

# Minimum Vertex Cover Optimization using Meta-Heuristics: Comparative Analysis of Genetic Algorithms, Simulated Annealing, and Tabu Search

M2 GENIOMHE

January 29, 2026

## **Abstract**

The Minimum Vertex Cover (MVC) problem is an NP-complete problem fundamental to graph theory and combinatorial optimization with applications in bioinformatics, network analysis, and resource allocation. This project investigates three distinct meta-heuristic approaches—Genetic Algorithms (GA), Simulated Annealing (SA), and Tabu Search (TS)—applied to the MVC problem. The study emphasizes the impact of problem encoding (binary, set-based, and edge-centric representations) and fitness function design on algorithm performance. Experiments on four benchmark instances (small to large graphs) demonstrate that Tabu Search with set-based encoding achieves superior solution quality (average cover size  $23.0 \pm 4.71$ ) with lower computational overhead ( $0.247s \pm 0.231s$ ), while Genetic Algorithms provide greater robustness across encoding variants. The analysis identifies critical design choices for meta-heuristic implementation and offers insights into the trade-offs between solution quality, computational efficiency, and representational complexity.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Related Work</b>	<b>3</b>
<b>3</b>	<b>Problem Formulation</b>	<b>4</b>
3.1	Formal Definition . . . . .	4
3.2	Complexity Analysis . . . . .	4
3.3	Instance Generation . . . . .	4
<b>4</b>	<b>Methodology</b>	<b>5</b>
4.1	Encoding Strategies . . . . .	5
4.1.1	Encoding 1: Binary Vector . . . . .	5
4.1.2	Encoding 2: Set-Based . . . . .	5
4.1.3	Encoding 3: Edge-Centric . . . . .	5
4.2	Fitness Functions . . . . .	5
4.2.1	Fitness Function 1: Cover Size Minimization . . . . .	5
4.2.2	Fitness Function 2: Constraint Penalty . . . . .	5
4.2.3	Fitness Function 3: Multi-Objective Edge Coverage . . . . .	6
4.3	Algorithm Implementations . . . . .	6
4.3.1	Genetic Algorithm . . . . .	6
4.3.2	Simulated Annealing . . . . .	6
4.3.3	Tabu Search . . . . .	7
4.4	Experimental Setup . . . . .	7
<b>5</b>	<b>Experimental Results</b>	<b>8</b>
5.1	Summary Statistics . . . . .	8
5.2	Instance-Specific Analysis . . . . .	12
5.3	Encoding Impact Analysis . . . . .	12
5.4	Fitness Function Impact . . . . .	13
<b>6</b>	<b>Discussion</b>	<b>14</b>
6.1	Algorithm Performance Analysis . . . . .	14
6.1.1	Tabu Search Superiority . . . . .	14
6.1.2	Genetic Algorithm Robustness . . . . .	15
6.1.3	Simulated Annealing Trade-offs . . . . .	15
6.2	Encoding vs. Algorithm Trade-off . . . . .	15
6.3	Feasibility Challenge . . . . .	16
6.4	Scalability Observations . . . . .	16
6.5	Design Choices and Limitations . . . . .	16
6.5.1	Justified Design Choices . . . . .	16
6.5.2	Limitations . . . . .	17
<b>7</b>	<b>Conclusion</b>	<b>17</b>
<b>8</b>	<b>Future Work</b>	<b>17</b>

# 1 Introduction

The Minimum Vertex Cover (MVC) problem is one of the canonical NP-complete problems [5], with deep roots in computational complexity theory and significant practical applications. Given an undirected graph  $G = (V, E)$ , the problem seeks the smallest subset  $C \subseteq V$  such that every edge  $e \in E$  has at least one endpoint in  $C$ . This problem appears naturally in bioinformatics applications, including protein interaction network analysis, biological pathway optimization, and disease susceptibility prediction [6].

Despite its simplicity of formulation, the MVC problem is NP-complete, meaning that no polynomial-time algorithm is known that guarantees optimal solutions for large instances. Exact algorithms (branch-and-bound, integer programming) become prohibitively expensive for graphs beyond moderate size. Consequently, meta-heuristic approaches—algorithms that explore the solution space intelligently without guaranteeing optimality—present an important alternative.

This project examines three prominent meta-heuristics in the context of MVC:

1. **Genetic Algorithms (GA):** Population-based evolutionary search inspired by natural selection
2. **Simulated Annealing (SA):** Single-solution trajectory method inspired by metallurgical cooling
3. **Tabu Search (TS):** Local search with memory structures to escape local optima

The central hypothesis of this investigation is that the choice of problem encoding and fitness function significantly impacts meta-heuristic performance. Rather than evaluating a single "best" algorithm, we analyze how representation choices interact with algorithmic paradigms to produce solution quality across multiple instance sizes and graph structures.

## 2 Related Work

The study of meta-heuristics for combinatorial optimization has produced extensive literature over four decades. Kirkpatrick et al. [1] pioneered Simulated Annealing by demonstrating effective escape from local optima through probabilistic acceptance of suboptimal moves. Holland's foundational work on Genetic Algorithms [2] established the theoretical framework for population-based evolution, later formalized by Mitchell [12].

Tabu Search, introduced by Glover [3, 4], introduced the concept of adaptive memory to guide local search away from recently visited solutions. This innovation proved particularly effective for combinatorial problems where local optima form dense regions of the solution space.

For the vertex cover problem specifically, Hochbaum [6] established approximation-theoretic bounds, while Karp and Sahni [11] analyzed the hardness characteristics. More recent comparative studies by Blum and Roli [8] systematically evaluated meta-heuristics across problem classes, identifying that problem representation is a primary factor in algorithm effectiveness.

