

EVALUATION METRICS

Metric	Definition	Value ranges	System's Value	Interpretation	Verdict
Problem Success Rate	Fraction of problems for which the system successfully completes: 1. formal specification 2. test generation 3. reference solution construction 4. golden test suite generation	0.0 → complete failure 1.0 → end-to-end success on all problems	0.71	~71 out of 100 vague problems successfully completed the full pipeline. ~29% failed due to: <ul style="list-style-type: none">specification ambiguity,invalid test generation,brute/optimize d solution disagreement,or circular healing exhaustion.	Moderate, not excellent For a system operating entirely on LLM reasoning with no human curation , this is acceptable but clearly improvable .
Test Survival Rate	Fraction of generated candidate tests that survive into the final golden suite.	0.0 → all tests rejected 1.0 → no tests discarded	0.9749	Only ~2.5% of generated tests are discarded. Indicates: <ul style="list-style-type: none">strong constraint adherence,effective test validation,minimal over-generati on of invalid inputs.	Excellent. The test oracle is disciplined and conservative,
Consensus Converg	Fraction of successful problems where	0.0 → no agreement	0.9718	When the pipeline succeeds, ~97% of problems reach	Very strong

ence Rate	optimized and brute-force solutions converge to consistent outputs.	1.0 → full agreement		consensus. Disagreements are rare and typically resolved via healing.	Shows that treating brute force as “truth” is operationally effective.
Circular Healing Frequency	Fraction of successful problems where circular healing (Stage 5 rejecting Stage 4 tests) was triggered.	0.0 → never needed Higher is worse	0.0563	Circular healing was triggered in ~5.6% of successful runs. Indicates occasional mismatch between: <ul style="list-style-type: none">• test generator assumptions, and• reference solution feasibility.	Low but non-zero The system catches its own mistakes , which is far better than silently accepting bad tests.
Healing Yield	Fraction of healing attempts that successfully corrected the optimized solution	0.0 → healing useless 1.0 → healing always succeeds	0.4348	Only ~43% of healing attempts successfully fix errors. More than half fail or stall.	Weak. LLM self-correction is unreliable , BUT with user injected counterexamples, the reference solutions would increase the yield multi-folds.

LATENCY ANALYSIS

Stage	Avg Latency (s)	Assessment
Stage 3	3.29	Reasonable
Stage 4	9.50	Heavy but expected
Stage 5	13.92	Expensive
Stage 5.2	7.09	Moderate
End-to-end	37.37	Slow but acceptable, because the user receives real-time feedback on which

		stage our model is in.
--	--	------------------------

GOLD BENCHMARK - SYSTEM SOLUTION CORRECTNESS

Vague Query	Persona	Analysis	Verdict	Question ID
<p>Python solution to fix a binary search tree by identifying and swapping the values of the two incorrect nodes back to their correct positions.</p> <p>Input: root = [1,3,null,null,2] Output: [3,1,null,null,2]</p>	implementation_specific	<p>Test suite failure. Can't create a binary tree that's just two nodes away from a BST.</p> <p>Good syntax understanding, but lacks semantics!</p>	Failed	1
<p>Input: \$T\$: binary search tree. Output: \$T'\$: corrected binary search tree with original node values restored. There are exactly two node values misplaced.</p> <p>Input: root = [1,3,null,null,2] Output: [3,1,null,null,2]</p>	technical_short_hand	<p>Test suite failure. Can't create a binary tree that's just two nodes away from a BST.</p> <p>Good syntax understanding, but lacks semantics!</p>	Failed	1
<p>Input: \$nums: array\$. Output: \$int\$. Count continuous subarrays where \$ nums[i1] - nums[i2] <= 2\$ for all \$i1, i2\$ in the subarray.</p> <p>Input: nums = [5,4,2,4] Output: 8</p>	technical_short_hand	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	2
<p>Python solution to count the number of continuous segments in an array where the difference between the largest and smallest element is 2 or less.</p> <p>Input: nums = [5,4,2,4] Output: 8</p>	implementation_specific	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	2
<p>Input: \$T: array\$ of characters, \$n: integer\$. Output: \$int\$. Return minimum intervals</p>	technical_short_hand	<p>Formal Specification: Correct</p>	Passed	3

