

Reinforcement-Learning-based Executive & Hierarchical Models for LLM Code Generation

Final Project Report

Amir Hosseinpour - Mani Vafapour

February 13, 2026

Project Notebooks (Google Colab)

`RL_3Episode.ipynb`
`RL_Hierarchical.ipynb`
`Combination_of_two_models.ipynb`

Abstract

This report documents a pipeline that trains two complementary RL agents to improve reasoning for LLM-based code generation. The first agent (“Executive”) is a task-level DQN trained to pick the correct sequence of actions to generate, test, debug, search, and stop on unit-test style programming tasks (HumanEval-like). The second agent (“Hierarchical / Planner”) is a single-episode DQN that learns the higher-level planning episode: planning → execution → integration → evaluation → stop. We describe the dataset, environment, state and reward design, model architectures (double + dueling DQN with action masking), training procedure, hyperparameters, results, and show the code used for all experiments.

1 Introduction

Large language models (LLMs) are powerful at code generation but often need iterative reasoning—generate, test, debug, and sometimes search for help. This project frames the iterative reasoning process as a reinforcement learning problem where agents learn to take high-level actions that orchestrate LLM code generation and correction. The pipeline has two models:

- **Executive model (DQN):** learns to pick actions within a small finite-state episode for a single task (generate, test, debug, search, stop).
- **Hierarchical / Planner model (DQN):** learns a single planning episode (planning, execution, integration, evaluation, stop), and uses the Executive model during the execution phase to solve subtasks.

We use HumanEval-style tasks (a small set of program synthesis tasks with unit tests) as our dataset for training and evaluation.

2 Dataset

2.1 Structure

The dataset used in experiments is a compact, HumanEval-like set of tasks. Each task entry contains:

- `task_id`: unique identifier
- `prompt`: the function signature and docstring given to the model
- `tests`: a short list of unit tests (assertions)
- `correct_code`: a correct implementation (used in simulation / seeding)
- `incorrect_code`: a representative incorrect output (to simulate LLM failures)
- `debugged_code`: the code after a successful debug (or after search + debug)

A minimal excerpt of the dataset loading and utilities used in experiments is included below (complete dataset appears in the appendix with all code cells). The dataset includes three training tasks that intentionally exercise the three episode types in the executive environment:

1. `task1_correct` – generate and test pass immediately.
2. `task2_incorrect` – generate, test fails, debug, test passes.
3. `task3_search` – requires a search step before the second debug can succeed.

3 Environment and Episode Design (Executive)

This section explains the precise finite-state environment used to train the Executive DQN.

3.1 Action Space

We use a small discrete action space:

$$\mathcal{A} = \{\text{generate, test, debug, search, stop}\}$$

implemented in code as an `IntEnum`. Each action corresponds to a high-level operation controlling (simulated) LLM behaviors.

3.2 Episode Types

The environment is designed so that each task’s `correct` trajectory belongs to one of three canonical episodes:

1. **Simple success:** `generate → test(pass) → stop`
2. **Debug once:** `generate → test(fail) → debug → test(pass) → stop`
3. **Search + double debug:** `generate → test(fail) → debug → test(fail) → search → debug → test(pass) → stop`

This enforces temporal preconditions and models realistic workflows: some failures can be fixed by a simple debug, others require an external hint/search.

3.3 State Definition

The state is represented as a dense vector composed of:

1. **Context embeddings:** embedding of the task prompt, the last produced code, and last test feedback. We used a sentence-transformers model (`all-MiniLM-L6-v2`) to obtain semantic embeddings and concatenated them.
2. **Flag vector:** small numeric flags appended to the embedding to explicitly encode temporal conditions. For the executive environment these flags are:

`flags = [has_generated, debug_count, has_searched, test_count, invalid_action, last_test_passed]`
where `has_generated ∈ {0, 1}`, `debug_count ∈ {0, 1, 2}`, etc.

Formally the state vector is:

$$s = \text{concat}(\text{embed}(\text{prompt}), \text{embed}(\text{code}), \text{embed}(\text{feedback}), \text{flags})$$

This explicit flagging makes it much easier for a tabular/ML-based agent to learn the temporal preconditions (e.g., "don't test before generating").

3.4 Reward Design

Reward structure is carefully shaped to encourage correct sequences and penalize illegal actions:

- Small step penalty for each step to encourage brevity: $r_{\text{step}} = -0.1$.
- Action-specific rewards:
 - Successful `generate`: +1.0
 - `test`: +1.0 if pass, -1.0 if fail
 - `debug`: +2.0 (successful debug transitions toward passing tests)
 - `search`: +1.5 (useful when required)
- Illegal action penalty (soft fail): a larger negative reward (-5.0) and the action is marked as `invalid_action`; episodes may continue but the final `stop` will be considered failure if any illegal action occurred.
- Terminal reward on `stop`:

$$r_{\text{stop}} = \begin{cases} +30.0 & \text{if success and no illegal actions} \\ -10.0 & \text{otherwise} \end{cases}$$

This mix combines dense shaping (small positive for progress) with sparse large rewards for successful termination, which stabilizes learning and encourages correct sequences.

4 Executive Agent: Network and Training

4.1 Network Architecture (Dueling Q-Network)

We use a **dueling Q-network**. The forward pass computes:

$$\begin{aligned} x &= \text{trunk}(s) \\ V(s) &= \text{value_head}(x) \quad (\text{scalar}) \end{aligned}$$

$$A(s, a) = \text{adv_head}(x) \quad (\text{vector over actions})$$

$$Q(s, a) = V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right)$$

Why dueling? it separates the estimation of state value $V(s)$ and advantages $A(s, a)$. This is helpful when many actions have similar value, improving learning stability.

4.2 Double DQN Targets

To reduce overestimation bias we use **Double DQN**: the online network selects the argmax next action while the target network evaluates it. Concretely:

$$a^* = \arg \max_{a'} Q_{\text{online}}(s', a')$$

$$y = r + \gamma Q_{\text{target}}(s', a^*)(1 - \text{done})$$

and the loss is MSE between $Q_{\text{online}}(s, a)$ and y .

4.3 Action Masking

We implement **action masking** derived from the flag vector so the agent does not choose impossible/illegal actions (e.g., test before generate). Masking is computed deterministically from the flag part of the state. Masking is applied both during policy selection (greedy) and exploration sampling (exploration respects mask when enabled). Masking reduces the effective action space and prevents the agent from wasting samples on obviously illegal moves.

4.4 Replay buffer, target network, and other details

The agent uses:

- Replay buffer (deque) with capacity 8000
- Target network updated every `target_update` steps
- Epsilon-greedy exploration with exponential decay to $\epsilon_{\min} = 0.05$
- Double-DQN loss and gradient clipping

4.5 Why standard DQN may fail and the motivation for our choices

A plain DQN without dueling, double updates, action masking, and careful state flags is unlikely to learn the required temporal preconditions. Reasons:

- **Sparse terminal reward:** Without reward shaping, the agent rarely reaches success and gets poor learning signal.
- **Illegal actions:** Without masking the agent will explore many obviously illegal actions, wasting samples.
- **Overestimation bias:** Single-network DQN tends to overestimate Q-values which harms stability.
- **State representation:** Without explicit flags the network must infer the temporal logic purely from embeddings a harder function to learn.

The architecture and state design directly address these failure modes.

5 Executive: Code and Explanation

Below is the core code used to define the dataset, the environment, the embedding and the DQN agent (executive). Each code block is followed by a brief explanation. The full code is preserved as-is so the experiment is reproducible.

