

The 2D Shortest Superstring Problem

Dat Thanh Tran^a, Khai Quang Tran^a and Van Khu Vu^{a,*}

^aVinUniversity, Hanoi, Vietnam

ARTICLE INFO

Keywords:

Shortest superstring problem
Two-dimensional strings
Genetic algorithm
Integer linear programming
Combinatorial optimization
Bounded-offset tree encoding

ABSTRACT

We introduce the Two-Dimensional Shortest Superstring Problem (2D-SSP): arrange rectangular symbol arrays on the integer lattice with symbol-consistent overlaps to minimize bounding-box area or square side length. We prove NP-hardness for both objectives and APX-hardness for area via L-reduction from 1D-SSP. A *Bounded-Offset Tree Representation* reduces 2D-SSP to search over spanning trees, justified by a connectivity lemma ensuring optimal solutions can be made 4-connected. We design a Tree-Based Genetic Algorithm with locality-preserving crossover that recombines subtrees rather than coordinates. Experiments show near-optimality (gap $\leq 2.6\%$) against exact solvers and 6–12% improvement over greedy baselines. The theoretical framework extends to d dimensions.

1. Introduction

The *Two-Dimensional Shortest Superstring Problem* (2D-SSP) generalizes the classical Shortest Superstring Problem (Gallant, Maier and Storer, 1980) from 1D strings to 2D arrays. Given rectangular 2D strings T_1, \dots, T_n over a finite alphabet, the goal is to place them in the integer plane with *symbol-consistent overlaps*, overlapping cells must contain identical symbols, while minimizing the bounding-box cost. We study two objectives: *area* ($W \cdot H$) and *balanced* ($\max\{W, H\}$). Figure 1 illustrates a small optimal placement.

The transition from 1D to 2D introduces three qualitative changes:

1. *Multi-directional overlap*. In 1D-SSP, strings overlap from two directions (left/right). In 2D-SSP, overlap can occur from four directions, dramatically increasing the number of pairwise configurations.
2. *Holes*. A 2D placement can have interior cells not covered by any string (see “*” in Figure 1). These “holes” inflate the bounding box and have no 1D analogue.
3. *Sequencing meets packing*. Optimal placement depends on global geometric configuration, not just local overlaps. This creates a hybrid problem combining aspects of stringology (overlap exploitation) and bin packing (spatial arrangement).

Contributions. We make the following contributions:

1. **Problem formalization** (Section 3). We define 2D-SSP precisely, including placement feasibility, symbol consistency, and the two bounding-box objectives.
2. **Complexity results** (Section 4). We prove NP-hardness for both objectives and APX-hardness for the area objective via L-reduction from 1D-SSP.
3. **Compaction theorem** (Theorem 5). Every optimal placement can be transformed into a 4-connected placement without increasing cost. This is analogous to compaction in VLSI floorplanning but requires a new proof to handle symbol-consistency constraints.
4. **Bounded-Offset Tree Representation** (Section 5.2). We show that connected placements correspond to spanning trees with bounded edge labels (Lemma 3), and that at least one optimal placement admits such a representation (Corollary 10). Combined with a global coordinate bound (Lemma 4), this yields a finite (though exponential) combinatorial search space.
5. **Tree-based Genetic Algorithm** (Section 6.4). We design a GA whose individuals are placement trees and whose crossover transplants subtrees, preserving spatially coherent clusters. On small instances ($n \leq 10$), T-GA achieves optimality gaps $\leq 2.6\%$ versus CPLEX; on larger instances, it outperforms greedy baselines by 6–12%.

*Corresponding author

✉ dat.tt3@vinuni.edu.vn (D.T. Tran); khai.tq@vinuni.edu.vn (K.Q. Tran); khu.vv@vinuni.edu.vn (V.K. Vu)
ORCID(s):

Scope and limitations. The tree encoding yields a *finite* search space in the sense that spanning-tree topologies and edge labels are bounded. However, this space remains *exponential* in n ; our exact MIP baseline does not scale beyond $n \approx 10$. The GA provides a practical heuristic but with no worst-case guarantees. Our experiments use binary alphabets; behavior on higher-entropy instances is unexplored.

Paper organization. Section 2 reviews related work. Section 3 defines 2D-SSP. Section 4 establishes complexity. Section 5 develops structural properties and the tree representation. Section 6 presents algorithms. Section 7 reports experiments. Section 8 concludes.

2. Background and Related Work

2.1. Shortest Superstring Problem

In the classical SSP, the input is a set of strings

$$S = \{s_1, \dots, s_n\}$$

over an alphabet Σ . A superstring is a string S in which each s_i appears as a substring. The objective is to minimize $|S|$. SSP is NP-hard, and there is a substantial literature on constant-factor approximations (e.g., greedy maximum-overlap merging, cycle-cover-based algorithms) and heuristic implementations used in practice.

2.2. 2D Covers and 2D Covering Sequences

Two-dimensional generalizations of string concepts appear in several areas:

- *2D covers and 2D strings.* Work on covers of 2D arrays considers how a small pattern can cover a larger 2D string with overlaps, generalizing the notion of a cover in 1D (Charalampopoulos, Pissis, Radoszewski, Waleń and Zuba, 2021).
- *Covering sequences and 2D covering sequences.* Recent research introduces covering sequences and covering 2D-sequences, where all $m \times n$ windows of a large 2D array form a covering code for patterns of that size up to a given radius. These provide natural sources of structured test instances (Yehezkeally and Schwartz, 2025).

These works focus on covering combinatorial spaces, whereas we focus on overlapping a *given finite set* of 2D strings with exact symbol consistency.

2.3. 2D Bin Packing and Cutting Stock

Classical two-dimensional bin packing problems ask how to place a collection of rectangles into one or more rectangular bins so as to minimize, for example, the number of bins used or the height of a single strip, under strict *non-overlap* constraints. The items are unlabeled shapes. A closely related problem is the *two-dimensional cutting stock problem* (Gilmore and Gomory, 1965; Lodi, Martello and Monaci, 2002), where the goal is to cut rectangular pieces from large stock sheets while minimizing waste. Both problems have been extensively studied using exact methods (branch-and-bound, column generation) and metaheuristics (genetic algorithms, simulated annealing) (Bennell and Oliveira, 2009).

Our setting is similar in that we also optimize a global bounding box for a family of rectangular pieces, but differs in two key ways. First, each string is a *discrete symbol array* rather than an unlabeled rectangle; second, *overlaps are allowed* as long as they are *symbol-consistent*. Thus a solution is not simply a packing of shapes, but a combinatorial “gluing” of patterns in which overlaps can reduce the effective occupied area, a phenomenon absent from standard 2D bin packing and cutting stock formulations. In the language of cutting stock, 2D-SSP allows “pieces” to share material when their patterns match, a constraint that transforms the problem from pure geometry to a hybrid of sequencing and packing.

2.4. Related Geometric and Assembly Problems

2D-SSP shares structural similarities with several problems in combinatorial optimization. In VLSI floorplanning, topological representations such as Sequence Pairs (Murata, Fujiyoshi, Nakatake and Kajitani, 1996a) and B*-trees (Chang, Chang, Wu and Wu, 2000) encode relative module positions; our placement trees adapt this approach to content-based adjacency. Patch-based texture synthesis (Efros and Freeman, 2001; Kwatra, Schödl, Essa, Turk and

Bobick, 2003) places patches to minimize visual error in overlaps; 2D-SSP is the discrete, lossless limit requiring exact symbol consistency. The Tile Assembly Model (Winfree, Liu, Wenzler and Seeman, 1998) studies self-assembly of Wang tiles with edge-matching constraints; 2D-SSP asks the inverse question of finding the most compact configuration containing a given set of patterned tiles. Unlike jigsaw puzzle assembly (Sholomon, David and Netanyahu, 2013) or polyomino packing (Golomb, 1994), 2D-SSP permits symbol-consistent overlaps that enable compression beyond pure geometric packing.

3. Problem Definition

This section formalizes the Two-Dimensional Shortest Superstring Problem (2D-SSP) and introduces the two objective variants.

Let Σ be a finite alphabet over which the 2D strings are defined. A *2D-string* over Σ is a finite $m \times n$ array $T \in \Sigma^{m \times n}$, for some $m, n \in \mathbb{N}$. For indices $1 \leq i \leq j \leq m$ and $1 \leq i' \leq j' \leq n$, we write $T[i..j, i'..j']$ for the corresponding subarray and call this a *2D-substring* of T .

We identify a 2D string T with a function on a finite index set $C_T \subset \mathbb{Z}^2$, its set of *local cell coordinates*. We write cells as pairs $(u, v) \in \mathbb{Z}^2$ and use the same coordinate system for both local and global positions: a translation by an offset $p(i) = (x_i, y_i)$ sends a local cell (u, v) of T_i to the global cell $p(i) + (u, v) = (x_i + u, y_i + v)$. The choice of which axis is drawn horizontally or vertically is irrelevant for our arguments; we only rely on coordinate-wise addition in \mathbb{Z}^2 .

Let P be an $m' \times n'$ 2D-string. We denote the set of its occurrences in T by

$$\text{Occ}(P, T) = \{(i, j) : T[i..i + m' - 1, j..j + n' - 1] = P\}.$$

We will derive cost functions from the dimensions of the minimal axis-aligned bounding rectangle of a placement, and consider two variants: one that minimizes the area and one that minimizes the maximum side length (square objective).

Definition 1. Let $\mathcal{T} = \{T_1, \dots, T_n\}$ be a finite set of 2D strings over Σ , which we call *2D strings*. An $m \times n$ 2D-string S is a *2D-superstring* of \mathcal{T} if each string T_i occurs as a 2D-substring of S , i.e.,

$$\text{Occ}(T_i, S) \neq \emptyset \quad \text{for all } i \in \{1, \dots, n\}.$$

Note that, under our bounding-box objectives, an optimal arrangement may leave cells inside the bounding box that are not covered by any input string (“holes”). Formally this causes no mismatch with $S \in \Sigma^{m \times n}$: any such uncovered cell is *unconstrained* by the feasibility conditions and may be assigned an arbitrary symbol from Σ without affecting whether the T_i occur as substrings. We denote by

$$|S|_{\text{area}} := m \cdot n \quad \text{and} \quad |S|_{\text{sq}} := \max\{m, n\}$$

the *area* and the *square side length* of S , respectively. The area measure is the natural 2D analogue of string length in 1D-SSP, while the square measure seeks the smallest enclosing square. Both are derived from the minimal axis-aligned rectangle containing S and penalize any empty cells inside that rectangle.

Thus $|S|_{\text{area}}$ is the primary objective, the natural 2D analogue of superstring length (total cells in the bounding box), while $|S|_{\text{sq}}$ is the side length of the smallest enclosing square, ensuring neither dimension dominates.

Remark 1 (Holes inside the bounding box). The “*” marker in Figure 1 highlights a fundamental distinction between 1D-SSP and 2D-SSP. In 1D-SSP, every position in an optimal superstring must be covered by some input string: if any position were uncovered, we could delete that character to obtain a shorter superstring, contradicting optimality.

In 2D-SSP, this property fails. The two-dimensional geometry permits *holes*: cells within the bounding box that no input string covers. These holes cannot simply be “deleted” as in 1D, because removing a row or column would disrupt the geometric arrangement of strings in other parts of the superstring. Hole cells contribute to the bounding-box cost but carry no information: they are unconstrained and can be filled arbitrarily with symbols from Σ without affecting which strings are embedded.

Formally, given a placement p , let $R(p) = \bigcup_{i=1}^n \text{footprint}(T_i, p)$ denote the set of occupied cells, and let $B(p)$ denote the bounding box. The *hole set* is $B(p) \setminus R(p)$, cells inside the bounding box but not covered by any string. In Figure 1, this set contains two cells (the “*” positions). These holes represent “wasted” area that inflates the objective value, and minimizing them is part of the optimization challenge unique to 2D-SSP.

1	1	1	0	0	0	*
0	1	0	1	0	1	*
1	0	0	1	1	0	0
1	0	1	1	1	1	1
1	0	0	1	0	1	0

Figure 1: A 2D-superstring S (5 rows \times 7 columns) containing six 2D strings T_1, \dots, T_6 as 2D-substrings (highlighted rectangles). The cells printed as “*” are *uncovered* positions (holes) inside the bounding box that are not covered by any input string; “*” is a *diagrammatic placeholder* (not an alphabet symbol). In any formal $S \in \Sigma^{m \times n}$ these cells can be filled arbitrarily with symbols from Σ . The area is $|S|_{\text{area}} = 35$, and the square side length is $|S|_{\text{sq}} = 7$.

Definition 2. Given a finite set \mathcal{T} of 2D strings over Σ , we define two variants of the Two-Dimensional Shortest Superstring Problem:

- *Area-based 2D-SSP* (2D-SSP_{area}): find a 2D-superstring S of \mathcal{T} minimizing $|S|_{\text{area}}$. This is the primary variant, directly generalizing the 1D-SSP objective.
- *Square 2D-SSP* (2D-SSP_{sq}): find a 2D-superstring S of \mathcal{T} minimizing $|S|_{\text{sq}} = \max\{m, n\}$, the side length of the smallest enclosing square.

Both objectives depend only on the minimal axis-aligned bounding rectangle of S and penalize all empty cells inside it.

4. Computational Complexity

We establish that 2D-SSP is NP-hard. Both proofs are straightforward reductions; we state them for completeness.

Theorem 1 (NP-hardness of 2D-SSP_{area}). *2D-SSP_{area} is NP-hard, even for binary alphabets.*

Proof. Reduction from 1D-SSP (NP-hard for $|\Sigma| \geq 2$ (Gallant et al., 1980)). Given 1D strings $S = \{s_1, \dots, s_n\}$, create height-1 2D strings $T_i = 1 \times |s_i|$. Any 2D placement with k rows and maximum row width W has area $\geq k \cdot W$. Concatenating rows yields a 1D superstring of length $\leq k \cdot W$. Conversely, the optimal 1D superstring gives a 1-row 2D placement. Thus $\text{OPT}_{2D} = \text{OPT}_{1D}$. \square

Theorem 2 (NP-hardness of 2D-SSP_{sq}). *2D-SSP_{sq} is NP-hard.*

Proof. Reduction from minimum enclosing square packing (NP-hard (Leung, Tam, Wong, Young and Chin, 1990)). Given rectangles R_1, \dots, R_n , assign each a unique symbol σ_i from alphabet Σ with $|\Sigma| = n$. The resulting 2D strings cannot overlap (symbol conflict), so the problem reduces to non-overlapping rectangle packing. \square

Remark 2 (Complexity summary and open problems). *Area objective:* NP-hard and APX-hard for any $|\Sigma| \geq 2$ (inherited from 1D-SSP via L-reduction with $\alpha = \beta = 1$; see Appendix A.1).