required to complete all tasks with at least \$n\$ intervals between identical tasks. Input: tasks = ["A", "A", "A", "B", "B", "B"], n = 2 Output: 8		# Attempts: 1 # Circular Heals: 1 # Test cycles: 1/1 # Solution cycles: 3/1 # User feedback: 0		
Python implementation to calculate the minimum number of intervals needed to finish a list of CPU tasks, ensuring a gap of at least 'n' intervals between tasks with the same label. Input: tasks = ["A", "A", "A", "B", "B", "B"], n = 2 Output: 8	implementation_specific	Formal Specification: Correct # Attempts: 1 # Circular Heals: 3 # Test cycles: 1/1/1 # Solution cycles: 3/3/1 # User feedback: 0	Passed	3
Input: \$S: string\$. Output: \$int\$. Return maximum product of lengths of two disjoint palindromic subsequences of \$\$\$. Input: s = "leetcodecom" Output: 9	technical_short_hand	Formal Specification: Correct # Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0	Passed	4
Python code to find the maximum product of lengths of two non-overlapping palindromic subsequences in a given string. Input: s = "leetcodecom" Output: 9	implementation_specific	Formal Specification: Correct # Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 1 (At first gave TLE, but upon user feedback, gave more optimized example)	Passed	4
Input: \$S: string\$. Output: \$List\$ of \$substrings\$. Return all 10-letter-long substrings of \$\$ that occur more than once. Input: s = "AAAAACCCCCAAAAACCCCCC AAAAAGGGTTT" Output: ["AAAAACCCCC", "CCCCCAAAA A"]	technical_short_hand	Formal Specification: Correct # Attempts: 1 # Circular Heals: 0 # Test cycles: 2 # Solution cycles: 1 # User feedback: 0	Passed	5
Python solution to identify all repeated 10-letter sequences in a given DNA string. Input: s = "AAAAACCCCCAAAAACCCCCC AAAAAGGGTTT" Output: ["AAAAACCCCC", "CCCCCAAAA A"]	implementation_specific	Formal Specification: Correct # Attempts: 1 # Circular Heals: 0 # Test cycles: 2 # Solution cycles: 1 # User feedback: 0	Passed	5

<p>Input: \$edges: array\$. \$distanceThreshold: int\$. Output: \$city: int\$. Return city with the least reachable cities under distanceThreshold, break ties by choosing the larger city number.</p> <p>Input: n = 4, edges = [[0,1,3],[1,2,1],[1,3,4],[2,3,1]], distanceThreshold = 4 Output: 3</p>	technical_short hand	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	6
<p>Python code to determine the city with the fewest reachable neighbors within a given distance, returning the largest city number in case of ties.</p> <p>Input: n = 4, edges = [[0,1,3],[1,2,1],[1,3,4],[2,3,1]], distanceThreshold = 4 Output: 3</p>	implementation_specific	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	6
<p>Input: \$nums: array\$, \$k: int\$. Output: \$int\$. Return the minimum number of k-bit flips to make all elements in \$nums\$ equal to 1, or -1 if not possible.</p> <p>Input: nums = [0,1,0], k = 1 Output: 2</p>	technical_short hand	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	7
<p>Python solution to calculate the least number of flips needed to turn an array of 0s and 1s into all 1s by flipping segments of length k.</p> <p>Input: nums = [0,1,0], k = 1 Output: 2</p>	implementation_specific	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	Passed	7
<p>Input: \$T\$: binary tree\$, \$L\$: linked list\$. Output: \$bool\$. Return \$True\$ if path from \$T\$ matches all elements of \$L\$ from head downwards.</p> <p>Input: head = [4,2,8], root = [1,4,4,null,2,2,null,1,null,6,8,null,null,null,1,3] Output: true</p>	technical_short hand	<p>Formal Specification: Wrong</p> <p># Attempts: 3 # Circular Heals: 3 # Test cycles: 3/3/3 # Solution cycles: 3/3/3 # User feedback: 0</p> <p>DIFFICULTY IN INTERPRETING THE LINKED-LIST and BINARY TREE INPUTS</p>	Failed	8