5.1 Dataset + Utilities

Listing 1: Dataset, utilities, and run_tests_locally

```
1 # Standard imports, dataset (TRAINING_TASKS) and dirs
2 import os, random, copy, math, json, time
3 from typing import List, Dict, Any, Tuple
4 from pathlib import Path
5 from collections import deque
6 import numpy as np
7 import matplotlib.pyplot as plt
8 from tqdm.auto import tqdm
9 from rich.console import Console
10 from rich.progress import track
11 import torch
12 import torch.nn as nn
13 import torch.nn.functional as F
14 from torch.optim import Adam
15
16 console = Console()
17 os.makedirs("logs", exist_ok=True)
18 os.makedirs("saved_models", exist_ok=True)
19
20 # TRAINING_TASKS (modified to add a third "search-required" task based on HumanEval style)
21 TRAINING_TASKS = [
22     {
23         "task_id": "task1_correct",
24         "prompt": "def is_palindrome(s: str) -> bool:\n    \"\"\"Return True if string is\n    palindrome, else False.\"\n    \"\"\"",
25         "tests": [
26             "assert is_palindrome('racecar') == True",
27             "assert is_palindrome('hello') == False",
28             "assert is_palindrome('') == True",
29             "assert is_palindrome('a') == True",
30             "assert is_palindrome('madam') == True"
31         ],
32         "correct_code": """def is_palindrome(s: str) -> bool:\n    \"\"\"Return True if string is palindrome, else False.\"\n    \"\"\"",
33         "incorrect_code": """def is_palindrome(s: str) -> bool:\n    \"\"\"Return True if string is palindrome, else False.\"\n    \"\"\"",
34         "debugged_code": """def is_palindrome(s: str) -> bool:\n    \"\"\"Return True if string is palindrome, else False.\"\n    \"\"\"",
35     },
36     {
37         "task_id": "task2_incorrect",
38         "prompt": "def find_max(nums: List[int]) -> int:\n    \"\"\"Return the maximum number in a\n    list.\"\n    \"\"\"",
39         "tests": [
40             "assert find_max([1, 2, 3, 4, 5]) == 5",
41             "assert find_max([-1, -2, -3]) == -1",
42             "assert find_max([10]) == 10",
43             "assert find_max([5, 3, 9, 1, 7]) == 9",
44             "assert find_max([0, 0, 0]) == 0"
45     }
```

```

51     ],
52     "correct_code": """def find_max(nums: List[int]) -> int:
53         \"""Return the maximum number in a list.\"""
54     if not nums:
55         raise ValueError("List cannot be empty")
56     return max(nums)""",
57     "incorrect_code": """def find_max(nums: List[int]) -> int:
58         \"""Return the maximum number in a list.\"""
59     return nums[0]""",
60     "debugged_code": """def find_max(nums: List[int]) -> int:
61         \"""Return the maximum number in a list.\"""
62     if not nums:
63         raise ValueError("List cannot be empty")
64     return max(nums)"""
65 },
66 {
67     "task_id": "task3_search",
68     "prompt": "def sum_unique(nums: List[int]) -> int:\n    \"""Return the sum of elements\n        that appear exactly once in the list.\\"\">\n",
69     "tests": [
70         "assert sum_unique([1,2,2,3,4]) == 1+3+4",
71         "assert sum_unique([]) == 0",
72         "assert sum_unique([5,5,5]) == 0",
73         "assert sum_unique([1,2,3]) == 6",
74         "assert sum_unique([0,1,0,2,3,2]) == 1+3"
75     ],
76     "correct_code": """def sum_unique(nums: List[int]) -> int:
77         \"""Return the sum of elements that appear exactly once in the list.\"""
78     from collections import Counter
79     c = Counter(nums)
80     return sum(x for x, cnt in c.items() if cnt == 1)""",
81     "incorrect_code": """def sum_unique(nums: List[int]) -> int:
82         \"""Return the sum of elements that appear exactly once in the list.\"""
83     # Wrong: sums unique set rather than only elements with count == 1
84     return sum(set(nums))""",
85     "debugged_code": """def sum_unique(nums: List[int]) -> int:
86         \"""Return the sum of elements that appear exactly once in the list.\"""
87     from collections import Counter
88     c = Counter(nums)
89     return sum(x for x, cnt in c.items() if cnt == 1)"""
90 }
91 ]
92
93 # Utility local test runner (copied from your code to run unit-tests on strings)
94 def run_tests_locally(code_str: str, tests: List[str]) -> Tuple[bool, Dict[str, Any]]:
95     import traceback
96     feedback = {"number_passed": 0, "total": len(tests), "failures": []}
97     exec_globals = {}
98     try:
99         exec("from typing import List", exec_globals)
100        exec(code_str, exec_globals)
101    except Exception as e:
102        return False, {"error": str(e), "trace": traceback.format_exc()}
103    for test in tests:
104        try:
105            exec(test, exec_globals)
106            feedback["number_passed"] += 1
107        except AssertionError:
108            feedback["failures"].append({"test": test, "error": "AssertionError"})
109        except Exception as e:
110            feedback["failures"].append({"test": test, "error": str(e)})
111    return feedback["number_passed"] == feedback["total"], feedback
112

```

```
113 | console.log("Dataset and basic utilities loaded.")
```

Explanation: This block defines the simulated HumanEval-style tasks and a local test runner that executes code strings and checks their unit tests. The simulated tasks include correct/incorrect/debugged code variants so we can simulate the LLM behavior deterministically during training and seeding.

5.2 Action enum and small execution demo

Listing 2: Action enum and print demo

```
1 from enum import IntEnum
2
3 class Action(IntEnum):
4     GENERATE = 0
5     TEST = 1
6     DEBUG = 2
7     SEARCH = 3 # NEW action
8     STOP = 4
9
10 # Human-readable names used by the environment + logging
11 Action.NAMES = {
12     Action.GENERATE: "generate",
13     Action.TEST: "test",
14     Action.DEBUG: "debug",
15     Action.SEARCH: "search",
16     Action.STOP: "stop",
17 }
18
19 print("Action space:", Action.NAMES)
```

Explanation: Defines a compact action space. Human-readable names are stored in a dictionary for logging.

5.3 Environment (CodeGenEnv)

Listing 3: CodeGenEnv: finite-state environment for executive

```
1 #cell-1 # STRICT finite-state environment illegal actions used to terminate previously
2 # MODIFIED: add 'invalid_action' flag and 'last_test_passed'. illegal actions set
3 # invalid_action True.
4 # last_test_passed included in state so agent can mask test-after-pass.
5
6 class CodeGenEnv:
7     def __init__(self, task: Dict, max_steps: int = 10):
8         self.task = task
9         self.max_steps = max_steps
10        self.reset()
11
12    def reset(self):
13        self.steps = 0
14        self.done = False
15        self.success = False
16        self.current_code = ""
17        self.test_feedback = ""
18        self.has_generated = False
19        self.test_results = [] # list of bools
20        self.debug_count = 0
21        self.has_searched = False
22        self.invalid_action = False # NEW: any illegal action sets this True
23        return self._state()
```

```

23
24     def _fail(self):
25         # SOFT FAIL: give a strong negative reward but do NOT force immediate termination.
26         # This enables exploration to continue and the agent to learn the correct preconditions.
27         return -5.0
28
29     def step(self, action: int):
30         if self.done:
31             raise RuntimeError("Episode already finished")
32         self.steps += 1
33         reward = -0.1
34         info = {"action": Action.NAMES[action]}
35
36         # ----- GENERATE -----
37         if action == Action.GENERATE:
38             if self.has_generated:
39                 reward += self._fail()
40                 self.invalid_action = True
41             else:
42                 self.current_code = sim_llm.generate_code(self.task["task_id"])
43                 self.has_generated = True
44                 reward += 1.0
45
46         # ----- TEST -----
47         elif action == Action.TEST:
48             # illegal: testing before generation
49             if not self.has_generated:
50                 reward += self._fail()
51                 self.invalid_action = True
52             # illegal: too many tests
53             elif len(self.test_results) >= 4:
54                 reward += self._fail()
55                 self.invalid_action = True
56             # illegal: re-testing immediately after a passed test (prevents extra test shortcuts)
57             elif self.test_results and self.test_results[-1] is True:
58                 reward += self._fail()
59                 self.invalid_action = True
60             else:
61                 passed, fb = run_tests_locally(self.current_code, self.task["tests"])
62                 self.test_results.append(passed)
63                 self.test_feedback = f"Tests {fb['number_passed']}/{fb['total']}"
64                 reward += 1.0 if passed else -1.0
65
66         # ----- DEBUG -----
67         elif action == Action.DEBUG:
68             # Only allowed when there was at least one test, last test failed, debug_count < 2,
69             # and if this is the 2nd debug it must either have searched where required (handled
70             # later)
71             if (
72                 not self.test_results
73                 or self.test_results[-1] is True
74                 or self.debug_count >= 2
75                 or (self.debug_count == 1 and not self.has_searched and len(self.test_results) != 2)
76             ):
77                 reward += self._fail()
78                 self.invalid_action = True
79             else:
80                 self.debug_count += 1
81                 self.current_code = sim_llm.debug_code(
82                     self.task["task_id"],
83                     self.current_code,
84                     searched=self.has_searched

```

```

84         )
85         reward += 2.0
86
87     # ----- SEARCH -----
88     elif action == Action.SEARCH:
89         if (
90             self.has_searched
91             or self.debug_count != 1
92             or self.test_results != [False, False]
93         ):
94             reward += self._fail()
95             self.invalid_action = True
96         else:
97             self.has_searched = True
98             reward += 1.5
99
100    # ----- STOP -----
101    elif action == Action.STOP:
102        self.done = True
103        tid = self.task["task_id"]
104        # Only allow success if no illegal actions occurred during the episode
105        if tid != "task3_search" and self.test_results == [True] and not self.invalid_action:
106            self.success = True
107        elif (
108            tid != "task3_search"
109            and self.test_results == [False, True]
110            and self.debug_count == 1
111            and not self.invalid_action
112        ):
113            self.success = True
114        elif (
115            tid == "task3_search"
116            and self.test_results == [False, False, True]
117            and self.debug_count == 2
118            and self.has_searched
119            and not self.invalid_action
120        ):
121            self.success = True
122        reward += 30.0 if self.success else -10.0
123
124    # ----- STEP LIMIT -----
125    if self.steps >= self.max_steps and not self.done:
126        self.done = True
127        reward -= 10.0
128
129        info["success"] = self.success
130        info["invalid_action"] = self.invalid_action # debugging helper
131        return self._state(), reward, self.done, info
132
133    def _state(self):
134        # include flags in the state embedding so the agent can easily learn the temporal
135        # preconditions
136        last_test_passed = bool(self.test_results[-1]) if self.test_results else False
137        return make_state_embedding(
138            self.task["prompt"],
139            self.current_code,
140            self.test_feedback,
141            has_generated=self.has_generated,
142            debug_count=self.debug_count,
143            has_searched=self.has_searched,
144            test_count=len(self.test_results),
145            invalid_action=self.invalid_action,
146            last_test_passed=last_test_passed

```

Explanation: The environment enforces legal action preconditions and returns shaped rewards. Illegal actions do not immediately terminate but set an invalid flag and apply a heavy negative reward so the agent learns to avoid them. The termination success check compares the sequence of test results and flags to the expected canonical sequences.