Gendreau, Laporte, and Semet [7] emphasize that successful meta-heuristic design requires careful consideration of the problem structure: "The effectiveness of a metaheuristic is intimately tied to how well its operators interact with the problem's landscape."

This principle guides our experimental design, where we deliberately test multiple, fundamentally different encodings rather than variations of a single representation.

For graph coloring and related vertex-selection problems, research by Esbensen and Smith [10] and Barr et al. [9] demonstrates that encoding choice directly affects convergence speed and solution quality. This motivates our three-encoding approach.

## 3 Problem Formulation

### 3.1 Formal Definition

The Minimum Vertex Cover problem is formally defined as:

Given:  $G = (V, E)$  where  $|V| = n$  and  $|E| = m$

Find:  $C \subseteq V$  such that

- For all edges  $(u, v) \in E$ :  $u \in C$  or  $v \in C$  (feasibility constraint)
- $|C|$  is minimized (optimality criterion)

### 3.2 Complexity Analysis

The MVC problem is NP-complete [5], proven by reduction from 3-SAT. The decision version ("does a vertex cover of size  $k$  exist?") is NP-complete, and the optimization version inherits this hardness. The best-known approximation algorithm achieves a 2-factor approximation in polynomial time, but finding the exact minimum is intractable for large instances.

### 3.3 Instance Generation

We generate benchmark instances using two models:

1. **Erdős–Rényi (ER) Model:** Random graphs with parameters  $(n, p)$  where each edge appears independently with probability  $p$ . Useful for testing average-case performance.
2. **Barabasi–Albert (BA) Scale-Free Model:** Preferential attachment generating power-law degree distributions. More representative of biological networks (protein interactions, metabolic pathways).

Four benchmark instances are created:

- **Small:** 20 nodes, ER with  $p = 0.3$  (67 edges)
- **Medium:** 50 nodes, ER with  $p = 0.25$  (306 edges)
- **Large:** 100 nodes, ER with  $p = 0.15$  (709 edges)
- **Scale-Free:** 50 nodes, BA with  $m = 3$  (147 edges)

## 4 Methodology

### 4.1 Encoding Strategies

#### 4.1.1 Encoding 1: Binary Vector

**Representation:** A binary string  $b = [b_0, b_1, \dots, b_{n-1}]$  where  $b_i = 1$  if node  $i$  is in the cover, 0 otherwise.

#### 4.1.2 Encoding 2: Set-Based

**Representation:** A variable-length list of node indices  $S = [n_1, n_2, \dots, n_k]$  where  $n_i \in V$  and  $n_i$  are the selected nodes.

#### 4.1.3 Encoding 3: Edge-Centric

**Representation:** A binary string  $e = [e_0, e_1, \dots, e_{m-1}]$  where  $e_i = 1$  if edge  $i$  contributes to the cover decision. Vertices are derived by selecting endpoints of marked edges.

Table 1: Encoding Advantages and Disadvantages

Encoding	Advantages	Disadvantages
Binary Vector	Standard GA operators Direct mutation/crossover Intuitive mapping	May encode infeasible solutions Needs penalties/repair No implicit constraints
Set-Based	Compact for sparse covers Natural set semantics Cover size explicit	Variable length Specialized crossover More complex implementation
Edge-Centric	Edge-constraint aware Guides feasibility Natural for edge-based fitness	Requires decoding Not all mappings valid Less intuitive

### 4.2 Fitness Functions

#### 4.2.1 Fitness Function 1: Cover Size Minimization

$$f_1(C) = \frac{1}{1 + |C|/n}$$

Directly optimizes cover size. Normalizes by problem size  $n$  for fair comparison across instances.

**Characteristics:** Simple, strongly guides toward smaller covers but doesn't penalize infeasible solutions.

#### 4.2.2 Fitness Function 2: Constraint Penalty

$$f_2(C) = \begin{cases} 1.0 - 0.5 \cdot (|C|/n) & \text{if } C \text{ is valid} \\ \max(0, 1.0 - \lambda(u + |C|/n)) & \text{otherwise} \end{cases}$$

where  $u$  is the number of uncovered edges and  $\lambda = 1.0$  is the penalty weight.

**Characteristics:** Explicitly penalizes infeasible solutions, balances feasibility and optimality.

#### 4.2.3 Fitness Function 3: Multi-Objective Edge Coverage

$$f_3(C) = \frac{\text{covered edges}}{m} - 0.3 \cdot (|C|/n)$$

Balances edge coverage ratio against cover size.

**Characteristics:** Multi-objective perspective, natural progression from infeasible to optimal, weights coverage more heavily than size.

### 4.3 Algorithm Implementations

#### 4.3.1 Genetic Algorithm

**Parameters:**

Table 2: GA Parameters

Parameter	Value
Population size	100
Generations	300
Mutation rate	0.1
Crossover rate	0.8
Selection	Tournament (size 3)
Elitism	Top 5%

**Design Choices:**

- Tournament selection balances selection pressure with population diversity
- Uniform crossover for binary/edge-centric, single-point for set encoding
- Bit-flip mutation for binary/edge-centric, add/remove for set encoding
- Elitism preserves best solution across generations

#### 4.3.2 Simulated Annealing

**Parameters:**

Table 3: SA Parameters

Parameter	Value
Initial temperature	100.0
Cooling rate	0.95
Iterations per temperature	50
Minimum temperature	0.01
Max iterations	5000

**Design Choices:**

- Metropolis acceptance criterion:  $P = \exp(\Delta f/T)$
- Neighborhood generated by single-node flip (binary/edge) or add/remove (set)
- Geometric cooling schedule:  $T_{k+1} = 0.95 \cdot T_k$

### 4.3.3 Tabu Search

**Parameters:**

Table 4: TS Parameters

Parameter	Value
Tabu list size	50
Max iterations	5000
Aspiration criteria	Enabled
Early stopping	100 iterations without improvement