Square objective: NP-hard, but Theorem 2 requires a large alphabet ($|\Sigma| = n$). This places the proof in the *packing regime* where no overlaps occur.

Open problem: *The complexity of 2D-SSP_{sq} over a fixed constant-size alphabet (e.g., binary) remains open.* In the sequencing regime where overlaps are common, neither the 1D-SSP reduction nor the packing reduction applies directly.

4.0.1. Experimental scope

Our experiments focus on the *sequencing regime*: binary alphabets ($|\Sigma| = 2$). Extended discussion of the sequencing–packing spectrum and entropy effects is deferred to Appendix A.2.

Whenever the distinction between the two variants is not important, we simply refer to either as *2D-SSP*. Hereafter, we use *string* to mean *2D string* unless otherwise specified.

The 2D Shortest Superstring Problem

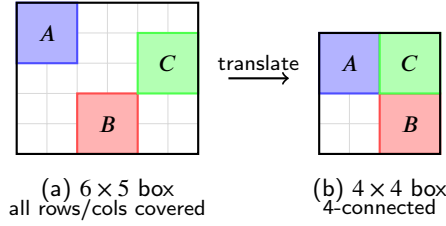


Figure 2: Row/column coverage is necessary but not sufficient for 2D optimality. (a) Both strings cover all rows and columns of the 6×5 bounding box, yet they are disconnected. (b) Translating B diagonally yields a 4-connected placement with a smaller 4×4 bounding box.

We assume throughout that strings are axis-aligned rectangles and cannot be rotated or reflected. Following standard 1D-SSP convention, we assume the input is *substring-free*: no string T_i is a 2D-substring of another string T_j . Redundant strings can be removed in polynomial-time preprocessing without affecting the optimal solution, since any 2D-superstring containing T_j automatically contains T_i .

Rather than working directly with the superstring S , we use placements on the integer grid.

Definition 3. A *placement* of \mathcal{T} is a function

$$p : \{1, \dots, n\} \rightarrow \mathbb{Z}^2, \quad p(i) = (x_i, y_i),$$

assigning an integer offset to each string T_i . A cell of T_i with local coordinates (u, v) (row and column indices) is mapped to global coordinates $(x_i + u, y_i + v) \in \mathbb{Z}^2$.

A placement p is symbol-consistent if for every global coordinate $(x, y) \in \mathbb{Z}^2$, all strings covering (x, y) under p write the same symbol. Note that 4-adjacent contact imposes no symbol constraints; only overlapping cells must agree. We denote by

$$R(p) := \{(x, y) \in \mathbb{Z}^2 : (x, y) \text{ is covered by some } T_i \text{ under } p\}$$

the union of occupied global cells, and let $B(p)$ be the minimal axis-aligned rectangle containing $R(p)$. Let $W(p)$ and $H(p)$ be the width and height of $B(p)$, and define

$$\text{cost}_{\text{area}}(p) := W(p) \cdot H(p), \quad \text{cost}_{\text{sq}}(p) := \max\{W(p), H(p)\}.$$

Restricting the symbol map to $B(p)$ yields an $m \times n$ array T_p ; by construction T_p is a 2D-superstring of \mathcal{T} , and its area and balanced side satisfy

$$|T_p|_{\text{area}} = \text{cost}_{\text{area}}(p), \quad |T_p|_{\text{sq}} = \text{cost}_{\text{sq}}(p).$$

In particular, empty cells of $B(p)$ that are not covered by any strings still contribute to both objectives. Thus every symbol-consistent placement defines a feasible solution to both $\text{2D-SSP}_{\text{area}}$ and $\text{2D-SSP}_{\text{sq}}$, with cost equal to the chosen bounding-box functional.

Remark 3. We minimize bounding box area $W \cdot H$ rather than union area $|R(p)|$ for two reasons: (i) it is the natural 2D generalization of 1D superstring length, and (ii) minimizing $|R(p)|$ would permit fragmented layouts with “holes,” whereas bounding-box cost penalizes such fragmentation and encourages compact, connected arrangements.

Remark 4. In 1D-SSP, every position in an optimal superstring must be covered by some input string. In 2D, the analogous condition, every row and column intersects some string, is necessary but not sufficient (Figure 2). The connectivity/compaction theorem (Theorem 5) provides the appropriate 2D analogue by ensuring optimal placements can be made 4-connected.

Conversely, let S be a 2D-superstring of \mathcal{T} and fix an arbitrary occurrence $(i, j) \in \text{Occ}(T_k, S)$ for each string T_k . Placing T_k with offset (i, j) then yields a symbol-consistent placement whose induced array is S up to a global translation. Hence optimizing over 2D-superstrings is equivalent to optimizing over symbol-consistent placements, modulo a global shift of all coordinates. We therefore work with placements from now on.

A subset $R \subseteq \mathbb{Z}^2$ is *4-connected* if its adjacency graph under the 4-neighbourhood ($\|x - y\|_1 = 1$) is connected. Two cells $x, y \in \mathbb{Z}^2$ overlap if $x = y$.

Definition 4. Given the string set \mathcal{T} , the *placement graph* G^{pl} has vertex set $\{1, \dots, n\}$. Its edges are triples

$$e = (i, j, \delta) \quad \text{with} \quad i \neq j, \delta \in \mathbb{Z}^2,$$

where δ is a relative offset such that placing T_j at position $p(j) = p(i) + \delta$ yields symbol-consistent contact between T_i and T_j , meaning they either overlap with matching symbols, or are 4-adjacent (share a boundary edge). Multiple edges may exist for the same unordered pair $\{i, j\}$, corresponding to different valid offsets.

Intuitively, each edge (i, j, δ) in G^{pl} specifies a way of gluing T_j next to T_i . The placement graph is determined entirely by \mathcal{T} and can be precomputed before searching for solutions.

Lemma 3. Let T_i and T_j have bounding boxes of dimensions $w_i \times h_i$ and $w_j \times h_j$. If (i, j, δ) is an edge in G^{pl} , then $\delta = (\Delta x, \Delta y)$ satisfies:

$$|\Delta x| \leq w_i + w_j - 1, \quad |\Delta y| \leq h_i + h_j - 1.$$

Proof. For T_i and T_j to be in contact (overlap or 4-adjacent), their bounding boxes must intersect or share an edge. The x-projections $[0, w_i - 1]$ and $[\Delta x, \Delta x + w_j - 1]$ intersect or are adjacent if $|\Delta x| \leq w_i + w_j - 1$. The bound for Δy follows symmetrically. \square

Lemma 4 (A polynomially describable global bound). Let $w_i \times h_i$ be the dimensions of T_i and define

$$W_0 := \sum_{i=1}^n w_i, \quad H_0 := \sum_{i=1}^n h_i.$$

There exists a symbol-consistent placement p_0 such that, after translating its bounding box to the origin, $B(p_0) \subseteq [0, W_0 - 1] \times [0, H_0 - 1]$. In particular,

$$\text{OPT}_{\text{area}} \leq W_0 H_0 \quad \text{and} \quad \text{OPT}_{\text{sq}} \leq \max\{W_0, H_0\}.$$

Consequently, there exists an optimal placement p^* whose bounding box (after translating to the origin) lies inside $[0, U - 1] \times [0, U - 1]$ for the polynomially describable bound

$$U := \begin{cases} W_0 H_0, & \text{for } 2D\text{-SSP}_{\text{area}}, \\ \max\{W_0, H_0\}, & \text{for } 2D\text{-SSP}_{\text{sq}}. \end{cases}$$

Proof. Construct p_0 by placing the strings without overlap, e.g., set $p_0(i) = (\sum_{t < i} w_t, \sum_{t < i} h_t)$. Then the translated footprints are pairwise disjoint, so symbol-consistency holds vacuously. The union of footprints is contained in $[0, W_0 - 1] \times [0, H_0 - 1]$. Thus $\text{cost}_{\text{area}}(p_0) \leq W_0 H_0$ and $\text{cost}_{\text{sq}}(p_0) \leq \max\{W_0, H_0\}$, implying the stated upper bounds on the optimum. Finally, for any objective, translating an optimal placement so that its bounding box has lower-left corner at the origin yields a congruent optimal placement; since its cost is at most the corresponding bound above, both width and height are at most U , hence the translated bounding box lies within $[0, U - 1] \times [0, U - 1]$. \square

Remark 5 (Graph size vs. tree space). Lemma 3 ensures that the *placement graph* G^{pl} has polynomial size: for each pair (T_i, T_j) , there are $O((w_i + w_j)(h_i + h_j))$ candidate offsets to check. For uniform $w \times h$ strings, the placement graph has $O(n^2 wh)$ edges.

However, the *search space of spanning trees* remains exponential. By Cayley's formula, a complete graph on n vertices has n^{n-2} spanning trees; while our placement graph is typically sparser, the number of feasible placement trees can still grow exponentially with n .

The main significance of the tree view is symmetry-breaking and structure: it represents solutions via *relative offsets* along a connected backbone rather than absolute coordinates, and enables structure-aware search operators (Section 6). Formally, one can also bound absolute coordinates a priori by a polynomially describable U (Lemma 4), so the set of placements up to translation is finite even in a coordinate encoding; the tree encoding remains valuable because it removes translational redundancy and exposes useful combinatorial structure.

We defer two performance-oriented discussions to the appendix: how string entropy affects edge density in G^{pl} (Appendix A.3) and how to enumerate offsets efficiently in practice (Appendix A.4).

Definition 5. Let p be a symbol-consistent placement of \mathcal{T} . The *contact graph* $G^{\text{ct}}(p)$ is the subgraph of G^{pl} induced by p : it has vertex set $\{1, \dots, n\}$, and edge (i, j, δ) is present if $\delta = p(j) - p(i)$ and this edge exists in G^{pl} .

Equivalently, strings i and j are adjacent in $G^{\text{ct}}(p)$ if they are in contact under p (overlapping or 4-adjacent).

The contact graph $G^{\text{ct}}(p)$ records which edges of the placement graph are “realized” by a given placement. Different placements of the same instance may realize different subsets of edges.

Remark 6. Our cost functionals depend only on the bounding rectangle $B(p)$, not directly on the cardinality of $R(p)$. In particular, two placements with the same bounding box have the same cost for both 2D-SSP_{area} and 2D-SSP_{sq}. The area-based objective directly generalizes the 1D shortest superstring objective (string length), while the square objective seeks a compact, near-square layout. Both differ from geometric covering formulations that minimize the area of the union $R(p)$.

Assumption 1. In the structural discussion below we work with placements whose occupied region $R(p)$ is 4-connected. The following theorem shows this is without loss of optimality.

5. Structural Properties and Tree Representation

This section develops the *Bounded-Offset Tree Representation* that underlies our algorithms. We first prove a compaction theorem (Section 5.1) showing that optimal placements can be made 4-connected, then introduce the placement tree encoding (Section 5.2) and summarize algorithmic implications (Section 5.3).

5.1. Compaction and Connectivity

The following theorem shows that optimal solutions can be restricted to connected placements, a result analogous to compaction in VLSI floorplanning (Murata, Fujiyoshi, Nakatake and Kajitani, 1996b).

Theorem 5 (Connectivity/compaction). *Let p be an optimal symbol-consistent placement for 2D-SSP under either objective. Then there exists an optimal symbol-consistent placement p' with $\text{cost}(p') \leq \text{cost}(p)$ such that $R(p')$ is 4-connected.*

The proof relies on two lemmas. The key observation is that symbol-consistency constrains only *overlapping* cells, not 4-adjacent ones. Thus we can slide disconnected components closer until they become 4-adjacent without creating overlaps, hence without risking symbol conflicts.

Lemma 6 (Unit shifts are safe). *Let p be a symbol-consistent placement with occupied region $R = R(p)$, and let C be a maximal 4-connected component of R . If $\text{dist}(C, R \setminus C) \geq 2$, then for any unit vector $e \in \{(\pm 1, 0), (0, \pm 1)\}$, shifting all strings in C by e preserves symbol-consistency.*

Proof. Suppose $(C + e) \cap (R \setminus C) \neq \emptyset$. Then there exist $c \in C$ and $d \in R \setminus C$ with $c + e = d$, so $\|c - d\|_1 = 1$, contradicting $\text{dist}(C, R \setminus C) \geq 2$. Thus the shifted component does not overlap any other component, and since the shift is rigid, symbol-consistency within C is preserved. \square

Lemma 7 (Merge move exists). *Let p be a symbol-consistent placement whose occupied region $R(p)$ has $k \geq 2$ maximal 4-connected components. Then there exist two components C_a, C_b and a unit vector e such that:*

- (i) *shifting C_b by e keeps $C_b + e$ within the same bounding box $B(R(p))$, and*
- (ii) *$\text{dist}(C_a, C_b + e) = \text{dist}(C_a, C_b) - 1$.*

Proof. Pick a closest pair (C_a, C_b) minimizing $\text{dist}(C_i, C_j)$ among all component pairs, with $d := \text{dist}(C_a, C_b) \geq 2$. Choose $u \in C_a$ and $v \in C_b$ with $\|u - v\|_1 = d$. Let e be a unit step from v toward u along a shortest Manhattan path, so $\|u - (v + e)\|_1 = d - 1$.

We show e can be chosen so that $C_b + e \subseteq B(R(p))$. Consider the case $e = (-1, 0)$ (left step). If $x_{\min}(C_b) > x_{\min}(B)$, then shifting C_b left by one keeps it within B . If $x_{\min}(C_b) = x_{\min}(B)$, then C_b is flush with the left boundary. But since $u \in C_a \subseteq B$ and $v \in C_b \subseteq B$, there exists at least one coordinate direction where u differs from v and the corresponding unit step from v toward u remains within B . The same reasoning applies symmetrically to other directions.

Finally, Lemma 6 ensures the move creates no overlaps (all components are at distance $\geq d \geq 2$ from C_b), so symbol-consistency is preserved. \square

Proof of Theorem 5. Among all optimal placements, choose one q that minimizes the number k of 4-connected components of $R(q)$.

Suppose $k \geq 2$. By Lemma 7, there exist components C_a, C_b and a unit shift e that keeps $C_b + e$ within the same bounding box (so cost does not increase) and reduces the distance between C_a and C_b by one. Repeating this shift at most $\text{dist}(C_a, C_b) - 1$ times makes C_a and C_b become 4-adjacent, at which point they merge into a single component.

