

# QHyer: Q-conditioned Hybrid Attention-Mamba Transformer for Offline Goal-conditioned RL

Xing Lei<sup>1</sup> Jincheng Wang<sup>2</sup> Xuetao Zhang<sup>1\*</sup> Donglin Wang<sup>3</sup>  
<sup>1</sup>Xian Jiaotong University <sup>2</sup>University College London <sup>3</sup>Westlake University  
 Correspondence: xuzetaozhang@xjtu.edu.cn



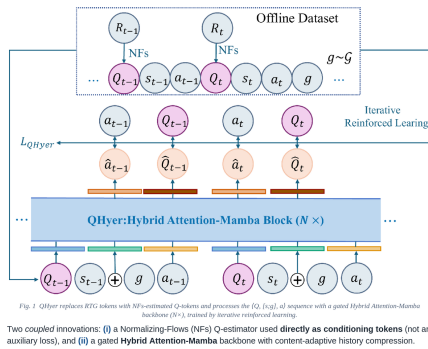
## 1 Background & Motivation

**Offline GCRL.** Learn a goal-reaching policy  $\pi(a|s, g)$  purely from a static dataset, with no environment interaction — a practical paradigm when online exploration is costly or unsafe.  
**The non-Markovian challenge.** Many real datasets are gathered by *history-dependent* behavior policies  $\{P_k | s_t, h_t\}$ . Value-based methods that assume Markovian transitions break down, while Decision Transformers (DT) fit such data via sequence modeling.  
**Yet DT has two limitations for GCRL:**  
 1. **Return-to-go (RTG) is trajectory-dependent** — under sparse rewards it gives no signal for stitching promising segments out of failed trajectories.  
 2. **Fixed-window local extractors** (convolution in LSDT / DMixer) cannot adapt the effective memory to variable-length temporal dependencies.

## 2 Limitation 1 — RTG Fails Under Sparse Rewards

**Root cause.** RTG answers "did this trajectory succeed?" not "how good is this state for the goal?" The same state gets RTG=1 on a success and RTG=0 on a failure, making cross-trajectory comparison — and stitching — impossible.  
**Our fix.** Replace RTG with a goal-reaching Q-value  $Q^\beta(s, a, g) = \mathbb{E}_{g \sim \mathcal{G}} [R_t | s_t, a_t, g]$ , measured per state-action and independent of which trajectory it came from — so high-value segments from failed demos can be identified and composed toward the goal.

## 3 QHyer Architecture



Two coupled innovations: (i) a Normalizing-Flows (NFs) Q-estimator used directly as conditioning tokens (not an auxiliary loss), and (ii) a gated Hybrid Attention-Mamba backbone with content-adaptive history compression.

## 4 NFs-based Q-Value Conditioning

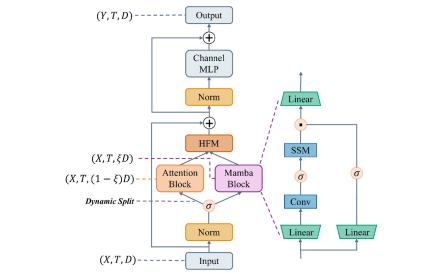
The goal-conditioned Q equals the discounted state-occupancy — the probability of reaching g:  

$$Q^\pi(s, a, g) = p_+^\pi(s_+ = g | s_0 = s, a_0 = a)$$
  
 Coupling-based NFs estimate it with an exact, properly-normalized log-density:  

$$Q_\theta^\beta(s, a, g) = \log p_\theta(f_\theta(g; z)) + \log \left| \det \frac{\partial f_\theta(g; z)}{\partial g} \right|$$
  
**Why NFs?** The transformer reads Q-tokens spanning *multiple* goals, so goal-independent normalization is essential. CVAE (ELBO ggg), contrastive RL (goal-dependent offset) and diffusion (ODE / trace-estimator variance) cannot provide it — coupling NFs uniquely can.  
 Explicit regression then extracts the in-distribution maximum Q ( $Q > 0.5$  penalizes underestimation;  $t - 1$  — optimal Q):  

$$\mathcal{L}_Q = \mathbb{E}_{(s, a, g) \sim \mathcal{D}} [L_\tau^2(Q_\theta^\beta(s_t, a_t, g) - \hat{Q}_\phi(s_t, g))]$$

## 5 Hybrid Attention-Mamba Block



Effective memory is set by **input-dependent SSM dynamics**, not a fixed kernel:  

$$\hat{A}_t = \exp(\Delta_t \cdot A), \quad \Delta_t = \text{softplus}(\text{Linear}_\Delta(x'_t))$$
  
 Small  $\Delta_t \rightarrow \hat{A}_t \approx 1$  (long memory, play); large  $\Delta_t \rightarrow \hat{A}_t \approx 0$  (short memory, noisy). A learnable gate fuses both branches:  $y = \alpha \cdot y_{\text{attn}} + (1 - \alpha) \cdot y_{\text{mamba}}$ .

## 6 Content-Adaptive Memory

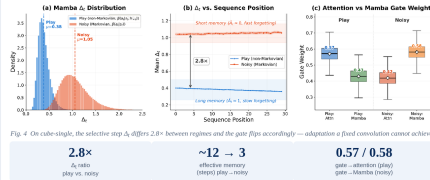


Fig. 4 On cube-stitch, the selective step  $\Delta_t$  differs 2.8<sup>x</sup> between registers and the gate ffn accordingly — adaptation a fixed convolution cannot achieve.

## 7 State-of-the-Art Results

SOTA on both non-Markovian (OGBench play, D4RL Maze) and Markovian (OGBench noisy) data.

Benchmark / setting	Metric	Best prior	QHyer
OGBench cube-play (total)	SR%	111 <small>100%</small>	152
OGBench puzzle-play (total)	SR%	147 <small>100%</small>	169
OGBench visual-scene-play	SR%	62 <small>0%</small>	96
D4RL AntMaze-v2 (total)	score	378 <small>0%</small>	483
D4RL Maze2D (total)	score	277 <small>0%</small>	292

Averaged over 8 seeds (OGBench) / 5 seeds (D4RL). Largest gains on big mazes that demand extensive stitching.

## 8 Ablations: Both Innovations Matter

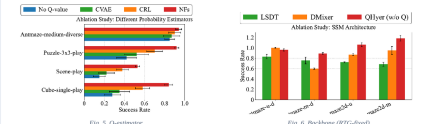


Fig. 5 Ablation study (D4RL-Play)

## 9 Conclusion

- First sequence-modelling framework to **jointly** resolve trajectory-dependent conditioning and fixed-window memory for non-Markovian offline GCRL.
- NFs Q-tokens (exact density + explicit regression) enable trajectory stitching from failed demonstrations — no bootstrapping, no policy projection.
- Hybrid Attention-Mamba delivers content-adaptive memory, giving SOTA across OGBench manipulation and D4RL navigation.