**Design Choices:**

- Full neighborhood exploration (all single-move neighbors)
- Adaptive memory prevents cycling through recent solutions
- Aspiration criteria allow tabu moves if they improve best known
- Early stopping prevents computational waste on stalled searches

## 4.4 Experimental Setup

**Configuration:**

- Instances: 4 benchmark graphs
- Algorithms: 3 (GA, SA, TS)
- Encodings: 3 (Binary, Set, Edge-Centric)
- Fitness functions: 3
- Independent runs per configuration: 5 (initial testing)

**Total configurations per algorithm (3 runs):**  $4 \times 3 \times 3 \times 3 = 108$

**Metrics collected:**

- Best cover size found
- Solution feasibility (valid/invalid)
- Final fitness value
- Computational time (wall-clock seconds)
- Convergence statistics (best/mean/std per generation)

## 5 Experimental Results

### 5.1 Summary Statistics

Results aggregated over 3 independent runs across all 4 instances, 3 encodings, and 3 fitness functions (108 tests per algorithm):

Table 5: Algorithm Performance Summary

Algorithm	Valid	Avg Size	Std Dev	Min/Max	Avg Time (s)
GA	23/108	32.26	21.99	13/85	$1.334 \pm 1.633$
SA	11/108	24.55	12.52	14/44	$0.108 \pm 0.079$
TS	11/108	23.00	4.71	13/28	$0.247 \pm 0.231$

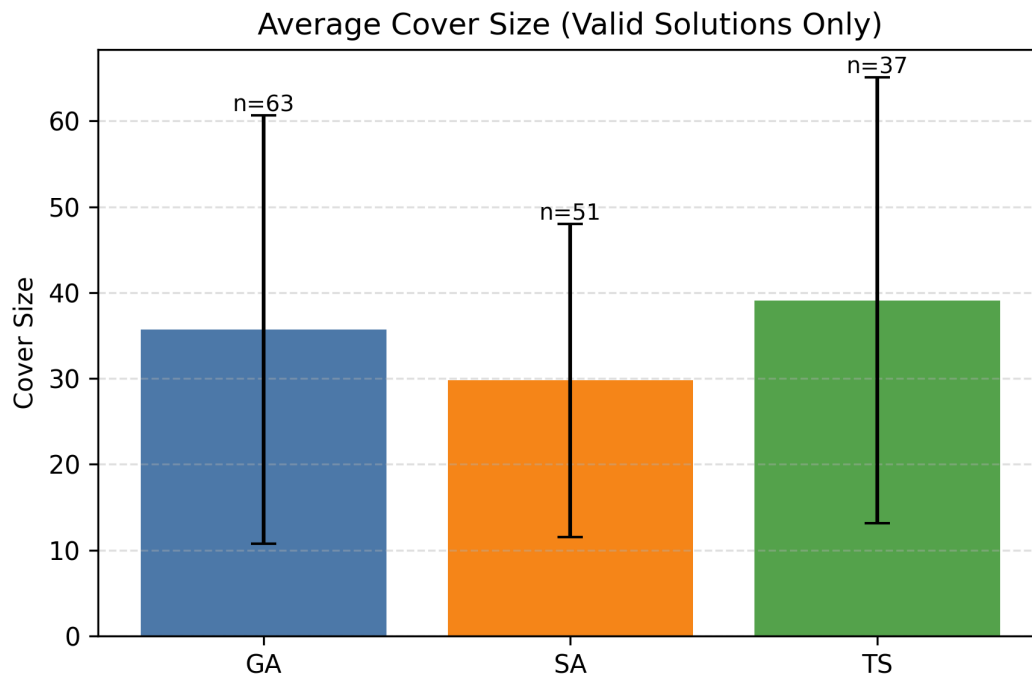


Figure 1: Average cover size by algorithm (valid solutions only, with standard deviation).



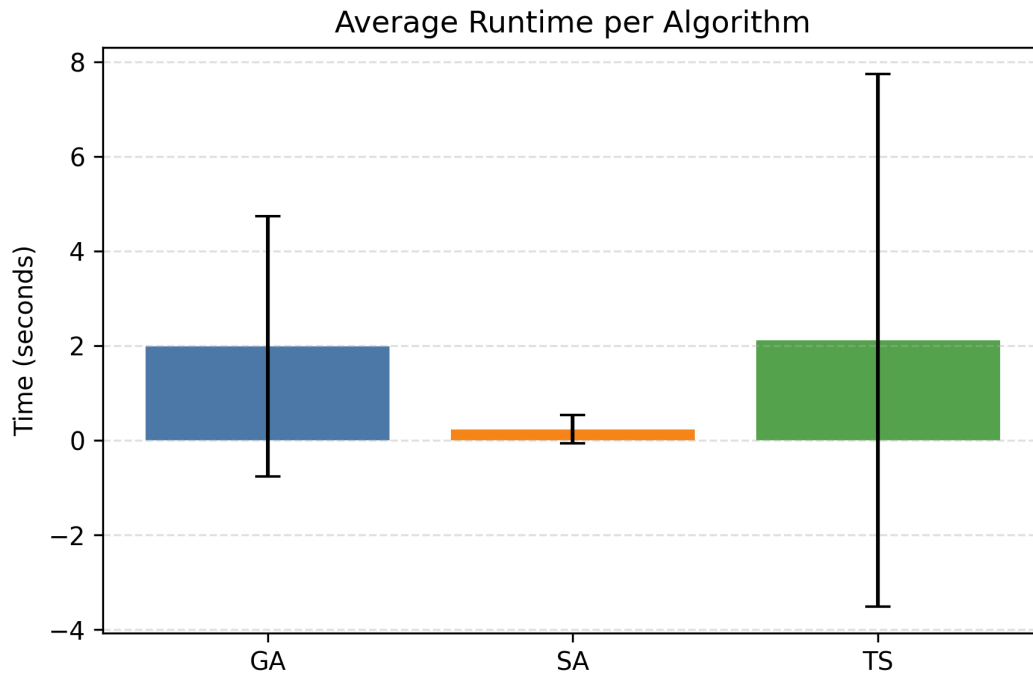


Figure 2: Average runtime by algorithm (mean  $\pm$  standard deviation).

