cdots \cup A_{n+1}$. Finally, let us define $S = X_0$, $Z = X_{n+1} \cup \Phi$, and set

(5.2) $$\mathcal{B}_\Phi = \{S, Z, T\} \cup \{X_i, Y_i \mid i = 1, ..., n\}.$$

Our aim is to show that with these definitions and appropriate choices for parameters $\alpha$, $\beta$ and $\tau$, the quantity $price_L(S, T)$, with respect the family $\mathcal{B}_\Phi$, attains its minimum possible value if and only if $\Phi$ is a satisfiable formula.

We plan to choose $\tau \gg \beta \gg \alpha \gg \max\{n, m\}$ such that we have

$$|S| > |X_1| = |Y_1| > \cdots > |X_n| = |Y_n| > |Z|.$$

Given this, let us recall that an optimal solution realizing $price_L(S, T)$ with respect to the family $\mathcal{B}_\Phi$ involves sets from $\mathcal{B}_\Phi$ in strictly decreasing order of their size by Proposition 4.6. To handle such sequences of sets, we introduce $\mathcal{P}(\sigma) = (P_1^{\sigma_1}, P_2^{\sigma_2}, ..., P_n^{\sigma_n})$ for $\sigma \in \{0, 1, *\}^{[n]}$, where for $1 \le i \le n$ and $\xi \in \{0, 1, *\}$ we have

$$P_i^\xi = \begin{cases} X_i & \text{if} \quad \xi = 1, \\ Y_i & \text{if} \quad \xi = 0. \end{cases}$$

Note that for index $i$ with $\sigma_i = *$ the corresponding sequence $\mathcal{P}(\sigma)$ simply skips both $X_i$ and $Y_i$. For instance, for $n = 4$ and $\sigma = (1, *, 0, *)$ we have $\mathcal{P}(\sigma) = (X_1, Y_3)$.

Note that an optimal sequence realizing $price_L(S, T)$, with respect to $\mathcal{B}_\Phi$, has the form $(S, \mathcal{P}(\sigma), T)$ or $(S, \mathcal{P}(\sigma), Z, T)$ for some $\sigma \in \{0, 1, *\}^{[n]}$. For this reason, we also use the notation $\sigma_0 = 1$ and $P_0^{\sigma_0} = P_0^1 = S = X_0 = Y_0$. For such sequences we also introduce

$$W_i(\sigma) = S \cup \left( \bigcup_{\substack{j \le i \\ \sigma_j \ne *}} P_j^{\sigma_j} \right)$$

for $i = 1, ..., n$ to denote the initial segments covered by the sequence.

In the rest of this section, we shall show that with the right choice of the parameters $\tau$, $\beta$, and $\alpha$, any optimal sequence realizing $price_L(S, T)$ has the form $(S, \mathcal{P}(\sigma), Z, T)$ for some $\sigma \in \{0, 1\}^{[n]}$. In particular, we will show the following properties of optimal sequences:

(I) $Z$ is a part of any optimal sequence, and

(II) for every $i$, $X_i$ or $Y_i$ is a part of the sequence.

Later, we will show that in any optimal sequence $\sigma$ minimizes the number of unsatisfied clauses in $\Phi$. In particular, there is a quantity $f$ which depends on the structure of formula $\Phi$ such that $price_L(S, T) = f + (|Z| + 1) \cdot |T|$ if $\Phi$ is satisfiable

and $price_L(S,T) > f + (|Z| + 1) \cdot |T|$ if $\Phi$ is not satisfiable. The reason for this is that if $P_i^{\sigma_i}$ is a part of the sequence and a clause $C$ in $\Phi$ is satisfied by literal $x_i$ or $\bar{x}_i$ (depending on the value of $\sigma_i$), then $C$ is already added to the forward chaining closure when reaching $P_i^{\sigma_i}$. Thus when adding $Z$ to the sequence, we do not have to add a clause with head $C$.

For simplicity, introduce $\delta_0 = \delta_{n+1} = 0$ and recall that $\delta_i = |M(x_i)| = |M(\bar{x}_i)|$ for $i = 1, ..., n$. We start by observing the following easy to see relations that we will rely on in our proof sometimes without mentioning them explicitly:

(i) $\delta_i = |X_i \cap M| = |Y_i \cap M| \le 16$ for $i = 0, ..., n+1$,

(ii) $|Z| = (n+1)\alpha + m$,

(iii) $|X_i| = |Y_i| = (n - i + 1)\beta + i\alpha + \delta_i$ for $i = 0, ..., n+1$,

In what follows, we show first that, with a right choice of parameters, such an optimal solution must include $Z$, thus proving statement (I).