5.4 State embedding

Listing 4: State embedding using sentence-transformers + flags

```

1 #cell-2
2 # Embedding (same approach as you used) - create state vector from prompt, code, and feedback
3 # MODIFIED: append small numeric flag vector (has_generated, debug_count, has_searched,
4 # to make the temporal aspects explicit to the agent.
5
6 from sentence_transformers import SentenceTransformer
7 embedder = SentenceTransformer("sentence-transformers/all-MiniLM-L6-v2")
8
9 def make_state_embedding(task_prompt: str, last_code: str, last_test_feedback: str,
10                         has_generated: bool = False, debug_count: int = 0,
11                         has_searched: bool = False, test_count: int = 0,
12                         invalid_action: bool = False, last_test_passed: bool = False) -> np.
13                         ndarray:
14     task_emb = embedder.encode([task_prompt], show_progress_bar=False)
15     code_emb = embedder.encode([last_code or ""], show_progress_bar=False)
16     feedback_emb = embedder.encode([last_test_feedback or ""], show_progress_bar=False)
17     state_vec = np.concatenate([task_emb[0], code_emb[0], feedback_emb[0]])
18     # append simple numeric flags (small in magnitude)
19     flags = np.array([1.0 if has_generated else 0.0,
20                      float(debug_count),
21                      1.0 if has_searched else 0.0,
22                      float(test_count),
23                      1.0 if invalid_action else 0.0,
24                      1.0 if last_test_passed else 0.0], dtype=np.float32)
25
26 # Quick smoke
27 sample_state = make_state_embedding(TRAINING_TASKS[0]["prompt"], "", "", False, 0, False, 0,
28                                     False)
29 console.log(f"State dim {sample_state.shape[0]}")
STATE_DIM = sample_state.shape[0]

```

Explanation: We use semantic embeddings for prompt/code/feedback so the network sees content features, and append a concise numeric flag vector for temporal preconditions. This hybrid semantic + symbolic state has proven effective for this setting.

5.5 Simulated LLM

Listing 5: Simulated LLM used in environment

```

1 #cell-3
2 # SimulatedLLM with enforced SEARCH dependency
3
4 class SimulatedLLM:
5     def __init__(self, tasks: List[Dict]):
6         self.tasks = {task["task_id"]: task for task in tasks}
7
8     def generate_code(self, task_id: str) -> str:
9         task = self.tasks[task_id]

```

```

10     return task["correct_code"] if task_id == "task1_correct" else task["incorrect_code"]
11
12 def debug_code(self, task_id: str, current_code: str, searched: bool = False) -> str:
13     task = self.tasks[task_id]
14
15     if task_id == "task1_correct":
16         return task["correct_code"]
17
18     if task_id == "task2_incorrect":
19         return task["debugged_code"]
20
21     if task_id == "task3_search":
22         if not searched:
23             # STILL WRONG guaranteed to fail tests
24             return """def sum_unique(nums: List[int]) -> int:
25 # Almost right but intentionally wrong
26 from collections import Counter
27 c = Counter(nums)
28 return sum(x for x, cnt in c.items() if cnt >= 1) # wrong condition
29 """
30         # Only after SEARCH can it be correct
31         return task["debugged_code"]
32
33     return current_code
34
35 def search(self, task_id: str) -> bool:
36     return task_id == "task3_search"
37
38
39 # REQUIRED: instantiate the simulated LLM
40 sim_llm = SimulatedLLM(TRAINING_TASKS)

```

Explanation: A deterministic simulated LLM simplifies training and debugging: the model's generate and debug behaviors are deterministic and designed to exercise the three episode types. For `task3_search` the debug step will only produce correct code after a `search` has occurred.

5.6 Dueling Double DQN Agent (Executive)

Listing 6: Dueling Double DQN agent with action masking

```

1 #cell-4
2 # ----- DQN agent (Double + Dueling) -----
3 # DQN with target network, replay buffer, epsilon-greedy, and prioritized-ish uniform sampling.
4 # MODIFIED: select_action now applies action masking derived from the flag vector at the end of
5 # state.
6
6 class DuelingQNetwork(nn.Module):
7     def __init__(self, input_dim, hidden=[512, 256], output_dim=len(Action.NAMES)):
8         super().__init__()
9         # shared trunk
10        layers = []
11        prev = input_dim
12        for h in hidden:
13            layers.append(nn.Linear(prev, h))
14            layers.append(nn.ReLU())
15            prev = h
16        self.trunk = nn.Sequential(*layers)
17        # value stream
18        self.value_head = nn.Sequential(
19            nn.Linear(prev, prev//2 if prev//2>0 else 32),
20            nn.ReLU(),

```

```

21         nn.Linear(prev//2 if prev//2>0 else 32, 1)
22     )
23     # advantage stream
24     self.adv_head = nn.Sequential(
25         nn.Linear(prev, prev//2 if prev//2>0 else 32),
26         nn.ReLU(),
27         nn.Linear(prev//2 if prev//2>0 else 32, output_dim)
28     )
29
30     def forward(self, x):
31         x = self.trunk(x)
32         value = self.value_head(x)
33         adv = self.adv_head(x)
34         # combine into Q-values:  $Q(s,a) = V(s) + (A(s,a) - \text{mean}_a A(s,a))$ 
35         q = value + (adv - adv.mean(dim=1, keepdim=True))
36         return q
37
38     class ReplayBuffer:
39         def __init__(self, capacity=8000):
40             self.buffer = deque(maxlen=capacity)
41         def push(self, s, a, r, ns, done):
42             self.buffer.append((s, a, r, ns, done))
43         def sample(self, batch_size):
44             batch = random.sample(self.buffer, min(batch_size, len(self.buffer)))
45             s,a,r,ns,d = zip(*batch)
46             return (np.stack(s), np.array(a), np.array(r, dtype=np.float32), np.stack(ns), np.array(
47                 d, dtype=np.float32))
48         def __len__(self):
49             return len(self.buffer)
50
51     class DQNAgent:
52         def __init__(self, state_dim, action_dim=len(Action.NAMES), hidden=[512,256], lr=1e-4, gamma
53             =0.99,
54             buffer_size=8000, batch_size=64, target_update=500, device=None, mask_actions=
55             True):
56             self.device = device or ("cuda" if torch.cuda.is_available() else "cpu")
57             self.qnet = DuelingQNetwork(state_dim, hidden, action_dim).to(self.device)
58             self.target = DuelingQNetwork(state_dim, hidden, action_dim).to(self.device)
59             self.target.load_state_dict(self.qnet.state_dict())
60             self.opt = Adam(self.qnet.parameters(), lr=lr)
61             self.gamma = gamma
62             self.buffer = ReplayBuffer(capacity=buffer_size)
63             self.batch_size = batch_size
64             self.action_dim = action_dim
65             self.eps = 1.0
66             self.eps_min = 0.05
67             self.eps_decay = 0.995
68             self.learn_steps = 0
69             self.target_update = target_update
70             self.mask_actions = mask_actions
71
72         def _compute_action_mask(self, state: np.ndarray):
73             # state ends with six flags: [has_generated, debug_count, has_searched, test_count,
74             # invalid_action, last_test_passed]
75             # assume state is numpy array
76             flags = state[-6:]
77             has_generated = bool(flags[0])
78             debug_count = int(flags[1])
79             has_searched = bool(flags[2])
80             test_count = int(flags[3])
81             invalid_action = bool(flags[4])
82             last_test_passed = bool(flags[5])

```

```

80     mask = np.ones(self.action_dim, dtype=bool) # True = allowed
81
82     # GENERATE allowed only if not generated yet
83     if has_generated:
84         mask[Action.GENERATE] = False
85
86     # TEST allowed only if generated, not too many tests, and last test was not a pass
87     if (not has_generated) or test_count >= 4 or last_test_passed:
88         mask[Action.TEST] = False
89
90     # DEBUG allowed only if there is at least one test and last test failed, debug_count < 2,
91     # and second debug requires search (handled conservatively)
92     if test_count == 0 or last_test_passed or debug_count >= 2:
93         mask[Action.DEBUG] = False
94     else:
95         # If this is the second debug (debug_count==1), require that test_count == 2 (we
96         # conservatively require two tests before second debug)
97         if debug_count == 1 and test_count != 2:
98             mask[Action.DEBUG] = False
99
100    # SEARCH allowed only if we have exactly two tests, both failed (we approximate using
101    # test_count==2 and last_test_passed==False), and debug_count==1 and not already
102    # searched
103    if not (debug_count == 1 and test_count == 2 and (not last_test_passed) and (not
104        has_searched)):
105        mask[Action.SEARCH] = False
106
107    # STOP allowed only if last test passed and no invalid action (conservative)
108    if not (last_test_passed and (not invalid_action)):
109        mask[Action.STOP] = False
110
111    return mask
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
def select_action(self, state):
    # state: np.array
    if random.random() < self.eps:
        # when exploring, respect mask by sampling only legal actions when possible
        if self.mask_actions:
            mask = self._compute_action_mask(state)
            legal_indices = np.flatnonzero(mask)
            if len(legal_indices) > 0:
                return int(np.random.choice(legal_indices))
            # fallthrough to random full action
            return random.randrange(self.action_dim)

    s = torch.FloatTensor(state).unsqueeze(0).to(self.device)
    with torch.no_grad():
        qvals = self.qnet(s).cpu().numpy().squeeze(0)
    if self.mask_actions:
        mask = self._compute_action_mask(state)
        legal_q = np.where(mask, qvals, -1e9) # very low for illegal actions
        if legal_q.max() <= -1e8:
            # no legal action (should be rare) -> fallback to random
            return random.randrange(self.action_dim)
        return int(int(legal_q.argmax()))
    else:
        return int(int(qvals.argmax()))

def push_transition(self, s,a,r,ns,done):
    self.buffer.push(s,a,r,ns,done)

def train_step(self):
    if len(self.buffer) < 32:

```

```

138     return 0.0
139     s,a,r,ns,d = self.buffer.sample(self.batch_size)
140     s = torch.FloatTensor(s).to(self.device)
141     a = torch.LongTensor(a).to(self.device)
142     r = torch.FloatTensor(r).to(self.device)
143     ns = torch.FloatTensor(ns).to(self.device)
144     d = torch.FloatTensor(d).to(self.device)
145
146     # current Q-values
147     qvals = self.qnet(s).gather(1, a.unsqueeze(1)).squeeze(1)
148
149     # ---- Double DQN target calculation ----
150     # use online network to select best next action, use target network to evaluate its Q
151     with torch.no_grad():
152         next_actions = self.qnet(ns).argmax(dim=1, keepdim=True) # shape (batch,1)
153         next_q_target = self.target(ns).gather(1, next_actions).squeeze(1)
154         target = r + self.gamma * next_q_target * (1 - d)
155
156     loss = F.mse_loss(qvals, target)
157     self.opt.zero_grad()
158     loss.backward()
159     nn.utils.clip_grad_norm_(self.qnet.parameters(), 0.5)
160     self.opt.step()
161
162     self.learn_steps += 1
163     if self.learn_steps % self.target_update == 0:
164         self.target.load_state_dict(self.qnet.state_dict())
165
166     # decay epsilon
167     self.eps = max(self.eps_min, self.eps * self.eps_decay)
168     return float(loss.item())
169
170 console.log("DQN agent (Double + Dueling + Action masking) ready.")

```

Explanation: This block defines a dueling Q-network, replay buffer, and agent. The agent uses action masking derived from the appended flags to avoid illegal actions. The train step uses Double-DQN target computation and gradient clipping. Hyperparameters (hidden sizes, learning rate, buffer size, batch size, gamma) are all visible and easily adjustable.

5.7 Training Loop and Results (Executive)

Listing 7: Training loop and seeding with expert episodes

```

1 #cell-5
2 # Training loop for DQN and logging / metrics (plots)
3 def train_dqn(agent: DQNAgent, tasks: List[Dict], num_episodes=1200, max_steps=10, log_every=50,
4             seed_expert_episodes_per_task=8):
5     metrics = {
6         "episode_rewards": [],
7         "success_rate_window": [],
8         "episode_steps": [],
9         "losses": [],
10        "sample_eff": []
11    }
12    success_history = deque(maxlen=100)
13    total_env_steps = 0
14
15    # --- Seed replay buffer with expert episodes (scripted correct trajectories) ---
16    ACTION_NAME_TO_INT = {v:k for k,v in Action.NAMES.items()}
17
18    expert_seqs = {

```

```

19     "task1_correct": ["generate", "test", "stop"],  

20     "task2_incorrect": ["generate", "test", "debug", "test", "stop"],  

21     # For task3, the correct sequence is: gen -> test(fail) -> debug -> test(fail) -> search  

22     # -> debug -> test(pass) -> stop  

23     "task3_search": ["generate", "test", "debug", "test", "search", "debug", "test", "stop"]  

24   }  

25  

26   # Push several expert episodes into the buffer  

27   for task in tasks:  

28     seq = expert_seqs.get(task["task_id"], None)  

29     if seq is None:  

30       continue  

31     for _ in range(seed_expert_episodes_per_task):  

32       env = CodeGenEnv(task, max_steps=max_steps)  

33       s = env.reset()  

34       for action_name in seq:  

35         a = ACTION_NAME_TO_INT[action_name]  

36         ns, r, done, info = env.step(a)  

37         agent.push_transition(s, a, r, ns, float(done))  

38         s = ns  

39         if done:  

40           break  

41  

42   best_success = -1.0 # track best "success(last100)" to save best model  

43   best_avg_reward = -1e9  

44  

45   for ep in tqdm(range(1, num_episodes+1), desc="DQN Training"):  

46     task = random.choice(tasks)  

47     env = CodeGenEnv(task, max_steps=max_steps)  

48     state = env.reset()  

49     ep_reward = 0.0  

50     ep_loss = 0.0  

51  

52     for t in range(max_steps):  

53       action = agent.select_action(state)  

54       next_state, reward, done, info = env.step(action)  

55       agent.push_transition(state, action, reward, next_state, float(done))  

56       loss = agent.train_step()  

57       state = next_state  

58       ep_reward += reward  

59       ep_loss += loss  

60       total_env_steps += 1  

61       if done:  

62         break  

63  

64     metrics["episode_rewards"].append(ep_reward)  

65     metrics["losses"].append(ep_loss / (t+1) if (t+1)>0 else 0.0)  

66     metrics["episode_steps"].append(t+1)  

67     success_history.append(1.0 if done and info.get("success", False) else 0.0)  

68  

69     metrics["success_rate_window"].append(np.mean(success_history))  

70     metrics["sample_eff"].append(sum(metrics["episode_rewards"])/ (total_env_steps + 1e-8))  

71  

72     # check for best model (prefer success_rate over avg reward)  

73     current_success = metrics["success_rate_window"][-1]  

74     current_avg_reward = np.mean(metrics["episode_rewards"][-log_every:])  

75     if current_success > best_success or (current_success == best_success and  

76       current_avg_reward > best_avg_reward):  

77       best_success = current_success  

78       best_avg_reward = current_avg_reward  

79       # save best model  

80       torch.save(agent.qnet.state_dict(), "saved_models/dqn_best.pt")

```

```

79     console.log(f"[green]Ep {ep:4d} | New best model saved: success(last100)={
80         best_success:.3%}, avg_reward={best_avg_reward:.3f}[/green]")
81
82     if ep % log_every == 0 or ep == 1:
83         console.log(f"[blue]Ep {ep:4d} | AvgReward {np.mean(metrics['episode_rewards'])[-log_every:]):.3f} | "
84             f"Success(last100) {metrics['success_rate_window'][-1]:.3%} | Eps {agent.
85             eps:.3f}")
86
87 # Quick run for smoke/training (lower episodes). Increase num_episodes to train more.
88 agent = DQNAgent(state_dim=STATE_DIM, action_dim=len(Action.NAMES), hidden=[512,256], lr=3e-4,
89     buffer_size=8000, batch_size=64, target_update=200, mask_actions=True)
90 metrics = train_dqn(agent, TRAINING_TASKS, num_episodes=500, max_steps=10, log_every=50)
91
92 # Save final model (also keep best which was saved during training)
93 torch.save(agent.qnet.state_dict(), "saved_models/dqn_final.pt")
94 console.log("Training complete and final model saved.")

```

Explanation: The training loop seeds the replay buffer with multiple expert demonstrations per task (this speeds up learning and ensures the agent sees reasonable trajectories). It then runs episodes with epsilon-greedy exploration, trains on minibatches, periodically updates the target network, and saves the best model based on recent success rate. The ‘metrics‘ object captures reward, loss, steps and rolling success rate for plotting.

5.8 Training Outputs / Visuals

Insert the print output image that shows the console print from the state dimension print or other console logs here:

Figure 1: Console output (example): state dimension and basic checks. (Place the file p1.png in the same folder.)

Note: ‘p1.png‘ is the requested image showing the print output after the action enum and state-dimension shows. Place your actual screenshot with that file name next to this .tex.