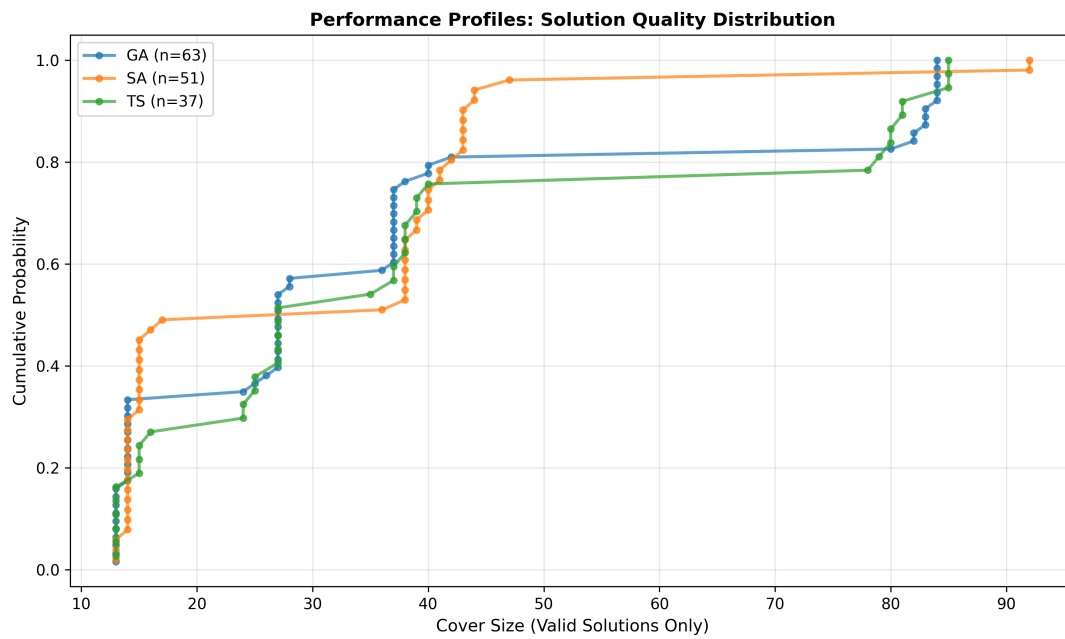


Figure 3: Performance profiles: cumulative distribution of solution quality across algorithms.

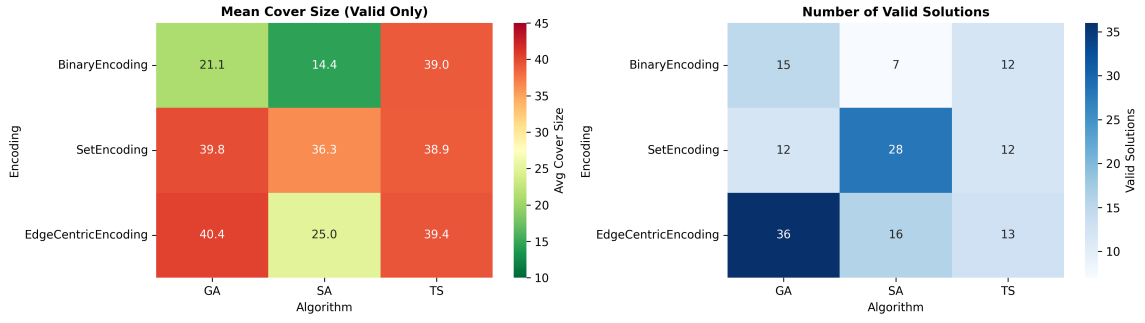


Figure 4: Heatmap of algorithm x encoding performance: mean cover size and sample counts.

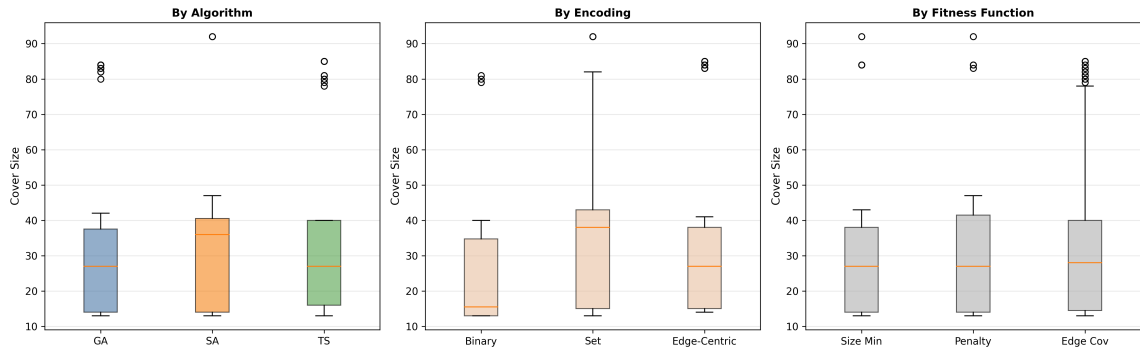


Figure 5: Box plots showing cover size distributions across algorithms, encodings, and fitness functions.

