
On the Role of Proposal Support in Diffusion-Based Offline RL for Sequential Decision-Making

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Proposal-selection frameworks for offline reinforcement learning (RL) enable
2 decision-making through candidate generation and value-based selection. However,
3 it remains unclear where performance variations arise within this process, particu-
4 larly across intermediate steps. We study a fundamental limitation of proposal-
5 based decision-making: the extent to which decision quality is constrained by
6 the support of the proposal distribution. Using ARC-AGI as a controlled analysis
7 environment, we construct Synthesized Offline Learning data for Abstraction and
8 Reasoning (SOLAR), a trajectory-based dataset that converts tasks into step-by-step
9 solution trajectories. This enables step-level analysis of decision-making across
10 controlled trajectory distributions. Through experiments with Latent Diffusion-
11 Constrained Q-Learning (LDCQ), we find that decision quality is closely tied to
12 proposal coverage: while the selection mechanism remains reliable when suitable
13 candidates are present, performance degrades sharply when the proposal distri-
14 bution fails to cover solution-relevant behaviors. These results identify proposal
15 coverage as a key bottleneck in decision-making with diffusion-based offline RL,
16 pointing to directions for improving robustness and generalization.

17 1 Introduction

18 Sequential decision-making in complex environments requires constructing solutions through a
19 series of intermediate steps. Recent approaches to offline reinforcement learning (RL) adopt a
20 proposal-selection structure, where candidate behaviors are generated and evaluated to guide decision-
21 making [2, 14]. In diffusion-based methods, this structure is typically implemented using generative
22 models for proposal and value functions for selection.

23 It remains unclear what governs decision quality when performance varies across tasks or data condi-
24 tions. In particular, it is difficult to attribute such variation to the quality of candidate behaviors versus
25 the selection mechanism, as these components are tightly coupled during inference. Disentangling
26 their contributions is essential for diagnosing failure modes and improving the robustness of the
27 proposal-selection-based offline RL.

28 We study this question using ARC-AGI [1] as a controlled environment in which intermediate state
29 transitions can be explicitly verified, enabling step-level analysis beyond final-outcome evaluation.
30 We construct Synthesized Offline Learning data for Abstraction and Reasoning (SOLAR), a trajectory-
31 based dataset that converts tasks into step-by-step solving trajectories, and evaluate Latent Diffusion-
32 Constrained Q-Learning (LDCQ) [14] within this framework.

33 Our experiments show that LDCQ constructs multi-step solutions effectively when the trajectory
34 distribution is well-aligned with the task, but performance degrades when the distribution contains sub-
35 optimal trajectories or when the task requires generalization beyond observed ones. We examine this

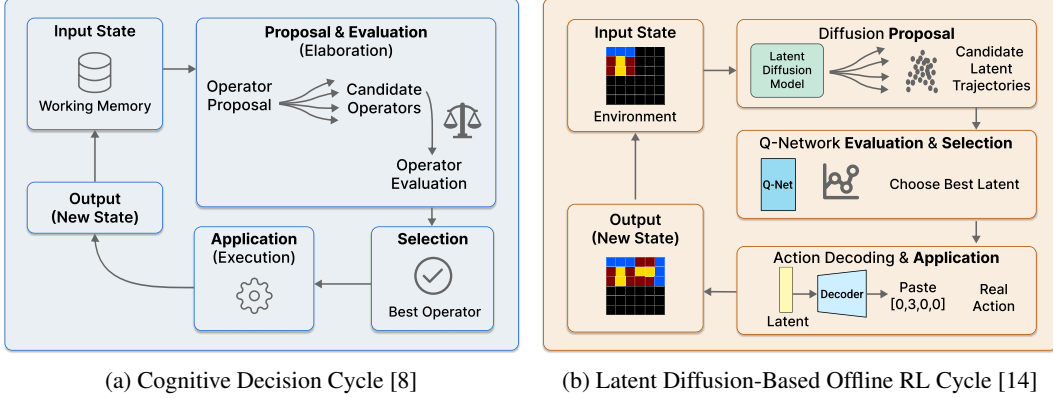


Figure 1: A staged decision motif studied in this paper. Both cycles can be decomposed into the phases of proposal, evaluation, selection, and application (execution).

36 behavior across multiple controlled settings, including single-task expert training, mixed-quality tra-
 37 jectory distributions, and compositional generalization. We find that this degradation stems primarily
 38 from insufficient proposal coverage: when the proposal distribution fails to include solution-relevant
 39 behaviors, the selection mechanism cannot recover regardless of its accuracy. These findings identify
 40 proposal support as the primary bottleneck in diffusion-based offline RL.

41 Our contributions are as follows:

- 42 • We study the role of proposal support in proposal–selection frameworks, isolating the
 43 contribution of candidate generation from value-based selection at the step level.
- 44 • We introduce SOLAR, a trajectory-based dataset that enables systematic, controlled variation
 45 in trajectory quality for fine-grained, step-level analysis of offline RL.
- 46 • We provide empirical evidence that proposal coverage is a key bottleneck in diffusion-based
 47 offline RL, identifying a concrete target for improving robustness.

48 2 Background

49 2.1 Diffusion-based Offline RL

50 LDCQ [14] models multi-step decision-making using a latent proposal–selection pipeline. This
 51 pipeline follows a staged decision cycle consisting of proposal, evaluation, selection, and execution,
 52 as illustrated in Figure 1. Trajectory segments $\tau_t = (s_t, a_t, \dots, s_{t+H-1}, a_{t+H-1})$ are encoded into
 53 latent variables z_t , which represent H -step behaviors.

54 **Latent Representation Learning.** A β -VAE is trained to encode trajectory segments into latent
 55 variables $z_t \sim q_\phi(z_t | \tau_t)$ and decode actions via a policy $\pi_\theta(a | s, z_t)$. The model is trained to learn
 56 compact representations of multi-step behaviors.

57 **Latent Proposal via Diffusion.** A conditional diffusion model $p_\psi(z_t | s_t)$ is trained over latent
 58 variables to generate candidate behaviors given the current state. This mainly serves as the proposal
 59 distribution in the decision pipeline.

