

# Filling crosswords is very hard

Laurent Gourvès, Ararat Harutyunyan, Michael Lampis, Nikolaos Melissinos

### ▶ To cite this version:

Laurent Gourvès, Ararat Harutyunyan, Michael Lampis, Nikolaos Melissinos. Filling crosswords is very hard. Theoretical Computer Science, 2024, 982, pp.114275. 10.1016/J.TCS.2023.114275 . hal-04487531

## HAL Id: hal-04487531 https://hal.science/hal-04487531

Submitted on 3 Mar 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

### Filling Crosswords is Very Hard<sup>\*</sup>

Laurent Gourvès<sup>a</sup>, Ararat Harutyunyan<sup>a</sup>, Michael Lampis<sup>a,1,2</sup>, Nikolaos Melissinos<sup>a,3</sup>

<sup>a</sup> Université Paris-Dauphine, Université PSL, CNRS, LAMSADE, 75016, Paris, France

#### Abstract

We revisit a classical crossword filling puzzle which already appeared in Garey& Jonhson's book. We are given a grid with n vertical and horizontal slots and a dictionary with m words. We are asked to place words from the dictionary in the slots so that shared cells are consistent. We attempt to pinpoint the source of intractability of this problem by carefully taking into account the structure of the grid graph, which contains a vertex for each slot and an edge if two slots intersect. Our main approach is to consider the case where this graph has a tree-like structure. Unfortunately, if we impose the common rule that words cannot be reused, we discover that the problem remains NP-hard under very severe structural restrictions, namely, if the grid graph is a union of stars and the alphabet has size 2, or the grid graph is a matching (so the crossword is a collection of disjoint crosses) and the alphabet has size 3. The problem does become slightly more tractable if word reuse is allowed, as we obtain an  $m^{tw}$ algorithm in this case, where tw is the treewidth of the grid graph. However, even in this case, we show that our algorithm cannot be improved to obtain fixed-parameter tractability. More strongly, we show that under the ETH the problem cannot be solved in time  $m^{o(k)}$ , where k is the number of horizontal slots of the instance (which trivially bounds tw).

Motivated by these mostly negative results, we also consider the much more restricted case where the problem is parameterized by the number of slots n. Here, we show that the problem does become FPT (if the alphabet has constant size), but the parameter dependence is exponential in  $n^2$ . We show that this dependence is also justified: the existence of an algorithm with running time  $2^{o(n^2)}$ , even for binary alphabet, would contradict the randomized ETH. After that, we consider an optimization version of the problem, where we seek to place as many words on the grid as possible. Here it is easy to obtain a

ararat.harutyunyan@dauphine.fr (Ararat Harutyunyan), michail.lampis@dauphine.fr

<sup>\*</sup>A preliminary version of this work appears in the proceedings of ISAACS 2021 [8]. Email addresses: laurent.gourves@dauphine.fr (Laurent Gourvès),

<sup>(</sup>Michael Lampis ), nikolaos.melissinos@dauphine.eu (Nikolaos Melissinos)

<sup>&</sup>lt;sup>1</sup>Michael Lampis was partially supported by ANR project ANR-21-CE48-0022 (S-EX-AP-PE-AL).

<sup>&</sup>lt;sup>2</sup>ORCID:0000-0002-5791-0887

<sup>&</sup>lt;sup>3</sup>ORCID:0000-0002-0864-9803

 $\frac{1}{2}$ -approximation, even on weighted instances, simply by considering only horizontal or only vertical slots. We show that this trivial algorithm is also likely to be optimal, as obtaining a better approximation ratio in polynomial time would contradict the Unique Games Conjecture. The latter two results apply whether word reuse is allowed or not.

Finally, we present some special cases where the problem is decidable in polynomial time. In particular, we present three reductions, one to 2-SAT, the second to MAXIMUM MATCHING and the third to EXACT MATCHING.

Keywords: Crossword Puzzle, Treewidth, ETH

#### 1. Introduction

Crossword puzzles are one-player games where the goal is to fill a (traditionally two-dimensional) grid with words. Since their first appearance more than 100 years ago, crossword puzzles have rapidly become popular. Nowadays, they can be found in many newspapers and magazines around the world like the *New York Times* in the USA, or *Le Figaro* in France. Besides their obvious recreational interest, crossword puzzles are valued tools in education [2] and medicine. In particular, crossword puzzles participation seems to delay onset of accelerated memory decline [21]. They are also helpful for developing and testing computational techniques; see for example [23]. In fact, both the design and the completion of a puzzle are challenging. In this article, we are interested in the task of solving a specific type of crossword puzzle.

There are different kinds of crossword puzzles. In the most famous ones, some clues are given together with the place where the answers should be located. A solution contains words that must be consistent with the given clues, and the intersecting pairs of words are constrained to agree on the letter they share. *Fill-in* crossword puzzles do not come with clues. Given a list of words and a grid in which some slots are identified, the objective is to fill all the slots with the given words. The list of words is typically succinct and provided explicitly.

In a variant of fill-in crossword puzzle currently proposed in a French TV magazine [16], one has to find up to 14 words and place them in a grid (the grid is the same for every instance, see Figure 1 for an illustration). The words are not explicitly listed but they must be *valid* (for instance, belong to the French language). In an instance of the game, some specified letters have a positive weight; the other letters have weight zero. The objective is to find a solution whose weight – defined as the total sum of the letters written in the grid – is at least a given threshold.

The present work deals with a theoretical study of this fill-in crossword puzzle (the grid is not limited to the one of Figure 1). We are mainly interested in two problems: Can the grid be entirely completed? How can the weight of a solution be maximized? Hereafter, these problems are called CROSSWORD PUZZLE DECISION and CROSSWORD PUZZLE OPTIMIZATION (CP-DEC and CP-OPT in short), respectively.

Figure 1: Place valid words in this grid. In a possible instance, letters S, U, I, V, R, E, and T have weight 7, 5, 4, 2, 6, 1, and 3, respectively. Any other letter has null weight. Try to obtain at least 330 points.

CP-DEC is not new; see GP14 in [6]. The proof of NP-completeness is credited to a personal communication with Lewis and Papadimitriou. Thereafter, an alternative NP-completeness proof appeared in [4] (see also [14]). Other articles on crossword puzzles exist and they are mostly empirically validated techniques coming from Artificial Intelligence and Machine Learning; see for example [7, 17, 15, 1, 23, 22] and references therein.

**Our Results** Our goal in this paper is to pinpoint the relevant structural parameters that make filling crossword puzzles intractable. We begin by examining the structure of the given grid. It is natural to think that, if the structure of the grid is tree-like, then the problem should become easier, as the vast majority of problems are tractable on graphs of small treewidth. We only partially confirm this intuition: by taking into account the structure of a graph that encodes the intersections between slots (the grid graph) we show in Section 3 that CP-OPT can be solved in polynomial time on instances of constant treewidth. However, our algorithm is not fixed-parameter tractable and, as we show, this cannot be avoided, even if one considers the much more restricted case where the problem is parameterized by the number of horizontal slots, which trivially bounds the grid graph's treewidth (Theorem 3). More devastatingly, we show that if we also impose the natural rule that words cannot be reused, the problem already becomes NP-hard when the grid graph is a matching for alphabets of size 3 (Theorem 5), or a union of stars for a binary alphabet (Theorem 4). Hence, a tree-like structure does not seem to be of much help in rendering crosswords tractable.

We then go on to consider CP-OPT parameterized by the total number of slots n. This is arguably a very natural parameterization of the problem, as in real-life crosswords, the size of the grid can be expected to be significantly smaller than the size of the dictionary. We show that in this case the problem does become fixed-parameter tractable (Corollary 9), but the running time of our algorithm is exponential in  $n^2$ . Our main result is to show that this disappointing dependence is likely to be best possible: even for a binary alphabet, an algorithm solving CP-DEC in time  $2^{o(n^2)}$  would contradict the randomized ETH (Theorem 13). Note that all our positive results up to this point work for the more general CP-OPT, while our hardness results apply to CP-DEC.

Afterwards, in Section 5 we consider the approximability of CP-OPT. Here,

it is easy to obtain a  $\frac{1}{2}$ -approximation by only considering horizontal or vertical slots. We are only able to slightly improve upon this, giving a polynomial-time algorithm with ratio  $\frac{1}{2} + O(\frac{1}{n})$ . Our main result in this direction is to show that this is essentially best possible: obtaining an algorithm with ratio  $\frac{1}{2} + \epsilon$  would falsify the Unique Games Conjecture (Theorem 16).

Before concluding, we explore in Section 6 the cases where CP-DEC can be resolved in polynomial time. We propose reductions from CP-DEC to some well-known problems that belong to P.

#### 2. Problem Statement, Encoding and Preliminaries

We are given a dictionary  $\mathcal{D} = \{d_1, \ldots, d_m\}$  whose words are constructed on an alphabet  $\mathcal{L} = \{l_1, \ldots, l_\ell\}$  of constant size (i.e.,  $\ell = O(1)$ ), and a twodimensional grid consisting of horizontal and vertical slots. A slot is composed of consecutive cells. Horizontal slots do not intersect each other; the same goes for vertical slots. However horizontal slots can intersect vertical slots. A cell is *shared* if it lies at the intersection of two slots. Unless specifically stated, n, m and  $\ell$  denote the total number of slots, the size of  $\mathcal{D}$ , and the size of  $\mathcal{L}$ , respectively.

In a feasible solution, each slot S receives either a word of  $\mathcal{D}$  of length |S|, or nothing (we sometimes say that a slot receiving nothing gets an *empty word*). Each cell gets at most one letter, and the words assigned to two intersecting slots must agree on the letter placed in the shared cell. All filled horizontal slots get words written from left to right (across) while all vertical slots get words written from top to bottom (down).

There is a weight function  $w : \mathcal{L} \to \mathbb{N}$ . The weight of a solution is the total sum of the weights of the letters placed in the grid. Observe that the weight of a solution is smaller than the total sum of the weights of its words because, in the former, the letters of the shared cells are counted only once.

The two main problems studied in this article are the following. Given a grid, a dictionary  $\mathcal{D}$  on alphabet  $\mathcal{L}$ , and a weight function  $w : \mathcal{L} \to \mathbb{N}$ , the objective of CROSSWORD PUZZLE OPTIMIZATION (CP-OPT in short) is to find a feasible solution of maximum weight. Given a grid and a dictionary  $\mathcal{D}$  on alphabet  $\mathcal{L}$ , the question posed by CROSSWORD PUZZLE DECISION (CP-DEC in short) is whether the grid can be completely filled or not?

Two cases will be considered: whether each word is used at most once, or if words can be assigned multiple times. In this article, we will sometimes suppose that some cells are pre-filled with some elements of  $\mathcal{L}$ . In this case, a solution is feasible if it is consistent with the pre-filled cells. Below we propose a first result when all the shared cells are pre-filled.

**Proposition 1.** CP-DEC and CP-OPT can be solved in polynomial time if all the shared cells in the grid are pre-filled, whether word reuse is allowed or not.

*Proof.* If word reuse is allowed, then for each combination of letters placed in these cells, we greedily fill out the rest of each slot with the maximum value

word that can still be placed there. This is guaranteed to produce the optimal solution. On the other hand, if word reuse is not allowed, we construct a bipartite graph, with elements of  $\mathcal{D}$  on one side and the slots on the other, and place an edge between a word and a slot if the word can still be placed in the slot. If we give each edge weight equal to the value of its incident word reduced by the weight of the letters imposed by the shared cells of the slot, then an optimal solution corresponds to a maximum weight matching.

One can associate a bipartite graph, hereafter called the *grid graph*, with each grid: each slot is a vertex and two vertices share an edge if the corresponding slots overlap. The grid (and then, the grid graph) is not necessarily connected.

Let us also note that so far we have been a bit vague about the encoding of the problem. Concretely, we could use a simple representation which lists for each slot the coordinates of its first cell, its size, and whether the slot is horizontal or vertical; and then supplies a list of all words in the dictionary and an encoding of the weight function. Such a representation would allow us to perform all the basic operations needed by our algorithms in polynomial time, such as deciding if it is possible to place a word d in a slot S, and which letter would then be placed in any particular cell of S. However, one drawback of this encoding is that its size may not be polynomially bounded in n + m, as some words may be exponentially long. We can work around this difficulty by using a more succinct representation: we are given the same information as above regarding the n slots; for each word we are given its total weight; and for each slot S and word d, we are told whether d fits exactly in S, and if yes, which letters are placed in the cells of S which are shared with other slots. Since the number of shared cells is  $O(n^2)$  this representation is polynomial in n+mand it is not hard to see that we are still able to perform any reasonable basic operation in polynomial time and that we can transform an instance given in the simple representation to this more succinct form. Hence, in the remainder, we will always assume that the size of the input is polynomially bounded in n+m.

We will rely on the Exponential Time Hypothesis (ETH) of Impagliazzo, Paturi, and Zane [11], which states the following:

**Conjecture 1.** Exponential Time Hypothesis: there exists an  $\epsilon > 0$ , such that 3-SAT on instances with n variables and m clauses cannot be solved in time  $2^{\epsilon(n+m)}$ .

Note that it is common to use the slightly weaker formulation which states the ETH as the assumption that 3-SAT cannot be solved in time  $2^{o(n+m)}$ . This is known to imply that k-INDEPENDENT SET cannot be solved in time  $n^{o(k)}[3]$ . We use this fact in Theorem 3. In Section 4 we will rely on the randomized version of the ETH, which has the same statement as Conjecture 1 but for randomized algorithms with expected running time  $2^{\epsilon(n+m)}$ .

#### 3. When the Grid Graph is Tree-like

In this section we are considering instances of CP-DEC and CP-OPT where the grid graph is similar to a tree. First, we give an algorithm for both problems in cases where the grid graph has bounded treewidth and we are allowed to reuse words. We show that this algorithm is essentially optimal. Then, we show that CP-DEC and CP-OPT are much harder to deal with, in the case where we are not allowed to reuse words, by proving that the problems are NP-hard even for instances where the grid graph is just a matching. For the instances such that CP-DEC is NP-hard, we know that CP-OPT is NP-hard. That happens because we can assume that all the letters have weight equal to 1. Hence, a solution for CP-DEC is an optimal solution for CP-OPT.

#### 3.1. Word Reuse

We propose a dynamic programming algorithm for CP-OPT and hence also for CP-DEC. Note that it can be extended to the case where some cells of the instance are pre-filled.

**Theorem 2.** If we allow word reuse, then CP-OPT can be solved in time  $(m + 1)^{tw}(n+m)^{O(1)}$  on inputs where tw is the treewidth of the grid graph.