This yields an optimal placement with strictly fewer than k components, contradicting minimality of k . Hence $k = 1$ and $R(q)$ is 4-connected. Set $p' := q$. \square

Remark 7. This proof does not rely on planar topology (crossing lemmas). The argument extends directly to d dimensions: components separated by grid distance ≥ 2 can be slid together one unit step at a time without creating overlaps, since $2d$ -adjacency (contact) is distinct from overlap.

Corollary 8. *For any instance of 2D-SSP, there exists an optimal placement whose occupied region is 4-connected. Consequently, Assumption 1 is justified for all instances.*

This connectivity/compaction theorem establishes that restricting attention to connected placements loses no optimal solutions. The assumption matches our experimental focus and enables the tree-based structural perspective developed in this section. It is not required for the correctness of the algorithms in Section 6, which operate on arbitrary symbol-consistent placements.

Corollary 9 (Connected occupied region implies connected contact graph). *Let p be a symbol-consistent placement such that $R(p)$ is 4-connected. Then the contact graph $G^{\text{ct}}(p)$ is connected.*

Proof. If $G^{\text{ct}}(p)$ were disconnected, we could partition the strings into two nonempty sets A and B with no edge between them. For each string T_i , let R_i be the set of global cells covered by T_i under p . Define

$$R_A = \bigcup_{i \in A} R_i, \quad R_B = \bigcup_{j \in B} R_j,$$

so that $R(p) = R_A \cup R_B$ and $R_A \cap R_B = \emptyset$. By construction there is no pair of 4-adjacent cells across R_A and R_B , which contradicts 4-connectivity of $R(p)$. \square

By Theorem 5 and Corollary 9, every optimal connected placement has a connected contact graph that admits a spanning tree. We now formalize this tree representation.

5.2. Placement Trees

A placement tree encodes a solution as a spanning tree with labeled edges specifying relative offsets between strings. The key insight is that edge labels are drawn from a *bounded* set determined by string dimensions, which, combined with the finite number of spanning tree topologies, yields a finite search space.

Definition 6. A *placement tree* for \mathcal{T} is a rooted tree $F = (V, E)$ with vertex set $V = \{1, \dots, n\}$ together with, for each edge $\{i, j\} \in E$, a label $\delta_{ij} \in \mathbb{Z}^2$ interpreted as the relative offset from i to j . For each oriented edge (i, j) we store δ_{ij} and require $\delta_{ji} = -\delta_{ij}$.

By Lemma 3, valid edge labels are bounded: if T_i and T_j have dimensions $w_i \times h_i$ and $w_j \times h_j$, then any symbol-consistent contact requires $|\Delta x| \leq w_i + w_j - 1$ and $|\Delta y| \leq h_i + h_j - 1$. This bounds the number of candidate labels per edge to $O((w_i + w_j)(h_i + h_j))$, for uniform $w \times h$ strings, $O(wh)$ per edge, ensuring the search space of placement trees is finite.

The *realization* of F with root r and root position $p(r) \in \mathbb{Z}^2$ is the placement p_F defined by

$$p_F(i) = p(r) + \sum_{t=0}^{\ell-1} \delta_{v_t, v_{t+1}},$$

where $r = v_0, v_1, \dots, v_\ell = i$ is the unique simple path from r to i in F .

We call F *feasible* if the realization p_F is symbol-consistent for some (equivalently, every) choice of root position.

Remark 8 (Placements vs. Trees). Not every symbol-consistent placement corresponds to a connected tree. If the occupied region $R(p)$ is disconnected, the contact graph $G^{\text{ct}}(p)$ is a forest of multiple components, not a single spanning tree. However, by Theorem 5, *at least one optimal placement* has a connected occupied region, and thus corresponds to a spanning tree. This is the crux of the search space reduction: we need not enumerate all placements (including disconnected ones), but only spanning trees of the placement graph.

Remark 9. Feasibility is a global property: even if all adjacent pairs (i, j) are locally consistent, collisions may occur between distant parts of the tree when their footprints overlap after summing the offsets. In particular, closed walks in the underlying contact graph induce non-trivial “loop-closure” constraints on the offsets. We treat feasibility algorithmically: given a placement tree F , we realize it via p_F and explicitly check symbol-consistency.

The following theorem establishes that optimal solutions can always be represented as placement trees, a crucial fact that justifies restricting our search to tree-based encodings.

Corollary 10 (Tree representation of optimal placements). *The following statements hold:*

- (i) **(Existence of tree-optimal solutions.)** *For any instance of 2D-SSP, there exists an optimal placement p^* that corresponds to a feasible placement tree. Specifically, by Theorem 5 we may assume $R(p^*)$ is 4-connected; any spanning tree F of $G^{\text{ct}}(p^*)$, equipped with edge labels $\delta_{ij} := p^*(j) - p^*(i)$, is then a feasible placement tree whose realization coincides with p^* up to global translation.*
- (ii) **(Completeness.)** *Conversely, any feasible placement tree F induces a symbol-consistent placement p_F that is a valid solution to 2D-SSP.*

Proof. (i) *Optimal Placement \rightarrow Tree.* Let p^* be any optimal placement. By Theorem 5, there exists an optimal placement p with the same cost such that $R(p)$ is 4-connected. By Corollary 9, $G^{\text{ct}}(p)$ is connected, so it admits a spanning tree F . For any edge (i, j) of F , we have $p(j) = p(i) + \delta_{ij}$ by definition of the labels. For any vertex i with path $r = v_0, \dots, v_\ell = i$ from root r :

$$p(i) = p(r) + \sum_{t=0}^{\ell-1} \delta_{v_t, v_{t+1}} = p_F(i).$$

Thus p equals the realization p_F (for root position $p(r)$), which is symbol-consistent, so F is feasible and optimal.

(ii) *Tree \rightarrow Placement.* By feasibility, p_F is symbol-consistent. Every string T_i is placed exactly once, so T_i occurs as a 2D-substring of the induced array T_{p_F} , making it a 2D-superstring of \mathcal{T} . \square

The significance of this corollary is that *no optimal solution is lost* by restricting to the Bounded-Offset Tree Representation. While not every placement corresponds to a tree, disconnected placements correspond to forests of multiple trees, Theorem 5 ensures that at least one optimal solution has a connected contact graph, and hence corresponds to a single spanning tree. Moreover, the bounded edge labels (Lemma 3) and finite tree topologies yield a finite set of candidate *encodings*, and the objective implies an explicit global bound on the size of an optimal bounding box (Lemma 4). This justifies designing algorithms that search exclusively over placement trees: the representation is *complete* (contains at least one optimum), symmetry-reduced, and it exposes a useful combinatorial structure for search.

Remark 10. All results in this section concern only symbol-consistency, connectivity, and combinatorial structure. They apply verbatim to both 2D-SSP_{area} and 2D-SSP_{sq}.

Remark 11 (Extension to higher dimensions). The theoretical framework (connectivity lemma, bounded offsets, tree representation) extends directly to d -dimensional SSP for any $d \geq 1$, with $2d$ -adjacency replacing 4-adjacency. See Appendix A.5 for details.

Table 1

Summary of main notation.

Symbol	Meaning
Σ	Alphabet
$\mathcal{T} = \{T_1, \dots, T_n\}$	Input set of strings
$C_T \subset \mathbb{Z}^2$	Local cell coordinates of string T
P	Pattern 2D-string (for occurrences)
$\text{Occ}(P, T)$	Set of occurrences of P in T
S	2D-superstring of \mathcal{T}
$ S _{\text{area}}$	Area of minimal bounding box of S
$ S _{\text{sq}}$	Square side length $\max\{m, n\}$ of S
$p(i) = (x_i, y_i)$	Placement (offset) of string T_i
$R(p) \subset \mathbb{Z}^2$	Union of occupied cells under placement p
$B(p)$	Minimal axis-aligned bounding box of $R(p)$
$W(p), H(p)$	Width and height of $B(p)$
$\text{cost}_{\text{area}}(p)$	Area-based cost $W(p)H(p)$
$\text{cost}_{\text{sq}}(p)$	Square cost $\max\{W(p), H(p)\}$
$G^{\text{ct}}(p)$	Contact graph induced by placement p
G^{pl}	Placement graph of symbol-consistent offsets
$F = (V, E)$	Placement tree on $\{1, \dots, n\}$
$\delta_{ij} \in \mathbb{Z}^2$	Relative offset from i to j in F
p_F	Realization (tree-induced placement) of F

5.3. Algorithmic Implications

Corollary 10 has two key algorithmic consequences. First, it guarantees that *optimal solutions are always reachable* via tree-based search: since at least one optimal placement corresponds to a feasible placement tree, any algorithm that exhaustively searches over placement trees is guaranteed to find an optimum. Second, the bounded offset property (Lemma 3) together with the objective-implied global bound (Lemma 4) justify focusing on a finite, discrete solution family: the tree encoding has finitely many topologies and finitely many admissible edge labels, and (up to translation) an optimal placement can be assumed to lie within a finite coordinate window. While the number of spanning trees can be exponential in n (see Remark 5), this replaces an unbounded coordinate description by a structured combinatorial one that is amenable to exact methods on small instances and structure-aware search on larger ones. Our GA searches only a restricted subspace of trees induced by greedy completion; the resulting existence-vs.-reachability gap is discussed in Appendix A.6.

The tree structure has three properties making it suitable for structure-aware search: (1) *locality preservation*, subtrees correspond to spatially coherent clusters; (2) *incremental realizability*, trees can be grown one string at a time; and (3) *symmetry breaking*, relative offsets factor out global translation, yielding a non-redundant search space.

These properties motivate our choice to design a genetic algorithm whose individuals are placement trees and whose crossover operates on subtrees, rather than using a more conventional coordinate-based representation.

6. Algorithms

Having established the structural foundation linking placements and trees, we present three algorithmic approaches for solving 2D-SSP. These methods span the spectrum from exact to structure-aware, offering different trade-offs between solution quality and computational cost:

- An *exact CPLEX MIP formulation* (Section 6.1) that enumerates all candidate placements on a discrete grid and optimizes over them using mixed-integer programming (MLP for the balanced objective and MIQP for the area objective). This approach guarantees optimal solutions but is limited to small instances.
- A *merge-based greedy heuristic* (Section 6.2) adapted from the classical 1D Shortest Superstring algorithm. This baseline repeatedly merges pairs of partial superstrings with maximum overlap.

- A *tree-growing greedy heuristic* (Section 6.3) that builds a placement tree incrementally, motivated by Corollary 10. At each step it attaches a new string to minimize the bounding-box cost.
- A *tree-based genetic algorithm* (Section 6.4) that represents individuals as placement trees and uses crossover operators that recombine subtrees. By exploiting the relative-offset encoding, this approach aims to combine the quality of exact methods with the scalability of structure-aware search.

All three methods can be instantiated with either of the two bounding-box cost variants, 2D-SSP_{area} or 2D-SSP_{sq}.

The genetic algorithm uses placement trees rather than the more obvious coordinate-based representation (where each individual is a vector of (x, y) coordinates for each string). This choice is motivated by several considerations. First, in the coordinate representation, an offspring produced by crossover rarely inherits good local structure from its parents: if parent 1 places strings A and B in a well-overlapping configuration, and parent 2 places strings B and C similarly, a crossover that takes A from parent 1 and C from parent 2 will likely place them far apart, destroying both favorable overlaps. In contrast, our tree-based crossover transplants entire subtrees, preserving the relative offsets among all strings in the subtree. Second, the coordinate representation is highly redundant: any global translation of a placement yields the same objective value, so the search space contains infinitely many equivalent solutions. The tree representation eliminates this redundancy by encoding only the pairwise offsets that matter. Third, the tree structure aligns naturally with the connectivity requirement: a spanning tree automatically ensures that all strings are geometrically connected, whereas the coordinate representation requires additional constraints or repair operators to enforce connectivity.

6.1. Exact Verification via CPLEX MIP

To validate the solution quality of our heuristic approaches, we formulate a direct grid-based mixed-integer program (Wolsey, 1998). We emphasize that this formulation is *not* intended as a scalable solver for general instances, but strictly as a *ground-truth oracle* for small-scale verification ($N \leq 10$). This allows us to measure exactly how close our genetic algorithm comes to the global optimum.

The model works with discrete candidate placements of each 2D string on a finite grid: it enumerates a finite set of allowed origins \mathcal{O}_i for each 2D string T_i , uses binary decision variables $b_{io} \in \{0, 1\}$ to choose exactly one origin per string, precomputes pairwise conflict indicators $\kappa_{ijoo'}$ to forbid symbol-inconsistent placements, and minimizes the bounding-box cost via big- M constraints that track the enclosing rectangle. For the balanced objective, we minimize $L = \max\{W, H\}$ via linear constraints. For the area objective, we use CPLEX to solve the resulting Mixed-Integer Quadratic Program (MIQP) directly with quadratic objective $W \cdot H$.

The formulation incorporates several refinements to improve tractability: (i) greedy-based grid bounds that accommodate all optimal aspect ratios while remaining much smaller than naïve worst-case bounds; (ii) symmetry breaking by fixing the first string at the origin; and (iii) instance-adaptive big- M constants that tighten the LP relaxation. The complete mathematical formulation, including all decision variables, constraints, and objective functions for both the balanced and area variants, is provided in Appendix B.

Despite these optimizations, the number of variables and conflict constraints grows quickly with the number of strings. Consequently, this exact MIP is practically limited to small instances ($N \leq 10$) and serves exclusively as a verification tool to certify the optimality gap of our heuristic methods.

6.2. Merge-Based Greedy Heuristic

Before presenting our tree-growing heuristic, we describe a natural baseline adapted from the classical 1D Shortest Superstring Problem: the *merge-based greedy* algorithm.

In the 1D setting, the Greedy Superstring Algorithm (Blum, Jiang, Li, Tromp and Yannakakis, 1994) is remarkably effective: it repeatedly merges the pair of strings with maximum overlap until a single superstring remains. This simple strategy achieves a 4-approximation for 1D-SSP (Blum et al., 1994), later improved to 2.5 (Kaplan and Shafrir, 2005). Empirically, the algorithm performs far better than these worst-case bounds suggest: extensive experiments on both random and biological sequences show approximation ratios consistently below 1.05 (Cazaux and Rivals, 2018), making it a strong practical baseline.