60 **Value-based Selection.** A Q-function $Q(s_t, z_t)$ is trained to evaluate latent behaviors over an H -
 61 step horizon. Given a transition from s_t to s_{t+H} , the target value is computed by sampling candidate
 62 next latents from the diffusion prior and selecting the one with the highest Q-value:

$$Q(s_t, z_t) \leftarrow r_{t:t+H} + \gamma^H Q \left(s_{t+H}, \underset{z \sim p_\psi(z | s_{t+H})}{\operatorname{argmax}} Q(s_{t+H}, z) \right), \quad (1)$$

63 At inference time, multiple candidates are sampled from $p_\psi(z | s_t)$ and the highest-valued one is
 64 selected. This enables sequential decision-making through iterative proposal and selection.

65 2.2 ARC-AGI as a Structured Decision Environment

66 ARC-AGI [1] is a benchmark where agents must infer transformation rules from a small number of
67 examples and apply them to unseen inputs. Although it is originally formulated as an input–output
68 mapping problem, solving ARC-AGI tasks often involves applying a sequence of transformations to
69 construct the final solution.

70 We formulate ARC-AGI as a Markov Decision Process, where the state s_t represents the current grid,
71 the operation opr_t specifies the transformation to apply, and the selection sel_t determines where the
72 transformation is applied. The resulting state transition is given by:

$$s_{t+1} = f(s_t, opr_t, sel_t), \quad (2)$$

73 A reward is provided only upon correct submission, resulting in sparse feedback:

$$r_t = \begin{cases} 1 & \text{if the submitted grid matches the target,} \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

74 To enable execution and trajectory-level analysis, we use ARCLE [9], a reinforcement learning
75 environment for ARC-AGI that provides explicit state transitions and a discrete action space. ARCLE
76 allows us to simulate sequential decision-making and validate action sequences within an environment.

77 However, standard ARC-AGI provides only input–output pairs without explicit trajectories, which
78 limits the analysis of intermediate decision processes. This motivates the need for trajectory-level
79 representations to enable step-by-step analysis.

80 3 SOLAR: A Trajectory-Based Dataset for Analysis

81 To enable fine-grained analysis of sequential decision-making in ARC-AGI, we introduce the Synthe-
82 sized Offline Learning data for Abstraction and Reasoning (SOLAR), a trajectory-based dataset that
83 provides explicit state–action transitions for ARC-AGI tasks. SOLAR augments standard ARC-AGI
84 formulations by incorporating executable trajectories, allowing us to observe intermediate decision
85 processes and analyze how solutions are constructed step by step.

86 This enables systematic investigation of decision-making behavior, including how candidate behaviors
87 are generated, how actions are selected, and how errors emerge across sequential steps. Furthermore,
88 it allows us to analyze not only final outcomes but also the intermediate processes that lead to success
89 or failure. By decoupling task structure from trajectory composition, SOLAR provides a controlled
90 setting for studying how different trajectory distributions affect decision-making behavior.

91 3.1 Dataset Construction

92 Each SOLAR task consists of demonstration input–output pairs and a test example, for which an
93 executable episode is provided. The demonstrations define the transformation rule, while the episode
94 specifies step-by-step transitions for solving the test instance.

95 An episode is represented as $\tau = \{(s_t, a_t, r_t, s_{t+1})\}_{t=0}^T$, where s_t is the grid state, a_t is an action,
96 and r_t is a sparse reward provided only upon correct submission. All episodes are executed and
97 validated using ARCLE [9], ensuring that each transition is consistent with a well-defined action
98 space and environment dynamics. For compatibility with diverse types of methods, SOLAR supports
99 both whole-episode and fixed-length segmented formats, accommodating a variety of offline RL
100 methods with temporal abstraction requirements.

101 3.2 Trajectory Generation

102 SOLAR is constructed using a rule-based synthesis procedure implemented in SOLAR-Generator.
103 For each task, a solving action sequence is first synthesized to correctly transform the input into the
104 target output, and then validated in ARCLE.

105 To increase diversity, additional trajectories are generated by varying action orderings or using
106 alternative combinations of operations. The generation process consists of three steps: (i) constructing

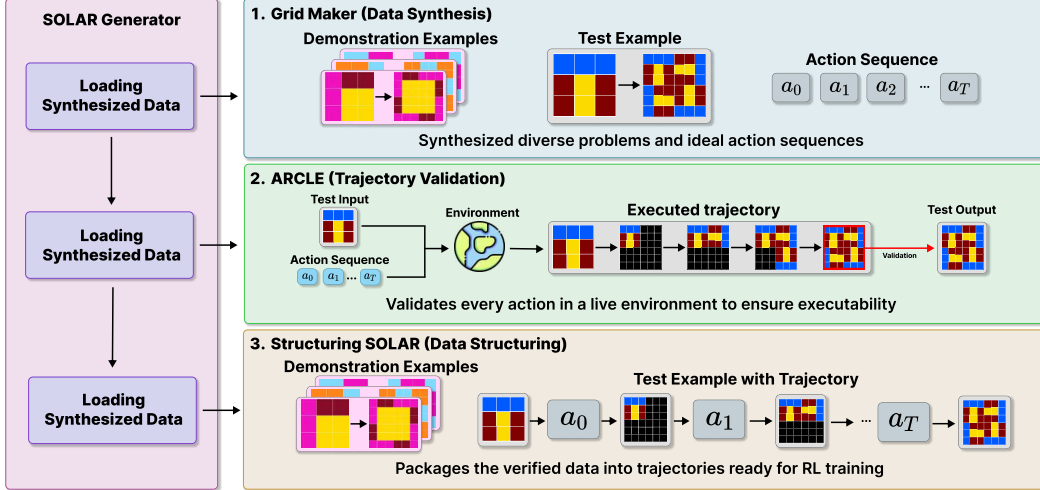


Figure 2: Dataset construction process with SOLAR-Generator. (i) Input-output pairs and action sequences are synthesized, (ii) validated through execution in ARCLE, and (iii) converted into JSON-style episodes for offline RL.

107 candidate solutions under task constraints, (ii) validating them through execution in ARCLE, and (iii)
 108 converting them into transition tuples for offline RL. This results in multiple valid trajectories per
 109 task, ranging from direct solutions to longer or suboptimal paths, exposing diverse decision patterns.

110 Importantly, SOLAR supports flexible construction of trajectory distributions by combining solution-
 111 oriented and suboptimal trajectories. This allows controlled analysis of how decision-making behavior
 112 changes under different trajectory conditions while keeping the underlying task fixed. Additional
 113 implementation details are provided in Appendix B.