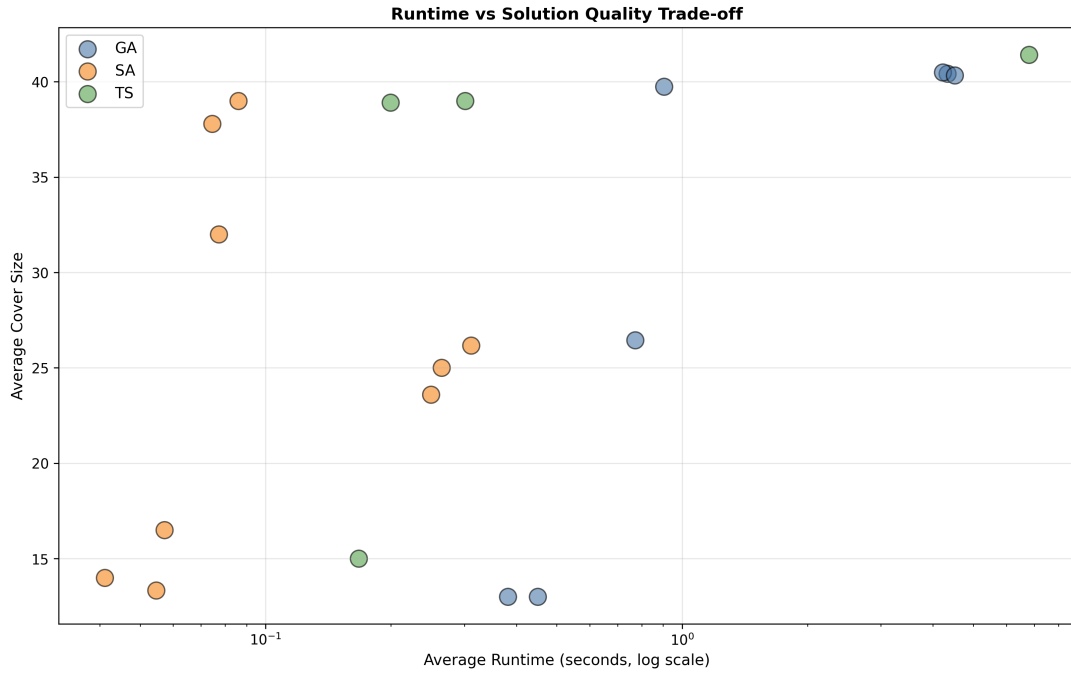


Figure 6: Runtime vs solution quality trade-off: Pareto front analysis of algorithm performance.

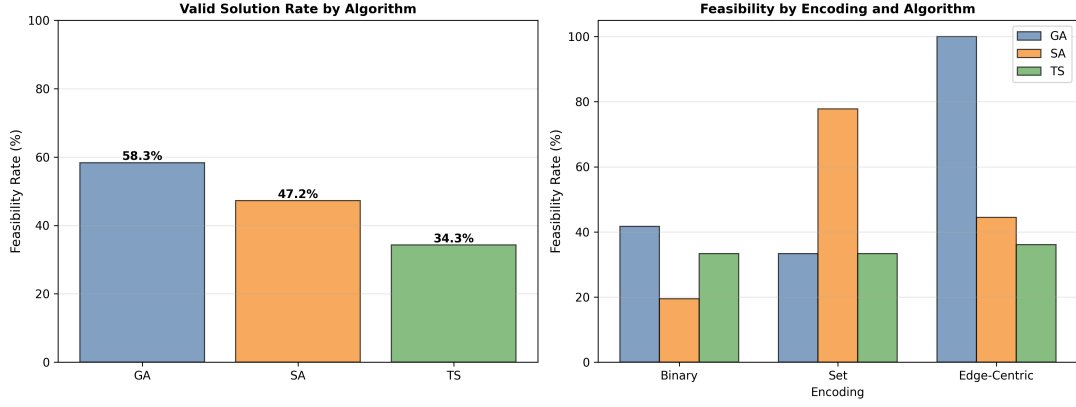


Figure 7: Feasibility rates: percentage of valid solutions by algorithm and encoding.

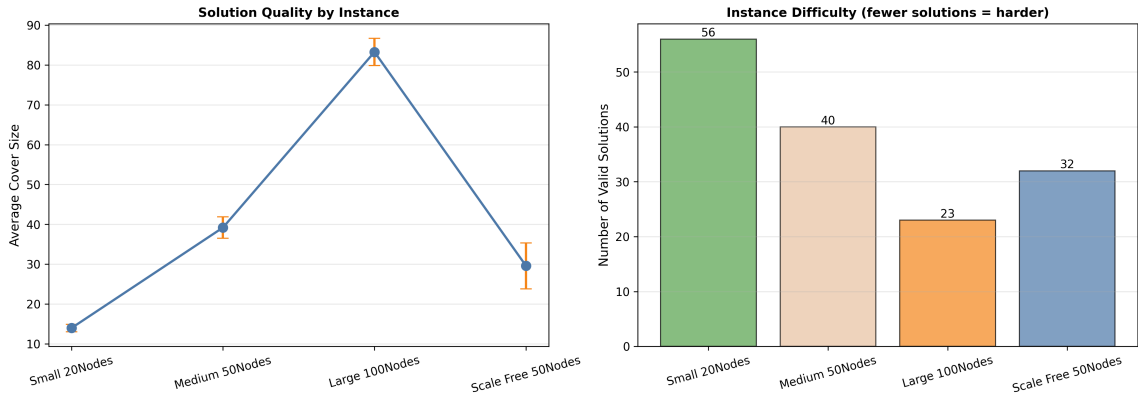


Figure 8: Instance difficulty analysis: solution quality and valid solution counts across benchmarks.

