

1 Best-case and Worst-case Sparsifiability of 2 Boolean CSPs

3 **Hubie Chen**

4 Birkbeck, University of London

5 hubie@dcs.bbk.ac.uk

6 **Bart M.P. Jansen**¹

7 Eindhoven University of Technology

8 b.m.p.jansen@tue.nl

9  <http://orcid.org/0000-0001-8204-1268>

10 **Astrid Pieterse**²

11 Eindhoven University of Technology

12 a.pieterse@tue.nl

13  <http://orcid.org/0000-0003-3721-6721>

14 — Abstract —

15 We continue the investigation of polynomial-time sparsification for NP-complete Boolean Con-
16 straint Satisfaction Problems (CSPs). The goal in sparsification is to reduce the number of
17 constraints in a problem instance without changing the answer, such that a bound on the num-
18 ber of resulting constraints can be given in terms of the number of variables n . We investigate
19 how the worst-case sparsification size depends on the types of constraints allowed in the problem
20 formulation (the constraint language). Two algorithmic results are presented. The first result
21 essentially shows that for any arity k , the only constraint type (up to negations) for which no
22 nontrivial sparsification is possible, is the logical OR. Our second result concerns linear sparsi-
23 fication, i.e., a reduction to an equivalent instance with $\mathcal{O}(n)$ constraints. Using linear algebra
24 over rings of integers modulo prime powers, we give an elegant necessary and sufficient condi-
25 tion for a constraint type to be captured by a degree-1 polynomial over such a structure, which
26 yields linear sparsifications. The combination of these algorithmic results allows us to prove two
27 characterizations that capture the optimal sparsification sizes for a range of Boolean CSPs. For
28 NP-complete Boolean CSPs whose constraints are *symmetric* (the satisfaction depends only on
29 the number of 1 values in the assignment, not on their positions), we give a complete charac-
30 terization of which constraint languages allow for a linear sparsification. For Boolean CSPs in
31 which every constraint has arity at most three, we characterize the optimal size of sparsifications
32 in terms of the largest OR that can be expressed by the constraint language.

33 **2012 ACM Subject Classification** Computing methodologies → Symbolic and algebraic algo-
34 rithms, Theory of computation → Problems, reductions and completeness, Theory of computa-
35 tion → Parameterized complexity and exact algorithms

36 **Keywords and phrases** constraint satisfaction problems; kernelization; sparsification; lower bounds

37 **Digital Object Identifier** 10.4230/LIPIcs.IPEC.2018.23

38 **Acknowledgements** We would like to thank Emil Jeřábek for the proof of Lemma 4.2.

¹ Supported by NWO Gravitation grant “Networks”.

² Supported by NWO Gravitation grant “Networks”.



39 **1** Introduction40 **Background**

41 The framework of constraint satisfaction problems (CSPs) provides a unified way to study
 42 the computational complexity of a wide variety of combinatorial problems such as CNF-
 43 SATISFIABILITY, GRAPH COLORING, and NOT-ALL-EQUAL SAT. The framework uncovers
 44 algorithmic approaches that simultaneously apply to several problems, and also identifies
 45 common sources of intractability. For the purposes of this discussion, a CSP is specified using
 46 a *constraint language*, which is a set of relations; the problem is to decide the satisfiability of
 47 a set of constraints, where each constraint has a relation coming from the constraint language.
 48 The fact that many problems can be viewed as CSPs motivates the following investigation:
 49 how does the complexity of a CSP depend its constraint language? A key result in this area
 50 is Schaefer's dichotomy theorem [20], which classifies each CSP over the Boolean domain as
 51 polynomial-time solvable or NP-complete.

52 Continuing a recent line of investigation [12, 14, 17], we aim to understand for which
 53 NP-complete CSPs an instance can be *sparsified* in polynomial time, without changing the
 54 answer. In particular, we investigate the following questions. Can the number of constraints
 55 be reduced to a small function of the number of variables n ? How does the sparsifiability of a
 56 CSP depend on its constraint language? We utilize the framework of kernelization [5, 8, 18],
 57 originating in parameterized complexity theory, to answer such questions.

58 The first results concerning polynomial-time sparsification in terms of the number n
 59 of variables or vertices were mainly negative. Under the assumption that $\text{NP} \not\subseteq \text{coNP/poly}$
 60 (which we tacitly assume throughout this introduction), Dell and van Melkebeek [7] proved a
 61 strong lower bound: For an integer $d \geq 3$ and positive real ε , there cannot be a polynomial-
 62 time algorithm that compresses any instance φ of d -CNF-SAT on n variables, into an
 63 equivalent SAT instance φ' of bitsize $\mathcal{O}(n^{d-\varepsilon})$. In fact, there cannot even be an algorithm
 64 that transforms such φ into equivalent small instances ψ of an *arbitrary* decision problem.
 65 Since an instance of d -CNF-SAT has at most $2^d n^d \in \mathcal{O}(n^d)$ distinct clauses, it can *trivially*
 66 be sparsified to $\mathcal{O}(n^d)$ clauses by removing duplicates, and can be compressed to size $\mathcal{O}(n^d)$
 67 by storing it as a bitstring indicating for each possible clause whether or not it is present.
 68 The cited lower bound therefore shows that the trivial sparsification for d -CNF-SAT cannot
 69 be significantly improved; we say that the problem does not admit nontrivial (polynomial-
 70 time) sparsification. Following these lower bounds for SAT, a number of other results were
 71 published [6, 11, 16] proving other problems do not admit nontrivial sparsification either.

72 This pessimistic state of affairs concerning nontrivial sparsification algorithms changed
 73 several years ago, when a subset of the authors [14] showed that the d -NOT-ALL-EQUAL SAT
 74 problem *does* have a nontrivial sparsification. In this problem, clauses have size at most d
 75 and are satisfied if the literals do not all evaluate to the same value. While there can be $\Omega(n^d)$
 76 different clauses in an instance, there is an efficient algorithm that finds a subset of $\mathcal{O}(n^{d-1})$
 77 constraints that preserves the answer, resulting in a compression of bitsize $\mathcal{O}(n^{d-1} \log n)$. The
 78 first proof of this result was based on an ad-hoc application of a theorem of Lovász [19]. Later,
 79 the underlying proof technique was extracted and applied to a wider range of problems [12].
 80 This led to the following understanding: if each relation in the constraint language can be
 81 represented by a polynomial of degree at most d , in a certain technical sense, then this
 82 allows the number of constraints in an n -variable instance of such a CSP to be reduced
 83 to $\mathcal{O}(n^d)$. The sparsification for d -NOT-ALL-EQUAL SAT is then explained by noting that
 84 such constraints can be captured by polynomials of degree $d - 1$. It is therefore apparent
 85 that finding a low-degree polynomial to capture the constraints of a CSP is a powerful tool

86 to obtain sparsification algorithms for it. Finding such polynomials of a certain degree d , or
 87 determining that they do not exist, proved a challenging and time-intensive task (cf. [13]).

88 The polynomial-based framework [12] also resulted in some *linear* sparsifications. Since
 89 “1-in- d ” constraints (to satisfy a clause, *exactly one* out of its $\leq d$ literals should evaluate to
 90 true) can be captured by *linear* polynomials, the 1-IN- d -SAT problem has a sparsification
 91 with $\mathcal{O}(n)$ constraints for each constant d . This prompted a detailed investigation into
 92 linear sparsifications for CSPs by Lagerkvist and Wahlström [17], who used the toolkit of
 93 universal algebra in an attempt to obtain a characterization of the Boolean CSPs with a
 94 linear sparsification. Their results give a necessary and sufficient condition on the constraint
 95 language of a CSP for having a so-called Maltsev embedding over an infinite domain. They
 96 also show that when a CSP has a Maltsev embedding over a *finite* domain, then this
 97 can be used to obtain a linear sparsification. Alas, it remains unclear whether Maltsev
 98 embeddings over infinite domains can be exploited algorithmically, and a characterization of
 99 the linearly-sparsifiable CSPs is currently not known.

100 Our contributions

101 We analyze and demonstrate the power of the polynomial-based framework for sparsifying
 102 CSPs using universal algebra, linear algebra over rings, and relational analysis. We present
 103 two new algorithmic results. These allow us to characterize the sparsifiability of Boolean
 104 CSPs in two settings, wherein we show that the polynomial-based framework yields *optimal*
 105 sparsifications. In comparison to previous work [12], our results are much more fine-grained
 106 and based on a deeper understanding of the reasons why a certain CSP *cannot* be captured
 107 by low-degree polynomials.

108 **Algorithmic results** Our first result (Section 3) shows that, contrary to the pessimistic
 109 picture that arose during the initial investigation of sparsifiability, the phenomenon of
 110 nontrivial sparsification is widespread and occurs for almost all Boolean CSPs! We prove
 111 that if Γ is a constraint language whose largest constraint has arity k , then the *only* reason
 112 that $\text{CSP}(\Gamma)$ does not have a nontrivial sparsification, is that it contains an arity- k relation
 113 that is essentially the k -ary OR (up to negating variables). When $R \subseteq \{0, 1\}^k$ is a relation
 114 with $|\{0, 1\}^k \setminus R| \neq 1$ (the number of assignments that fail to satisfy the constraints is not
 115 equal to 1), then it can be captured by a polynomial of degree $k - 1$. This yields a nontrivial
 116 sparsification compared to the $\Omega(n^k)$ distinct applications of this constraint that can be in
 117 such an instance.

118 Our second algorithmic result (Section 4) concerns the power of the polynomial-based
 119 framework for obtaining linear sparsifications. We give a necessary and sufficient condition
 120 for a relation to be captured by a degree-1 polynomial. Say that a Boolean relation $R \subseteq$
 121 $\{0, 1\}^k$ is *balanced* if there is *no* sequence of vectors $s_1, \dots, s_{2n}, s_{2n+1} \in R$ for $n \geq 1$ such
 122 that $s_1 - s_2 + s_3 \dots - s_{2n} + s_{2n+1} = u \in \{0, 1\}^k \setminus R$. (The same vector may appear
 123 multiple times in this sum.) In other words: R is balanced if one cannot find an odd-length
 124 sequence of vectors in R for which alternating between adding and subtracting these vectors
 125 component-wise results in a 0/1-bitvector u that is outside R . For example, the binary OR
 126 relation $2\text{-OR} = \{0, 1\}^2 \setminus \{(0, 0)\}$ is *not* balanced, since $(0, 1) - (1, 1) + (1, 0) = (0, 0) \notin 2\text{-OR}$,
 127 but the 1-in-3 relation $R_{=1} = \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ is. We prove that if a Boolean
 128 relation R is balanced, then it can efficiently be captured by a degree-1 polynomial and the
 129 number of constraints that are applications of this relation can be reduced to $\mathcal{O}(n)$. Hence
 130 when *all* relations in a constraint language Γ are balanced—we call such a constraint language
 131 *balanced*—then $\text{CSP}(\Gamma)$ has a sparsification with $\mathcal{O}(n)$ constraints. We also show that, on

132 the other hand, if a Boolean relation R is not balanced, then there does not exist a degree-1
 133 polynomial (over any ring $\mathbb{Z}/q\mathbb{Z}$) that captures R in the sense required for application of the
 134 polynomial framework. The property of being balanced is (as defined) a universal-algebraic
 135 property; these results thus tightly bridge universal algebra and the polynomial framework.

136 **Characterizations** The property of being balanced gives an easy way to prove that certain
 137 Boolean CSPs admit linear sparsifications. But perhaps more importantly, this character-
 138 ization constructively exhibits a certain *witness* when a relation can *not* be captured by
 139 a degree-1 polynomial, in the form of the alternating sum of satisfying assignments that
 140 yield an unsatisfying assignment. In several scenarios, we can turn this witness structure
 141 against degree-1 polynomials into a lower bound proving that the problem does not have a
 142 linear sparsification. As a consequence, we can prove two fine-grained characterizations of
 143 sparsification complexity.

144 **Characterization of symmetric CSPs with a linear sparsification (Section 5).**

145 We say that a Boolean relation is *symmetric* if the satisfaction of a constraint only depends
 146 on the *number* of 1-values taken by the variables (the *weight* of the assignment), but does
 147 not depend on the *positions* where these values appear. For example, “1-in- k ”-constraints
 148 are symmetric, just as “not-all-equal”-constraints, but the relation $R_{a \rightarrow b} = \{(0, 1), (1, 1)\}$
 149 corresponding to the truth value of $a \rightarrow b$ is not. We prove that if a symmetric Boolean
 150 relation R is not balanced, then it can implement a binary OR using constants and negations
 151 but without having to introduce fresh variables. Building on this, we prove that if such
 152 an unbalanced symmetric relation R occurs in a constraint language Γ for which $\text{CSP}(\Gamma)$
 153 is NP-complete, then $\text{CSP}(\Gamma)$ does *not* admit a sparsification of size $\mathcal{O}(n^{2-\varepsilon})$ for any
 154 $\varepsilon > 0$. Consequently, we obtain a characterization of the sparsification complexity of NP-
 155 complete Boolean CSPs whose constraint language consists of symmetric relations: there is
 156 a linear sparsification if and only if the constraint language is balanced. This yields linear
 157 sparsifications in several new scenarios that were not known before.

158 **Characterization of sparsification complexity for CSPs of low arity (Section 6).**

159 By combining the linear sparsifications guaranteed by balanced constraint languages with the
 160 nontrivial sparsification when the largest-arity relations do not have exactly one falsifying
 161 assignment, we obtain an exact characterization of the optimal sparsification size for all
 162 Boolean CSPs where each relation has arity at most three. For a Boolean constraint language Γ
 163 consisting of relations of arity at most three, we characterize the sparsification complexity
 164 of Γ as an integer $k \in \{1, 2, 3\}$ that represents the largest OR that Γ can implement using
 165 constants and negations, but without introducing fresh variables. Then we prove that $\text{CSP}(\Gamma)$
 166 has a sparsification of size $\mathcal{O}(n^k)$, but no sparsification of size $\mathcal{O}(n^{k-\varepsilon})$ for any $\varepsilon > 0$, giving
 167 matching upper and lower bounds. Hence for all Boolean CSPs with constraints of arity at
 168 most three, the polynomial-based framework gives provably *optimal* sparsifications in all
 169 cases.

170 **2 Preliminaries**

171 For a positive integer n , define $[n] := \{1, 2, \dots, n\}$. For an integer q , we let $\mathbb{Z}/q\mathbb{Z}$ denote
 172 the integers modulo q . These form a field if q is prime, and a ring otherwise. We will use
 173 $x \equiv_q y$ to denote that x and y are congruent modulo q , and $x \not\equiv_q y$ to denote that they are
 174 incongruent modulo q . For statements marked with a star (\star), the (full) proof can be found
 175 in Appendix A.