114 4 Experiments

115 4.1 Experimental Setup

116 All experiments are conducted using the SOLAR dataset, which provides executable state-action
 117 trajectories for ARC-AGI tasks. The agent interacts with the environment by applying actions
 118 sequentially, receiving a reward only upon correct submission.

119 To control the trajectory distribution, we construct offline datasets by mixing expert trajectories \mathcal{D}_{exp}
 120 and suboptimal trajectories \mathcal{D}_{sub} :

$$\mathcal{D}(\alpha) = \alpha \mathcal{D}_{\text{exp}} + (1 - \alpha) \mathcal{D}_{\text{sub}} \quad (4)$$

121 where $\alpha \in [0, 1]$ controls the proportion of expert trajectories. Suboptimal trajectories may either
 122 reach the correct solution via a longer path or fail to reach it entirely. We consider $\alpha = 1$ (expert-only)
 123 and $\alpha = 1/2$ (mixed-quality) settings.

124 For each task, we synthesize 5,000 training episodes and 100 test episodes using SOLAR-Generator,
 125 with non-overlapping test examples. We report task success rate, defined as the proportion of episodes
 126 that terminate with a correct final output.

127 4.2 Behavior under Expert Trajectories

128 We evaluate LDCQ in a controlled, single-task setting in which the trajectory distribution is fully
 129 aligned with the task, using expert-only trajectories. In this setting, all training episodes success-
 130 fully reach the correct solution, allowing us to isolate the agent’s ability to learn valid multi-step
 131 transformations.

132 Table 1 summarizes performance across representative ARC-AGI tasks. While LDCQ achieves
 133 near-perfect accuracy on several tasks, performance varies depending on task structure. Notably,
 134 some tasks such as 0d3d703e exhibit near-zero performance despite expert-only training.

135 These tasks require inferring global transformation rules (e.g., color permutations) that are not
 136 explicitly reflected in intermediate state transitions, making them difficult to capture through step-
 137 level trajectory representations. In contrast, tasks involving spatial transformations (e.g., coloring
 138 specific regions) yield more direct state changes that are easier to model.

139 Overall, LDCQ performs reliably when solutions can be decomposed into explicit sequences of
 140 intermediate transformations, confirming that the model can learn effective multi-step behaviors when
 141 the trajectory distribution provides sufficient coverage.

Table 1: Single-task performance with expert-only trajectories.

Task ID	Accuracy	Task ID	Accuracy
5c0a986e-colorfix	100%	a65b410d	100%
4c4377d9	98%	6f8cd79b	91%
4258a5f9	79%	007bbfb7	77%
46442a0e	76%	5c0a986e-colordiff	68%
74dd1130	24%	0d3d703e	0%

142 4.3 Selecting Efficient Solution Paths

143 We examine whether LDCQ can identify more efficient behaviors when multiple solutions exist. To
 144 study this, we construct a controlled SOLAR dataset for task a65b410d, where all trajectories reach
 145 the correct solution but differ in efficiency. As illustrated in Figure 3, the two training trajectories
 146 diverge before a shared intermediate state and reconverge after it: one takes an optimal path in the
 147 first half but an inefficient path in the second, and the other does the reverse. As a result, no single
 148 training trajectory is fully optimal end-to-end, and a fully efficient solution can only be constructed
 by combining efficient segments from different trajectories.

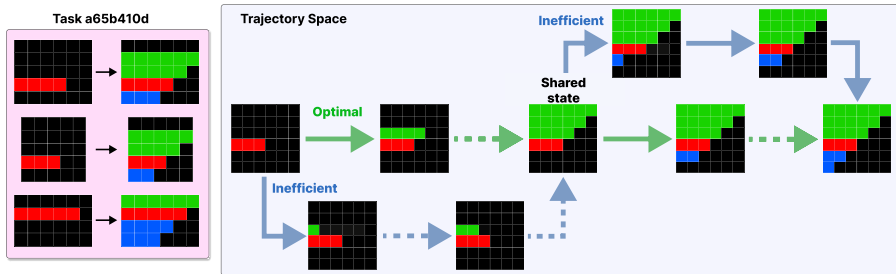


Figure 3: Two ways of solving Task a65b410d. The two trajectories share an intermediate state, with each being optimal in one half and inefficient in the other. A fully efficient solution requires stitching the optimal segments across trajectories, which no single training episode achieves end-to-end.

149

150 We compare LDCQ against two ablated baselines that omit selection: VAE Prior and Diffusion Prior,
 151 which use a single candidate sampled from the VAE and from the diffusion model, respectively. As
 152 shown in Table 2, all methods achieve high accuracy, but LDCQ achieves with fewer steps.

Table 2: Comparison of policies on a task where all trajectories reach correct solutions but differ in efficiency. LDCQ significantly reduces the number of steps, highlighting the role of selection.

Policy	Accuracy	Avg. Steps
VAE Prior	90%	19.4
Diffusion Prior	99%	15.4
LDCQ (w/ Q)	100%	7.0

153 While generative policies (VAE prior, Diffusion prior) tend to follow training trajectories and produce
 154 longer solutions, LDCQ achieves the same accuracy with significantly fewer steps. This suggests that
 155 value-based selection enables the model to prefer more efficient action sequences rather than simply
 156 reproducing observed trajectories.

157 **4.4 Behavior under Mixed-quality Trajectories**

158 We next consider a multi-task setting where the dataset contains a mixture of expert and suboptimal
 159 trajectories. We vary α to control the proportion of expert trajectories and evaluate how performance
 160 changes across methods.

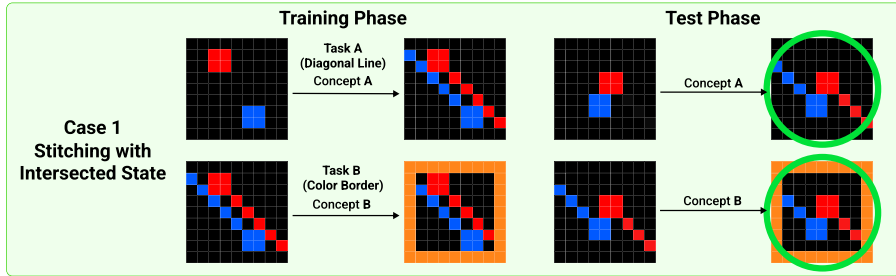
Table 3: Baseline comparison on multi-task learning under expert-only ($\alpha = 1$) and mixed-quality ($\alpha = 1/2$) training.