6 Executive: Test Rollouts and Results

After training we load the best model and run demonstration rollouts across the three tasks. The code saving/loading and demo rollouts are shown next:

Listing 8: Load, test rollouts and log examples

```

1 #cell-7
2 # Save final agent (and quick test on training tasks)
3 # Modified to load the best model saved during training (falls back to final if not present)
4
5 best_path = "saved_models/dqn_best.pt"
6 final_path = "saved_models/dqn_final.pt"
7
8 # load best if available
9 if os.path.exists(best_path):
10     try:
11         agent.qnet.load_state_dict(torch.load(best_path, map_location=agent.device))
12         console.log(f"Loaded best model from {best_path}")
13     except Exception as e:

```

```

14     console.log(f"[red]Failed to load best model ({e}), loading final model if present.[/red")
15         ])
16     if os.path.exists(final_path):
17         agent.qnet.load_state_dict(torch.load(final_path, map_location=agent.device))
18         console.log(f"Loaded final model from {final_path}")
19     else:
20         if os.path.exists(final_path):
21             agent.qnet.load_state_dict(torch.load(final_path, map_location=agent.device))
22             console.log(f"No best model found; loaded final model from {final_path}")
23         else:
24             console.log("[yellow]No saved model found; running with current agent weights.[/yellow]")
25
26 torch.save(agent.qnet.state_dict(), "saved_models/dqn_final.pt")
27 console.log("Final model saved to saved_models/dqn_final.pt")
28
29 # Quick run: show three example rollouts per task with printed paths
30 for task in TRAINING_TASKS:
31     console.log(f"[bold]Demo rollouts for {task['task_id']}[/bold]")
32     for i in range(3):
33         env = CodeGenEnv(task, max_steps=10)
34         s = env.reset()
35         path = []
36         for _ in range(12):
37             a = agent.select_action(s)
38             path.append(Action.NAMES[a])
39             s, r, done, info = env.step(a)
40             if done: break
41         console.log(f" Path {i+1}: {path} | success={info.get('success', False)}")
```

Explanation: Demonstrates that the trained Executive agent reliably follows the canonical correct episodes for each task (shown in console logs and saved screenshots).

Insert the training plot / success log screenshot (your reported best result) here:

```

[13:14:37] Ep 331 | New best model saved: success(last100)=97.000%, avg_reward=32.164
...
[13:14:39] Ep 332 | New best model saved: success(last100)=97.000%, avg_reward=32.288
[13:14:49] Ep 330 | New best model saved: success(last100)=97.000%, avg_reward=32.684
Ep 330 | Avgreward 32.684 | Success(last100) 97.000% | Eps 0.050
[13:14:52] Ep 331 | New best model saved: success(last100)=97.000%, avg_reward=32.728
[13:14:57] Ep 361 | New best model saved: success(last100)=97.000%, avg_reward=32.772
[13:14:58] Ep 362 | New best model saved: success(last100)=98.000%, avg_reward=32.788
[13:15:03] Ep 368 | New best model saved: success(last100)=98.000%, avg_reward=32.832
[13:15:04] Ep 369 | New best model saved: success(last100)=98.000%, avg_reward=32.892
[13:15:11] Ep 381 | New best model saved: success(last100)=99.000%, avg_reward=32.800
[13:15:15] Ep 385 | New best model saved: success(last100)=99.000%, avg_reward=32.844
[13:15:16] Ep 386 | New best model saved: success(last100)=99.000%, avg_reward=32.884
[13:15:18] Ep 388 | New best model saved: success(last100)=99.000%, avg_reward=32.932
[13:15:19] Ep 389 | New best model saved: success(last100)=99.000%, avg_reward=32.976
[13:15:20] Ep 390 | New best model saved: success(last100)=99.000%, avg_reward=33.036
[13:15:23] Ep 396 | New best model saved: success(last100)=99.000%, avg_reward=33.052
[13:15:25] Ep 400 | New best model saved: success(last100)=100.000%, avg_reward=33.036
Ep 400 | Avgreward 33.036 | Success(last100) 100.000% | Eps 0.050
[13:16:03] Ep 450 | Avgreward 31.100 | Success(last100) 98.000% | Eps 0.050
[13:16:46] Ep 500 | Avgreward 30.584 | Success(last100) 95.000% | Eps 0.050
Training complete and final model saved.
```

Figure 2: Executive training log: best model saved and performance metrics (example screenshot).

Reported best result (from your run):

[13:15:25] Ep 400 | New best model saved: success(last100)=100.000%, avg_reward=33.036

This indicates the executive agent reached a consistent perfect success rate over a last-100 window in that run.

7 Hierarchical / Planner Model

The planner learns a single canonical planning episode:

planning → execution → integration → evaluation → stop

During **execution**, the planner spawns multiple subtasks (as produced by the planning action). Each subtask is then handled by the Executive agent (we either run the Executive network or simulate it, depending on the experiment) to produce helper functions. Integration merges the helpers and evaluation runs tests on the integrated code.

7.1 Planner State and Flags

Planner state is a concatenation of semantic embeddings and a small symbolic flag vector. Concretely, the planner state contains:

- the task prompt, for example: `Write a Python function compute(a,b) that returns (a+b) to the power 2/3.`
- the last code produced (the integration output).
- the last test feedback (a short textual summary like “Tests 2/3”).
- a compact flag vector that encodes temporal and structural information. We represent it as:
`[has_planned, executed_count, integrated, evaluated, invalid_action, last_test_passed]`.

This encoding lets the planner learn the single-episode temporal workflow.

7.2 Planner Environment and Rewards

Rewards are shaped to encourage correct progress through the planning episode:

- planning: +1.0
- execution: +1.5
- integration: +2.0
- evaluation: +5.0 if pass else -1.0
- terminal success: +30.0 else -10.0
- step penalty: -0.1
- illegal/invalid actions: -5.0 (soft fail)

These are analogous to the Executive environment but targeted to the planner’s single-episode workflow.

7.3 Planner DQN

We reuse the dueling Double-DQN architecture and replay buffer, with the planner-specific state dimension and action set:

$$\mathcal{A}_{\text{planner}} = \{\text{planning, execution, integration, evaluation, stop}\}.$$

During early training exploration, the planner samples from the full action space (no mask) to discover failure modes; exploitation uses a mask derived from flags to avoid impossible choices. This gives the agent the ability to discover temporal dependencies but prevents catastrophic exploitation errors later.

8 Planner: Code and Explanation

The planner code closely mirrors the executive code; here are the core parts and the training loop (complete listing in appendix):

Listing 9: Planner training task, env, network, and training loop

```
1 # planner_dqn_final.py
2 # Single-episode DQN trainer: planning -> execution -> integration -> evaluation -> stop
3 # Final version: exploration samples from full action space (no mask) so agent must learn.
4
5 import os, random, json, math, time
6 from collections import deque
7 from typing import List, Dict, Any, Tuple
8 import numpy as np
9 from tqdm.auto import tqdm
10 from rich.console import Console
11 import torch
12 import torch.nn as nn
13 import torch.nn.functional as F
14 from torch.optim import Adam
15
16 console = Console()
17 os.makedirs("logs", exist_ok=True)
18 os.makedirs("saved_models", exist_ok=True)
19
20 # ----- Training task (single planner task) -----
21 PLANNER_TASK = {
22     "task_id": "planner_compute_2o3",
23     "prompt": "Write a Python function `compute(a, b)` that returns  $(a + b)^{(2/3)}$ .",
24     "tests": [
25         "import math",
26         "assert math.isclose(compute(1,2), (1+2)**(2/3), rel_tol=1e-9)",
27         "assert math.isclose(compute(0,0), (0+0)**(2/3), rel_tol=1e-9)",
28         "assert math.isclose(compute(8,1), (8+1)**(2/3), rel_tol=1e-9)",
29         "assert math.isclose(compute(-1,8), (-1+8)**(2/3), rel_tol=1e-9)",
30     ],
31     "correct_code": """def compute(a, b):
32     # compute  $(a + b)^{(2/3)}$ 
33     return (a + b) ** (2.0/3.0)""",
34     "incorrect_code": """def compute(a, b):
35     # wrong: uses integer division or wrong exponent
36     return (a + b) ** (1/2)""",
37 }
38
39 # ----- Local test runner -----
40 def run_tests_locally(code_str: str, tests: List[str]) -> Tuple[bool, Dict[str, Any]]:
41     import traceback
42     feedback = {"number_passed": 0, "total": len(tests), "failures": []}
43     exec_globals = {}
44     try:
45         exec(code_str, exec_globals)
46     except Exception as e:
47         return False, {"error": str(e), "trace": traceback.format_exc()}
48     for test in tests:
49         try:
50             exec(test, exec_globals)
51         feedback["number_passed"] += 1
52     except AssertionError:
53         feedback["failures"].append({"test": test, "error": "AssertionError"})
54     except Exception as e:
55         feedback["failures"].append({"test": test, "error": str(e)})
56     return feedback["number_passed"] == feedback["total"], feedback
57
```

```

58 console.log("Planner training task and test runner ready.")
59
60 # ----- Planner action space -----
61 from enum import IntEnum
62
63 class PlannerAction(IntEnum):
64     PLANNING = 0
65     EXECUTION = 1
66     INTEGRATION = 2
67     EVALUATION = 3
68     STOP = 4
69
70 PlannerAction.NAMES = {
71     PlannerAction.PLANNING: "planning",
72     PlannerAction.EXECUTION: "execution",
73     PlannerAction.INTEGRATION: "integration",
74     PlannerAction.EVALUATION: "evaluation",
75     PlannerAction.STOP: "stop",
76 }
77 console.log("Planner Action space:", PlannerAction.NAMES)
78
79 # ----- Simulated Planner LLM (deterministic) -----
80 class SimulatedPlannerLLM:
81     def __init__(self, task: Dict):
82         self.task = task
83
84     def planning(self, prompt: str) -> Dict[str, str]:
85         return {
86             "task_a": "Write function add(a,b) that returns a + b",
87             "task_b": "Write function square(x) that returns x ** 2",
88             "task_c": "Write function cbrt(x) that returns x ** (1/3)",
89         }
90
91     def execute(self, subtask_prompts: Dict[str, str]) -> Dict[str, str]:
92         return {
93             "task_a": "def add(a, b):\n    return a + b\n",
94             "task_b": "def square(x):\n    return x ** 2\n",
95             "task_c": "def cbrt(x):\n    return x ** (1.0/3.0)\n",
96         }
97
98     def integrate(self, exec_outputs: Dict[str, str], original_prompt: str) -> str:
99         integrated = """def compute(a, b):
100     # integrated result computing (a + b) ** (2/3)
101     return (a + b) ** (2.0/3.0)
102 """
103         helpers = """
104     def add(a, b):
105         return a + b
106
107     def square(x):
108         return x ** 2
109
110     def cbrt(x):
111         return x ** (1.0/3.0)
112 """
113         return helpers + "\n" + integrated
114
115     def evaluate(self, integrated_code: str, tests: List[str]) -> Tuple[bool, Dict[str, Any]]:
116         return run_tests_locally(integrated_code, tests)
117
118     sim_planner = SimulatedPlannerLLM(PLANNER_TASK)
119
120 # ----- Planner environment -----