176 **Parameterized complexity** A parameterized problem \mathcal{Q} is a subset of $\Sigma^* \times \mathbb{N}$, where Σ is
 177 a finite alphabet. Let $\mathcal{Q}, \mathcal{Q}' \subseteq \Sigma^* \times \mathbb{N}$ be parameterized problems and let $h: \mathbb{N} \rightarrow \mathbb{N}$ be a
 178 computable function. A *generalized kernel for \mathcal{Q} into \mathcal{Q}' of size $h(k)$* is an algorithm that,
 179 on input $(x, k) \in \Sigma^* \times \mathbb{N}$, takes time polynomial in $|x| + k$ and outputs an instance (x', k')
 180 such that: (i) $|x'|$ and k' are bounded by $h(k)$, and (ii) $(x', k') \in \mathcal{Q}'$ if and only if $(x, k) \in \mathcal{Q}$.
 181 The algorithm is a *kernel* for \mathcal{Q} if $\mathcal{Q}' = \mathcal{Q}$.

182 Since a polynomial-time reduction to an equivalent sparse instance yields a generalized
 183 kernel, lower bounds against generalized kernels can be used to prove the non-existence of
 184 such sparsification algorithms. To relate the sparsifiability of different problems to each
 185 other, the following notion is useful.

186 ► **Definition 2.1.** Let $\mathcal{P}, \mathcal{Q} \subseteq \Sigma^* \times \mathbb{N}$ be two parameterized problems. A *linear-parameter*
 187 *transformation* from \mathcal{P} to \mathcal{Q} is a polynomial-time algorithm that, given an instance $(x, k) \in$
 188 $\Sigma^* \times \mathbb{N}$ of \mathcal{P} , outputs an instance $(x', k') \in \Sigma^* \times \mathbb{N}$ of \mathcal{Q} such that the following holds:

- 189 1. $(x, k) \in \mathcal{Q}$ if and only if $(x', k') \in \mathcal{P}$, and
- 190 2. $k' \in \mathcal{O}(k)$.

191 It is well-known (cf. [1, 2]) that the existence of a polynomial-time linear-parameter
 192 transformation from problem \mathcal{P} to \mathcal{Q} implies that any generalized kernelization lower bound
 193 for \mathcal{P} , also holds for \mathcal{Q} .

194 **Operations, relations, and preservation** A *Boolean operation* is a mapping from $\{0, 1\}^k$
 195 to $\{0, 1\}$, where k , a natural number, is said to be the *arity* of the operation; we assume
 196 throughout that operations have positive arity. From here, we define a *partial Boolean*
 197 *operation* in the usual way, that is, it is a mapping from a subset of $\{0, 1\}^k$ to $\{0, 1\}$. We
 198 say that a partial Boolean operation f of arity k is *idempotent* if $f(0, \dots, 0) = 0$ and
 199 $f(1, \dots, 1) = 1$; and, *self-dual* if for all $(a_1, \dots, a_k) \in \{0, 1\}^k$, when $f(a_1, \dots, a_k)$ is defined,
 200 it holds that $f(\neg a_1, \dots, \neg a_k)$ is defined and $f(a_1, \dots, a_k) = \neg f(\neg a_1, \dots, \neg a_k)$.

201 ► **Definition 2.2.** A partial Boolean operation $f: \{0, 1\}^k \rightarrow \{0, 1\}$ is *balanced* if there exist
 202 integer values $\alpha_1, \dots, \alpha_k$, called the *coefficients* of f , such that

- 203 ■ $\sum_i \alpha_i = 1$,
- 204 ■ (x_1, \dots, x_k) is in the domain of f if and only if $\sum_{i \in [k]} \alpha_i x_i \in \{0, 1\}$, and
- 205 ■ $f(x_1, \dots, x_k) = \sum_{i \in [k]} \alpha_i x_i$ for all tuples in its domain.

206 A *relation* over the set D is a subset of D^k ; here, k is a natural number called the *arity*
 207 of the relation. Throughout, we assume that each relation is over a finite set D . A *Boolean*
 208 *relation* is a relation over $\{0, 1\}$.

209 ► **Definition 2.3.** For each $k \geq 1$, we use k -OR to denote the relation $\{0, 1\}^k \setminus \{(0, \dots, 0)\}$.

210 A *constraint language over D* is a finite set of relations over D ; a *Boolean constraint*
 211 *language* is a constraint language over $\{0, 1\}$. For a Boolean constraint language Γ , we define
 212 $\text{CSP}(\Gamma)$ as follows.

CSP(Γ)

Parameter: The number of variables $|V|$.

213 **Input:** A tuple (\mathcal{C}, V) , where \mathcal{C} is a finite set of constraints, V is a finite set of variables,
 and each constraint is a pair $R(x_1, \dots, x_k)$ for $R \in \Gamma$ and $x_1, \dots, x_k \in V$.

Question: Does there exist a *satisfying assignment*, that is, an assignment $f: V \rightarrow \{0, 1\}$
 such that for each constraint $R(x_1, \dots, x_k) \in \mathcal{C}$ it holds that $(f(x_1), \dots, f(x_k)) \in R$?

214 Let $f : \{0, 1\}^k \rightarrow \{0, 1\}$ be a partial Boolean operation, and let $T \subseteq \{0, 1\}^n$ be a Boolean
 215 relation. We say that T is *preserved* by f when, for any tuples $t^1 = (t_1^1, \dots, t_n^1), \dots, t^k =$
 216 $(t_1^k, \dots, t_n^k) \in T$, if all entries of the tuple $(f(t_1^1, \dots, t_1^k), \dots, f(t_n^1, \dots, t_n^k))$ are defined, then
 217 this tuple is in T . We say that a Boolean constraint language Γ is *preserved* by f if each
 218 relation in Γ is preserved by f . We say that a Boolean relation is *balanced* if it is preserved by
 219 all balanced operations, and that a Boolean constraint language is *balanced* if each relation
 220 therein is balanced.

221 Define an *alternating operation* to be a balanced operation $f : \{0, 1\}^k \rightarrow \{0, 1\}$ such
 222 that k is odd and the coefficients alternate between $+1$ and -1 , so that $\alpha_1 = +1, \alpha_2 = -1,$
 223 $\alpha_3 = +1, \dots, \alpha_k = +1$. We have the following.

224 ► **Proposition 2.4 (★).** *A relation is balanced if and only if it is preserved by all alternating*
 225 *operations.*

226 We will use the following straightforwardly verified fact tacitly, throughout.

227 ► **Observation 2.5.** *Each balanced operation is idempotent and self-dual.*

228 For $b \in \{0, 1\}$, let $u_b : \{0, 1\} \rightarrow \{0, 1\}$ be the unary operation defined by $u_b(0) = u_b(1) = b$;
 229 let **major** : $\{0, 1\}^3 \rightarrow \{0, 1\}$ to be the operation defined by **major** $(x, y, z) = (x \wedge y) \vee (x \wedge z) \vee (y \wedge$
 230 $z)$; and, let **minor** : $\{0, 1\}^3 \rightarrow \{0, 1\}$ to be the operation defined by **minor** $(x, y, z) = x \oplus y \oplus z,$
 231 where \oplus denotes exclusive OR. We say that a Boolean constraint language Γ is *tractable* if it
 232 is preserved by one of the six following operations: $u_0, u_1, \wedge, \vee, \text{minor}, \text{major}$; we say that
 233 Γ is *intractable* otherwise. It is known that, in terms of classical complexity, the problem
 234 $\text{CSP}(\Gamma)$ is polynomial-time decidable when Γ is tractable, and that the problem $\text{CSP}(\Gamma)$ is
 235 NP-complete when Γ is intractable (see [4] for a proof; in particular, refer there to the proof
 236 of Theorem 3.21).

237 Constraint Satisfaction and Definability

238 ► **Assumption 2.6.** *By default, we assume in the sequel that the operations, relations, and*
 239 *constraint languages under discussion are Boolean, and that said operations and relations are*
 240 *of positive arity. We nonetheless sometimes describe them as being Boolean, for emphasis.*

241 ► **Definition 2.7.** Let us say that a Boolean relation T of arity m is *cone-definable* from a
 242 Boolean relation U of arity n if there exists a tuple (y_1, \dots, y_n) where:

- 243 ■ for each $j \in [n]$, it holds that y_j is an element of $\{0, 1\} \cup \{x_1, \dots, x_m\} \cup \{\neg x_1, \dots, \neg x_m\}$;
- 244 ■ for each $i \in [m]$, there exists $j \in [n]$ such that $y_j \in \{x_i, \neg x_i\}$; and,
- 245 ■ for each $f : \{x_1, \dots, x_m\} \rightarrow \{0, 1\}$, it holds that $(f(x_1), \dots, f(x_m)) \in T$ if and only
 246 if $(\hat{f}(y_1), \dots, \hat{f}(y_n)) \in U$. Here, \hat{f} denotes the natural extension of f where $\hat{f}(0) = 0,$
 247 $\hat{f}(1) = 1$, and $\hat{f}(\neg x_i) = \neg f(x_i)$.

248 (The prefix *cone* indicates the allowing of **constants** and **negation**.)

249 ► **Example 2.8.** Let $R = \{(0, 0), (0, 1)\}$ and let $S = \{(0, 1), (1, 1)\}$. We have that R is
 250 cone-definable from S via the tuple $(\neg x_2, \neg x_1)$; also, S is cone-definable from R via the same
 251 tuple.

252 When Γ is a constraint language over D , we use Γ^* to denote the expansion of Γ where
 253 each element of D appears as a relation, that is, we define Γ^* as $\Gamma \cup \{\{(d)\} \mid d \in D\}$.

254 The following is a key property of cone-definability; it states that relations that are
 255 cone-definable from a constraint language Γ may be simulated by the constraint language,
 256 and thus used to prove hardness results for $\text{CSP}(\Gamma)$.

257 ▶ **Proposition 2.9 (★)**. *Suppose that Γ is an intractable constraint language, and that Δ*
 258 *is a constraint language such that each relation in Δ is cone-definable from a relation in Γ .*
 259 *Then, there exists a linear-parameter transformation from $\text{CSP}(\Gamma^* \cup \Delta)$ to $\text{CSP}(\Gamma)$.*

260 3 Trivial versus non-trivial sparsification

261 It is well known that k -CNF-SAT allows no non-trivial sparsification, for each $k \geq 3$ [7]. This
 262 means that we cannot efficiently reduce the number of clauses in such a formula to $\mathcal{O}(n^{k-\varepsilon})$.
 263 The k -OR relation is special, in the sense that there is exactly one k -tuple that is not contained
 264 in the relation. We show in this section that when considering k -ary relations for which there
 265 is more than one k -tuple not contained in the relation, a non-trivial sparsification is always
 266 possible. In particular, the number of constraints of any input can efficiently be reduced to
 267 $\mathcal{O}(n^{k-1})$. Using Lemmas 3.4 and 3.6, we can completely classify the constraint languages
 268 that allow a non-trivial sparsification.

269 ▶ **Theorem 3.1 (★)**. *Let Γ be an intractable constraint language. Let k be the maximum*
 270 *arity of any relation $R \in \Gamma$. The following dichotomy holds.*

271 ■ *If for all $R \in \Gamma$ it holds that $|R| \neq 2^k - 1$, then $\text{CSP}(\Gamma)$ has a kernel with $\mathcal{O}(n^{k-1})$*
 272 *constraints that can be stored in $\mathcal{O}(n^{k-1} \log n)$ bits.*

273 ■ *If there exists $R \in \Gamma$ with $|R| = 2^k - 1$, then $\text{CSP}(\Gamma)$ has no generalized kernel of bitsize*
 274 *$\mathcal{O}(n^{k-\varepsilon})$ for any $\varepsilon > 0$, unless $\text{NP} \subseteq \text{coNP}/\text{poly}$.*

275 To obtain the kernels given in this section, we will heavily rely on the following procedure
 276 based on representing constraints by polynomials.

277 ▶ **Definition 3.2**. Let R be a k -ary Boolean relation. We say that a polynomial p_u captures
 278 an unsatisfying assignment $u \in \{0, 1\}^k \setminus R$ with respect to R over $\mathbb{Z}/q_u\mathbb{Z}$ or over \mathbb{Z} , if the
 279 following two conditions hold (over $\mathbb{Z}/q_u\mathbb{Z}$, or over \mathbb{Z} respectively).

$$280 \quad p_u(x_1, \dots, x_k) \equiv 0 \text{ for all } (x_1, \dots, x_k) \in R, \text{ and} \quad (1)$$

$$281 \quad p_u(u_1, \dots, u_k) \not\equiv 0. \quad (2)$$

283 The following Theorem is a generalization of Theorem 16 in [15]. The main improvement
 284 is that we now allow the usage of different polynomials, over different rings, for each $u \notin R$.
 285 Previously, all polynomials had to be given over the same ring, and each constraint was
 286 captured by a single polynomial.

287 ▶ **Theorem 3.3 (★)**. *Let $R \subseteq \{0, 1\}^k$ be a fixed k -ary relation, such that for every $u \in$*
 288 *$\{0, 1\}^k \setminus R$ there exists a polynomial p_u of degree at most d (and optionally, a prime power q_u)*
 289 *that captures u with respect to R (over $\mathbb{Z}/q_u\mathbb{Z}$ or over \mathbb{Z}). Then there exists a polynomial-time*
 290 *algorithm that, given a set of constraints \mathcal{C} over $\{R\}$ over n variables, outputs $\mathcal{C}' \subseteq \mathcal{C}$ with*
 291 *$|\mathcal{C}'| = \mathcal{O}(n^d)$, such that any Boolean assignment satisfies all constraints in \mathcal{C} if and only if it*
 292 *satisfies all constraints in \mathcal{C}' .*

293 The next lemma states that any k -ary relation R with $|\{0, 1\}^k \setminus R| > 1$ admits a non-
 294 trivial sparsification. To prove the lemma, we show that such relations can be represented
 295 by polynomials of degree at most $k - 1$, such that the sparsification can be obtained using
 296 Theorem 3.3. Since relations with $|R| = 2^k$ have a sparsification of size $\mathcal{O}(1)$, as constraints
 297 over such languages are satisfied by any assignment, it will follow that k -ary relations with
 298 $|\{0, 1\}^k \setminus R| \neq 1$ always allow a non-trivial sparsification.