Method	Expert-only	Mixed-quality
BC	38.8±1.86%	9.67±2.44%
CQL [7]	61.9±0.31%	37.5±0.31%
IQL [6]	58.3±0.31%	32.4±1.11%
LDCQ [14]	66.9±6.39%	23.7±6.04%

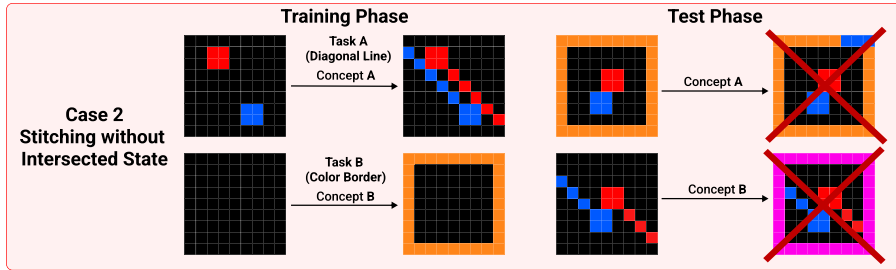
161 As shown in Table 3, LDCQ achieves the highest performance under expert-only data, but its
 162 performance drop under mixed-quality trajectories is notably larger than that of other methods. This
 163 higher sensitivity to distributional variation may be attributed to the nature of latent-based proposal
 164 generation: unlike CQL and IQL, which apply explicit action-level constraints or regularization to
 165 stabilize learning, LDCQ relies on a diffusion prior over latent trajectory segments, making it more
 166 susceptible to degradation when the training distribution contains suboptimal behaviors. Qualitatively,
 167 the agent often reaches states close to the target but fails to complete the task, producing inconsistent
 168 action sequences.

169 These results suggest that decision quality is closely tied to the availability of suitable candidate
 170 behaviors in the proposal stage, which we analyze in detail in Section 5.

171 **4.5 Compositional Generalization**



(a) Applying learned concepts to seen intermediate states (Case 1)



(b) Applying learned concepts to unseen intermediate states (Case 2)

Figure 4: Compositional generalization with and without intermediate state support. (a) When the intermediate state at the composition boundary is covered by training data, the agent successfully composes the two learned transformations. (b) When it is absent, composition fails despite each transformation having been learned individually, highlighting that generalization is constrained by trajectory-level support rather than concept-level learning.

172 We examine whether LDCQ can generalize to compositional tasks that require combining independently
 173 learned transformations. The agent is trained separately on two individual transformations,

Table 4: Effect of trajectory support on compositional generalization.

Condition	State Seen	Transition Seen	Accuracy
Seen task	Yes	Yes	100%
Unseen state	No	Partial	0–33%
Composition (with shared states)	Yes	Yes	100%
Composition (without shared states)	No	No	0%

174 Task A (diagonal line) and Task B (color border), and evaluated on a compositional task that requires
 175 executing them sequentially. We consider two cases depending on whether the training episodes share
 176 an intermediate state at the composition boundary: in Case 1, the intermediate state after applying
 177 Task A appears in the training data for Task B, while in Case 2 it does not. As a baseline, we also
 178 evaluate each transformation individually on unseen intermediate states (Case 2 without composition),
 179 finding that performance already degrades (0–33%) even before attempting composition, confirming
 180 that the agent struggles to generalize to states outside the training distribution.

181 As shown in Figure 4 and Table 4, composition succeeds in Case 1 where the intermediate state is
 182 covered by training data, but fails in Case 2 where it is not, even though each transformation has
 183 been learned individually. When transitions are only partially supported, performance varies across
 184 concepts, with Concept A failing completely (0%) and Concept B achieving limited success (33%).

185 These results indicate that compositional generalization is fundamentally constrained by the availabil-
 186 ity of compatible trajectory segments in the proposal distribution: even when individual transforma-
 187 tions are learned, the agent cannot compose them if the necessary intermediate states and transitions
 188 are absent from the training data.

189 5 Analysis

190 In this section, we analyze the failure modes observed across our experiments to identify the underly-
 191 ing factors that govern the success and limitations of proposal-based decision-making.

192 5.1 Failure Decomposition

193 As observed in Section 4.4, LDCQ exhibits significant performance degradation under mixed-quality
 194 data. To better understand this behavior, we analyze how failures occur during the decision process.
 195 We observe that the agent often reaches states that are close to the target but fails to consistently
 196 complete the task. In many cases, early-stage transitions are reasonable, but errors accumulate in later
 197 steps, leading to incorrect final outputs.

198 Moreover, generated episodes frequently contain locally valid transitions but fail to form coherent
 199 sequences toward the goal. This suggests that failures are not caused by isolated incorrect actions, but
 200 by inconsistencies across sequential decisions. These observations motivate a separate analysis of the
 201 proposal and selection stages, which we examine next.

202 5.2 Diagnosing the Source of Failure

203 To identify the source of failure, we analyze the proposal and selection stages separately. LDCQ
 204 follows a proposal–selection pipeline, where candidate behaviors are generated from a diffusion prior
 205 and evaluated using a learned Q-function.

206 We first examine whether the Q-network fails to correctly select among candidate behaviors. To isolate
 207 this effect, we perform an oracle proposal analysis, where the correct action is explicitly included
 208 in the candidate set. Since SOLAR provides ground-truth trajectories for test instances, we encode
 209 the corresponding trajectory segments using the trained β -VAE to obtain the latent representation
 210 of the correct behavior, and inject this latent into the candidate set during inference. For each of
 211 500 randomly sampled test states, we generate 100 latent candidates from the diffusion prior, check
 212 whether the oracle latent is included, and evaluate whether the Q-network selects it when present.

Table 5: Proposal coverage and conditional selection accuracy.

	Expert-only	Mixed-quality
$P(\text{oracle} \in \text{proposals})$	91.0%	62.2%
$P(\text{correct action} \mid \text{oracle})$	95.6%	88.4%

213 As shown in Table 5, the probability that the correct candidate is included in the proposal set drops
 214 notably under mixed-quality data (91.0% \rightarrow 62.2%). However, when the correct candidate is available,
 215 the Q-network selects it with high accuracy (88.4%), even under mixed-quality data.