LEMMA 5.1. *For all $\sigma \in \{0, 1, *\}^{[n]}$, we have*

$$cost_L(S, \mathcal{P}(\sigma), T) > cost_L(S, \mathcal{P}(\sigma), Z, T)$$

*whenever*

(5.3)          $(\beta - \alpha - m) \cdot \tau > ((n+1)\beta + 17) \cdot ((n+1)\alpha + 8m).$

*Proof.* Define $i \in \{0, 1, ..., n\}$ to be the largest index such that $\sigma_i \in \{0, 1\}$. Since we defined $\sigma_0 = 1$ above, such an $i$ exists. Then we can write

$$cost_L(S, \mathcal{P}(\sigma), T) = \mu + (|P_i^{\sigma_i}| + 1)|T \setminus W_i(\sigma)|, \text{ and}$$
$$cost_L(S, \mathcal{P}(\sigma), Z, T) = \mu + (|P_i^{\sigma_i}| + 1)|Z \setminus W_i(\sigma)| + (|Z| + 1)|T \setminus (Z \cup W_i(\sigma))|,$$

where $\mu = cost_L(S, \mathcal{P}(\sigma))$ denotes the sum contributed to $cost_L$ by the common initial sequence.

Since $T \cap W_i(\sigma) = T \cap Z = \emptyset$, we have $|T \setminus W_i(\sigma)| = |T \setminus (Z \cup W_i(\sigma))| = \tau$. Since $Z \setminus W_i(\sigma) \subseteq M \cup A_{i+1} \cup ... \cup A_{n+1}$, we have $|Z \setminus W_i(\sigma)| \le (n+1-i)\alpha + 8m$. Thus we can write

$$
\begin{aligned}
cost_L(S, \mathcal{P}(\sigma), T) &- cost_L(S, \mathcal{P}(\sigma), Z, T) \\
&\ge (|P_i^{\sigma_i}| - |Z|)\tau - (|P_i^{\sigma_i}| + 1)((n+1-i)\alpha + 8m) \\
&= ((n-i+1)(\beta - \alpha) + \delta_i - m)\tau \\
&\quad - ((n-i+1)\beta + i\alpha + \delta_i + 1)((n-i+1)\alpha + 8m),
\end{aligned}
$$

where the last equality follows by (ii) and (iii). Since $(n-i+1)(\beta - \alpha) + \delta_i - m \ge \beta - \alpha - m$, $(n-i+1)\beta + i\alpha + \delta_i + 1 \le (n+1)\beta + 17$, and $(n-i+1)\alpha + 8m \le (n+1)\alpha + 8m$, our claim, that is, the positivity of the above difference, is implied by our assumption (5.3).                                                                         □

For $\sigma \in \{0, 1, *\}^{[n]}$ with $\sigma_j = *$, let us denote by $\sigma^{j \to 0}$ and $\sigma^{j \to 1}$ the sequences obtained by switching the $j$th entry in $\sigma$ to 0 and 1, respectively. Next we show that, with a right choice of parameters, an optimal solution must include exactly one of $X_i$ and $Y_i$ for all $i = 1, \ldots, n$, thus proving statement (II).

LEMMA 5.2. *For every $\sigma \in \{0, 1, *\}^{[n]}$ with $\sigma_j = *$, and for every $\epsilon \in \{0, 1\}$ we have*

$$cost_L(S, \mathcal{P}(\sigma), Z, T) > cost_L(S, \mathcal{P}(\sigma^{j \to \epsilon}), Z, T),$$

*if the following inequality holds:*

(5.4) $$(\beta - \alpha) \cdot (\alpha + 8m) \; > \; 16m(n\beta + 17).$$

*Proof.* Similarly as before, let us define $\sigma_0 = \sigma_{n+1} = 1$, $P_0^{\sigma_0} = S$ and $P_{n+1}^{\sigma_{n+1}} = Z$. Let us set $i$ to be the largest index $i < j$ with $\sigma_i \neq *$ while $k$ to be the smallest index $j < k$ with $\sigma_i \neq *$. As $\sigma_0 = \sigma_{n+1} = 1$, such $i$ and $k$ exist, and $i < j < k$.