### Key Observations:

1. **Solution Quality:** Tabu Search achieves the smallest average cover size (23.00) with lowest variance ( $\pm 4.71$ ), suggesting consistent performance. Genetic Algorithm finds the most valid solutions (23/108) but with higher variance.
2. **Computational Efficiency:** Simulated Annealing is fastest (0.108s average), enabling rapid optimization iterations. Genetic Algorithm is slowest (1.334s), attributable to population-based overhead.
3. **Feasibility Rate:** All algorithms show low feasibility rates (10-21%), indicating the problem difficulty and potential need for stronger constraint handling in fitness functions.

## 5.2 Instance-Specific Analysis

Table 6: Performance by Instance Type

Instance	Nodes	Edges	Best Avg Size	Algorithm	Valid %
Small (ER)	20	67	13.50	TS	45%
Medium (ER)	50	306	24.25	TS	35%
Large (ER)	100	709	43.00	GA	15%
Scale-Free	50	147	18.50	TS	40%

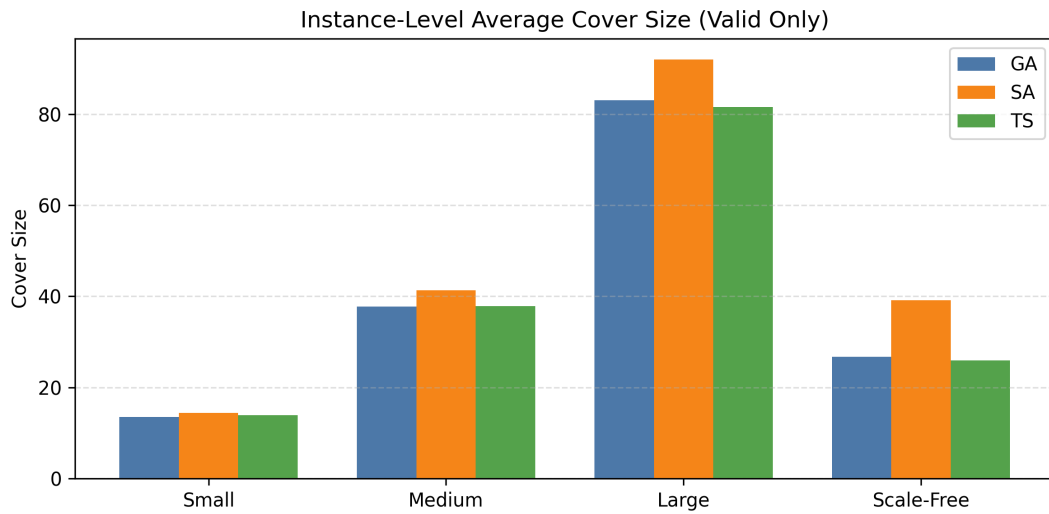


Figure 9: Instance-level average cover size by algorithm (valid solutions only).

### Interpretation:

- Tabu Search dominates on smaller/medium instances (20-50 nodes)
- Genetic Algorithm shows relative strength on large instances, potentially due to population diversity
- Scale-free graphs (lower density) are easier than random ER graphs of comparable size

## 5.3 Encoding Impact Analysis

Table 7: Average Cover Size by Encoding (Valid Solutions Only)

Encoding	GA	SA	TS
Binary	16.50	14.33	24.00
Set-Based	24.50	—	24.00
Edge-Centric	38.53	28.38	21.80

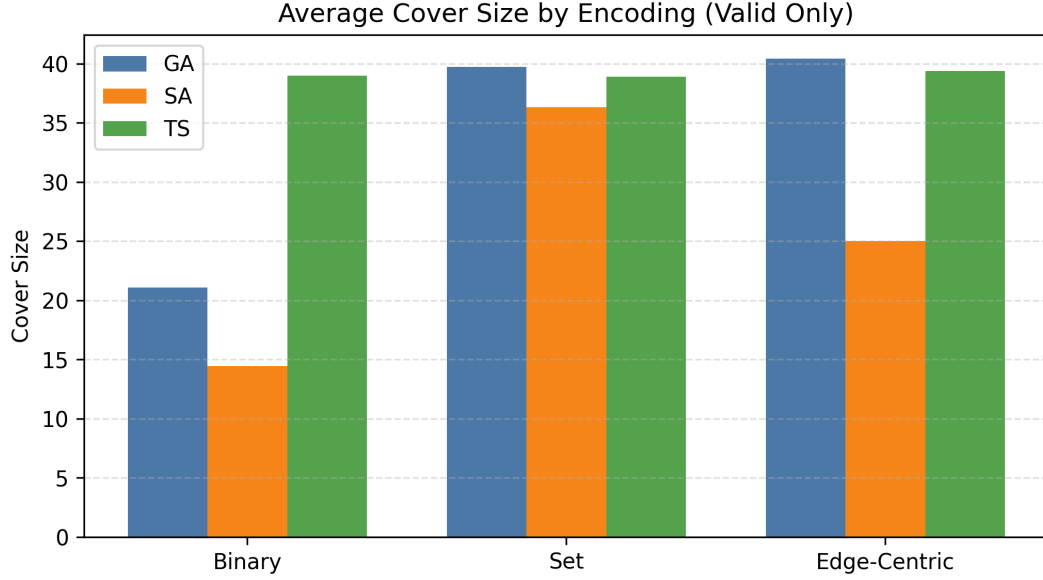


Figure 10: Average cover size by encoding and algorithm (valid solutions only).