216 These results indicate that the selection mechanism remains reliable when suitable candidates are
 217 present. However, this reliability is conditioned on the quality of the proposal distribution, and the
 218 primary source of failure is the limited coverage of the proposal distribution rather than inaccuracies
 219 in value-based selection.

220 5.3 Proposal Support Limitation

221 We analyze how often correct candidates appear in the proposal set across states. As shown in Table 6,
 222 under mixed-quality data, a large proportion of states (37.8%) contain no correct candidates at all.
 223 This indicates that the correct behavior is entirely absent from the proposal set, forcing the agent to
 224 select among suboptimal options regardless of the accuracy of the Q-function.

Table 6: Distribution of oracle candidates across states (100 proposals), showing that many states lack any correct candidate under mixed-quality data.

# Oracle proposals	Expert-only	Mixed-quality
0	9.0%	37.8%
1–20	1.4%	37.2%
20–50	2.4%	8.4%
50–99	10.8%	6.8%
100	76.4%	9.8%

225 Increasing the number of proposals improves the chance of including correct candidates, but also
 226 introduces more out-of-distribution behaviors. As shown in Table 7, performance improves as
 227 the number of proposals increases, but degrades beyond a certain point. These out-of-distribution
 228 candidates are not well calibrated by the Q-network, causing selection to become biased toward
 229 unreliable candidates and leading to degraded performance despite increased proposal diversity.

Table 7: Effect of the number of proposals on performance ($\alpha = 1/2$).

# Proposals	Accuracy (%)
1	5.6
10	18.2
50	26.2
100	25.4
500	20.0
1000	19.8

230 This indicates that proposal limitations arise not only from insufficient coverage, but also from the
 231 quality of the proposal distribution. When correct candidates are absent, the agent is forced to select
 232 among suboptimal options; when too many out-of-distribution candidates are present, selection itself
 233 becomes less reliable.

234 Together, these results reveal a structural limitation of the LDCQ framework: decision quality is
 235 fundamentally constrained by the support and quality of the proposal distribution. Even a well-
 236 trained value function cannot recover correct behavior when suitable candidates are absent or poorly
 237 represented.

238 6 Discussion: Towards Better Proposal Support

239 Our analysis reveals that decision quality in diffusion-based offline RL is fundamentally constrained
240 by the support of the proposal distribution. This manifests in three failure modes: insufficient coverage
241 under mixed-quality data, degraded selection under out-of-distribution proposals, and failure of
242 compositional generalization. We discuss each and outline directions for addressing these limitations.

243 **Improving proposal coverage under data limitations.** A substantial fraction of states lack correct
244 candidates under mixed-quality data, as the diffusion prior $p_\psi(z_t | s_t)$ is trained to replicate the
245 offline dataset, which may itself be incomplete or noisy. This suggests that proposal diversity is
246 fundamentally constrained by data support. One direction is to decouple proposal generation from
247 raw data quality, for example, through latent-space augmentation, retrieval-augmented proposals, or
248 filtering of suboptimal trajectories. More broadly, this highlights the need for mechanisms that can
249 enrich the proposal distribution beyond the empirical dataset.

250 **Calibrating selection and decoding under unreliable proposals.** Increasing the number of
251 proposals introduces out-of-distribution latents that are not reliably evaluated by the Q-network,
252 leading to overestimation bias. This reflects a fundamental tension between proposal coverage and
253 selection reliability. More broadly, the current formulation tightly couples latent representations with
254 action decoding, meaning that even directionally correct latents may fail to produce valid action
255 sequences during execution. This suggests that latent representations are better interpreted as high-
256 level directional guidance rather than precise action specifications, and should be coupled with more
257 expressive decoding mechanisms that ground them into executable actions conditioned on the current
258 state.

259 **Compositional generalization and latent structure.** Our results show that composition fails when
260 intermediate states or transitions are not covered by the training data. This limitation is closely tied to
261 the current trajectory-level latent representation, which encodes fixed segments of behavior rather
262 than reusable primitives. Structuring the latent space to capture modular or compositional behaviors,
263 rather than monolithic trajectory segments, may enable more flexible recombination and improve
264 generalization beyond the training distribution.

265 **Towards hybrid offline-to-online adaptation.** The limitations identified above reflect a fundamen-
266 tal constraint of purely offline learning: the proposal distribution is inherently bounded by the support
267 of the dataset, and correct behaviors cannot be recovered when they are absent from the offline data,
268 regardless of the quality of the value function. Our results show that the model performs well on seen
269 distributions but struggles to generalize, highlighting the need for mechanisms that actively expand
270 proposal support beyond what offline data provides. Limited online interaction offers a principled
271 way to address this gap, by exploring regions where offline coverage is sparse and augmenting the
272 proposal distribution with behaviors that cannot be inferred from static datasets alone.

273 7 Conclusion

274 We studied diffusion-based offline RL with a focus on proposal–selection decision-making, using
275 ARC-AGI as a controlled analysis environment. Using the SOLAR dataset, we conducted a systematic
276 analysis across single-task, mixed-quality, and compositional settings. Our results show that LDCQ
277 constructs multi-step solutions effectively when the trajectory distribution is well-aligned with the
278 task, but performance degrades when the distribution contains suboptimal trajectories or when the task
279 requires generalization beyond observed ones. Through detailed analysis, we find that this degradation
280 is not primarily due to inaccuracies in value-based selection, but is closely tied to the availability of
281 candidate behaviors, identifying proposal coverage as a key bottleneck in proposal–selection-based
282 offline RL. The SOLAR dataset and diagnostic methodology introduced in this work provide a
283 controlled testbed for future approaches targeting improved proposal coverage, uncertainty-aware
284 selection, and compositional generalization.