**Proof.** As the techniques we are going to use are standard we are sketching some details. For more details on tree decomposition (definition and terminology) see [3, Chap. 7]. Assuming that we have a rooted nice tree decomposition of the grid graph, we are going to perform dynamic programming on the nodes of this tree decomposition. For a node  $B_t$  of the given tree decomposition of the grid graph we denote by  $B_t^{\downarrow}$  the set of vertices of the grid graph that appears in the nodes of the subtree with  $B_t$  as a root. Since each vertex of the grid graph and its corresponding slot. In particular, we say that a solution  $\sigma$  assigns words to the vertices of the grid graph, and  $\sigma(v)$  denotes the word assigned to v.

For each node  $B_t$  of the tree decomposition we are going to keep all the triplets  $(\sigma, W, W_t)$  such that:

- $\sigma$  is an assignment of words to the vertices of  $B_t$ ;
- W is the weight of  $\sigma$  restricted to the vertices appearing in  $B_t$ ;
- and  $W_m$  is the maximum weight, restricted to the vertices appearing in  $B_t^{\downarrow}$ , of an assignment consistent with  $\sigma$ .

In order to create all the possible triplets for all the nodes of the tree decomposition we are going to explore the nodes from leaves to the root. Therefore, each time we visit a node we assume that we have already created the triplets for all its children. Let us explain how we deal with the different types of nodes.

In the Leaf nodes we have no vertices so we keep an empty assignment ( $\sigma$  does not assign any word) and the weights W and  $W_m$  are equal to 0.

For an Introduce node  $B_t$  we need to take in consideration its child node. Assume that u is the introduced vertex; for each triplet  $(\sigma, W, W_m)$  of the child node we are going to create all the triplets  $(\sigma', W', W'_m)$  for the new node as follows. First we find all the words  $d \in \mathcal{D}$  that fit in the corresponding slot of u and respect the assignment  $\sigma$  (i.e., if there are cells that are already filled under  $\sigma$  and d uses these cells then it must have the same letters). We create one triplet  $(\sigma', W', W'_m)$  for each such a d as follows:

- We set  $\sigma'(u) := d$  and  $\sigma'(v) := \sigma(v)$  for all  $v \in B_t \setminus \{u\}$ .
- We can easily calculate the total weight, W', of the words in  $B_t$  where the shared letters are counted only once under the assignment  $\sigma'$ .
- For the maximum weight  $W'_m$  we know that it is increased by the same amount as W; so we set  $W'_m = W_m + W' W$ .

Observe that we do not need to consider the intersection with slots whose vertices appear in  $B_t^{\downarrow} \setminus B_t$  as each node of a tree decomposition is a cut set.

Finally, we need to take in consideration that we can leave a slot empty. For this case we create a new word  $d_*$  which, we assume that, fits in all slots and  $d_*$ has weight 0. Because the empty word has weight 0, W' and  $W'_m$  are identical to W and  $W_m$  so for each triplet of the child node, we only need to extend  $\sigma$ by assigning  $d_*$  to u. In the case we assign the empty word somewhere we will consider that the cells of this slot are empty unless another word  $d \neq d_*$  uses them.

For the Forget nodes we need to restrict the assignments of the child node to the vertex set of the Forget node, as it has been reduced by one vertex (the forgotten vertex), and reduce the weight W (which we can calculate easily). The maximum weight is not changed by the deletion.

However, if we restrict the assignments we may end up with several triplets  $(\sigma, W, W_m)$  with identical assignments  $\sigma$ . In that case we are keeping only the triplet with maximum  $W_m$ . Observe that we are allowed to keep only triplets with the maximum  $W_m$  because each node of a tree decomposition is a cut set so the same holds for the Forget nodes. Specifically, the vertices that appear in the nodes higher than a Forget node  $B_t$  of the tree decomposition do not have edges incident to vertices in  $B_t^{\downarrow} \setminus B_t$  so we only care for the assignment in  $B_t$ .

Finally, we need to consider the Join nodes. Each Join node has exactly two children. For each possible assignment  $\sigma$  on the vertices of this Join node, we create a triplet iff this  $\sigma$  appears in a triplet of both children of the Join node.

Because W is related only to the assignment  $\sigma$ , it is easy to see that it will be the same as in the children of the Join node. So we need to find the maximum weight  $W_m$ . Observe that between the vertices that appear in the subtrees of two children of a Join node there are no edges except those incident to the vertices of the Join node. Therefore, we can calculate the maximum weight  $W_m$ as follows: first we consider the maximum weight of each child of the Join node reduced by W, we add all these weights and, in the end, we add again the W. It is easy to see that this way we consider the weight of the cells appearing in each subtree without those of the slots of the Join node and we add the weight of the words assigned to the vertices of the Join node in the end.

For the running time we need to observe that the number of nodes of a nice tree decomposition is  $O(\mathsf{tw} \cdot n)$  and all the other calculations are polynomial in n + m so we only need to consider the different assignments for each node. Because for each vertex we have  $|\mathcal{D}| + 1$  choices, the number of different assignments for a node is at most  $(|\mathcal{D}| + 1)^{\mathsf{tw}+1}$ .

It seems that the algorithm we propose for CP-DEC is essentially optimal, even if we consider a much more restricted case.

**Theorem 3.** CP-DEC with word reuse is W[1]-hard parameterized by the number of horizontal slots of the grid, even for alphabets with two letters. Furthermore, under the ETH, no algorithm can solve this problem in time  $m^{o(k)}$ , where k is the number of horizontal slots.

*Proof.* We perform a reduction from k-INDEPENDENT SET, where we are given a graph G = (V, E) with |V| vertices and |E| edges and are looking for an independent set of size k. This problem is well-known to be W[1]-hard and not solvable in  $|V|^{o(k)}$  time under the ETH [3]. We assume without loss of generality that  $|E| \neq k$ . Furthermore, we can safely assume that G has no isolated vertices.

We first describe the grid of our construction which fits within an area of 2k - 1 lines and 2|E| - 1 columns. We construct:

- 1. k horizontal slots, each of length 2|E| 1 (so each of these slots is as long horizontally as the whole grid). We place these slots in the unique way so that no two of these slots are in consecutive lines. We number these horizontal slots  $1, \ldots, k$  from top to bottom.
- 2. |E| vertical slots, each of length 2k-1 (so each of these slots is long enough to cover the grid top to bottom). We place these slots in the unique way so that no two of them are in consecutive columns. We number them  $1, \ldots, |E|$  from left to right.

Before we describe the dictionary, let us give some intuition about the grid. The main idea is that in the k horizontal slots we will place k words that signify which vertices we selected from the original graph. Each vertical slot represents an edge of E, and we will be able to place a word in it if and only if we have not placed words representing two of its endpoints in the horizontal slots.

Our alphabet has two letters, say 0, 1. In the remainder, we assume that the edges of the original graph are numbered, that is,  $E = \{e_1, \ldots, e_{|E|}\}$ . The dictionary is as follows:

1. For each vertex v we construct a word of length 2|E| - 1. For each  $i \in \{1, \ldots, |E|\}$ , if the edge  $e_i$  is incident on v, then the letter at position 2i-1 of the word representing v is 1. All other letters of the word representing v are 0. Observe that this means that if  $e_i$  is incident on v and we place the word representing v on a horizontal slot, the letter i will appear on the *i*-th vertical slot. Furthermore, the word representing v has a number of 1s equal to the degree of v.

2. We construct k + 1 words of length 2k - 1. One of them is simply  $0^{2k-1}$ . The remaining are  $0^{2j-2}10^{2k-2j}$ , for  $j \in \{1, \ldots, k\}$ , that is, the words formed by placing a 1 in an odd-numbered position and 0s everywhere else. Observe that if we place one of these k words on a vertical slot, a 1 will be placed on exactly one horizontal slot.



Figure 2: In the left, we have the given graph where we want to find an independent set of size 4. In the right, we have the created grid. ( $\alpha$ ): We have placed the word that represents  $u_1$  (in the first horizontal slot) and the word that represents  $v_1$  (in the second horizontal slot). Because both forces 1s in the vertical slot that represents the edge  $u_1v_1$ , we cannot fill this slot. This happens because  $S = \{u_1, v_1\}$  is not an independent set of the starting graph. ( $\beta$ ): We have placed the words that represent  $v_1, \ldots, v_4$  into the horizontal slots. Now, we can fill all vertical slots as we have a maximum of one 1 per vertical slot. This happens because  $S = \{v_1, \ldots, v_4\}$  is an independent set of the starting graph.

This completes the construction. We now observe that we have constructed exactly k horizontal slots, therefore, if the reduction preserves the answer, the hardness results for k-INDEPENDENT SET transfer to our problem, since we preserve the value of the parameter.

We claim that if there exists an independent set of size k in G, then it is possible to fill the grid. Indeed, take such a set S and for each  $v \in S$  we place the word representing v in a horizontal slot. Consider the *i*-th vertical slot. We will place in this slot one of the k + 1 words of length 2k - 1. We claim that the vertical slot at this moment contains the letter 1 at most once, and if 1 appears it must be at an odd position (since these are the positions shared with the horizontal slots). If this is true, clearly there is a word we can place. To see that the claim is true, recall that since S is an independent set of k distinct vertices, there exists at most one vertex in S incident on  $e_i$ . For the converse direction, recall that  $|E| \neq k$ . This implies that if there is a way to fill out the whole grid, then words representing vertices must go into horizontal slots and words of length 2k - 1 must go into vertical slots. By looking at the words that have been placed in the horizontal slots we obtain a collection of k (not necessarily distinct) vertices of G. We will prove that these vertices must actually be an independent set of size exactly k. To see this, consider the *i*-th vertical slot. If our collection of vertices contained two vertices incident on  $e_i$ , it would have been impossible to fill out the *i*-th vertical slot, since we would need a word with two 1s. Observe that the same argument rules out the possibility that our collection contains the same vertex v twice, as the column corresponding to any edge  $e_i$  incident on v would have been impossible to fill.

Notice that the number of horizontal slots is an upper bound on the size of a minimum vertex cover of the grid graph which also bounds the treewidth of the grid graph. This shows that the algorithm of Theorem 2 for CP-DEC is essentially optimal.

#### 3.2. No Word Reuse

If a word cannot be reused, then CP-DEC looks more challenging. Indeed, in the following theorem we prove that if reusing words is not allowed, then the problem becomes NP-hard even if the grid graph is acyclic and the alphabet size is 2. (Note that if the alphabet size is 1, the problem is trivial, independent of the structure of the graph.)

**Theorem 4.** CP-DEC is NP-hard, even for instances where all of the following restrictions apply: (i) the grid graph is a union of stars (ii) the alphabet contains only two letters (iii) words cannot be reused.

*Proof.* We show a reduction from 3-PARTITION. Recall that in 3-PARTITION we are given a collection of 3n distinct positive integers  $x_1, \ldots, x_{3n}$  and are asked if it is possible to partition these integers into n sets of three integers (triples), such that all triples have the same sum. This problem has long been known to be strongly NP-hard [6] and NP-hardness when the integers are distinct was shown by Hulett et al. [10]. We can assume that  $\sum_{i=1}^{3n} x_i = nB$  and that if a partition exists each triple has sum B. Furthermore, we can assume without loss of generality that  $x_i > 6n$  for all  $i \in \{1, \ldots, 3n\}$  (otherwise, we can simply add 6n to all numbers and adjust B accordingly without changing the answer).

Given an instance of 3-PARTITION as above, we construct a crossword instance as follows. First, the alphabet only contains two letters, say the letters \* and !. To construct our dictionary we do the following:

1. For each  $i \in \{1, \ldots, 3n\}$ , we add to the dictionary one word of length  $x_i$  that begins with ! and n - 1 words of length  $x_i$  that begin with \*. The remaining letters of these words are chosen in an arbitrary way so that all words remain distinct.

2. For each  $i, j, k \in \{1, ..., 3n\}$  with i < j < k we check if  $x_i + x_j + x_k = B$ . If this is the case, we add to the dictionary the word  $*^{2i-2}!*^{2j-2i-1}!*^{2k-2j-1}!*^{6n-2k}$ . In other words, we constructed a word that has \* everywhere except in positions 2i - 1, 2j - 1, and 2k - 1. The length of this word is 6n - 1. Let f be the number of words added to the dictionary in this step. We have  $f \leq {3 \choose 3} = O(n^3)$ .

We now also need to specify our grid. We first construct f horizontal slots, each of length 6n - 1. Among these f slots, we select n, which we call the "interesting" horizontal slots. For each interesting horizontal slot, we construct 3n vertical slots, such that the *i*-th of these slots has length  $x_i$  and its first cell is the cell in position 2i - 1 of the interesting horizontal slot. This completes the construction, which can clearly be carried out in polynomial time. Observe that the first two promised restrictions are satisfied as we have an alphabet with two letters and each vertical slot intersects at most one horizontal slot (so the grid graph is a union of stars).

We claim that if there exists a partition of the original instance, then we can place all the words of the dictionary on the grid. Indeed, for each  $i, j, k \in \{1, \ldots, 3n\}$  such that  $\{x_i, x_j, x_k\}$  is one of the triples of the partition, we have constructed a word of length 6n - 1 corresponding to the triple (i, j, k), because  $x_i + x_j + x_k = B$ . We place each of these n words on an interesting horizontal slot and we place the remaining words of length 6n - 1 on the non-interesting horizontal slots. Now, for every  $i \in \{1, \ldots, 3n\}$  we have constructed n words, one starting with ! and n - 1 starting with \*. We observe that among the interesting horizontal slots, there is one that contains the letter ! at position 2i - 1 (the one corresponding to the triple containing  $x_i$  in the partition) and n - 1 containing the letter \* at position 2i - 1. By construction, the vertical slots that begin in these positions have length  $x_i$ . Therefore, we can place all n words corresponding to  $x_i$  on these vertical slots. Proceeding in this way we fill the whole grid, fulfilling the third condition.

For the converse direction, suppose that there is a way to fill the whole grid. Then, vertical slots must contain words that were constructed in the second step and represent integers  $x_i$ , while horizontal slots must contain words constructed in the first step (this is a consequence of the fact that  $x_i > 6n$  for all  $i \in \{1, \ldots, 3n\}$ ). We consider the n interesting horizontal slots. Each such slot contains a word that represents a triple (i, j, k) with  $x_i + x_j + x_k = B$ . We therefore collect these n triples and attempt to construct a partition from them. To do this, we must prove that each  $x_i$  must belong to exactly one of these triples. However, recall that we have exactly n words of length  $x_i$  (since all integers of our instance are distinct) and exactly n vertical slots of this length. We conclude that exactly one vertical slot must have ! as its first letter, therefore  $x_i$  appears in exactly one triple and we have a proper partition.

Actually, the problem remains NP-hard even in the case where the grid graph is a matching and the alphabet contains three letters. This is proved for grid graphs composed of  $\mathcal{T}$ s, where a  $\mathcal{T}$  is a horizontal slot solely intersected by the first cell of a vertical slot.

**Theorem 5.** CP-DEC is NP-hard, even for instances where all of the following restrictions apply: (i) each word can be used only once (ii) the grid is consisted only by  $\mathcal{T}s$  and (iii) the alphabet contains only three letters.