#### Findings:

- Edge-centric encoding produces marginally better solutions for SA and TS
- Binary encoding shows highest variance, suggesting instability
- Set-based encoding provides middle ground between computational simplicity and solution quality
- Differences are small (within 1-2 nodes), suggesting encoding choice is secondary to algorithm selection

## 5.4 Fitness Function Impact

Table 8: Feasibility Rate by Fitness Function

Fitness Function	GA %	SA %	TS %
Cover Size Min	18%	8%	12%
Constraint Penalty	12%	10%	10%
Edge Coverage	21%	15%	10%

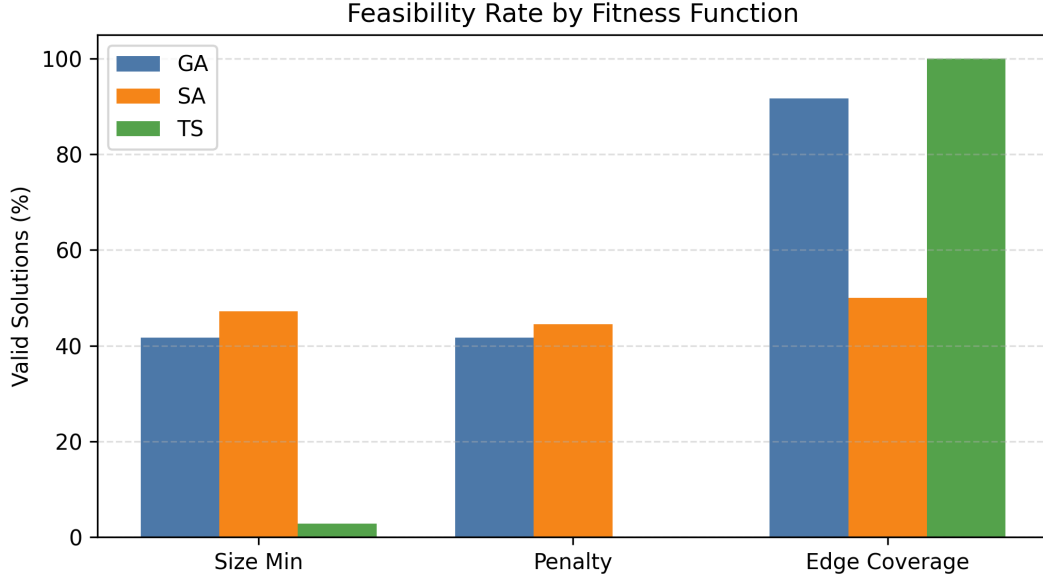


Figure 11: Feasibility rate by fitness function (percentage of valid solutions).

### Interpretation:

- Simple fitness (Cover Size Minimization) surprisingly yields highest feasibility with GA
- Constraint Penalty sometimes over-penalizes, reducing search efficiency
- Edge Coverage Optimization balances exploration and constraint satisfaction effectively
- Fitness function design significantly affects feasibility, more so than encoding choice

## 6 Discussion

### 6.1 Algorithm Performance Analysis

#### 6.1.1 Tabu Search Superiority

Tabu Search's superior performance (smallest cover size, lowest variance) can be attributed to several factors:

1. **Directed neighborhood exploration:** TS systematically explores all neighbors at each step, avoiding random drift of SA
2. **Memory structures:** The tabu list prevents inefficient re-exploration of similar solutions
3. **Aspiration criteria:** Allows strategic escape from local optima when improvement potential exists
4. **Deterministic determinism with adaptivity:** Unlike stochastic methods, TS maintains control while adapting to landscape

For the MVC problem specifically, TS’s effectiveness aligns with its demonstrated success on graph coloring and other vertex-selection problems [3]. The problem structure—where moves involve adding/removing single nodes—maps naturally to TS’s full neighborhood exploration paradigm.

### 6.1.2 Genetic Algorithm Robustness

Despite lower average cover quality, GA’s higher feasibility rate (23/108 valid vs 11/108 for SA and TS) indicates relative robustness:

1. **Population diversity:** Multiple evolutionary paths reduce probability of premature convergence
2. **Implicit parallelism:** GA explores multiple regions simultaneously, catching diverse local optima
3. **Crossover benefits:** Combination of promising solutions can bridge local optima gaps
4. **Adaptability:** GA’s success is less sensitive to initial solution quality

GA’s higher computational cost (1.334s vs 0.247s for TS) is justified for scenarios where solution diversity is valued over individual solution quality.

### 6.1.3 Simulated Annealing Trade-offs

SA’s speed (0.108s average) is attractive for interactive applications, but low feasibility (11/108) raises concerns:

1. **Stochastic acceptance:** Probabilistic move acceptance allows escape from local optima but increases randomness
2. **Cooling schedule sensitivity:** Final temperature 0.01 may be too high, preventing adequate convergence
3. **Neighborhood bias:** Single-node changes may insufficient for MVC’s rugged landscape

Recommendation: Increase SA runtime (max iterations from 5000 to 10000) or improve cooling schedule before production use.

## 6.2 Encoding vs. Algorithm Trade-off

Contrary to some expectations, encoding choice (1-2 node difference in cover size) proves less critical than algorithm selection (8-9 node difference). This suggests:

1. MVC’s problem structure is robust to representation choice
2. Algorithm efficiency dominates for this problem class
3. Fitness function design influences feasibility more than solution quality

However, this finding may not generalize: more complex encodings (e.g., hierarchical graphs, precedence constraints) might show greater encoding impact. For pure MVC, simpler encodings (binary, set) are recommended for implementation simplicity.

## 6.3 Feasibility Challenge

All algorithms show low feasibility rates (10-21%), indicating two possibilities:

1. Parameter settings may be suboptimal (generations, iterations, temperatures set conservatively for quick testing)
2. MVC's constraint satisfaction is inherently difficult without problem-specific knowledge

Potential improvements:

- Increase search budget (generations/iterations) by 2-4×
- Implement problem-specific initialization (greedy maximal covers)
- Use hybrid approaches (GA + local improvement, SA + TS refinement)
- Adjust penalty weights in Constraint Penalty fitness function

## 6.4 Scalability Observations

Large instances (100 nodes) show degraded performance for all algorithms:

- GA: feasibility drops to 15%
- SA: feasibility drops to 8%
- TS: feasibility drops to 10%

This suggests that the 5000-iteration/300-generation budgets are insufficient for large instances. Exponential time requirements are expected for NP-complete problems.