Let us introduce the notation $I = P_i^{\sigma_i}$, $J = P_j^{\epsilon}$ and $K = P_k^{\sigma_k}$. Let us further denote by $\mathcal{Q}$ the initial and by $\mathcal{R}$ the terminating subsequence of $\mathcal{P}(\sigma)$ such that $\mathcal{P}(\sigma) = (\mathcal{Q}, I, \mathcal{R})$. Finally, set $U = W_i(\sigma)$, $\mu = cost_L(S, \mathcal{Q})$, $\nu = cost_L(S, \mathcal{P}(\sigma), Z, T) - cost_L(S, \mathcal{Q}, I, K)$, and $\nu' = cost_L(S, \mathcal{P}(\sigma^{j \to \epsilon}), Z, T) - cost_L(S, \mathcal{Q}, I, J, K)$.

Note that by the definition of $cost_L$, the expression for $\nu$ and $\nu'$ are almost the same. The only difference is the sum defining $\nu'$ has $J$ added to the unions which are taken away in each term and thus the corresponding cardinalities cannot increase. We thus have that $\nu' \leq \nu$.

Note also that with the above notation we can write

$$\begin{aligned} cost_L(S, \mathcal{P}(\sigma), Z, T) &= \mu + (|I| + 1)|K \setminus U| + \nu, \text{ and} \\ cost_L(S, \mathcal{P}(\sigma^{j \to \epsilon}), Z, T) &= \mu + (|I| + 1)|J \setminus U| + (|J| + 1)|K \setminus (U \cup J)| + \nu'. \end{aligned}$$

Thus, the difference between the above two left hand sides is at least $(|I| + 1)(|K \setminus U| - |J \setminus U|) - (|J| + 1)|K \setminus (U \cup J)|$. By our definitions of these sets we have

$$\begin{aligned} K \setminus U &\supseteq A_{i+1} \cup \cdots \cup A_k, \\ J \setminus U &\subseteq M \cup A_{i+1} \cup \cdots \cup A_j, \text{ and} \\ K \setminus (U \cup J) &\subseteq M \cup A_{j+1} \cup \cdots \cup A_k. \end{aligned}$$

Hence, using our notation and (iii) we get

$$\begin{aligned} cost_L(S, \mathcal{P}(\sigma), Z, T) \; &- \; cost_L(S, \mathcal{P}(\sigma^{j \to \epsilon}), Z, T) \\ &\geq \; ((n - i + 1)\beta + i\alpha + \delta_i + 1)((k - i)\alpha - (j - i)\alpha - 8m) \\ &\quad - \; ((n - j + 1)\beta + j\alpha + \delta_j + 1)(8m + (k - j)\alpha) \\ \\ &= \; (k - j)\alpha((j - i)(\beta - \alpha) + \delta_j - \delta_i) \\ &\quad - \; 8m((2n - i - j + 2)\beta + (i + j)\alpha + \delta_i + \delta_j + 2). \end{aligned}$$

Using (i), $j - i \geq 1$, $k - j \geq 1$, and $i + j \geq 1$ we can conclude that

$$\begin{aligned} cost_L(S, \mathcal{P}(\sigma), Z, T) \; &- \; cost_L(S, \mathcal{P}(\sigma^{j \to \epsilon}), Z, T) \\ &\geq \; (\beta - \alpha)(\alpha + 8m) - 16m(n\beta + 17), \end{aligned}$$

where the right hand side is positive by our assumption (5.4), hence completing our proof. □

It is easy to see that we can choose $\alpha$, $\beta$, and $\tau$ such that (5.3) and (5.4) hold, and all of these parameters are $O(m^2 n^3)$. Indeed, set $\alpha := 16mn$. Then (5.4) simplifies to $(\beta - 16mn) \cdot (16mn + 8m) > 16m(n\beta + 17)$, which holds if we set $\beta := 32mn^2 + 16mn + 35$. Now (5.3) transforms into $(32mn^2 - m + 35) \cdot \tau > ((16mn + 1)(32mn^2 + 16mn + 35) + 17) \cdot ((n + 1)16mn + 8m)$, therefore setting $\tau := \lceil ((16mn + 1)(32mn^2 + 16mn + 35) + 17) \cdot ((n + 1)16mn + 8m)/(32mn^2 - m + 35) + 1 \rceil$ gives a proper choice of the parameters. In this way our construction above has polynomial size in the size of $\Phi$. Let us assume for the rest of our proof that we fix such a choice for $\alpha$, $\beta$ and $\tau$.

In what follows we show that $price_L(S, T)$ is the smallest if and only if $\Phi$ is satisfiable.