```

```

121 class PlannerEnv:
122     def __init__(self, task: Dict, max_steps: int = 10):
123         self.task = task
124         self.max_steps = max_steps
125         self.reset()
126
127     def reset(self):
128         self.steps = 0
129         self.done = False
130         self.success = False
131         self.has_planned = False
132         self.has_executed = False
133         self.has_integrated = False
134         self.has_evaluated = False
135         self.invalid_action = False
136         self.plan_output = None
137         self.exec_output = None
138         self.integrated_code = ""
139         self.test_feedback = ""
140         self.test_results = []
141         return self._state()
142
143     def _fail(self):
144         return -5.0
145
146     def step(self, action: int):
147         if self.done:
148             raise RuntimeError("Episode already finished")
149         self.steps += 1
150         reward = -0.1
151         info = {"action": PlannerAction.NAMES[PlannerAction(action)]}
152
153         # PLANNING
154         if action == PlannerAction.PLANNING:
155             if self.has_planned:
156                 reward += self._fail()
157                 self.invalid_action = True
158             else:
159                 self.plan_output = sim_planner.planning(self.task["prompt"])
160                 self.has_planned = True
161                 reward += 1.0
162
163         # EXECUTION
164         elif action == PlannerAction.EXECUTION:
165             if (not self.has_planned) or self.has_executed:
166                 reward += self._fail()
167                 self.invalid_action = True
168             else:
169                 self.exec_output = sim_planner.execute(self.plan_output)
170                 self.has_executed = True
171                 reward += 1.5
172
173         # INTEGRATION
174         elif action == PlannerAction.INTEGRATION:
175             if (not self.has_executed) or self.has_integrated:
176                 reward += self._fail()
177                 self.invalid_action = True
178             else:
179                 self.integrated_code = sim_planner.integrate(self.exec_output, self.task["prompt"])
180                 self.has_integrated = True
181                 reward += 2.0
182
183         # EVALUATION

```

```

184 elif action == PlannerAction.EVALUATION:
185     if (not self.has_integrated) or self.has_evaluated:
186         reward += self._fail()
187         self.invalid_action = True
188     else:
189         passed, fb = sim_planner.evaluate(self.integrated_code, self.task["tests"])
190         self.test_results.append(passed)
191         self.test_feedback = f"Tests {fb['number_passed']}/{fb['total']}]"
192         self.has_evaluated = True
193         reward += 5.0 if passed else -1.0
194
195 # STOP
196 elif action == PlannerAction.STOP:
197     self.done = True
198     if self.has_evaluated and self.test_results and self.test_results[-1] and (not self.
199         invalid_action):
200         self.success = True
201         reward += 30.0 if self.success else -10.0
202
203 # step limit
204 if self.steps >= self.max_steps and not self.done:
205     self.done = True
206     reward -= 10.0
207
208 info["success"] = self.success
209 info["invalid_action"] = self.invalid_action
210 return self._state(), reward, self.done, info
211
212 def _state(self):
213     last_code = self.integrated_code if self.integrated_code else ""
214     last_feedback = self.test_feedback
215     return make_state_embedding(
216         self.task["prompt"],
217         last_code,
218         last_feedback,
219         has_planned=self.has_planned,
220         executed_count=1 if self.has_executed else 0,
221         integrated=1 if self.has_integrated else 0,
222         evaluated=1 if self.has_evaluated else 0,
223         invalid_action=self.invalid_action,
224         last_test_passed=bool(self.test_results[-1]) if self.test_results else False,
225     )
226
227 # ----- Embedding helper -----
228 from sentence_transformers import SentenceTransformer
229 embedder = SentenceTransformer("sentence-transformers/all-MiniLM-L6-v2")
230
231 def make_state_embedding(task_prompt: str, last_code: str, last_test_feedback: str,
232     has_planned: bool = False, executed_count: int = 0,
233     integrated: bool = False, evaluated: bool = False,
234     invalid_action: bool = False, last_test_passed: bool = False) -> np.ndarray:
235     task_emb = embedder.encode([task_prompt], show_progress_bar=False)
236     code_emb = embedder.encode([last_code or ""], show_progress_bar=False)
237     feedback_emb = embedder.encode([last_test_feedback or ""], show_progress_bar=False)
238     state_vec = np.concatenate([task_emb[0], code_emb[0], feedback_emb[0]])
239     flags = np.array([
240         1.0 if has_planned else 0.0,
241         float(executed_count),
242         1.0 if integrated else 0.0,
243         1.0 if evaluated else 0.0,
244         1.0 if invalid_action else 0.0,
245         1.0 if last_test_passed else 0.0,
246     ], dtype=np.float32)

```

```

246     return np.concatenate([state_vec, flags]).astype(np.float32)
247
248 sample_state = make_state_embedding(PLANNER_TASK["prompt"], "", "", False, 0, False, False,
249                                     False, False)
250 STATE_DIM = sample_state.shape[0]
251 console.log(f"Planner state dim: {STATE_DIM}")
252
253 # ----- Dueling DQN + Replay Buffer -----
254 class DuelingQNetwork(nn.Module):
255     def __init__(self, input_dim, hidden=[512, 256], output_dim=len(PlannerAction.NAMES)):
256         super().__init__()
257         layers = []
258         prev = input_dim
259         for h in hidden:
260             layers.append(nn.Linear(prev, h))
261             layers.append(nn.ReLU())
262             prev = h
263         self.trunk = nn.Sequential(*layers)
264         self.value_head = nn.Sequential(
265             nn.Linear(prev, prev//2 if prev//2>0 else 32),
266             nn.ReLU(),
267             nn.Linear(prev//2 if prev//2>0 else 32, 1)
268         )
269         self.adv_head = nn.Sequential(
270             nn.Linear(prev, prev//2 if prev//2>0 else 32),
271             nn.ReLU(),
272             nn.Linear(prev//2 if prev//2>0 else 32, output_dim)
273         )
274
275     def forward(self, x):
276         x = self.trunk(x)
277         value = self.value_head(x)
278         adv = self.adv_head(x)
279         q = value + (adv - adv.mean(dim=1, keepdim=True))
280         return q
281
282     class ReplayBuffer:
283         def __init__(self, capacity=8000):
284             self.buffer = deque(maxlen=capacity)
285         def push(self, s,a,r,ns,done):
286             self.buffer.append((s,a,r,ns,done))
287         def sample(self, batch_size):
288             batch = random.sample(self.buffer, min(batch_size, len(self.buffer)))
289             s,a,r,ns,d = zip(*batch)
290             return (np.stack(s), np.array(a), np.array(r, dtype=np.float32), np.stack(ns), np.array(d,
291                                         dtype=np.float32))
292         def __len__(self):
293             return len(self.buffer)
294
295     class DQNAgent:
296         def __init__(self, state_dim, action_dim=len(PlannerAction.NAMES), hidden=[512, 256], lr=1e-4,
297                      gamma=0.99,
298                      buffer_size=8000, batch_size=64, target_update=500, device=None, mask_actions=True,
299                      eps_decay=0.9995):
300             self.device = device or ("cuda" if torch.cuda.is_available() else "cpu")
301             self.qnet = DuelingQNetwork(state_dim, hidden, action_dim).to(self.device)
302             self.target = DuelingQNetwork(state_dim, hidden, action_dim).to(self.device)
303             self.target.load_state_dict(self.qnet.state_dict())
304             self.opt = Adam(self.qnet.parameters(), lr=lr)
305             self.gamma = gamma
306             self.buffer = ReplayBuffer(capacity=buffer_size)
307             self.batch_size = batch_size
308             self.action_dim = action_dim