285 **References**

- 286 [1] François Chollet. On the Measure of Intelligence. *arXiv:1911.01547*, 2019.
- 287 [2] Scott Fujimoto, David Meger, and Doina Precup. Off-Policy Deep Reinforcement Learning
288 without Exploration. In *ICML*, 2019.
- 289 [3] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing Function Approximation Error
290 in Actor-Critic Methods. In *ICML*, 2018.
- 291 [4] Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and Baining
292 Guo. Efficient Diffusion Training via Min-SNR Weighting Strategy. In *ICCV*, 2023.
- 293 [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In
294 *NeurIPS*, 2020.
- 295 [6] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline Reinforcement Learning with Implicit
296 Q-Learning. In *ICLR*, 2022.
- 297 [7] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-Learning for
298 Offline Reinforcement Learning. In *NeurIPS*, 2020.
- 299 [8] John E Laird, Allen Newell, and Paul S Rosenbloom. Soar: An Architecture for General
300 Intelligence. *Artificial Intelligence*, 33(1):1–64, 1987.
- 301 [9] Hosung Lee, Sejin Kim, Seungpil Lee, Sanha Hwang, Jihwan Lee, Byung-Jun Lee, and Sundong
302 Kim. ARCLE: The Abstraction and Reasoning Corpus Learning Environment for Reinforcement
303 Learning. In *CoLLAs*, 2024.
- 304 [10] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
305 Text-Conditional Image Generation with CLIP Latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 306 [11] Tom Schaul. Prioritized Experience Replay. In *ICLR*, 2016.
- 307 [12] Suyeon Shim, Dohyun Ko, Hosung Lee, Seokki Lee, Doyoon Song, Sanha Hwang, Sejin Kim,
308 and Sundong Kim. O2ARC 3.0: A Platform for Solving and Creating ARC Tasks. In *IJCAI*,
309 2024.
- 310 [13] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. 2020.
- 311 [14] Siddarth Venkatraman, Shivesh Khaitan, Ravi Tej Akella, John Dolan, Jeff Schneider, and Glen
312 Berseth. Reasoning with Latent Diffusion in Offline Reinforcement Learning. In *ICLR*, 2024.

313 **A Training Details**

314 **Training Latent Encoder and Policy Decoder** The first stage in training with LDCQ is to train a
315 β -VAE that learns latent representations. In this stage, the β -VAE learns how actions are executed over
316 multiple steps to change the state. With H -horizon latents, it becomes easier to capture longer-term
317 changes in the state. We use SOLAR as the training dataset \mathcal{D} , which contains H -length segmented
318 trajectories τ_t . Each τ_t consists of state sequences $\mathbf{s}_{t:t+H} = [\mathbf{s}_t, \dots, \mathbf{s}_{t+H-1}]$ and action sequences
319 $\mathbf{a}_{t:t+H} = [\mathbf{a}_t, \dots, \mathbf{a}_{t+H-1}]$, along with additional information such as demonstration examples. As
320 shown in Figure 5a, during the β -VAE training stage, the encoder q_ϕ is trained to encode τ_t into the
321 latent representation \mathbf{z}_t , and the low-level policy decoder π_θ is trained to decode actions based on
322 the given state and latent. For example, given the latent \mathbf{z}_t and a state from the segment trajectory,
323 \mathbf{s}_{t+h} where $h \in [0, H)$, the policy decoder decodes the action \mathbf{a}_{t+h} for \mathbf{s}_{t+h} . The β -VAE is trained
324 by maximizing the evidence lower bound (ELBO), minimizing the loss in Eq. 5. The loss consists
325 of the reconstruction loss from the low-level policy decoder and the KL divergence between the
326 approximate posterior $q_\phi(\mathbf{z}_t|\tau_t)$ and the prior $p_\omega(\mathbf{z}_t|\mathbf{s}_t)$.

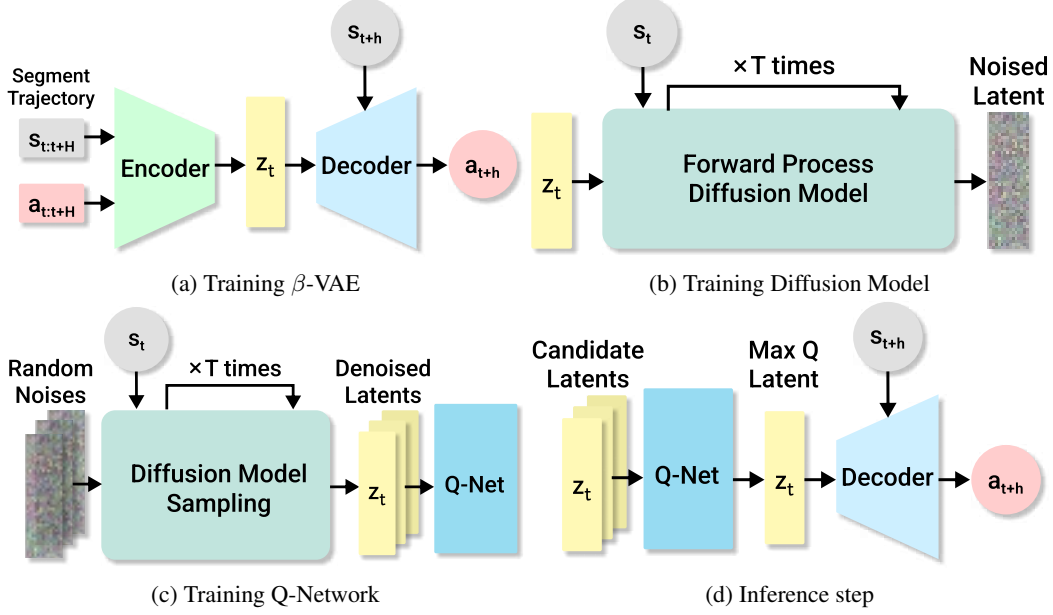


Figure 5: (a)–(c) Training stages of LDCQ. (a) Training a β -VAE with an encoder that encodes H -horizon segment trajectories into latents z_t , and a policy decoder that decodes actions based on z_t and state s_{t+h} where $h \in [0, H)$ contained in the latent. (b) Training a diffusion model based on z_t and the s_t . (c) Training a Q-network using latents sampled through the diffusion model. (d) LDCQ inference step at s_{t+h} . Possible latents at s_t are sampled through the diffusion model, and the agent executes actions resulting from decoding the latent with the highest Q-value.

$$\mathcal{L}_{\text{VAE}}(\theta, \phi, \omega) = -\mathbb{E}_{\tau_t \sim \mathcal{D}} \left[\mathbb{E}_{q_\phi(z_t | \tau_t)} \left[\sum_{l=t}^{t+H-1} \log \pi_\theta(a_l | s_l, z_t) \right] - \beta \cdot D_{KL} \left(q_\phi(z_t | \tau_t) \parallel p_\omega(z_t | s_t) \right) \right] \quad (5)$$