LEMMA 5.3. *There exists a function $d : [n] \to \mathbb{Z}_+$ such that*

$$|X_{i+1} \setminus W_i(\sigma)| = |Y_{i+1} \setminus W_i(\sigma)| = d(i).$$

*for every $i = 0, \ldots, n$ and $\sigma \in \{0,1\}^{[n]}$.*

*Proof.* To see the claim, let us consider a clause $C$ of $\Phi$ that contains variable $x_{i+1}$ or its negation. Recall that $\mathcal{C}(C) \subseteq M$ denotes the set of eight possible clauses included in $M$ consisting of the three variables in C. Let us further denote by $I(C) = \{u, v, w\}$ the indices $u < v < w$ of the variables that are involved (with or without a complementation) in C. Let us then observe that if $i + 1 = u$ is the smallest index in $I(C)$, then both $X_{i+1} \setminus W_i(\sigma)$ and $Y_{i+1} \setminus W_i(\sigma)$ contain exactly four elements of $\mathcal{C}(C)$. This is because $W_i(\sigma)$ contains clauses from $M$ that contain variables $x_j$ or its negation, depending on $\delta_j$, for $j \leq i$. Thus none of the eight clauses of $\mathcal{C}(C)$ are contained in $W_i(\sigma)$, and exactly four of those contain $x_{i+1}$ and four contain its negation. If $i + 1 = v$ is the second smallest index in $I(C)$, then both $X_{i+1} \setminus W_i(\sigma)$ and $Y_{i+1} \setminus W_i(\sigma)$ contain exactly two elements of $\mathcal{C}(C)$. This is because $\mathcal{C}(C) \setminus W_i(\sigma)$ contains now exactly the four clauses that contain either $x_v$ or its negation, depending on $\sigma_v$, and thus two of those four contain $x_{i+1}$ and two contain $\overline{x}_{i+1}$. Finally, if $i + 1 = w$ is the largest index in $I(C)$, then both $X_{i+1} \setminus W_i(\sigma)$ and $Y_{i+1} \setminus W_i(\sigma)$ contain exactly 1 element of $\mathcal{C}(C)$. This is because in this case $\mathcal{C}(C) \setminus W_i(\sigma)$ contains only the two clauses that do not contain the particular combination of $x_u$ or its negation and $x_v$ or its negation that corresponds to $(\sigma_u, \sigma_v)$, and of those two one contains $x_{i+1}$ and one contains its negation.

Note that these counts do not depend on $\sigma \in \{0,1\}^{[n]}$, and hence the claim follows. □

LEMMA 5.4. *There exists an integer $g \in \mathbb{Z}_+$ such that*

$$cost_L(S, \mathcal{P}(\sigma)) = g$$

*for every $\sigma \in \{0,1\}^{[n]}$.*

*Proof.* The claim follows by Lemma 5.3 and the fact that $|X_i| = |Y_i|$ for $i = 1, \ldots, n$. □

LEMMA 5.5. *There exists an integer $f$ such that for all $\sigma \in \{0,1\}^{[n]}$ we have*

$$cost_L(S, \mathcal{P}(\sigma), X_{n+1}) = f.$$

*Proof.* By (iii) we have $|P_n^{\sigma_n}| = |X_n| = |Y_n| = \beta + n\alpha + \delta_n$, and by our construction we have $X_{n+1} \setminus W_n(\sigma) = A_{n+1}$. Thus, by Lemma 5.4 we get $f = g + (\beta + n\alpha + \delta_n + 1)|A_{n+1}| = g + (\beta + n\alpha + \delta_n + 1)\alpha$ and the statement follows. □

LEMMA 5.6. *For $\sigma \in \{0,1\}^{[n]}$ we have*

$$cost_L(S, \mathcal{P}(\sigma), Z, T) = f + (\beta + n\alpha + \delta_n + 1) \cdot |\Phi(\sigma)| + (|Z| + 1) \cdot |T|,$$

*where $|\Phi(\sigma)|$ denotes the number of clauses of $\Phi$ that are not satisfied by $\sigma$.*