```

```

306 self.eps = 1.0
307 self.eps_min = 0.05
308 self.eps_decay = eps_decay
309 self.learn_steps = 0
310 self.target_update = target_update
311 self.mask_actions = mask_actions
312
313 def _compute_action_mask(self, state: np.ndarray):
314     flags = state[-6:]
315     has_planned = bool(flags[0])
316     executed_count = int(flags[1])
317     integrated = bool(flags[2])
318     evaluated = bool(flags[3])
319     invalid_action = bool(flags[4])
320     last_test_passed = bool(flags[5])
321
322     mask = np.ones(self.action_dim, dtype=bool)
323     if has_planned:
324         mask[PlannerAction.PLANNING] = False
325     if (not has_planned) or (executed_count >= 1):
326         mask[PlannerAction.EXECUTION] = False
327     if (not executed_count) or integrated:
328         mask[PlannerAction.INTEGRATION] = False
329     if (not integrated) or evaluated:
330         mask[PlannerAction.EVALUATION] = False
331     if not (last_test_passed and (not invalid_action)):
332         mask[PlannerAction.STOP] = False
333     return mask
334
335 def select_action(self, state):
# EXPLORATION: sample uniformly from full action space (no mask) -> allows discovering failures
336     if random.random() < self.eps:
337         return random.randrange(self.action_dim)
338
# EXPLOIT: use greedy Q with masking to avoid selecting impossible terminal actions
339     s = torch.FloatTensor(state).unsqueeze(0).to(self.device)
340     with torch.no_grad():
341         qvals = self.qnet(s).cpu().numpy().squeeze(0)
342         if self.mask_actions:
343             mask = self._compute_action_mask(state)
344             legal_q = np.where(mask, qvals, -1e9)
345             if legal_q.max() <= -1e8:
346                 return random.randrange(self.action_dim)
347             return int(int(legal_q.argmax()))
348         else:
349             return int(int(qvals.argmax()))
350
351     def push_transition(self, s,a,r,ns,done):
352         self.buffer.push(s,a,r,ns,done)
353
354     def train_step(self):
355         if len(self.buffer) < 32:
356             return 0.0
357         s,a,r,ns,d = self.buffer.sample(self.batch_size)
358         s = torch.FloatTensor(s).to(self.device)
359         a = torch.LongTensor(a).to(self.device)
360         r = torch.FloatTensor(r).to(self.device)
361         ns = torch.FloatTensor(ns).to(self.device)
362         d = torch.FloatTensor(d).to(self.device)
363
364         qvals = self.qnet(s).gather(1, a.unsqueeze(1)).squeeze(1)
365         with torch.no_grad():
366             next_actions = self.qnet(ns).argmax(dim=1, keepdim=True)

```

```

369 next_q_target = self.target(ns).gather(1, next_actions).squeeze(1)
370 target = r + self.gamma * next_q_target * (1 - d)
371
372 loss = F.mse_loss(qvals, target)
373 self.opt.zero_grad()
374 loss.backward()
375 nn.utils.clip_grad_norm_(self.qnet.parameters(), 0.5)
376 self.opt.step()
377
378 self.learn_steps += 1
379 if self.learn_steps % self.target_update == 0:
380     self.target.load_state_dict(self.qnet.state_dict())
381
382 self.eps = max(self.eps_min, self.eps * self.eps_decay)
383 return float(loss.item())
384
385 console.log("DQN agent for planner ready.")
386
387 # ----- Training loop -----
388 def train_planner(agent: DQNAgent, task: Dict, num_episodes=1000, max_steps=10, log_every=50,
389 seed_expert_episodes=0, save_after=100):
390     metrics = {"episode_rewards": [], "losses": [], "episode_steps": [], "success_rate_window": [],
391     "sample_eff": []}
392     success_history = deque(maxlen=100)
393     total_env_steps = 0
394
395 ACTION_NAME_TO_INT = {v:k for k,v in PlannerAction.NAMES.items()}
396 expert_seq = ["planning", "execution", "integration", "evaluation", "stop"]
397
398 for _ in range(seed_expert_episodes):
399     env = PlannerEnv(task, max_steps=max_steps)
400     s = env.reset()
401     for action_name in expert_seq:
402         a = ACTION_NAME_TO_INT[action_name]
403         ns, r, done, info = env.step(a)
404         agent.push_transition(s, a, r, ns, float(done))
405         s = ns
406         if done: break
407
408     best_success = 0.0
409     best_avg_reward = -1e9
410
411     for ep in tqdm(range(1, num_episodes+1), desc="Planner DQN Training"):
412         env = PlannerEnv(task, max_steps=max_steps)
413         state = env.reset()
414         ep_reward = 0.0
415         ep_loss = 0.0
416
417         for t in range(max_steps):
418             action = agent.select_action(state)
419             next_state, reward, done, info = env.step(action)
420             agent.push_transition(state, action, reward, next_state, float(done))
421             loss = agent.train_step()
422             state = next_state
423             ep_reward += reward
424             ep_loss += loss
425             total_env_steps += 1
426             if done:
427                 break
428
429         metrics["episode_rewards"].append(ep_reward)
430         metrics["losses"].append(ep_loss / (t+1) if (t+1)>0 else 0.0)
431         metrics["episode_steps"].append(t+1)

```

```

431 success_history.append(1.0 if done and info.get("success", False) else 0.0)
432 metrics["success_rate_window"].append(np.mean(success_history))
433 metrics["sample_eff"].append(sum(metrics["episode_rewards"]) / (total_env_steps + 1e-8))
434
435 current_success = metrics["success_rate_window"][-1]
436 current_avg_reward = np.mean(metrics["episode_rewards"][-log_every:])
437
438 if ep > save_after and (current_success > best_success or (current_success == best_success and
439     current_avg_reward > best_avg_reward)):
440     best_success = current_success
441     best_avg_reward = current_avg_reward
442     torch.save(agent.qnet.state_dict(), "saved_models/planner_dqn_best.pt")
443     console.log(f"[green]Ep {ep:4d} | New best model saved: success(last100)={best_success:.3%},
444                 avg_reward={best_avg_reward:.3f}[/green]")
445
446 if ep % log_every == 0 or ep == 1:
447     console.log(f"[blue]Ep {ep:4d} | AvgReward {np.mean(metrics['episode_rewards'])[-log_every:]:.3
448 f} | "
449 f"Success(last100) {metrics['success_rate_window'][-1]:.3%} | Eps {agent.eps:.3f}")
450
451 return metrics
452
453 # ----- Run training -----
454 if __name__ == '__main__':
455     agent = DQNAgent(state_dim=STATE_DIM, action_dim=len(PlannerAction.NAMES), hidden=[512, 256], lr
456         =3e-4,
457         buffer_size=8000, batch_size=64, target_update=200, mask_actions=True, eps_decay=0.9995)
458
459     metrics = train_planner(agent, PLANNER_TASK, num_episodes=3000, max_steps=10, log_every=50,
460         seed_expert_episodes=0, save_after=100)
461
462     torch.save(agent.qnet.state_dict(), "saved_models/planner_dqn_final.pt")
463     console.log("Planner training complete. Models saved to saved_models/planner_dqn_*.pt")
464
465 for i in range(5):
466     env = PlannerEnv(PLANNER_TASK, max_steps=10)
467     s = env.reset()
468     path = []
469     for _ in range(12):
470         a = agent.select_action(s)
471         path.append(PlannerAction.NAMES[PlannerAction(a)])
472         s, r, done, info = env.step(a)
473         if done: break
474     console.log(f"Demo {i+1}: {path} | success={info.get('success', False)}")
```

Explanation: The planner's role is to produce subtask prompts (planning), call the executive to execute each subtask (execution), combine helpers (integration), and validate the integrated code (evaluation). The training loop follows similar patterns: replay buffer, double-dueling q-networks, target updates, epsilon decay, and occasional saving of best models.

Include an example planner training screenshot and demo results in the report:

```

[23:04:00] Ep 1002 | New best model saved: success(last100)=87.000%, avg_reward=35.484
...
[23:04:04] Ep 1003 | New best model saved: success(last100)=88.000%, avg_reward=35.484
[23:04:07] Ep 1009 | New best model saved: success(last100)=89.000%, avg_reward=33.516
[23:04:08] Ep 1012 | New best model saved: success(last100)=89.000%, avg_reward=33.826
[23:04:08] Ep 1013 | New best model saved: success(last100)=90.000%, avg_reward=35.028
[23:04:08] Ep 1014 | New best model saved: success(last100)=91.000%, avg_reward=35.028
[23:04:36] Ep 1050 | AvgReward 26.374 | Success(last100) 85.000% | Eps 0.050
[23:05:12] Ep 1100 | AvgReward 25.924 | Success(last100) 77.000% | Eps 0.050
[23:05:47] Ep 1150 | AvgReward 29.900 | Success(last100) 80.000% | Eps 0.050
[23:06:24] Ep 1200 | AvgReward 29.144 | Success(last100) 84.000% | Eps 0.050
[23:07:01] Ep 1250 | AvgReward 27.252 | Success(last100) 82.000% | Eps 0.050
[23:07:36] Ep 1300 | AvgReward 32.388 | Success(last100) 84.000% | Eps 0.050
[23:08:11] Ep 1350 | AvgReward 30.026 | Success(last100) 86.000% | Eps 0.050
[23:08:47] Ep 1400 | AvgReward 28.546 | Success(last100) 83.000% | Eps 0.050
[23:09:22] Ep 1450 | AvgReward 28.208 | Success(last100) 81.000% | Eps 0.050
[23:09:56] Ep 1500 | AvgReward 31.154 | Success(last100) 82.000% | Eps 0.050
[23:10:31] Ep 1550 | AvgReward 33.334 | Success(last100) 87.000% | Eps 0.050
[23:11:08] Ep 1600 | AvgReward 25.556 | Success(last100) 83.000% | Eps 0.050
[23:11:44] Ep 1650 | AvgReward 29.256 | Success(last100) 79.000% | Eps 0.050
[23:12:21] Ep 1700 | AvgReward 24.054 | Success(last100) 78.000% | Eps 0.050
[23:12:57] Ep 1750 | AvgReward 27.204 | Success(last100) 76.000% | Eps 0.050
[23:13:33] Ep 1800 | AvgReward 30.718 | Success(last100) 82.000% | Eps 0.050
[23:14:08] Ep 1850 | AvgReward 30.116 | Success(last100) 85.000% | Eps 0.050
[23:14:44] Ep 1900 | AvgReward 30.668 | Success(last100) 84.000% | Eps 0.050
[23:15:18] Ep 1950 | AvgReward 32.436 | Success(last100) 86.000% | Eps 0.050
[23:15:55] Ep 2000 | AvgReward 22.190 | Success(last100) 79.000% | Eps 0.050
[23:16:30] Ep 2050 | AvgReward 32.896 | Success(last100) 79.000% | Eps 0.050
[23:17:06] Ep 2100 | AvgReward 33.166 | Success(last100) 89.000% | Eps 0.050
[23:17:41] Ep 2150 | AvgReward 24.040 | Success(last100) 82.000% | Eps 0.050
[23:18:16] Ep 2200 | AvgReward 32.816 | Success(last100) 81.000% | Eps 0.050

```

Figure 3: Planner training success-rate plot / console log.