In order to prove this theorem we need first to define a restricted version of EXACTLY-1 3-SAT.

**Definition 1 (Restricted Exactly 1 (3,2)-SAT).** Assume that  $\phi$  is a CNF formula where each clause has either three or two literals and each variable appears at most three times. We want to determine whether there exists a satisfying assignment so that each clause has exactly one true literal.

Lemma 6. The RESTRICTED EXACTLY-1 (3,2)-SAT is NP-complete.

*Proof.* We show a reduction from EXACTLY-1 3-SAT which is known to be NP-complete [6] (LO4, ONE-IN-THREE 3SAT).

Let  $I = (\phi, X)$  be an instance of EXACTLY-1 3-SAT with |X| = n variables and m clauses. If there exists a variable x with k > 3 appearances, we replace each appearance with a fresh variable  $x_i, i \in [k]$  and add to the formula the clauses  $(\neg x_1 \lor x_2) \land (\neg x_2 \lor x_3) \ldots (\neg x_k \lor x_1)$ . We repeat this for all variables that appear more than three times. Let  $I' = (\phi', X')$  be this new instance.

We claim that  $I = (\phi, X)$  is a yes instance of EXACTLY-1 3-SAT iff  $I' = (\phi', X')$  is a yes instance of RESTRICTED EXACTLY-1 (3,2)-SAT.

Let  $S: X \to \{T, F\}$  be a satisfying assignment for  $\phi$  such that each clause of  $\phi$  has exactly one true literal. It is not hard to see that  $S': X' \to \{T, F\}$  such that S'(x) = S(x) if  $x \in X$  and  $S'(x_i) = S(x)$  if  $x_i$  replaces one appearance of  $x \in X$ , is a satisfying assignment for  $\phi'$  such that each clause of  $\phi'$  has exactly one true literal.

Conversely, let  $S': X' \to \{T, F\}$  be a satisfying assignment for  $\phi'$  such that each clause of  $\phi'$  has exactly one true literal. Let  $x_i, i \in [k]$ , be the variables replacing x. Because we have clauses  $(\neg x_1 \lor x_2) \land (\neg x_2 \lor x_3) \ldots (\neg x_k \lor x_1)$  we know that all the  $x_i, i \in [k]$ , must have the same value in order to guarantee that all of these clauses have exactly one true literal. Furthermore, is not hard to see that  $S: X \to \{T, F\}$  where S(x) = S'(x) if  $x \in X'$  and  $S(x) = S'(x_1)$ if  $x_1$  replaces one appearance of x, then S is a satisfying assignment for  $\phi$  such that each clause of  $\phi$  has exactly one true literal.  $\Box$ 

Now, let us give a construction that we are going to use. **Construction.** 

Let  $\phi$  be an instance of RESTRICTED EXACTLY 1 (3,2)-SAT with variables  $X = \{x_1, \ldots, x_n\}$  and clauses  $C = \{c_1, \ldots, c_m\}$ . We will construct an instance of the crossword problem with alphabet  $\mathcal{L} = \{s_1, s_2, s_3\}$  where each letter has weight 1. The dictionary  $\mathcal{D}$  is as follows.

Let  $nl_j \in \{2,3\}$  be the number of literals in  $c_j$ . For each variable  $x_i$ , let  $a_i \leq 3$  be the number of its appearances in  $\phi$ . Then, we create  $3a_i$  words,  $d_{i,k,T}$ ,  $d_{i,k,F}$  and  $d_{i,k}$ , for each  $k \in [a_i]$  as follows.

•  $d_{i,k,T}$  and  $d_{i,k,F}$  have length m + n + 3i + k,

- the last letter of  $d_{i,k,T}$  is  $s_k$ ,
- the last letter of  $d_{i,k,F}$  is  $s_{k'}$  where k' := k + 1 when  $k < a_i$ , otherwise k' := 1,
- if the k-th appearance of  $x_i$  is positive then,  $d_{i,k,T}$  starts with  $s_1$  and  $d_{i,k,F}$  starts with  $s_2$ ,
- if the k-th appearance of  $x_i$  is negative then,  $d_{i,k,T}$  starts with  $s_2$  and  $d_{i,k,F}$  starts with  $s_1$ ,
- the word  $d_{i,k}$  has length m + i + 1 and starts with  $s_k$ , and
- all the other letters of these words can be chosen arbitrarily.

Observe that the above process gives three words for each literal in  $\phi$ .

For each clause  $c_j$ ,  $j \in [m]$ , we construct  $nl_j$  distinct words  $d_j^t$ ,  $t \in [nl_j]$  of length 1+j such that one of them starts with the letter  $s_2$ , the other  $nl_j-1$  words start with  $s_1$ , and the unspecified letters can be chosen arbitrarily. Observe that we have enough positions in order to create  $nl_j - 1$  distinct words starting with  $s_1$ , which indicates that we can create  $nl_j$  pairwise distinct words for each  $c_j$ .

In order to finish our construction we have to specify the grid. For each clause  $c_j$  and each literal l in  $c_j$  we construct two pairs of slots as follows. Let l be the k-th appearance of variable  $x_i, k \in [a_i]$ . The first pair of slots (type 1) consists of one horizontal slot  $hSlot_{j,1}^{i,k}$  of length m + n + 3i + k, and one vertical slot  $vSlot_{j,1}^{i,k}$  of length m + i + 1 such that, the last cell of the horizontal slot and the first cell of the vertical slot is the shared cell. The second pair of slots (type 2) consists of one horizontal slot  $hSlot_{j,2}^{i,k}$  of length m + n + 3i + k, and one vertical slot  $vSlot_{j,2}^{i,k}$  of length j + 1, that share their first cells. Here let us mention that the grid we constructed is consisted only by  $\mathcal{T}$ s.

Before we continue with the proof let us observe that in the instance of crossword puzzle we created the number of slots in the grid is equal to the number of words in the dictionary. Furthermore, we can specify in which slots each word can be assigned by considering the size of the words and slots. For any  $i \in [n]$  and  $k \in [a_i]$  the word  $d_{i,k}$  can be assigned only to the vertical slots of the type 1 pairs of slots. For any  $j \in [m]$  and  $t \in [nl_j]$  the word  $d_j$  can be assigned only to the vertical slots of the type 2 pairs of slots. The rest of the words can be assigned to horizontal slots of any type.

Let us first prove the following property where j(i, k) denotes the index of the clause where the k-th occurrence of  $x_i$  appears.

**Property 1.** For any given  $i \in [n]$ , slots  $hSlot_{j(i,k),1}^{i,k}$  and  $vSlot_{j(i,k),1}^{i,k}$  for  $k \in [a_i]$  are all filled iff we have assigned either all the words of  $\{d_{i,k,T} : k \in [a_i]\}$ , or all the words of  $\{d_{i,k,F} : k \in [a_i]\}$ , to the slots  $hSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ .

*Proof.* In one direction, if we have assigned to slots  $hSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ , all the words of  $\{d_{i,k,T} : k \in [a_i]\}$  or all the words of  $\{d_{i,k,F} : k \in [a_i]\}$ , then all the

letters  $s_1, \ldots, s_{a_i}$  appear exactly once in the end of these  $a_i$  slots. Because the words of  $\{d_{i,k}: k \in [a_i]\}$  start exactly with this set of letters, there is a unique way to assign them properly to the slots  $vSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ .

Conversely, assume that all the type 1 pairs of slots of  $x_i$  are filled. Because the only words that have the same length as slots  $vSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ , are the words of  $\{d_{i,k} : k \in [a_i]\}$ , we know that in the end of slots  $hSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ , each letter of  $\{s_1, \ldots, s_{a_i}\}$  appears exactly once. It is not hard to see that no combination of words except  $\{d_{i,k,T} : k \in [a_i]\}$  or  $\{d_{i,k,F} : k \in [a_i]\}$ , gives the same letters in the shared positions.

Now we are ready to present the proof of Theorem 5.

*Proof.* We show a reduction from RESTRICTED EXACTLY 1 (3,2)-SAT. We claim that  $\phi$  is a yes instance of RESTRICTED EXACTLY 1 (3,2)-SAT iff we can fill all the slots of the grid.

Suppose  $f: X \to \{T, F\}$  is a truth assignment so that each clause of  $\phi$  has exactly one true literal that satisfies  $\phi$ .

We are going to show a way to fill all the slots of the grid. Each variable  $x_i$  appears in  $a_i$  literals; let  $l(i, k), k \in [a_i]$ , be these literals and  $j(i, k) \in [m]$ ,  $k \in [a_i]$ , be the indices of the clauses  $c_{j(i,k)}$  that contain the corresponding literals.

For each variable  $x_i$ , fill the  $3a_i$  slots  $hSlot_{j(i,k),1}^{i,k}$ ,  $hSlot_{j(i,k),2}^{i,k}$  and  $vSlot_{j(i,k),1}^{i,k}$ for all  $k \in [a_i]$  as follows. If  $f(x_i) = T$ , then:

- assign  $d_{i,k,T}$  to  $hSlot_{i(i,k),1}^{i,k}$  for all  $k \in [a_i]$  and
- assign  $d_{i,k,F}$  to  $hSlot_{i(i,k),2}^{i,k}$  for all  $k \in [a_i]$ .

Otherwise  $(f(x_i) = F)$ :

- assign  $d_{i,k,F}$  to  $hSlot_{j(i,k),1}^{i,k}$  for all  $k \in [a_i]$  and
- assign  $d_{i,k,T}$  to  $hSlot_{i(i,k),2}^{i,k}$  for all  $k \in [a_i]$ .

Finally, in both cases, we assign the words of  $\{d_{i,k} : k \in [a_i]\}$  to the slots  $vSlot_{i(i,k),1}^{i,k}$  for  $k \in [a_i]$  in any way they fit.

In order to fill the grid completely, for each  $j \in [m]$ , we assign to the  $nl_j$  slots,  $vSlot_{j,2}^{i,k}$ , the words  $d_j^{k'}$  for  $k' \in [nl_j]$  in any way they fit.

It is not hard to see that we have assigned words to slots of the same length. It remains to prove that the words we have assigned have the same letters in the shared positions.

First observe that for a variable  $x_i$  and the slots  $hSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ , we have put either  $\{d_{i,k,T} : k \in [a_i]\}$  or  $\{d_{i,k,F} : k \in [a_i]\}$ . Therefore, we know by Property 1 that we can use the words of  $\{d_{i,k} : k \in [a_i]\}$  in the slots  $vSlot_{j(i,k),1}^{i,k}$ ,  $k \in [a_i]$ .

In the  $nl_j$  slots,  $vSlot_{j,2}^{i,k}$ , related to clause  $c_j$ , we have put the words  $d_j^{k'}$ ,  $k' \in [nl_j]$ . One of these words starts with  $s_2$  and the  $nl_j - 1$  others start with  $s_1$ . We will show that the same holds for the words we have assigned in the  $nl_j$  slots  $hSlot_{j,2}^{i,k}$ .

Observe that each literal  $l \in c_j$  can be described by a unique triplet (j, i, k)where  $j \in [m]$  is the index of the clause,  $i \in [n]$  is the index of the variable  $x_i$ on which l is built, and  $k \in [a_i]$  is the number of times that  $x_i$  has appeared in  $\phi$  until now. We claim that if the literal l described by (j, i, k) satisfies  $c_j$ , then the word assigned to  $hSlot_{j,2}^{i,k}$  starts with  $s_2$ , otherwise it starts with  $s_1$ . If l satisfies  $c_j$ , then either  $l = x_i$  and  $f(x_i) = T$  or  $l = \neg x_i$  and  $f(x_i) = F$ .

If l satisfies  $c_j$ , then either  $l = x_i$  and  $f(x_i) = T$  or  $l = \neg x_i$  and  $f(x_i) = F$ . If  $l = x_i$  (resp.,  $l = \neg x_i$ ), then we have assigned  $d_{i,k,F}$  (resp.,  $d_{i,k,T}$ ) to  $hSlot_{j,2}^{i,k}$  which starts with  $s_2$  because  $f(x_i) = T$  (resp.,  $f(x_i) = F$ ). If l does not satisfy  $c_j$ , then we used  $d_{i,k,T}$  (resp.,  $d_{i,k,F}$ ) which starts with  $s_1$ .

Finally, because we assumed that each clause is satisfied by exactly one literal, we know that one of the clause words starts with  $s_2$  and the other  $nl_j - 1$  clause words start with  $s_1$ .

Conversely, we claim that if we can fill the whole grid, then we can construct a truth assignment  $f: X \to \{T, F\}$  such that each clause of  $\phi$  has exactly one true literal. Furthermore, one such assignment is the following:

$$f(x_i) = \begin{cases} T, \text{ if } d_{i,1,T} \text{ is assigned to } hSlot_{j(i,1),1}^{i,1}, \\ F, \text{ otherwise.} \end{cases}$$
(1)

We first prove the following claim.

**Claim.** Let *l* be the literal of a clause  $c_j$  corresponding to the *k*-th appearance of some variable  $x_i$ . *l* is true under the truth assignment (1) iff the word in  $hSlots_{j,2}^{i,k}$  starts with  $s_2$ .

PROOF. Due to its length,  $hSlots_{j,2}^{i,k}$  receives either  $d_{i,k,T}$  or  $d_{i,k,F}$ , and one of these words starts with  $s_2$  whereas the other starts with  $s_1$ . Therefore, we have two cases. In the first case  $d_{i,k,F}$  starts with  $s_2$ , then  $d_{i,k,T}$  starts with  $s_1$  and  $l = x_i$ . In the second case,  $d_{i,k,T}$  starts with  $s_2$ ,  $d_{i,k,F}$  starts with  $s_1$  and  $l = \neg x_i$ .

Assume that  $d_{i,k,F}$  (resp.,  $d_{i,k,T}$ ) starts with  $s_2$ . By construction, we have that  $l = x_i$  (resp.,  $l = \neg x_i$ ).

If  $d_{i,k,F}$  (resp.,  $d_{i,k,T}$ ) is assigned to  $hSlots_{j,2}^{i,k}$ , then  $d_{i,k,T}$  (resp.,  $d_{i,k,F}$ ) is assigned to  $hSlots_{j,1}^{i,k}$ . By Property 1 we know that  $hSlots_{j,1}^{i,1}$  must contain  $d_{i,1,T}$  (resp.,  $d_{i,1,F}$ ) so  $f(x_i) = T$  (resp.,  $f(x_i) = F$ ). So, if  $d_{i,k,F}$  (resp.,  $d_{i,k,T}$ ) is assigned to  $hSlots_{j,2}^{i,k}$ , then we know that  $f(x_i) = T$  (resp.,  $f(x_i) = F$ ) and  $l = x_i$  (resp.,  $l = \neg x_i$ ) which means that l must be true under the truth assignment (1).