*Proof.* The lemma follows by the construction and by Lemma 5.5. □

LEMMA 5.7. *For $\mathcal{B}_\Phi$ defined in (5.2) we have*

$$price_L(S, T) = f + (|Z| + 1) \cdot |T|$$

*if and only if $\Phi$ is satisfiable.*

*Proof.* The construction of $\Phi^{(1)}$ in Section 4.2 shows that there exists a pure Horn CNF attaining the minimum in $price_L(S, T)$ that can be written in form (5.1) for some sequence $\{S_0, \ldots, S_s\} \subseteq \mathcal{B}_\Phi$ where $|S_0| > |S_1| > \ldots > |S_s|$. By Lemmas 5.1 and 5.2, we may assume that $\mathcal{S} = \{S, \mathcal{P}(\sigma), Z, T\}$ for some truth assignment $\sigma \in \{0, 1\}^{[n]}$. Lemma 5.6 implies that $price_L(S, T) = cost_L(S, \mathcal{P}(\sigma), Z, T) = f + (|Z| + 1) \cdot |T|$ if and only if $|\Phi(\sigma)| = 0$, that is, if $\sigma$ is a true point of $\Phi$.    □

THEOREM 5.8. *Computing $price_L$ is NP-hard.*

*Proof.* Let $\Phi$ be a 3-CNF in which every variable appears at most 4 times. Recall that SAT is NP-complete even when restricted to this class of CNF formulas [26]. By Lemma 5.7, $\Phi$ is satisfiable if and only if $price_L(S, T) = f + (|Z| + 1) \cdot |T|$ that is if and only if there exists a $\sigma \in \{0, 1\}^{[n]}$ such that $|\Phi(\sigma)| = 0$. This shows that computing $price_L$ is NP-hard.    □

**6. Clause minimization and strongly connected subgraphs.** For a strongly connected graph $D = (V, E)$ and non-negative weights $w : E \to \mathbb{Z}_+$, we denote by $\mathtt{MWSCS}(D, w)$ the problem of finding a minimum weight subset $F \subseteq E$ of the arcs such that $(V, F)$ is also strongly connected. We denote by $\mathtt{mwscs}(D, w) = w(F)$ the weight of such a minimum weight arc subset. $\mathtt{MWSCS}$ is an NP-hard problem, for which polynomial time approximation algorithms are known. For the case of uniform weights a 1.61-approximation was given by Khuller et al. [21]. For general weights a simple 2-approximation is due to Fredericson and Jájá [17]. Note that in the case of general weights, we can assume that $D$ is a complete directed graph.

As already observed in the beginning of Section 4, there is a natural relation of the above problem to the minimization of a key Horn function. Let us consider a Sperner hypergraph $\mathcal{B} \subseteq 2^V \setminus \{V\}$ and the corresponding Horn function

$$(6.1) \qquad h_\mathcal{B} \;=\; \bigwedge_{B \in \mathcal{B}} B \to (V \setminus B).$$

The body graph of $\mathcal{B}$ was a complete directed graph $D_\mathcal{B}$ where $V(D_\mathcal{B}) = \mathcal{B}$. Define a weight function $w$ on the arcs of this graph by setting $w(B, B') = price_*(B, B')$ for all $B, B' \in \mathcal{B}$, $B \neq B'$, where $price_*$ is defined in (3.2). Then any solution $H \subseteq E(D_\mathcal{B}) = \mathcal{B} \times \mathcal{B}$ of problem $\mathtt{MWSCS}(D_\mathcal{B}, w)$ defines a representation of $h_\mathcal{B}$:

$$(6.2) \qquad \Phi(H) \;=\; \bigwedge_{(B, B') \in H} \Phi_*(B, B'),$$

where $\Phi_*(B, B')$ is a formula for which $B' \subseteq F_{\Phi_*(B, B')}(B)$, $\mathcal{B}_{\Phi_*(B, B')} \subseteq \mathcal{B}$ and $|\Phi_*(B, B')|_* = price_*(B, B')$. It is immediate to see that $OPT_*(\mathcal{B}) \leq w(H)$ holds. Thus, it is natural to expect that a polynomial time approximation of problem $\mathtt{MWSCS}(D_\mathcal{B}, w)$ provides also a good approximation for $OPT_*(\mathcal{B})$. This however turns out to be false for the case of $* = C$. Our construction uses finite projective spaces $PG(d, q)$ where $d$ is the dimension and $q$ is the order.