299 ► **Lemma 3.4 (★)**. *Let R be a k -ary relation with $|R| < 2^k - 1$. Let \mathcal{C} be a set of constraints*
 300 *over $\{R\}$, using n variables. Then there exists a polynomial-time algorithm that outputs*
 301 *$\mathcal{C}' \subseteq \mathcal{C}$ with $|\mathcal{C}'| = \mathcal{O}(n^{k-1})$, such that a Boolean assignment satisfies all constraints in \mathcal{C}' if*
 302 *and only if it satisfies all constraints in \mathcal{C} .*

303 **Proof sketch.** We will show that for every $u \in \{0, 1\}^k \setminus R$, there exists a degree- $(k - 1)$
 304 polynomial p_u over the integers that captures u , such that the result follows from Theorem 3.3.
 305 We will prove the existence of such a polynomial by induction on k . For $k = 1$, the lemma
 306 statement implies that $R = \emptyset$. Thereby, for any $u \notin R$, we simply choose $p_u(x_1) := 1$. This
 307 polynomial satisfies the requirements, and has degree 0. Let $k > 1$ and let $u = (u_1, \dots, u_k) \in$
 308 $\{0, 1\}^k \setminus R$. Since $|R| < 2^k - 1$, we can choose $w = (w_1, \dots, w_k)$ such that $w \in \{0, 1\}^k \setminus R$
 309 and $w \neq u$. We distinguish two cases, depending on whether u and w agree on some position.

310 Suppose $u_i \neq w_i$ for all i , and assume for concreteness that $u = (0, \dots, 0)$ and $w =$
 311 $(1, \dots, 1)$. Then the polynomial $p_u(x_1, \dots, x_k) := \prod_{i=1}^{k-1} (i - \sum_{j=1}^k x_j)$ suffices: $p_u(0, \dots, 0) =$
 312 $\prod_{j=1}^{k-1} j \neq 0$, while for any $(x_1, \dots, x_k) \in R$, it holds that $\sum_{i=1}^k x_i \in [k - 1]$ and thereby
 313 $p_u(x_1, \dots, x_k) = 0$; the product has a 0-term. Other values of u and w are handled similarly.

314 Now suppose $u_i = w_i$ for some $i \in [k]$, and assume for concreteness that $u_1 = w_1 = 1$.
 315 Define $R' := \{(x_2, \dots, x_k) \mid (1, x_2, \dots, x_k) \in R\}$ and let $u' := (u_2, \dots, u_k)$. Since (u_2, \dots, u_k)
 316 and (w_2, \dots, w_k) are distinct tuples not in R' , by induction there is a polynomial $p_{u'}$ of
 317 degree $k - 2$ that captures u' with respect to R' . Then the polynomial $p_u(x_1, \dots, x_k) :=$
 318 $x_1 \cdot p_{u'}(x_2, \dots, x_k)$ has degree $k - 1$ and captures u with respect to R . ◀

319 To show the other part of the dichotomy, we will need the following theorem.

320 ► **Theorem 3.5 (★)**. *Let Γ be an intractable constraint language, and let $k \geq 1$. If there*
 321 *exists $R \in \Gamma$ such that R cone-defines k -OR, then $\text{CSP}(\Gamma)$ does not have a generalized kernel*
 322 *of size $\mathcal{O}(n^{k-\varepsilon})$, unless $\text{NP} \subseteq \text{coNP/poly}$.*

323 The next lemma formalizes the idea that any k -ary relation with $|\{0, 1\}^k \setminus R| = 1$ is
 324 equivalent to k -OR, up to negation of variables. The proof of the dichotomy given in Theorem
 325 3.1 will follow from Lemma 3.4, together with the next lemma and Theorem 3.5.

326 ► **Lemma 3.6 (★)**. *Let R be a k -ary relation with $|R| = 2^k - 1$. Then R cone-defines k -OR.*

327 4 From balanced operations to linear sparsification

328 The main result of this section is the following theorem, which we will prove later in this
 329 section.

330 ► **Theorem 4.1**. *Let Γ be a balanced constraint language. Then $\text{CSP}(\Gamma)$ has a kernel with*
 331 *$\mathcal{O}(n)$ constraints that are a subset of the original constraints. The kernel can be stored using*
 332 *$\mathcal{O}(n \log n)$ bits.*

333 To prove the theorem, we will use two additional technical lemmas. To state them, we
 334 introduce some notions from linear algebra. Given a set $S = \{s_1, \dots, s_n\}$ of k -ary vectors in
 335 \mathbb{Z}^k , we define $\text{span}_{\mathbb{Z}}(S)$ as the set of all vectors y in \mathbb{Z}^d for which there exist $\alpha_1, \dots, \alpha_n \in \mathbb{Z}$
 336 such that $y = \sum_{i \in [n]} \alpha_i s_i$. Similarly, we define $\text{span}_q(S)$ as the set of all k -ary vectors y over
 337 $\mathbb{Z}/q\mathbb{Z}$, such that there exist $\alpha_1, \dots, \alpha_n$ such that $y \equiv_q \sum_{i \in [n]} \alpha_i s_i$. For an $m \times n$ matrix S ,
 338 we use s_i for $i \in [m]$ to denote the i 'th row of S .

339 ► **Lemma 4.2 (★)**. *Let S be an $m \times n$ integer matrix. Let $u \in \mathbb{Z}^n$ be a row vector. If*
 340 *$u \in \text{span}_q(\{s_1, \dots, s_m\})$ for all prime powers q , then $u \in \text{span}_{\mathbb{Z}}(\{s_1, \dots, s_m\})$.*

341 ▶ **Lemma 4.3 (★)**. *Let q be a prime power. Let A be an $m \times n$ matrix over $\mathbb{Z}/q\mathbb{Z}$. Suppose*
 342 *there exists no constant $c \not\equiv_q 0$, such that the system $Ax \equiv_q b$ has a solution, where*
 343 *$b := (0, \dots, 0, c)^T$ is the vector with c on the last position and zeros in all other positions.*
 344 *Then $a_m \in \text{span}_q(\{a_1, \dots, a_{m-1}\})$.*

345 **Proof of Theorem 4.1.** We show that for all $R \in \Gamma$, for all $u \notin R$, there exists a linear
 346 polynomial p_u that captures u with respect to R . It will then follow from Theorem 3.3, that
 347 given an instance for $\text{CSP}(\Gamma)$ on n variables, we can reduce it to an equivalent instance with
 348 $\mathcal{O}(n)$ constraints. Since each constraint uses constantly many variables, we can store the
 349 resulting instance using $\mathcal{O}(n \log n)$ bits, as desired.

350 Suppose for contradiction that there exists $R \in \Gamma$ and $u \notin R$, such that no prime power
 351 q and polynomial p exists that satisfy conditions (1) and (2). We can view the process of
 352 finding such a linear polynomial, as solving a set of linear equations whose unknowns are
 353 the coefficients of the polynomial. We have a linear equation for each evaluation of the
 354 polynomial for which we want to enforce a certain value.

355 Let $R = \{r_1, \dots, r_\ell\}$. By the non-existence of p and q , the following system has no
 356 solution for any prime power q

$$357 \begin{pmatrix} 1 & r_{1,1} & r_{1,2} & \dots & r_{1,k} \\ 1 & r_{2,1} & r_{2,2} & \dots & r_{2,k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & r_{\ell,1} & r_{\ell,2} & \dots & r_{\ell,k} \\ 1 & u_1 & u_2 & \dots & u_k \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_k \end{pmatrix} \equiv_q \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ c \end{pmatrix}$$

358 and for any $c \not\equiv_q 0$. Otherwise, it is easy to verify that q is the desired prime power and
 359 $p(x_1, \dots, x_k) := \alpha_0 + \sum_{i=1}^k \alpha_i x_i$ is the desired polynomial.

360 The fact that no solution exists, implies that $(1, u_1, \dots, u_k)$ is in the span of the remaining
 361 rows of the matrix, by Lemma 4.3.

362 But this implies that for any prime power q , there exist coefficients $\beta_1, \dots, \beta_\ell$ over $\mathbb{Z}/q\mathbb{Z}$
 363 such that $u \equiv_q \sum \beta_i r_i$. Furthermore, since the first column of the matrix is the all-ones
 364 column, we obtain that $\sum \beta_i \equiv_q 1$. By Lemma 4.2, it follows that there exist integer
 365 coefficients $\gamma_1, \dots, \gamma_\ell$ such that $\sum \gamma_i = 1$ and furthermore $u = \sum \gamma_i r_i$. But it immediately
 366 follows that R is not preserved by the balanced operation given by $f(x_1, \dots, x_\ell) := \sum \gamma_i x_i$,
 367 which is a contradiction. ◀

368 The kernelization result above is obtained by using the fact that when Γ is balanced, the
 369 constraints in $\text{CSP}(\Gamma)$ can be replaced by linear polynomials. We show in the next theorem
 370 that this approach fails when Γ is not balanced.

371 ▶ **Theorem 4.4 (★)**. *Let R be a k -ary relation that is not balanced. Then there exists*
 372 *$u \in \{0, 1\}^k \setminus R$ such that there exists no polynomial p_u such that p_u captures u with respect*
 373 *to R over \mathbb{Z} , and no integer q such that p_u captures u over $\mathbb{Z}/q\mathbb{Z}$.*

374 5 Characterization of symmetric CSPs with linear sparsification

375 In this section, we characterize constraint languages Γ such that $\text{CSP}(\Gamma)$ has a linear
 376 sparsification, under the assumption that Γ is symmetric.

377 ▶ **Definition 5.1.** We say a k -ary Boolean relation R is *symmetric*, if there exists $S \subseteq$
 378 $\{0, 1, \dots, k\}$ such that a tuple $x = (x_1, \dots, x_k)$ is in R if and only if $\text{weight}(x) \in S$. We call
 379 S the set of *satisfying weights* for R .

23:10 Best-case and Worst-case Sparsifiability of Boolean CSPs

380 We will say that a constraint language Γ is symmetric, if it only contains symmetric relations.
 381 We will prove the following theorem at the end of this section.

- 382 ► **Theorem 5.2.** *Let Γ be a Boolean symmetric intractable constraint language.*
 383 ■ *If Γ is balanced, $\text{CSP}(\Gamma)$ has a kernel with $\mathcal{O}(n)$ constraints that can be stored in $\mathcal{O}(n \log n)$*
 384 *bits.*
 385 ■ *If Γ is not balanced, then $\text{CSP}(\Gamma)$ does not have a generalized kernel of size $\mathcal{O}(n^{2-\varepsilon})$ for*
 386 *any $\varepsilon > 0$, unless $\text{NP} \subseteq \text{coNP}/\text{poly}$.*

387 To show this, we use the following lemma.

388 ► **Lemma 5.3 (★).** *Let R be a symmetric relation with satisfying weights $S \subseteq \{0, 1, \dots, k\}$.*
 389 *Let $U := \{0, 1, \dots, k\} \setminus S$. If there exist $a, b, c \in S$ and $d \in U$ such that $a - b + c = d$, then R*
 390 *cone-defines 2-OR.*

391 **Proof sketch.** We will demonstrate the result in the case that $b \leq a$, $b \leq c$, and $b \leq d$; the
 392 other cases are similar. We use the following tuple to express $x_1 \vee x_2$.

$$393 \quad \underbrace{(\neg x_1, \dots, \neg x_1)}_{(a-b) \text{ copies}}, \underbrace{(\neg x_2, \dots, \neg x_2)}_{(c-b) \text{ copies}}, \underbrace{(1, \dots, 1)}_{b \text{ copies}}.$$

394 Let $f : \{x_1, x_2\} \rightarrow \{0, 1\}$, then $(\neg f(x_1), \dots, \neg f(x_1), \neg f(x_2), \dots, \neg f(x_2), 1, \dots, 1)$ has weight
 395 $d \notin S$ when $f(x_1) = f(x_2) = 0$. It is easy to verify that in all other cases, the weight is either
 396 a , b , or c . In these cases the tuple is contained in R , as the weight is contained in S . ◀

397 We now give the main lemma that is needed to prove Theorem 5.2. It shows that if a
 398 relation is symmetric and not balanced, it must cone-define 2-OR.

399 ► **Lemma 5.4.** *Let R be a symmetric relation of arity k . If R is not balanced, then R*
 400 *cone-defines 2-OR.*

401 **Proof.** Let f be a balanced operation that does not preserve R . Since f has integer
 402 coefficients, it follows that there exist (not necessarily distinct) $r_1, \dots, r_m \in R$, such that
 403 $r_1 - r_2 + r_3 - r_4 \cdots + r_m = u$ for some $u \in \{0, 1\}^k \setminus R$. Thereby, $\text{weight}(r_1) - \text{weight}(r_2) +$
 404 $\text{weight}(r_3) - \text{weight}(r_4) \cdots + \text{weight}(r_m) = \text{weight}(u)$. Let S be the set of satisfying weights
 405 for R and let $U := \{0, \dots, k\} \setminus S$. Define $s_i := \text{weight}(r_i)$ for $i \in [m]$, and $t = \text{weight}(u)$,
 406 such that $s_1 - s_2 + s_3 - s_4 \cdots + s_m = t$, and furthermore $s_i \in S$ for all i , and $t \in U$. We show
 407 that there exist $a, b, c \in S$ and $d \in U$ such that $a - b + c = d$, such that the result follows
 408 from Lemma 5.3. We do this by induction on the length of the alternating sum.

409 If $m = 3$, we have that $s_1 - s_2 + s_3 = t$ and define $a := s_1$, $b := s_2$, $c := s_3$, and $d := t$.

410 If $m > 3$, we will use the following claim.

411 ► **Claim 5.5 (★).** *Let $s_1, \dots, s_m \in S$ and $t \in U$ such that $s_1 - s_2 + s_3 - s_4 \cdots + s_m = t$.*
 412 *There exist distinct $i, j, \ell \in [m]$ with i, j odd and ℓ even, such that $s_i - s_\ell + s_j \in \{0, \dots, k\}$.*

413 Use Claim 5.5 to find i, j, ℓ such that $s_i - s_\ell + s_j \in \{0, \dots, k\}$. We consider two options.
 414 If $s_i - s_\ell + s_j \in U$, then define $d := s_i - s_\ell + s_j$, $a := s_i$, $b := s_\ell$, and $c := s_j$ and we
 415 are done. The other option is that $s_i - s_\ell + s_j = s \in S$. Replacing $s_i - s_j + s_k$ by s in
 416 $s_1 - s_2 + s_3 - s_4 \cdots + s_m$ gives a shorter alternating sum with result t . We obtain a, b, c , and
 417 d by the induction hypothesis.

418 Thereby, we have obtained $a, b, c \in S$, $d \in U$ such that $a - b + c = d$. It now follows from
 419 Lemma 5.3 that R cone-defines 2-OR. ◀

420 Using the lemma above, we can now prove Theorem 5.2.

421 **Proof of Theorem 5.2.** If Γ is balanced, it follows from Theorem 4.1 that $\text{CSP}(\Gamma)$ has a
 422 kernel with $\mathcal{O}(n)$ constraints that can be stored in $\mathcal{O}(n \log n)$ bits. Note that the assumption
 423 that Γ is symmetric is not needed in this case.