Role in our evaluation. We include merge-greedy not as a strawman but for two principled reasons. First, it provides *scalability comparison*: merge-greedy runs in polynomial time and scales to arbitrarily large instances, whereas our exact ILP baseline is limited to $n \leq 10$ strings. Second, its near-optimal performance on 1D-SSP ($< 5\%$ gap) makes it a

calibration baseline: if our 2D algorithms cannot substantially outperform a method that is nearly optimal in 1D, then the 2D structure is not being exploited. Our primary validation of solution quality, however, comes from comparison against exact ILP solutions on tractable instances (Table 4), where convergence to the global optimum is verifiable.

We adapt this approach to the 2D setting as follows. Given a set \mathcal{T} of 2D strings, we maintain a set of *partial superstrings* (initially, each string is its own partial superstring). At each step, we identify the pair (S_i, S_j) of partial superstrings that can be merged with maximum symbol-consistent overlap, merge them into a single partial superstring, and repeat until only one remains.

Crucially, the overlap function measures *geometrical overlap*, the area of the overlapping region, rather than the number of matching non-wildcard characters:

$$\text{overlap}(S_i, S_j) = \max_{\delta \in D_{ij}} |R_i \cap (R_j + \delta)|,$$

where R_i, R_j are the cell sets of S_i, S_j , the offset δ ranges over all translations that yield symbol-consistent overlap, and $|\cdot|$ denotes the cardinality of the intersection (i.e., the number of overlapping cells, regardless of whether those cells contain alphabet symbols or wildcards).

Geometrical vs. character overlap. In the 1D setting, maximizing geometrical overlap and maximizing character overlap are equivalent: every overlapping position contains exactly one character. In the 2D setting, however, these objectives diverge. Consider two 3×3 strings that can overlap in two ways: (a) a 2×2 region with 4 matching characters, or (b) a 1×3 strip with 3 matching characters. A character-based criterion would prefer (a), but this choice might force the merged component into an unfavorable shape that increases the final bounding box.

We use geometrical overlap because it directly controls the *union area* of occupied cells when merging two partial superstrings. Although our optimization objective is the *bounding-box* area (which is not determined solely by pairwise overlap), overlap provides a fast, geometry-aware proxy that often correlates with smaller bounding boxes.

Lemma 11 (Overlap maximization minimizes union-area increase). *Let $R_i, R_j \subset \mathbb{Z}^2$ be the occupied cell sets of two partial superstrings S_i and S_j . For any translation (offset) $\delta \in \mathbb{Z}^2$, the occupied set after placing S_j at offset δ is*

$$R_{\text{merged}}(\delta) := R_i \cup (R_j + \delta).$$

Then

$$|R_{\text{merged}}(\delta)| = |R_i| + |R_j| - |R_i \cap (R_j + \delta)|.$$

Consequently, for fixed R_i and R_j , maximizing the geometrical overlap $|R_i \cap (R_j + \delta)|$ is equivalent to minimizing the union size $|R_{\text{merged}}(\delta)|$, and equivalently to minimizing the union-area increase

$$\Delta_{\text{union}}(\delta) := |R_{\text{merged}}(\delta)| - |R_i|.$$

Proof. The identity follows from inclusion–exclusion: $|A \cup B| = |A| + |B| - |A \cap B|$ applied to $A = R_i$ and $B = R_j + \delta$. Since $|R_i|$ and $|R_j|$ do not depend on δ , minimizing $|R_{\text{merged}}(\delta)|$ (or $\Delta_{\text{union}}(\delta)$) is equivalent to maximizing $|R_i \cap (R_j + \delta)|$. \square

This observation motivates our choice of overlap function in merge-greedy. In contrast, for the bounding-box objectives $W \cdot H$ and $\max\{W, H\}$, overlap is only a heuristic signal: two merges with the same overlap can yield different bounding boxes due to shape effects.

Hypothesis: 1D vs. 2D performance. We hypothesize that the merge-based greedy performs well when the input strings are *nearly one-dimensional* (i.e., have aspect ratios close to $1 \times m$ or $m \times 1$), since this regime closely resembles the classical 1D setting where the algorithm has provable guarantees. However, as strings become *genuinely two-dimensional* (aspect ratios closer to 1), the merge-based approach may struggle: the “best overlap” criterion optimizes locally for overlap size but ignores how the merge affects the global bounding-box shape. In contrast, our tree-growing heuristic (Section 6.3) explicitly optimizes the bounding-box cost at each step, which we expect to yield better solutions for 2D instances.

The procedure is summarized in Algorithm 1.

We test this hypothesis experimentally in Section 7: Table ?? compares merge-based greedy against our methods on instances with 1×8 strings (nearly 1D), while Tables 2 and ?? evaluate performance on 3×3 strings (genuinely 2D).

Algorithm 1 Merge-Based Greedy for 2D-SSP

Require: Set of 2D strings $\mathcal{T} = \{T_1, \dots, T_n\}$
Ensure: 2D superstring S

```

1:  $S \leftarrow \{T_1, \dots, T_n\}$  ▷ Set of partial superstrings
2: while  $|S| > 1$  do
3:    $(S_i^*, S_j^*, \delta^*) \leftarrow \arg \max_{S_i, S_j \in S, \delta} \text{overlap}(S_i, S_j, \delta)$  ▷ Find best symbol-consistent overlap
4:    $S_{\text{merged}} \leftarrow \text{MERGE}(S_i^*, S_j^*, \delta^*)$  ▷ Create merged superstring
5:    $S \leftarrow (S \setminus \{S_i^*, S_j^*\}) \cup \{S_{\text{merged}}\}$ 
6: end while
7: return the single element of  $S$ 
    
```

6.3. Tree-Growing Greedy Heuristic

Motivated by Corollary 10, which establishes that every connected placement corresponds to a spanning tree of relative offsets, we design a greedy heuristic that directly constructs a *placement tree*. Rather than merging pairs of partial superstrings as in the merge-based approach, we grow a single tree by iteratively attaching strings to the current structure.

The algorithm maintains:

- A *placement tree* $F = (V_F, E_F)$ with $V_F \subseteq \{1, \dots, n\}$ representing the strings placed so far;
- A *canvas* of occupied global cells with their symbols;
- The bounding-box coordinates $(x_{\min}, x_{\max}, y_{\min}, y_{\max})$ of the current placement.

At each iteration, we select an unplaced string $T_j \notin V_F$ and attach it to some string $T_i \in V_F$ via an edge (i, j, δ) from the placement graph G^{pl} , choosing the pair (i, j, δ) that minimizes the resulting bounding-box cost while maintaining symbol-consistency. This directly mirrors the structural insight of Corollary 10: we are constructing an optimal placement by “growing” a spanning tree one edge at a time.

Crucially, by Lemma 3, we need only consider offsets δ within the bounded window $|\Delta x| \leq w_i + w_j - 1$ and $|\Delta y| \leq h_i + h_j - 1$. This transforms what might seem like an unbounded search into a tractable enumeration: for each candidate parent T_i in the current tree and each unplaced string T_j , we examine $O((w_i + w_j)(h_i + h_j))$ candidate attachment positions.

The current width and height are

$$w(C) = x_{\max} - x_{\min} + 1, \quad h(C) = y_{\max} - y_{\min} + 1,$$

and we define either

$$\text{size}_{\text{sq}}(C) = \max\{w(C), h(C)\} \quad \text{or} \quad \text{size}_{\text{area}}(C) = w(C) \cdot h(C)$$

depending on the chosen objective.

For any candidate attachment of a string, we can (i) check whether overlapping cells agree symbolically, and (ii) compute the resulting bounding-box cost.

For a given target side length L in the square objective case, we enumerate translations $(\Delta x, \Delta y)$ of a string T that would result in a bounding box of side length exactly L when combined with the current canvas. Intuitively:

- if the canvas is empty, we place the first string so that the bounding box matches its own width/height;
- otherwise, when L equals the current size s , we slide the string in all ways that keep the bounding box within a virtual $s \times s$ square;
- when $L > s$, we consider placements that extend this square along one of its four sides so that the new size becomes exactly L .

Symbol-consistency of these candidate placements is checked against the canvas. For the area-based variant we similarly enumerate candidate positions and evaluate them according to the area-based cost.

6.3.1. Deterministic Tree-Growing

The deterministic variant starts from a chosen root string T_r (placed at the origin) and iteratively attaches the remaining strings. At each step, it considers all candidate edges (i, j, δ) where T_i is already in the tree and T_j is not, and selects the edge that:

- (i) achieves the smallest possible increase in the chosen bounding-box cost, and
- (ii) among those, maximizes the *character overlap*, the number of overlapping cells containing matching non-wildcard symbols, as a tie-breaker.

Rationale for the two-level criterion. The primary criterion explicitly minimizes the increase in the chosen bounding-box cost (area or balanced side length). This is the objective we ultimately care about, and it is *not* determined solely by pairwise overlap because global shape effects (e.g., creating long thin layouts) matter.

Among candidate attachments that yield the same bounding-box cost, we use *character overlap* as a tie-breaker. This reflects the fact that, while overlap maximization is exactly equivalent to minimizing the *union-area* increase (Lemma 11), it is only a secondary signal for bounding-box objectives. Preferring higher character overlap among cost-ties encourages tighter interlocking and tends to reduce holes in practice.

The algorithm terminates when all n strings have been added to the tree. By construction, it always increases the cost threshold until a consistent attachment exists for some string, and therefore returns a complete placement tree spanning all strings.

This procedure directly implements the “growing” perspective of Corollary 10: we construct a spanning tree of the contact graph by adding one vertex at a time, always maintaining a valid partial placement.

6.3.2. Stochastic Tree-Growing

The stochastic variant follows the same tree-growing structure but introduces randomization in the selection among equally good candidate edges. Instead of deterministically picking the single best attachment, it uses *roulette wheel selection* among all candidates that achieve the minimal cost increase, with selection probabilities proportional to their *character overlap* (number of matching non-wildcard symbols). Specifically, for candidates $\{(i_1, j_1, \delta_1), \dots, (i_k, j_k, \delta_k)\}$ all achieving cost c^* , the probability of selecting candidate ℓ is

$$P(\ell) = \frac{\text{char_overlap}(i_\ell, j_\ell, \delta_\ell)}{\sum_{m=1}^k \text{char_overlap}(i_m, j_m, \delta_m)},$$

where $\text{char_overlap}(i, j, \delta)$ denotes the number of overlapping cells containing matching non-wildcard symbols when T_j is placed at offset δ relative to T_i . This weighting biases the selection toward attachments with more character matches (favoring compression) while still allowing exploration of alternative placements.

This stochastic heuristic often yields slightly worse single solutions than the deterministic variant, but it produces a diverse set of placement trees, which is useful for initializing the genetic algorithm population. The diversity arises because different runs may grow the tree in different orders, exploring distinct regions of the solution space.

The tree-growing procedure is summarized in Algorithm 2.

6.4. Tree-based genetic algorithm

We now describe the tree-based genetic algorithm (GA) used in our experiments. More precisely, our approach is a *Memetic Algorithm* (Moscato, 1989), a hybrid that couples global evolutionary search with a problem-specific local completion operator. This design reflects a modern consensus in combinatorial optimization: “pure” GAs rarely compete with state-of-the-art methods, but GAs hybridized with domain-specific heuristics consistently achieve top performance on structured problems (Merz and Freisleben, 2000). In our setting, the evolutionary process explores the *high-level topology*, which strings cluster together and how subtrees combine, while the greedy completion operator handles the *low-level boundary repair*, attaching the small fraction of strings that cannot be consistently inherited from parents. This division of labor is deliberate: the expensive evolutionary search focuses on the critical structural backbone of the placement tree, while trivial leaf attachments are delegated to fast, deterministic (or stochastic) greedy. Recall that a solution is represented as a *placement tree* $F = (V, E, r)$ whose vertices are strings, whose directed edges $(u \rightarrow v)$ are annotated by integer offsets $(\Delta x, \Delta y)$, and whose root r is the string placed at the origin. Decoding such a tree yields a concrete placement and an objective value.

As a preprocessing step we build a *placement graph* $G = (V, E_G)$ on the strings, whose directed edges encode all locally valid relative placements between string pairs. For each ordered pair of distinct strings (u, v) we enumerate

Algorithm 2 Tree-Growing Greedy for 2D-SSP**Require:** Set of 2D strings $\mathcal{T} = \{T_1, \dots, T_n\}$, root index r , objective cost $\in \{\text{area}, \text{sq}\}$, mode $\in \{\text{DET}, \text{STOCH}\}$ **Ensure:** Placement tree F with placement p_F

```

1:  $V_F \leftarrow \{r\}; E_F \leftarrow \emptyset; p_F(r) \leftarrow (0, 0)$ 
2:  $\text{canvas} \leftarrow \text{cells of } T_r \text{ at origin}$ 
3: while  $|V_F| < n$  do
4:    $\text{candidates} \leftarrow \emptyset$ 
5:   for each  $i \in V_F$ , each  $j \notin V_F$ , each valid offset  $\delta$  do
6:     if attaching  $T_j$  at  $p_F(i) + \delta$  is symbol-consistent with  $\text{canvas}$  then
7:        $c \leftarrow \text{cost}(\text{canvas} \cup \text{cells of } T_j \text{ at } p_F(i) + \delta)$ 
8:        $\text{candidates} \leftarrow \text{candidates} \cup \{(i, j, \delta, c)\}$ 
9:     end if
10:  end for
11:   $c^* \leftarrow \min\{c : (i, j, \delta, c) \in \text{candidates}\}$ 
12:   $\text{best} \leftarrow \{(i, j, \delta) : (i, j, \delta, c^*) \in \text{candidates}\}$ 
13:  if mode = DET then
14:     $(i^*, j^*, \delta^*) \leftarrow \text{element of } \text{best} \text{ maximizing overlap}$ 
15:  else
16:    // Roulette wheel selection weighted by overlap
17:     $w_\ell \leftarrow \text{overlap}(i_\ell, j_\ell, \delta_\ell)$  for each  $(i_\ell, j_\ell, \delta_\ell) \in \text{best}$ 
18:     $(i^*, j^*, \delta^*) \leftarrow \text{sample from } \text{best} \text{ with } P(\ell) \propto w_\ell$ 
19:  end if
20:   $V_F \leftarrow V_F \cup \{j^*\}; E_F \leftarrow E_F \cup \{(i^*, j^*, \delta^*)\}$ 
21:   $p_F(j^*) \leftarrow p_F(i^*) + \delta^*$ 
22:  Update  $\text{canvas}$  with cells of  $T_{j^*}$  at  $p_F(j^*)$ 
23: end while
24: return  $(F, p_F)$ 

```

translations of v relative to u within the bounded search window established by Lemma 3: $|\Delta x| \leq w_u + w_v - 1$ and $|\Delta y| \leq h_u + h_v - 1$. We collect all offsets within this window that yield symbol-consistent overlaps or 4-adjacent contact. This bounded enumeration is crucial for efficiency: for uniform $w \times h$ strings, we examine $O(wh)$ candidate offsets per pair rather than an unbounded search space.