Reported best planner run excerpt (from your logs):

```
Ep 1014 | New best model saved: success(last100)=91.000%, avg_reward=35.028
```

```

... [13:17:14] Loaded best model from saved_models/dqn_best.pt
Final model saved to saved_models/dqn_final.pt
Demo rollouts for task1_correct
[13:17:15] Path 1: ['generate', 'test', 'stop'] | success=True
Path 2: ['generate', 'test', 'stop'] | success=True
Path 3: ['generate', 'test', 'stop'] | success=True
Demo rollouts for task2_incorrect
[13:17:16] Path 1: ['generate', 'test', 'debug', 'test', 'stop'] | success=True
Path 2: ['generate', 'test', 'debug', 'test', 'stop'] | success=True
Path 3: ['generate', 'test', 'debug', 'test', 'stop'] | success=True
Demo rollouts for task3_search
[13:17:17] Path 1: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
Path 2: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
Path 3: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
[13:17:18] Path 1: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
Path 2: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
Path 3: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'search', 'stop'] | success=True
[13:17:19] Path 1: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'test', 'stop'] | success=True
Path 2: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'test', 'stop'] | success=True
Path 3: ['generate', 'test', 'debug', 'test', 'search', 'debug',
'test', 'stop'] | success=True

```

Figure 4: Planner demo rollouts showing successful canonical episodes.

9 Hierarchical Inference: Combining Planner + Executive

After training both models, the full hierarchical inference works as follows:

1. Load the planner and executive networks (the trained DQN weights).
2. The planner receives the top-level prompt (for example, `compute (a+b) to the power 2/3`), and then executes its policy, which proceeds through the following sub-steps:
 - (a) **Planning:** the planner produces subtask prompts (for example, `add`, `square`, `cbrt`).
 - (b) **Execution:** for each subtask the Executive is invoked. The Executive follows the local protocol (for example, `generate`, `test`, optionally `debug`, then `stop`) to produce a helper function and verifies it against unit tests.
 - (c) **Integration:** the planner combines the produced helper functions into a single integrated code body (for example, the final `compute()` implementation).
 - (d) **Evaluation:** the integrated code is executed and tested (either by the Executive or by a direct test runner). Based on the evaluation outcome the planner decides whether to `stop` or to attempt further corrective actions such as re-planning or re-execution.

3. If the integrated evaluation fails, the planner can in principle learn to re-plan or to trigger additional corrective flows. This extension was not explored in depth in the current experiments.

9.1 Pretty-logged Inference and Final Output

You provided a script that loads networks and pretty-prints the hierarchical inference flow. The script prints only the final `compute` function to be delivered to the user. The provided run demonstrates success: subtasks (add, square, cbrt) were produced, tested, integrated, and the final `compute` passed planner tests.

Final compute function printed by the script:

Listing 10: Final `compute()` output printed to the user

```

1
2 =====
3 [18:29:08] Hierarchical inference start
4 =====
5 [18:29:08] Query: write me a code for computing the third root of (a+b) ^ 2
6
7 -----
8 [18:29:08] Executive: run subtasks (simulated LLM actions)
9 -----
10
11 [18:29:08] Planner I will ask the Executive to solve subtask:
12 compute_add_001
13 [18:29:08] Executive generate:
14 Producing candidate code...
15
16 def add(a: int, b: int) -> int:
17     return a + b
18 [18:29:08] Executive test:
19 Running unit tests for the generated code
20 [18:29:08] Test results: 3/3 passed
21 [18:29:08] Executive stop:
22 Episode finished for subtask
23
24 [18:29:08] Planner I will ask the Executive to solve subtask:
25 square_001
26 [18:29:08] Executive generate:
27 Producing candidate code...
28
29 def square(x: int) -> int:
30     return x ** 2
31 [18:29:08] Executive test:
32 Running unit tests for the generated code
33 [18:29:08] Test results: 3/3 passed
34 [18:29:08] Executive stop:
35 Episode finished for subtask
36
37 [18:29:08] Planner I will ask the Executive to solve subtask:
38 cbrt_001
39 [18:29:08] Executive generate:
40 Producing candidate code...
41
42 def cbrt(x: int) -> float:
43     if x == 0:
44         return 0
45     sign = -1 if x < 0 else 1
46     x_abs = abs(x)
47     r = round(x_abs ** (1.0 / 3.0))
48     if r ** 3 == x_abs:

```

```

49 return sign * r
50 return sign * (x_abs ** (1.0 / 3.0))
51 [18:29:08] Executive test:
52 Running unit tests for the generated code
53 [18:29:08] Test results: 4/4 passed
54 [18:29:08] Executive stop:
55 Episode finished for subtask
56
57 -----
58 [18:29:08] Integration: combining helper functions into final compute()
59 -----
60 [18:29:08] Integrated code assembled. Full content below:
61
62 def add(a: int, b: int) -> int:
63     return a + b
64
65 def square(x: int) -> int:
66     return x ** 2
67
68 def cbrt(x: int) -> float:
69     if x == 0:
70         return 0
71     sign = -1 if x < 0 else 1
72     x_abs = abs(x)
73     r = round(x_abs ** (1.0 / 3.0))
74     if r ** 3 == x_abs:
75         return sign * r
76     return sign * (x_abs ** (1.0 / 3.0))
77
78 def compute(a, b):
79     """Compute the cube root of (a+b)^2 (i.e. (a+b)^(2/3))."""
80     s = add(a, b)
81     s2 = square(s)
82     return cbrt(s2)
83
84
85 -----
86 [18:29:08] Final Evaluation: validate integrated compute()
87 -----
88 [18:29:08] TEST: compute(1,2) -> 2.080083823051904 expected 2.080083823051904 ok=True
89 [18:29:08] TEST: compute(0,0) -> 0 expected 0.0 ok=True
90 [18:29:08] TEST: compute(8,1) -> 4.3267487109222245 expected 4.3267487109222245 ok=True
91 [18:29:08] TEST: compute(-1,8) -> 3.6593057100229713 expected 3.6593057100229713 ok=True
92 [18:29:08] Final tests passed: 4/4
93
94 =====
95 [18:29:08] FINAL CODE (compute function only)
96 =====
97
98 def compute(a, b):
99     """Compute the cube root of (a+b)^2 (i.e. (a+b)^(2/3))."""
100    s = add(a, b)
101    s2 = square(s)
102    return cbrt(s2)

```

10 Interpretation of Results

- The Executive DQN successfully learned the three canonical episode types and achieved perfect or near-perfect success rates in the reported runs (example: $\text{success}(\text{last100})=100\%$ at one checkpoint).
- The Planner DQN learned the single planning episode, producing sensible decomposition (3 subtasks) and integrating the helper functions to pass the final tests (reported $\text{success}(\text{last100})=91\%$ best checkpoint).
- The hierarchical system composed the subtask solutions and produced a correct final function `compute` that passed all planner-level tests in the provided demonstration.

10.1 Why the design works

- **Explicit flags** in the state make temporal preconditions explicit, which is critical for finite-state workflows.
- **Semantic embeddings** for prompt/code/feedback provide content-awareness so the network can generalize across prompts and code.
- **Action masking** prevents the agent from learning trivial but illegal shortcuts and reduces wasted exploration.
- **Dueling + Double DQN** increases stability (separating state value from advantage) and reduces overestimation bias.
- **Seeding with expert episodes** gives early guidance and accelerates learning of the correct sequences.

11 Hyperparameters (selected)

Parameter	Value
Network hidden sizes	[512, 256]
Learning rate	3e-4 (training code uses 3e-4; agent default 1e-4)
Batch size	64
Replay buffer	8000
Target update (steps)	200 (example)
Discount factor (γ)	0.99
Epsilon start / min / decay	1.0 / 0.05 / 0.995 (exec)
Episode max steps	10
Seeding expert episodes	8 per task (exec) / configurable for planner

Table 1: Key hyperparameters used in experiments.