424 If Γ is not balanced, it follows that Γ contains a relation R that is not balanced. It
 425 follows from Lemma 5.4, that R cone-defines the 2-OR relation. Thereby, we obtain from
 426 Theorem 3.5 that $\text{CSP}(\Gamma)$ has no generalized kernel of size $\mathcal{O}(n^{2-\varepsilon})$ for any $\varepsilon > 0$, unless
 427 $\text{NP} \subseteq \text{coNP}/\text{poly}$. ◀

428 6 Low-arity classification

429 In this section, we will give a full classification of the sparsifiability for constraint languages
 430 that consist only of low-arity relations. The next results will show that in this case, if the
 431 constraint language is not balanced, it can cone-define the 2-OR relation.

432 ▶ **Observation 6.1.** *Each relation of arity 1 is balanced.*

433 ▶ **Theorem 6.2 (★).** *A relation of arity 2 is balanced if and only if it is not cone-interdefinable
 434 with the 2-OR relation.*

435 ▶ **Theorem 6.3 (★).** *Suppose that $U \subseteq \{0, 1\}^3$ is an arity 3 Boolean relation that is not
 436 balanced. Then, the 2-OR relation is cone-definable from U .*

437 Combining the results in this section with the results in previous sections, allows us
 438 to give a full classification of the sparsifiability of constraint languages that only contain
 439 relations of arity at most three. Observe that any k -ary relation R such that $R \neq \emptyset$ and
 440 $\{0, 1\}^k \setminus R \neq \emptyset$ cone-defines the 1-OR relation. Since we assume that Γ is intractable in the
 441 next theorem, it follows that k is always defined and $k \in \{1, 2, 3\}$.

442 ▶ **Theorem 6.4.** *Let Γ be an intractable Boolean constraint language such that each relation
 443 therein has arity ≤ 3 . Let $k \in \mathbb{N}$ be the largest value for which k -OR can be cone-defined from
 444 a relation in Γ . Then $\text{CSP}(\Gamma)$ has a kernel with $\mathcal{O}(n^k)$ constraints that can be encoded in
 445 $\mathcal{O}(n^k \log k)$ bits, but for any $\varepsilon > 0$ there is no kernel of size $\mathcal{O}(n^{k-\varepsilon})$, unless $\text{NP} \subseteq \text{coNP}/\text{poly}$.*

446 **Proof.** To show that there is a kernel with $\mathcal{O}(n^k)$ constraints, we do a case distinction on k .

447 ■ ($k = 1$) If $k = 1$, there is no relation in Γ that cone-defines the 2-OR relation. It follows
 448 from Observation 6.1 and Theorems 6.2 and 6.3 that thereby, Γ is balanced. It now
 449 follows from Theorem 4.1 that $\text{CSP}(\Gamma)$ has a kernel with $\mathcal{O}(n)$ constraints that can be
 450 stored in $\mathcal{O}(n \log n)$ bits.

451 ■ ($k = 2$) If $k = 2$, there is no relation $R \in \Gamma$ with $|R| = 2^3 - 1 = 7$, as otherwise by Lemma
 452 3.6 such a relation R would cone-define 3-OR which is a contradiction. Thereby, it follows
 453 from Theorem 3.1 that $\text{CSP}(\Gamma)$ has a sparsification with $\mathcal{O}(n^{3-1}) = \mathcal{O}(n^2)$ constraints
 454 that can be encoded in $\mathcal{O}(n^2 \log n)$ bits.

455 ■ ($k = 3$) Given an instance (\mathcal{C}, V) , it is easy to obtain a kernel of with $\mathcal{O}(n^3)$ constraints by
 456 simply removing duplicate constraints. This kernel can be stored in $\mathcal{O}(n^3)$ bits, by storing
 457 for each relation $R \in \Gamma$ and for each tuple $(x_1, x_2, x_3) \in V^3$ whether $R(x_1, x_2, x_3) \in \mathcal{C}$.
 458 Since $|\Gamma|$ is constant and there are $\mathcal{O}(n^3)$ such tuples, this results in using $\mathcal{O}(n^3)$ bits.

459 It remains to prove the lower bound. By definition, there exists $R \in \Gamma$ such that R
 460 cone-defines the k -OR relation. Thereby, the result follows immediately from Theorem 3.5.
 461 Thus, $\text{CSP}(\Gamma)$ has no kernel of size $\mathcal{O}(n^{k-\varepsilon})$ for any $\varepsilon > 0$, unless $\text{NP} \subseteq \text{coNP}/\text{poly}$. ◀

462 **7 Conclusion**

463 In this paper, we studied the best-case and worst-case sparsifiability of $\text{CSP}(\Gamma)$ for intractable
 464 constraint languages. First of all, we studied which constraint languages allow for non-trivial
 465 sparsification. If k is the largest arity of any relation in Γ , $\text{CSP}(\Gamma)$ has a trivial sparsification
 466 of size $\mathcal{O}(n^k)$. We introduced the notion of cone-definability and have shown that when no
 467 relation in Γ cone-defines k -OR, a non-trivial sparsification is always possible. By showing
 468 that if a relation cone-defines k -OR, it has no $\mathcal{O}(n^{k-\varepsilon})$ sparsification, we obtained a full
 469 classification of Boolean CSPs that admit a non-trivial sparsification.

470 Secondly, we introduced the notion of balanced constraint languages and have shown
 471 that $\text{CSP}(\Gamma)$ allows a sparsification with $\mathcal{O}(n)$ constraints whenever Γ is balanced. The
 472 constructive proof of this statement can be transformed into an effective algorithm to find a
 473 series of low-degree polynomials to capture the constraints, which earlier had to be done by
 474 hand. By combining the resulting upper and lower bound framework, we fully classified the
 475 symmetric constraint languages for which $\text{CSP}(\Gamma)$ allows a linear sparsification. Furthermore,
 476 we fully classified the sparsifiability of $\text{CSP}(\Gamma)$ when Γ contains relations of arity at most
 477 three. To obtain these results, we combined the above results with the additional insight that
 478 low-arity or symmetric relations that are not balanced, always cone-define the 2-OR relation.

479 The goal of this line of research is to fully classify the sparsifiability of $\text{CSP}(\Gamma)$, depending
 480 on Γ . In particular, we would like to classify those Γ for which $\mathcal{O}(n)$ sparsifiability is possible.
 481 In this paper, we have shown that Γ being balanced is a sufficient condition to obtain a linear
 482 sparsification. We conjecture that for intractable Boolean constraint languages Γ , $\text{CSP}(\Gamma)$
 483 admits no $\mathcal{O}(n)$ sparsification if Γ is not balanced. If true, this would give a full classification
 484 of the constraint languages that allow for linear sparsification.

485 One way to prove the conjecture, would be to show that if Γ is not balanced, it cone-defines
 486 2-OR and therefore has no subquadratic sparsification. As we have seen, this is true for both
 487 the low-arity and the symmetric case. Unfortunately, this does not hold in general, as can
 488 be shown by an example inspired by a construction of Lagerkvist and Wahlström [17, p.169].
 489 Let Γ consist of a single relation R , we start by defining the balanced relation that will not
 490 preserve R . Let $f(x_1, \dots, x_5) := x_1 - x_2 + x_3 - x_4 + x_5$ and let the domain D_f of f consist of
 491 those tuples $(x_1, \dots, x_5) \in \{0, 1\}^5$ for which $x_1 - x_2 + x_3 - x_4 + x_5 \in \{0, 1\}$. We now define
 492 R as a $|D_f|$ -ary relation with $R = \{r_1, r_2, r_3, r_4, r_5\}$, as follows. Let $D_f = \{d_1, \dots, d_m\}$ and
 493 let $d_i = \{d_{i,1}, \dots, d_{i,5}\}$ for $i \in [m]$. Define $r_{i,j} := d_{j,i}$ for $i \in [5]$ and $j \in [m]$. It can be shown
 494 that R is not balanced by observing that $f(r_1, \dots, r_5) \notin R$. However, R does not cone-define
 495 the 2-OR relation: by construction it is preserved by all alternating operations of arity three
 496 (because its outcome is never defined when applied to three distinct tuples from R), while a
 497 cone-definition of 2-OR would yield a violation by an arity-three alternating operation.

498 Luckily, cone-defining the 2-OR relation is not strictly necessary to obtain a non-linear
 499 lower bound. For example, a polynomial parameter transformation from 4-CNF-SAT to
 500 $\text{CSP}(\Gamma)$ using at most $\mathcal{O}(n^2)$ additional variables, also shows that no $\mathcal{O}(n^{2-\varepsilon})$ sparsification
 501 exists, unless $\text{NP} \subseteq \text{coNP/poly}$. We thereby leave as an open question whether there exist
 502 non-balanced constraint languages for which $\text{CSP}(\Gamma)$ allows a linear sparsification. Our
 503 suggested characterization in terms of balanced operations can be compared to that of
 504 Lagerkvist and Wahlström [17]. They suggest a weaker condition on constraint languages,
 505 and show that it is necessary and sufficient for $\text{CSP}(\Gamma)$ having a Maltsev embedding over
 506 an infinite domain (which does not have immediate consequences for sparsification). Our
 507 condition of being balanced is a stronger requirement, but as we have shown, leads directly
 508 to linear sparsification.

509 — **References** —

- 510 **1** Hans L. Bodlaender, Bart M. P. Jansen, and Stefan Kratsch. Kernelization lower bounds
511 by cross-composition. *SIAM J. Discrete Math.*, 28(1):277–305, 2014. doi:10.1137/
512 120880240.
- 513 **2** Hans L. Bodlaender, Stéphan Thomassé, and Anders Yeo. Kernel bounds for disjoint cycles
514 and disjoint paths. *Theor. Comput. Sci.*, 412(35):4570–4578, 2011. doi:10.1016/j.tcs.
515 2011.04.039.
- 516 **3** A. Bulatov, P. Jeavons, and A. Krokhin. Classifying the Complexity of Constraints using
517 Finite Algebras. *SIAM Journal on Computing*, 34(3):720–742, 2005. doi:10.1137/
518 S0097539700376676.
- 519 **4** Hubie Chen. A rendezvous of logic, complexity, and algebra. *ACM Computing Surveys*,
520 42(1), 2009. doi:10.1145/1189056.1189076.
- 521 **5** Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin
522 Pilipczuk, Michal Pilipczuk, and Saket Saurabh. *Parameterized Algorithms*. Springer, 2015.
523 doi:10.1007/978-3-319-21275-3.
- 524 **6** Holger Dell and Dániel Marx. Kernelization of packing problems. In *Proc. 23rd SODA*,
525 pages 68–81, 2012. doi:10.1137/1.9781611973099.6.
- 526 **7** Holger Dell and Dieter van Melkebeek. Satisfiability allows no nontrivial sparsification
527 unless the polynomial-time hierarchy collapses. *J. ACM*, 61(4):23:1–23:27, 2014. doi:
528 10.1145/2629620.
- 529 **8** Rodney G. Downey and Michael R. Fellows. *Fundamentals of Parameterized Complexity*.
530 Texts in Computer Science. Springer, 2013. doi:10.1007/978-1-4471-5559-1.
- 531 **9** Andrew Drucker. New limits to classical and quantum instance compression. *SIAM J.*
532 *Comput.*, 44(5):1443–1479, 2015. doi:10.1137/130927115.
- 533 **10** M.S. Gockenbach. *Finite-Dimensional Linear Algebra*. Discrete Mathematics and Its Ap-
534 plications. Taylor & Francis, 2011.
- 535 **11** Bart M. P. Jansen. On sparsification for computing treewidth. *Algorithmica*, 71(3):605–635,
536 2015. doi:10.1007/s00453-014-9924-2.
- 537 **12** Bart M. P. Jansen and Astrid Pieterse. Optimal sparsification for some binary CSPs
538 using low-degree polynomials. In *Proc. 41st MFCS*, pages 71:1–71:14, 2016. doi:10.4230/
539 LIPIcs.MFCS.2016.71.
- 540 **13** Bart M. P. Jansen and Astrid Pieterse. Optimal data reduction for graph coloring using
541 low-degree polynomials. In *Proc. 12th IPEC*, pages 22:1–22:12, 2017. doi:10.4230/LIPIcs.
542 IPEC.2017.22.
- 543 **14** Bart M. P. Jansen and Astrid Pieterse. Sparsification upper and lower bounds for
544 graph problems and not-all-equal SAT. *Algorithmica*, 79(1):3–28, 2017. doi:10.1007/
545 s00453-016-0189-9.
- 546 **15** Bart M. P. Jansen and Astrid Pieterse. Optimal sparsification for some binary CSPs using
547 low-degree polynomials. *CoRR*, abs/1606.03233, 2018. URL: [http://arxiv.org/abs/
548 1606.03233v2](http://arxiv.org/abs/1606.03233v2).
- 549 **16** Stefan Kratsch, Geevarghese Philip, and Saurabh Ray. Point line cover: The easy kernel is
550 essentially tight. *ACM Trans. Algorithms*, 12(3):40:1–40:16, 2016. doi:10.1145/2832912.
- 551 **17** Victor Lagerkvist and Magnus Wahlström. Kernelization of constraint satisfaction prob-
552 lems: A study through universal algebra. In *Proc. 23rd CP*, pages 157–171, 2017.
553 doi:10.1007/978-3-319-66158-2_11.
- 554 **18** Daniel Lokshtanov, Neeldhara Misra, and Saket Saurabh. Kernelization - Preprocessing
555 with a guarantee. In *The Multivariate Algorithmic Revolution and Beyond*, pages 129–161,
556 2012. doi:10.1007/978-3-642-30891-8_10.
- 557 **19** László Lovász. Chromatic number of hypergraphs and linear algebra. In *Studia Scientiarum*
558 *Mathematicarum Hungarica 11*, pages 113–114, 1976.

23:14 Best-case and Worst-case Sparsifiability of Boolean CSPs

- 559 20 Thomas J. Schaefer. The complexity of satisfiability problems. In *Proc. 10th STOC*, pages
560 216–226, 1978. doi:10.1145/800133.804350.