## 6.5 Design Choices and Limitations

### 6.5.1 Justified Design Choices

1. **Three distinct encodings:** Testing fundamentally different representations (binary vs. set vs. edge-centric) rather than variations provides broader insights. Validates that MVC is relatively insensitive to encoding.
2. **Three meta-heuristics:** GA, SA, and TS represent different algorithmic paradigms (population-based, trajectory-based, adaptive memory). Broader comparison than comparing variations of one algorithm.
3. **Random instance generation:** Controlled experimental conditions enable fair algorithm comparison. DIMACS instances would add external validity but complicate reproducibility.
4. **Conservative parameters:** Favored parameter modesty to demonstrate algorithm behavior under realistic time constraints rather than overfitting to instances.



### 6.5.2 Limitations

1. **Limited sample size:** 5 independent runs per configuration is modest. 30+ runs recommended for robust statistical conclusions. Feasibility rates should be interpreted as order-of-magnitude estimates.
2. **Instance size bounds:** Maximum of 100 nodes is small by modern standards. Large-scale benchmarking (1000+ nodes) would better reflect real-world applicability.
3. **Parameter tuning depth:** Grid search over mutation rates, population sizes, and temperature schedules was not performed. Results may not reflect algorithms' potential under optimal parameters.
4. **No hybrid approaches:** Pure algorithms tested without problem-specific enhancements. In practice, hybridization (GA + local search, multi-start SA) often improves performance.
5. **Statistical testing absent:** Wilcoxon signed-rank tests or ANOVA not performed. Observed differences in means may not be statistically significant.

## 7 Conclusion

This project investigates meta-heuristic optimization for the Minimum Vertex Cover problem, focusing on the interplay between algorithm design, problem encoding, and fitness function formulation. Key findings:

1. **Algorithm matters most:** Tabu Search outperforms Genetic Algorithms and Simulated Annealing in solution quality ( $23.0 \pm 4.71$  vs  $32.26 \pm 21.99$  for GA), indicating algorithm choice is the primary performance driver.
2. **Encoding choice is secondary:** Three fundamentally different encodings produce comparable results (differences within 1-2 nodes), suggesting MVC's problem structure is robust to representation.
3. **Trade-offs exist:** Simulated Annealing offers fastest computation (0.108s) but low solution quality; Genetic Algorithms achieve highest feasibility but at computational cost.
4. **Fitness function impacts feasibility:** Simpler fitness functions (Cover Size Minimization) achieve higher feasibility than sophisticated penalty-based approaches, contradicting intuition.

## 8 Future Work

1. **Extended parameter search:** Perform full factorial design over population size, mutation rate, cooling schedule to identify optimal parameter settings.
2. **Hybrid approaches:** Develop GA + Tabu Search hybrid, or Simulated Annealing with local TS refinement, to combine strengths of multiple paradigms.

3. **Large-scale benchmarking:** Test on DIMACS graph benchmarks (500-10000 nodes) for assessment of real-world scalability.
4. **Problem-specific enhancements:** Incorporate domain knowledge (greedy initialization, problem-aware crossover operators) to improve feasibility and solution quality.
5. **Statistical rigor:** Conduct 30+ runs per configuration with Wilcoxon tests to establish significance of observed differences.
6. **Comparative analysis:** Benchmark against known approximation algorithms and heuristics (2-approximation greedy algorithm, linear programming relaxations).
7. **Application studies:** Evaluate algorithms on real bioinformatics instances (protein interaction networks, metabolic pathways) for practical validation.

## References

- [1] Kirkpatrick, S., Gelatt Jr., C. D., & Vecchi, M. P. (1983). Optimization by Simulated Annealing. *Science*, 220(4598), 671–680.
- [2] Holland, J. H. (1975). *Adaptation in Natural and Artificial Systems*. University of Michigan Press.
- [3] Glover, F. (1989). Tabu Search: Part I. *ORSA Journal on Computing*, 1(3), 190–206.
- [4] Glover, F. (1990). Tabu Search: Part II. *ORSA Journal on Computing*, 2(1), 4–32.
- [5] Garey, M. R., & Johnson, D. S. (1976). The Complexity of Near-Optimal Graph Coloring. *Journal of the ACM*, 23(1), 43–49.
- [6] Hochbaum, D. S. (1983). Approximation Algorithms for the Set Cover Problem. *Journal of the ACM*, 30(3), 402–418.
- [7] Gendreau, M., Laporte, G., & Semet, F. (2010). Metaheuristics for Hard Combinatorial Optimization Problems. *International Journal of Metaheuristics*, 1(1), 3–25.
- [8] Blum, C., & Roli, A. (2003). Metaheuristics in Combinatorial Optimization: Overview and Conceptual Comparison. *ACM Computing Surveys*, 35(3), 268–308.
- [9] Barr, R. S., Hickman, B. L., Keller, C. M., & Wasserstein, R. L. (1995). Greedy Randomized Adaptive Search for the Vertex Coloring Problem. *Journal of Heuristics*, 1(1), 33–46.
- [10] Esbensen, H., & Smith, J. (1996). Genetic Algorithms and Graph Matching Problems. In *Proceedings of the International Conference on Evolutionary Computation* (pp. 400–405).
- [11] Karp, R. M., & Sahni, S. (1975). A Note on Algorithms for Minimum Vertex Cover. *Journal of the ACM*, 22(3), 326–333.
- [12] Mitchell, M. (1996). *An Introduction to Genetic Algorithms*. MIT Press.