In reverse direction, if we have assigned  $d_{i,k,T}$  (resp.,  $d_{i,k,F}$ ) to  $hSlots_{j,2}^{i,k}$ , then we know that  $f(x_i) = F$  (resp.,  $f(x_i) = T$ ) and  $l = x_i$  (resp.,  $l = \neg x_i$ ) thus, l is false under the truth assignment (1). Based on the previous claim, we will show that each clause has exactly one true literal under the truth assignment f given in (1).

For any  $j \in [m]$  there are exactly  $nl_j$  pairs (i, k) where  $i \in [n]$  and  $k \in [a_i]$  such that the k-th appearance of  $x_i$  is in  $c_j$ . Let  $C_j$  be the set that contains contains all these pairs (i, k).

Observe that for each pair  $(i, k) \in C_j$  there exists a pair of slots  $hSlots_{j,2}^{i,k}$ ,  $vSlots_{j,2}^{i,k}$  which share their first cells. Because the grid is full, the  $nl_j$  vertical slots,  $vSlots_{j,2}^{i,k}$ , where  $(i,k) \in C_j$ , must contain the words  $d_j^t$ ,  $t \in [nl_j]$ . One of these words starts with  $s_2$  and  $nl_j - 1$  others start with  $s_1$ . Therefore, the same must hold for the words that have been assigned in the slots  $hSlots_{j,2}^{i,k}$  for  $(i,k) \in C_j$ .

Using the previous claim, we know that one of the literals in  $c_j$  is true and the other  $nl_j - 1$  are false under the truth assignment 1. Therefore, if we can fill the whole grid, then there exists a truth assignment such that exactly one literal of each clause of  $\phi$  is true.

**Remark 1.** In our construction each  $\mathcal{T}$  has unique shape<sup>4</sup> so the problem remains NP-hard even in this case.

**Remark 2.** Theorem 3 can be adjusted to work also for the case where word reuse is not allowed. We simply need to add a suffix of length  $\log m$  to all words of length 2k - 1 and add rows to the grid accordingly. Hence, under the ETH, no algorithm can solve this problem in time  $m^{o(k)}$ , where k is the number of horizontal slots.

Finally, we present two algorithms. The first is a modification of the algorithm presented in Theorem 2.

**Theorem 7.** If word reuse is not allowed, then CP-OPT can be solved in time  $4^m(m+1)^{tw}n^{O(1)}$  on inputs where tw is the treewidth of the grid graph.

*Proof.* The algorithm is similar to the one presented for Theorem 2. The main difference is that, for each node of the tree decomposition, we also need to keep a subset of the dictionary. This subset, S, indicates which words have been already used to fill the slots represented by vertices in the  $B_t^{\downarrow} \setminus B_t$ , where t is the considered node. We will say that S is the set of used words. Since we do not want to reuse words, we only keep the tuples where, for any pair  $(\sigma, S)$  appearing in a tuple it holds that:

- $\sigma(v) \neq \sigma(u)$  for all  $v, u \in B_t$  where  $v \neq u$  and
- $\sigma(v) \notin S$  for all  $v \in B_t$ .

 $<sup>^{4}</sup>$ Two crosses are of the same shape if they are identical: same number of horizontal cells, same number of vertical cells, and same shared cell.

This increases the number of tuples we need to keep by a factor of  $2^m$ . We also need to pay attention when we consider Join and Forget nodes. In a Forget node t, if c is the child node,  $v \in B_c \setminus B_t$  and we have a tuple for c that includes the pair  $(\sigma, S)$  then we create a tuple for t in the same way we created it in Theorem 2 and we include the set  $S \cup \sigma(v)$  as the set of used words. Finally in the Join nodes we keep a tuple for a pair  $(\sigma, S)$  if there exist tuples stored for the children of the considered node such that:

- the assignments of both tuples is  $\sigma$  and
- the sets  $S_1$ ,  $S_2$  stored in the tuples of the children nodes are disjoint, i.e.  $S_1 \cap S_2 = \emptyset$ .

This resolves to an algorithm with running time  $4^m(m+1)^{tw}n^{O(1)}$ .

For the second algorithm, observe that by filling the slots represented by a vertex cover of the grid graph, all the shared cells are pre-filled. Since there are at most  $(m+1)^k$  (where k is the size of the vertex cover) ways to assign words to these slots, by Proposition 1, we get the following corollary.

**Corollary 8.** Given a vertex cover of size k of the grid graph we can solve CP-DEC and CP-OPT in time  $(m+1)^k(n+m)^{O(1)}$ . Furthermore, as vertex cover we can take the set of horizontal slots.

Therefore, the bound given in Remark 2 for the parameter vertex cover is tight.

#### 4. Parameterized by Total Number of Slots

In this section we consider a much more restrictive parameterization of the problem: we consider instances where the parameter is n, the total number of slots. Recall that in Theorem 3 (and Remark 2) we already considered the complexity of the problem parameterized by the number of *horizontal* slots of the instance. We showed that this case of the problem cannot be solved in  $m^{o(k)}$  and that an algorithm with running time roughly  $(m + 1)^k$  is possible whether word reuse is allowed or not.

Since parameterizing by the number of horizontal slots is not sufficient to render the problem FPT, we therefore consider our parameter to be the total number of slots. This is, finally, sufficient to obtain a simple FPT algorithm.

**Corollary 9.** There is an algorithm that solves CP-DEC and CP-OPT in time  $O^*((\ell+1)^{n^2/4})$ , where n is the total number of slots and  $\ell$  the size of the alphabet, whether word reuse is allowed or not.

*Proof.* Since there are n slots in the instance, even if the grid is a complete bipartite graph, the instance contains at most  $n^2/4$  cells which are shared between two slots. In time  $(\ell+1)^{n^2/4}$  we consider all possible letters that could be placed in these cells. Finally, as we have shown in Proposition 1, each of these instances can be solved in polynomial time.

Even though the running time guaranteed by Corollary 9 is FPT for parameter n, we cannot help but observe that the dependence on n is rather disappointing, as our algorithm is exponential *in the square* of n. It is therefore a natural question whether an FPT algorithm for this problem can achieve complexity  $2^{o(n^2)}$ , assuming the alphabet size is bounded. The main result of this section is to establish that this is likely to be impossible.

**Overview** Our hardness proof consists of two steps. In the first step we reduce 3-SAT to a version of the same problem where variables and clauses are partitioned into  $O(\sqrt{n+m})$  groups, which we call SPARSE 3-SAT. The key property of this intermediate problem is that interactions between groups of variables and groups of clauses are extremely limited. In particular, for each group of variables  $V_i$  and each group of clauses  $C_j$ , at most one variable of  $V_i$  appears in a clause of  $C_j$ . We obtain this rather severe restriction via a randomized reduction that runs in expected polynomial time. The second step is to reduce SPARSE 3-SAT to CP-DEC. Here, every horizontal slot will represent a group of variables and every vertical slot a group of clauses, giving  $O(\sqrt{n+m})$  slots in total. Hence, an algorithm for CP-DEC whose dependence on the total number of slots is subquadratic in the exponent will imply a sub-exponential time (randomized) algorithm for 3-SAT. The limited interactions between groups of clauses and variables will be key in allowing us to execute this reduction using a *binary* alphabet.

Let us now define our intermediate problem.

**Definition 2.** In SPARSE 3-SAT we are given an integer n which is a perfect square and a 3-SAT formula  $\phi$  with at most n variables and at most n clauses, such that each variable appears in at most 3 clauses. Furthermore, we are given a partition of the set of variables V and the set of clauses C into  $\sqrt{n}$  sets  $V_1, \ldots, V_{\sqrt{n}}$  and  $C_1, \ldots, C_{\sqrt{n}}$  of size at most  $\sqrt{n}$  each, such that for all  $i, j \in [\sqrt{n}]$  the number of variables of  $V_i$  which appear in at least one clause of  $C_j$  is at most one.

Now, we are going to prove the hardness of SPARSE 3-SAT, which is the first step of our reduction.

**Lemma 10.** Suppose the randomized ETH is true. Then, there exists an  $\epsilon > 0$  such that SPARSE 3-SAT cannot be solved in time  $2^{\epsilon n}$ .

The first step of our reduction will be to prove that SPARSE 3-SAT cannot be solved in sub-exponential time (in n) under the randomized ETH, via a reduction from 3-SAT. To do this, we will need the following combinatorial lemma.

**Lemma 11.** For each  $\epsilon > 0$  there exists C > 0 such that for sufficiently large n we have the following. There exists a randomized algorithm running in expected polynomial time which, given a bipartite graph G = (A, B, E) such that |A| = |B| = n and the maximum degree of G is 3, produces a set  $V' \subseteq A \cup B$  with

 $|V'| \ge 2(1-\epsilon)n$  and a coloring  $c: V' \to [k]$  of the vertices of V' with k colors, where  $k \le C\sqrt{n}$ , such that for all  $i \in [k]$  we have  $|c^{-1}(i)| \le \sqrt{n}$  and for all  $i, j \in [k]$  the graph induced by  $c^{-1}(i) \cup c^{-1}(j)$  contains at most one edge.

*Proof.* Let  $k = C\lceil \sqrt{n} \rceil$ , where C is a sufficiently large constant (depending only on  $\epsilon$ ) to be specified later. We color each vertex of the graph uniformly at random from a color in [k], call this coloring c. Let  $X_{i,j}$  be the set of edges which have as endpoints a vertex of color i and a vertex of color j.

Our algorithm is rather simple: initially, we set V' = V. Then, for each  $i, j \in [k]$  we check whether  $X_{i,j}$  contains at most one edge. If yes, we do nothing; if not, we select for each edge  $e \in X_{i,j}$  an arbitrary endpoint and remove that vertex from V'. In the end we return the set V' that remains and its coloring. It is clear that this satisfies the property that  $c^{-1}(i) \cup c^{-1}(j)$  contains at most one edge for the graph induced by V' for all  $i, j \in [k]$ , so what we need to argue is that (i)  $|c^{-1}(i)| \leq \sqrt{n}$  for all i with high probability and (ii) that V' has the promised size with at least constant probability. If we achieve this it will be sufficient to repeat the algorithm a polynomial number of times to obtain the claimed properties with high probability, hence we will have an expected running time polynomial in n.

For the first part, fix an  $i \in [k]$  and observe that  $E[|c^{-1}(i)|] \leq \frac{2\sqrt{n}}{C}$ . To prove that all  $|c^{-1}(i)|$  are of size at most  $4\sqrt{n}/C$  with high probability (and hence also at most  $\sqrt{n}$  for C sufficiently large), we will use Chernoff's Inequality.

**Proposition 12** (Chernoff's Inequality). Let X be a binomial random variable and  $\epsilon > 0$ . Then  $P[|X - E[X]| > \epsilon E[X]| < 2e^{-\epsilon^2 E[X]/3}$ 

We take  $\epsilon = 1$ . It follows that  $P[|c^{-1}(i)| > 4\sqrt{n}/C] \leq 2e^{-2\sqrt{n}/3C}$ . Now, taking the union bound, we obtain that almost surely for all color i,  $|c^{-1}(i)| < 4\sqrt{n}/C$ 

The more interesting part of this proof is to bound the expected size of V'. Let e be an edge whose endpoints are colored with colors i and j. We say that e is good if no other edge in G has one endpoint colored i and the other colored j by the coloring c. Let u and v be the endpoints of e. The probability of another edge having endpoints of colors i and j in the graph  $G - \{u, v\}$  is at most  $\frac{2|E|}{C^2n} \leq \frac{6}{C^2}$ . The probability that at least one of the at most four edges incident to e has endpoints colored i and j is at most  $\frac{4}{C\sqrt{n}}$ . Thus, the probability that e is good is at least  $1 - \frac{6}{C^2} - \frac{4}{C\sqrt{n}} > 1 - \frac{7}{C^2}$ , if n is sufficiently large. Let X be the number of edges which are not good. Then,  $E[X] \leq 7C^{-2}|E|$ . By Markov's Inequality  $P[X > 21C^{-2}|E|] < 1/3$ . Thus, with probability at least 2/3, our algorithm will remove at most  $2\epsilon n$  vertices, it suffices to select any value  $C \geq \frac{8}{\sqrt{\epsilon}}$ .

Now, we present the proof of Lemma 10

*Proof.* Suppose that the statement is false, therefore for any  $\epsilon > 0$  we can solve SPARSE 3-SAT in which the number of variables and clauses can be upperbounded by N in expected time  $2^{\epsilon N}$  using some supposed algorithm. Fix an arbitrary  $\epsilon' > 0$ . We will show how to solve an arbitrary instance of 3-SAT with n variables and m clauses in expected time  $2^{\epsilon'(n+m)}$  using this supposed algorithm for SPARSE 3-SAT. If we can do this for any arbitrary  $\epsilon'$ , this will contradict the randomized ETH.

Start with an arbitrary 3-SAT instance  $\phi$  with *n* variables and *m* clauses. We first edit  $\phi$  to ensure that each variable appears at most three times. In particular, if *x* appears k > 3 times, we replace each appearance of *x* with a fresh variable  $x_i$ ,  $i \in [k]$ , and add the clauses  $(\neg x_1 \lor x_2) \land (\neg x_2 \lor x_3) \land \ldots \land (\neg x_k \lor x_1)$ .

The number of variables in the new instance is at most n+3m. The number of clauses is at most 4m. This is because every new clause and every new variable corresponds to an occurrence of an original variable in an original clause and there are at most 3m such occurrences.

We now have an instance  $\phi'$  equivalent to  $\phi$  with at most n + 3m variables and at most 4m clauses, such that each variable appears at most 3 times. Let Nbe the smallest perfect square such that  $N \ge n + 4m$ . We have N < 10(n+m). What we need now is to produce a partition of the vertices and clauses of  $\phi'$ .

In order to produce this partition we invoke Lemma 11 on the incidence graph of  $\phi'$ , that is, the bipartite graph where we have variables on one side and clauses on the other, and edges signify that a variable appears in a clause. Add some dummy isolated vertices on each side so that both sides of the incidence graph contain N vertices. We invoke Lemma 11 by setting  $\epsilon$  to be  $\epsilon'/80$ . We obtain a coloring of all but at most  $\frac{\epsilon' N}{40} \leq \frac{\epsilon'(n+m)}{4}$  of the vertices of the incidence graph.

Let U be the set of variables and clauses that correspond to uncolored vertices of the incidence graph. Then, for each such variable we produce two formulas (one by setting it to True and one by setting it to False), and for each such clause, at most 3 formulas (one by setting each of the literals of the clause to True). We thus construct at most  $3^{\epsilon'(n+m)/4} \leq 2^{\epsilon'(n+m)/2}$  new formulas, such that one of them is satisfiable if and only if  $\phi$  was satisfiable. We will then use the supposed algorithm for SPARSE 3-SAT to decide each of these formulas one by one.