To evaluate a tree $F = (V, E, r)$ we traverse it from the root r , assign absolute coordinates $p_F(i)$ to every 2D string T_i by summing the edge offsets along the unique path from r to i , and construct the global canvas, rejecting any edge that would create a symbol conflict. The fitness is then the bounding-box cost of the resulting placement, either

$$\text{cost}_{\text{area}}(p_F) = W(p_F) \cdot H(p_F) \quad \text{or} \quad \text{cost}_{\text{sq}}(p_F) = \max\{W(p_F), H(p_F)\},$$

depending on which objective variant is being optimized.

Each individual in the initial population is obtained by running a greedy constructive heuristic from a given start string, producing a full placement with absolute coordinates; we then extract a spanning tree of relative offsets. Thus every individual faithfully encodes the relative structure of a greedy solution.

We consider two variants that differ in their use of deterministic versus stochastic greedy:

- **T-GA** (Tree-based GA): Uses *deterministic* tree-growing greedy (DET mode) for both population initialization and greedy completion. Since deterministic greedy produces the same tree for a given root, initial population members differ *only* in their choice of starting root. This limits initial diversity to n distinct individuals.
- **ST-GA** (Stochastic Tree-based GA): Uses *stochastic* tree-growing greedy (STOCH mode) for both population initialization and greedy completion. Each initial individual is generated by an independent stochastic run, introducing diversity at the population level. Greedy completion during crossover is also stochastic, allowing different repair trajectories for incomplete offspring.

Algorithm 3 Tree-Based Genetic Algorithm for 2D-SSP

Require: Strings \mathcal{T} , population size N , generations G , crossover rate ρ , elite fraction ϵ , mode $\in \{\text{DET}, \text{STOCH}\}$
Ensure: Best placement tree F^*

```

1: // Initialization
2: for  $k = 1$  to  $N$  do
3:    $r_k \leftarrow$  random root from  $\{1, \dots, n\}$ 
4:    $P[k] \leftarrow \text{TREEGROWINGGREEDY}(\mathcal{T}, r_k, \text{mode})$  ▷ DET for T-GA, STOCH for ST-GA
5: end for
6: // Main loop
7: for  $g = 1$  to  $G$  do
8:   Evaluate fitness  $f[k] \leftarrow \text{cost}(P[k])$  for all  $k$ 
9:   Sort population by fitness (ascending)
10:   $P' \leftarrow$  copy of top  $\lfloor \epsilon N \rfloor$  individuals ▷ Elitism
11:  while  $|P'| < N$  do
12:    if  $\text{rand}() < \rho$  then
13:      Select parents  $F_1, F_2$  via tournament selection
14:       $F_{\text{child}} \leftarrow \text{TREECROSSOVER}(F_1, F_2, \text{mode})$ 
15:    else
16:       $F_{\text{child}} \leftarrow$  copy of tournament-selected individual
17:    end if
18:     $P' \leftarrow P' \cup \{F_{\text{child}}\}$ 
19:  end while
20:   $P \leftarrow P'$ 
21: end for
22: return  $\arg \min_{F \in P} \text{cost}(F)$ 

```

Our *locality-preserving crossover* operator combines two parent trees by *alternating* their local tree structures while maintaining geometric feasibility. This design directly implements the building block hypothesis: rather than mixing raw coordinates, we transplant entire subtrees, thereby preserving the *schema*, the pattern of relative offsets that makes a cluster of strings fit together well. Starting from the root, we expand a child tree by re-using parent edges whenever they can be realized without conflicts on the canvas. If some strings cannot be connected using parent edges alone, we perform a final *greedy completion* step to attach all remaining strings. The greedy completion ensures that every crossover produces a complete tree containing all strings, even when parental structures are incompatible.

The full GA is summarized in Algorithm 3, and the crossover operator in Algorithm 4. We maintain a population of trees, initialized from greedy placements with different starting strings (T-GA) or independent stochastic runs (ST-GA). In each generation we decode all individuals, rank them by fitness, copy the best few (elitism (De Jong, 1975)), and fill the remaining slots by crossover or by copying fit parents.

Our GA deliberately omits a dedicated mutation operator. This design choice is motivated by two considerations. First, when crossover is carefully designed to recombine meaningful building blocks, as our subtree-based crossover does, explicit mutation often provides diminishing returns or can even be counterproductive by disrupting well-structured solutions. This observation aligns with findings in other combinatorial optimization domains where problem-specific crossover operators dominate the search dynamics (Sholomon et al., 2013). Second, maintaining feasibility under traditional mutation is non-trivial in our setting: a random perturbation of edge offsets in a placement tree can easily introduce symbol conflicts or break connectivity.

However, we observe that the *greedy completion* step in our crossover operator (Algorithm 4, lines 18–22) implicitly serves as a *feasibility-preserving mutation mechanism*. When crossover produces an incomplete offspring, one that does not contain all strings, the greedy completion step attaches the missing strings using fresh, locally-optimal placements that differ from both parents. In the stochastic variant, this completion step uses *stochastic greedy* placement, which randomly samples among equally-good candidate positions rather than deterministically choosing one. This randomization introduces genuine variation: even when the same set of strings must be completed, different runs may attach them at different positions, exploring alternative regions of the solution space.

This design offers several advantages over traditional mutation:

Algorithm 4 Locality-Preserving Tree Crossover with Greedy Completion**Require:** Parent trees $F_1 = (V, E_1, r_1)$, $F_2 = (V, E_2, r_2)$, mode $\in \{\text{DET}, \text{STOCH}\}$ **Ensure:** Child tree F_c

```

1:  $r_c \leftarrow r_1$  ▷ Inherit root from first parent
2:  $V_c \leftarrow \{r_c\}$ ;  $E_c \leftarrow \emptyset$ ;  $p_c(r_c) \leftarrow (0, 0)$ 
3:  $canvas \leftarrow$  cells of  $T_{r_c}$  at origin
4:  $Q \leftarrow [r_c]$  ▷ Queue for BFS expansion
5: while  $Q \neq \emptyset$  do
6:    $i \leftarrow Q.\text{DEQUEUE}()$ 
7:   for each child  $j$  of  $i$  in  $F_1$  or  $F_2$  with  $j \notin V_c$  do
8:      $\delta \leftarrow$  offset of edge  $(i, j)$  in the parent tree
9:     if attaching  $T_j$  at  $p_c(i) + \delta$  is symbol-consistent with  $canvas$  then
10:        $V_c \leftarrow V_c \cup \{j\}$ ;  $E_c \leftarrow E_c \cup \{(i, j, \delta)\}$ 
11:        $p_c(j) \leftarrow p_c(i) + \delta$ 
12:       Update  $canvas$ ;  $Q.\text{ENQUEUE}(j)$ 
13:     end if
14:   end for
15: end while
16: // Greedy completion for missing strings
17: if  $|V_c| < n$  then
18:    $missing \leftarrow V \setminus V_c$ 
19:   for each  $j \in missing$  do
20:      $(i^*, \delta^*) \leftarrow$  best attachment via greedy (mode) ▷ DET or STOCH
21:      $V_c \leftarrow V_c \cup \{j\}$ ;  $E_c \leftarrow E_c \cup \{(i^*, j, \delta^*)\}$ 
22:      $p_c(j) \leftarrow p_c(i^*) + \delta^*$ ; update  $canvas$ 
23:   end for
24: end if
25: return  $F_c = (V_c, E_c, r_c)$ 

```

1. *Guaranteed feasibility.* Unlike random perturbations of tree edges, greedy completion always produces symbol-consistent placements by construction.
2. *Adaptive intensity.* The amount of “mutation” adapts to the compatibility of the parents: when parent structures are highly compatible, few strings require completion and offspring closely resemble parents; when parents are incompatible, many strings are completed afresh, introducing substantial variation.
3. *Local optimization.* Newly attached strings are placed greedily, ensuring that the mutated portion of the solution is locally reasonable rather than random.

Table 5 confirms that crossover dominates solution construction (95–97.5% of placements), with greedy completion providing lightweight repair for 3–5% of placements. See Appendix A.7 for detailed analysis of population dynamics.

Our experimental analysis (Section 7, Table 5) shows that approximately 3–5% of string placements arise from greedy completion, providing a consistent but moderate level of exploration that complements the exploitation performed by crossover.

In all experiments we fix the crossover rate $\rho = 0.7$ and an elite fraction of 10% of the population. Unless otherwise stated, the objective used in selection is the area-based cost of the decoded placement. The same GA can be run with the square cost by simply changing the fitness function to cost_{sq} , and we report results for both objective variants.

7. Experiments

This section presents our experimental evaluation. We first describe the experimental setup (Section 7.1), then report results comparing algorithms across different problem scales and string geometries (Section 7.2), and finally analyze the internal dynamics of the genetic algorithm (Section 7.3).

7.1. Experimental Setup

Hardware and implementation. All experiments were conducted on a system with an Intel Core i5-13400F processor and 32 GB RAM. The exact verification formulation was solved using IBM ILOG CPLEX 22.1 (MILP for the balanced objective and MIQP for the area objective). All heuristic algorithms were implemented in Python.

Instance generation. For each configuration, we generated 10 synthetic instances consisting of random binary 2D strings (alphabet $\Sigma = \{0, 1\}$), where each cell is assigned 0 or 1 uniformly at random with probability 0.5. The same 10 instances were used across all algorithms within each configuration to ensure fair comparison.

Table 2 presents results for configurations mixing 1D-like and 2D string shapes with the *area* objective. Table 3 reports results for genuinely 2D strings with the *square* objective.

Experimental configurations. We designed four experimental configurations to systematically evaluate algorithm performance across different scales and string geometries:

1. **1D-like instances** (1d): Strings with extreme aspect ratios (1×2 , 1×4 , 1×8) that closely resemble classical 1D strings. Instance sizes: $n \in \{10, 20, 50, 100, 200, 300\}$ strings. This configuration uses the *area* objective ($H \cdot W$) and tests whether merge-based greedy retains its 1D effectiveness. GA parameters: population size 150, 300 generations.
2. **Small instances** (small): Genuinely 2D strings (3×3 , 2×4 , 5×5 , 4×6) with $n \in \{6, 8, 10\}$ strings. This configuration includes CPLEX as a baseline (time limit: 300 s) to validate heuristic quality against optimal solutions. Uses the *square* objective ($\max\{H, W\}$). GA parameters: population size 100, 200 generations.
3. **Medium instances** (medium): Same string shapes as small, but with $n \in \{20, 30, 50\}$ strings. CPLEX is excluded due to computational intractability. Uses the *square* objective. GA parameters: population size 150, 300 generations.
4. **Large instances** (large): Same string shapes, with $n \in \{60, 80, 100\}$ strings for scalability testing. Uses the *square* objective. GA parameters: population size 200, 400 generations.

Algorithms compared. We evaluate five heuristic algorithms plus an exact solver:

- CPLEX: Exact ILP solver (small instances only, 300 s time limit).
- M-Greedy: Merge-based greedy (Algorithm 1).
- T-Greedy: Tree-growing greedy with deterministic tie-breaking.
- ST-Greedy: Tree-growing greedy with stochastic tie-breaking.
- T-GA: Genetic algorithm using deterministic T-Greedy for both population initialization and greedy completion. Initial population members differ only in the choice of starting root.
- ST-GA: Genetic algorithm using stochastic ST-Greedy for both population initialization and greedy completion. This introduces diversity at both stages: each initial individual is generated by an independent stochastic greedy run, and incomplete offspring are completed stochastically.

All GA variants use tournament selection with size 3 and crossover probability 0.7.

Performance metrics. For each algorithm and configuration, we report the objective value (mean \pm standard deviation over 10 instances) and mean runtime. Best results per configuration are shown in bold.

7.2. Results

Key observations (area objective, Table 2).

- Tree-based methods outperform merge-greedy by 9–15% on genuinely 2D instances.
- Merge-greedy matches tree-based methods on 1D-like (1×8) instances.
- GA variants improve over greedy by 6–12% across all 2D configurations.
- Stochastic tie-breaking (ST-GA) consistently achieves best objective values.