327 **Training Latent Diffusion Model** In the second stage, latent diffusion model is trained to generate
328 latents based on the latent representations encoded by the β -VAE. The training data consists of
329 (s_t, z_t) pairs, which are used to train a conditional latent diffusion model $p_\psi(z_t | s_t)$ by learning
330 the denoising function $\mu_\psi(z_t^j, s_t, j)$, where $j \in [0, T]$ is diffusion timestep. This allows the model
331 to capture the distribution of trajectory latents conditioned on s_t . $q(z_t^j | z_t^0)$ denotes the forward
332 Gaussian diffusion process that noising the original data. Following previous research [10, 14], we
333 predict the original latent rather than the noise, balancing the loss across diffusion timesteps using the
334 Min-SNR- γ strategy [4]. The loss function used to train the diffusion model is shown in Eq. 6. Here,
335 z_t^j , $j \in [0, T]$ represents noised latent on j -th diffusion time step, when $j = 0$ then $z_t^0 = z_t$ and z_t^T
336 is Gaussian noise.

$$\mathcal{L}(\psi) = \mathbb{E}_{j \sim [1, T], \tau_H \sim \mathcal{D}, \substack{z_t \sim q_\phi(z_t | \tau_t), \\ z_t^j \sim q(z_t^j | z_t^0)}} \left[\min\{\text{SNR}(j), \gamma\} \cdot \|z_t^0 - \mu_\psi(z_t^j, s_t, j)\|^2 \right] \quad (6)$$

337 **Training Q-Network** Finally, the latent vectors sampled by the latent diffusion model are used
338 for Q-learning. For latent sampling, while the original LDCQ framework uses DDPM sampling [5],
339 we instead adopt DDIM sampling [13] to accelerate latent generation with fewer denoising steps.
340 The trained diffusion model samples latents by denoising random noise conditioned on the state
341 information s_t . We use the data consisting of $(s_t, z_t, r_{t:t+H}, s_{t+H})$ for training Q-network, where

342 $r_{t:t+H} = \sum_{l=t}^{t+H-1} \gamma^l r_l$ denotes the discounted sum of rewards. Here, DDIM sampling is used to
 343 sample z_{t+H} conditioned on s_{t+H} . For Q-learning, we use Clipped Double Q-learning [3] as shown
 344 in Eq. 7 with Prioritized Experience Replay buffer [11] to improve learning stability and mitigate
 345 overestimation. The trained Q-network $Q(s_t, z_t)$ evaluates the expected return of performing various
 346 H -length actions, with z_t sampled based on s_t . This allows the network to efficiently calculate the
 347 value of actions over H -steps to estimate future returns. Additionally, since ARC-AGI tasks involve
 348 inferring analogies from demonstration pairs, the embedded representation of the demonstration pair,
 349 p_{emb} , is also used in the Q-function calculation.

$$Q(s_t, z_t, p_{emb}) \leftarrow \left(r_{t:t+H} + \gamma^H Q\left(s_{t+H}, \underset{z \sim p_\psi(z_{t+H}|s_{t+H})}{\operatorname{argmax}} Q(s_{t+H}, z, p_{emb}), p_{emb}\right) \right) \quad (7)$$

350 A.1 Hyperparameters

351 We used a horizon length of 5 for encoding skill latents, meaning the model plans and evaluates
 352 actions over a five-step lookahead.

353 We trained the diffusion model with 500 diffusion steps. If the number of diffusion steps is too small,
 354 it can lead to high variance in the sampling process, potentially causing errors during the decoding of
 355 operations or selections in ARCLE. To minimize these errors, we set the number of diffusion steps to
 356 500, ensuring more accurate operation and selection decoding from the sampled latents.

357 We set the discount factor to 0.5 to ensure the model appropriately balances immediate and future
 358 rewards. Since the total steps required to reach the correct answer in ARCLE are usually fewer than
 359 20, a high discount factor could cause the agent to struggle in distinguishing between submitting at
 360 the correct state and continuing with additional steps, which could lead to episode failure.

361 The hyperparameters that we used for training three stages of LDCQ are shown in Tables 8, 9 and 10.

Table 8: Hyperparameters for training β -VAE

Parameter	Value
Learning rate	5e-5
Batch size	128
Epochs	400
Horizon (H)	5
Latent dimension (z)	256
KL loss ratio (β)	0.1
Hidden layer dimension	512

Table 9: Hyperparameters for training latent diffusion model

Parameter	Value
Learning rate	1e-4
Batch size	64
Epochs	400
Diffusion steps (T)	500
Variance schedule	linear
Sampling algorithm	DDIM
γ (For Min-SNR- γ weighting)	5

362 A.2 Hardware

363 We used an NVIDIA A100-SXM4-40GB GPU to train the model. Training the β -VAE took about 7
 364 hours, while training the diffusion model and Q-network each took around 6 to 10 hours.

Table 10: Hyperparameters for training DQN

Parameter	Value
Learning rate	5e-4
Batch size	64
Discount factor (γ)	0.5
Target net update rate (ρ)	0.995
PER buffer α	0.7
PER buffer β range	[0.3, 1)
PER buffer β increment	+0.03 per 2,000 steps
Diffusion samples for batch argmax	100

365 B Details of SOLAR-Generator

366 B.1 Design Rationale for SOLAR

367 We develop SOLAR as the trajectory-based dataset for this study for three reasons.

368 **Controlled mixed-quality trajectory generation based on realistic human error patterns.**

369 SOLAR-Generator synthesizes expert trajectories based on ARCTraj, a dataset of real human task-
 370 solving trajectories collected via O2ARC 3.0 [12]. From this data collection process, we observed that
 371 humans typically follow a correct trajectory up to an intermediate step and then suddenly deviate with
 372 an incorrect action. SOLAR-Generator models this error pattern by branching from expert trajectories
 373 at random intermediate steps and executing plausible but incorrect actions. Importantly, the injected
 374 actions are not uniformly random; instead, they are constrained to semantically plausible alternatives,
 375 such as changing only the spatial extent of a selection or choosing a different color from the set of
 376 valid colors. This design produces suboptimal trajectories that reflect realistic human mistakes rather
 377 than arbitrary noise.