Each new formula we have contains at most N variables and at most N clauses, and by Lemma 11 we have partitions of the variables and clauses into  $C\sqrt{N}$  groups, where C is a constant (that depends on  $\epsilon'$ ). By setting  $N' = \lceil C \rceil^2 N$  we can view these instances as instances of SPARSE 3-SAT, because then the number of groups becomes equal to the square root of the upper bound on the number of variables and clauses, and by the properties of Lemma 11 there is at most one edge between each group of variables and each group of clauses. Since we suppose that for all  $\epsilon > 0$  such instances can be solved in time  $2^{\epsilon N'}$ , by setting  $\epsilon = \epsilon'/50\lceil C \rceil^2$  we can solve each formula in  $2^{\epsilon'(n+m)/5}$ . The total expected running time of our algorithm is at most  $2^{\epsilon'(n+m)/2} \cdot 2^{\epsilon'(n+m)/5} \cdot (n+m)^{O(1)} < 2^{\epsilon'(n+m)}$ , so we contradict the ETH.

We are now ready to prove the main theorem of this section.

**Theorem 13.** Suppose the randomized ETH is true. Then, there exists an  $\epsilon > 0$  such that CP-DEC on instances with a binary alphabet cannot be solved in time  $2^{\epsilon n^2} \cdot m^{O(1)}$ . This holds also for instances where all slots have distinct sizes (so words cannot be reused).

Proof. Suppose for the sake of contradiction that for any fixed  $\epsilon > 0$ , CP-DEC on instances with a binary alphabet can be solved in time  $2^{\epsilon n^2} \cdot m^{O(1)}$ . We will then contradict Lemma 10. In particular, we will show that for any  $\epsilon'$  we can solve SPARSE 3-SAT in time  $2^{\epsilon'N}$ , where N is the upper bound on the number of variables and clauses. Fix some  $\epsilon' > 0$  and suppose that  $\phi$  is an instance of SPARSE 3-SAT with at most N variables and at most N clauses, where N is a perfect square. Recall that the variables are given partitioned into  $\sqrt{N}$  sets,  $V_1, \ldots, V_{\sqrt{N}}$  and the clauses partitioned into  $\sqrt{N}$  sets  $C_1, \ldots, C_{\sqrt{N}}$ . In the remainder, when we write  $V(C_j)$  we will denote the set of variables that appear in a clause of  $C_j$ . Recall that the partition satisfies the property that for all  $i, j \in [\sqrt{N}]$  we have  $|V_i \cap V(C_j)| \leq 1$ . Suppose that the variables of  $\phi$  are ordered  $x_1, x_2, \ldots, x_N$ .

We construct a grid as follows: for each group  $V_i$  we construct a horizontal slot and for each group  $C_j$  we construct a vertical slot, in a way that all slots have distinct lengths. More precisely, the *i*-th horizontal slot, for  $i \in [\sqrt{N}]$  is placed on row 2i - 1, starts in the first column and has length  $2\sqrt{N} + 2i$ . The *j*-th vertical slot is placed in column 2j - 1, starts in the first row and has length  $5\sqrt{N} + 2j$ . (As usual, we number the rows and columns top-to-bottom and left-to-right). Observe that all horizontal slots intersect all vertical slots; in particular, the cell in row 2i - 1 and column 2j - 1 is shared between the *i*-th horizontal and *j*-th vertical slot, for  $i, j \in [\sqrt{N}]$ . We define  $\mathcal{L}$  to contain two letters  $\{0, 1\}$ .

What remains is to describe the dictionary.

- For each  $i \in [\sqrt{N}]$  and for each assignment function  $\sigma : V_i \to \{0, 1\}$  we construct a word  $w_{\sigma}$  of length  $2\sqrt{N} + 2i$ . The word  $w_{\sigma}$  has the letter 0 in all positions, except positions 2j 1, for  $j \in [\sqrt{N}]$ . For each such j, we consider  $\sigma$  restricted to  $V_i \cap V(C_j)$ . By the properties of SPARSE 3-SAT, we have  $|V_i \cap V(C_j)| \leq 1$ . If  $V_i \cap V(C_j) = \emptyset$  then we place letter 0 in position 2j 1; otherwise we set in position 2j 1 the letter that corresponds to the value assigned by  $\sigma$  to the unique variable of  $V_i \cap V(C_j)$ .
- For each  $j \in [\sqrt{N}]$  and for each satisfying assignment function  $\sigma : V(C_j) \rightarrow \{0, 1\}$ , that is, every assignment function that satisfies all clauses of  $C_j$ , we construct a word  $w'_{\sigma}$  of length  $5\sqrt{N} + 2j$ . The word  $w'_{\sigma}$  has the letter 0 in all positions, except positions 2i 1, for  $i \in [\sqrt{N}]$ . For each such i, we consider  $\sigma$  restricted to  $V_i \cap V(C_j)$ . If  $V_i \cap V(C_j) = \emptyset$  then we place letter 0 in position 2i 1; otherwise we set in position 2i 1 the letter that corresponds to the value assigned by  $\sigma$  to the unique variable of  $V_i \cap V(C_j)$ .

The construction is now complete. We claim that if  $\phi$  is satisfiable, then

it is possible to fill out the grid we have constructed. Indeed, fix a satisfying assignment  $\sigma$  to the variables of  $\phi$ . For each  $i \in [\sqrt{N}]$  let  $\sigma_i$  be the restriction of  $\sigma$  to  $V_i$ . We place in the *i*-th horizontal slot the word  $w_{\sigma_i}$ . Similarly, for each  $j \in [\sqrt{N}]$  we let  $\sigma'_j$  be the restriction of  $\sigma$  to  $V(C_j)$  and place  $w'_{\sigma'_j}$  in the *j*-th vertical slot. Now if we examine the cell shared by the *i*-th horizontal and *j*-th vertical slot, we can see that it contains a letter that represents  $\sigma$  restricted to (the unique variable of)  $V_i \cap V(C_j)$  or 0 if  $V_i \cap V(C_j) = \emptyset$ , and both the horizontal and vertical word place the same letter in that cell.

For the converse direction, if the grid is filled, we can extract an assignment  $\sigma$  for the variables of  $\phi$  as follows: for each  $x \in V_i$  we find a  $C_j$  such that x appears in some clause of  $C_j$  (we can assume that every variable appears in some clause). We then look at the cell shared between the *i*-th horizontal and the *j*-th vertical slot. The letter we have placed in that cell gives an assignment for the variable contained  $V_i \cap V(C_j)$ , that is x. Having extracted an assignment to all the variables, we claim it must satisfy  $\phi$ . If not, there is a group  $C_j$  that contains an unsatisfied clause. Nevertheless, in the *j*-th vertical slot we have placed a word that corresponds to a *satisfying* assignment for the clauses of  $C_j$ , call it  $\sigma_j$ . Then  $\sigma_j$  must disagree with  $\sigma$  in a variable x that appears in  $C_j$ . Suppose this variable is part of  $V_i$ . Then, this would contradict the fact that we extracted an assignment for x from the word placed in the *i*-th horizontal slot.

Observe that the new instance has  $n = 2\sqrt{N}$  slots. If there exists an algorithm that solves CP-DEC in time  $2^{\epsilon n^2} m^{O(1)}$  for any  $\epsilon > 0$ , we set  $\epsilon = \epsilon'/8$  (so  $\epsilon$  only depends on  $\epsilon'$ ) and execute this algorithm on the constructed instance. We observe that  $m \leq 2\sqrt{N} \cdot 7^{\sqrt{N}}$ , and that  $2^{\epsilon n^2} \leq 2^{\epsilon' N/2}$ . Assuming that N is sufficiently large, using the supposed algorithm for CP-DEC we obtain an algorithm for SPARSE 3-SAT with complexity at most  $2^{\epsilon' N}$ . Since we can do this for arbitrary  $\epsilon'$ , this contradicts the randomized ETH.

#### 5. Approximability of CP-Opt

This section begins with a  $\left(\frac{1}{2}+O(\frac{1}{n})\right)$ -approximation algorithm which works when words can, or cannot, be reused. After that, we prove that under the unique games conjecture, an approximation algorithm with a significantly better ratio is unlikely.

**Theorem 14.** CP-OPT is  $(\frac{1}{2} + \frac{1}{2(\varepsilon n+1)})$ -approximable in polynomial time, for all  $\varepsilon \in (0, 1]$ .

Proof. Fix some  $\varepsilon \in (0, 1]$ . Let  $k_v := \min(\lceil \frac{1}{\varepsilon} \rceil, n-h)$  and  $r_v := \lceil \frac{n-h}{k_v} \rceil$ , where h is the number of horizontal slots in the grid. Create  $r_v$  groups of vertical slots  $G_1, \ldots, G_{r_v}$  such that  $|G_i| \leq k_v$  for all  $i \in [r_v]$  and  $G_1 \cup \ldots \cup G_{r_v}$  covers the entire set of vertical slots. For each  $G_i$ , guess an optimal choice of words, i.e., identical to a global optimum, and complete this partial solution by filling the horizontal slots (use the aforementioned matching technique where the words selected for  $G_i$  are excluded from  $\mathcal{D}$ ). Each slot of  $\bigcup_{i \neq i} G_j$  gets the empty word.

Since  $|G_i| \leq k_v$ , guessing an optimal choice of words for  $G_i$  by brute force requires at most  $(m+1)^{k_v}$  combinations. This is done  $r_v$  times (once for each  $G_i$ ). The maximum matching runs in time  $\mathcal{O}((m+n)^2 \cdot mn)$ . In all, the time complexity of the algorithm is  $\mathcal{O}((m+1)^{k_v} \cdot r_v \cdot (m+n)^2 \cdot mn) \leq \mathcal{O}((m+1)^{1/\varepsilon} \cdot mn)$  $\varepsilon n \cdot (m+n)^2 \cdot mn$ ).

Assume that, given an optimal solution,  $W_H^*$  and  $W_V^*$  are the total weight of the words assigned to the horizontal and vertical slots, respectively, both including the shared cells. Furthermore, let  $W_S^*$  be the weight of the letters assigned to the shared cells in the optimal solution. Observe that the weight of the optimal solution is  $W_H^* + W_V^* - W_S^*$  and the weight of our solution is at

least  $W_H^* + \frac{1}{r_v}(W_V^* - W_S^*)$ . We repeat the same process, but the roles of vertical and horizontal slots are interchanged. Fix a parameter  $k_h := \min(\lceil \frac{1}{\varepsilon} \rceil, h)$ . Create  $r_h := \lceil \frac{h}{k_h} \rceil$  groups of horizontal slots  $G_1, \ldots, G_{r_h}$  such that  $|G_i| \leq k_h$  for all  $i \in [r_h]$  and  $G_1 \cup \ldots \cup G_{r_h}$ covers the entire set of horizontal slots. For each  $G_i$ , guess an optimal choice of words and complete this partial solution by filling the vertical slots. Each slot of  $\bigcup_{j\neq i} G_j$  gets the empty word.

Using the same arguments as above, we can conclude that the time complexity is  $O((m+1)^{1/\varepsilon} \cdot \varepsilon n \cdot (m+n)^2 \cdot mn)$  and that we return a solution of weight at least  $W_V^* + \frac{1}{r_h}(W_H^* - W_S^*)$ .

Finally, between the two solutions, we return the one with the greater weight. It remains to argue about the approximation ratio. We need to consider two cases:  $W_H^* \ge W_V^*$  and  $W_V^* > W_H^*$ . Suppose  $W_H^* \ge W_V^*$ . The first approximate solution has value  $W_H^* + \frac{1}{r_v}(W_V^* - W_V^*)$ .

 $W_{S}^{*} \geq \frac{1+1/r_{v}}{2}(W_{H}^{*}+W_{V}^{*}-W_{S}^{*})$ . If  $k_{v} = n-h$  then  $r_{v} = 1$  and our approximation ratio is 1. Otherwise,  $k_v = \lceil \frac{1}{\varepsilon} \rceil$  and  $r_v = \lceil \frac{n-h}{\lceil 1/\varepsilon \rceil} \rceil \le \frac{n-h}{\lceil 1/\varepsilon \rceil} + 1 = \frac{n-h+\lceil 1/\varepsilon \rceil}{\lceil 1/\varepsilon \rceil}$ . It follows that  $\frac{1}{r_v} \ge \frac{\lceil l/\varepsilon \rceil}{n-h+\lceil l/\varepsilon \rceil}$ . Use  $n-h+\lceil l/\varepsilon \rceil \le n+\frac{1}{\varepsilon}$  and  $\lceil l/\varepsilon \rceil \ge l/\varepsilon$  to get that  $\frac{1}{r_v} \ge \frac{l/\varepsilon}{n+1/\varepsilon} = \frac{1}{\varepsilon n+1}$ . Our approximation ratio is at least  $\frac{1+l/(\varepsilon n+1)}{2}$ . Suppose  $W_V^* > W_H^*$ . The second approximate solution has value  $W_V^* + \frac{1}{r_h}(W_H^* - W_S^*) > \frac{1+l/r_h}{2}(W_H^* + W_V^* - W_S^*)$ . If  $k_h = h$ , then our approximation

ratio is 1. Otherwise,  $k_h = \lceil \frac{1}{\varepsilon} \rceil$  and, using the same arguments, our approximation ratio is at least  $\frac{1+1/(\varepsilon n+1)}{2}$ . Note that  $\frac{1+1/(\varepsilon n+1)}{2} \leq 1$ . In all, we have a  $\frac{1+1/(\varepsilon n+1)}{2}$ -approximate solution in  $\mathcal{O}((m+1)^{1/\varepsilon} \cdot \varepsilon n \cdot (m+n)^2 \cdot mn)$  for all  $\varepsilon \in (0,1]$ .

The previous approximation algorithm only achieves an approximation ratio of  $\frac{1}{2} + O(\frac{1}{n})$ , which tends to  $\frac{1}{2}$  as n increases. At first glance this is quite disappointing, as someone can observe that a ratio of  $\frac{1}{2}$  is achievable simply by placing words only on the horizontal or the vertical slots of the instance<sup>5</sup>. Nevertheless, we are going to show that this performance is justified, as im-

 $<sup>^{5}</sup>$ This placement is done in a way that maximizes the weight, using the matching technique as in the proof of Proposition 1.

proving upon this trivial approximation ratio would falsify the Unique Games Conjecture (UGC).