A

 Omitted proofs

A.1 Proofs omitted from Section 2

Proof of Proposition 2.4. It suffices to show that if a relation T is not balanced, then there exists an alternating operation that does not preserve T . Let f be a k -ary balanced operation that does not preserve T . Then there exist tuples t^1, \dots, t^k in T such that $\alpha_1 t^1 + \dots + \alpha_k t^k$ is not in T , where the sum of the α_i is equal to 1 (and where we may assume that no α_i is equal to 0). For each positive α_i , replace $\alpha_i t^i$ in the sum with $t^i + \dots + t^i$ (α_i times); likewise, for each negative α_i , replace $\alpha_i t^i$ in the sum with $-t^i - \dots - t^i$ ($-\alpha_i$ times). Each tuple then has coefficient $+1$ or -1 in the sum; since the sum of coefficients is $+1$, by permuting the sum's terms, the coefficients can be made to alternate between $+1$ and -1 . ◀

For the proof of Proposition 2.9, we will need the following additional theorem.

► **Theorem A.1.** (follows from [3]) Let Γ be a constraint language over a finite set D such that each unary operation $u : D \rightarrow D$ that preserves Γ is a bijection. Then, there exists a linear-parameter transformation from $\text{CSP}(\Gamma^*)$ to $\text{CSP}(\Gamma)$.

Note that in particular, an *intractable* Boolean constraint language can only be preserved by unary operations that are bijections. Hence for intractable Boolean Γ , there is a linear-parameter transformation from $\text{CSP}(\Gamma^*)$ to $\text{CSP}(\Gamma)$.

Proof of Theorem A.1. The desired transformation is the final polynomial-time reduction given in the proof of Theorem 4.7 of [3]. This reduction translates an instance of $\text{CSP}(\Gamma^*)$ with n variables to an instance of $\text{CSP}(\Gamma \cup \{=_D\})$ with $n + |D|$ variables; here, $=_D$ denotes the equality relation on domain D . Each constraint of the form $=_D(v, v')$ may be removed (while preserving satisfiability) by taking one of the variables v, v' , and replacing each instance of that variable with the other. The resulting instance of $\text{CSP}(\Gamma)$ has $\leq n + |D|$ variables. ◀

► **Definition A.2.** A relation $T \subseteq D^k$ is *pp-definable* (short for *primitive positive definable*) from a constraint language Γ over D if there exists an instance (\mathcal{C}, V) of $\text{CSP}(\Gamma)$ and there exist pairwise distinct variables $x_1, \dots, x_k \in V$ such that, for each map $f : \{x_1, \dots, x_k\} \rightarrow \{0, 1\}$, it holds that f can be extended to a satisfying assignment of the instance if and only if $(f(x_1), \dots, f(x_k)) \in T$.

The following is a known fact; for an exposition, we refer the reader to Theorems 3.13 and 5.1 of [4].

► **Proposition A.3.** If Γ is an intractable Boolean constraint language, then every Boolean relation is pp-definable from Γ^* .

Proof of Proposition 2.9. It suffices to give a linear-parameter transformation from $\text{CSP}(\Gamma^* \cup \Delta)$ to $\text{CSP}(\Gamma^*)$, by Theorem A.1. Let (\mathcal{C}, V) be an instance of $\text{CSP}(\Gamma^* \cup \Delta)$, and let n denote $|V|$. We generate an instance (\mathcal{C}', V') of $\text{CSP}(\Gamma^*)$ as follows.

■ For each variable $v \in V$, introduce a primed variable v' . By Proposition A.3, the relation \neq (that is, the relation $\{(0, 1), (1, 0)\}$) is pp-definable from Γ^* . Fix such a pp-definition, and let d be the number of variables in the definition. For each $v \in V$, include in \mathcal{C}' all constraints in the pp-definition of \neq , but where the variables are renamed so that v and v' are the distinguished variables, and the other variables are fresh.

The number of variables used so far in \mathcal{C}' is nd .

■ For each $b \in \{0, 1\}$, introduce a variable z_b , and include the constraint $\{(b)\}(z_b)$ in \mathcal{C}' .

603 ■ For each constraint $T(v_1, \dots, v_k)$ in \mathcal{C} such that $T \in \Gamma^*$, include the constraint in \mathcal{C}' .
604 ■ For each constraint $T(v_1, \dots, v_k)$ in \mathcal{C} such that $T \in \Delta \setminus \Gamma^*$, we use the assumption that
605 T is cone-definable from a relation in Γ to include a constraint in \mathcal{C}' that has the same
606 effect as $T(v_1, \dots, v_k)$. In particular, assume that T is cone-definable from $U \in \Gamma$ via
607 the tuple (y_1, \dots, y_ℓ) , and that U has arity ℓ . Include in \mathcal{C}' the constraint $U(w_1, \dots, w_\ell)$,
608 where, for each $i \in [\ell]$, the entry w_i is set to v_j if $y_i = x_j$, v'_j if $y_i = \neg x_j$, z_0 if $y_i = 0$, and
609 z_1 if $y_i = 1$.
610 The set V' of variables used in \mathcal{C}' is the union of $V \cup \{v' \mid v \in V\} \cup \{z_0, z_1\}$ with the
611 other variables used in the copies of the pp-definition of \neq . We have $|V'| = nd + 2$. It is
612 straightforward to verify that an assignment $f : V \rightarrow \{0, 1\}$ satisfies \mathcal{C} if and only if there
613 exists an assignment $f' : V' \rightarrow \{0, 1\}$ of f that satisfies \mathcal{C}' . ◀

614 A.2 Proofs omitted from Section 3

615 We start by proving the main Theorem of this section, using the other lemmas in the section.

616 **Proof of Theorem 3.1.** Suppose that for all $R \in \Gamma$, it holds that $|R| \neq 2^k - 1$. We give
617 the following kernelization procedure. Suppose we are given an instance of $\text{CSP}(\Gamma)$, with
618 set of constraints \mathcal{C} . We show how to define $\mathcal{C}' \subseteq \mathcal{C}$. For each constraint $R(x_1, \dots, x_\ell) \in \mathcal{C}$
619 where R is a relation of arity $\ell < k$, add one such constraint to \mathcal{C}' (thus removing duplicate
620 constraints). Note that this adds at most $\mathcal{O}(n^\ell)$ constraints for each ℓ -ary relation $R \in \Gamma$.

621 For a k -ary relation $R \in \Gamma$, let \mathcal{C}_R contain all constraints of the form $R(x_1, \dots, x_k)$. For
622 all k -ary relations R with $|R| < 2^k - 1$, apply Lemma 3.4 to obtain $\mathcal{C}'_R \subseteq \mathcal{C}_R$ such that
623 $|\mathcal{C}'_R| = \mathcal{O}(n^{k-1})$ and any Boolean assignment satisfying \mathcal{C}'_R also satisfies \mathcal{C}_R . Add \mathcal{C}'_R to
624 \mathcal{C}' . This concludes the definition of \mathcal{C}' . Note that the procedure removes constraints of the
625 form $R(x_1, \dots, x_k)$ with $|R| = 2^k$, as these are always satisfied. It is easy to verify that
626 $|\mathcal{C}'| \leq |\Gamma| \cdot \mathcal{O}(n^{k-1}) = \mathcal{O}(n^{k-1})$. Since each constraint can be stored in $\mathcal{O}(\log n)$ bits, this
627 gives a kernel of bitsize $\mathcal{O}(n^{k-1} \log n)$.

628 Suppose that there exists $R \in \Gamma$ with $|R| = 2^k - 1$. It follows from Lemma 3.6 that R
629 cone-defines k -OR. Since Γ is intractable, it now follows from Theorem 3.5 that $\text{CSP}(\Gamma)$ has
630 no generalized kernel of size $\mathcal{O}(n^{k-\varepsilon})$, unless $\text{NP} \subseteq \text{coNP}/\text{poly}$. ◀

631 To prove Theorem 3.3, we use the following theorem, that was proven by a subset of the
632 current authors [15]. We recall the required terminology. Define d -POLYNOMIAL ROOT CSP
633 over $\mathbb{Z}/m\mathbb{Z}$ as the problem whose input consists of a set L of polynomial equalities over $\mathbb{Z}/m\mathbb{Z}$
634 of degree at most d , over a set of variables V . Each equality is of the form $p(x_1, \dots, x_k) \equiv_m 0$.
635 The question is whether there exists a Boolean assignment to the variables in V that satisfies
636 all equalities in L .

637 ▶ **Theorem A.4** ([15, Theorem 16]). *There is a polynomial-time algorithm that, given an*
638 *instance (L, V) of d -POLYNOMIAL ROOT CSP over $\mathbb{Z}/m\mathbb{Z}$ for some fixed integer $m \geq 2$*
639 *with r distinct prime divisors, outputs an equivalent instance (L', V) of d -POLYNOMIAL*
640 *ROOT CSP over $\mathbb{Z}/m\mathbb{Z}$ with at most $r \cdot (n^d + 1)$ constraints such that $L' \subseteq L$.*

641 Note that if m is a prime power, m has only one distinct prime divisor and thereby $r = 1$
642 in the above theorem statement.

643 **Proof of Theorem 3.3.** Let \mathcal{C} be a set of constraints over R and let V be the set of variables
644 used. We will create $|\{0, 1\}^k \setminus R|$ instances of d -POLYNOMIAL ROOT CSP over V . For each
645 $u \in R \setminus \{0, 1\}^k$, we create an instance (L_u, V) of d -POLYNOMIAL ROOT CSP over $\mathbb{Z}/q_u\mathbb{Z}$,
646 as follows. Choose q_u and p_u such that they satisfy (1) and (2) for u . For each constraint

647 $(x_1, \dots, x_k) \in \mathcal{C}$, add the equality $p_u(x_1, \dots, x_k) \equiv_{q_u} 0$ to the set L_u . Let $L := \bigcup_{u \notin R} L_u$ be
 648 the union of all created sets of equalities. From this construction, we obtain the following
 649 claim.

650 ► **Claim A.5.** *Any Boolean assignment f that satisfies all equalities in L , satisfies all*
 651 *constraints in \mathcal{C} .*

652 **Proof.** Let f be a Boolean assignment that satisfies all equalities in L . Suppose f does not
 653 satisfy all equalities in \mathcal{C} , thus there exists $(x_1, \dots, x_k) \in \mathcal{C}$, such that $(f(x_1), \dots, f(x_k)) \notin R$.
 654 Let $u := (f(x_1), \dots, f(x_k))$. Since $u \notin R$, the equation $p_u(x_1, \dots, x_k) \equiv_{q_u} 0$ was added to
 655 $L_u \subseteq L$. However, it follows from (2) that $p_u(f(x_1), \dots, f(x_k)) \not\equiv_{q_u} 0$, which contradicts the
 656 assumption that f satisfies all equalities in L . ◻

657 For each instance (L_u, V) of d -POLYNOMIAL ROOT CSP, apply Theorem A.4 to obtain
 658 an equivalent instance (L'_u, V) with $L'_u \subseteq L_u$ and $|L'_u| = \mathcal{O}(n^d)$. Let $L' := \bigcup L'_u$. By this
 659 definition, any Boolean assignment satisfies all equalities in L , if and only if it satisfies all
 660 equalities in L' . Construct \mathcal{C}' as follows. For any $(x_1, \dots, x_k) \in \mathcal{C}$, add (x_1, \dots, x_k) to \mathcal{C}'
 661 if there exists $u \in \{0, 1\}^k \setminus R$ such that $p_u(x_1, \dots, x_k) \equiv_{q_u} 0 \in L'$. Hereby, $\mathcal{C}' \subseteq \mathcal{C}$. The
 662 following two claims show the correctness of this sparsification procedure.

663 ► **Claim A.6.** *Any Boolean assignment f satisfies all constraints in \mathcal{C}' , if and only if it*
 664 *satisfies all constraints in \mathcal{C} .*

665 **Proof.** Since $\mathcal{C}' \subseteq \mathcal{C}$, it follows immediately that any Boolean assignment satisfying the
 666 constraints in \mathcal{C} also satisfies all constraints in \mathcal{C}' . It remains to prove the opposite direction.

667 Let f be a Boolean assignment satisfying all constraints in \mathcal{C}' . We show that f satisfies
 668 all equalities in L' . Let $p_u(x_1, \dots, x_k) \equiv_{q_u} 0 \in L'$. Thereby, $(x_1, \dots, x_k) \in \mathcal{C}'$ and since
 669 f is a satisfying assignment, $(f(x_1), \dots, f(x_k)) \in R$. It follows from property (1) that
 670 $p_u(f(x_1), \dots, f(x_k)) \equiv_{q_u} 0$ as desired.

671 Since f satisfies all equalities in L' , it satisfies all equalities in L by the choice of L' . It
 672 follows from Claim A.5 that thereby f satisfies all constraints in \mathcal{C} . ◻

673 ► **Claim A.7.** $|\mathcal{C}'| = \mathcal{O}(n^d)$.

674 **Proof.** By the construction of \mathcal{C}' , it follows that $|\mathcal{C}'| \leq |L'|$. We know $|L'| = \sum_{u \notin R} |L'_u| \leq$
 675 $\sum_{u \notin R} \mathcal{O}(n^d) \leq 2^k \mathcal{O}(n^d) = \mathcal{O}(n^d)$, as k is considered constant. ◻

676 The lemma statement now follows from Claims A.6 and A.7. ◀

677 **Proof of Lemma 3.4.** We will prove this by showing that for every $u \in \{0, 1\}^k \setminus R$, there
 678 exists a k -ary polynomial p_u over the integers of degree at most $k - 1$ satisfying (1) and (2),
 679 such that the result follows from Theorem 3.3.

680 We will prove the existence of such a polynomial by induction on k . For $k = 1$, the lemma
 681 statement implies that $R = \emptyset$. Thereby, for any $u \notin R$, we simply choose $p_u(x_1) := 1$. This
 682 polynomial satisfies the requirements, and has degree 0.

683 Let $k > 1$ and let $u = (u_1, \dots, u_k) \in \{0, 1\}^k \setminus R$. Since $|R| < 2^k - 1$, we can choose
 684 $w = (w_1, \dots, w_k)$ such that $w \in \{0, 1\}^k \setminus R$ and $w \neq u$. Choose such w arbitrarily, we now
 685 do a case distinction.

686 **(There exists no $i \in [k]$, such that $u_i = w_i$)** This implies $u_i = \neg w_i$ for all i . One
 687 may note that for $u = (0, \dots, 0)$ and $w = (1, \dots, 1)$ this situation corresponds to MONOTONE
 688 k -NAE-SAT. We show that there exists a polynomial p_u such that $p_u(u_1, \dots, u_k) \neq 0$, and
 689 $p_u(x_1, \dots, x_k) = 0$ for all $(x_1, \dots, x_k) \in R$. Hereby p_u satisfies conditions (1) and (2) for u .