378 **Flexible multi-task training and evaluation.** SOLAR provides approximately 20 task variants
 379 derived from ARC-AGI, and SOLAR-Generator can synthesize trajectories for these tasks on demand.
 380 Our experiments use 5–10 tasks, enabling controlled multi-task training and evaluation of interference
 381 effects.

382 **Explicit action-level state transitions for stage-wise analysis.** SOLAR represents task-solving
 383 trajectories at the action level, and we integrate ARCLE [9] to obtain explicit state transitions. This
 384 setup allows us to attribute performance degradation to specific stages (proposal, evaluation, selection,
 385 application) of the decision pipeline.

386 B.2 Operations in SOLAR

387 The operations from 0 to 34 are identical to those used in ARCLE [9]. Since `Submit` is an operation
 388 that receives a reward, it should only be used when the state is considered correct and not excessively.
 389 Due to LDCQ’s fixed horizon, and to ensure that the agent only uses `Submit` when the state is
 390 definitively correct, we added a `None` operation that fills all subsequent states after `Submit` with the
 391 11th color (10), which does not exist in the original ARC (0–9). In other words, during training, the
 392 `None` action emphasizes that the episode ends after `Submit`.

393 B.3 Grid Maker

394 To generate SOLAR, we develop a SOLAR-Generator that synthesizes a large volume of data for a
 395 given rule. Grid Maker is a hard-coded program specific to each task. Grid Maker contains the rules
 396 for synthesizing demonstration examples and test examples, and the synthesized solution action path
 397 consists of operations and selections. In Grid Maker, data is formatted to be compatible with ARCLE.
 398 The Grid Maker constructs analogies with the same problem semantics but with various attributes,
 399 such as the shape, color, size, and position of objects. SOLAR-Generator can generate intermediate
 400 trajectories by interacting with ARCLE. The SOLAR-Generator algorithm is designed to augment
 401 specific tasks using the Grid Maker, which can be primarily divided into three parts.



Figure 6: All operations compatible with SOLAR, 0–34 operations follow ARCLE, and only in SOLAR, 35 (None) is for terminated episode. This indicates that the episode ends after Submit.

402 Grid Maker was built as a data loader used in ARCLE. In the original ARCLE environment, there was
 403 no need to load operations and selections. Only the grid was loaded with the original ARC. To change
 404 this structure, the entire environment would need to be recreated. Instead, operations and selections
 405 are now loaded from the data loader’s description, thereby preserving the original environment.
 406 Therefore, the process of creating input-output examples and generating action sequences works
 407 within a single file.

408 **Specifying Common Parts** Each task in the ARC dataset usually contains 3 demonstration exam-
 409 ples, with common elements observed across these pairs. In the common parts, attributes such as
 410 color, task type, and object presence are determined by random values before pair generation.

411 **Synthesizing Examples** In the example synthesis phase, the input of the original task is augmented
 412 in a way that ensures diversity while preserving the integrity of the problem-solving method. A random
 413 input grid is generated under conditions that satisfy the task’s required analogy. A solution grid is
 414 created using a hard-coded algorithm. For tasks involving pattern-based problems, as experimented in
 415 the paper, selections are made to fit the grid size, and various operations are executed either randomly
 416 or in a predetermined order. For object-based problems, the solution grid is generated by an algorithm
 417 that identifies the required objects in the input grid and processes them in accordance with the task
 418 requirements.

419 **Converting to ARCLE Trajectories** This stage involves the creation of an ARCLE-based trajectory
 420 that meticulously adheres to the problem-solving schema of the synthesized examples. The entire
 421 process is implemented using a hard-coded algorithm. During the example synthesis process, object
 422 locations may already be known, or they can be identified using a search algorithm. The information
 423 obtained is used to make appropriate selections, and the trajectory is converted into an ARCLE
 424 trajectory via an algorithm that yields the correct solution.

425 If all steps are properly coded, it is possible to generate the operations and selections that yield the
 426 correct solution for any randomly generated input grid. These are then fed into ARCLE to obtain
 427 intermediate states, rewards, and other information, and to verify whether the correct result is reached.
 428 Once steps 1) to 3) are correctly implemented, SOLAR-Generator can continuously and automatically
 429 generate as much data for the given task as the user desires, using the Grid Maker.

Algorithm 1: SOLAR-Generator

Input: Task set T , grid size (H, W) , samples N , examples E , horizon H_{seg} **Output:** Training and test trajectory datasets $\{\tau^{train}, \tau^{test}\}$

```
for  $task \in T$  do
   $\mathcal{H}_{test} \leftarrow \emptyset$ ; // Test set hash tracking
  // Generate test data first for independent evaluation
2  for  $type \in \{test, train\}$  do
    // Generate grids using GridMaker
3     $\mathcal{D}_s \leftarrow \text{GridMaker}(task, (H, W), N_{type}, E, seed_{type})$ 
4    for  $(ex_{in}, ex_{out}, pr_{in}, pr_{out}, desc) \in \mathcal{D}_s$  do
      // Duplicate detection and filtering
5       $hash \leftarrow \text{SHA256}(pr_{in} || pr_{out})$ 
6      if duplicate or ( $type = train$  and  $hash \in \mathcal{H}_{test}$ ) then
7        continue
8      end
      // Execute action sequence in ARCLE environment
9       $s_0 \leftarrow \text{ARCLE.reset}(pr_{in})$ 
10     for  $(op_t, sel_t) \in (desc.ops, desc.sels)$  do
11        $s_{t+1}, r_t, done_t \leftarrow \text{ARCLE.step}(s_t, a_t)$ 
12        $\tau \leftarrow \tau \cup \{s_t, a_t, r_t, done_t\}$ 
13     end
      // Validate expert trajectories
14     if  $is\_expert$  and  $s_t.grid \neq pr_{out}$  then
15       continue
16     end
      // Save trajectory (whole and segmented)
17     Save  $\tau$  and  $\{\tau_i^{seg}\}$  with horizon  $H_{seg}$ 
18   end
19   if  $type = test$  then
20      $\mathcal{H}_{test} \leftarrow$  current hashes
21   end
22 end
end
return  $\{\tau^{train}, \tau^{test}\}$ 
```