Config	M-Greedy	T-Greedy	ST-Greedy	T-GA	ST-GA
T10_n3_m3	71.80 ± 7.90 (0.019s)	69.20 ± 5.10 (0.005s)	69.10 ± 5.72 (0.005s)	62.30 ± 4.80 (0.195s)	59.00 ± 2.72 (0.266s)
T10_n4_m6	249.80 ± 21.80 (0.255s)	225.80 ± 6.95 (0.100s)	225.40 ± 4.20 (0.121s)	222.20 ± 5.40 (1.773s)	221.00 ± 5.60 (4.040s)
T10_n5_m5	277.00 ± 19.61 (0.202s)	240.00 ± 5.92 (0.089s)	241.00 ± 5.83 (0.089s)	237.00 ± 5.57 (1.530s)	236.00 ± 5.39 (4.976s)
T20_n1_m2	5.10 ± 0.54 (0.013s)	5.20 ± 0.60 (0.001s)	5.10 ± 0.54 (0.001s)	4.90 ± 0.30 (0.162s)	4.90 ± 0.30 (0.154s)
T20_n1_m4	18.10 ± 1.87 (0.046s)	18.30 ± 1.10 (0.003s)	19.10 ± 1.97 (0.003s)	16.80 ± 0.75 (0.331s)	16.50 ± 0.50 (0.509s)
T20_n1_m8	71.40 ± 6.41 (0.415s)	77.00 ± 7.97 (0.025s)	77.30 ± 5.92 (0.045s)	74.00 ± 5.67 (0.880s)	71.40 ± 6.23 (1.416s)
T20_n2_m4	94.30 ± 6.10 (0.147s)	98.90 ± 10.06 (0.013s)	100.60 ± 14.70 (0.014s)	84.60 ± 5.94 (1.196s)	77.50 ± 4.61 (1.427s)
T20_n3_m3	113.40 ± 8.89 (0.130s)	120.30 ± 8.26 (0.010s)	122.60 ± 10.76 (0.013s)	104.20 ± 5.06 (1.041s)	96.10 ± 3.81 (1.628s)
T20_n4_m6	483.00 ± 31.12 (4.002s)	444.20 ± 7.61 (0.737s)	444.20 ± 8.96 (0.740s)	437.60 ± 5.12 (18.391s)	432.40 ± 7.68 (34.618s)
T20_n5_m5	518.40 ± 21.12 (3.390s)	470.00 ± 11.40 (0.562s)	467.50 ± 9.01 (0.552s)	459.00 ± 8.31 (18.089s)	452.00 ± 6.00 (32.016s)
T30_n2_m4	118.80 ± 10.60 (0.526s)	134.80 ± 17.28 (0.031s)	132.00 ± 14.83 (0.037s)	109.90 ± 8.19 (2.329s)	109.30 ± 7.93 (2.883s)
T30_n3_m3	143.10 ± 12.51 (0.396s)	154.70 ± 9.95 (0.021s)	164.80 ± 11.70 (0.024s)	138.40 ± 7.14 (2.821s)	131.90 ± 7.02 (3.518s)
T30_n4_m6	710.30 ± 30.83 (17.240s)	656.80 ± 16.35 (1.786s)	656.00 ± 22.11 (1.839s)	638.80 ± 8.40 (45.850s)	633.20 ± 8.77 (83.731s)
T30_n5_m5	748.67 ± 24.68 (15.296s)	692.78 ± 8.53 (1.696s)	690.00 ± 11.30 (1.682s)	681.25 ± 8.93 (52.411s)	673.75 ± 11.66 (91.774s)
T50_n1_m2	5.50 ± 0.50 (0.148s)	5.20 ± 0.40 (0.003s)	5.40 ± 0.49 (0.003s)	5.00 ± 0.00 (0.653s)	5.00 ± 0.00 (0.656s)
T50_n1_m4	19.30 ± 1.10 (0.626s)	21.10 ± 1.92 (0.004s)	20.90 ± 1.70 (0.004s)	18.50 ± 0.67 (1.258s)	18.40 ± 0.49 (1.716s)
T50_n1_m8	123.60 ± 13.51 (5.470s)	135.10 ± 14.82 (0.148s)	136.50 ± 13.03 (0.191s)	131.30 ± 13.84 (7.273s)	127.10 ± 13.40 (6.853s)
T50_n2_m4	164.50 ± 11.86 (2.593s)	181.40 ± 10.00 (0.078s)	181.00 ± 8.64 (0.069s)	156.50 ± 7.34 (6.798s)	153.10 ± 6.32 (8.555s)
T50_n3_m3	211.10 ± 16.83 (1.806s)	227.30 ± 14.45 (0.086s)	236.60 ± 15.92 (0.097s)	197.90 ± 6.70 (13.097s)	191.33 ± 9.01 (18.018s)
T50_n5_m5	1166.44 ± 37.06 (110.392s)	1122.22 ± 17.18 (6.525s)	1127.22 ± 18.12 (6.606s)	1100.56 ± 12.35 (234.643s)	1091.00 ± 7.76 (302.926s)
T100_n1_m2	5.00 ± 0.00 (1.188s)	5.20 ± 0.40 (0.004s)	5.40 ± 0.49 (0.004s)	5.00 ± 0.00 (1.628s)	5.00 ± 0.00 (1.619s)
T100_n1_m4	20.60 ± 1.43 (4.948s)	20.00 ± 1.10 (0.007s)	21.30 ± 1.55 (0.006s)	19.30 ± 0.46 (3.078s)	19.00 ± 0.00 (4.048s)
T100_n1_m8	178.20 ± 11.41 (37.071s)	198.50 ± 11.72 (0.399s)	196.10 ± 10.15 (0.380s)	185.10 ± 8.94 (25.846s)	183.40 ± 9.32 (25.341s)
T200_n1_m2	5.33 ± 0.47 (9.418s)	5.00 ± 0.00 (0.008s)	5.33 ± 0.47 (0.009s)	5.00 ± 0.00 (4.230s)	5.00 ± 0.00 (4.374s)

Table 2

Objective type: area. Each cell shows objective (top) as mean ± std and runtime is mean only (bottom).

Key observations (square objective, Table 3).

- Even stronger relative performance: ST-GA reduces cost by > 50% vs. merge-greedy on some configurations.
- Merge-greedy's overlap-maximization is poorly suited for aspect-ratio control.
- All heuristics scale to $n = 50$ strings; GA runtimes grow to ~230 s on large instances.

Config	M-Greedy	T-Greedy	ST-Greedy	T-GA	ST-GA
T6_n3_m3	7.70 ± 0.90 (0.007s)	7.33 ± 0.54 (0.001s)	7.47 ± 0.56 (0.001s)	7.00 ± 0.37 (0.165s)	6.77 ± 0.42 (0.225s)
T6_n5_m5	15.30 ± 1.27 (0.030s)	14.90 ± 1.04 (0.001s)	14.80 ± 0.60 (0.001s)	13.50 ± 0.50 (0.476s)	13.10 ± 0.30 (1.376s)
T6_n10_m10	33.00 ± 1.73 (0.681s)	31.80 ± 1.25 (0.009s)	31.10 ± 1.22 (0.009s)	29.60 ± 0.66 (5.311s)	28.90 ± 0.30 (14.190s)
T8_n3_m3	8.60 ± 0.92 (0.013s)	8.23 ± 0.56 (0.001s)	8.40 ± 0.55 (0.001s)	7.77 ± 0.42 (0.304s)	7.37 ± 0.48 (0.285s)
T8_n5_m5	17.40 ± 1.28 (0.088s)	16.30 ± 0.64 (0.002s)	16.60 ± 0.66 (0.002s)	15.10 ± 0.30 (1.561s)	15.00 ± 0.00 (2.793s)
T8_n10_m10	42.40 ± 3.93 (2.267s)	34.30 ± 1.27 (0.016s)	34.10 ± 1.51 (0.016s)	31.70 ± 0.90 (24.422s)	30.60 ± 0.49 (40.423s)
T10_n3_m3	9.20 ± 0.40 (0.020s)	8.97 ± 0.41 (0.001s)	9.13 ± 0.43 (0.001s)	8.17 ± 0.37 (0.237s)	8.03 ± 0.18 (0.380s)
T10_n5_m5	19.50 ± 2.91 (0.202s)	17.90 ± 0.70 (0.002s)	18.20 ± 0.40 (0.003s)	16.70 ± 0.46 (1.786s)	16.50 ± 0.50 (6.692s)
T10_n10_m10	51.50 ± 5.00 (5.939s)	39.60 ± 1.20 (0.024s)	39.50 ± 0.50 (0.027s)	37.10 ± 0.54 (58.876s)	36.90 ± 0.30 (111.969s)
T20_n3_m3	11.80 ± 0.75 (0.127s)	11.33 ± 0.47 (0.002s)	11.77 ± 0.56 (0.002s)	10.27 ± 0.44 (1.280s)	10.13 ± 0.34 (2.314s)
T20_n5_m5	26.40 ± 3.32 (3.331s)	24.00 ± 0.45 (0.008s)	24.00 ± 0.77 (0.010s)	22.60 ± 0.49 (16.813s)	22.80 ± 0.40 (27.894s)
T20_n10_m10	105.50 ± 2.50 (129.009s)	53.60 ± 1.11 (0.113s)	53.40 ± 0.92 (0.135s)	50.30 ± 0.64 (739.766s)	50.00 ± 0.00 (1088.163s)
T30_n3_m3	13.80 ± 1.72 (0.391s)	12.83 ± 0.52 (0.003s)	13.57 ± 0.56 (0.004s)	11.93 ± 0.36 (2.897s)	11.90 ± 0.30 (5.085s)
T30_n5_m5	30.60 ± 2.76 (15.057s)	28.70 ± 1.00 (0.019s)	29.10 ± 0.70 (0.030s)	27.10 ± 0.30 (75.373s)	27.40 ± 0.49 (108.513s)
T50_n3_m3	16.30 ± 1.19 (1.793s)	15.00 ± 0.52 (0.006s)	15.77 ± 0.80 (0.008s)	14.13 ± 0.34 (12.208s)	14.20 ± 0.48 (18.474s)
T50_n5_m5	35.50 ± 0.50 (109.531s)	35.00 ± 0.00 (0.055s)	36.00 ± 0.00 (0.076s)	33.50 ± 0.50 (236.194s)	35.00 ± 0.00 (358.765s)

Table 3

Objective type: square. Each cell shows objective (top) as mean ± std and runtime is mean only (bottom).

Comparison with CPLEX (Table 4). On small instances where CPLEX finds provably optimal solutions, ST-GA achieves optimality gaps of 0.4–2.6%. On larger instances (T20–T30), ST-GA outperforms CPLEX (which hits the 300 s time limit), validating the GA’s effectiveness.

Runtime trade-offs. Greedy methods complete in milliseconds; GA variants require seconds to minutes. This defines a Pareto frontier: greedy for interactive or exploratory use, GA for offline optimization where solution quality justifies additional computation time.

Visual comparison. Figure 3 illustrates the qualitative difference on a $N = 50$ instance: the GA solution achieves a ~10% smaller bounding box through tight interlocking that greedy’s myopic decisions cannot discover.

7.3. Genetic Algorithm Dynamics

7.3.1. Crossover statistics

To better understand how the genetic algorithms construct their solutions, we instrument T-GA and ST-GA with additional counters. For each run we record:

- the total number of crossover operations, C (total_crossovers);
- the number of crossovers that produce an incomplete placement and therefore require greedy completion, C_{repair} (crossovers_needing_completion);
- the total number of strings that are finally placed by the greedy completion step across all incomplete offspring, T_{fix} (total_strings_completed).

Config	CPLEX	M-Greedy	T-Greedy	ST-Greedy	beam_search	T-GA	ST-GA
T6_n2_m4	6.40 ± 0.49 (0.553s)	8.30 ± 1.10 (0.005s)	6.90 ± 0.30 (0.001s)	7.20 ± 0.40 (0.001s)	6.40 ± 0.49 (0.015s)	6.50 ± 0.50 (0.131s)	6.40 ± 0.49 (0.126s)
T6_n3_m3	6.60 ± 0.49 (1.194s)	7.70 ± 0.90 (0.007s)	7.33 ± 0.54 (0.001s)	7.47 ± 0.56 (0.001s)	7.00 ± 0.45 (0.020s)	7.00 ± 0.37 (0.165s)	6.77 ± 0.42 (0.225s)
T6_n4_m6	12.00 ± 0.00 (21.692s)	15.80 ± 2.99 (0.035s)	14.40 ± 0.49 (0.002s)	14.10 ± 0.70 (0.001s)	12.90 ± 0.54 (0.069s)	13.10 ± 0.30 (0.400s)	12.40 ± 0.49 (0.571s)
T6_n5_m5	13.00 ± 0.00 (292.930s)	15.30 ± 1.27 (0.030s)	14.90 ± 1.04 (0.001s)	14.80 ± 0.60 (0.001s)	13.80 ± 0.40 (0.074s)	13.50 ± 0.50 (0.476s)	13.10 ± 0.30 (1.376s)
T8_n2_m4	7.00 ± 0.00 (3.270s)	10.00 ± 1.73 (0.011s)	7.90 ± 0.30 (0.001s)	8.00 ± 0.45 (0.001s)	7.20 ± 0.40 (0.038s)	7.20 ± 0.40 (0.194s)	7.00 ± 0.00 (0.225s)
T8_n3_m3	7.30 ± 0.46 (18.737s)	8.60 ± 0.92 (0.013s)	8.23 ± 0.56 (0.001s)	8.40 ± 0.55 (0.001s)	7.80 ± 0.40 (0.043s)	7.77 ± 0.42 (0.304s)	7.37 ± 0.48 (0.285s)
T8_n4_m6	14.40 ± 0.49 (302.317s)	19.50 ± 3.17 (0.099s)	16.00 ± 1.00 (0.002s)	16.20 ± 0.75 (0.002s)	15.20 ± 0.40 (0.176s)	14.90 ± 0.30 (1.099s)	14.70 ± 0.46 (1.962s)
T8_n5_m5	14.70 ± 0.46 (302.312s)	17.40 ± 1.28 (0.088s)	16.30 ± 0.64 (0.002s)	16.60 ± 0.66 (0.002s)	15.20 ± 0.40 (0.180s)	15.10 ± 0.30 (1.561s)	15.00 ± 0.00 (2.793s)
T10_n2_m4	7.40 ± 0.49 (38.092s)	11.00 ± 1.84 (0.023s)	8.70 ± 0.46 (0.001s)	8.80 ± 0.60 (0.001s)	7.80 ± 0.40 (0.076s)	7.70 ± 0.46 (0.248s)	7.50 ± 0.50 (0.340s)
T10_n3_m3	8.00 ± 0.00 (127.740s)	9.20 ± 0.40 (0.020s)	8.97 ± 0.41 (0.001s)	9.13 ± 0.43 (0.001s)	8.40 ± 0.49 (0.078s)	8.17 ± 0.37 (0.237s)	8.03 ± 0.18 (0.380s)
T10_n4_m6	16.67 ± 0.67 (303.327s)	24.50 ± 3.11 (0.249s)	17.40 ± 0.66 (0.002s)	18.00 ± 0.77 (0.002s)	16.60 ± 0.49 (0.376s)	16.50 ± 0.50 (1.637s)	16.10 ± 0.30 (3.626s)
T10_n5_m5	16.50 ± 0.50 (302.949s)	19.50 ± 2.91 (0.202s)	17.90 ± 0.70 (0.002s)	18.20 ± 0.40 (0.003s)	16.80 ± 0.40 (0.391s)	16.70 ± 0.46 (1.786s)	16.50 ± 0.50 (6.692s)

Table 4

Objective type: square. Each cell shows objective (top) as mean ± std and runtime is mean only (bottom).

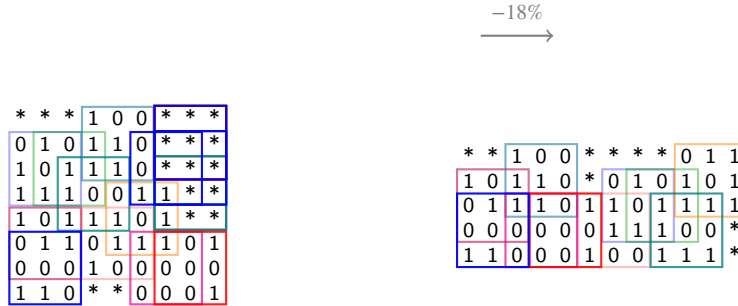


Figure 3: Visual comparison of placements on a representative instance (10 strings of size 3×3). *Left:* ST-GA discovers a tightly interlocking arrangement with overlapping strings, achieving a 12×5 bounding box. *Right:* M-Greedy places strings more linearly with minimal overlaps, yielding a larger 9×8 bounding box. Colored rectangles represent individual strings; darker regions indicate symbol-consistent overlaps. The GA's global search over the tree representation finds configurations that greedy's myopic decisions cannot reach.