THEOREM 6.1. *Let $d \geq 4$ be a positive integer, $n$ be the number of points of $PG(d, 2)$ and $V = \mathbb{Z}_n$. Then we have*

$$(6.3) \qquad \max_{\mathcal{B} \subseteq 2^V \setminus \{V\}} \frac{\mathtt{mwscs}(D_\mathcal{B}, price_C)}{OPT_C(\mathcal{B})} \geq \frac{n + 1}{8}.$$

Before proving the theorem, let us recall first some basic facts on finite projective spaces from the book [14]. The finite projective space $PG(d, q)$ of dimension $d$ over

a finite field $GF(q)$ of order $q$ (prime power) has $n = q^d + q^{d-1} + \cdots + q + 1$ points.
Subspaces of dimension $k$ are isomorphic to $PG(k, q)$ for $0 \leq k < d$, where 0-dimension
subspaces are the points themselves. The number of subspaces of dimension $k < d$ is

$$N_k(d, q) \;=\; \prod_{i=0}^{k} \frac{q^{d+1-i} - 1}{q^{i+1} - 1},$$

and the number of points of such a subspace is $q^k + q^{k-1} + \cdots + q + 1$. In particular,
the number of subspaces of dimension $d - 1$ is $N_{d-1}(d, q) = n$. If $F$ and $F'$ are two
distinct subspaces of dimension $k$, then

$$2k - d \leq dim(F \cap F') \leq k - 1.$$

Furthermore, any $k + 1$ points belong to at least one subspace of dimension $k$.

Let us also recall that $PG(d, q)$ has a cyclic automorphism. In other words the
points of $PG(d, q)$ can be identified with the integers of the cyclic group $\mathbb{Z}_n$ of modulo
$n$ addition such that if $F \subseteq \mathbb{Z}_n$ is a subspace of dimension $k$, then $F + i = \{f + i$
$\mod n \mid f \in F\}$ is also a subspace of dimension $k$. Furthermore, two subspaces $F$ and
$F + i$ are distinct if $i \neq 0 \pmod n$.

Let us consider a subspace $Q \subseteq \mathbb{Z}_n$ of dimension $d - 1$. Then the family defined
as $\mathcal{Q} = \{Q + i \mid i \in \mathbb{Z}_n\}$ contains all subspaces of $PG(d, q)$ of dimension $d - 1$ and
the size of $Q$ is $n$. In the rest of this section we use $+$ for the modulo $n$ addition of
integers.

LEMMA 6.2. *For every $k = 0, ..., d-1$ there exists a unique subspace of dimension
$k$ that contains $\{0, 1, ..., k\}$.*

*Proof.* By the properties we recalled above it follows that there is at least one
such subspace for every $0 \leq k < d$. We prove that there is at most one by induction
on $k$. For $k = 0$ this is obvious, since the points are the only subspaces of dimension
0. Assume next that the claim is already proved for all $k' < k$, and assume that there
are two distinct subspaces, $R$ and $R'$, of dimension $k$ both of which contains the set
$\{0, 1, ..., k\}$. Then $R \cap R'$ and $(R - 1) \cap (R' - 1) = (R \cap R') - 1$ are two distinct
subspaces of dimension $k' < k$ and both contain $\{0, 1, ..., k - 1\}$, contradicting our
inductive assumption, and thus proving our claim.                                   □

Thus, by Lemma 6.2, there exists a unique subspace $Q \subseteq \mathbb{Z}_n$ of dimension $d - 1$
that contains $\{0, 1, ..., d - 1\}$.

LEMMA 6.3. $d \notin Q$.

*Proof.* Assume to the contrary that $d \in Q$. Then the set $\{0, 1, ..., d - 1\}$ is
contained by both $Q$ and $Q - 1 = Q + (n - 1)$, contradicting Lemma 6.2, since $Q$ and
$Q - 1$ are distinct subspaces of dimension $d - 1$.                                 □

Let us also introduce the set $D = \{0, 1, ..., d\}$. Now we are in the position to prove
the theorem.

*Proof of Theorem 6.1.* Let us define $\mathcal{B} := \mathcal{Q} \cup \{D + i \mid i \in \mathbb{Z}_n\}$, and observe that
for any distinct pair $B \in \mathcal{Q}$ and $B' \in \mathcal{B}$ we have $|B \setminus B'| \geq 2^{d-1}$. This is obvious if
$B' \in \mathcal{Q}$ by properties of subspaces, and follows easily for $B' \in \mathcal{B} \setminus \mathcal{Q}$ because $d$ is at
least four. Since in any solution $H \subseteq \mathcal{B} \times \mathcal{B}$ we must have an arc entering $B$ for all
$B \in \mathcal{Q}$, and for each such arc $(B', B) \in H$ the CNF $\Phi_C(B', B)$ must contain a clause
with head $x$ for each $x \in B \setminus B'$, we get

(6.4)                        $\mathtt{mwscs}(D_\mathcal{B}, price_C) \;\geq\; n \cdot 2^{d-1}.$