23:18 Best-case and Worst-case Sparsifiability of Boolean CSPs

690 For $i \in [k]$, define $r_i(x) := (1 - x)$ if $u_i = 1$ and $r_i(x) := x$ if $u_i = 0$. It follows immediately
 691 from this definition, that $r_i(u_i) = 0$ for all $i \in [k]$. Define

$$692 \quad p_u(x_1, \dots, x_k) := \prod_{i=1}^{k-1} \left(i - \sum_{j=1}^k r_j(x_j) \right).$$

693 By this definition, p_u has degree $k - 1$. It remains to verify that p_u has the desired properties.
 694 First of all, since $\sum_{j=1}^k r_j(u_j) = 0$ by definition, it follows that

$$695 \quad p_u(u_1, \dots, u_k) = \prod_{i=1}^{k-1} i \neq 0,$$

696 as desired. By the assumption that $u_i \neq w_i$ for all i , we obtain $p_u(w_1, \dots, w_k) = \prod_{i=1}^{k-1} (i - k) \neq$
 697 0 , which is allowed since $w \notin R$. It is easy to verify that in all other cases, $\sum_{j=1}^k r_j(x_j) \in$
 698 $\{1, 2, \dots, k-1\}$ and thereby one of the terms of the product is zero, implying $p_u(x_1, \dots, x_k) =$
 699 0 .

700 **(There exists $i \in [k]$, such that $u_i = w_i$)** Let u' and w' be defined as the results of
 701 removing coordinate i from u and w respectively. Note that $u' \neq w'$. Define

$$702 \quad R' := \{(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k) \mid (x_1, \dots, x_{i-1}, u_i, x_{i+1}, x_k) \in R\}.$$

703 By this definition, $u', w' \notin R'$ and thereby R' is a $(k - 1)$ -ary relation with $|R'| < 2^{k-1} - 1$.
 704 By the induction hypothesis, there exists a polynomial $p_{u'}$ of degree at most $k - 1$, such that
 705 $p_{u'}(u'_1, \dots, u'_{k-1}) \neq 0$ and $p_{u'}(x'_1, \dots, x'_{k-1}) = 0$ for all $x' \in R'$. Now define

$$706 \quad p_u(x_1, \dots, x_k) := (1 - x_i - u_i) \cdot p_{u'}(x_1, \dots, x_{i-1}, x_{i+1}, x_k).$$

707 We show that p_u has the desired properties. By definition, p_u has the degree of $p_{u'}$ plus one.
 708 Since $p_{u'}$ has degree $k - 2$ by the induction hypothesis, it follows that p_u has degree $k - 1$.
 709 Let $(x_1, \dots, x_k) \in R$. We do a case distinction on the value taken by x_i .

710 ■ $x_i \neq u_i$. In this case, $(1 - x_i - u_i) = 0$, and thereby $p_u(x_1, \dots, x_k) = 0$, thus satisfying
 711 equation (1).

712 ■ $x_i = u_i$. Since $x = (x_1, \dots, x_{i-1}, u_i, x_{i+1}, \dots, x_k) \in R$, it follows that $(x_1, \dots, x_{i-1}, x_{i+1},$
 713 $\dots, x_k) \in R'$. By definition of $p_{u'}$, it follows that $p_{u'}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k) = 0$ and
 714 thus $p_u(x_1, \dots, x_k) = 0$, showing (1).

715 It remains to show that $p_u(u_1, \dots, u_k) \neq 0$. This follows from $(1 - u_i - u_i) \in \{-1, 1\}$, and
 716 $p_{u'}(u_1, \dots, u_{i-1}, u_{i+1}, \dots, u_k) \neq 0$, showing that (2) holds.

717 Since we have shown for all $u \in \{0, 1\}^k \setminus R$, that there exists a polynomial p_u satisfying
 718 (1) and (2), the result follows from Theorem 3.3. ◀

719 **Proof of Theorem 3.5.** We do a case distinction on k .

720 ($k = 1$) Suppose that there exists $\varepsilon > 0$ such that $\text{CSP}(\Gamma)$ has a (generalized) kernel of
 721 size $\mathcal{O}(n^{1-\varepsilon})$. Using this hypothetical generalized kernel, one could obtain a polynomial-time
 722 algorithm that takes as input a series of inputs $(\mathcal{C}_1, V_1), \dots, (\mathcal{C}_t, V_t)$ of $\text{CSP}(\Gamma)$, and outputs
 723 in polynomial time an instance x^* of some fixed decision problem L such that:

724 ■ $x^* \in L$ if and only if all (\mathcal{C}_i, V_i) are YES-instances of $\text{CSP}(\Gamma)$, and

725 ■ x^* has bitsize $\mathcal{O}(N^{1-\varepsilon})$, where $N := \sum_{i=1}^t |V_i|$.

726 To obtain such an AND-compression algorithm from a hypothetical generalized kernel of
 727 $\text{CSP}(\Gamma)$ into a decision problem L , it suffices to do the following:

- 728 1. On input a series of instances $(\mathcal{C}_1, V_1), \dots, (\mathcal{C}_t, V_t)$ of $\text{CSP}(\Gamma)$, form a new instance $(\mathcal{C}^* :=$
729 $\bigcup_{i=1}^t \mathcal{C}_i, V^* := \bigcup_{i=1}^t V_i)$ of $\text{CSP}(\Gamma)$. Hence we take the disjoint union of the sets of
730 variables and the sets of constraints, and it follows that the new instance has answer YES
731 if and only if all the inputs (\mathcal{C}_i, V_i) have answer YES.
- 732 2. Run the hypothetical generalized kernel on (\mathcal{C}^*, V^*) , which has $|V^*| = N$ variables and is
733 therefore reduced to an equivalent instance x^* of L with bitsize $\mathcal{O}(N^{1-\varepsilon})$.

734 If we apply this AND-compression scheme to a sequence of $t_1(m) := m^\alpha$ instances of m bits each
735 (which therefore have at most m variables each), the resulting output has $\mathcal{O}(|V^*|^{1-\varepsilon}) = \mathcal{O}((m \cdot$
736 $m^\alpha)^{1-\varepsilon}) = \mathcal{O}(m^{(1+\alpha)(1-\varepsilon)})$ bits. By picking α large enough that it satisfies $(1+\alpha)(1-\varepsilon) \leq \alpha$,
737 we therefore compress a sequence of $t_1(m)$ instances of bitsize m into an equivalent instance
738 of size at most $t_2(m) \leq \mathcal{O}(m^{(1+\alpha)(1-\varepsilon)}) \leq C \cdot t_1(m)$, for some suitable constant C . Drucker [9,
739 Theorem 5.4] has shown that an error-free deterministic AND-compression algorithm with
740 these parameters for an NP-complete problem into a fixed decision problem L , implies
741 $\text{NP} \subseteq \text{coNP/poly}$. Hence the lower bound for $k = 1$ follows since $\text{CSP}(\Gamma)$ is NP-complete.

742 ($k \geq 2$) For $k \geq 2$, we prove the lower bound using a linear-parameter transformation
743 (recall Definition 2.1). Let Δ be the set of k -ary relations given by $\Delta := \{\{0, 1\}^k \setminus \{u\} \mid u \in$
744 $\{0, 1\}^k\}$. In particular, note that Δ contains the k -OR relation. Since R cone-defines k -OR, it
745 is easy to see that by variable negations, R cone-defines all relations in Δ . Thereby, it follows
746 from Proposition 2.9 that there is a linear-parameter transformation from $\text{CSP}(\Gamma^* \cup \Delta)$ to
747 $\text{CSP}(\Gamma)$. Thus, to prove the lower bound for $\text{CSP}(\Gamma)$, it suffices to prove the desired lower
748 bound for $\text{CSP}(\Gamma^* \cup \Delta)$.

749 ($k = 2$) If $k = 2$, we do a linear-parameter transformation from VERTEX COVER to
750 $\text{CSP}(\Gamma^* \cup \Delta)$. Since it is known that VERTEX COVER parameterized by the number of
751 vertices n has no generalized kernel of size $\mathcal{O}(n^{2-\varepsilon})$ for any $\varepsilon > 0$, unless $\text{NP} \subseteq \text{coNP/poly}$ [7],
752 the result follows.

753 Suppose we are given a graph $G = (V, E)$ and integer k for VERTEX COVER. We create
754 an instance (\mathcal{C}, V') of $\text{CSP}(\Gamma^*)$ as follows. We introduce a new variable x_v for each $v \in V$.
755 For each edge $\{u, v\} \in E$, we add the constraint $2\text{-OR}(x_u, x_v)$ to \mathcal{C} .

756 At this point, any vertex cover in G corresponds to a satisfying assignment, and vice versa.
757 It remains to ensure that the size is bounded by k . Let $H_{n,k}$ be the n -ary relation given
758 by $H_{n,k} = \{(x_1, \dots, x_n) \mid x_i \in \{0, 1\} \text{ for all } i \in [n] \text{ and } \sum_{i \in [n]} x_i = k\}$. By Proposition A.3,
759 we obtain that Γ^* pp-defines all boolean relations. It follows from [17, Lemma 17] that Γ^*
760 pp-defines $H_{n,k}$ using $\mathcal{O}(n+k)$ constraints and $\mathcal{O}(n+k)$ existentially quantified variables.
761 We add the constraints from this pp-definition to \mathcal{C} , and remove the existential quantification.
762 This concludes the construction of \mathcal{C} . It is easy to see that \mathcal{C} has a satisfying assignment if
763 and only if G has a vertex cover of size k . Furthermore, we used $\mathcal{O}(n)$ variables and thereby
764 this is a linear-parameter transformation from VERTEX COVER to $\text{CSP}(\Gamma^* \cup \Delta)$.

765 ($k \geq 3$) In this case there is a trivial linear-parameter transformation from $\text{CSP}(\Delta)$ to
766 $\text{CSP}(\Gamma^* \cup \Delta)$. It is easy to verify that $\text{CSP}(\Delta)$ is equivalent to k -CNF-SAT. The result now
767 follows from the fact that for $k \geq 3$, k -CNF-SAT has no kernel of size $\mathcal{O}(n^{k-\varepsilon})$ for any $\varepsilon > 0$,
768 unless $\text{NP} \subseteq \text{coNP/poly}$ [7]. \blacktriangleleft

769 **Proof of Lemma 3.6.** Let $u = (u_1, \dots, u_k)$ be the unique k -tuple not contained in R . Define
770 the tuple (y_1, \dots, y_k) as follows. Let $y_i := x_i$ if $u_i = 0$, and let $y_i := \neg x_i$ otherwise.
771 Clearly, this satisfies the first two conditions of cone-definability. It remains to prove the
772 last condition. Let $f: \{x_1, \dots, x_m\} \rightarrow \{0, 1\}$. Suppose $(f(x_1), \dots, f(x_k)) \in k\text{-OR}$. We
773 show $(\hat{f}(y_1), \dots, \hat{f}(y_k)) \in R$. Since $(f(x_1), \dots, f(x_k)) \in k\text{-OR}$, there exists at least one
774 $i \in [k]$ such that $f(x_i) \neq 0$. Thereby, $\hat{f}(y_i) \neq u_i$ and thus $(\hat{f}(y_1), \dots, \hat{f}(y_k)) \neq u$, implying
775 $(\hat{f}(y_1), \dots, \hat{f}(y_k)) \in R$.

776 Suppose $(f(x_1), \dots, f(x_k)) \notin k\text{-OR}$, implying $f(x_i) = 0$ for all $i \in [k]$. But this implies
 777 $\hat{f}(y_i) = u_i$ for all $i \in [k]$ and thus $(\hat{f}(y_1), \dots, \hat{f}(y_k)) = u \notin R$. ◀

778 A.3 Proofs omitted from Section 4

779 To give the proofs that were omitted from Section 4, we need the following additional
 780 definitions.

781 ▶ **Definition A.8.** We say an $m \times n$ matrix A is a *diagonal matrix*, if all entries $a_{i,j}$ with
 782 $i \neq j$ are zero. Thus, all non-zero elements occur on the diagonal.

783 Note that by the above definition of diagonal matrices, a matrix can be diagonal even if it is
 784 not a square matrix.

785 We denote the greatest common divisor of two integers x and y as $\gcd(x, y)$. Recall that
 786 by Bézout's lemma, if $\gcd(x, y) = z$ there exist integers a and b such that $ax + by = z$. We
 787 will use $x \mid y$ to indicate that x divides y (over the integers) and $x \nmid y$ to indicate that it does
 788 not.

789 **Proof of Lemma 4.2.** Suppose $u \notin \text{span}_{\mathbb{Z}}(\{s_1, \dots, s_m\})$, thus u cannot be written as a
 790 linear combination of the rows of S over \mathbb{Z} . This implies that the system $yS = u$ has no
 791 solutions for y over \mathbb{Z} . We will show that there exists a prime power q , such that $yS \equiv_q u$
 792 has no solutions over $\mathbb{Z}/q\mathbb{Z}$ and thus $u \notin \text{span}_q(\{s_1, \dots, s_m\})$.

793 There exist an $m \times m$ matrix M and an $n \times n$ matrix N over \mathbb{Z} , such that M and N are
 794 invertible over \mathbb{Z} and furthermore $S' := MSN$ is in Smith Normal Form (c.f. [10, Theorem
 795 368]). In particular, this implies that all non-zero elements of S' are on the diagonal. Define
 796 $u' := uN$.

797 ▶ **Claim A.9.** $yS = u$ is solvable for y over \mathbb{Z} , if and only if $y'S' = u'$ is solvable for y'
 798 over \mathbb{Z} .

799 **Proof.** Let y be such that $yS = u$. Define $y' := yM^{-1}$. We verify that $y'S' = u'$ as follows.

$$800 \quad y'S' = y'MSN = yM^{-1}MSN = ySN = uN = u'.$$

801 For the other direction, consider y' such that $y'S' = u'$. In this case one can verify that
 802 $y := y'M$ solves $yS = u$, as

$$803 \quad yS = y'MS = y'MSN N^{-1} = y'S' N^{-1} = u' N^{-1} = u N N^{-1} = u. \quad \lrcorner$$

804 It remains to show that if $y'S' = u'$ is unsolvable over \mathbb{Z} , then there exists a prime power q
 805 such that $y'S' \equiv_q u'$ is unsolvable. By Claim A.9, this implies the lemma statement.

806 Suppose $y'S' = u'$ has no solutions. Since all non-zero elements of S' are on the diagonal,
 807 this implies that either there exists $i \in [n]$, such that u'_i is not divisible by $s'_{i,i}$, or $s'_{i,i}$ is zero
 808 while $u'_i \neq 0$. We finish the proof by a case distinction.

809 ■ Suppose there exists $i \in [n]$ such that $s'_{i,i} = 0$, while $u'_i \neq 0$. Choose a prime power q
 810 such that $q \nmid u'_i$. It is easy to see that thereby, $u'_i \not\equiv_q 0$. Since $s'_{i,i} \equiv_q 0$ holds trivially in
 811 this case, the system has $y'S' \equiv_q u'$ no solution.