From these quantities we derive:

$$r_{\text{repair}} = \frac{C_{\text{repair}}}{C}$$

repair rate: fraction of crossovers that need greedy fix

$$\rho_{\text{greedy}} = \frac{T_{\text{fix}}}{Cn}$$

share of string placements coming from greedy completion

$$\rho_{\text{direct}} = 1 - \rho_{\text{greedy}}$$

“completion rate”: share of strings placed directly by crossover

$$\bar{i}_{\text{repair}} = \frac{T_{\text{fix}}}{C_{\text{repair}}}$$

average #missing strings per repaired offspring

Algorithm	Scale	Mean #crossovers	Repair rate r_{repair}	Direct strings ρ_{direct}	Avg. strings/repair \bar{t}_{repair}
T-GA	small	3 152	9.2%	97.0%	3.0
T-GA	medium	12 594	6.7%	97.0%	22.6
T-GA	large	28 357	3.1%	97.5%	68.3
ST-GA	small	3 151	12.8%	95.9%	2.7
ST-GA	medium	12 602	8.9%	95.2%	22.0
ST-GA	large	28 354	6.3%	96.1%	60.8

Table 5

Internal statistics of the genetic algorithms. “Direct strings” is the fraction of string placements coming directly from crossover ($\rho_{\text{direct}} = 1 - \rho_{\text{greedy}}$); the remainder are filled by the greedy completion procedure. All experiments here use the area-based objective.

Several trends are apparent:

- The number of crossovers per run grows with instance size. On small instances each GA performs about 3.1×10^3 crossovers; this increases to roughly 1.26×10^4 on medium instances and 2.8×10^4 on large instances.
- The repair rate r_{repair} is modest: between 3% and 13% of crossovers produce offspring that are not complete placements. The ST-GA variant tends to trigger repairs slightly more often than T-GA.
- At the string level, only a very small fraction of the solution is delegated to the greedy completion phase. Across all scales and both GA variants, the mean ρ_{greedy} is about 3.6%, so approximately 96% of string placements come directly from crossover. In other words, the GA’s recombination operators are doing almost all of the constructive work.
- When a repair is needed, it can be substantial on large instances: on small instances, an incomplete offspring is missing on average 3 strings out of 8; on medium instances, around 22 out of 33 strings; and on large instances, around 61–68 out of roughly 80 strings. This reflects the increasing difficulty of producing fully consistent placements purely by recombination when the search space grows.

The role of greedy completion. Table 5 shows that 95–97.5% of string placements are inherited from parents via crossover; greedy completion acts as lightweight boundary repair for the remaining 3–5%. The evolutionary search over tree structures is the primary driver of solution quality.

Figure 4 focuses on the string-level interaction between crossover and greedy repair by plotting ρ_{greedy} (the fraction of strings placed by the greedy completion step) as a function of the number of strings. For both GA variants this fraction remains between roughly 2% and 7% across all sizes, confirming that the vast majority of strings in the final placements are produced by crossover rather than by the repair heuristic.

8. Conclusion

This paper has introduced the Two-Dimensional Shortest Superstring Problem (2D-SSP), establishing it as a rich combinatorial optimization problem unifying string sequencing with geometric packing.

Theoretical foundations. We established NP-hardness for both objectives and APX-hardness for area via L-reduction from 1D-SSP. The *Bounded-Offset Tree Representation* supports a structured, symmetry-reduced combinatorial view of 2D-SSP: bounded edge offsets and finite tree topologies yield a finite (though exponential) family of candidate encodings, and the objectives imply an explicit global bound on the size of an optimal bounding box (Lemma 4). The connectivity/compaction theorem (Theorem 5) proves optimal solutions can be made 4-connected without increasing cost.

Algorithmic innovation. Our Memetic Algorithm (T-GA) features locality-preserving crossover that preserves beneficial subtree structures. Crossover dominates solution construction ($> 95\%$ of placements), with greedy completion providing lightweight repair. The GA achieves optimality gaps $\leq 2.6\%$ on ILP-verifiable instances and outperforms greedy baselines by 6–12% on larger instances.

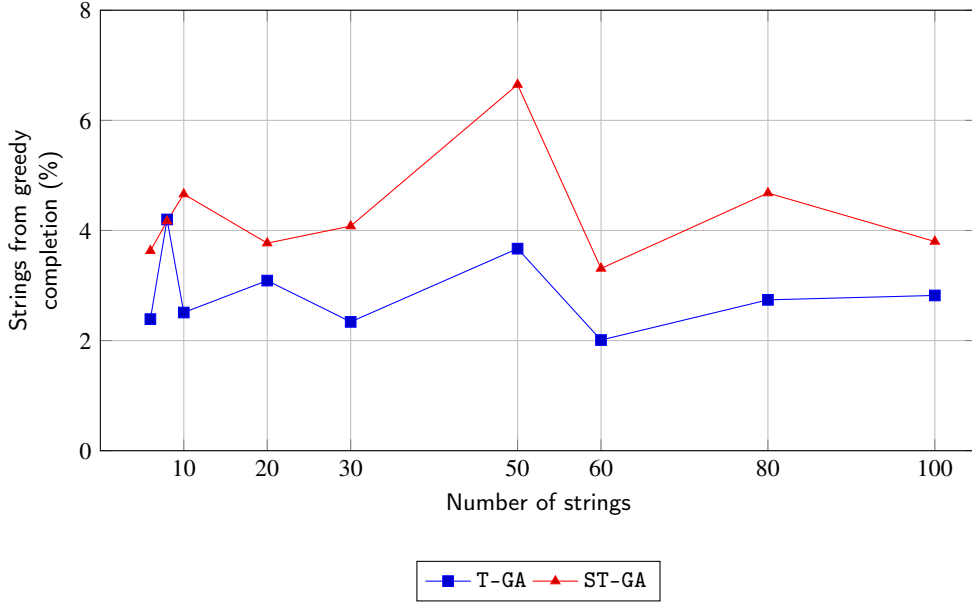


Figure 4: Share of string placements produced by the greedy completion step as a function of the number of strings. The complement to 100% can be interpreted as the “completion rate” of the crossover operators.

Extensions. The theoretical framework extends directly to d -dimensional SSP for any $d \geq 1$, with $2d$ -adjacency replacing 4-adjacency; see Appendix A.5.

Limitations. The ILP does not scale beyond $n \leq 10$. The gap between existence of optimal trees and the GA’s greedy-completion subspace remains theoretically open. Experiments focus on binary alphabets; behavior on higher-entropy instances is unexplored.

Future directions. Key open problems include: (1) approximation algorithms generalizing 1D-SSP’s 2.5-approximation; (2) APX-hardness of 2D-SSP_{sq} in the sequencing regime; (3) scaling exact methods via decomposition or column generation; and (4) extension to non-rectangular patterns.

A. Supplementary Technical Details

This appendix collects extended discussions moved from the main text for space reasons.

A.1. APX-Hardness Proof Details

Recall that an L-reduction from problem A to problem B requires constants $\alpha, \beta > 0$ such that:

1. $\text{OPT}_B \leq \alpha \cdot \text{OPT}_A$, and
2. for any solution to I_B with cost c_B , one can construct a solution to I_A with cost c_A satisfying $|c_A - \text{OPT}_A| \leq \beta \cdot |c_B - \text{OPT}_B|$.

Condition (1): From Theorem 1, $\text{OPT}_{2D} = \text{OPT}_{1D}$, so $\alpha = 1$.

Condition (2): Given any 2D placement with area A , the row-concatenation argument in Theorem 1 produces a 1D superstring of length $\leq A$. Thus $c_{1D} \leq c_{2D}$, which implies $c_{1D} - \text{OPT}_{1D} \leq c_{2D} - \text{OPT}_{2D}$, so $\beta = 1$.

A.2. Complexity Details: Alphabet Size and Sequencing–Packing Spectrum

Approximation complexity of 2D-SSP_{sq}. The L-reduction does *not* extend to 2D-SSP_{sq}. For height-1 strings, multiple rows can reduce the square cost: if $\text{OPT}_{1D} = 100$ and strings can be arranged in 10 rows of width 10 each, the square cost becomes $\max\{10, 10\} = 10 < 100$.

Our NP-hardness proof for 2D-SSP_{sq} (Theorem 2) operates in the *packing regime*: large alphabets force non-overlapping placements, reducing to rectangle packing. However, the *approximation complexity* in the *sequencing regime* (small $|\Sigma|$, frequent overlaps) remains open.

Alphabet size and hardness. The case $|\Sigma| = 1$ is trivial: all strings can be stacked at the origin, yielding bounding box equal to the largest string's dimensions. At the opposite extreme ($|\Sigma| \geq n$ with unique symbols), no overlaps are possible and 2D-SSP reduces to rectangle packing.

SSP vs. Packing: Opposing objectives. In packing problems, the goal is to fit rectangles into a container *without overlap*. In 2D-SSP, overlaps are *encouraged*, since they could potentially reduce the bounding box. For $|\Sigma| = 1$, 2D-SSP is trivial (maximal overlap permitted) while packing remains NP-hard (no-overlap constraint persists).

Information entropy and the sequencing–packing spectrum. The effective difficulty depends on string entropy:

- *High entropy (random strings, large $|\Sigma|$):* Probability of symbol-consistent overlap decreases exponentially with overlap size. 2D-SSP degenerates toward a packing problem.
- *Low entropy (repetitive strings, small $|\Sigma|$):* Many pairs admit large overlaps. The problem becomes a sequencing problem with combinatorial explosion of valid configurations.

Binary alphabets represent the sequencing regime where 2D-SSP is most distinct from pure geometric packing.

A.3. Edge Density and String Entropy

The bound $O(n^2 wh)$ on $|E(G^{\text{pl}})|$ is a worst-case geometric bound. The set of valid offsets C_{ij} decomposes as $C_{ij} = C_{ij}^{\text{adj}} \cup C_{ij}^{\text{ovl}}$.

Adjacency edges $|C_{ij}^{\text{adj}}|$: Always present, contributing $O(w_i + w_j + h_i + h_j)$ edges per pair (the perimeter of the contact region).

Overlap edges $|C_{ij}^{\text{ovl}}|$: For symbols drawn uniformly from Σ , $\Pr[k\text{-cell overlap is consistent}] = |\Sigma|^{-k}$. For random strings over large alphabets, $|C_{ij}^{\text{ovl}}| \approx 0$ and G^{pl} is sparse.

Periodic/low-entropy strings: Many overlaps become symbol-consistent. For $|\Sigma| = 1$, every offset in the contact region is valid, yielding $|C_{ij}| = O(w_i w_j + h_i h_j)$.

A.4. Optimized Offset Enumeration

The naïve approach enumerates all $O(wh)$ candidate offsets per pair. In the sequencing regime, we can improve performance by iterating offsets in order of *bounding-box increase*.

For two strings $T_i (w_i \times h_i)$ and $T_j (w_j \times h_j)$, placing T_j at offset $(\Delta x, \Delta y)$ increases the bounding box by a computable amount $\text{inc}(\Delta x, \Delta y)$. We enumerate offsets in non-decreasing order of inc :

- $\text{inc} = 0$: Full containment.
- $\text{inc} = 1$: Partial overlap leaving 1 row/column exposed.
- ... up to $\text{inc} = w_j + h_j - 2$ (4-adjacent contact, no overlap).

For fixed inc , valid offsets form a predictable geometric “frame.” A greedy algorithm can **stop at the first level containing a symbol-consistent offset**. Worst-case complexity remains $O(wh)$, but average-case improves when overlaps are frequent.

A.5. Extension to Higher Dimensions

The theoretical framework extends directly to d -dimensional SSP for any $d \geq 1$:

(i) *Geometry:* Objects are hyper-rectangles with dimensions $n_1 \times \dots \times n_d$. Cost becomes volume $\prod_{k=1}^d W_k$ or maximum side $\max_k W_k$. Offsets are vectors $\delta \in \mathbb{Z}^d$. The bounded offset property holds: valid offsets are bounded by the sum of dimensions along each axis.

(ii) *Connectivity lemma:* The sliding argument works in \mathbb{Z}^d : disconnected components A and B can be translated along a cardinal axis through empty space until they share a $(d-1)$ -dimensional face, without increasing any bounding-box dimension.

(iii) *Graph representation*: 4-connectivity generalizes to $2d$ -connectivity (sharing a $(d-1)$ -face). The contact graph definition is identical; Corollary 10 holds verbatim with $2d$ -adjacency replacing 4-adjacency.

Applications include 3D voxel assembly, volumetric data compression, and higher-dimensional tensor compression.

A.6. Existence vs. Search Reachability

Corollary 10 establishes *existence* of an optimal placement tree, but our GA searches only “greedily-completable” trees. A natural question: *can the optimal tree be “ungreedy”*? That is, might the optimal require a locally suboptimal edge to achieve minimum cost globally?

In principle, yes: the optimal tree might contain an edge (i, j, δ) dominated by (i, j, δ') with smaller local cost, yet δ enables a globally superior arrangement.

Two design choices mitigate this:

1. *Stochastic completion* randomly samples among equally-good candidates, exploring alternative attachment points.
2. *Crossover dominance*: Table 5 shows 95–97.5% of placements are inherited from parents via crossover, not constructed by greedy. The greedy-completion bias affects only a small fraction of solution structure.

Empirically, the GA matches ILP solutions on small instances, suggesting the greedy-completion subspace contains near-optimal solutions for typical inputs.

A.7. Population Diversity Analysis

Our greedy completion is *context-sensitive*: placement of each missing string depends on the current canvas state, which varies across offspring.

Table 5 confirms:

- Crossover dominates: 95–97.5% of placements come directly from crossover.
- Repair rate r_{repair} remains bounded at 3–13% even as instances scale.
- This indicates sustained population diversity rather than convergence to “super-individuals.”

The greedy operator acts as lightweight boundary repair, filling 2–5% of solution structure. This is analogous to mutation: essential for feasibility and diversity, but not the primary driver of solution quality.

B. MIP Formulation Details (MILP/MIQP)

This appendix provides the complete mathematical formulation of the mixed-integer program used for exact verification of small instances. The main text (Section 6.1) provides an overview; here we present the full technical details.