On the other hand, we have that

$$(6.5) \qquad \Phi \; = \; \left( \bigwedge_{i \in \mathbb{Z}_n} (Q+i) \to (d+i) \right) \wedge \left( \bigwedge_{i \in \mathbb{Z}_n} (D+i) \to (d+1+i) \right)$$

is a representation of $h_{\mathcal{B}}$ and $|\Phi|_C \leq 2n$. As $n = 2^d + \cdots + 2 + 1 = 2^{d+1} - 1$, we have $2^{d-1} = (n+1)/4$. Thus

$$(6.6) \qquad \mathtt{mwscs}(D_{\mathcal{B}}, price_C) \; \geq \; \frac{n+1}{8} \cdot OPT_C(\mathcal{B}),$$

completing the proof of the theorem. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \square$

Let us note that for such a negative result, we need to rely on Horn functions with large bodies. For the case when we limit the body sizes of the underlying Horn function by a constant $k$, we have already showed that there exists a solution which is a $k$-approximation for the CNF minimization problem as well as for the MWSCS problem, see Lemma 4.1.

**7. Conclusions.** In this paper we study the class of key Horn functions which is a generalization of a well-studied class of hydra functions [22, 25]. Given a CNF representing a key Horn function, we are interested in finding the minimum size logically equivalent pure Horn CNF, where the size of the output CNF is measured in several different ways. This problem is known to be NP-hard already for hydra CNFs for most common measures of the CNF size.

The main results of the paper are two approximation algorithms for key Horn CNFs – one for minimizing the number of clauses and the other for minimizing the total number of literals in the output CNF. Both algorithms achieve a logarithmic approximation bound with respect to the size of the largest body in the input CNF (denoted by $k$). This parameter can be also defined as the size of the largest clause in the input CNF minus one. Note that $k$ is a trivial lower bound on the number of variables (denoted by $n$).

These algorithms are (to the best of our knowledge) the first approximation algorithms for NP-hard Horn minimization problems that guarantee a sublinear approximation bound with respect to $k$. It follows that both algorithms also guarantee a sublinear approximation bound with respect to $n$. There are two approximation algorithms for Horn minimization known in the literature, one for general Horn CNFs [19], and one for hydra CNFs [25], but both of them guarantee only a linear (or higher) approximation bound with respect to $k$ (see Table 1 and the relevant text in the introduction section for details).

For a given pair of sets $S, T$ and set of bodies $\mathcal{B}$, we prove NP-hardness of the problem of finding a literal minimum pure Horn CNF $\Phi$ that uses bodies only from $\mathcal{B}$ and for which the forward chaining procedure starting from $S$ reaches all the variables in $T$.

In contrast to our approach which takes an in-branching in the body graph and extends it with a small number of additional edges, we show that no polynomial time approximation of the minimum weight strongly connected subgraph problem in the body graph may provide a good solution for the CNF minimization problem. The counterexample is based on a construction using finite projective spaces.

Our analysis of Procedure 1 provides an approximation factor of $\min\{\lceil \log n \rceil + 1, \lceil \log k \rceil + 2\}$ for (C) and (BC). However, we do not know whether our analysis provides the best bound in general. We actually believe that the proposed algorithm

(possibly with slight modifications) could be used to obtain a constant factor approximation for (C) and (BC). Similarly, no example is known for which the solution provided by Procedure 2 attains the proved approximation bound tightly. A better analysis of these procedures possibly leading to a constant factor approximation or a better lower bound than the one given in Lemma 3.6 is subject of future research.

## REFERENCES

[1] A. V. AHO, M. R. GAREY, AND J. D. ULLMAN, *The transitive reduction of a directed graph*, SIAM Journal on Computing, 1 (1972), pp. 131–137.

[2] G. AUSIELLO, A. D'ATRI, AND D. SACCA, *Minimal representation of directed hypergraphs*, SIAM Journal on Computing, 15 (1986), pp. 418–431.

[3] A. BHATTACHARYA, B. DASGUPTA, D. MUBAYI, AND G. TURÁN, *On approximate Horn formula minimization*, in International Colloquium on Automata, Languages, and Programming (ICALP), Springer, 2010, pp. 438–450.

[4] E. BOROS, O. ČEPEK, AND A. KOGAN, *Horn minimization by iterative decomposition*, Annals of Mathematics and Artificial Intelligence, 23 (1998), pp. 321–343.