Before we proceed, let us recall some relevant definitions regarding Unique Games. The UNIQUE LABEL COVER problem is defined as follows: we are given a graph G = (V, E), with some arbitrary total ordering  $\prec$  of V, an integer R, and for each  $(u,v) \in E$  with  $u \prec v$  a 1-to-1 constraint  $\pi_{(u,v)}$  which can be seen as a permutation on [R]. The vertices of G are considered as variables of a constraint satisfaction problem, which take values in [R]. Each constraint  $\pi_{(u,v)}$  defines for each value of u a unique value that must be given to v in order to satisfy the constraint. The goal is to find an assignment to the variables that satisfies as many constraints as possible. The Unique Games Conjecture states that for all  $\epsilon > 0$ , there exists R, such that distinguishing instances of UNIQUE LABEL COVER for which it is possible to satisfy a  $(1 - \epsilon)$ -fraction of the constraints from instances where no assignment satisfies more than an  $\epsilon$ fraction of the constraints is NP-hard. In this section we will need a slightly different version of this conjecture, which was defined by Khot and Regev as the Strong Unique Games Conjecture. Despite the name, Khot and Regev showed that this version is implied by the standard UGC. The precise formulation is the following:

**Theorem 15.** [Theorem 3.2 of [13]] If the Unique Games Conjecture is true, then for all  $\epsilon > 0$  it is NP-hard to distinguish between the following two cases of instances of UNIQUE LABEL COVER G = (V, E):

- (Yes case): There exists a set  $V' \subseteq V$  with  $|V'| \ge (1 \epsilon)|V|$  and an assignment for V' such that all constraints with both endpoints in V' are satisfied.
- (No case): For any assignment to V, for any set  $V' \subseteq V$  with  $|V'| \ge \epsilon |V|$ , there exists a constraint with both endpoints in V' that is violated by the assignment.

Using the version of the UGC given in Theorem 15 we are ready to present our hardness of approximation argument for the crossword puzzle.

**Theorem 16.** Suppose that the Unique Games Conjecture is true. Then, for all  $\epsilon$  with  $\frac{1}{4} > \epsilon > 0$ , there exists an alphabet  $\Sigma_{\epsilon}$  such that it is NP-hard to distinguish between the following two cases of instances of the crossword problem on alphabet  $\Sigma_{\epsilon}$ :

- (Yes case): There exists a valid solution that fills a  $(1 \epsilon)$ -fraction of all cells.
- (No case): No valid solution can fill more than a  $(\frac{1}{2} + \epsilon)$ -fraction of all cells.

Moreover, the above still holds if all slots have distinct lengths (and hence reusing words is trivially impossible). Proof. Fix an  $\epsilon > 0$ . We will later define an appropriately chosen value  $\epsilon' \in (0, \epsilon)$ whose value only depends on  $\epsilon$ . We present a reduction from a UNIQUE LABEL COVER instance, as described in Theorem 15. In particular, suppose we have an instance G = (V, E), with |V| = n, alphabet [R], such that (under UGC) it is NP-hard to distinguish if there exists a set V' of size  $(1 - \epsilon')n$  that satisfies all its induced constraints, or if all sets V' of size  $\epsilon'n$  induce at least one violated constraint for any assignment. Throughout this proof we assume that n is sufficiently large (otherwise the initial instance is easy). In particular, let  $n > \frac{20}{\epsilon}$ .

We construct an instance of the crossword puzzle that fits in an  $N \times N$  square, where  $N = 4n + n^2$ . We number the rows  $1, \ldots, N$  from top to bottom and the columns  $1, \ldots, N$  from left to right. The instance contains n horizontal and nvertical slots. For  $i \in [n]$ , the *i*-th horizontal slot is placed in row 2i, starting at column 1, and has length  $2n + n^2 + i$ . For  $j \in [n]$ , the *j*-th vertical slot is placed in column 2j, starts at row 1 and has length  $3n + n^2 + j$ . Observe that all horizontal slots intersect all vertical slots and in particular, for all  $i, j \in [n]$  the cell in row 2i, column 2j belongs to the *i*-th horizontal slot and the *j*-th vertical slot. Furthermore, each slot has a distinct length, as the longest horizontal slot has length  $3n + n^2$  while the shortest vertical slot has length  $3n + n^2 + 1$ .

We define the alphabet as  $\Sigma_{\epsilon} = [R] \cup \{*\}$ . Before we define our dictionary, let us give some intuition. Let  $V = \{v_1, \ldots, v_n\}$ . The idea is that a variable  $v_i \in V$  of the original instance will be represented by both the *i*-th horizontal slot and the *i*-th vertical slot. In particular, we will define, for each  $\alpha \in [R]$ a pair of words that we can place in these slots to represent the fact that  $v_i$ is assigned with the value  $\alpha$ . We will then ensure that if we place words on both the *i*-th horizontal slot and the *j*-th horizontal slot, where  $(v_i, v_j) \in E$ , then the assignment that can be extracted by reading these words will satisfy the constraint  $\pi_{(v_i, v_j)}$ . The extra letter \* represents an indifferent assignment (which we need if  $(v_i, v_j) \notin E$ ).

Armed with this intuition, let us define our dictionary.

- For each  $i \in [n]$ , for each  $\alpha \in [R]$  we define a word  $d_{(i,\alpha)}$  of length  $2n + n^2 + i$ . The word  $d_{(i,\alpha)}$  has the character \* everywhere except at position 2i and at positions 2j for  $j \in [n]$  and  $(v_i, v_j) \in E$ . In these positions the word  $d_{(i,\alpha)}$  has the character  $\alpha$ .
- For each  $j \in [n]$ , for each  $\alpha \in [R]$  we define a word  $d'_{(j,\alpha)}$  of length  $3n + n^2 + j$ . The word  $d'_{(j,\alpha)}$  has the character \* everywhere except at position 2j and at positions 2i for  $i \in [n]$  and  $(v_i, v_j) \in E$ . In position 2j we have the character  $\alpha$ . In position 2i with  $(v_i, v_j) \in E$ , we place the character  $\beta \in [R]$  such that the constraint  $\pi_{(v_i, v_j)}$  is satisfied by assigning  $\beta$  to  $v_i$  and  $\alpha$  to  $v_j$ . (Note that  $\beta$  always exists and is unique, as the constraints are permutations on [R], that is, for each value  $\alpha$  of  $v_j$  there exists a unique value  $\beta$  of  $v_i$  that satisfies the constraint).

This completes the construction. Suppose now that  $V = \{v_1, \ldots, v_n\}$  and that we started from the Yes case of UNIQUE LABEL COVER, that is, there



Figure 3: This is a part of the grid we construct. Here, we have assumed that  $(v_i, v_j) \in E$ and assigning  $\alpha$  to  $v_i$  and  $\beta$  to  $v_j$  satisfies the constraint  $\pi_{(v_i, v_j)}$ . The words that fit in the  $i^{th}$  horizontal slot forcing the same symbol in all the cells that correspond to existing edges. Here we have forced the letter  $\alpha$ . If both the  $i^{th}$  horizontal and vertical slots are filled, then the vertical words force letters in the rest of the slots that satisfy the constraints. Notice that filling the *i*th vertical and horizontal slots together with the *j*th vertical and horizontal slots is equivalent to assigning to  $v_i$  and  $v_j$  integers from [R] such that the constraint  $\pi_{(v_i, v_j)}$  is satisfied.

exists a set  $V' \subseteq V$  such that  $|V'| \ge (1 - \epsilon')n$  and all constraints induced by V'can be simultaneously satisfied. Fix an assignment  $\sigma : V' \to [R]$  that satisfies all constraints induced by V'. For each  $i \in [n]$  such that  $v_i \in V'$  we place in the *i*-th horizontal slot (that is, in row 2i) the word  $d_{(i,\sigma(v_i))}$ . For each  $j \in [n]$  such that  $v_j \in V'$  we place in the *j*-th vertical slot the word  $d'_{(j,\sigma(v_j))}$ . We leave all other slots empty. We claim that this solution is valid, that is, no shared cell is given different values from its horizontal and vertical slot. To see this, examine the cell in row 2i and column 2j. If both of the slots that contain it are filled, then  $v_i, v_j \in V'$ . If  $(v_i, v_j) \notin E$  and  $i \neq j$ , then the cell contains \* from both words. If i = j, then the cell contains  $\sigma(v_i)$  from both words. If  $i \neq j$  and  $(v_i, v_j) \in E$ , then the cell contains  $\sigma(v_i)$ . This is consistent with the vertical word, as the constraint  $\pi_{(v_i, v_j)}$  is assumed to be satisfied by  $\sigma$ . We now observe that this solution covers at least  $2(1 - \epsilon')n^3$  cells, as we have placed  $2(1 - \epsilon')n$ words, each of length at least  $n^2 + 2n$ , that do not pairwise intersect beyond their first 2n characters.

Suppose now we started our construction from a No instance of UNIQUE LABEL COVER. We claim that the optimal solution in the new instance cannot cover significantly more than half the cells. In particular, suppose a solution covers at least  $(1 + \epsilon')n^3 + 10n^2$  cells. We claim that the solution must have placed at least  $(1 + \epsilon')n$  words. Indeed, if we place at most  $(1 + \epsilon')n$  words, as the longest word has length  $n^2 + 4n$ , the maximum number of cells we can cover is  $(1 + \epsilon')n(n^2 + 4n) \leq (1 + \epsilon')n^3 + 4(1 + \epsilon')n^2 < (1 + \epsilon')n^3 + 10n^2$ . Let x be the number of indices  $i \in [n]$  such that the supposed solution has placed a word in both the *i*-th horizontal slot and the *i*-th vertical slot. We claim that  $x \geq \epsilon' n$ . Indeed, if  $x < \epsilon' n$ , then the total number of words we might have placed is at most  $(n - x) + 2x < (1 + \epsilon')n$ , which contradicts our previous observation that

we placed at least  $(1 + \epsilon')n$  words. Let  $V' \subseteq V$  be defined as the set of  $v_i \in V$ such that the solution places words in the *i*-th horizontal and vertical slot. Then  $|V'| \geq \epsilon'n$ . We claim that it is possible to satisfy all the constraints induced by V' in the original instance, obtaining a contradiction. Indeed, we can extract an assignment for each  $v_i \in V'$  by assigning to  $v_i$  value  $\alpha$  if the *i*-th horizontal slot contains the word  $d_{(i,\alpha)}$ . Note that the *i*-th horizontal slot must contain such a word, as these words are the only ones that have an appropriate length. Observe that in this case the *i*-th vertical slot must also contain  $d'_{(i,\alpha)}$ . Now, for  $v_i, v_j \in V'$ , with  $(v_i, v_j) \in E$  we see that  $\pi_{(v_i, v_j)}$  is satisfied by our assignment, otherwise we would have a conflict in the cell in position (2i, 2j). Therefore, in the No case, it must be impossible to fill more than  $(1 + \epsilon')n^3 + 10n^2$  cells.

The only thing that remains is to define  $\epsilon'$ . Let C be the total number of cells in the instance. Recall that we proved that in the Yes case we cover at least  $2(1-\epsilon')n^3$  cells and in the No case at most  $(1+\epsilon')n^3+10n^2$  cells. So we need to define  $\epsilon'$  such that  $2(1-\epsilon')n^3 \ge (1-\epsilon)C$  and  $(1+\epsilon')n^3+10n^2 \le (\frac{1}{2}+\epsilon)C$ . To avoid tedious calculations, we observe that  $2n^3 \le C \le 2n^3 + 8n^2$ . Therefore, it suffices to have  $2(1-\epsilon')n^3 \ge 2(1-\epsilon)(n^3+4n^2)$  and  $(1+\epsilon')n^3+10n^2 \le (1+2\epsilon)n^3$ . The first inequality is equivalent to  $(\epsilon-\epsilon')n \ge 4(1-\epsilon)$  and the second inequality is equivalent to  $(2\epsilon-\epsilon')n \ge 10$ . Since we have assumed that  $n \ge 20/\epsilon$ , it is sufficient to set  $\epsilon' = \epsilon/2$ .

#### 6. Special Cases Solvable in Polynomial Time

In this section we give some instances of CP-DEC which can be solved in polynomial time when word reuse is not allowed. Hereafter, we will always silently assume that word reuse is not allowed. We propose reductions from these instances to some well-known problems that belong to P. The first of these problems is 2-SAT which can be solved in linear time (see [5]). The second problem is the maximum matching problem. In a bipartite graph G = (V, E)this problem can be solved in  $O(\sqrt{|V|}(|E| + |V|))$  [9]. For general graphs, a much more involved algorithm by Micali and Vazirani matches the previous performance [18]. Finally, we will reduce some instances of CP-DEC to the Exact Matching problem. Karzanov in [12] proved that Exact Matching with 0-1 weights can be solved in polynomial time in complete balanced bipartite graphs. In general graphs, under the same weight restrictions, there exists an RNC algorithm by Mulmuley et al. [19], which implies that Exact Matching is solvable in randomized polynomial time – though we note that finding a deterministic polynomial time algorithm for Exact Matching is a notorious open problem.

If the grid graph of the crossword puzzle is a matching, then we will call crosses the pairs of slots that intersect. We will say that a pair of words  $(d_i, d_j)$ ,  $i \neq j$  and  $i, j \in [m]$ , indicated by their indices (i, j), can be assigned to a cross C if  $d_i$  fits into the horizontal slot of this cross,  $d_j$  fits into the vertical slot and they have the same letter in the shared cell. In the sequel, pairs (i, j) and (j, i) for any  $i \neq j$  and  $i, j \in [m]$  will count as different pairs.

**Proposition 17.** CP-DEC can solved in polynomial time for instances where all of the following conditions apply: (i) the grid consists of n/2 crosses, (ii) for any cross of the grid we have at most two pairs of words that can fit into it.

We can recognize this kind of instances in polynomial time by counting the pairs that fit in each cross as follows. For each cross  $C_k$ ,  $k \in [n/2]$ ,

- first we find the set  $\mathcal{D}_{k,h}$  of words that fit into the horizontal slot of  $\mathcal{C}_k$ ,
- for each word  $d_i \in \mathcal{D}_{k,h}$  we find the set of words  $d_j$ , different from  $d_i$ , that can fit into the vertical slot of  $\mathcal{C}_k$ , and such that  $d_i$  and  $d_j$  agree on the shared cell of  $\mathcal{C}_k$ ; let  $\mathcal{D}_{k,i,v}$  be this set,
- let  $P_k = \{(i, j) \mid d_i \in \mathcal{D}_{k,h} \text{ and } d_j \in \mathcal{D}_{k,i,v}\}$  be the set of pairs that can be assigned to  $\mathcal{C}_k$ ,
- the number of pairs that fit in  $C_k$  is:  $np(C_k) = \sum_{i:d_i \in \mathcal{D}_{k,h}} |\mathcal{D}_{k,i,v}|$ .

So we can recognize if we have this kind of instance by checking if  $np(\mathcal{C}_k) \leq 2$  for all  $k \in [n/2]$ . Furthermore, the above process is polynomial and can give us the pairs of words that can be assigned to each cross. Now, we are going to present the proof of Proposition 17.

*Proof.* First let  $\mathcal{D} = \{d_1, \ldots, d_m\}$  be the dictionary,  $\mathcal{C}_k, k \in [n/2]$ , be all the crosses of the grid and  $P_k, k \in [n/2]$ , the pairs of word indices that can be assigned to  $\mathcal{C}_k$ . Observe that if  $P_k = 0$  for some  $k \in [n/2]$ , then we know that we can not completely fill the grid. Moreover, if  $P_k = 1$ , for some  $k \in [n/2]$ , then there exists only one pair (i, j) for the cross  $\mathcal{C}_k$  so we can assign the only pair of words that fits into  $\mathcal{C}_k$  and reduce the instance by removing  $\mathcal{C}_k$  from the grid, and the words from the dictionary. Therefore, we can assume that  $P_k = 2$  for all  $k \in [n/2]$ .