B.1. Notation and Grid Setup

We reuse the notation from Section 3. We have a finite set of 2D strings $\mathcal{T} = \{T_1, \dots, T_n\}$, each represented by a finite set of local cells $C_i \subset \mathbb{Z}^2$ and a symbol function $T_i : C_i \rightarrow \Sigma$. For each i we denote the local bounding box of T_i by

$$\begin{aligned} x_i^{\min} &= \min_{(u,v) \in C_i} u, & x_i^{\max} &= \max_{(u,v) \in C_i} u, \\ y_i^{\min} &= \min_{(u,v) \in C_i} v, & y_i^{\max} &= \max_{(u,v) \in C_i} v, \end{aligned}$$

and its width and height by

$$w_i = x_i^{\max} - x_i^{\min} + 1, \quad h_i = y_i^{\max} - y_i^{\min} + 1.$$

We embed all strings into a common rectangular grid $[0, W_g] \times [0, H_g] \subset \mathbb{Z}^2$. To guarantee that the MIP is a *true exact solver*, the grid must be large enough to accommodate *any* possible optimal arrangement, including those with extreme aspect ratios.

Remark 12 (Grid bounds and exactness). A naïve approach would set $W_g = \sum_i w_i$ and $H_g = \sum_i h_i$ (the maximum possible dimensions if all strings are placed in a line with no overlap). This guarantees global optimality but creates an enormous grid.

A tempting optimization is to use greedy bounds: run a heuristic to obtain dimensions $W_{\text{greedy}} \times H_{\text{greedy}}$, then set $W_g := W_{\text{greedy}}$ and $H_g := H_{\text{greedy}}$. However, this is *not* sound for all objectives:

- *For 2D-SSP_{area}*: An optimal solution may have a different aspect ratio than the greedy solution. For example, if greedy yields a 10×10 placement (area 100), the optimal might be 2×40 (area 80). If $W_g = H_g = 10$, the MIP would exclude the 2×40 solution since $40 > 10$.
- *For 2D-SSP_{sq}*: The greedy bound *is* sound. If greedy achieves $\max\{W_{\text{greedy}}, H_{\text{greedy}}\} = L$, any optimal solution has $\max\{W^*, H^*\} \leq L$, so both dimensions are bounded by L .

To ensure global exactness for both objectives, we use:

$$W_g := \min\left(\sum_i w_i, \text{cost}_{\text{area}}^{\text{greedy}}\right), \quad H_g := \min\left(\sum_i h_i, \text{cost}_{\text{area}}^{\text{greedy}}\right).$$

This exploits the fact that if the greedy area is A_{greedy} , any optimal area satisfies $W^* \cdot H^* \leq A_{\text{greedy}}$, so $W^* \leq A_{\text{greedy}}$ (since $H^* \geq 1$). The grid remains bounded by the greedy area rather than the sum of all dimensions, while accommodating all aspect ratios.

For small verification instances, this grid size is manageable. For larger instances, the MIP becomes intractable regardless of grid bounds.

B.2. Symmetry Breaking and Candidate Origins

To eliminate redundant symmetric solutions arising from translation invariance, we apply *symmetry breaking*: we fix the first string T_1 at the origin by restricting its set of allowed origins to $\mathcal{O}_1 := \{(0, 0)\}$. This constraint removes all translated copies of any given solution from the search space without affecting optimality.

For each string $i \geq 2$ we define a finite set of allowed *origins* $\mathcal{O}_i \subset \mathbb{Z}^2$ such that translating T_i by o keeps all its local cells inside the global grid:

$$\mathcal{O}_i := \left\{ o = (x, y) \in \mathbb{Z}^2 \mid \begin{array}{l} (x + u, y + v) \in [0, W_g] \times [0, H_g] \\ \text{for all } (u, v) \in C_i \end{array} \right\}. \quad (1)$$

Recall that $\mathcal{O}_1 = \{(0, 0)\}$ by the symmetry-breaking constraint. When T_i is placed at an origin $o = (x, y)$, each local cell $(u, v) \in C_i$ is mapped to global coordinates $(x + u, y + v)$.

B.3. Conflict Precomputation

Two candidate placements (i, o) and (j, o') are incompatible if they assign different symbols to the same global coordinate. Formally, we precompute a Boolean conflict indicator

$$\kappa_{ijoo'} = \begin{cases} 1, & \begin{array}{l} \text{if there exist } (u, v) \in C_i, (u', v') \in C_j \\ \text{with } (x + u, y + v) = (x' + u', y' + v') \\ \text{and } T_i(u, v) \neq T_j(u', v'), \\ \text{for } o = (x, y), o' = (x', y'), \end{array} \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

for all $i < j$, $o \in \mathcal{O}_i$, and $o' \in \mathcal{O}_j$. This preprocessing reduces symbol consistency to simple pairwise constraints in the MIP.

B.4. Decision Variables

The model uses the following variables.

- For each string $i \in \{1, \dots, n\}$ and each origin $o \in \mathcal{O}_i$:

$$b_{io} \in \{0, 1\} \quad (1 \text{ if } T_i \text{ is placed at origin } o, 0 \text{ otherwise}).$$

- Integer coordinates of the global bounding box:

$$X_{\min}, X_{\max}, Y_{\min}, Y_{\max} \in \mathbb{Z},$$

which represent the minimum and maximum global row/column indices among all occupied cells.

- The width and height of the bounding box:

$$W, H \in \mathbb{Z}_{\geq 0}.$$

- A maximum side variable (for the square objective):

$$L \in \mathbb{Z}_{\geq 0},$$

representing the maximum of width and height.

- An area variable (for the area objective):

$$A \in \mathbb{Z}_{\geq 0},$$

used to model or approximate the bounding-box area $W \cdot H$.

In the implementation we bound these variables by a constant

$$M := \max\{W_g, H_g\} + \max_i \max\{w_i, h_i\},$$

so that $X_{\min}, X_{\max}, Y_{\min}, Y_{\max} \in [-M, M]$ and $W, H, L \in [0, M]$, $A \in [0, M^2]$.

Remark 13 (Big- M calibration). The choice of M involves a trade-off. If M is too small, valid placements may be incorrectly excluded; if M is excessively large, the LP relaxation becomes weak (the big- M constraints provide little tightening when $b_{io} = 0$), leading to slow branch-and-bound convergence.

Our choice $M = \max\{W_g, H_g\} + \max_i \max\{w_i, h_i\}$ is valid given the grid bounds from Remark 12: since the grid accommodates all optimal solutions, no coordinate can exceed $\max\{W_g, H_g\}$ plus the maximum string dimension. The bound is instance-adaptive and typically much smaller than a naïve bound like $n \cdot \max_i \max\{w_i, h_i\}$, improving LP relaxation quality.

B.5. Common Constraints

Both objective variants share the following groups of constraints.

(i) *Exactly one origin per string.* Each 2D string must be placed at exactly one origin:

$$\sum_{o \in \mathcal{O}_i} b_{io} = 1 \quad \forall i \in \{1, \dots, n\}. \quad (3)$$

(ii) *No symbol conflicts.* If two candidate placements (i, o) and (j, o') conflict, they cannot be chosen simultaneously:

$$b_{io} + b_{jo'} \leq 1 \quad \forall i < j, \forall o \in \mathcal{O}_i, \forall o' \in \mathcal{O}_j \text{ with } \kappa_{ijoo'} = 1. \quad (4)$$

(iii) *Bounding box must contain all placed strings.* Let $o = (x, y) \in \mathcal{O}_i$ be a candidate origin for T_i . If $b_{io} = 1$, then the global footprint of T_i is

$$[x + x_i^{\min}, x + x_i^{\max}] \times [y + y_i^{\min}, y + y_i^{\max}],$$

and this rectangle must lie inside $[X_{\min}, X_{\max}] \times [Y_{\min}, Y_{\max}]$. We encode these implications with big- M constraints:

$$X_{\min} \leq x + x_i^{\min} + M(1 - b_{io}), \quad (5)$$

$$X_{\max} \geq x + x_i^{\max} - M(1 - b_{io}), \quad (6)$$

$$Y_{\min} \leq y + y_i^{\min} + M(1 - b_{io}), \quad (7)$$

$$Y_{\max} \geq y + y_i^{\max} - M(1 - b_{io}), \quad (8)$$

for all i and all $o = (x, y) \in \mathcal{O}_i$. When $b_{io} = 1$ these reduce to the desired inequalities $X_{\min} \leq x + x_i^{\min}$, $X_{\max} \geq x + x_i^{\max}$, etc., and when $b_{io} = 0$ they are relaxed by the big- M terms.

(iv) *Definition of width and height.* The bounding box dimensions are defined by

$$W = X_{\max} - X_{\min} + 1, \quad H = Y_{\max} - Y_{\min} + 1. \quad (9)$$

In addition, W and H must be at least as large as the widest and tallest individual 2D string:

$$W \geq \max_i w_i, \quad H \geq \max_i h_i. \quad (10)$$

B.6. Square Objective

The square objective minimises the maximum side length $\max\{W, H\}$. We link L to W and H via

$$L \geq W, \quad L \geq H. \quad (11)$$

At optimality, $L = \max\{W, H\}$.

The *square* MILP is

$$\begin{aligned} \min \quad & L \\ \text{s.t.} \quad & (3) - (10), (11), \\ & b_{io} \in \{0, 1\} \forall i, o, \quad X_{\min}, X_{\max}, Y_{\min}, Y_{\max}, W, H, L \in \mathbb{Z}. \end{aligned} \quad (12)$$

B.7. Area Objective

The area-based variant uses the product $W \cdot H$ as its cost. Since the constraints remain linear, this yields a Mixed-Integer Quadratic Program (MIQP), which we solve directly with CPLEX.

The *bounding area* MIQP is

$$\begin{aligned} \min \quad & W \cdot H \\ \text{exts.t.} \quad & (3) - (10), \\ & b_{io} \in \{0, 1\} \forall i, o, \quad X_{\min}, X_{\max}, Y_{\min}, Y_{\max}, W, H \in \mathbb{Z}. \end{aligned} \quad (13)$$

The objective value W^*H^* coincides with the area of the minimal axis-aligned bounding rectangle $B(p)$ enclosing all occupied cells in the induced placement p .

This formulation is conceptually straightforward: each b_{io} encodes a specific choice of origin for string T_i ; the conflict constraints (4) enforce symbol consistency; the big- M constraints (5)–(8) define a global bounding box that contains all chosen placements; and depending on the variant, either (11) or the quadratic objective (13) expresses the chosen cost. The greedy-based grid bounds and symmetry breaking significantly reduce the search space, but the number of variables and conflict constraints still grows quickly with the number of strings.

References

- Bennell, J.A., Oliveira, J.F., 2009. A tutorial in irregular shape packing problems. *Journal of the Operational Research Society* 60, S93–S105.
 Blum, A., Jiang, T., Li, M., Tromp, J., Yannakakis, M., 1994. Linear approximation of shortest superstrings. *Journal of the ACM* 41, 630–647.
 Cazaux, B., Rivals, E., 2018. Hierarchical overlap graph. *Information Processing Letters* 136, 78–84.
 Chang, Y.C., Chang, Y.W., Wu, G.M., Wu, S.W., 2000. B*-trees: a new representation for non-slicing floorplans, in: *Proceedings of the 37th Design Automation Conference*, pp. 458–463.

- Charalampopoulos, P., Pissis, S.P., Radoszewski, J., Waleń, T., Zuba, W., 2021. Computing covers of 2d strings, in: 32nd Annual Symposium on Combinatorial Pattern Matching (CPM 2021), Schloss Dagstuhl–Leibniz-Zentrum für Informatik. pp. 12:1–12:20.
- De Jong, K.A., 1975. An Analysis of the Behavior of a Class of Genetic Adaptive Systems. Ph.D. thesis. University of Michigan.
- Efros, A.A., Freeman, W.T., 2001. Image quilting for texture synthesis and transfer, in: Proceedings of the 28th annual conference on Computer graphics and interactive techniques, pp. 341–346.
- Gallant, J., Maier, D., Storer, J.A., 1980. On finding minimal length superstrings. *Journal of Computer and System Sciences* 20, 50–58.
- Gilmore, P.C., Gomory, R.E., 1965. Multistage cutting stock problems of two and more dimensions. *Operations Research* 13, 94–120.
- Golomb, S.W., 1994. Polyominoes: Puzzles, Patterns, Problems, and Packings. 2nd ed., Princeton University Press.
- Kaplan, H., Shafir, N., 2005. The greedy algorithm for shortest superstrings. *Information Processing Letters* 93, 13–17.
- Kwatra, V., Schödl, A., Essa, I., Turk, G., Bobick, A., 2003. Graphcut textures: image and video synthesis using graph cuts. *ACM Transactions on Graphics (ToG)* 22, 277–286.
- Leung, J.Y.T., Tam, T.W., Wong, C.S., Young, G.H., Chin, F.Y.L., 1990. Packing squares into a square. *Journal of Parallel and Distributed Computing* 10, 271–275.
- Lodi, A., Martello, S., Monaci, M., 2002. Two-dimensional packing problems: A survey. *European Journal of Operational Research* 141, 241–252.
- Merz, P., Freisleben, B., 2000. Fitness landscape analysis and memetic algorithms for the quadratic assignment problem. *IEEE Transactions on Evolutionary Computation* 4, 337–352.
- Moscato, P., 1989. On evolution, search, optimization, genetic algorithms and martial arts: Towards memetic algorithms. Technical Report C3P Report 826. Caltech Concurrent Computation Program.
- Murata, H., Fujiyoshi, K., Nakatake, S., Kajitani, Y., 1996a. Vlsi module placement based on rectangle-packing by the sequence-pair. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 15, 1518–1524.
- Murata, H., Fujiyoshi, K., Nakatake, S., Kajitani, Y., 1996b. Vlsi module placement based on rectangle-packing by the sequence-pair. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 15, 1518–1524.
- Sholomon, D., David, O.E., Netanyahu, N.S., 2013. A genetic algorithm-based solver for very large jigsaw puzzles. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1767–1774.
- Winfree, E., Liu, F., Wenzler, L.A., Seeman, N.C., 1998. Algorithmic self-assembly of dna. *Nature* 394, 539–544.
- Wolsey, L.A., 1998. *Integer Programming*. Wiley-Interscience.
- Yehezkeally, Y., Schwartz, M., 2025. Constructions of covering sequences and 2d-sequences. *Designs, Codes and Cryptography* 93, 1–25.