[5] E. BOROS, O. ČEPEK, A. KOGAN, AND P. KUČERA, *A subclass of Horn CNFs optimally compressible in polynomial time*, Annals of Mathematics and Artificial Intelligence, 57 (2009), pp. 249–291.

[6] E. BOROS, O. ČEPEK, AND P. KUČERA, *A decomposition method for CNF minimality proofs*, Theoretical Computer Science, 510 (2013), pp. 111–126.

[7] E. BOROS AND A. GRUBER, *Hardness results for approximate pure Horn CNF formulae minimization*, Annals of Mathematics and Artificial Intelligence, 71 (2014), pp. 327–363.

[8] D. BUCHFUHRER AND C. UMANS, *The complexity of Boolean formula minimization*, Journal of Computer and System Sciences, 77 (2011), pp. 142–153.

[9] N. CASPARD AND B. MONJARDET, *The lattices of closure systems, closure operators, and implicational systems on a finite set: a survey*, Discrete Applied Mathematics, 127 (2003), pp. 241–269.

[10] Y.-J. CHU, *On the shortest arborescence of a directed graph*, Scientia Sinica, 14 (1965), pp. 1396–1400.

[11] S. A. COOK, *The complexity of theorem-proving procedures*, in Proceedings of the 3rd annual ACM Symposium on Theory of Computing, ACM, 1971, pp. 151–158.

[12] Y. CRAMA AND P. L. HAMMER, *Boolean functions: Theory, algorithms, and applications*, Cambridge University Press, 2011.

[13] R. DECHTER AND J. PEARL, *Structure identification in relational data*, Artificial Intelligence, 58 (1992), pp. 237 – 270.

[14] P. DEMBOWSKI, *Finite geometries*, Ergebnisse der Mathematik und ihrer Grenzgebiete, Springer-Verlag, 1968.

[15] W. F. DOWLING AND J. H. GALLIER, *Linear-time algorithms for testing the satisfiability of propositional Horn formulae*, The Journal of Logic Programming, 1 (1984), pp. 267 – 284.

[16] J. EDMONDS, *Optimum branchings*, Journal of Research of the National Bureau of Standards, B, 71 (1967), pp. 233–240.

[17] G. N. FREDERICKSON AND J. JÁJÁ, *Approximation algorithms for several graph augmentation problems*, SIAM Journal on Computing, 10 (1981), pp. 270–283.

[18] J.-L. GUIGUES AND V. DUQUENNE, *Familles minimales d'implications informatives résultant d'un tableau de données binaires*, Mathématiques et Sciences humaines, 95 (1986), pp. 5–18.

[19] P. L. HAMMER AND A. KOGAN, *Optimal compression of propositional Horn knowledge bases: complexity and approximation*, Artificial Intelligence, 64 (1993), pp. 131–145.

[20] P. L. HAMMER AND A. KOGAN, *Quasi-acyclic propositional Horn knowledge bases: optimal compression*, IEEE Transactions on Knowledge and Data Engineering, 7 (1995), pp. 751–762.

[21] S. KHULLER, B. RAGHAVACHARI, AND N. YOUNG, *On strongly connected digraphs with bounded cycle length*, Discrete Applied Mathematics, 69 (1996), pp. 281–289.

[22] P. KUČERA, *Hydras: Complexity on general graphs and a subclass of trees*, Theoretical Computer Science, 658 (2017), pp. 399–416.

[23] D. MAIER, *Minimum covers in relational database model*, Journal of the ACM (JACM), 27 (1980), pp. 664–674.

[24] D. MAIER, *The theory of relational databases*, vol. 11, Computer Science Press, Rockville, 1983.

[25] R. H. SLOAN, D. STASI, AND G. TURÁN, *Hydras: Directed hypergraphs and Horn formulas*, Theoretical Computer Science, 658 (2017), pp. 417–428.

[26] C. A. TOVEY, *A simplified NP-complete satisfiability problem*, Discrete Applied Mathematics, 8 (1984), pp. 85–89.

[27] J. D. ULLMAN, *Principles of database systems*, Galgotia Publications, 1984.

[28] C. UMANS, *Hardness of approximating/spl sigma//sub 2//sup p/minimization problems*, in 40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039), IEEE, 1999, pp. 465–474.

[29] C. UMANS, *The minimum equivalent dnf problem and shortest implicants*, Journal of Computer and System Sciences, 63 (2001), pp. 597–611.

[30] C. UMANS, T. VILLA, AND A. L. SANGIOVANNI-VINCENTELLI, *Complexity of two-level logic minimization*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 25 (2006), pp. 1230–1246.