<p>Python solution to verify if a linked list matches a downward path in a binary tree, starting from the tree's nodes.</p> <p>Input: head = [4,2,8], root = [1,4,4,null,2,2,null,1,null,6,8,null,null,null,1,3] Output: true</p>	implementation_specific	<p>Formal Specification: Wrong</p> <p># Attempts: 10</p> <p>DIFFICULTY IN INTERPRETING THE LINKED-LIST and BINARY TREE INPUTS</p>	Failed - Took me 8 attempts	8
<p>Input: \$nums: array of integers\$. Output: \$int\$. Count permutations of \$nums\$ such that \$nums[i] + nums[i+1]\$ is a perfect square for all valid \$i\$\$. Input: nums = [1,17,8] Output: 2</p>	technical_short_hand	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 1 (Pass HINT as solution to Strong Prior Bias)</p> <p>Suffers from GPT's Strong Prior Bias. 'Solution' object has no attribute 'solve'</p>	Passed	9
<p>Python implementation to count distinct arrangements of an integer array where the sum of each adjacent pair results in a perfect square. Input: nums = [1,17,8] Output: 2</p>	implementation_specific	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 1 (Pass HINT as solution to Strong Prior Bias)</p> <p>Suffers from GPT's Strong Prior Bias. 'Solution' object has no attribute 'solve'</p>	Passed	9
<p>You are given the root of a binary tree. We install cameras on the tree nodes where each camera at a node can monitor its parent, itself, and its immediate children. Return the minimum number of cameras needed to monitor all nodes of the tree. Input: root = [0,0,null,0,0] Output: 1</p>	Formal (Leetcode)	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 2 # Solution cycles: 1 # User feedback: 0</p>	Passed	10

<p>A binary tree is uni-valued if every node in the tree has the same value.</p> <p>Given the root of a binary tree, return true if the given tree is uni-valued, or false otherwise.</p> <p>Input: root = [1,1,1,1,1,null,1] Output: true</p>	<p>Formal (Leetcode)</p>	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 2 # Solution cycles: 1 # User feedback: 0</p>	<p>Passed</p>	<p>11</p>
<p>Input: \$nums: array\$. Output: \$int\$. Return minimum swaps to group all 1's in \$nums\$ considering it as a circular array.</p> <p>Input: nums = [0,1,0,1,1,0,0] Output: 1</p>	<p>technical_short hand</p>	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	<p>Passed</p>	<p>12</p>
<p>Python solution to count land cells in a binary matrix that are completely enclosed by water, meaning they cannot reach the grid's edge.</p>	<p>implementation_specific</p>	<p>Formal Specification: Correct</p> <p># Attempts: 1 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	<p>Passed</p>	<p>13</p>
<p>Python solution to manage two lists where you can flip values in the first list based on commands, update the second list using values from the first list multiplied by a given number, and finally calculate the total of the second list after all operations.</p>	<p>implementation_specific</p>	<p>Formal Specification: Wrong</p> <p># Attempts: 2</p>	<p>Failed</p>	<p>14</p>
<p>Python implementation to group integers into pairs from an array and calculate the maximum sum of the minimums of those pairs.</p> <p>Input: nums = [1,4,3,2] Output: 4</p>	<p>implementation_specific</p>	<p>Formal Specification: Correct</p> <p># Attempts: 2 # Circular Heals: 0 # Test cycles: 1 # Solution cycles: 1 # User feedback: 0</p>	<p>Passed</p>	<p>15</p>