812 ■ Otherwise, there exists $i \in [n]$ such that $s'_{i,i} \nmid u'_i$. Choose a prime power q such that
 813 $q \mid u'_i$ and $q \nmid s'_{i,i}$. Such a prime power can be chosen by letting $q := p^\ell$ for a prime p that
 814 occurs ℓ times in the prime factorization of u'_i , but less often in the prime factorization of
 815 $s'_{i,i}$. Thereby, $u'_i \not\equiv_q 0$, while $s'_{i,i} \equiv_q 0$. It again follows that the system $y'S' \equiv_q u'$ has no
 816 solutions. ◀

817 **Proof of Lemma 4.3.** Let A' be the $(m-1) \times n$ matrix consisting of the first $m-1$ rows
 818 of A . Find the Smith normal form [10] of A' over \mathbb{Z} , thus there exist an $(m-1) \times (m-1)$
 819 matrix M' and an $n \times n$ matrix N , such that $S' := M'A'N$ is in Smith Normal Form and
 820 M' and N are invertible over \mathbb{Z} .

821 We show that similar properties hold over $\mathbb{Z}/q\mathbb{Z}$. Let $(M')^{-1}$, N^{-1} be the inverses of
 822 M' and N over \mathbb{Z} , it is easy to verify that $NN^{-1} = I \equiv_q I$ and $M'(M')^{-1} = I \equiv_q I$, such
 823 that M' and N are still invertible over $\mathbb{Z}/q\mathbb{Z}$. Furthermore, $S' \pmod{q}$ remains a diagonal
 824 matrix.

825 Define M to be the following $m \times m$ matrix

$$826 \quad M := \left(\begin{array}{c|c} M' & \mathbf{0} \\ \hline \mathbf{0} & 1 \end{array} \right),$$

827 then M has an inverse over $\mathbb{Z}/q\mathbb{Z}$ that is given by the following matrix

$$828 \quad M^{-1} \equiv_q \left(\begin{array}{c|c} (M')^{-1} & \mathbf{0} \\ \hline \mathbf{0} & 1 \end{array} \right).$$

829 Define $S := MAN$ and verify that

$$830 \quad S := MAN \equiv_q \left(\begin{array}{c} S' \\ a_m N \end{array} \right), \quad (3)$$

831 meaning that the first $m-1$ rows of S are equal to the first $m-1$ rows of S' , and the last
 832 row of S is given by the row vector $a_m N$.

833 The following two claims will be used to show that proving the lemma statement for
 834 matrix S , will give the desired result for A .

835 **► Claim A.10.** *Let $b := (0, \dots, 0, c)$ for some constant c . The system $Sx' \equiv_q b$ has a solution,*
 836 *if and only if the system $Ax \equiv_q b$ has a solution.*

837 **Proof.** Let x be a solution for $Ax \equiv_q b$. Define $x' := N^{-1}x$. Then $MANx' \equiv_q MAx \equiv_q Mb$.
 838 Observe that by the definitions of M and b , $Mb \equiv_q b$, which concludes this direction of the
 839 proof.

840 For the other direction, let x' be a solution for $MANx' \equiv_q b$. Define $x := Nx'$. Then
 841 $M^{-1}MANx' \equiv_q M^{-1}b$ and thus $ANx' \equiv_q M^{-1}b$ and thereby $Ax \equiv_q M^{-1}b$. By the
 842 definition of M^{-1} and b , we again have $M^{-1}b \equiv_q b$. \square

843 **► Claim A.11.** *$s_m \in \text{span}_q(\{s_1, \dots, s_{m-1}\})$ if and only if $a_m \in \text{span}_q(\{a_1, \dots, a_{m-1}\})$.*

844 **Proof.** Suppose $s_m \in \text{span}_q(\{s_1, \dots, s_{m-1}\})$. This implies that there exist $\alpha_1, \dots, \alpha_{m-1}$
 845 such that $\sum_{i \in [m-1]} \alpha_i s_i \equiv_q s_m \equiv_q a_m N$. Thus, $\sum_{i \in [m-1]} \alpha_i s'_i \equiv_q a_m N$, and for $\alpha =$
 846 $(\alpha_1, \dots, \alpha_{m-1})$ we therefore have $\alpha S' \equiv_q a_m N$, implying $(\alpha M')A'N \equiv_q a_m N$. Since N is
 847 invertible, it follows that

$$848 \quad (\alpha M')A' \equiv_q a_m$$

849 and thus $a_m \in \text{span}_q(\{a_1, \dots, a_{m-1}\})$.

850 For the other direction, suppose $a_m \in \text{span}_q(\{a_1, \dots, a_{m-1}\})$. Thus, there exists $\alpha \equiv_q$
 851 $(\alpha_1, \dots, \alpha_{m-1})$ such that $\alpha A' \equiv_q a_m$. Let $\alpha' := \alpha(M')^{-1}$. Then

$$852 \quad \alpha' S' \equiv_q \alpha' M' A' N \equiv_q \alpha A' N \equiv_q a_m N \equiv_q s_m,$$

853 and now it follows from the definition of S given in (3) that $s_m \in \text{span}_q(\{s_1, \dots, s_{m-1}\})$. \square

854 It follows from Claims A.10 and A.11, that it suffices to show that if $Sx = (0, \dots, 0, c)^T$
 855 has no solutions for any $c \not\equiv_q 0$, then $s_m \in \text{span}_q(\{s_1, \dots, s_{m-1}\})$. So suppose $s_m \notin$
 856 $\text{span}_q(\{s_1, \dots, s_{m-1}\})$, we show that the system has a solution for some non-zero c . Observe
 857 that since S' (the first $m-1$ rows of S) is a diagonal matrix, there must exist $i \in [m-1]$
 858 such that there is no α_i such that $s_{i,i} \cdot \alpha_i \equiv_q s_{m,i}$. Otherwise, it is easy to see that
 859 $\sum_{i \in [m-1]} \alpha_i s_i \equiv_q s_m$, contradicting that $s_m \notin \text{span}_q(\{s_1, \dots, s_{m-1}\})$. We now do a case
 860 distinction.

861 Suppose there exists $i \in [m-1]$ such that $s_{i,i} \equiv_q 0$, while $s_{m,i} \not\equiv_q 0$. Let $x =$
 862 $(0, \dots, 0, 1, 0, \dots, 0)$ be the vector with 1 in the i 'th position and zeros in all other po-
 863 sitions. It is easy to verify that $Sx \equiv_q (0, \dots, 0, s_{m,i})^T$ and thereby the system $Sx \equiv_q b$ has
 864 a solution for $c = s_{m,i}$.

865 Otherwise, choose i such that there exists no integer α_i such that $s_{i,i} \cdot \alpha_i \equiv_q s_{m,i}$ and
 866 $s_{i,i} \not\equiv_q 0$. It is given that q is a prime power, let $q = p^k$ for prime p . Let $0 \leq \ell < k$ be the
 867 largest integer such that $p^\ell \mid s_{i,i}$ over the integers. We consider the following two cases.

868 ■ Suppose $p^\ell \mid s_{m,i}$. Let c such that $s_{i,i} \equiv_q cp^\ell$ and choose d such that $s_{m,i} \equiv_q d \cdot p^\ell$. Note
 869 that $\gcd(c, q) = 1$ and it follows from Bézout's lemma that c has an inverse c^{-1} such that
 870 $cc^{-1} \equiv_q 1$. Then

$$871 \quad s_{i,m} \equiv_q (d \cdot c^{-1})s_{i,i},$$

872 which is a contradiction with the assumption that no integer α_i exists such that
 873 $s_{i,i} \cdot \alpha_i \equiv_q s_{i,m}$.

874 ■ Suppose $p^\ell \nmid s_{m,i}$. Define $x := (0, \dots, 0, p^{k-\ell}, 0, \dots, 0)^T$ as the vector with $p^{k-\ell}$ in
 875 position i . Then

$$876 \quad Sx \equiv_q (0, \dots, 0, p^{k-\ell} \cdot s_{i,i}, 0, \dots, 0, p^{k-\ell} \cdot s_{m,i})^T.$$

877 Since $p^\ell \mid s_{i,i}$ it follows that $p^{k-\ell} \cdot s_{i,i} \equiv_q 0$. Furthermore, since $p^\ell \nmid s_{m,i}$, it follows that
 878 $p^{k-\ell} \cdot s_{m,i} \not\equiv_q 0$, and thereby the system $Sx \equiv_q b$ has a solution for $b := (0, \dots, 0, p^{k-\ell} s_{m,i})$.
 879 ◀

880 **Proof of Theorem 4.4.** Suppose R is not balanced, let f be a balanced operation that does
 881 not preserve R . Let the coefficients of f be $\alpha_1, \dots, \alpha_m$. Then, there exist $r_1, \dots, r_m \in$
 882 R and $u \in \{0, 1\}^k \setminus R$ such that $\sum_{i \in [m]} \alpha_i r_i = u$. Now suppose there exists a linear
 883 polynomial p_u (and optionally, a prime power q_u) such that p_u captures u (over \mathbb{Z} or $\mathbb{Z}/q_u\mathbb{Z}$
 884 respectively). Let $r_i = (r_{i,1}, \dots, r_{i,k})$ for $i \in [m]$. By definition, $p_u(r_{i,1}, \dots, r_{i,k}) \equiv 0$ (over \mathbb{Z}
 885 or $\mathbb{Z}/q_u\mathbb{Z}$ respectively). Since p_u is a linear polynomial, suppose $p_u(x_1, \dots, x_k)$ is given by
 886 $p_u(x_1, \dots, x_k) \equiv \beta_0 + \sum_{i \in [m]} \beta_i x_i$ for integer coefficients β_i . Let $u = (u_1, \dots, u_k)$. But then,

$$\begin{aligned}
 887 \quad p_u(u_1, \dots, u_k) &\equiv \beta_0 + \sum_{i \in [m]} \beta_i u_i \\
 888 &\equiv \beta_0 + \sum_{i \in [m]} \beta_i \left(\sum_{j \in [m]} \alpha_j r_{j,i} \right) \\
 889 &\equiv \beta_0 + \sum_{i \in [m]} \sum_{j \in [m]} \alpha_j \beta_i r_{j,i} \\
 890
 \end{aligned}$$

891 (Since $\sum_{i \in [m]} \alpha_i \equiv 1$)

$$\begin{aligned}
 892 \quad & \equiv \sum_{j \in [m]} \alpha_j \beta_0 + \sum_{j \in [m]} \sum_{i \in [m]} \alpha_j \beta_i r_{j,i} \\
 893 \quad & \equiv \sum_{j \in [m]} \alpha_j \beta_0 + \alpha_j \sum_{i \in [m]} \beta_i r_{j,i} \\
 894 \quad & \equiv \sum_{j \in [m]} p_u(r_{j,1}, \dots, r_{j,k}) \equiv 0, \\
 895 \quad &
 \end{aligned}$$

896 But this contradicts the fact that $p_u(u_1, \dots, u_k) \neq 0$. Thereby, there exists no linear
 897 polynomial that captures u with respect to R . \blacktriangleleft

898 A.4 Proofs omitted from Section 5

899 **Proof of Lemma 5.3.** We first show the result when $b \leq a$, $b \leq c$, and $b \leq d$. In this case,
 900 we use the following tuple to express $x_1 \vee x_2$.

$$901 \quad (\underbrace{\neg x_1, \dots, \neg x_1}_{(a-b) \text{ copies}}, \underbrace{\neg x_2, \dots, \neg x_2}_{(c-b) \text{ copies}}, \underbrace{1, \dots, 1}_{b \text{ copies}}).$$

902 Let $f : \{x_1, x_2\} \rightarrow \{0, 1\}$, then $(\neg f(x_1), \dots, \neg f(x_1), \neg f(x_2), \dots, \neg f(x_2), 1, \dots, 1)$ has weight
 903 $(a-b)(1-f(x_1)) + (c-b)(1-f(x_2)) + b$. It is easy to verify that for $f(x_1) = f(x_2) = 0$,
 904 this implies the tuple has weight $a + c - b = d \notin S$ and thus the tuple is not in R . Otherwise,
 905 the weight is either a , b , or c . In these cases the tuple is contained in R , as the weight is
 906 contained in S .

907 Note that the above case applies when b is the smallest of all four integers. We now
 908 consider the remaining cases. Suppose $a \leq b$, $a \leq c$, and $a \leq d$ (the case where c is smallest
 909 is symmetric). In this case, use the tuple

$$910 \quad (\underbrace{\neg x_1, \dots, \neg x_1}_{(d-a) \text{ copies}}, \underbrace{x_2, \dots, x_2}_{(b-a) \text{ copies}}, \underbrace{1, \dots, 1}_{a \text{ copies}}).$$

911 Consider an assignment f satisfying $x_1 \vee x_2$, verify that the weight of the above tuple under
 912 this assignment lies in $\{a, b, c\}$, and thus the tuple is contained in R . Assigning 0 to both x_1
 913 and x_2 gives weight d , such that the tuple is not in R .

914 Otherwise, we have $d \leq a$, $d \leq b$, and $d \leq c$ and use the tuple

$$915 \quad (\underbrace{x_1, \dots, x_1}_{(a-d) \text{ copies}}, \underbrace{x_2, \dots, x_2}_{(c-d) \text{ copies}}, \underbrace{1, \dots, 1}_{d \text{ copies}}).$$

916 It is again easy to verify that any assignment to x_1 and x_2 satisfies this tuple if and only if it
 917 satisfies $(x_1 \vee x_2)$, using the fact that $a - d + c = b \in S$. \blacktriangleleft

918 **Proof of Claim 5.5.** If there exist distinct $i, j, \ell \in [m]$ with i, j odd and ℓ even, such that
 919 $s_i \geq s_\ell \geq s_j$, then these i, j, ℓ satisfy the claim statement. Suppose these do not exist, we
 920 consider two options.

- 921 ■ Suppose $s_i \geq s_\ell$ for all $i, \ell \in [m]$ with i odd and ℓ even. It is easy to see that thereby,
 922 for any i, j, ℓ with i, j odd and ℓ even it holds that $s_i - s_\ell + s_j \geq 0$. Furthermore,
 923 $s_i - s_\ell + s_j \leq s_1 - s_2 + s_3 - s_4 \cdots + s_m = t$ and $t \leq k$ since $t \in U$. Thus, any distinct
 924 $i, j, \ell \in [m]$ with i, j odd and ℓ even satisfy the statement.