From the sets  $P_k$ ,  $k \in [n/2]$ , we will construct an instance of 2-SAT that is satisfiable if and only if we can completely fill the given grid without reusing the words of the dictionary  $\mathcal{D}$ . Before we continue the construction, let us mention that 2-SAT can be solved in linear time [5].

We start by creating  $\sum_{k \in [n/2]} |P_k| = n$  variables as follows; for each  $k \in [n/2]$ and each  $(i, j) \in P_k$  we create the variable  $x_{k,i,j}$ . Let X be the set of all the variables we created. Now, we will construct a CNF formula  $\phi$  such that each clause contains at most two variables. First we add in  $\phi n/2$  clauses  $c_k, k \in [n/2]$ as follows. For each cross  $\mathcal{C}_k, k \in [n/2]$ , we add the clause  $c_k = (x_{k,i,j} \lor x_{k,i',j'})$ where (i, j) and (i', j') are in  $P_k$ .

After that, for all  $i \in [m]$  let  $X_i$  be the set of variables related to the word  $d_i$ , i.e.,  $X_i := X \cap \{x_{k,i,j}, x_{k,j,i} \mid k \in [n/2] \text{ and } d_j \in \mathcal{D}\}$ . If  $|X_i| \ge 2$ , then for any pair of variables  $x, x' \in X_i, x \ne x'$ , we add in  $\phi$  a new clause  $(\neg x \lor \neg x')$ . That process creates  $\sum_{i \in [m]} {|X_i| \choose 2} < mn^2$  additional clauses as  $|X_i| \le n$ . Therefore, the number of clauses in  $\phi$  is  $O(mn^2)$ .

It remains to prove that  $\phi$  is satisfiable if and only if we can completely fill the grid. Assume that  $\phi$  is satisfiable and  $f: X \to \{T, F\}$  is a satisfying assignment. Then, we claim that for each  $k \in [n/2]$  there exists at least one variable  $x_{k,i,j}$  such that  $f(x_{k,i,j}) = T$ . For  $k \in [n/2]$ , we select one such a variable,  $x_{k,i,j}$ , and we assign  $d_i$  horizontally and  $d_j$  vertically to the cross  $C_k$ .

Let us argue why there is always such a variable for each  $k \in [n/2]$ . By construction, for each  $k \in [n/2]$  we have a clause  $c_k$  that can be satisfied only if a variable  $x_{k,i,j}$ , such that  $(i, j) \in P_k$ , has the value T. So, we know that for each  $k \in [n/2]$  we have one such a variable. Furthermore, since  $(i, j) \in P_k$ , the words  $d_i$  and  $d_j$  fits to the horizontal and vertical slots of the cross  $C_k$  respectively.

We need to prove that we have not assigned the same word to more than one crosses. Assume that we have assigned a word  $d_i$  to two different crosses. Then, there exist two variables x and x' in  $X_i$  such that f(x) = f(x') = T. This is a contradiction because for each such pair of variables we have a clause  $(\neg x \lor \neg x')$ .

For the reverse direction we will show the following. If the grid is completely filled then the truth assignment f such that

$$f(x_{k,i,j}) = \begin{cases} T, \text{ if we have assigned the pair } (i,j) \in P_k \text{ to the cross } \mathcal{C}_k, \\ F, \text{ otherwise,} \end{cases}$$

is a satisfying assignment for the formula  $\phi$ . First, observe that for each  $k \in [n/2]$  the clause  $c_k$  is satisfied because in  $\mathcal{C}_k$  we have assigned a pair  $(i, j) \in \mathcal{P}_k$  (these are the only pairs that can be assigned to  $\mathcal{C}_k$ ) so the variable  $x_{k,i,j}$  that appears positively in  $c_k$  takes the value T by f. In order to complete the proof we need to show that at most one of the variables in  $X_i$  is true. Observe that if two of them are true then we know that we have two pairs of words containing  $d_i$  and they are assigned to a cross. This is a contradiction because we can not reuse words.

If the shared cell of a cross is the first cell of the vertical slot, then the cross is called  $\mathcal{T}$  because the slots have the form of a capital T. The next proposition shows that under some restrictions, the crossword problem can be solved efficiently if the grid is only made of  $\mathcal{T}$ 's. We remark that the grid could have an unbounded number of different kinds of  $\mathcal{T}$ 's.

**Proposition 18.** CP-DEC can be solved in polynomial time if the alphabet  $\mathcal{L}$  has only 2 letters,  $l_1$  and  $l_2$ , and the grid has the following properties: (i) it only consists of  $\mathcal{T}$ 's and (ii) all the horizontal slots have length  $\ell_h$  and all the vertical slots have length  $\ell_v$ , where  $\ell_h \neq \ell_v$ .

*Proof.* Starting from the dictionary  $\mathcal{D}$ , we can create two disjoint sub-dictionaries  $\mathcal{D}_h$  and  $\mathcal{D}_v$  containing the words that can fit into the horizontal slots and vertical slots, respectively.

Now, let  $d_i$ ,  $i \in [|\mathcal{D}_h|]$  be the words in  $\mathcal{D}_h$  and  $d'_i$ ,  $i \in [|\mathcal{D}_v|]$  be the words in  $\mathcal{D}_v$ . Furthermore, let  $m_1$  and  $m_2$  be the number of words in  $\mathcal{D}_v$  that start with  $l_1$  and  $l_2$ , respectively. Finally, let us call  $\mathcal{T}_i$ , for each  $i \in [n/2]$ , the  $\mathcal{T}$ 's of the grid.

Observe that, for any completely filled grid, we know the minimum and the maximum possible number of appearances of the symbol  $l_1$  in the shared cells. In particular, since  $m_2$  words of  $\mathcal{D}_v$  do not start with  $l_1$  we have at least  $w_1 = \max\{0, n/2 - m_2\}$  appearances and since  $m_1$  words of  $\mathcal{D}_v$  start with  $l_1$ we have at most  $w_2 = \min\{n/2, m_1\}$  appearances. So we need to decide if there exists a way to fill completely the horizontal slots of the grid that forces w appearances of  $l_1$  in the shared cells for a  $w \in \{w_1, \ldots, w_2\}$ .

Based on  $\mathcal{D}_h$  and the  $\mathcal{T}$ 's we will construct an instance of the EXACT MATCHING problem. First, we create a complete balanced bipartite graph G = (V, U, E) where  $V = \{v_1, \ldots, v_{|\mathcal{D}_h|}\}$  represents the words in  $\mathcal{D}_h$ , vertices  $u_j \in U$  where  $j \in [n/2]$  represent the  $\mathcal{T}$ 's of the grid and vertices  $u_j \in U$  such that  $j \in \{n/2 + 1, \ldots, |\mathcal{D}_h|\}$  are added so that G is balanced.

It remains to assign weights to the edges. For each edge  $(v_i, u_j) \in E$ , if  $i \in [|\mathcal{D}_h|], j \in [n/2]$  and  $d_i$  contains symbol  $l_1$  in the position of the shared cell of the horizontal slot of  $\mathcal{T}_i$ , then  $(v_i, u_j)$  has weight 1, otherwise its weight is 0.

Now it is not difficult to observe that G has a perfect matching of weight exactly w if and only if we can fill all the horizontal slots of the grid with words of  $\mathcal{D}_h$  such that  $l_1$  appears exactly w times in the shared cells. Finally, we can decide if G has a perfect matching of weight w, for any  $w \in \{w_1, \ldots, w_2\}$ , in polynomial time by using Theorem 1 in [12].

**Remark 3.** The same proof works if instead of  $\mathcal{T}$ 's (where two slots intersect only in the first position of the vertical slot), one has crosses such that any two slots that intersect do so only at the same position of the vertical slot.

At first sight, the case covered by Proposition 18 looks restricted but it necessitates to match three objects (1 horizontal word, 1 vertical word, and a cross) and 3D-matchings are, in general, hard problems [6].

**Proposition 19.** CP-DEC can be solved in polynomial time if the instance has the following properties: (i) the grid graph is a matching and (ii) the number of different types of crosses<sup>6</sup> is a constant.

*Proof.* First, we will prove that if we know the number of appearances of each letter in the shared cells of each type of crosses, then we can find a way to fill the grid (if there exists one) using the maximum matching problem.

Let  $t \in \mathbb{N}$  be the number of different types of crosses and let  $\mathcal{L} = \{l_1, \ldots, l_\ell\}$ be the alphabet. Furthermore, suppose that the given instance has a solution and that we are given values  $a_{i,j}, i \in [\ell]$  and  $j \in [t]$  which indicate the number of appearances of the letter  $l_i$  in the shared cell of type j crosses in this solution. In the end we will repeat the algorithm for all combinations of such values, so we can assume that these values are given to us in the input. We will therefore look for a solution that agrees with the given values  $a_{i,j}$ .

Now, we construct the following bipartite graph.

 $<sup>^{6}</sup>$ Two crosses are of the same type if they are identical: same number of horizontal cells, same number of vertical cells, and same shared cell.

- For each pair  $(i, j) \in [\ell] \times [t]$  we create a set  $V_{i,j} = \{v_{i,j,k}, v'_{i,j,k} \mid k \in [a_{i,j}]\}$  of  $2a_{i,j}$  vertices. Furthermore, let V be the set  $\bigcup_{(i,j) \in [\ell] \times [t]} V_{i,j}$ .
- For each word  $d \in \mathcal{D}$  we create a vertex  $u_d$ . Let U be the set of these vertices.
- For all  $d \in \mathcal{D}$  and  $(i, j) \in [\ell] \times [t]$ , if the word d fits in the horizontal slot of a cross of type j and it forces the letter  $l_i$  in the shared cell of this slot, then we add all the edges  $(u_d, v_{i,j,k})$  where  $k \in [a_{i,j}]$ .
- Finally, for all  $d \in \mathcal{D}$  and  $(i, j) \in [\ell] \times [t]$ , if the word d fits in the vertical slot of a cross of type j and it forces the letter  $l_i$  in the shared cell of this slot, then we add all the edges  $(u_d, v'_{i,j,k})$  where  $k \in [a_{i,j}]$ .

Now, we claim that this graph has a matching that covers all the vertices in V if and only if we can completely fill the grid in a way that respects the given appearances of the letters. That is, for each letter  $i \in [\ell]$  and each type  $j \in [t]$  there are exactly  $a_{i,j}$  crosses of type j that contain i in the shared cell in the solution.

In one direction, assume that we have an assignment of words in the grid that respects the given appearances of the letters. For each  $(i, j) \in [\ell] \times [t]$  let  $\mathcal{D}_{i,j}$  be the set of words that have been assigned to the crosses of type j and these words force the letter  $l_i$  in the shared cell. Observe that  $\mathcal{D}_{i,j}$  has size  $2a_{i,j}$ as it respects the appearances of the letters and for each  $d \in \mathcal{D}_{i,j}$  the vertex  $u_d$ is adjacent to  $v_{i,j,k}$ , for all  $k \in [a_{i,j}]$ , if d has been assigned to a horizontal slot or is adjacent to  $v'_{i,j,k}$ , for all  $k \in [a_{i,j}]$  if d has been assigned to a vertical slot.

Now we are going to create a matching of the graph. Starting from an empty set S, for each  $(i, j) \in [\ell] \times [t]$  and each word d in  $\mathcal{D}_{i,j}$ , if d has been assigned to a horizontal slot, then we add in S an edge  $(u_d, v_{i,j,k})$  for some k in  $[a_{i,j}]$  that covers an uncovered vertex  $v_{i,j,k}$ ; otherwise we add an edge  $(u_d, v'_{i,j,k})$  for some k in  $[a_{i,j}]$  that covers an uncovered vertex  $v'_{i,j,k}$  (we know that there are enough vertices by the previous observations).

Observe that S is a matching because each time we add an edge incident to two uncovered vertices. Furthermore, because the size of  $\mathcal{D}_{i,j}$  is  $2a_{i,j}$  and we have exactly  $a_{i,j}$  horizontal and  $a_{i,j}$  vertical slots we know that we will cover all the  $v_{i,j,k}$  with the corresponding vertices of the horizontal slots and all the  $v'_{i,j,k}$ with the corresponding vertices of the vertical slots. Hence, S is a matching that covers all the vertices of V.

For the reverse direction, assume that S is a matching that covers all the vertices of V. Because V is an independent set we know that for all  $v \in V$  there exists an edge  $(u_d, v) \in S$  for some word  $d \in \mathcal{D}$  and  $u_d \in U$ . We will assign words to the slots of the grid as follows. For each edge  $(u_d, v_{i,j,k})$  we assign the word d to an empty horizontal slot of a cross of type j. Because all  $v_{i,j,k}$ ,  $k \in [a_{i,j}]$ , are adjacent to vertices  $u_d$  such that d fits in the horizontal slots of crosses of type j and forcing the letter  $l_i$  in the shared cells, we have filled all the horizontal slots in a way that respects the given appearances of the letters.

Now, we are going to fill the vertical slots. For each edge  $(u_d, v'_{i,j,k})$  we assign the word d to an empty vertical slot of a cross of type j where the shared cell has the letter  $l_i$ . Because for a given pair  $(i, j) \in [\ell] \times [t]$  we have  $a_{i,j}$  vertices  $v_{i,j,k}$  and  $v'_{i,j,k}$  we know that we have forced exactly  $a_{i,j}$  times the letter  $l_i$  in the shared cell of crosses of type j, so we have the exact number of vertical slots we want.

Finally, because the total number of vertices in V is the same as the number of slots in the grid, we know that we have completely filled the grid.

Now, note that we can decide in polynomial time if there exists such a matching (by finding a maximum matching).

In order to complete the proof we need to show that the different guesses for the number of appearances  $a_{i,j}$ ,  $(i,j) \in [\ell] \times [t]$  is polynomially bounded by the input size. Fix a type j and suppose it has  $\nu_j$  crosses. Enumerating all the possible  $a_{i,j}$ 's is equivalent to choosing  $\ell - 1$  positions in a row vector of size  $\nu_j + \ell - 1$  (the chosen cells "separate" the  $a_{i,j}$ 's). Then, there are  $\binom{\nu_j + \ell - 1}{\ell - 1}$ choices, which is upper bounded by  $\binom{n+\ell-1}{\ell-1}$ . Finally, because we have t types of crosses, the total number of guesses is

Finally, because we have t types of crosses, the total number of guesses is  $O(n^{t\ell})$  which is polynomial as t and  $\ell$  are constant.

So far we have assumed that the size of the alphabet is constant. In contrast, the next two propositions hold without this hypothesis. Therefore, there are independent of Proposition 19.