925 ■ Otherwise, $s_i \leq s_\ell$ for all $i, \ell \in [m]$ with i odd and ℓ even. It follows that for any i, j, ℓ
 926 with i, j odd and ℓ even $s_i - s_\ell + s_j \leq k$, as $s_i - s_\ell \leq 0$ and $s_j \leq k$. Furthermore,
 927 $s_i - s_\ell + s_j \geq s_1 - s_2 + s_3 - s_4 \cdots + s_m = t$ and $t \geq 0$ by definition. Thus, any distinct
 928 $i, j, \ell \in [m]$ with i, j odd and ℓ even satisfy the statement. \lrcorner

929 A.5 Proofs omitted from Section 6

930 To state the proofs that were omitted from this section, we will first give a number of relevant
 931 observations and propositions.

932 ► **Observation A.12.** *If a relation T of arity m is cone-definable from a relation U of arity*
 933 *n , then $m \leq n$.*

934 ► **Observation A.13.** *Suppose a relation T is cone-definable from a relation U , and that g*
 935 *is a partial operation that is idempotent and self-dual. If g preserves U , then g preserves T .*

936 ► **Observation A.14.** *(transitivity of cone-definability) Suppose that T_1, T_2, T_3 are relations*
 937 *such that T_2 is cone-definable from T_1 , and T_3 is cone-definable from T_2 . Then T_3 is*
 938 *cone-definable from T_1 .*

939 ► **Definition A.15.** Let us say that two Boolean relations T, U are *cone-interdefinable* if
 940 each is cone-definable from the other.

941 The following two propositions are consequences of Observations A.12 and A.13. We will
 942 tacitly use them in the sequel. Morally, they show that the properties of relations that we
 943 are interested in are invariant under cone-interdefinability.

944 ► **Proposition A.16.** *If relations T, U are cone-interdefinable, then they have the same arity.*

945 ► **Proposition A.17.** *Suppose that T and U are relations that are cone-interdefinable, and*
 946 *that g is a partial operation that is idempotent and self-dual. Then, g preserves T if and only*
 947 *if g preserves U .*

948 **Proof of Theorem 6.2.** Let $R \subseteq \{0, 1\}^2$ be a relation. We prove the two directions sepa-
 949 rately.

950 (\Leftarrow) Proof by contraposition. Suppose that R is cone-interdefinable with 2-OR. Then
 951 in particular, R cone-defines the 2-OR relation. Let (y_1, y_2) be a tuple witnessing cone-
 952 definability as in Definition 2.7. Since 2-OR is symmetric in its two arguments, we may assume
 953 without loss of generality that y_i is either x_i or $\neg x_i$ for $i \in [2]$. Define $g: \{0, 1\}^2 \rightarrow \{0, 1\}^2$
 954 by letting $g(i_1, i_2) := (\hat{i}_1, \hat{i}_2)$ where $\hat{i}_\ell = i_\ell$ if $y_i = x_i$ and $\hat{i}_\ell = 1 - i_\ell$ if $y_i = \neg x_i$. By
 955 definition of cone-definability we then have $g(1, 0), g(0, 1), g(1, 1) \in R$ while $g(0, 0) \notin R$.
 956 But $g(1, 0) - g(1, 1) + g(0, 1) = g(0, 0)$, showing that R is not preserved by all alternating
 957 operations and therefore is not balanced.

958 (\Rightarrow) We again use contraposition. Suppose R is not balanced; we will prove R is cone-
 959 interdefinable with 2-OR. Let $f: \{0, 1\}^k \rightarrow \{0, 1\}$ be a balanced partial Boolean operation of
 960 *minimum arity* that does not preserve R . Let $\alpha_1, \dots, \alpha_k \in \mathbb{Z}$ be the coefficients of f , as in
 961 Definition 2.2. Let $s^1, \dots, s^k \in R$ such that $f(s^1, \dots, s^k) = u \in \{0, 1\}^2 \setminus R$ witnesses that f
 962 does not preserve R . By Definition 2.2 we have $u = \sum_{i=1}^k \alpha_i s^i$ and $\sum_{i=1}^k \alpha_i = 1$. This shows
 963 that if $\alpha_i = 0$ for some coordinate i , then that position does not influence the value of f ,
 964 implying the existence of a smaller-arity balanced relation that does not preserve T . Hence
 965 our choice of f as a minimum-arity operation ensures that α_i is nonzero for all $i \in [k]$.

966 ► **Claim A.18.** *The tuples s^1, \dots, s^k are all distinct.*

967 **Proof.** Suppose that $s^i = s^j$ for some distinct $i, j \in [k]$, and assume without loss of generality
 968 that $i = k - 1$ and $j = k$. But then the balanced operation f' of arity $k - 1$ defined
 969 by the coefficients $(\alpha_1, \dots, \alpha_{k-2}, \alpha_{k-1} + \alpha_k)$ does not preserve T since $f'(s^1, \dots, s^{k-1}) =$
 970 $f(s^1, \dots, s^k) = u \notin R$, contradicting that f is a minimum-arity balanced operation that does
 971 not preserve R . \square

972 **► Claim A.19.** *The arity k of operation f is 3.*

973 **Proof.** Since $s^1, \dots, s^k \in R \subseteq \{0, 1\}^2$ are all distinct, while $u \in \{0, 1\}^2 \setminus R$, we have $k \leq 3$. We
 974 cannot have $k = 1$ since that would imply $f(s^1) = s^1 \in R$ and $f(s^1, \dots, s^k) = f(s^1) = u \notin R$.
 975 It remains to show that $k \neq 2$. So assume for a contradiction that $k = 2$. Since s^1 and s^2
 976 are distinct, there is a position $\ell \in [2]$ such that $s_\ell^1 \neq s_\ell^2$. Assume without loss of generality
 977 that $s_\ell^1 = 1$ while $s_\ell^2 = 0$. Since $f(s^1, \dots, s^k) = f(s^1, s^2) = \alpha_1 s^1 + \alpha_2 s^2 = u \in \{0, 1\}^2 \setminus R$,
 978 we find $u_\ell = \alpha_1 s_\ell^1 + \alpha_2 s_\ell^2 = \alpha_1 \cdot 1 + \alpha_2 \cdot 0 \in \{0, 1\}$. Since α_1 is a nonzero integer, we must
 979 have $\alpha_1 = 1$. But since $\alpha_1 + \alpha_2 = 1$ by definition of a balanced operation, this implies $\alpha_2 = 0$,
 980 contradicting that f is a minimum-arity balanced operation that does not preserve R . \square

981 The previous two claims show that there are at least three distinct tuples in $R \subseteq \{0, 1\}^2$.
 982 Since $u \in \{0, 1\}^2 \setminus R$ it follows that $|R| = 3$. Hence R and 2-OR are both Boolean relations of
 983 arity two that each have three tuples. To cone-define one from the other, one may easily verify
 984 that it suffices to use the tuple (y_1, y_2) , where $y_i = x_i$ if $u_i = 0$ and $y_i = \neg x_i$ otherwise. \blacktriangleleft

985 In order to prove Theorem 6.3, we first present some additional lemmas and definitions.

986 **► Lemma A.20.** *Suppose that $U \subseteq \{0, 1\}^n$ is a Boolean relation, and that there exist an
 integer a and a natural number $m > 1$ such that, for each $u \in U$, it holds that*

$$\text{weight}(u) \equiv a \pmod{m}.$$

987 *Then, if w is a witness for U , it holds that $\text{weight}(w) \equiv a \pmod{m}$.*

988 Let $U \subseteq \{0, 1\}^n$ be a relation. We say that $w \in \{0, 1\}^n$ is a *witness* for U if $w \notin R$, and
 989 there exists a balanced operation $f : \{0, 1\}^k \rightarrow \{0, 1\}$ and tuples $t^1, \dots, t^k \in U$ such that
 990 $w = f(t^1, \dots, t^k)$. Observe that U is not balanced if and only if there exists a witness for U .
 991 We will view tuples of arity n as maps $f : [n] \rightarrow \{0, 1\}$, via the natural correspondence where
 992 such a map f represents the tuple $(f(1), \dots, f(n))$. We freely interchange between these
 993 two representations of tuples. We say that $f : S \rightarrow \{0, 1\}$ is a *no-good* of U when: $S \subseteq [n]$;
 994 each extension $g : [n] \rightarrow \{0, 1\}$ of f is not an element of U ; and there exists an extension
 995 $h : [n] \rightarrow \{0, 1\}$ of f that is a witness for U . We say that $f : S \rightarrow \{0, 1\}$ is a *min-no-good* of
 996 f is a no-good, but no proper restriction of f is a no-good. Observe that the following are
 997 equivalent, for a relation: the relation is not balanced; it has a witness; it has a no-good; it
 998 has a min-no-good.

999 When $U \subseteq \{0, 1\}^n$ is a relation and $S \subseteq [n]$, let $s_1 < \dots < s_m$ denote the elements of S ;
 then, we use $U \upharpoonright S$ to denote the relation $\{(h(s_1), \dots, h(s_m)) \mid h \in U\}$.

1000 **► Proposition A.21.** *Let $U \subseteq \{0, 1\}^n$ be a relation, let $S \subseteq [n]$, and suppose that $f : S \rightarrow$
 1001 $\{0, 1\}$ is a min-no-good of U . Then f is a min-no-good of $U \upharpoonright S$.*

1002 **Proof.** Observe that f is not in $U \upharpoonright S$; since f has an extension that is a witness for U , it
 1003 follows that f is a witness for $U \upharpoonright S$. Thus, f is a no-good of $U \upharpoonright S$. In order to obtain that
 1004 f is a min-no-good of $U \upharpoonright S$, it suffices to establish that, for any restriction $f^- : S^- \rightarrow \{0, 1\}$
 1005 of f , it holds that f^- is a no-good of U if and only if f^- is a no-good of $U \upharpoonright S$. This

1006 follows from what we have established concerning f and the following fact: all extensions
 1007 $h : S \rightarrow \{0, 1\}$ of f^- are not in $U \upharpoonright S$ if and only if all extensions $h' : [n] \rightarrow \{0, 1\}$ of f^-
 1008 are not in U . ◀

Proof. Since w is a witness for U , there exist tuples $t^1 = (t_1^1, \dots, t_n^1), \dots, t^k = (t_1^k, \dots, t_n^k)$ and
 a balanced operation $f : \{0, 1\}^k \rightarrow \{0, 1\}$ such that $f(t^1, \dots, t^k) = w$. Let $\alpha_1, \dots, \alpha_k$ be the
 coefficients of f . From $f(t^1, \dots, t^k) = w$, we obtain that $\alpha_1 \text{weight}(t^1) + \dots + \alpha_k \text{weight}(t^k) =$
 $\text{weight}(w)$. As

$$\alpha_1 \text{weight}(t^1) + \dots + \alpha_k \text{weight}(t^k) \equiv \alpha_1 a + \dots + \alpha_k a \equiv a \pmod{m},$$

1009 the result follows. ◀

1010 **Proof of Theorem 6.3.** Let $f : S \rightarrow \{0, 1\}$ be a min-no-good of U .

1011 It cannot hold that $|S| = 0$, since then U would be empty and hence preserved by all
 1012 balanced operations. It also cannot hold that $|S| = 1$, since then f would be a min-no-good
 1013 of $U \upharpoonright S$ (by Proposition A.21), which is not possible since $U \upharpoonright S$ would have arity 1 and
 1014 hence would be preserved by all balanced operations (by Observation 6.1).

1015 For the remaining cases, by replacing U with a relation that is interdefinable with it, we
 1016 may assume that $f : S \rightarrow \{0, 1\}$ maps each $s \in S$ to 0.

1017 Suppose that $|S| = 2$, and assume for the sake of notation that $S = \{1, 2\}$ (this can be
 1018 obtained by replacing U with a relation that is interdefinable with it). By Proposition A.21,
 1019 f is a min-no-good of $U \upharpoonright S$. By Theorem 6.2, we obtain that $U \upharpoonright S$ contains all tuples
 1020 other than f , that is, we have $\{(0, 1), (1, 0), (1, 1)\} = U \upharpoonright S$. It follows that there exists
 1021 a *realization*, where we define a *realization* to be a tuple $(a_1, a_2, a_3) \in \{0, 1\}^3$ such that
 1022 $(0, 1, a_1), (1, 0, a_2), (1, 1, a_3) \in U$. Let us refer to $(0, 0, 1)$ and $(1, 1, 0)$ as *bad* tuples, and to all
 1023 other arity 3 tuples as *good* tuples. We claim that there exists a realization that is a good
 1024 tuple; this suffices, for then the 2-OR relation is cone-definable from U via a tuple of the
 1025 form (x_1, x_2, y) where $y \in \{0, 1, x_1, x_2, \neg x_1, \neg x_2\}$. In particular, let us state explicitly how y
 1026 can be chosen as a function of a realization that is a good tuple: chose $y = 0$ for $(0, 0, 0)$;
 1027 $y = 1$ for $(1, 1, 1)$; $y = x_1$ for $(0, 1, 1)$; $y = x_2$ for $(1, 0, 1)$; $y = \neg x_1$ for $(1, 0, 0)$; and, $y = \neg x_2$
 1028 for $(0, 1, 0)$. We prove the claim by contradiction. If there exists no realization that is a good
 1029 tuple, every realization is a bad tuple; moreover, there is a unique realization, for if there
 1030 were more than one, there would exist a realization that was a good tuple. We may assume
 1031 (up to interdefinability of U) that the unique realization is $(1, 1, 0)$. Then, U is the relation
 1032 $\{(0, 1, 1), (1, 0, 1), (1, 1, 0)\}$ containing exactly the weight 2 tuples; applying Lemma A.20 to
 1033 U with $a = 2$ and $m = 3$, we obtain that for any witness w for U , it holds that $\text{weight}(w) \equiv 2$
 1034 $\pmod{3}$. This implies that f has no extension w' that is a witness, since any such extension
 1035 must have $\text{weight}(w')$ equal to 0 or 1; we have thus contradicted that f is a no-good of U .

1036 Suppose that $|S| = 3$. Since f is both a min-no-good and a witness, it follows that
 1037 each of the weight 1 tuples $(1, 0, 0), (0, 1, 0), (0, 0, 1)$ is contained in U . We claim that U
 1038 contains a weight 2 tuple; if not, then U would contain only weight 1 and weight 3 tuples,
 1039 and by invoking Lemma A.20 with $a = 1$ and $m = 2$, we would obtain that $\text{weight}(f) \equiv 1$
 1040 $\pmod{2}$, a contradiction. Assume for the sake of notation that U contains the weight 2 tuple
 1041 $(0, 1, 1)$. Then U cone-defines the 2-OR relation via the tuple $(0, x_1, x_2)$, since $(0, 0, 0) \notin R$
 1042 and $(0, 1, 0), (0, 0, 1), (0, 1, 1) \in R$. ◀