**Proposition 20.** CP-DEC can be solved in polynomial time when the grid consists of identical crosses. This holds even if the size of the alphabet is not constant.

*Proof.* Let  $\mathcal{D}$  be the dictionary of size m and let n/2 be the number of crosses in the grid. In order to answer the question we will create a graph G on mvertices such that G has a matching of size at least n/2 if and only if we can completely fill the grid. We construct G = (V, E) as follows: V has exactly one vertex  $v_i$  for each word  $d_i \in \mathcal{D}$  and E contains the edge  $(v_i, v_j)$  if and only if at least one of the pairs of words (i, j) or (j, i) fits in a cross.

Now, observe that an edge of G gives us a pair of words that can fit into a cross. Therefore, each matching S represents pairs of words that (all of them) fit in crosses. Because the words are represented by vertices in the graph, the pairs we take from a matching are independent as each vertex is covered at most once by S. Finally, if the matching has at least n/2 edges then we know that we have enough pairs to fill the grid completely.

Conversely, if we can cover the grid completely, then there exist n/2 distinct pairs of words that fit in a cross. For each of these pairs we have an edge in G. Furthermore, the set of these edges is a matching as we have not used the same word twice so there are no two edges incident to the same vertex. Hence, G has a matching of size at least n/2.

Because we can find a maximum matching of a graph in polynomial time we know that we can decide if we can fill the whole grid in the same time.  $\Box$ 

**Proposition 21.** CP-DEC is polynomial time reducible to 0-1 Exact Matching if the grid has the following properties: (i) it is a matching and (ii) the number

of different types of crosses is a constant. This holds even if the size of the alphabet is not constant.

*Proof.* Let  $t \in \mathbb{N}$  be the number of different types of crosses. Let  $a_k$ , for  $k \in [t]$ , be the number of crosses of type t in the grid. We assume w.l.o.g. that the number of words  $|\mathcal{D}| = m$  is even. We first create a multi-graph where parallel edges are allowed, and we assign weights to the edges from the set  $\{0, m^0, \ldots, m^{t-1}\}$  as follows:

- Let G = (V, E) be a complete graph with m vertices and give weight 0 to each edge.
- For each  $k \in [t]$ , add a set of edges  $E_k$  such that  $(v_i, v_j) \in E_k$  if one of the pairs of words (i, j) or (j, i) fits in a cross of type  $k \in [t]$ . Give weight  $m^{k-1}$  to all the edges in  $E_k$ .

In order to distinguish the parallel edges, denote by  $(v_i, v_j)_k$  the edge  $(v_i, v_j) \in E_k$ . Now, we claim that this graph has a perfect matching of weight  $W = \sum_{k \in [t]} a_k m^{k-1}$  if and only if we can fill the grid completely.

Assume that the grid is filled and construct a perfect matching as follows. Start from an empty set S. For each cross, add to S the edge  $(v_i, v_j)_k$  where (i, j) is the pair of words allocated to the cross and k is the type of the cross. Since reusing words is not allowed, we know that S is a matching. Furthermore, because each cross of type k has an edge of  $E_k$ , we get that S has weight W. Finally, S can be extended to a perfect matching of the same weight because all the vertices have edges of weight 0 between them.

Conversely, assume that the graph has a perfect matching S of weight W. Observe that any perfect matching has exactly m/2 edges. Therefore we know that we have at most m/2 edges from each set  $E_k$ , for any  $k \in [t]$ . Now we are going to show that we need exactly  $a_k$  edges from each set  $E_k$ . In particular, we claim that for any  $k \in \{1, \ldots, t\}$  we have that  $|S \cap E_k| = a_k$ . We prove this by induction.

**Base Case** (t): Assume that there exists a perfect matching S with weight W such that  $|S \cap E_t| = b_t \neq a_t$ . Assume that  $b_t < a_t$ . Since  $0 \leq b_t$ , we get that  $1 \leq a_t$ . Even if we use a maximum number of edges of maximum weight (i.e., m/2 edges of weight  $m^{t-2}$ ), we cannot compensate, i.e., S cannot have weight W because  $m^{t-2} \cdot m/2 < m^{t-1}$ . Now, we have to check the case  $b_t > a_t$ . Since S has weight at least  $(a_t + 1)m^{t-1}$ , we have  $W \geq (a_t + 1)m^{t-1}$ . That gives us:

$$\sum_{k \in [t]} a_k m^{k-1} \ge (a_t + 1) m^{t-1} \Rightarrow \sum_{k \in [t-1]} a_k m^{k-1} \ge m^{t-1} \Rightarrow m^{t-2} \sum_{k \in [t-1]} a_k \ge m^{t-1}$$

which is a contradiction since  $\sum_{k \in [t-1]} a_k \leq m/2$ . Therefore, for any perfect matching of weight W it must be  $|S \cap E_t| = a_t$ .

**Induction Hypothesis**  $(\{k, k+1, \ldots, t\}, k > 1)$ : We assume that, for a given k > 1,  $|S \cap E_i| = a_i$  holds for all  $i \in \{k, k+1, \ldots, t\}$ .

**Induction Step** (k-1): For every  $l \in \{k, \ldots, t\}$  we know that  $|S \cap E_l| = a_l$ . Therefore, the set  $S' = S \setminus \bigcup_{l=k}^{t} (S \cap E_l)$  has weight exactly  $\sum_{l=1}^{k-1} a_l m^{l-1}$ . Assume that k - 1 = 1; then  $\sum_{l=1}^{k-1} a_l m^{l-1} = a_1$ . Since S' consists only of edges in  $E_1 \cup E$ , the edges in  $E_1$  have weight 1 and the edges in E have weight 0, so we can conclude that S' contains exactly  $a_1$  edges from  $E_1$ .

Now, we consider the case k-1 > 1. Similarly to the base case, assume that  $|S' \cap E_{k-1}| = b_{k-1}$  and  $b_{k-1} < a_{k-1}$ . In this case, even if we use a maximum number of edges (i.e., m/2) of maximum weight (i.e.,  $m^{k-2}$ ), we cannot compensate, i.e.,  $m^{k-2} \cdot m/2 < m^{k-1}$ . Now we prove that, assuming  $b_{k-1} > a_{k-1}$  leads to a contradiction. Indeed, in this case S' has weight at least  $(a_{k-1}+1)m^{k-2}$  which implies that:

$$\sum_{l \in [k-1]} a_l m^{l-1} \ge (a_{k-1} + 1) m^{k-2} \Rightarrow \sum_{l \in [k-2]} a_l m^{l-1} \ge m^{k-2} \Rightarrow m^{k-3} \sum_{l \in [k-2]} a_k \ge m^{t-1}$$

The last inequality contradicts the fact that  $\sum_{l \in [k-2]} a_k \leq \sum_{l \in [t]} a_t \leq m/2$ . Hence,  $b_{k-1} = a_{k-1}$ .

Finally, observe that for any  $k \in [t]$  each edge in  $S \cap E_k$  gives us a pair of words that fits into a cross of type k. Moreover, because S is a matching and  $|S \cap E_k| = a_k$  for all  $k \in [t]$ , we can completely fill the grid.

Now, starting from the graph we created in the first reduction, we will create an instance of 0-1 Exact Matching, using the same technique as Papadimitriou and Yannakakis in their proof of Proposition 1 of [20]. We build the new graph G' by replacing each edge  $(v_i, v_j)_k$ ,  $k \in [t]$ , with a path,  $p_{i,j,k} = \langle v_i u_{i,j,1}^k \dots u_{i,j,2m^{k-1}}^k v_j \rangle$ . Assign weight 0 to the edges  $(u_{i,j,2m^{k-1}}^k, v_j)$  and  $(u_{i,j,l}^k, u_{i,j,l+1}^k)$ , where  $l < 2m^{k-1}$  is odd. All the other edges of the path have weight 1.

Now, let  $E_{i,j,k}$  be the edges of a path  $p_{i,j,k}$ . Observe that for any perfect matching  $M \subseteq E'$  of G' and any path  $p_{i,j,k}$ , the set of edges  $M \cap E_{i,j,k}$ :

- either, contains edges of total weight 0 and it does not cover any endpoint of  $p_{i,j,k}$ ,
- or, it contains edges of total weight  $m^{k-1}$  and it covers both  $v_i$  and  $v_j$ .

Based on this observation we can transform any perfect matching of G into a perfect matching of G' with the same weight and vice versa.

Finally, since G' has  $O(m^t)$  vertices, the starting instance of CP-DEC is polynomial time reducible to 0-1 Exact Matching.

We can also use the technique of Papadimitriou and Yannakakis [20] in the proof of Proposition 18, giving us Corollary 22 which more general than Proposition 18.

**Corollary 22.** CP-DEC with alphabet of constant size is polynomial time reducible to 0-1 Exact Matching if the grid has the following properties: (i) it is a matching, (ii) all shared cells are at the same position of the vertical slots, (iii) all vertical slots have the same length  $l_v$  and (iv) there is no horizontal slot of length  $l_v$ .

The details are omitted but observe that we don't need to construct a complete bipartite graph and we can allow weights up to  $m^k$ , where  $m = O(|\mathcal{D}|)$ and  $k = O(|\mathcal{L}|)$ , before we apply the aforementioned technique.

#### 7. Conclusion

We studied the parameterized complexity of some crossword puzzles under several different parameters and we gave some positive results followed by proofs which show that our algorithms are essentially optimal. Based on our results the most natural questions that arise are:

- What is the complexity of CP-DEC when the grid graph is a matching and the alphabet has size 2?
- Can Theorem 13 be strengthened by starting from ETH instead of randomized ETH?
- Can we beat the  $(\frac{1}{2} + O(\frac{1}{n}))$ -approximation ratio of CP-OPT if we restrict our instances?
- Can Theorem 15 be strengthened by dropping the UGC?

Finally, as a future work, we could consider a variation of the crossword puzzle problems where each word can be used a given number of times. This would be an intermediate case between word reuse and no word reuse.

#### Acknowledgements

We thank Dominique Taton for communicating the puzzle to us.

#### References

- Anbulagan and Adi Botea. Crossword puzzles as a constraint problem. In Principles and Practice of Constraint Programming, 14th International Conference, CP 2008, Sydney, Australia, September 14-18, 2008. Proceedings, pages 550–554, 2008.
- [2] Edward K. Crossman and Sharyn M. Crossman. The crossword puzzle as a teaching tool. *Teaching of Psychology*, 10(2):98–99, 1983.
- [3] Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michal Pilipczuk, and Saket Saurabh. *Parame*terized Algorithms. Springer, 2015.
- [4] Jakob Engel, Markus Holzer, Oliver Ruepp, and Frank Sehnke. On computer integrated rationalized crossword puzzle manufacturing. In Fun with Algorithms - 6th International Conference, FUN 2012, Venice, Italy, June 4-6, 2012. Proceedings, pages 131–141, 2012.
- [5] Shimon Even, Alon Itai, and Adi Shamir. On the complexity of timetable and multicommodity flow problems. SIAM J. Comput., 5(4):691–703, 1976.
- [6] Michael R. Garey and David S. Johnson. Computers and Intractability: A Guide to the Theory of NP-Completeness. W. H. Freeman, 1979.

- [7] Matthew L. Ginsberg, Michael Frank, Michael P. Halpin, and Mark C. Torrance. Search lessons learned from crossword puzzles. In Proceedings of the 8th National Conference on Artificial Intelligence. Boston, Massachusetts, USA, July 29 - August 3, 1990, 2 Volumes, pages 210–215, 1990.
- [8] Laurent Gourvès, Ararat Harutyunyan, Michael Lampis, and Nikolaos Melissinos. Filling crosswords is very hard. In 32nd International Symposium on Algorithms and Computation, ISAAC 2021, December 6-8, 2021, Fukuoka, Japan, pages 36:1–36:16, 2021.
- [9] John E. Hopcroft and Richard M. Karp. An n<sup>5/2</sup> algorithm for maximum matchings in bipartite graphs. SIAM J. Comput., 2(4):225–231, 1973.
- [10] Heather Hulett, Todd G. Will, and Gerhard J. Woeginger. Multigraph realizations of degree sequences: Maximization is easy, minimization is hard. Oper. Res. Lett., 36(5):594–596, 2008.
- [11] Russell Impagliazzo, Ramamohan Paturi, and Francis Zane. Which problems have strongly exponential complexity? J. Comput. Syst. Sci., 63(4):512–530, 2001.
- [12] Alexander V. Karzanov. Maximum matching of given weight in complete and complete bipartite graphs. *Kibernetika*, 23(1):7–11, (1987) (English translation in CYBNAW23(1), 8–13 (1987)).
- [13] Subhash Khot and Oded Regev. Vertex cover might be hard to approximate to within 2-epsilon. J. Comput. Syst. Sci., 74(3):335–349, 2008.
- [14] Michael Lampis, Valia Mitsou, and Karolina Soltys. Scrabble is PSPACEcomplete. J. Inf. Process., 23(3):284–292, 2015.
- [15] Michael L. Littman, Greg A. Keim, and Noam M. Shazeer. Solving crosswords with PROVERB. In Proceedings of the Sixteenth National Conference on Artificial Intelligence and Eleventh Conference on Innovative Applications of Artificial Intelligence, July 18-22, 1999, Orlando, Florida, USA, pages 914–915, 1999.
- [16] Télé Magazine. Publications Grand Public.
- [17] Gary Meehan and Peter Gray. Constructing crossword grids: Use of heuristics vs constraints. In In: Proceedings of Expert Systems 97: Research and Development in Expert Systems XIV, SGES, pages 159–174, 1997.
- [18] Silvio Micali and Vijay V. Vazirani. An o(sqrt(|v|) |e|) algorithm for finding maximum matching in general graphs. In 21st Annual Symposium on Foundations of Computer Science, Syracuse, New York, USA, 13-15 October 1980, pages 17–27. IEEE Computer Society, 1980.
- [19] Ketan Mulmuley, Umesh V. Vazirani, and Vijay V. Vazirani. Matching is as easy as matrix inversion. Comb., 7(1):105–113, 1987.

- [20] Christos H. Papadimitriou and Mihalis Yannakakis. The complexity of restricted spanning tree problems. J. ACM, 29(2):285–309, 1982.
- [21] Jagan A. Pillai, Charles B. Hall, Dennis W. Dickson, Herman Buschke, Richard B. Lipton, and Joe Verghese. Association of crossword puzzle participation with memory decline in persons who develop dementia. *Journal* of the International Neuropsychological Society, 17(6):1006–1013, 2011.
- [22] Leonardo Rigutini, Michelangelo Diligenti, Marco Maggini, and Marco Gori. Automatic generation of crossword puzzles. Int. J. Artif. Intell. Tools, 21(3), 2012.
- [23] Christopher D. Rosin. Nested rollout policy adaptation for Monte Carlo tree search. In IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011, pages 649–654, 2011.